# Predicting age of acquisition for children's early vocabulary in five languages using language model surprisal

Eva Portelance<sup>1</sup>\*, Yuguang Duan<sup>2</sup>, Michael C. Frank<sup>3</sup>, Gary Lupyan<sup>2</sup>

<sup>1</sup>Department of Linguistics, McGill University

<sup>2</sup>Department of Psychology, University of Wisconsin-Madison

<sup>3</sup>Department of Psychology, Stanford University

July 2023

#### **Abstract**

What makes a word easy to learn? Early-learned words are frequent and tend to name concrete referents. But words typically do not occur in isolation. Some words are predictable from their contexts; others less so. Here, we investigate whether predictability relates to when children start producing different words (Age of Acquisition; AoA). We operationalised predictability in terms of a word's surprisal in child directed speech, computed using n-gram and long-short-term-memory (LSTM) language models. Predictability derived from LSTMs was generally a better predictor than predictability derived from n-gram models. Across five languages, average surprisal was positively correlated with the AoA of predicates and function words, but not nouns. Controlling for concreteness and word frequency, more predictable predicates and function words were learned earlier. Differences in predictability between languages were associated with cross-linguistic differences in AoA: the same word (when it was a predicate) was produced earlier in languages where the word was more predictable.

#### 1 Introduction

In the first 2 years of life, children's grammatical knowledge and lexicon grow in tandem (Bates et al., 1994; Brinchmann, Braeken, & Lyster, 2019; Frank, Braginsky, Marchman, & Yurovsky, 2021). In addition, the order in which children acquire their first words show remarkable consistency (Clark, 1993; Tardif et al., 2008; Goodman, Dale, & Li, 2008; Braginsky, Yurovsky, Marchman, & Frank, 2019). For example, the words 'ball', 'car', and 'nose' are, on average, learned almost a year earlier than words like 'drawer', 'green' and 'animal'. This general pattern holds across multiple languages (Braginsky et al., 2019). Modeling when words are acquired can therefore help us understand what factors drive language learning more generally.

One way to study these ordering effects is by attempting to predict a word's age of acquisition (AoA) from lexical properties such as its part of speech, frequency, length, and concreteness (Goodman et al., 2008; Kuperman, Stadthagen-Gonzalez, & Brysbaert, 2012; Braginsky et al., 2019). A common way is to estimate AoA is to use retrospective reports from adults of their own AoA for words (e.g., Kuperman et al., 2012). In contrast, here, AoA is determined using objective and timely parental reports of their children's linguistic productions. Though these methods are related, they are not the same and produce only moderately correlated AoA estimates for words (Łuniewska et al., 2016). Using child-based AoA estimates, studies find that nouns tend to be learned before verbs; more frequent words before less frequent words, shorter words before longer words, and words with concrete referents before those referring to more abstract entities. Controlling for part of speech, frequency tends to be one of the strongest predictors of AoA (Goodman et al., 2008; Roy, Frank, DeCamp, Miller, & Roy, 2015; Braginsky et al., 2019; Fourtassi, Bian, & Frank, 2020). Here, we go beyond using these word-level predictors, by considering the linguistic context in which words appear in speech directed at children.

One such contextual predictor, previously examined by Braginsky et al. (2019), is the mean length of utterance (MLU) in which the word occurs. MLU can be considered as a proxy for syntactic complexity. If word learning

<sup>\*</sup>Corresponding author: eva.portelance@mcgill.ca

is constrained by a child's ability to parse longer utterances (which are often more complex), we should find longer MLUs to be associated with later AoAs. This is precisely what Braginsky et al. (2019) found, though MLU was a significant predictor only for predicates and function words (a result that turns out to be relevant to the current investigation). Others have used contextual diversity – a measure of semantic co-occurrence – as a predictor of AoA (Hills, Maouene, Riordan, & Smith, 2010; Roy et al., 2015; Amatuni & Bergelson, 2017; Stella, Beckage, & Brede, 2017; Fourtassi et al., 2020). Words with higher contextual diversity – those appearing in many different semantic contexts in children's input – tend to be easier to learn and process by children (Rosa, Tapia, & Perea, 2017; Hsiao & Nation, 2018; Pagán, Bird, Hsiao, & Nation, 2020). Like with MLU, this predictor does not show the same effect across lexical categories. Hills et al. (2010) found that higher contextual diversity predicted earlier word acquisition, though this effect appeared to be strongest for predicates and function words when controlling for the effect of word frequency. Such a predictor can be considered a proxy for some semantic factors, approximating the overall semantic richness of children's learning environment. However, contextual diversity abstracts away from the semantic information encoded in specific sequential linguistic contexts – utterances or conversations – experienced by learners (a point we discuss in section 3.1 of this paper), and thus, does not directly measure syntactic complexity.

Here, we examine another contextual factor – the predictability of words across the linguistic contexts in which they appear. A word's predictability takes into account its linguistic context as a whole, including both syntactic and semantic information. It allows us to identify words that might be more or less difficult to learn as a function of the word sequences they tend to appear in.

There are multiple reasons for focusing on predictability. First, it is easy to compute, as we detail below. Second, it is psychologically real – we know that statistical factors such as transition probabilities play important roles in just about every aspect of early language learning (Saffran, 2020). Predictability (operationalized in terms of *surprisal*, see below) has also been shown to be a strong predictor of processing difficulty in adult psycholinguistic experiments (Levy, 2008; Demberg & Keller, 2008; Smith & Levy, 2013). More predictable words (i.e., those with lower surprisal) tend to be easier to process.<sup>1</sup>

Thus, to help us understand the role of sequential predictability on the complexity and growth of children's vocabulary across languages, we propose to study the relation that exists between word predictability and AoA, treating AoA as a proxy for word learning difficulty. We frame our experiments around the following two sets of concrete questions:

- 1. Does a word's predictability across linguistic contexts help explain how difficult a word is to learn beyond previously known predictors? If so, how much sequential context matters? And, is predictability equally important for all words or not and why?
- 2. Second, are the effects of predictability on AoA observable across many languages and linguistic communities or are they isolated to only some? And, do differences in word predictability predict differences in AoA between languages?

To determine a word's predictability, we use computational language models that consider a word's previous sequential linguistic context. Specifically, word predictability can be measured as the average surprisal of a word – its negative log probability in a given context, averaged across all contexts in which it appears – given a *language model* as our probability model.

Language models are sequential predictive models trained to generate linguistic output which are commonly used in natural language processing (NLP). They do so by learning probability distributions over strings of words in a corpus, where the length of these word strings may vary.<sup>2</sup> The additional information contained in these substrings in the form of preceding words is what allows us to consider a word's predictability given previous syntactic and semantic context. In this paper, we consider two types of language models: n-gram models and LSTM language models (Sundermeyer, Schlüter, & Ney, 2012). N-gram models are simple conditional probability models while LSTM models are more sophisticated neural network models.

There has been growing interest in training or evaluating neural network language models on corpora that resemble the linguistic input of children more than the standard newspaper, Wikipedia, or web corpora used in NLP. This interest is motivated by two main goals: first, determining how neural network models come to learn languages, and second, using these models as tools to better understand language learning in children. In one model-centered investigation,

<sup>&</sup>lt;sup>1</sup>For an overview of empirical evidence supporting the validity of surprisal as a predictor of processing difficulty, see Hale, 2016.

<sup>&</sup>lt;sup>2</sup>Context size can vary depending on the language model, ranging from a single previous token in n-gram models to a bounded ordered list of all previous or subsequent tokens in an LSTM model.

Huebner, Sulem, Cynthia, and Roth (2021) found that training language models on corpora of child-directed utterances can actually help these models perform better on grammatical knowledge evaluations than when trained on similar sized or larger corpora of traditional NLP data, suggesting that child-directed language may help grammar learning. Chang and Bergen (2022) have also proposed to use the average surprisal of words as a proxy for the AoA of words for language models in order to develop a new model evaluation task. They evaluated whether previously known predictors of the AoA of words in children – like frequency, concreteness, MLU, lexical category, and number of characters – also predicted the 'AoA of words' in language models, and found that there are clear differences between the orders in which children and language models acquire words.

Language models have also been proposed as tools to evaluate children's language development. For example, Sagae (2021) suggest using LSTMs trained on children's utterances to quantify children's syntactic development, finding that they perform as well or better than previous metrics.

All of this work – both model-based and child-focused – has been limited to English data. Here, we expand our analyses to cross-linguistic data, considering models trained on five different languages: English, German, French, Spanish, and Mandarin.

Our approach In the rest of this paper, we will take the following approach: for each language, we fit a set of language models on a corpus of child directed utterances and extract the average surprisal of words for which we have AoA estimates in children. We then compare regression models of children's AoA with average surprisal as one key predictor in concert with previous significant predictors, using cross-validation to estimate out-of-sample performance. We present two different modeling methods in two experiments. The first considers the effect of predictability on the AoA of words in each language individually, while the second considers its effect on all languages as a whole. We close this paper by considering how what we have learnt about the relation between word predictability and AoA can inform our understanding about the role predictability plays on both children's vocabulary growth and the complexity of languages as a whole.

### 2 Data

We relied on two types of data: (1) Corpora of child-directed utterances used to train the language models.<sup>3</sup> this These were taken from the CHILDES database (MacWhinney, 2000); (2) Age of Acquisition estimates. These were based on parental reports of children's language use, taken from the Wordbank repository (Frank, Braginsky, Yurovsky, & Marchman, 2016). We go into further detail in the following subsections about both these resources and the data they contain. Importantly, in order for a language to be considered in this cross-linguistic study, there had to be sufficient data in this language in both of these resources. This criterion narrowed down the list of languages we could consider to English, German, French, Spanish, and Mandarin. English was by far the most represented language in CHILDES, but there were still enough data for the other languages to be able to fit our language models (see Table 1 for the amounts of data available in each language for both CHILDES and Wordbank).

#### 2.1 CHILDES and child-directed utterances

The CHILDES database (MacWhinney, 2000) is a repository of child language data, containing text transcripts of child-caregiver interactions as well as video and sound recordings of some of these interactions. The data comes from many different studies conducted during the past 60 years, spanning multiple languages and countries. For the most part, the children in these studies range in age between nine months old and five and a half years old.

For this paper, we only considered text transcription data and no other modalities. For each of the five languages considered (English, German, French, Spanish, Mandarin), we collected all of the available transcripts across all corpora available through the childes-db API (Sanchez et al., 2018) in July 2021. We then removed all utterances spoken by the target child, leaving only the utterances said to the child or around the child. These utterances can be considered as an estimate of the linguistic input the children in these transcript have access to. We combined all of these child-directed utterances<sup>4</sup> into a corpus used to fit the language models presented in the next section and

<sup>&</sup>lt;sup>3</sup>All of the data, models, and experiment code presented in this paper are publicly available at www.qithub.com/evaportelance/multilingual-aoa-prediction.

<sup>&</sup>lt;sup>4</sup>We will use the term child-directed somewhat loosely here such that it may also refer to utterances that were directed to other adults or children present, but that the target child could still hear.

	Table 1:	Amount of data	available in	CHILDES	and Wordbar	nk by language
--	----------	----------------	--------------	---------	-------------	----------------

Language	Child-directed to- kens (CHILDES)	Dialect	Vocabulary reports (Wordbank)
English	25,659,263	American	7,955
		British	23,129
		Australian	1,497
German	5,663,294		1,181
French	3,183,037	European	863
		Quebecois	1,364
Spanish	2,267,707	European	1,005
		Mexican	1,934
Mandarin	2,369,896	Beijing	1,938
		Taiwanese	2,654

to calculate the relative frequency of words for the regression models presented in our experiments. The result was 5 corpora of child-directed utterances, one for each language. Unlike the Wordbank database (Frank et al., 2016) presented in the next subsection, childes-db does not explicitly distinguish data based on dialectal varieties of each language, so we use the same aggregated data for each language across all varieties when fitting our models.

#### 2.2 Wordbank and age of acquisition estimates

The Wordbank database (Frank et al., 2016) is a repository of parental reports about their children's vocabularies – essentially, a checklist of words where parents can check off words their child produces or understands. Most of these reports are versions of the MacArthur-Bates Communicative Development Inventories (CDI) (Fenson et al., 1993). The database is a collection of reports originating from different studies that were conducted across the world. These studies and the vocabulary checklists they use are dialect-specific, so for each of five languages we consider in this study, we collected data from all available dialects. We did not combine the data from different dialects into single languages as each dialect contains different word lists on their reports, leaving fewer words at their intersection.

Our predictive target is the age at which a word is acquired. Since not all children learn a given word at the same time, we instead follow prior work in quantifying AoA as the age at which 50% of children are reported to produce a word on the CDI (Goodman et al., 2008).<sup>5</sup>. There are a number of methods to estimate this 50% point from a group of binary responses for children of different ages. The simplest method is to determine the youngest age group at which the empirical proportion of children producing the word is > 50%, but this approach has several shortcomings. If words are very hard or very easy to learn, then it is possible that for the covered age range some words never reach the 50% point (e.g., *beside*), or have already surpassed the 50% point (e.g., *Mommy*) for even the youngest children. Such words would have to be discarded if we were to use this method. Another issue is that this approach is susceptible to bias AoA estimates towards ages for which more CDI instruments were available since the number of observations at each age is not equal (i.e., there may be more CDI instruments from 24-month-olds than with 20-month-olds in the dataset, but this density shouldn't lead to more words being acquired at exactly 24 months). For these reasons, we used Bayesian generalized linear models predicting acquisition as a function of age to estimate the AoA for each word, following the method suggested in Frank et al., 2021<sup>6</sup>.

From the reports available in each language, we narrowed down the list of items used in our experiments to all single word items on the forms that were classified as either nouns, predicates (verbs and adjectives), or function words (closed class words like pronouns, prepositions, question words, connectives, determiners). We excluded items that were multi-word expressions (e.g. 'all gone') or that were classified as being part of the "other" lexical category, which included animal sounds, onomatopoeia, and other non-word expressions. Words were also excluded if they weren't in the five thousand most common words in each language in our corpora of child-directed utterances from CHILDES,

<sup>&</sup>lt;sup>5</sup>We chose to use these AoA estimates from parental reports over those of Kuperman et al., 2012 which exist for a much larger vocabulary, because the latter are based on adult estimates of their own AoA, rather than timely reports of children's AoA. Additionally, Kuperman et al. (2012) report word comprehension AoA estimates, while the estimates we use here are word production AoA estimates.

<sup>&</sup>lt;sup>6</sup>For a more detailed description of this method see Appendix E of Frank et al., 2021.

Language	Number	Nouns	Predicates	<b>Function words</b>
	of items			
English (American)	563	301 (53%)	165 (29%)	97 (17%)
English (British)	333	196 (59%)	100 (30%)	37 (11%)
English (Australian)	368	230 (63%)	138 (37%)	0 (0%)
German	361	189 (52%)	108 (30%)	64 (18%)
French (European)	394	222 (56%)	119 (30%)	53 (13%)
French (Quebecois)	468	248 (53%)	159 (34%)	61 (13%)
Spanish (European)	430	212 (49%)	121 (28%)	97 (23%)
Spanish (Mexican)	339	181 (53%)	90 (27%)	68 (20%)
Mandarin (Beijing)	500	269 (54%)	185 (37%)	46 (9%)
Mandarin (Taiwanese)	412	257 (62%)	119 (29%)	36(8%)

as this was the vocabulary size used for the language models (described in the next section). Table 2 contains the exact number of items taken from Wordbank that we considered for each language, as well as their breakdown by lexical category.

# 3 Language models and predictability

To determine the predictability of words in the child-directed utterance corpora described above, we use language models as our probability models. Language models define probability distributions over subsequent words in a given context. Here, we consider the predictability of words solely based on linguistic contextual information. Specifically, we define the overall predictability of a word,  $w_i$ , as its average surprisal, or average negative log probability, across all its contexts of use, C (eq. 1).

$$predictability(w_i) = \sum_{C:w_i \in C} -\log P(w_i \mid w_1, ..., w_{i-1}) \times \frac{1}{|C|}$$
 (1)

where  $w_1, ..., w_{i-1}$  is a sequence of words of bounded length representing the preceding linguistic context.

There are many different types of language models. They vary in terms of the context sizes they consider, in how they represent words, and in how they come to calculate the overall probability of a word. As our definition in eq. 1 suggests, we only consider language models that take into account preceding linguistic context (and not following context) and do so in an incremental order. Specifically, the experiments that follow will contain average surprisal values obtained from two types of language models, n-gram models and LSTM models.

#### 3.1 N-gram language models

N-gram models are basic language models that consider contexts of sequence length n, such that a bi-gram model keeps track of two word sequences and the contextual probability of a word is based on a single preceding word, and a tri-gram keep track of 3 word sequences and contextual probability of a word is based on the two preceding words. Thus, given our formula in eq. 1, we simply need to replace i by n to determine the predictability measure, or average surprisal, of a word for a given n-gram model.

In an n-gram probability model, the probability of a word  $w_n$  in a given context,  $w_1, ..., w_{n-1}$ , or  $P_{n-gram}(w_n \mid w_1, ..., w_{n-1})$ , is simply its normalized count across all words that follow this context in the corpus. For example, if our probability model was a tri-gram model and we wanted to determine the probability of the word 'bird' in the context 'is that a bird', we would take  $P_{tri-gram}(\text{bird} \mid \text{that a})$ , which in practice is equal to : count(that a bird)/count(that a), where count() returns the number of instances of an expression in a corpus. In this study, we used four different n-gram models: uni-grams, bi-grams, tri-grams, and four-grams. Note that uni-gram models only track single word contexts, in other words they represent the normalized frequency counts of words, so the average surprisal of a word for a uni-gram model is equivalent to its negative log frequency. Uni-gram surprisal is also in effect equivalent to contextual

diversity. Though their definitions are different, when calculated on a large enough corpus, these two metrics converge, as we show in Appendix C, reaching a correlation score of -0.97 between uni-gram surprisal and log-transformed contextual diversity. For all intents and purposes, uni-gram surprisal, or the predictability of words irrespective of previous linguistic context, can be also be interpreted as contextual diversity in the experiments which follow.<sup>7</sup>

One downside of n-gram models is that the context size is fixed across the whole probability model. This means that though some words may be better predicted from a single preceding word while others may be better predicted by two preceding words, we can only consider one of these context sizes at a time. LSTM models can help us get around this problem.

#### 3.2 LSTM language models

Recurrent neural networks (RNNs) which use long-short term memory gating unit layers (Hochreiter & Schmidhuber, 1997), commonly known as LSTMs, are neural networks that can be trained on sequential data, such as sentences, up to some bounded maximum length n. These models can be used for language modeling (Sundermeyer et al., 2012) and have become a staple baseline that continue to be used in NLP because of their useful analytical properties, even though more recent model architectures outperform them (Vaswani et al., 2017). Furthermore, regular RNNs have previously been proposed as cognitive models for language learning (Elman, 1990, 1993; Christiansen, Allen, & Seidenberg, 1998), however, these earlier models were computationally limited and could be used only with small schematic datasets; in contrast, LSTMs can be applied to larger datasets. Thus, LSTM language models lend themselves well to our analyses.

LSTM language models process utterances incrementally and make use of nested layers of hidden units to learn abstract representations that can predict sequential dependencies between words across a range of dependency lengths (Linzen, Dupoux, & Goldberg, 2016). LSTM neural units use a gating system that allows them to 'forget' some of the previous states while 'remembering' others, thus learning to prioritize some dependencies in a sequence over others at each state. So, unlike n-gram models, when determining the predictability, or average surprisal, (eq. 1) of a word  $w_i$  for these models, the preceding contexts,  $w_1, ..., w_{i-1}$  in C, can vary in length, usually representing all of the previous words in the sentence. Further, the probability of a word in context can weigh the importance of preceding words differently based on the information encoded in the model's different layers. The added richness of these representations may lead to a better probability model overall.

For the experiments that follow, we use a two-layered LSTM language model (Figure 1). The model has randomly initialized 100-dimensional word embeddings as its input layer, which are updated during learning. Hidden states encode information about the preceding context. At each time-step, the current word embedding  $w_t$  and the hidden state from the previous time-step  $h_{t-1}^1$  are passed through a transformation function, resulting in a new hidden state  $h_t^1$ . This hidden state  $h_t^1$  and the hidden state from the previous time-step in the second hidden layer  $h_{t-1}^2$  are then also passed through a transformation function, resulting in a new hidden state  $h_t^2$ . This final hidden state is then resized through a linear layer to the size of our vocabulary before going through a softmax transformation to produce the output – a distribution over the whole vocabulary W representing a prediction about the upcoming word. We use a vocabulary size of 5,000, representing the most frequent words, because we found that including the 5,000 most common words usually resulted in the inclusion of almost all the words we had on our AoA word lists for all languages. Figure 1 shows how the probability of a word given its preceding context  $P_{LSTM}$  ( $w_i \mid w_1, ..., w_{i-1}$ ) can be extracted from the model.

The model's objective during training is to maximize the likelihood of the next word at each step – in other words, the model updates its parameters to minimize the surprisal of words in context. We performed cross-validation tests to find the hyperparameter settings for the LSTM language models that best minimized overall word surprisal. The hyperparameters tested were the vocabulary size (2000, 4000, or 5000), the word embedding size (100, 150), the hidden state dimension size (100, 150), the batch size (128, 256, 512), and the number of epochs (up to 50). We found that the optimal parameter combination was a vocabulary size of 5000, 100 dimensional word embeddings, a 100 hidden state size, a batch size of 256, and about 20 epochs of training.

We trained the models on all of the child-directed utterances for each language since the models were to be used as probability models and not predictive models. We were therefore not concerned with overfitting to the training data. Utterances were shuffled at each epoch of training. For further details on the model implementation see Appendix A.

<sup>&</sup>lt;sup>7</sup>Though contextual diversity is strongly correlated with uni-gram surprisal, importantly, it is not related to higher-order surprisal metrics which take into account contextualized sequential predictability, like LSTM average surprisal, as shown in Appendix C.

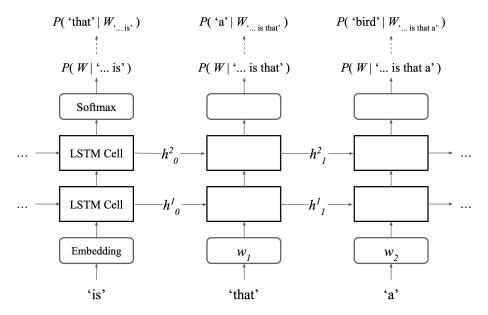


Figure 1: The LSTM model architecture incrementally processing the utterance 'is that a bird'. At each step, the model generates a conditional probability distribution over the vocabulary *W* representing the likelihood of being the next word. During training, the model updates its weights to maximize the probability of the actual next word. We can estimate the probability of words in context by retrieving their probability from the model output at each state.

# 4 Experiment 1: The role of word predictability beyond log frequency

In this first experiment, we consider how word predictability beyond frequency predicts the AoA of words. Previous work (Goodman et al., 2008; Kuperman et al., 2012; Roy et al., 2015; Braginsky et al., 2019; Fourtassi et al., 2020) has found that log frequency – or, equivalently, uni-gram surprisal – is an important predictor of the AoA of words; here, we evaluate the explanatory power of larger context sizes by using their residualised effect beyond uni-gram surprisal as predictors. We compare models with different versions of average surprisal, obtained using different language models: bi-gram, tri-gram, four-gram, or LSTM average surprisals. We do so using leave-one-out (LOO) cross-validation.

## 4.1 Predictors

There are two main types of predictors considered in our models. First, we consider several methods for computing average surprisal using language models conditioned on different sizes and types of previous linguistic context. Second, we include other predictors that have been found to be informative in previous work: concreteness and lexical category (Goodman et al., 2008; Kuperman et al., 2012; Braginsky et al., 2019). All predictors are scaled by centering their mean at zero and dividing by their standard deviation so that their magnitudes can be compared.

**uni-gram surprisal** is computed as the negative logarithm of frequency. Log frequency has been found to explain substantial variance in AoA in previous work.

**Residualised n-gram surprisal** represents the residual variance left after fitting a linear model which predicts n-gram average surprisal as a function of uni-gram surprisal, n-gram average surprisal  $\sim$  uni-gram surprisal. n-gram average surprisal is the predictability of a word given all contexts of size n in which it appears in the n-th position. We compare the average surprisal of bi-gram, tri-gram, and four-gram language models. If an item had multiple word forms associated to it on the parental report instrument (e.g. 'inside/in'), we used the weighted mean average surprisal across all forms — in other words, if one form was overall more frequent, then it was weighted it

accordingly.8

Residualised LSTM surprisal represents the residual variance left after fitting a linear model which predicts LSTM average surprisal as a function of uni-gram surprisal, LSTM average surprisal  $\sim$  uni-gram surprisal. LSTM average surprisal is calculated using the LSTM language models described above. We compute the average surprisal of each word across all of the child-directed utterances available in each language. We trained three LSTM language models in each language using different random seeds and then used the mean average surprisal across these three runs as our measure of word predictability in each language. As with n-gram surprisal, we used the weighted mean average surprisal across all word forms.

**Concreteness** is a rating score ranging between 1 and 5 for each word representing some measure along the abstract to concrete scale. These scores are taken from (Brysbaert, Warriner, & Kuperman, 2014). In order to obtain them for languages other than English, Wordbank data associates each item to a 'unilemma', which is an equivalent English concept across all languages for that item in the database. The concreteness score for the equivalent English concept was then used for each word. This practice follows previous work (Braginsky et al., 2019).

**Lexical category** is included via contrast coding of three lexical categories: nouns (common nouns), predicates (verbs and adjectives), and function words (closed-class words) following Bates et al., 1994. Word categories were derived from the categories on the CDI forms (e.g., verbs are listed as 'action words'). Lexical category serves as interacting variable with all predictors.

## 4.2 Regression models

Our approach involves a nested model comparison between the model containing previously known predictors, the uni-gram model, and augmented models that additionally contain residualised average surprisal from either LSTM or n-gram language models. This comparison allows us to assess whether adding information from more dynamic context sizes using LSTMs or from fixed n-gram context sizes beyond log frequency helps predict AoA. We consider the uni-gram model as our base model:

1. The uni-gram model: AoA  $\sim$  lexical category \* (uni-gram surprisal + concreteness)

We then compare it to the following augmented models:

- 2. The residualised n-gram model: AoA  $\sim$  lexical category \*(uni-gram surprisal + n-gram residual average surprisal + concreteness), where n-grams are either bi-grams, tri-grams, or four-grams.
- 3. The residualised LSTM model: AoA  $\sim$  lexical category \*(uni-gram surprisal + LSTM residual average surprisal + concreteness)

#### 4.3 Results

We compare the uni-gram models in each language to their augmented versions which additionally contains residualised LSTM or n-gram average surprisals as predictors. We analyse the difference between the base uni-gram models and the augmented models using both Leave-one-out (LOO) cross-validation and an ANOVA nested model comparison across languages. We report the mean absolute deviation (MAD) across all LOO model fits in each language: the lower the MAD, the better the model fit.

<sup>&</sup>lt;sup>8</sup>We also tried considering different forms as separate items; these results are available in the Appendix D.

#### 4.3.1 Predictability overall

Table 3: Model comparison results augmenting uni-gram surprisal model with residualised surprisals by language and model. Numbers indicate leave-one-out cross validation mean absolute deviation (in months) with 95% CIs.

Language	uni-gram	2-gram	3-gram	4-gram	LSTM
English (American)	$2.0_{[1.86,2.14]}$	$2.02_{[1.88,2.16]}$	$2.01_{[1.87,2.15]}$	$2.01_{[1.87,2.14]}$	2.01 <sub>[1.87,2.15]</sub> *
English (British)	$2.19_{[2.0,2.39]}$	$2.22_{[2.02,2.41]}$	$2.21_{[2.01,2.4]}$	$2.21_{[2.01,2.4]}$	<b>2.18</b> <sub>[1.99,2.37]</sub> *
English (Australian)	<b>1.93</b> <sub>[1.76,2.1]</sub>	<b>1.93</b> <sub>[1.76,2.1]</sub>	$1.94_{[1.77,2.11]}$	$1.94_{[1.77,2.11]}$	
German	$2.2_{[1.99,2.41]}$	$2.22_{[2.02,2.43]}$	$2.22_{[2.01,2.43]}$	$2.21_{[2,2.42]}$	$2.3_{[2.02,2.58]}$
French (European)	$2.3_{[2.1,2.5]}$	$2.33_{[2.12,2.53]}$	$2.31_{[2.11,2.52]}$	$2.3_{[2.1,2.51]}$	$2.46_{[2.12,2.8]}$
French (Quebecois)	$2.54_{[2.35,2.73]}$	$2.56_{[2.37,2.75]}$	<b>2.54</b> <sub>[2.31,2.73]</sub>		
Spanish (European)	<b>2.48</b> <sub>[2.29,2.67]</sub>	$2.49_{[2.3,2.69]}$	$2.5_{[2.31,2.69]}$	$2.51_{[2.31,2.7]}$	$2.49_{[2.3,2.68]}$
Spanish (Mexican)			$2.1_{[1.93,2.28]}$	$2.09_{[1.92,2.27]}$	<b>2.06</b> <sub>[1.89,2.23]</sub> *
Mandarin (Beijing)			* <b>1.86</b> <sub>[1.73,2.0]</sub> * *	1.87 <sub>[1.74,2.01]</sub>	
Mandarin (Taiwanese)			$3.0_{[2.78,3.23]}$	$3.0_{[2.77,3.22]}*$	<b>2.98</b> <sub>[2.75,3.2]</sub> *

<sup>\*</sup> (p < 0.05) and \*\* (p < 0.01) indicate that the nested ANOVA is significant.

The overall results are available in Table 3, where we report MAD and 95% confidence intervals across LOO cross-validation folds, as well as *p* values from our ANOVA nested model comparison. The models with the smallest MAD, or best fits, are bolded.

The nested ANOVA results suggest that adding residualised LSTM surprisal as a predictor significantly increases model fit in four of the ten datasets. However, given the cross-validation results which show that there is very little difference in MAD values between our base models and augmented models, we may want to be cautious about these results and suggest instead that the overall effects are likely small.

The large majority of items across all languages are nouns. Average surprisal using more linguistic context may not be such a good predictor of the AoA of nouns, at least not as much as simple frequency (Portelance, Degen, & Frank, 2020). For this reason, the fact that we do not see a large difference between our base and augmented models here may be expected. If we consider the interaction between lexical category and residualised LSTM surprisal, we find that the interaction terms for predicates are significant (p < 0.05) in five of the languages (*English* (*American*), *English* (*British*), *English* (*Australian*), *French* (*Quebecois*), *Spanish* (*Mexican*)).

#### 4.3.2 Predictability by lexical category

Adding residualised surprisal beyond uni-gram surprisal as a predictor of AoA does not substantially improve overall prediction, but it may make a difference for predicting the AoA of specific words. We next consider how effect sizes may differ by lexical category. Here, we do so by plotting the estimated coefficients by lexical category for the different surprisal predictors. Here, we do so by plotting the estimated coefficients by lexical category for the

The plots in Figure 2 show the estimated coefficients for variables across LOO folds by lexical category. In each of these graphs, a point represents the estimated coefficient of a predictor for one fitted model run from the LOO cross-validation, so for example in *English* (*American*) there are 563 items and therefore 563 folds all of which have been plotted here.

For each language in Figure 2, we see that residualised bi-gram, tri-gram, four-gram, and LSTM surprisals generally have little to no effect for nouns, but have a positive effect for function words and predicates in most languages. In other words, the harder function words and predicates are to predict in their linguistic contexts, the later they are

<sup>&</sup>lt;sup>9</sup>We used contrast coding for our lexical categories, such that interaction terms between a lexical category and residualised surprisal can be interpreted as main effects, indicating a difference from the overall mean.

<sup>&</sup>lt;sup>10</sup>Since lexical categories differ significantly in their concreteness ratings, we also explored whether the by-lexical-category differences in surprisal effects seen in this subsection can be better explained by an interaction between residualised surprisal and concreteness in Appendix F.

<sup>&</sup>lt;sup>11</sup>In Appendix E, we also fit models to items in each lexical category individually and consider their MAD; Table 5 shows that nouns are generally best predicted by simply using the uni-gram base model, while predicates and function words often benefit from the augmented versions which include some form of higher order residualised surprisal, LSTM and tri-gram residualised surprisal generally doing best.

acquired. The exception to this rule is Mandarin, where both function words and predicates show a negative effect for all types of residual surprisal, meaning predicates with higher average surprisal, or less predictability, seem to be learned earlier and have a lower AoA.

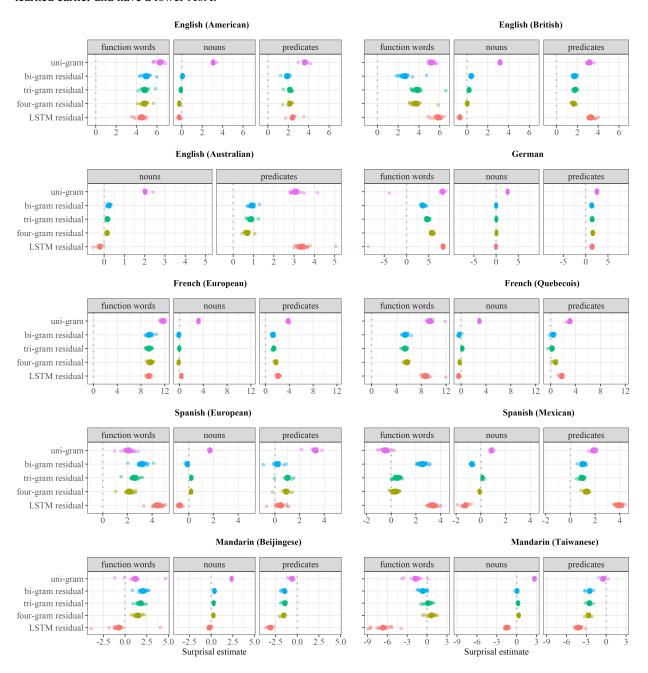


Figure 2: Coefficient estimates for surprisal values by lexical category in each language

Several explanations for the disparate results in Mandarin are possible. First, Mandarin has been known to pattern differently to other languages in other early word learning studies. For example, it has been said that Mandarin learners do not show the same 'noun bias' – the observations that learners tend to initially learn to produce more nouns before increasing their productions of verbs – that English learners seem to have (Tardif, Gelman, & Xu, 1999), an observation that has been reproduced using the Wordbank parental report data (Frank et al., 2021; Yee, 2020). Instead, Mandarin learners have been found to produce more predicates early on during learning, even though both English and Mandarin speaking parents produce relatively more predicates than nouns (Tardif et al., 1999). Another

second possible explanation for this difference may however lie in some of the Wordbank data itself. As Frank et al. (ch.11; 2019) note, there seem to be some discrepancies with some of the Mandarin data in the repository, specifically forms collected for *Mandarin* (*Beijing*) from the Tardif, Fletcher, Liang, & Kaciroti, 2009 study seem to show a much stronger predicate bias than other forms available for *Mandarin* (*Beijing*) or other languages. This data imbalance may also be contributing to this effect for predicates. On the other hand, the *Mandarin* (*Taiwanese*) data is not known to have this same issue, yet it shows a similarly negative effect for average surprisal with predicates albeit weaker than that of *Mandarin* (*Beijing*).<sup>12</sup>

Finally, uni-gram surprisal has a positive effect for all lexical categories across almost all language: words that are more predictable based on their frequency are generally learned earlier, here again the exception being Mandarin function words and predicates.

#### 4.4 Interim conclusions

In this first experiment, we addressed our first set of research questions: Does a word's predictability across linguistic contexts help explain how difficult a word is to learn beyond previously known predictors? If so, how much sequential context matters? And, is predictability equally important for all words or not and why? Since average surprisal and log frequency are correlated (see Appendix B), we elected to remove any variance explained by frequency by using residualised surprisal values as our predictor. We found that including word predictability as a predictor of AoA helped explain some of the variance unaccounted for by previous predictors like log frequency. However, these effects were not universal, showing differences across lexical categories. Specifically, more predictable predicates and function words were found to generally be acquired earlier than their less predictable counterparts. As for how much previous sequential context mattered when determining word predictability, based on our plots in Figure 2, residual LSTM surprisal had a notably greater effect size then all other n-gram surprisals, suggesting that dynamic context sizes which vary from word to word may be most useful in measuring predictability.<sup>13</sup>

# 5 Experiment 2: The role of word predictability across languages

In the previous section, we fit models for different languages using different word lists (table 2). It was therefore hard to determine whether the differences in effect sizes we observed across languages were caused by variation within each individual language word lists or by real distinctions in effect sizes. To remedy this issue, we needed to use the same word list for all languages. We achieved this by unifying words with the same concept across languages using their unilemmas, and then taking their intersection.

First, we aggregated our data across languages to find the intersection of all unilemmas. There were a total of 89 unilemmas for which we have AoA estimates in all languages. These included 64 nouns and 25 predicates; no function words were left because one language, *English* (*Australian*), did not have function words in its word lists. Although we were left with very few unilemmas overall, since we have 10 language groups, we still had 640 noun items and 250 predicate items. Our previous results suggested that the effects of residualised surprisal differed by lexical category, for this reason we split our data into nouns and predicates, testing each category separately. We then fitted a mixed-effects models on each category with by-language random effects to compare effect sizes across languages.

#### 5.1 Regression models

The models in this section differ from the those used in the previous section as they include an additional random effect term: (1 + uni-gram surprisal + residual average surprisal | language). This term means that the models consider by-language differences in coefficient estimates for intercepts, the effects of uni-gram surprisal and LSTM or n-gram residual average surprisal, taking those differences as a source of variance in the data. We did not include a separate random effect by language for concreteness ratings, since these were based on the unilemmas for words and were therefore identical in all languages.

<sup>&</sup>lt;sup>12</sup>It is also possible that these differences are simply due to strength and nature of the resources we had at our disposal for Mandarin. For example, it is possible that different corpora in CHILDES data used different character systems or word segmentation norms, since there exists different standards for this language, which in turn could have led to noisier data.

<sup>&</sup>lt;sup>13</sup>The importance of context size has previously been studied in the case of contextual diversity (Hills et al., 2010), but as we note in Appendix C, contextual diversity and predictability measure different properties and may therefore require different types of contexts.

- 1. The residualised n-gram mixed-effects model: AoA ~ uni-gram surprisal + n-gram residual average surprisal + concreteness + (1 + uni-gram surprisal + n-gram residual average surprisal | language), where n-grams are either bi-grams, tri-grams, or four-grams.
- 2. The residualised LSTM mixed-effects model: AoA  $\sim$  uni-gram surprisal + LSTM residual average surprisal + concreteness + (1 + uni-gram surprisal + LSTM residual average surprisal | language)

#### 5.2 Results

All predictors were scaled by language, centering predictors at zero and setting the standard deviation to one<sup>14</sup>.

#### 5.2.1 Predictability overall

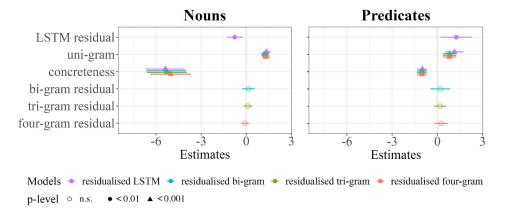


Figure 3: Mean coefficient estimates for surprisal values by lexical category for mixed-effects model using the intersection of items in all languages

As shown in Figure 3, when we examine words occurring in all languages, we find that the effects of frequency (unigram surprisal) and concreteness are the same as before. Words with higher uni-gram surprisal are learned later, i.e., more frequent words are learned earlier. Controlling for frequency and concreteness, none of the residual n-gram surprisal predictors showed significant effects overall, but LSTM surprisal did. The reason may be that n-grams only include information from a fixed length of context, while LSTM surprisal contains information from more dynamic context sizes in the course of word learning. For predicates, higher LSTM average surprisal predicts later AoA, so less predictable predicates across linguistic contexts are harder to learn and therefore learned later. However, for nouns, residual LSTM surprisal instead shows a negative effect, meaning that less predictable nouns tend to be learned earlier.

<sup>&</sup>lt;sup>14</sup>Doing so for each language individually allows us to compare the effects across languages without worrying that any differences in variation that are due to our predictive n-gram or LSTM models having different surprisal distributions in different languages.

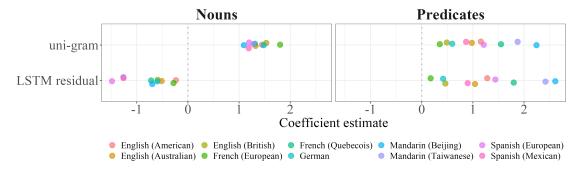


Figure 4: Coefficient estimates for surprisal values for nouns and predicates in each language. Only words with data for all languages are included in this analysis.

A possible explanation for the case of nouns is that cross-situational learning for nouns across diverse social and visual contexts in theory may also lead to more diverse linguistic contexts and therefore higher surprisal (Hills, 2013; Roy et al., 2015). For predicates, however, children rely heavily on linguistic context to learn these words. Therefore, higher predictability across linguistic contexts may indicate that there are more linguistic cues helping children learn certain predicates earlier than others.

#### 5.2.2 Surprisal and AoA across languages

Figure 4 plots effect sizes of uni-gram surprisal and residual LSTM surprisal for each individual language by lexical category <sup>15</sup>. Comparing the predictive power of surprisal across languages, we find that uni-gram surprisal (word frequency) is equally, and positively predictive of noun AoA. LSTM surprisal for nouns shows some variation, being somewhat more negative – greater surprisal is associated with earlier AoA – in Mandarin and Spanish than in other tested languages. In the case of predicates, we see more variation. This result may simply be due to there being fewer items used to fit the predicate mixed-effects model, with only 250 words that were predicates compared to 640 words in the nouns mixed-effects model. The positive effect of both uni-gram surprisal and LSTM surprisal is strongest for Mandarin (Taiwanese and Beijing). This is an interesting result because our analyses from our first experiment using separate models for each language (section 4) found that Mandarin was the only language to show a negative effect for surprisal when predicting the AoA of predicates. It is likely that this difference in effect polarity is due to items included the mixed-effects models being a small sub-sample of those included in the previous experiment, suggesting that the specific words used to fit the models may be introducing a bias towards one or the other effect direction. We test and confirm this hypothesis in Appendix H.

What makes the items from experiment 2 different from the rest in experiment 1 specifically in Mandarin is unclear (the item lists are provided in the appendix as well). However, we note that in experiment 1 uni-gram surprisal also seemed to have negative estimates, albeit smaller, for function words and predicates in this language. Less frequent function words and predicates were supposedly acquired earlier contradicting the pattern seen in all other languages. Given that the discrepancy in average surprisal estimate polarity between experiment 1 and 2 extends to uni-gram surprisal and not only LSTM average surprisal, the issue is unlikely to be with our predictability measure. Instead, it likely follows from issues with overall data quality in this language specifically. There are two places where data quality may be eroded. First, as we explain in the results section of experiment 1, the AoA estimates for Mandarin may be questionable. Second, the corpus built from CHILDES data that we used to calculate both log frequency and probability model surprisal values could be the issue. It could simply be noisier and of poorer quality than the data in other languages, for example, because the Mandarin data was much sparser in CHILDES than many of the other languages (see Table 1).

Although children learning different languages follow broadly similar learning trajectories (e.g., learning frequent concrete nouns before less common abstract verbs), there *are* differences between the AoA of words that mean roughly the same thing across languages. Taking just the words available in all the languages, we find that for 45% of the word/language pairs, the mean difference in AoA is greater than 2 months. For 16% of the pairs, it is greater than 4 months. Are these differences predictable by differences in average surprisal for the same word across languages?

<sup>&</sup>lt;sup>15</sup>The effect sizes by language for residualised n-gram surprisal values are available in Appendix G.

To answer this question, we first computed differences in AoA, uni-gram surprisal, and LSTM residual average surprisal (i.e., LSTM average surprisal controlling for uni-gram surprisal) for each lemma attested in each pair of languages (E.g., American English and Spanish). Uni-gram surprisal and LSTM residual average surprisal values were scaled within each language before computing the differences. We then used a mixed-effects model to predict cross-linguistic differences in AoA from differences in uni-gram surprisal and differences in LSTM residual surprisal. Since some pairs of languages belong to the same base language, e.g., English (American) and English (Australian), we also added a binary variable to the model indicating whether the two languages have the same base language. Moreover, we also included by-unilemma, by-the-first-base-language, and by-the-second-base-language random intercepts in the model, i.e., (1 | unilemma) + (1 | base language 1) + (1 | base language 2) because those random effects could also be a source of variance in the data. Concreteness was not included as a predictor because we did not have separate concreteness estimates for each language.

• The AoA differences mixed-effects model: AoA differences  $\sim$  differences in uni-gram surprisal + differences in LSTM residual average surprisal + whether the base languages are the same +  $(1 \mid \text{unilemma}) + (1 \mid \text{base language 1}) + (1 \mid \text{base language 2})$ 

From our final analysis (Figure 5), we see that for both nouns and predicates, if a word occurs more frequently in language 1 than its equivalent <sup>16</sup> in language 2, then its AoA is likely to be earlier in language 1. The same pattern is true with higher-order LSTM surprisal for predicates, though for nouns we see an opposite effect, which aligns with our previous finding that LSTM surprisal has opposite effects on predicates and nouns. A greater difference in surprisal between two languages is associated with earlier AoAs for nouns and later AoAs for predicates.

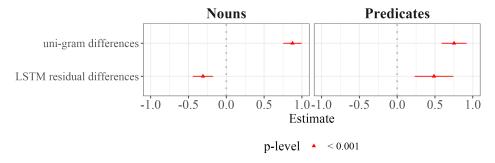


Figure 5: Differences in frequency (uni-gram surprisal) and LSTM-based word surprisal are associated with differences in AoA. A positive coefficient indicates that if the same word has higher surprisal in language 1 than language 2, it is learned later in language 1 than language 2.

#### 5.3 Interim conclusions

In this second experiment, we tried to answer our second set of questions: are the effects of predictability on AoA observable across many languages and linguistic communities or are they isolated to only some? And, do differences in word predictability predict differences in AoA between languages? We found that word predictability in context was a good predictor across language, showing little variation in effect size. The effects for nouns and predicates had polarity, however, such that less predictable nouns in linguistic contexts were learnt earlier, while less predictable predicates were learnt later. Additionally, in our secondary analysis we found that for predicates, there was a positive relationship between differences in surprisal and differences in AoA: predicate lemmas that have a greater surprisal in language 1 than language 2 tend to have a later AoA in language 1.

#### 6 Discussion

Are words that are more predictable given their linguistic context (i.e., words with lower surprisal), easier to learn by young children? For predicates and function words, the answer seems to be yes. Although uni-gram surprisal (i.e., word log frequency) was overall a better predictor of AoA than higher-order surprisal, surprisal derived from an LSTM

<sup>&</sup>lt;sup>16</sup>Equivalent words refer to word with the same unilemma in Wordbank.

neural network predicted AoA beyond uni-gram surprisal and the words' frequency. This effect was largely restricted to predicates and function words – words whose meaning is especially dependent on their context.

One route by which linguistic context is known to affect learning is via *syntactic bootstrapping*. Introduced by Brown (1957) and coined by Gleitman (1990), syntactic bootstrapping refers to the use of syntactic context to help determine meaning, especially for predicates. A wealth of evidence from individual experiments and computational models supports the ability of children to use linguistic information to make inferences about meaning (Fisher, Gertner, Scott, & Yuan, 2010; Fisher, Jin, & Scott, 2020). Although our experiments do not directly test this idea, our results are certainly congruent with syntactic bootstrapping accounts.

A somewhat puzzling finding is that when predicting the AoA of nouns, greater LSTM surprisal was associated with *earlier* age of acquisition – the opposite of what we observed for predicates and function words. One possible explanation for this pattern is that noun learning is more dependent on the perception of concrete referents, the salience of which is often obvious from the extra-linguistic context (Gillette, Gleitman, Gleitman, & Lederer, 1999). Frequent nouns are likely to appear in a broader set of contexts, which would *increase* their higher-order surprisal while at the same time making it more semantically interconnected and improve learning (Hills et al., 2010; Hills, 2013; Roy et al., 2015). The results in Appendix C which show how frequency and predictability relate to contextual diversity support this hypothesis. Frequency and contextual diversity are strongly correlated, thus, more frequent words tend to also be more contextually diverse. Furthermore, contextual diversity is negatively correlated with predictability, in other words more predictable words also tend to have higher contextual diversity.<sup>17</sup>

We found the effects of surprisal to be mostly consistent across the tested languages. The one exception was Mandarin. In our first experiment using separate models for each language, Mandarin showed the opposite effects for surprisal on predicates and function words compared to other languages. However, in experiment 2 using a single-mixed effects model over common unilemmas, Mandarin showed the largest effect of LSTM-based surprisal on AoA. We suggested that this discrepancy may follow from data quality issues in Mandarin specifically.

#### **6.1** Broader theoretical implications

Beyond helping us understand the role played by sequential word predictability in determining which words are easier or harder to learn, this research can also offer insights about other theoretical questions.

First, a broad literature has explored the so-called 'AoA effect' which refers to the finding that earlier learnt words are both easier to process for adults (Elsherif, Preece, & Catling, 2023) and more robustly accessible in aphasia patients (Brysbaert & Ellis, 2016). Most studies have looked at the effect specifically with referential nouns; a few exceptions have also considered it with predicates and none have looked at function words, i.e. Colombo & Burani, 2002; Morrison, Hirsh, & Duggan, 2003; Bogka et al., 2003; Bonin, Boyer, Méot, Fayol, & Droit, 2004; Boulenger, Décoppet, Roy, Paulignan, & Nazir, 2007. The first four mentioned studies found an AoA-effect for both nouns and verbs, all however used an image naming task and AoA-effects are known to be strongly modulated by the imageability of words (Elsherif et al., 2023). The most recent study, Boulenger et al., 2007, used a lexical decision task which requires participants to read words and decided whether they are a real or nonce word. The AoA effect was only found with nouns and not verbs. The authors suggest that verbs and nouns may be learnt and processed differently resulting in different AoA effects depending on task demands. Our results confirm that nouns and predicates may not be learnt using the same cues and that their AoAs are dependent on different predictive factors, including predictability, or average surprisal, in the case of predicates. These learning differences could then result in different AoA effects, depending on task demands. Surprisal is known to predict sentence processing difficulty in incremental reading tasks (Hale, 2016); Reading task demands – as opposed to image naming task demands – may be different for words whose AoAs are more dependent on surprisal, like predicates, than those that are not, like nouns. 18

Second, our results shed light on a puzzle of linguistic complexity. Languages spoken by larger communities, those with many nonnative speakers, and those that have undergone substantial language contact (factors that are often, but not always positively correlated), tend to be morphologically simpler (Wray & Grace, 2007; Lupyan & Dale, 2010; Trudgill, 2011; Winters, Kirby, & Smith, 2015; Bentz, Dediu, Verkerk, & Jäger, 2018; Trudgill, 2011). On the other hand, languages spoken by smaller groups are more complex and less compressible (Lupyan & Dale, 2010; Lupyan,

<sup>&</sup>lt;sup>17</sup>We note that when controlling for the effect of log frequency, predictability and contextual diversity are no longer as correlated. Additionally, nouns pattern differently from predicates and function words in terms of the relationship between predictability and frequency.

<sup>&</sup>lt;sup>18</sup>We note the caveat that Boulenger et al., 2007 used subjective AoA measures based on retrospective adult ratings; Our research used objective ratings based of parental reports of children's productions. As mentioned in the introduction these measures are related but not always strongly correlated.

2019; Koplenig, 2019). But why? One possibility – a prediction made by Lupyan and Dale's (2010) Linguistic Niche Hypothesis – is that there may be a trade-off between optimizing languages for use by larger/more diverse groups, and optimizing languages for maximally efficient learning by young children. Greater compressibility is intimately linked to predictability: a more compressible language is one that contains utterances where one part is better predicted from other parts, making them informationally redundant. An intriguing possibility is that the reason smaller languages are more compressible – i.e., have overall higher predictability, or lower average surprisal – is that predictability is especially important to young children learning their first language (Lupyan & Dale, 2010, 2016). Our results show initial support for this idea: at least for predicates, a decrease in surprisal is associated with faster word learning suggesting that redundancy (lower surprisal) may provide more effective learning opportunities for young children. Further tests of this hypothesis would require comparison across a much wider range of languages, however.

## 7 Conclusion and Limitations

We investigated the relationship between word predictability – formulated as average surprisal in linguistic contexts – and Age of Acquisition – a proxy for how difficult learning a given word is for children. Word predictability seems to be especially important for the acquisition of predicates and function words, rather than nouns. Less predictable verbs and adjective in linguistic contexts tend to be acquired later by children and are therefore harder to learn. We found this effect to be true is multiple languages, and that differences in surprisal predict differences in AoA. This finding is broadly in line with the prediction that decreasing surprisal (i.e., increasing redundancy) may – all else being equal – facilitate child language learning, particularly of more relational words.

Theories of language learning must explain not just words like 'ball' and 'dog' but also words whose meaning in context is almost entirely dependent on other words. Sequential models like the LSTM we used here may be a promising avenue to help explain the acquisition of these 'hard words'.

A major limitation of our analyses is that it is based on the data available in CHILDES. These corpora of child-directed utterances were – by necessity – assembled from many sub-corpora from different children and studies, and may not be the best representations of the regularity, idiosyncrasies, and contextual diversity found in the language targeted to a single child. Furthermore, our corpora contained utterances directed at children who were older in some cases than those surveyed for the AoA estimates. Ideally, these sentences would span the exact same developmental stages. A second limitation is that the words and AoA estimates we use are based on parents' reports of their children's language. Although validated and highly reliable, these data cannot capture the full richness of individual children's vocabularies and language use. Finally, our sample of language is restricted to languages with a large digital footprint. We hope future studies can begin to overcome these limitations.

#### References

- Amatuni, A., & Bergelson, E. (2017). Semantic networks generated from early linguistic input. bioRxiv, 157701.
- Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P., Reznick, J. S., ... Hartung, J. (1994). Developmental and stylistic variation in the composition of early vocabulary. *Journal of Child Language*, 21(1), 85–123.
- Bentz, C., Dediu, D., Verkerk, A., & Jäger, G. (2018). The evolution of language families is shaped by the environment beyond neutral drift. *Nature Human Behaviour*, 2(11), 816–821. doi: 10.1038/s41562-018-0457-6
- Bogka, N., Masterson, J., Druks, J., Fragkioudaki, M., Chatziprokopiou, E.-S., & Economou, K. (2003). Object and action picture naming in english and greek. *European Journal of Cognitive Psychology*, 15(3), 371–403.
- Bonin, P., Boyer, B., Méot, A., Fayol, M., & Droit, S. (2004). Psycholinguistic norms for action photographs in french and their relationships with spoken and written latencies. *Behavior Research Methods, Instruments, & Computers*, 36, 127–139.
- Boulenger, V., Décoppet, N., Roy, A. C., Paulignan, Y., & Nazir, T. A. (2007). Differential effects of age-of-acquisition for concrete nouns and action verbs: Evidence for partly distinct representations? *Cognition*, 103(1), 131–146.
- Braginsky, M., Yurovsky, D., Marchman, V. A., & Frank, M. C. (2019). Consistency and variability in children's word learning across languages. *Open Mind*, 1–16.
- Brinchmann, E. I., Braeken, J., & Lyster, S.-A. H. (2019). Is there a direct relation between the development of vocabulary and grammar? *Developmental Science*, 22(1), e12709.
- Brown, R. W. (1957). Linguistic determinism and the part of speech. *The Journal of Abnormal and Social Psychology*, 55, 1–5.

- Brysbaert, M., & Ellis, A. W. (2016). Aphasia and age of acquisition: are early-learned words more resilient? *Aphasiology*, 30(11), 1240–1263.
- Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014). Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior Research Methods*, 46(3), 904–911.
- Chang, T. A., & Bergen, B. K. (2022). Word acquisition in neural language models. *Transactions of the Association of Computational Linguistics*. Retrieved from https://arxiv.org/abs/2110.02406
- Christiansen, M. H., Allen, J., & Seidenberg, M. S. (1998). Learning to segment speech using multiple cues: A connectionist model. *Language and cognitive processes*, 13(2-3), 221–268.
- Clark, E. V. (1993). Early lexical development. In *The lexicon in acquisition* (p. 21–42). Cambridge University Press.
- Colombo, L., & Burani, C. (2002). The influence of age of acquisition, root frequency, and context availability in processing nouns and verbs. *Brain and Language*, 81(1-3), 398–411.
- Demberg, V., & Keller, F. (2008). Data from eye-tracking corpora as evidence for theories of syntactic processing complexity. *Cognition*, 109(2), 193–210.
- Elman, J. L. (1990). Finding structure in time. Cognitive Science, 14(2), 179–211.
- Elman, J. L. (1993). Learning and development in neural networks: The importance of starting small. *Cognition*, 48(1), 71–99.
- Elsherif, M., Preece, E., & Catling, J. (2023). Age-of-acquisition effects: A literature review. *Journal of Experimental Psychology: Learning, Memory, and Cognition*.
- Fenson, L., Dale, P., Reznick, J. S., Thal, D., Bates, E., Hartung, J., ... Reilly, J. (1993). MacArthur Communicative Inventories: User's guide and technical manual. *San Diego*.
- Fisher, C., Gertner, Y., Scott, R. M., & Yuan, S. (2010). Syntactic bootstrapping. Wiley Interdisciplinary Reviews: Cognitive Science, 1(2), 143–149.
- Fisher, C., Jin, K.-s., & Scott, R. M. (2020). The developmental origins of syntactic bootstrapping. *Topics in Cognitive Science*, *12*(1), 48–77.
- Fourtassi, A., Bian, Y., & Frank, M. C. (2020). The growth of children's semantic and phonological networks: Insight from 10 languages. *Cognitive Science*, 44(7), e12847.
- Frank, M. C., Braginsky, M., Marchman, V., & Yurovsky, D. (2021). *Variability and Consistency in Early Language Learning: The Wordbank Project*. MIT Press.
- Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2016). Wordbank: An open repository for developmental vocabulary data. *Journal of Child Language*, 44(3), 677–694.
- Gillette, J., Gleitman, H., Gleitman, L., & Lederer, A. (1999). Human simulations of vocabulary learning. *Cognition*, 73(2), 135–176. doi: 10.1016/s0010-0277(99)00036-0
- Gleitman, L. (1990). The structural sources of verb meanings. Language Acquisition, 1(1), 3-55.
- Goodman, J. C., Dale, P. S., & Li, P. (2008). Does frequency count? Parental input and the acquisition of vocabulary. *Journal of Child Language*, 35(3), 515–531.
- Hale, J. (2016). Information-theoretical complexity metrics. Language and Linguistics Compass, 10(9), 397–412.
- Hills, T. T. (2013). The company that words keep: comparing the statistical structure of child-versus adult-directed language. *Journal of child language*, 40(3), 586–604.
- Hills, T. T., Maouene, J., Riordan, B., & Smith, L. B. (2010). The associative structure of language: Contextual diversity in early word learning. *Journal of Memory and Language*, 63(3), 259–273.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9(8), 1735–1780.
- Hsiao, Y., & Nation, K. (2018). Semantic diversity, frequency and the development of lexical quality in children's word reading. *Journal of Memory and Language*, 103, 114-126. Retrieved from https://www.sciencedirect.com/science/article/pii/S0749596X18300718 doi: https://doi.org/10.1016/j.jml.2018.08.005
- Huebner, P. A., Sulem, E., Cynthia, F., & Roth, D. (2021). BabyBERTa: Learning more grammar with small-scale child-directed language. In *Proceedings of the 25th conference on computational natural language learning* (pp. 624–646).
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. Retrieved from https://arxiv.org/pdf/1412.6980.pdf
- Koplenig, A. (2019). Language structure is influenced by the number of speakers but seemingly not by the proportion of non-native speakers. *Royal Society Open Science*, 6(2), 181274. Retrieved 2022-11-08, from https://royalsocietypublishing.org/doi/full/10.1098/rsos.181274 (Publisher: Royal Society) doi: 10.1098/rsos.181274

- Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for 30,000 English words. *Behavior Research Methods*, 44(4), 978–990.
- Levy, R. (2008). Expectation-based syntactic comprehension. Cognition, 106(3), 1126–1177.
- Linzen, T., Dupoux, E., & Goldberg, Y. (2016). Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *TACL*, *4*, 521–535.
- Łuniewska, M., Haman, E., Armon-Lotem, S., Etenkowski, B., Southwood, F., Andelković, D., ... Ünal Logacev, O. (2016). Ratings of age of acquisition of 299 words across 25 languages: Is there a cross-linguistic order of words? *Behavior research methods*, 48, 1154–1177.
- Lupyan, G. (2019). Larger languages have higher entropy.. Retrieved from https://cle.ppls.ed.ac.uk/index
  .php/ielc2019/
- Lupyan, G., & Dale, R. (2010). Language structure is partly determined by social structure. *PLoS ONE*, 5(1), e8559. Retrieved from http://dx.doi.org/10.1371/journal.pone.0008559 doi: 10.1371/journal.pone.0008559
- Lupyan, G., & Dale, R. (2016). Why are there different languages? the role of adaptation in linguistic diversity. *Trends in Cognitive Sciences*, 20(9), 649–660. doi: http://dx.doi.org/10.1016/j.tics.2016.07.005
- MacWhinney, B. (2000). The CHILDES Project: Tools for analyzing talk. Third Edition. Erlbaum.
- Morrison, C. M., Hirsh, K. W., & Duggan, G. B. (2003). Age of acquisition, ageing, and verb production: Normative and experimental data. *The Quarterly Journal of Experimental Psychology Section A*, 56(4), 1–26.
- Pagán, A., Bird, M., Hsiao, Y., & Nation, K. (2020). Both semantic diversity and frequency influence children's sentence reading. *Scientific Studies of Reading*, 24(4), 356-364. Retrieved from https://doi.org/10.1080/10888438.2019.1670664 doi: 10.1080/10888438.2019.1670664
- Portelance, E., Degen, J., & Frank, M. C. (2020). Predicting Age of Acquisition in Early Word Learning Using Recurrent Neural Networks. In *Cogsci*.
- Rosa, E., Tapia, J. L., & Perea, M. (2017). Contextual diversity facilitates learning new words in the classroom. *PLoS One*, 12(6), e0179004.
- Roy, B. C., Frank, M. C., DeCamp, P., Miller, M., & Roy, D. (2015). Predicting the birth of a spoken word. *Proceedings of the National Academy of Sciences*, 112(41), 12663–12668.
- Saffran, J. R. (2020). Statistical language learning in infancy. *Child development perspectives*, *14*(1), 49–54. Retrieved from https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8078161/ doi: 10.1111/cdep.12355
- Sagae, K. (2021). Tracking child language development with neural network language models. *Frontiers in Psychology*, 12.
- Sanchez, A., Meylan, S., Braginsky, M., MacDonald, K., Yurovsky, D., & Frank, M. C. (2018). childes-db: a flexible and reproducible interface to the Child Language Data Exchange System.
- Smith, N. J., & Levy, R. (2013). The effect of word predictability on reading time is logarithmic. *Cognition*, 128, 302–319.
- Stella, M., Beckage, N. M., & Brede, M. (2017). Multiplex lexical networks reveal patterns in early word acquisition in children. *Scientific reports*, 7(1), 1–10.
- Sundermeyer, M., Schlüter, R., & Ney, H. (2012). LSTM neural networks for language modeling. In 13th annual conference of the international speech communication association.
- Tardif, T., Fletcher, P., Liang, W., & Kaciroti, N. (2009). Early vocabulary development in Mandarin (Putonghua) and Cantonese. *Journal of child language*, *36*, 1115–1144.
- Tardif, T., Fletcher, P., Liang, W., Zhang, Z., Kaciroti, N., & Marchman, V. A. (2008). Baby's first 10 words. *Developmental Psychology*, 44(4), 929.
- Tardif, T., Gelman, S. A., & Xu, F. (1999). Putting the "noun bias" in context: A comparison of English and Mandarin. *Child development*, 70, 620–635.
- Trudgill, P. (2011). Sociolinguistic typology: Social determinants of linguistic complexity. Oxford University Press.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998–6008).
- Winters, J., Kirby, S., & Smith, K. (2015). Languages adapt to their contextual niche. Language and Cognition, 7(3), 415-449. Retrieved from https://www.cambridge.org/core/journals/language-and-cognition/article/abs/languages-adapt-to-their-contextual-niche/83E9F516875C340E0A9263B4A7C38F43 (Publisher: Cambridge University Press) doi: 10.1017/langcog.2014.35
- Wray, A., & Grace, G. (2007). The consequences of talking to strangers: Evolutionary corollaries of socio-cultural influences on linguistic form. *Lingua*, 117(3), 543–578. Retrieved from http://cat.inist.fr/?aModele=

afficheN&cpsidt=18462025 Yee, S. (2020). Is noun bias universal? Evidence from Chinese and Korean compared with French and English. In Studies in the linguistic sciences: Illinois working papers (pp. 32–44).

## A LSTM model details

Models were implemented in Python 3.8 using the Pytorch 1.8.1 and trained on a single Nvidia RTX 3080 GPU. They are composed of an embedding layer, two sequential and fully connected LSTM layers, and a linear layer. We used an Adam optimizer across all models during training (Kingma & Ba, 2014) and simple cross-entropy across the outputs of all states as our loss function.

## A.1 Hyperparameters

Here are the hyperparameters tested, bolded ones are the settings we used for the experiments reported in the main paper.

number of epochs: up to 50, **20** used

learning rate: **0.0001**batch size: 128 / **256** / 512
hidden size: **100** / 150
embedding size: **100** / 150

vocabulary size: 2000 / 4000 / 5000

random seed: [0:2]

# B Relationship between uni-gram and higher order surprisals

Table 4: Pearson correlation between LSTM, bi-gram, tri-gram, four-gram average surprisal and uni-gram surprisal in each language

Language	1gm-LSTM	1gm-2gm	1gm-3gm	1gm-4gm
English (American)	0.65	0.91	0.80	0.56
English (British)	0.68	0.89	0.77	0.5
English (Australian)	0.72	0.85	0.68	0.34
German	0.63	0.81	0.51	0.18
French (European)	0.65	0.89	0.59	0.23
French (Quebecois)	0.69	0.87	0.55	0.21
Spanish (European)	0.97	0.87	0.42	0.16
Spanish (Mexican)	0.97	0.83	0.44	0.12
Mandarin (Beijing)	0.96	0.81	0.55	0.43
Mandarin (Taiwanese)	0.97	0.81	0.55	0.41

# C Relationship between contextual diversity and contextual predictability

Contextual diversity was calculated using the procedure outlined in Hills et al., 2010. Using a moving window of size 5, we created co-occurrence matrices for words. We used the same corpora of child-directed utterances in 5 languages that were described in the main paper and used to extract average surprisal. Like with average surprisal, we limited our co-occurrence matrices to the 5000 most common words in each language. We then counted the number of unique words each word appeared with, creating a contextual diversity score (CD). We also calculated the log transformed CD since both uni-gram surprisal (negative log frequency) and average LSTM surprisal are in log scale. Table 5 shows the correlation scores between these four variables in addition to frequency across languages. Both average LSTM surprisal and CD, especially log transformed CD, are highly correlated with uni-gram surprisal. Average LSTM surprisal and CD are not however as highly correlated with one another. In fact, if we control for the effect of log frequency by taking the residual variance of both predictors once we have accounted for uni-gram surprisal and compare them again in Table 6, we see that they are not correlated. Though both average surprisal and CD depend on log frequency, beyond that they encode very different information; CD is a word-level property, it abstracts away from

individual ordered contexts in which words are used, while predictability is a contextual property, taking into account the sequences of words a given word appears in.

Table 5: Pearson correlation between frequency, uni-gram surprisal, LSTM average surprisal, contextual diversity (CD), and log transformed contextual diversity (log CD) across all languages

	uni-gram	LSTM	CD	log CD
frequency	-0.60	-0.31	0.81	0.54
uni-gram		0.79	-0.84	-0.97
LSTM			-0.51	-0.78
CD				0.83

Table 6: Pearson correlation between uni-gram surprisal, residualised LSTM average surprisal, and residualised log transformed contextual diversity by lexical category across all languages

Lexical category		LSTM-rd	log CD-rd
overall	LSTM-rd		-0.14
nouns	uni-gram	0.43	-0.12
	LSTM-rd		-0.14
predicates	uni-gram	-0.20	0.04
	LSTM-rd		-0.05
function words	uni-gram	-0.52	0.37
	LSTM-rd		-0.35

# D Approach 1 - All words as separate items

The following section reproduces the analyses from our first approach, but instead of combining items with multiple lexical forms (e.g. in/inside), it considers each form as a separate item when fitting models. This method may be problematic if one form is much more frequent then another since it will give them equal importance when fitting regression models, as each word form would be considered independently from its related forms.

Table 7 shows the results overall of our model comparison using LOO cross-validation and reporting MAD. The results are similar to those reported in the main paper.

Table 7: Approach 1 model comparison results augmenting uni-gram surprisal model with residualised surprisal by language

			LOO MAD <sub>[95%</sub>	CI]	
Language	1gm-base	2gm-rd	3gm-rd	4gm-rd	LSTM-rd
English (American)	<b>2.01</b> <sub>[1.87,2.15]</sub>	$2.03_{[1.89,2.17]}$	<b>2.01</b> <sub>[1.88,2.15]</sub>	<b>2.01</b> <sub>[1.8,2.15]</sub>	2.02 <sub>[1.88,2.16]</sub>
English (British)	$2.19_{[2.0,2.38]}$	$2.21_{[2.02,2.41]}$	$2.21_{[2.01,2.4]}$	$2.20_{[2.01,2.4]}$	<b>2.17</b> <sub>[1.87,2.15]</sub> *
English (Australian)	$1.94_{[1.77,2.11]}$	$1.93_{[1.77,2.1]}$	$1.94_{[1.77,2.11]}$	$1.94_{[1.77,2.11]}$	<b>1.93</b> <sub>[1.76,2.1]</sub>
German	$2.2_{[1.99,2.41]}$	$2.22_{[2.02,2.43]}$	$2.22_{[2.01,2.43]}$	$2.21_{[2,2.42]}$	$2.3_{[2.02,2.58]}$
French (European)	$2.32_{[2.13,2.52]}$	$2.35_{[2.15,2.54]}$	$2.32_{[2.13,2.53]}$	$2.32_{[2.12,2.51]}$	$2.34_{[2.15,2.54]}$
French (Quebecois)	<b>2.56</b> <sub>[2.37,2.74]</sub>	$2.57_{[2.38,2.76]}$	$2.56_{[2.37,2.75]}$	$2.56_{[2.38,2.75]}$	$2.57_{[2.38,2.76]}$
Spanish (European)	$2.53_{[2.34,2.71]}$	<b>2.52</b> <sub>[2.34,2.71]</sub> >	* 2.54 <sub>[2.36,2.73]</sub>	$2.55_{[2.37,2.73]}$	$2.53_{[2.35,2.71]}$
Spanish (Mexican)	$2.09_{[1.92,2.26]}$	$2.1_{[1.93,2.27]}$	$2.11_{[1.93,2.28]}$	$2.10_{[1.92,2.27]}$	<b>2.06</b> <sub>[1.89,2.34]</sub> **
Mandarin (Beijing)	$1.88_{[1.75,2.01]}$	1.88 <sub>[1.74,2.01]</sub>	* <b>1.86</b> <sub>[1.73,2.99]</sub> * *	1.87 <sub>[1.73,2]</sub> *	$1.89_{[1.75,2.02]}$
Mandarin (Taiwanese)	$3.01_{[2.78,3.23]}$	$3.03_{[2.8,3.26]}$	$3.03_{[2.8,3.26]}$		<b>3</b> <sub>[2.78,3.22]</sub> **

<sup>\* (</sup>p <0.05) and \*\* (p <0.01) indicate that the nested ANOVA is significant.

Figure 6 shows the estimated coefficients for surprisal and residualised surprisal by language and lexical category. These results are similar to those in the main paper.

# E Approach 1 - By lexical category MAD

Table 8 shows the LOO MAD results for models fit on each lexical category separately.

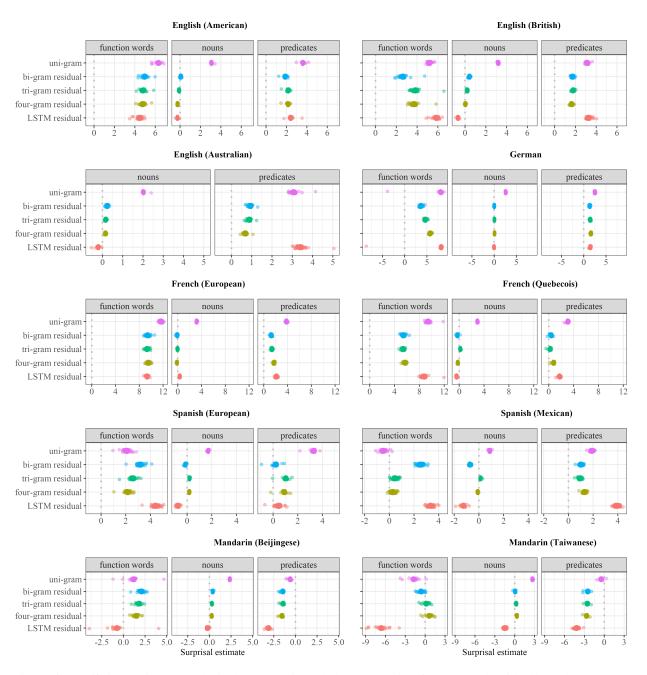


Figure 6: Coefficient estimates by lexical category in each language using first approach with all words as separate items.

Table 8: Approach 1 model comparison results augmenting uni-gram surprisal model with residualised surprisal by language and lexical category

			LOO MA	D <sub>[95% CI]</sub>			
Language	Category	1gm-base	2gm-rd	3gm-rd	4gm-rd	LSTM-rd	
English	Nouns	1.95 <sub>[1.76,2.14]</sub>	1.96 <sub>[1.77,2.15]</sub>	1.96 <sub>[1.77,2.15]</sub>	$1.96_{[1.77,2.15]}$	1.96 <sub>[1.77,2.15]</sub>	
(American)	Predicates	$1.79_{[1.58,2.01]}$	$1.81_{[1.59,2.03]}$	$1.78_{[1.55,2.0]}$	$1.77_{[1.55,2.0]}$	$1.81_{[1.58,2.04]}$	
	Function	$2.52_{[2.13,2.91]}$	$2.56_{[2.16,2.97]}$	$2.55_{[2.15,2.94]}$	$2.54_{[2.15,2.94]}$	$2.54_{[2.15,2.93]}$	
	words						
English	Nouns	$2.21_{[1.96,2.47]}$	$2.22_{[1.97,2.48]}$	$2.23_{[1.97,2.48]}$	$2.23_{[1.97,2.49]}$	$2.2_{[1.94,2.45]}$	
(British)	Predicates	$1.92_{[1.62,2.22]}$	$1.93_{[1.62,2.23]}$	$1.91_{[1.61,2.21]}$	$1.92_{[1.62,2.22]}$	$1.91_{[1.6,2.21]}$	
	Function words	$2.83_{[2.17,3.49]}$	$2.97_{[2.29,3.66]}$	$2.92_{[2.24,3.6]}$	$2.88_{[2.21,3.55]}$	<b>2.17</b> <sub>[1.87,2.15]</sub>	
English	Nouns	$2.06_{[1.83,2.29]}$	$2.06_{[1.83,2.29]}$	$2.06_{[1.83,2.29]}$	<b>2.06</b> <sub>[1.83,2.29]</sub>	$2.07_{[1.83,2.3]}$	
(Australian)	Predicates	$1.73_{[1.49,1.96]}$	$1.72_{[1.48,1.96]}$	$1.72_{[1.49,1.96]}$	$1.74_{[1.5,2.97]}$	$1.71_{[1.48,1.95]}$	
German	Nouns	$1.86_{[1.66,2.05]}$	$1.86_{[1.67,2.06]}$	$1.86_{[1.67,2.06]}$	$1.86_{[1.67,2.06]}$	$1.87_{[1.67, 2.06]}$	
	Predicates	<b>1.86</b> <sub>[1.58,2.14]</sub>	$1.88_{[1.6,2.17]}$	$1.88_{[1.59,2.16]}$	$1.87_{[1.59,2.15]}$	$1.89_{[1.6,2.18]}$	
	Function words	<b>3.77</b> <sub>[2.94,4.6]</sub>	3.86 <sub>[3.06,4.65]</sub>	3.84 <sub>[3.02,4.66]</sub>	3.81 <sub>[2.97,4.65]</sub>	4.27 <sub>[2.99,5.55]</sub>	
French	Nouns	$2.03_{[1.76,2.3]}$	$2.05_{[1.77,2.32]}$	$2.05_{[1.77,2.32]}$	$2.04_{[1.77,2.31]}$	$2.05_{[1.78,2.32]}$	
(European)	Predicates	$2.45_{[2.12,2.79]}$	$2.48_{[2.15,2.81]}$	<b>2.43</b> <sub>[2.1,2.76]</sub>	$2.44_{[2.11,2.78]}$	$2.48_{[2.14,2.82]}$	
	Function words	<b>3.06</b> <sub>[2.46,3.67]</sub>	$3.16_{[2.51,3.82]}$	$3.16_{[2.52,3.81]}$	$3.07_{[2.43,3.71]}$	4.14[2.06,6.21]	
French	Nouns	<b>2.47</b> <sub>[2.21,2.72]</sub>	$2.49_{[2.23,2.75]}$	<b>2.47</b> <sub>[2.21,2.74]</sub>	$2.48_{[2.22,2.74]}$	$2.48_{[2.22,2.73]}$	
(Quebecois)	Predicates	$2.48_{[2.17,2.79]}$	$2.48_{[2.16,2.79]}$	<b>2.46</b> <sub>[2.14,2.77]</sub>	$2.49_{[2.17,2.8]}$	$2.45_{[2.13,2.76]}$	
	Function words	$3.01_{[2.43,3.6]}$	$3.07_{[2.46,3.67]}$	$3.05_{[2.46,3.64]}$	<b>3.0</b> <sub>[2.42,3.59]</sub>	$3.04_{[2.46,3.62]}$	
Spanish	Nouns	<b>2.11</b> <sub>[1.86,2.37]</sub>	$2.12_{[1.86,2.37]}$	$2.13_{[1.88,2.38]}$	$2.13_{[1.88,2.38]}$	<b>2.11</b> <sub>[1.85,2.36]</sub>	
(European)	Predicates	$2.16_{[1.81,2.5]}$	$2.17_{[1.83,2.52]}$	$2.18_{[1.83,2.53]}$	$2.19_{[1.84,2.54]}$	$2.18_{[1.83,2.52]}$	
	Function words	<b>3.69</b> <sub>[3.29,4.09]</sub>	$3.71_{[3.3,4.13]}$	$3.71_{[3.3,4.13]}$	$3.74_{[3.32,4.15]}$	3.71 <sub>[3.31,4.1]</sub>	
Spanish	Nouns	<b>2.04</b> <sub>[1.79,2.29]</sub>	<b>2.04</b> <sub>[1.79,2.29]</sub>	$2.05_{[1.79,2.3]}$	$2.05_{[1.8,2.3]}$	$2.05_{[1.8,2.3]}$	
(Mexican)	Predicates	$1.79_{[1.56,2.03]}$	$1.83_{[1.59,2.07]}$	$1.82_{[1.58,2.06]}$	$1.8_{[1.56,2.05]}$	<b>1.72</b> <sub>[1.47,1.96]</sub>	
	Function words	$2.58_{[2.15,3.0]}$	$2.63_{[2.23,3.04]}$	$2.62_{[2.19,3.06]}$	$2.59_{[2.16,3.03]}$	<b>2.54</b> <sub>[2.13,2.95]</sub>	
Mandarin	Nouns	$1.91_{[1.74, 2.09]}$	$1.9_{[1.72,2.08]}$	$1.9_{[1.72,2.08]}$	$1.9_{[1.73,2.08]}$	$1.92_{[1.74,2.1]}$	
(Beijing)	Predicates	$1.67_{[1.47,1.87]}$	$1.69_{[1.48,1.89]}$	$1.66_{[1.46,1.87]}$	$1.67_{[1.46,1.87]}$	<b>1.65</b> <sub>[1.45,1.85]</sub>	
	Function words	$2.59_{[1.94,3.25]}$	$2.6_{[1.97,3.34]}$	<b>2.46</b> <sub>[1.82,3.09]</sub>	$2.51_{[1.87,3.15]}$	$2.76_{[2.02,3.49]}$	
Mandarin	Nouns	$2.88_{[2.6,3.17]}$	$2.9_{[2.61,3.18]}$	$2.9_{[2.61,3.18]}$	$2.89_{[2.61,3.18]}$	$2.9_{[2.61,3.18]}$	
(Taiwanese)	Predicates	<b>3.07</b> <sub>[2.65,3.48]</sub>	$3.09_{[2.65,3.52]}$	$3.08_{[2.67,3.49]}$	$3.08_{[2.68,3.48]}$	<b>3.07</b> <sub>[2.66,3.47]</sub>	
	Function words	3.39[2.56,4.22]	3.48 <sub>[2.64,4.33]</sub>	3.5 <sub>[2.65,4.34]</sub>	$3.47_{[2.63,4.3]}$	<b>3.25</b> <sub>[2.46,4.04]</sub>	

# **F** Approach 1 - Interactions with concreteness

Since lexical categories differ significantly in their concreteness ratings (Figure 7), it is possible that the differences in effect sizes for surprisal by lexical categories are caused by concreteness. In this set of model comparisons, we replace the interactions between surprisal and lexical category with interactions between surprisal and concreteness. We want to know whether interactions with concreteness may better predict AoA than interactions with lexical category.

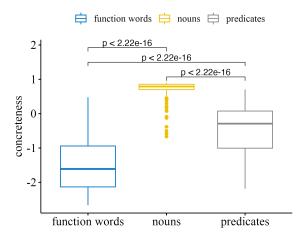


Figure 7: Concreteness distribution for different lexical categories

# F.1 Regression models

The regression models used here resemble those from the second approach in the main paper, though this time the interaction terms are concreteness, while lexical category is just another predictor.

- 1. The uni-gram concreteness model: AoA  $\sim$  concreteness \*(uni-gram surprisal + lexical category)
- 2. The residualised n-gram concreteness model: AoA  $\sim$  concreteness \*(uni-gram surprisal + n-gram residual average surprisal + lexical category), where n-grams are either bi-grams, tri-grams, or four-grams
- 3. The residualised LSTM concreteness model: AoA  $\sim$  concreteness \* (uni-gram surprisal + LSTM residual average surprisal + lexical category)

#### F.2 Results

Table 9 shows the LOO MAD results for our approach 1 model comparison using interactions with concreteness ratings instead of lexical category. Residual surprisal seems to have less of an effect overall with these interaction terms than with those reported in the main paper.

Table 9: Approach 1 model comparison results using interaction with concreteness by language

	LOO MAD <sub>[95% CI]</sub>					
Language	1gm-base	2gm-rd	3gm-rd	4gm-rd	LSTM-rd	
English (American)	<b>2.0</b> <sub>[1.87,2.14]</sub>	2.01 <sub>[1.87,2.14]</sub>	2.01 <sub>[1.87,2.14]</sub>	$2.01_{[1.88,2.15]}$	<b>2.0</b> <sub>[1.86,2.14]</sub>	
English (British)	$2.23_{[2.04,2.42]}$	$2.24_{[2.05,2.43]}$	$2.23_{[2.04,2.42]}$	$2.23_{[2.04,2.42]}$	$2.24_{[2.05,2.43]}$	
English (Australian)	$1.92_{[1.76, 2.09]}$	$1.93_{[1.76,2.09]}$	$1.93_{[1.77,2.1]}$	$1.94_{[1.77,2.1]}$	$1.93_{[1.77,2.1]}$	
German	$2.2_{[1.99,2.42]}$	$2.23_{[2.01,2.44]}$	$2.22_{[2.01,2.43]}$	$2.2_{[1.99,2.41]}$	$2.23_{[2.02,2.44]}$	
French (European)	$2.36_{[2.17,2.56]}$	$2.38_{[2.19,2.58]}$	$2.38_{[2.18,2.57]}$	$2.38_{[2.18,2.57]}$	$2.36_{[2.16,2.55]}$	
French (Quebecois)	$2.55_{[2.36,2.74]}$	$2.57_{[2.38,2.75]}$	$2.56_{[2.37,2.75]}$	$2.56_{[2.37,2.75]}$	$2.57_{[2.38,2.76]}$	
Spanish (European)	$2.52_{[2.34,2.71]}$	$2.53_{[2.34,2.71]}$	$2.53_{[2.35,2.72]}$	$2.54_{[2.36,2.72]}$	$2.52_{[2.34,2.71]}$	
Spanish (Mexican)	<b>2.09</b> <sub>[1.92,2.26]</sub>	$2.1_{[1.92,2.27]}$	$2.1_{[1.93,2.27]}$	$2.10_{[1.93,2.28]}$	$2.1_{[1.93,2.27]}$	
Mandarin (Beijing)	$1.91_{[1.78,2.04]}$	$1.9_{[1.76,2.03]}$	<b>1.9</b> <sub>[1.76,2.03]</sub>	$1.9_{[1.77,2.04]}$	$1.9_{[1.77,2.04]}$	
Mandarin (Taiwanese)	$2.99_{[2.77,3.22]}$	$3.01_{[2.78,3.23]}$	$3.01_{[2.78,3.23]}$	$3.01_{[2.78,3.23]}$	<b>2.95</b> <sub>[2.74,3.16]</sub>	

# G Approach 2 - N-gram predictability by language

Figure 8 shows the by language coefficient estimates for the mixed-effects models of approach 2 for residualised bi-gram, tri-gram, and four-gram models.

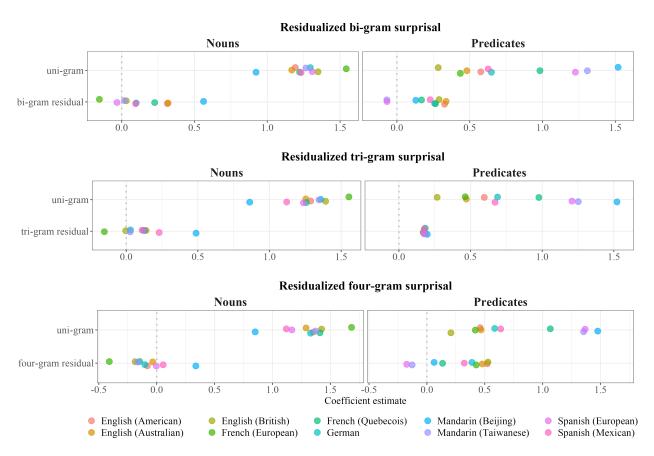


Figure 8: By-language random effects in Residualised n-gram mixed-effects models for second approach

# **H** Mandarin predicates

In experiment 1, we found that Mandarin predicates showed a negative effect for surprisal, less predictable predicates were learnt earlier; in experiment 2, we found the opposite effect. Experiment 2 used a subset (25) of the 185 Mandarin (Beijing) predicates and 119 (Mandarin (Taiwanese) ones. We hypothesized that this discrepancy was do to the specific items that were retained in experiment 2 and that there may be item specific effects. Here we test this hypothesis by fitting the model from experiment 1 to the item lists from experiment 2. As shown in Figure 9, the effect directions of surprisal are the same to experiment 2 results when using the model from experiment 1 on the Mandarin data from experiment 2, confirming this hypothesis.

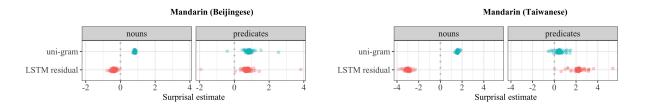


Figure 9: Coefficient estimates for surprisal values using the model from experiment 1 on the dataset from experiment 2 for Mandarin

We additionally provide a list of the unilemmas that were included and excluded from experiment 2.

**Included predicates from experiment 2:** big, bite, blue, clean (description), close, cold, cry, dirty, draw, eat, give, green, hit, hot, open (action), play, put, red, sad, take, walk, wash, write, yellow, rain.

Excluded Mandarin (Beijing) unilemmas (predicates) from experiment 2: a lot, angry, asleep, bad, black, blow, break, bring, broken, brush (action), bump, buy, can (auxiliary), cap, careful, catch, clap, clean (action), climb, comb (action), come, count, cover (action), cute, dance, down, drink (action), drive, drop, dry (description), dump, enter, fall, fast, feed, find, finished, flip, fly (action), full, go, good, happy, hard, have, hear, heavy, help, hide, high, hop, hug, hungry, hurry, hurt (description), jump, kick, kiss, lick, lift, light (description), little (amount), little (description), look, love, much, naughty, need, new, no, old, pee, press, pretend, pretty, pull, push, read, ride, run, salty, say, scared, shake, shout, sick, sing, sit, sleep, sleepy, slow, smart, smell, smile, soft, spill, spit, splash, stand, stick, stop, swallow, sweep, sweet, swim, tasty, tell, thirsty, throw, tickle, tired, touch, turn, ugly, wait, wake, wear, wet (description), white, wind, wipe, work (action).

Excluded Mandarin (Taiwanese) unilemmas (predicates) from experiment 2: angry, black, blow, bright, bring, call, careful, catch, clap, climb, come, cook, count, cut, cute, dark, drink (action), drop, feed, first, go, good, hang, hate, help, hide, high, hug, hungry, hurry, jump, kick, kiss, knock, listen, look, loud, love, new, no, pick up, pinch, poor, pour, press, pretend, pretty, pull, push, quickly, quiet, remember, ride, run, same, say, scared, scold, shake, sick, sit, slide (action), small, smell, smile (action), sour, squat, stand, stop, swallow, sweet, tear, tell, think, throw, tie, tired, touch, wanna, wear, wet, white, wipe, yes, yum yum.