RESEARCH ARTICLE

# The role of statistical learning in the L2 acquisition and use of nonadjacent predicate-argument constructions

Jiaqi Feng Guo[1] and Pascual Pérez-Paredes[2]

[1]Department of Chinese, School of Languages and Translation, University of Turku, Turku and [2]Department of Applied Linguistics, University of Murcia, Murcia, Spain
**Corresponding author:** Jiaqi Guo; Email: jiaqi.guo@utu.fi.

## Abstract

While statistical learning of adjacent constructions is well-documented in SLA, our knowledge of this cognitive mechanism concerning nonadjacent constructions remains limited. To address this, we investigated the acquisition of Mandarin predicate-argument constructions containing the preposition *duì*. Specifically, via a corpus-based approach, we probed whether learners' core predicate use within these nonadjacent constructions mirrors the patterns of frequency and contingency in their natural language input. Our findings show that learners' usage aligns with target language distributional regularities, which is consistent with statistical learning. However, our study underscores the necessity of going beyond a sole focus on distributional factors within learners' input to more fully comprehend L2 production choices and the intricacies of statistical learning. This includes examining variables that shape learners' exposure to input, such as input accessibility, proficiency, and prototypicality. Finally, we demonstrate the suitability of mixed-effects negative binomial regression to effectively address non-normality and overdispersion in linguistic data.

## Introduction

Statistical learning, which is the cognitive ability to implicitly discern and internalize regularities in one's environment, has been shown to play a vital role in second language acquisition (SLA) (Erickson & Thiessen, 2015; Rebuschat & Williams, 2012). Specifically, usage-based research, which focuses on grammatical constructions (*i.e.*, conventionalized form/meaning pairings) as the basic unit of linguistic analysis, has demonstrated that L2 learners are sensitive to distributional patterns in their input. These patterns include frequency (*i.e.*, how many times a construction is encountered),

contingency (*i.e.*, how predictable/dependent the occurrence of one construction is on another), and contextual dispersion (*i.e.*, how evenly/widely a construction is distributed across different input contexts such as reading or listening) (Gries, 2010a).

However, knowledge of this vital cognitive mechanism in the field of SLA remains subject to some limitations. A particularly important issue relates to a lack of research regarding nonadjacent constructions. These are constructions with collocating elements, which jointly create meaning but are separated by intervening linguistic items, such as a noun, noun phrase with modifiers, or complex clause. Although such intricate constructions appear likely to impose a greater cognitive burden which might interrupt learners' discernment of distributional patterns due to the challenges of processing long-distance dependencies (Isbilen et al., 2022), any impeding effects remain unattested (Hopp, 2023; Marinis et al., 2005). Furthermore, other limitations evident in prior statistical learning research relate both to the scope of the factors investigated (which are overly narrow), as well as the data analysis methods employed (which have tended to lack sufficient granularity to capture the complex patterns that typically characterize linguistic data). This study aims to address these issues by using corpus linguistic methods and mixed-effects regression to investigate how natural language input affects the acquisition of a type of L2 non-adjacent prepositional predicate-argument construction (PAC) in Mandarin Chinese.

## Statistical learning in SLA

The concept of statistical learning emphasizes that the input regularities which language learners are exposed to are crucial in shaping learners' ability to recognize, internalize, and use constructions (Bybee, 2010; Hoey, 2012). Specifically, the distributional factors of frequency, contextual dispersion, and contingency are seen to be of utmost importance. This review will commence by analyzing the literature on these factors, as this contextualizes the current study's focus.

Firstly, substantial evidence indicates that input frequency is vital for L2 acquisition. This evidence has been gathered in relation to various constructions such as morphemes (Ellis et al., 2016), words (Pellicer-Sánchez, 2016), collocations (Wolter & Gyllstad, 2013), idiomatic expressions (Martinez & Murphy, 2011), phraseology (Edmonds & Gudmestad, 2023), and schematic patterns (*e.g.*, verb-argument constructions) (Ellis & Ferreira–Junior, 2009). Data have been collected via numerous experimental methods, including eye-tracking (Winke et al., 2013), judgment tasks (Wolter & Gyllstad, 2013), and multimedia input (Peters, 2018), as well as via corpus-based studies, which have found strong correlations between high-frequency features in reference corpora and learners' output as derived from learner corpora. Finally, previous research has considered input frequency across different types of learners, such as those from different L1 backgrounds (Römer & Garner, 2019; Tono, 2004) and proficiency levels (Crossley et al., 2019). These studies have found that higher input frequency is consistently associated with better recognition, faster processing, and more accurate usage. However, it is important to note that where studies have looked beyond the isolated impact of frequency, results suggest other factors can moderate its impact. For example, Uchihara et al. (2019) found that variables like learner age and spaced learning influence the effects of frequency on vocabulary learning. Similarly, Stutterheim et al. (2021) found that attention allocation and L1 conceptual frames (*i.e.*, the mental structures through which speakers organize and interpret information) may also appreciably moderate frequency effects, while Eckerth and Tavakoli (2012) found that L2 word retention is more strongly influenced by the depth of word processing

rather than by input frequency. These findings highlight that frequency effects are nuanced and may interact with other learner factors.

Another apparently important variable in construction learning is contingency, which is the degree of reliability or strength of the association between a linguistic cue (such as a verb) and a specific grammatical structure or construction (Ellis & Ferreira–Junior, 2009; Gries & Ellis, 2015). It is typically measured through corpus-based collexeme analysis, which uses contingency tables to calculate the strength of association between linguistic elements. This analysis considers both the frequency of co-occurrence between elements, such as a verb and a construction, and their independent frequencies in the corpus. This helps to determine whether elements are significantly "attracted to" (more naturally associated) or "repelled from" (less frequently associated) each other (Gries, 2022a). Various statistical measures have been used to quantify contingency, including log-likelihood, Fisher's exact test, and delta P scores (Gries, 2010a).

Numerous studies have produced robust evidence of the importance of contingency, but, similar to frequency, have shown that its impact can be moderated. For instance, Murakami and Ellis (2022) explored the relationship between accuracy of grammatical morphemes in L2 writing and several distributional factors. Their findings indicate that contingency, defined as the token frequency relative to other forms that include the same lemma, is a robust predictor of morpheme accuracy. Edmonds and Gudmestad (2023) investigated the effect of immersion on L2 learners' sensitivity to cue contingency in phraseological development. Through studying anglophone Spanish and French learners during a study abroad, they discovered that at the end of their stay, those learners used more strongly associated combinations. This suggests that immersion enhances learners' sensitivity to cue contingency, leading to increased use of more strongly associated phraseological units. However, other studies have indicated a more limited effect. For example, Boone et al.'s (2022) longitudinal study of productive collocation knowledge found that while there was a significant effect of congruency (*i.e.*, when a multi-word construction has an equivalent in learners' L1) and prior productive vocabulary knowledge on collocation learning, association strength did not demonstrate a significant impact. Similar findings were also reported by Alzahrani (2021) and Siyanova and Schmitt (2008).

Finally, the third main distributional factor emphasized in statistical learning is contextual dispersion, which is the distribution of linguistic items across different texts/contexts. While fewer studies have examined dispersion compared to frequency and contingency, those that have corroborate its significance. Specifically, findings indicate that greater input dispersion enhances understanding and retention. For instance, Gries (2010b) found that wider distribution of linguistic features in different contexts leads to more effective and appropriate language use. Similarly, Çandarli (2021) found evidence suggesting that the dispersion of multiword constructions (MWCs) in an input corpus significantly affects MWC frequency in learners' essays.

Collectively, these studies suggest that distributional factors in learners' input play a crucial role in the acquisition of constructions via statistical learning. Specifically, it appears that frequency establishes familiarity, dispersion strengthens learners' ability to recognize patterns, and contingency assists the comprehension of constructions' form-meaning pairings and collocational tendencies. However, studies also show the effects are complex and subject to modulation. Moreover, some notable gaps in the literature remain.

An especially important issue relates to the construction types examined. For instance, while the extant research probes statistical learning across morphemes,

collocations, phrases, MWCs, and certain VACs, nearly all studies have investigated adjacent constructions (*i.e.*, where the collocating elements are contiguous with no open slot). Although it is important to acknowledge that some research has been conducted on non-adjacent dependencies in an L1 context or via artificial languages, given the differences between L1 and L2 acquisition and the limitations of artificial experiments, the findings cannot be assumed to have full applicability in SLA. Consequently, understanding of L2 statistical learning involving non-adjacent constructions remains limited. Specifically, it is unclear how much long-distance dependencies between collocating elements and the presence of intervening words/phrases affect learners' acquisition of these constructions and their lexical choices when planning for production, or if they can rely on explicit grammar knowledge to form associations between nonadjacent elements. The lack of research is an issue because nonadjacent constructions are crucial in many languages (Elder et al., 2017), often enhancing linguistic complexity and flexibility by facilitating the expression of nuanced functions, ideas, and relationships (Croft & Cruse, 2004; Goldberg, 1995).

Regarding why there is a lack of robust research on nonadjacent constructions, Gries (2022a) suggests that a key reason is methodological. Specifically, until recently, the off-the-shelf natural language processing (NLP) programs and corpus tools which researchers depend on to extract construction usage patterns have tended to rely on target items being near each other within the data. Therefore, they have only been able to extract adjacent constructions with sufficient accuracy to ensure reliable analysis (McEnery & Brezina, 2022; Pérez-Paredes, 2020). This limitation has affected both corpus and experimental studies, which generally also derive their reference frequency from corpora. Consequently, scholars (*e.g.*, Kyle, 2021; Meurers, 2015; Yan & Liu, 2022) have called for more advanced customized NLP programs.

Another weakness of much L2 statistical learning research is that while studies have started to move beyond a sole focus on distributional factors in learners' input by examining how such factors can be modulated by other variables, the range of additional factors investigated remains limited (Deshors & Gries, 2022). The interaction between various factors has also been underexplored (Wulff, 2020). This is problematic because, as has been shown in numerous other areas of SLA, the tendency to focus on a limited number of isolated factors risks overlooking the multifaceted reality of language acquisition by producing findings which oversimplify or overstate the effects of target variables (Ibbotson, 2013; McManus, 2024).

An area especially deserving of analytical attention concerns how variables that shape learners' exposure to input and their ability to process it possibly interact with features of input such as frequency, contingency, and contextual dispersion. Exposure-shaping factors include: 1) total duration of L2 exposure (which can be inferred from proficiency level); 2) the timing of a learner's initial contact with a construction; and 3) the complexity of the constructions they encounter. Investigating these variables will help illuminate potentially important nuances in the operation of statistical learning such as: do high-frequency items in the target language have the same influence on L2 acquisition if encountered at different times in the learning journey; does sensitivity to distributional factors change as learners' proficiency progresses; and does the timing of initial contact with a construction lead to the formation of strong prototypes that can either facilitate or override distributional effects. Addressing these currently less-well understood questions will help refine our understanding of whether L2 statistical learning is a relatively stable capacity, or if it exerts a more complex/dynamic effect.

Finally, in recent years, SLA research has increasingly used generalized linear mixed models (GLMMs) in order to take account of variability across participants and

linguistic items, with mixed-effects logistic regression commonly used for binary outcomes (*e.g.*, Baayen et al., 2008; Johnson, 2009). However, little use has been made of two types of regression—Poisson regression and negative binomial regression—that were developed to handle types of non-normal count data (*e.g.*, corpus frequencies of words), which have been of interest to SLA researchers for decades (Winter & Bürkner, 2021). Basic distinctive characteristics of count data are as follows: the data points are discrete (*i.e.*, there are no decimals); and they cannot be negative, but they can be very large—potentially infinitely large, at least in theory (McElreath, 2020).

Given the need to explore multiple interacting factors and to employ more appropriate statistical methods for linguistic count data, several influential scholars have argued that claims regarding statistical learning in SLA must be treated cautiously and more detailed investigation is required (Eskildsen, 2022; Gries, 2022b). We respond to this call. In this study, we examine L2 statistical learning of a type of non-adjacent predicate-argument construction (PAC) in Mandarin Chinese, namely, PACs containing the preposition *duì*. Henceforth, we refer to them as *duì*-constructions.

The goal of our study is to probe whether learners' core predicate use within *duì*-constructions correlates with the usage patterns found in their natural language input. Alongside the distributional factors of input frequency and contingency, we also probe how several factors that shape learners' exposure to input co-determine their usage. The additional factors are fully elaborated in the Method section. Finally, by using mixed-effects count data regression, our study demonstrates the applicability of this method in SLA research.

The study addresses two research questions:

1. Do L2 learners of Chinese display evidence of statistical learning in relation to their core predicate choices in *duì*-constructions?
2. How do input exposure-shaping factors and distributional factors co-determine L2 learners' core predicate choices in *duì*-constructions?

## Method

This section details the study's linguistic focus, data sources, extraction processes, and analytical approach.

### Duì-*constructions*

The linguistic focus of this study is *duì*-constructions, a type of prepositional predicate-argument construction. *Duì*-constructions are semischematic and combine one fixed component (the preposition *duì*, which can be translated as "to/toward," "at/on," "about/in terms of," or "for/as for," depending on the context) with three open slots: a subject argument (X), an oblique argument (Y), and a core predicate. Thus, the form of a *duì*-construction is:

$$X + duì + Y + \text{core predicate}$$

The core predicate slot of a *duì*-construction may be filled by either a verb (*e.g.*, *shuō* "say", *pīpíng* "criticize") or an adjective (*e.g.*, *hǎo* "to be kind"/"to be beneficial", *shīwàng* "to be disappointed"). Unlike English, Chinese adjectives can function directly as core predicates without requiring a copula verb (*e.g.*, "is," "am," or "are") (Huang, 2006;

Thompson & Tao, 2010). For example, the Chinese sentence *wǒ duì tā hěn hǎo*, which can be literally translated as "I toward him very good", contains the adjective *hǎo* ("good") that functions as the core predicate without requiring a copula verb "am."

A defining characteristic of *duì*-constructions is that the preposition *duì* combines with the construction's core predicate to establish the semantic relationship between X and Y and to collectively create the construction's meaning/function. This conceptualization of *duì*-constructions is informed by Goldberg's (1995, 2015) construction grammar approach and the predicate-argument structure framework (Langacker, 1987; Levin & Hovav, 2005), which provide the theoretical tools to analyze how specific predicates constrain and license particular arguments within these constructional patterns.

Our previous comprehensive analysis of *duì*-constructions (Guo, 2023) provides empirical evidence supporting this constructional approach and demonstrates that *duì* combines with different core predicates to produce six distinct types of constructional meaning/function. To develop this function taxonomy, we analyzed over 8000 instances of *duì*-constructions used by native speakers of Chinese derived from a native Chinese reference corpus (Guo, 2023). This is the same corpus used in this study (see Section 3.2). The taxonomy was subsequently verified by two independent Chinese linguistics experts. The six identified *duì*-construction functions are detailed in Table 1. As will be explained later, these functions relate to a vital part of this study's analysis as we seek to determine whether L2 learners use *duì*-constructions to express similar functions as native Chinese speakers.

There are two further points that it is important to discuss regarding *duì*-constructions. The first relates to the core predicates that can be used in the construction, and the second to the non-adjacent nature of the construction.

Firstly, as noted above, *duì*-constructions are semischematic. Like other semi- or fully schematic constructions, they allow speakers to use different lexis (Kay & Fillmore, 1999). However, in contrast to the "open-choice principle" (Sinclair, 1991) that suggests the predicate slot in a construction could theoretically accept any semantically

**Table 1.** Six duì-construction functions

| Function | Description | Example |
|---|---|---|
| Target-Action | Involves X (an agent) performing an action toward Y (a target or goal). | S1: 我对他说 (*wǒ duì tā shuō*) "I said to him". |
| Gesture-Attitude | Describes X (an experiencer) expressing an emotion or gesture toward Y (a recipient). | S2: 我对他很好 (*wǒ duì tā hěn hǎo*) "I am kind to him". |
| Transformative-Scope | Delineates X (an agent) performing an action that targets a defined group/scope Y. | S3: 政府对旧城区进行了改造 (*zhèngfǔ duì jiù chéngqū jìnxíngle gǎizào*) "The government carried out renovations on the old urban areas". |
| Evaluative-Perspective | Describes X (a theme) being evaluated as having an effect on Y (a beneficiary or affected party). | S4: 吸烟对身体有害 (*xīyān duì shēntǐ yǒuhài*) "Smoking is harmful to one's health". |
| Psychological-Reaction | Describes X (an experiencer) experiencing a psychological response to Y (a theme). | S5: 我对他很失望 (*wǒ duì tā hěn shīwàng*) "I am very disappointed in him". |
| Thematic-Relation | Represents X (an agent) performing a non-transformative action in relation to topic Y. | S6: 他们对这件事进行讨论 (*tāmen duì zhè jiàn shì jìnxíng tǎolùn*) "They conducted a discussion about this matter". |

compatible verb or adjective, *duì*-constructions exhibit systematic constraints on which verbs and adjectives can be used in the predicate slot. These constraints manifest as collocational preferences (Stefanowitsch & Gries, 2003)—statistical tendencies for certain lexis to occur more frequently in the construction than others. These collocational preferences are a central concept in usage-based approaches to language learning (Bybee, 2008), where such statistical patterns represent conventionalized associations formed through repeated language exposure and use in the speech community, rather than being predetermined by strict grammatical rules. For example, while both *shuō* ("say") and *tán* ("talk/discuss") are speech verbs and are grammatically acceptable in *duì*-constructions, native speakers show a systematic preference for *shuō* (Guo, 2023). This notion of collocational preferences is crucial to this study's analysis as we seek to examine whether L2 learners mirror native speakers' tendency to use certain verbs and adjectives in *duì*-constructions more than other available lexical choices.

The second key point to discuss is *duì*-constructions are characterized by a key structural feature: nonadjacency between *duì* and its collocating predicate. Specifically, *duì* and the construction's core predicate are always separated by the oblique argument —either a noun, a noun phrase, or a complex clause, and sometimes additional adverbial modifiers. Therefore, *duì*-constructions, which follow the word order [subject argument + preposition + oblique argument + core predicate], differ from many prepositional PACs in English examined in the field of SLA (*e.g.*, "look at the boy" or "depend on your friends") that follow the word order [subject argument + core predicate + preposition + oblique argument], where the preposition and predicate are generally adjacent (Kyle et al., 2021; Römer, 2019).

Figure 1 below illustrates the nonadjacency of *duì*-constructions. It shows a *duì*-construction extracted from the native Chinese reference corpus used in this study (see Section 3.2). The construction exemplifies the Target-Action function, where *tā* ("she") as agent X performs the action of speaking (*shuō*) directed toward *zhèng yào chūmén de érzi* ("the son who was about to leave") as target Y. Notably, *duì* and its collocating verb *shuō* are separated not only by this complex noun phrase, but also by a nonargument element—the adverbial modifier *qīngshēng de* ("softly"). Such a gap is not unusual. Our analysis of all *duì*-constructions extracted from this study's reference corpus shows a median gap of 3 characters, with the longest being 137 characters. Such gaps underscore the processing challenges of nonadjacent constructions, which may hinder statistical learning (Hopp, 2023; Marinis et al., 2005). Thus, they make *duì*-constructions a suitable case study for exploring L2 statistical learning of nonadjacent constructions.

| Chinese | 她 | 对 | 正 | 要 | 出门 | 的 | 儿子 | 轻声 | 地 | 说 |
|---|---|---|---|---|---|---|---|---|---|---|
| Pinyin | tā | duì | zhèng | yào | chūmén | de | érzi | qīngshēng | de | shuō |
| POS | pron | prep | adv | aux | verb | part | noun | adv | part | verb |
| English | she | to | just | about | leave | (mod) | son | softly | (mod) | speak |

**Translation:** She softly said to her son who was just about to leave.

**Figure 1.** Example duì-construction.

### Data sources

This study draws *duì*-construction data from two corpora: a native Chinese reference corpus and an L2 learner corpus that was derived from students living in China for an extended period. We now describe these corpora and discuss their suitability for investigating statistical learning.

The reference corpus is a balanced multi-genre L1 Chinese corpus compiled by the Beijing Language and Culture University Corpus Centre (Link). It consisted of 1.06 billion characters ≈ 606 million words/tokens at the time of our research. According to the compilers (Xun et al., 2016), it aims to reflect the entire range of actual language use in contemporary Chinese society. It is derived from both formal sources such as newspaper articles (28% of the total) and literature (23%), and also from informal everyday language use drawn from micro-blogs and discussions (24%). This corpus serves as a proxy for the natural language input that L2 learners of Chinese residing in China are likely to encounter in daily life, an input source which has been shown to be vital to L2 development, especially for those undertaking immersive study (Monteiro et al., 2020). Henceforth, it is the "reference corpus."

Although we acknowledge that using a reference corpus to represent learners' natural language input has limitations because the data it contains is not actual input and cannot fully replicate the limitless variety of learners' experiences (Gass et al., 2020), it nonetheless remains one of the most feasible means of investigating the effects of natural language input on L2 acquisition. This is because it circumvents the practical impossibility of tracking a learner's real input, as well as some of the main weaknesses of language experiments conducted in controlled environments (Bley-Vroman, 2002; Kartal & Sarigul, 2017). In particular, by providing a systematic collection of authentic language data from a wide variety of sources and genres, an adequately sized and representative reference corpus can approximate the diverse types and regularities of input that learners are likely to encounter to a sufficient degree for analytical purposes such as ours (Paquot & Gries, 2020; Wolter & Gyllstad, 2013).

The second corpus is the 1.2 million-word Guangwai-Lancaster Chinese Learner Corpus (Chen & Xu, 2019), available via Sketch Engine (Link). It includes data from 1,473 learners from 106 countries who resided in China while undertaking a degree program at Guangdong University of Foreign Studies (Guangwai). The data is split across three learner proficiency levels—beginner, intermediate, and advanced—established by Guangwai. The written data (comprising 52% of the total) and spoken data (48%) were collected from exams and tutorial sessions. Henceforth, it is the "learner corpus."

It is important to note that limited access to the full corpora data impacted our analysis. Specifically, this study does not incorporate contextual dispersion as a factor. This is because we could only access concordance lines retrievable through the corpora's online platforms, which do not support the automatic calculation of dispersion metrics.

### Data extraction

To accurately extract the core predicate within *duì*-constructions from both corpora, we followed several steps. Initially, all instances of *duì*-constructions were extracted using the string-searching techniques provided by the corpora platforms. Specifically, regular expressions (*e.g.*, [对/p * v]) were employed in the reference corpus, and corpus query language (CQL) was used in Sketch Engine for the learner corpus, with the query

[word="对" & tag="P"] [tag!="PU"]{0,20} [tag="V.*"] within <s/>. However, using the off-the-shelf analysis tools provided by the two corpora resulted in a significant percentage of errors, including incorrect collocations (67% of the collocating verbs and adjectives were incorrectly identified), and irrelevant constructions and duplicates (22% of all extracted *duì*-constructions). The errors arose because generic tools are not designed to handle nonadjacent structures. To address these challenges, we developed a custom NLP script in Python using dependency parsing, semantic dependency parsing, and semantic role labeling modules from Harbin Institute of Technology's Language Technology Platform (Che et al., 2021). This approach allowed us to correctly identify the collocational relationships between *duì* and the core predicate of each construction based on their semantic connections rather than just their surface proximity. This script (LTP_Verb_Extraction.py) is available on IRIS (Link). With these sophisticated techniques, we removed non-*duì*-constructions and duplicates, extracted the correct core predicates (97% accuracy rate), and calculated the overall token and type frequency of the unique verbs and adjectives across different learner levels and among native speakers.

### Data selection and annotation

In total, we extracted 8,364 unique verbs and adjectives as the core predicates from *duì*-constructions in the reference corpus and 135 unique predicates from the learner corpus. Of the 135 learner core predicates, 7 have multiple (*i.e.*, two or more) functions in *duì*-constructions, resulting in 143 entries. All of these 143 learner core predicates were also present in the reference corpus. To ensure direct comparability, we focused our analysis on these 143 common core predicates. Although we considered including all the unique core predicate usage in the reference corpus that were not present in the learner corpus, we excluded them due to potential model convergence issues and bias from zero-inflated data points (Winter & Grice, 2021). However, we added four high-frequency verbs from the reference corpus that were absent from the learner corpus, resulting in 147 entries. These four additional high-frequency verbs were included to ensure balanced analysis and investigate potential learner avoidance.

These 147 core predicates were then annotated. This process involved annotating six key factors: frequency in both corpora, contingency, function, predicate semantics, proficiency, and accessibility. The first two of these are distributional factors, which are traditionally investigated in research into statistical learning:

Frequency in our study refers to the occurrence rate of each core predicate in *duì*-constructions across the reference corpus, and the learner corpus split by three proficiency levels. To account for corpus size differences, corpus sizes for each proficiency level and the reference corpus size were included as offset terms.

Contingency, which refers to the association strength between core predicates and *duì*-constructions, was calculated using a 2-by-2 simple collexeme analysis contingency table (Table 2) adapted from Hilpert (2014). This table compares core predicate frequencies in *duì*-constructions versus other constructions in the reference corpus, with log-likelihood scores determining the statistical significance of these associations.

The next four factors were included to probe how learners' exposure to input and their ability to process it affect their acquisition, and potentially modulate the impact of the distributional factors. Each of these variables is now explained.

The first input-exposure related factor is L2 Proficiency, which directly links to the quantity and type of linguistic input that learners are likely to encounter and assimilate.

**Table 2.** Core predicate association contingency table for *duì*-constructions

|  | In *duì*-constructions | Not in *duì*-constructions | Row Totals |
|---|---|---|---|
| **Core Predicate X** | Frequency of predicate X in *duì*-constructions | Frequency of predicate X not in *duì*-constructions | Total frequency of predicate X |
| **Not Core predicate X** | Frequency of other predicates in *duì*-constructions | Frequency of other predicates not in *duì*-constructions | Total frequency of other predicates |
| **Column Totals** | Total frequency of *duì*-constructions | Total frequency of other constructions | Total corpus frequency |

The expectation is that as they gain proficiency, learners are more likely to access and benefit from varied and complex language input, which affects how frequently they encounter specific language features and structures in authentic contexts (DeKeyser, 2007). This, in turn, may facilitate their discernment of input patterns. In this study, proficiency levels (beginner, intermediate, and advanced) were based on the Guangwai classification. These levels reflect cumulative exposure to Chinese in a structured educational setting at Guangwai. While these categories offer a useful framework, they should be interpreted cautiously, as the Guangwai classification may not fully correspond to international standards such as the CEFR.

The second input-exposure related factor is constructional Function, which refers to the meaning that a speaker tries to convey with a construction by employing different semantic/thematic relations within it. The reason functions are particularly relevant to input-exposure is because they are not rules, but rather are generalized cognitive structures (mental representations) that are time-sensitive and emerge from speakers' interactions with language (Bybee, 2010). Studies have shown that when learners initially associate a specific word or structure with a particular meaning or function, they tend to form a prototypical association with this initial form/function and favor it even when alternatives are available—a phenomenon known as the pre-emption effect (Trahey & White, 1993). Similarly, the prototypical or most commonly encountered meaning of a multifunctional construction is heavily influenced by learners' timing of initial contact and most frequent interactions with it (Granena & Long, 2013). This conceptualization can significantly influence learners' understanding and usage of a construction.

In this study, we adopt the six distinct communicative functions of *duì*-constructions, outlined in Section 3.1, which include: Target-Action, Gesture-Attitude, Transformative-Scope, Evaluative-Perspective, Psychological-Reaction, and Thematic-Relation. For annotation, each unique core predicate was manually assigned a function, by the first author and an independent expert in Chinese linguistics. A small number of these core predicates were annotated with multiple functions to reflect their ability to express different semantic relationships when combined with *duì* and different arguments. When assigning functions, each verb or adjective was examined in context by analyzing 25 random concordance lines (50 for multifunctional lexical items) from the reference corpus to understand its typical usage and semantic relations within the *duì*-constructions. Inter-rater reliability for primary raters reached Cohen's Kappa = 0.798 ($p < .001$), showing good agreement. Discrepancies were resolved through discussion, resulting in 94% agreement. For the remaining 6%, a third expert provided external validation, achieving a final Fleiss' Kappa of 0.864 ($p < .001$) among all raters.

The third input-exposure factor, Predicate Semantics, categorizes each item according to its primary meaning(s) to explore the range of contexts in which both native

speakers and learners use *duì*-constructions. Using *A Thesaurus of Modern Chinese* (Su, 2013), we classified verbs and adjectives into seven distinct semantic classes: Communication, Social, Functional, Psych, Manner, Attribute, and Physical. This classification was designed to examine whether early exposure to a specific verb or adjective within certain semantic classes might lead learners to prototype and favor these lexical items. However, we excluded this factor from our final model due to collinearity issues (see Section 3.5).

The fourth input-exposure related factor is what we conceptualize as Accessibility. This is a hybrid notion of the timing of initial encounter with a specific linguistic item (Nation, 2013) and the concept of word difficulty (Hashimoto & Egbert, 2019). Together, these aspects not only reflect the complexity of a word but also how familiar learners are with it and the typical proficiency level at which they first encounter it. The examination of accessibility alongside distributional factors was designed to clarify whether learners more readily encode and acquire simpler, more familiar words in nonadjacent constructions compared to more complex, less familiar ones.

We have operationalized Accessibility through a 6-level scale by evaluating each of the 147 verbs and adjectives using 4 criteria: (1) HSK level (the 9-level standardized Mandarin Chinese proficiency framework created by China's Ministry of Education); (2) timing of first appearance in teaching materials (as derived from the Xiamen University Textbook Corpus, which comprises data from many of the most widely used textbooks for teaching Chinese as a foreign language in mainland China); (3) native speaker usage frequency derived from the reference corpus used in this study (described in Section 3.2); and (4) expert evaluation by three specialists—the first author (with over 15 years of Chinese teaching experience) and two independent Chinese linguistics researchers. Through this comprehensive process, we have assigned each verb or adjective to one of six accessibility levels. The resulting scale organizes these core predicates from the easiest and earliest encountered, to more difficult, typically later encountered ones. We provide the full classification of all 147 core predicates in Appendix A. This details the rationale for every classification decision. While we are confident our systematic procedure has resulted in a robust classification, we acknowledge that measuring Accessibility is a complex challenge and that this remains a partly subjective process.

To illustrate our annotation approach, Table 3 presents an example core predicate with annotations for all six factors:

**Table 3.** Example core predicate annotation

| Factor | Annotation |
|---|---|
| Core Predicate | 道歉 |
| Translation | Apologize |
| Function | Target-Action |
| Semantics | Communication |
| Accessibility | 3 |
| Contingency | Attracted |
| Frequency | |
|   **Beginner Frequency** | 0 |
|   **Intermediate Frequency** | 5 |
|   **Advanced Frequency** | 1 |
|   **Native Frequency** | 135 |

*Modeling and data analysis*

Our modelling and analysis process followed a systematic sequence. First, using RStudio (Version 2023.06.0+421) with ggplot2, we examined the distribution of predicative verb and adjective frequencies in the reference corpus and the learner corpus split by proficiency level to help inform our choice of model. Figure 2-A depicts that core predicate frequency in *duì*-constructions in both the reference and learner corpora follows a Zipfian distribution. This means a small number of verbs and adjectives has a very high usage frequency, while the rest exhibit a low usage frequency, resulting in long-tail distributions. Figure 2-B illustrates substantial overdispersion in the learner corpus across all three proficiency levels. That is, learners at the same proficiency level use core predicates within *duì*-constructions very differently from each other. The long error bars indicate that even within each proficiency level, some learners use certain verbs or adjectives frequently while others use them rarely. If core predicate usage followed a normal distribution, we would expect most learners at the same level to show similar patterns, with fewer extreme differences. However, the data show much more individual variation.

Given the non-normal and highly skewed nature of our count data, we adopted a Generalized Linear Model (GLM) framework. This approach allowed us to effectively capture and analyze the observed frequency patterns (Winter & Bürkner, 2021). While a Poisson model is often used for count data, it assumes equal mean and variance. We therefore decided it was unsuitable for our dataset. Instead, we chose a negative binomial (NB) model, which includes an additional dispersion parameter to better capture the variability in core predicate frequencies. This allows it to better accommodate highly skewed frequency patterns with long tails—characteristics that align well with our dataset.

As a second step, before validating our proposed factors in model testing, we assessed multicollinearity among the six factors detailed in Section 3.4 using Cramér's V. This analysis revealed high collinearity between Function and Predicate Semantics ($V = 0.75$) and moderate collinearity between Predicate Semantics and Contingency ($V = 0.40$). To address these issues, we tested the model with all factors and evaluated the impact of removing Predicate Semantics on model stability and performance. Based on these evaluations, Predicate Semantics was excluded from the final model to reduce redundancy and resolve multicollinearity issues (Gelman et al., 2012).

Thirdly, to further validate our choice of NB models with our refined factors (*i.e.*, excluding Semantics), we systematically tested several models using glmmTMB (Brooks et al., 2017). Comparing AIC and BIC values, the NB-GLMM (*AIC*=1230.15, *BIC*=1299.67) offered the best fit and interpretability (Figure 3), effectively managing high dispersion without risking model misfit (McElreath, 2020).

Fourthly, after selecting the optimal NB-GLMM, we tested for key two-way interactions (*e.g.*, Native_Frequency × Proficiency) to explore potential interaction effects. However, none of the interaction models converged successfully, likely due to limited data support for complex interaction structures. Consequently, we retained the base model without interactions. Model diagnostics conducted using DHARMa confirmed that the base model reliably captures the underlying patterns in L2 core predicate choice. Specifically, as shown in Figure 4, the QQ and residual plots revealed no misfit ($p = 0.12$) or over-/under-dispersion ($p = 0.88$), which suggests that the model fits the data well.

To implement the final model, we used the R syntax shown in Figure 5, specifying key predictors and accounting for corpus size through offset terms.
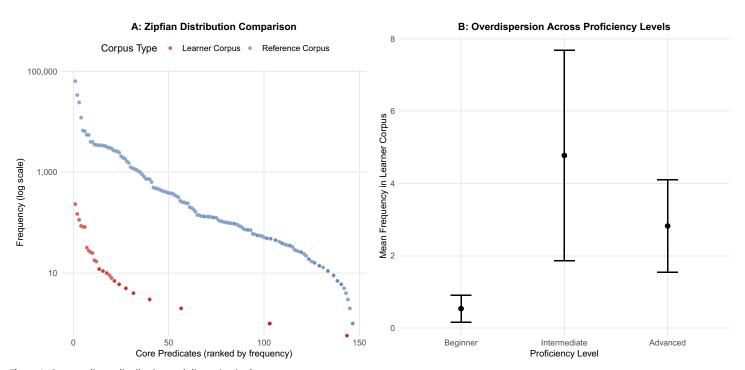
**A: Zipfian Distribution Comparison**

Corpus Type    ● Learner Corpus    ● Reference Corpus

**B: Overdispersion Across Proficiency Levels**



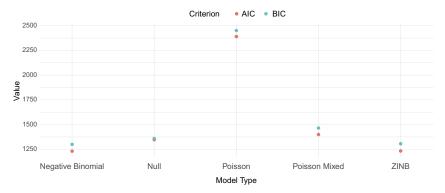**Figure 2.** Core predicate distribution and dispersion in the two corpora.

**Figure 3.** Model comparison.

Finally, to examine potential post hoc interaction effects, we applied postestimation Expected Marginal Means (EMMs) for three-way interactions. As Gelman et al. (2012) argue, multilevel models reduce the need for multiple comparison adjustments by utilizing partial pooling, which balances interpretability and statistical power. To address potential Type I errors due to multiple testing, we applied the Benjamini-Hochberg (BH) correction to the $p$ values.

## Results

This section presents our results concerning *duì*-construction usage by native speakers and L2 learners. We start with descriptive analysis of core predicate frequencies and their distribution in the reference and learner corpora. Subsequently, we examine the influence of various factors on L2 core predicate choice, as determined by our NB-GLMM analysis. Finally, we assess how the distributional and input-exposure factors interact to affect L2 production.

### *Predicate usage within duì-constructions*

Table 4 summarizes the *duì*-construction usage frequency and the unique verbs and adjectives in the predicate slot in both the reference and learner corpora. The learner corpus data is also split by proficiency level.

As described above, a notable feature of the core predicates used in *duì*-constructions in both the reference and learner corpora is that they follow a Zipfian distribution. In the reference corpus, 8364 unique core predicates were extracted from *duì*-constructions. However, only three—*shuō* ("to speak"), *yǒu* ("to have"), and *jìnxíng* ("to conduct")—account for over one third of the total core predicate usage. Their frequencies are 61,540 (101.45 pmw), 39,361 (64.89 pmw), and 33,797 (55.71 pmw), respectively. The learner corpus shows a similar Zipfian trend, with its top four core predicates representing 62% of the total core predicate usage. Specifically, the four core predicates are *shuō* ("to speak"), *hǎo* ("to be nice/beneficial"), *yǒu* ("to have"), and *gǎn* ("to feel"), with frequencies of 231 (179.87 pmw), 196 (152.62 pmw), 174 (135.49 pmw), and 148 (115.24 pmw), respectively. This underscores that a small subset of verbs and adjectives is heavily favored.
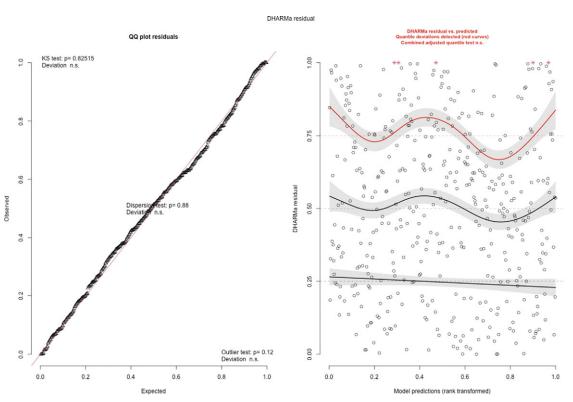
DHARMa residual

**QQ plot residuals**

**Figure 4.** DHARMa diagnostic plots for model fit and residual analysis.

```
Final_model <- glmmTMB(L2_Core_Predicate_Use ~ Native_Frequency + Accessibility + Function + Contingency +
                       Proficiency + offset(log_corpus_size) + offset(log_native_size) +
                       (1 | Core_Predicate), family = nbinom2, data = model_data_long)
```

**Figure 5.** Final NB-GLMM model specification in R.

**Table 4.** Comparison of duì-construction usage frequency and unique predicates

| Corpus | Corpus Size | Raw Frequency | *Duì*-constructions (pmw) | Unique Predicates |
|---|---|---|---|---|
| **Reference** | 606,060,606 | 392,926 | 648 | 8,364 |
| **Learner** | 1,294,228 | 1201 | 928 | 146 |
| Beginner | 288,534 | 80 | 277 | 21 |
| Intermediate | 627,227 | 699 | 1,114 | 89 |
| Advanced | 378,467 | 422 | 1,115 | 91 |

*Note*: pmw refers to normalized frequency per million words.

### Factors that influence L2 core predicate choice

Our NB-GLMM model analyzed the factors which influence learners' core predicate choices in *duì*-constructions, including the fixed effects: Native (frequency in reference corpus), Contingency, Accessibility, Function, and Proficiency, along with random intercepts for Core_Predicate (Figure 5). The model fit was assessed using the dispersion parameter ($\theta = 1.29$), and random effects variance for Core_Predicate was estimated at 0.653 ($SD = 0.808$). The model identified all factors as statistically significant at the group level ($p < .05$), based on Type III ANOVA with BH correction. However, to provide a more detailed view of how specific categories within these group factors (*e.g.*, the six Accessibility levels) influence predicate choices, we used the tab_model function in R for our model.

As shown in Table 5, the model uses Incidence Rate Ratios (IRRs) to show how specific categories within each factor influence core predicate choices. Each IRR for a categorical factor compares the frequency of the use of a core predicate relative to its respective reference category, including Accessibility Level 1, Proficiency (Beginner), Function (Thematic-Relation), and Contingency (Attracted). An IRR greater than 1 indicates an increased probability of a core predicate being used compared to its reference category, while an IRR less than 1 suggests a decreased probability.

The results in Table 5 suggest that most factor levels demonstrate statistically significant associations with learners' core predicate choices in *duì*-constructions. Turning first to the distributional factors. Native Frequency has statistical significance ($p < .001$), which indicates a reliable relationship with L2 predicate usage. While the IRR of 1.00 might suggest no effect, the model summary statistics reveal a small but reliable positive relationship. This is indicated by the coefficient for Native Frequency ($7.295 \times 10^{-5}$, $SE = 1.290 \times 10^{-5}$). This suggests that although small changes in Native Frequency (*e.g.*, one additional occurrence of a verb or an adjective in the native usage) have minimal impact on learners' core predicate choices, if Native Frequency increases by a larger amount, such as 10,000 occurrences, this small effect becomes meaningful. The second distributional factor, Contingency, is also statistically significant ($p = .042$), with core predicates categorized as Repelled being 44% less likely to be chosen by learners than core predicates categorized as Attracted.

Among input exposure-related factors, Accessibility is a highly significant group-level predictor of learners' predicate choices ($p < .001$). However, specific levels within Accessibility show distinct effects. For instance, verbs and adjectives at Accessibility level 6 (*IRR = 0.14*, 95% confidence interval [*CI*] [*0.05–0.39*]) are 86% less likely to be

**Table 5.** Negative binomial GLMM output

| Predictor | IRR (95% CI) | Adjusted p value |
|---|---|---|
| (Intercept) | 0 (0–0) | < .001 *** |
| Native Frequency | 1 (1–1) | < .001 *** |
| Accessibility (Reference: Level 1) | | |
|   Accessibility 1 (ref) | 1 | — |
|   Accessibility 2 | 0.70 (0.31–1.61) | .408 |
|   Accessibility 3 | 0.33 (0.13–0.80) | .023 * |
|   Accessibility 4 | 0.20 (0.09–0.45) | < .001 *** |
|   Accessibility 5 | 0.20 (0.07–0.58) | .006 ** |
|   Accessibility 6 | 0.14 (0.05–0.39) | < .001 *** |
| Proficiency (Reference: Beginner) | | |
|   Beginner (ref) | 1 | — |
|   Intermediate | 4.17 (2.71–6.43) | < .001 *** |
|   Advanced | 6.25 (4.01–9.75) | < .001 *** |
| Function (Reference: Thematic-Relation) | | |
|   Thematic-Relation (ref) | 1 | — |
|   Evaluative-Perspective | 2.17 (1.11–4.23) | .035 * |
|   Psychological-Reaction | 2.38 (1.28–4.42) | .012 * |
|   Target-Action | 0.68 (0.34–1.36) | .313 |
|   Transformative-Scope | 0.63 (0.26–1.56) | .341 |
|   Gesture-Attitude | 1.50 (0.77–2.92) | .293 |
| Contingency (Reference: Attracted) | | |
|   Attracted (ref) | 1 | — |
|   Repelled | 0.56 (0.33–0.95) | .042 * |

Note: IRR = Incidence Rate Ratio; CI = Confidence Interval. $p < .05$ = *; $p < .01$ = **; $p < .001$ = ***. Reference levels (shown as IRR = 1) were determined by R's default factor handling for ordered variables (Level 1 for Accessibility, Beginner for Proficiency) and alphabetical ordering for unordered variables (Function and Contingency).

used than those at level 1. Function is also statistically significant ($p < .001$, based on the overall model), with certain functions preferred by learners relative to others. Specifically, verbs or adjectives that collocate with *duì* to form the Psychological-Reaction (*IRR* = 2.38) and Evaluative-Perspective (*IRR* = 2.17) functions are used more than twice as frequently by learners than verbs or adjectives that form the Thematic-Relation function baseline. This underscores the importance of incorporating constructional function into our analysis, as learners' preferences reflect an internalized sensitivity to prototypical usage patterns. Finally, Proficiency is statistically significant for learners' core predicate choice. The results suggest that intermediate learners (*IRR* = 4.17) are nearly four times more likely to use *duì*-constructions than beginners are, while advanced learners (*IRR* = 6.25) are over six times as likely.

Figure 6 shows the effects of Native Frequency, Accessibility, Function, Contingency, and Proficiency on learners' likelihood of choosing a particular core predicate in *duì*-constructions. It depicts the IRR findings, highlighting the strength and direction of each factor's effect.

Although the results in Table 5 highlight the significance of the investigated factors in shaping L2 core predicate choice, how these factors interact requires further investigation. For instance, although Native Frequency shows a positive association with L2 core predicate choice (see Figure 6, top left), the wide confidence intervals suggest considerable variability. For example, when the Native Frequency is 50,000, the model predicts an L2 frequency of 18.78, with a CI ranging from 5.54 to 63.71 (Table 6). This indicates that while high native frequency correlates with increased L2 core predicate use, the substantial variability in predictions suggests that its influence is both cumulative (becoming more meaningful at higher frequency values) and is
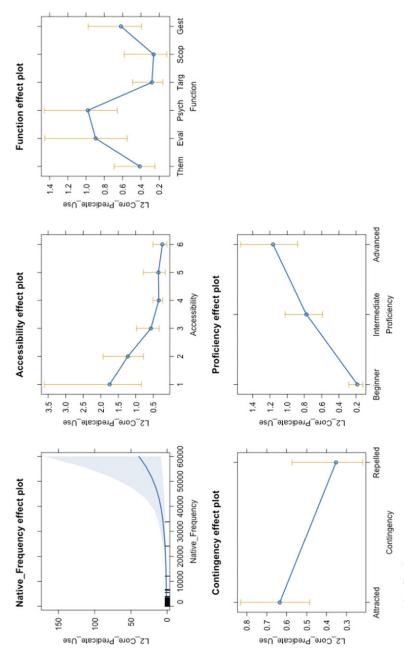
**Figure 6.** NB-GLMM visualization.

**Table 6.** Predicted L2 core predicate usage based on native frequency

| Native Frequency | Predicted | 95% CI |
|---|---|---|
| 0 | 0.49 | [0.38, 0.62] |
| 10,000 | 0.49 | [0.38, 0.62] |
| 20,000 | 2.11 | [1.28, 3.47] |
| 30,000 | 4.37 | [2.10, 9.07] |
| 40,000 | 4.37 | [2.10, 9.07] |
| 50,000 | 18.78 | [5.54, 63.71] |
| 60,000 | 38.95 | [9.47, 169.58] |

modulated by other factors. This means the impact of native frequency is best understood when considered in combination with these additional variables.

### Interactions via post hoc EMMs

To investigate interaction effects among the five factors, which could not be directly modelled in the NB-GLMM due to convergence issues, we conducted post hoc analysis using estimated marginal means (EMMs) (Lenth, 2022). These EMMs were derived from the NB-GLMM model. While we modelled every possible interaction, due to space constraints, we only present two interaction plots (Figures 7 and 8). Both plots illustrate interactions between three factors to maintain visual clarity and interpretability.
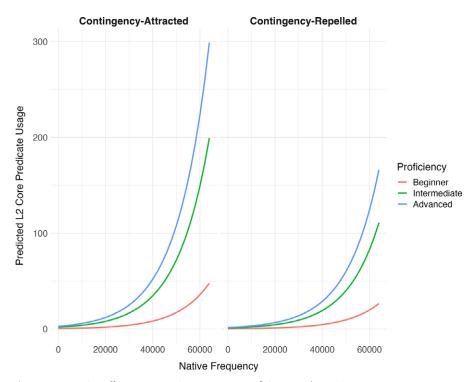


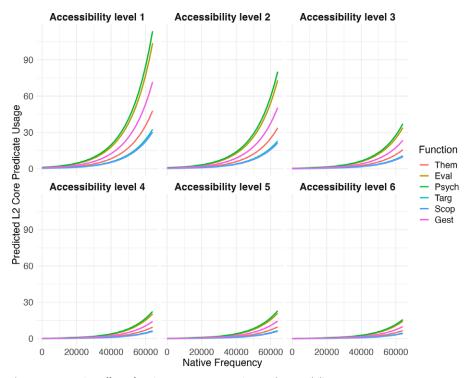**Figure 7.** Interaction effect among Native Frequency, Proficiency, and Contingency.

**Figure 8.** Interaction effect of Native Frequency, Function, and Accessibility.

Figure 7 demonstrates how Proficiency, Native Frequency, and Contingency collectively influence L2 core predicate choice. The figure is divided into two panels: the left for core predicates with high contingency with *duì*-constructions (Attracted), and the right for those with low contingency (Repelled). The horizontal axis shows predicate frequency in *duì*-constructions from the reference corpus (ranging from 0 to 60,000). Different proficiency levels are indicated by color-coded lines, with the vertical axis showing L2 learners' expected predicate usage frequency. The figure shows that L2 core predicate choice is shaped by all three factors. In both the Attracted and Repelled panels, *duì*-construction usage increases with learner proficiency. In addition, verbs and adjectives used more frequently by native speakers in the predicate slot are more likely to be used by L2 learners, with this trend particularly notable for verbs and adjectives attracted to *duì*-constructions.

Figure 8 illustrates another three-way interaction—Native Frequency, Function, and Accessibility—through six sub-plots. Each subplot corresponds to one of the six Accessibility levels, which shows how core predicates of different frequencies and functions are used by learners at that Accessibility level. The horizontal axis represents reference predicate frequency (0–60,000), the vertical axis shows predicted L2 usage frequency (0–100), and the color-coded lines indicate different functions.

Through looking at each panel of Figure 8, it becomes clear that it is not only the function that a verb or an adjective predicate is used to express within *duì*-constructions and how often it appears in the native language that matters for learners' core predicate choices, but also how accessible the verb or the adjective is. For instance, at Accessibility level 1, learners more frequently use verbs or adjectives associated with Psychological-

Reaction, Gesture-Attitude, and Evaluative-Perspective functions compared to verbs that are linked with Thematic-Relation, Target-Action, and Transformative-Scope functions. Moreover, at Accessibility level 1, the probability of a verb or an adjective being used by learners increases with its prevalence in native speech, and this tendency holds irrespective of its function. This pattern persists across Accessibility levels 1 to 3, suggesting the more readily accessible a verb or an adjective is, the more frequently learners use it, conditional on Native Frequency and Function. However, as predicates become less accessible (levels 4 to 6), their usage by learners exhibits a marked decline, highlighting how Accessibility, alongside Native Frequency and Function, collectively influence statistical learning.

## Discussion and conclusion

This study's first research question relates to whether L2 learners display evidence of statistical learning in their usage of non-adjacent *duì*-constructions, that is, does learners' usage mirror the distributional patterns in the target language? To address this question, we examined two factors—input frequency and contingency. We found that learners' usage of core predicates in *duì*-constructions generally aligns with the frequency and contingency patterns seen in the reference corpus, suggesting statistical learning of these complex linguistic structures. In this respect, our study accords with findings in broader linguistic research and the study of cognitive processes in language acquisition (Cadierno & Eskildsen, 2015; Gass & Mackey, 2002; Wulff, 2020), reinforcing the idea that L2 learners are sensitive to distributional patterns in their input.

However, the main value of our results lies in extending knowledge of statistical learning beyond adjacent constructions to nonadjacent constructions. By showing that learners' usage aligns with input regularities concerning these challenging constructions, our study casts some doubt on the speculation that L2 learners must engage in shallow processing due to attentional biases, risk of cognitive overload, and working memory limitations (Clahsen & Felser, 2006). We therefore add a new dimension to earlier research which suggests that learners are mainly sensitive to word-level elements and multiword sequences with no open internal slots in their input (Ellis et al., 2008; Felser et al., 2003; Sorace, 2006). As a result, our study contributes evidence that statistical learning operates on a wider scale than previously expected by showing sensitivity to non-adjacent dependencies. Since statistical learning is a crucial component of usage-based theory, our findings strengthen its empirical foundations.

Our study's investigation of its second research question underscores the significance of analyzing a broader range of factors than those typically considered in statistical learning research. In particular, we show that it is not just the distribution of constructions in learners' input that matters for their output, but also the accessibility and timing of initial exposure to specific linguistic items, the learner's proficiency level, and the cumulative duration of exposure to the target language. The relevance of examining these variables alongside the distributional factors of input frequency and contingency is underlined by the interaction effects shown in Figures 7 and 8. For instance, Figure 7 shows that intermediate and advanced learners are significantly more sensitive to contingency between core predicates and *duì*-constructions than beginners, suggesting that higher proficiency level learners are better able to extract patterns from their input, likely due to their prolonged L2 exposure.

Moreover, as Figure 8 shows, learners across all proficiency levels are more sensitive to collocating predicates in *duì*-constructions that are less complex and encountered

earlier in the learning journey. This sensitivity is especially pronounced when these early-encountered verbs and adjectives appear frequently in *duì*-constructions, compared to more complex verbs and adjectives typically encountered at later learning stages. In other words, the repeated encounter of earlier-learned and less complex lexis makes them more readily retrievable in memory and deployable in *duì*-constructions than later-acquired, more advanced lexis. Our findings thus conform with previous studies, which have shown that earlier learned items are often processed more quickly than those learned later because they are more strongly represented in a learner's mental lexicon than other alternative lexical representations (Sebastián-Gallés et al., 2005). This shows that the concept of Accessibility that we have introduced is a highly useful metric to understand L2 statistical learning, alongside input frequency and contingency. The explanatory power of this concept is likely due to the fact that Accessibility levels are partly guided by an L2 syllabus crafted to align with learners' developmental stages, prior knowledge, and learning objectives (Hulstijn et al., 2015).

Another area where we can see that distributional effects are modulated by L2 exposure is the formation of prototypical functions for *duì*-constructions. The formation of these prototypes is also closely tied to when initial input occurs, which means an earlier encountered function is more likely to form a prototype in a learner's repertoire. Specifically, we have found that all learners demonstrated a pronounced preference for *duì*-constructions that express Psychological-Reaction, Gesture-Attitude, and Evaluative-Perspective functions compared to Target-Action and Transformative-Scope functions. This inclination toward emotionally charged semantic functions, which resemble learner-generated prototypes, aligns with research on L2 epistemic stance. Such studies have suggested that L2 learners, especially those at lower proficiency levels, often engage more deeply in contexts requiring reactions, attitudes, and expressions of personal judgment or emotional states (Bucholtz & Hall, 2005; Gablasova et al., 2017; Lim et al., 2024). Such preferences for the Psychological-Reaction, Gesture-Attitude, and Evaluative-Perspective *duì*-construction functions are also likely influenced by the early introduction of these three functions within teaching materials (Guo, 2023). This early exposure could lead to a pre-emption effect (Trahey & White, 1993), wherein initial associations with the preposition *duì* are reinforced. Given these results, the regularities learners extract from their input are clearly influenced not only by raw frequency counts and association strengths but also by when and how a certain linguistic feature is encountered, with earlier encountered meanings or functions of a polysemous structure often becoming entrenched as prototypes. Our findings therefore suggest that investigations of statistical learning should broaden their traditional strong focus on frequency, contingency, and contextual dispersion to additionally incorporate other factors which shape exposure to input.

The importance of these exposure-related factors is particularly evident when examining apparent anomalies to the overall trend of statistical learning we have observed. In particular, while we found the most frequently used core predicates in *duì*-constructions are similar in both the reference and learner corpora, there were some exceptions. Take, for example, the verb *jìnxíng* ("to conduct"). This is one of the top three most frequently employed core predicates in *duì*-constructions in the reference corpus (55.71 pmw). Given the evidence of statistical learning, the expectation is that the verb *jìnxíng* would also be used frequently by learners. However, this is not the case. In fact, it was not used at all by the beginner or intermediate learners and only twice (5.28 pmw) by advanced learners. In seeking to account for this seeming disruption of the frequency effect, the study has found that the factors Accessibility and Function have strong explanatory power. Specifically, *jìnxíng* is a hard-to-learn verb

(Accessibility level 5) that learners would typically encounter later during their learning journeys. This verb is also normally used to express the Transformative-Scope function, which has the lowest IRR in the NB-GLMM model. In other words, from a learner's perspective, Transformative-Scope is the least prototypical or favored *duì*-construction function. Consequently, although learners are likely regularly exposed to *jìnxíng* via natural language input, the frequency effect and its contingency with *duì*-constructions are moderated.

Our findings also have important methodological implications. The study reinforces the importance of selecting statistical methods that can handle the typical characteristics of language data, particularly its Zipfian distribution and overdispersion. These are key features of both L1 and L2 language production (see Figure 2 and Section 4.1). By employing an NB-GLMM, which is well-suited for analyzing count data, we were able to capture both the systematic patterns and individual variation in core predicate usage, providing insights into how learners acquire and use verbal and adjectival predicates in *duì*-constructions.

Finally, several limitations to our study must be acknowledged. The first is that while we have broadened the range of factors examined in relation to L2 statistical learning relative to many existing studies by including three input exposure-related variables, our analysis still only captures a small proportion of the multitude of variables that, in reality, can affect SLA and statistical learning. The second issue concerns using a reference corpus as an input proxy. As described earlier, while an adequately sized and representative corpus can provide a valuable approximation of learners' natural language input, it cannot fully capture the diversity of authentic language exposure that learners encounter. Therefore, research which follows this approach must acknowledge the risk of glossing over the nuanced ways in which actual input shapes language acquisition. It is also vital to note that learners are exposed to input in a broad range of settings and while reference corpus data offer insights into general language patterns, the input that learners receive in the classroom and from teachers is also vital. Therefore, future studies should investigate how such input correlates to L2 usage patterns regarding nonadjacent constructions to further unpack the intricacies of statistical learning in SLA. Finally, while the learner corpus used in this study provides valuable learner-generated data, it suffers certain limitations. Specifically, although substantial in size, it only represents a controlled sample of language production from exams and tutorials, which may not reflect the full range of communicative experiences of L2 learners. Moreover, the learner corpus data may exhibit patterns aligned with academic and instructional settings rather than everyday language use, potentially limiting our findings' generalizability.

To conclude, this study reveals that statistical learning in L2 acquisition operates through a complex ecology where distributional patterns interact with various learner-specific and input exposure-shaping variables. By demonstrating that L2 learners not only discern nonadjacent collocational patterns that mirror native speaker usage but also show distinctive usage characteristics influenced by initial exposure timing, linguistic accessibility, and functional prototypicality, our findings support a conceptualization of statistical learning as a multidimensional process shaped by input sequence and learner experience—not just input frequency. These insights suggest promising directions for both theoretical models, particularly within the usage-based approach to language learning, and practical pedagogical approaches, where carefully sequenced input might leverage these interactional effects to enhance learning outcomes. Future research might explore how deliberate manipulation of input accessibility and timing could optimize

statistical learning across diverse linguistic constructions, potentially bridging the gap between classroom instruction and natural language acquisition.

**Author notes/acknowledgements.** Jiaqi Feng Guo conducted this research while affiliated with the Faculty of Education, University of Cambridge, UK. She is now affiliated with the Department of Chinese, School of Languages and Translation, University of Turku, Finland.

**CRediT author statement**
- Jiaqi Feng Guo: conceptualization; methodology; investigation; writing—original draft preparation; writing—review and editing.
- Pascual Perez-Paredes: conceptualization; formal analysis; writing—review and editing.

**Accessible summary.** All data files, R code for statistical analysis, and Python scripts used in this study have been deposited in the IRIS repository (https://www.iris-database.org) and are freely available for access, replication, and further research purposes.

# References

Alzahrani, A. (2021). The effects of two association measures on L2 collocation processing. *International Journal of English Linguistics*, *11*(5), 28. https://doi.org/10.5539/ijel.v11n5p28

Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, *59*(4), 390–412. https://doi.org/10.1016/j.jml.2007.12.005

Bley-Vroman, R. (2002). Frequency in production, comprehension, and acquisition. *Studies in Second Language Acquisition*, *24*(2), 209–213.

Boone, G., Wilde, V. D., & Eyckmans, J. (2022). A longitudinal study into learners' productive collocation knowledge in L2 German and factors affecting the learning. *Studies in Second Language Acquisition*. https://doi.org/10.1017/s0272263122000377

Brooks, P. J., Kwoka, N., & Kempe, V. (2017). Distributional effects and individual differences in L2 morphology learning: Determinants of L2 morphology learning. *Language Learning*, *67*(1), 171–207. https://doi.org/10.1111/lang.12204

Bucholtz, M., & Hall, K. (2005). Identity and interaction: A sociocultural linguistic approach. *Discourse Studies*, *7*(4–5), 585–614.

Bybee, J. L. (2008). Usage-based grammar and second language acquisition. In P. Robinson & N. C. Ellis (Eds.), *Handbook of Cognitive Linguistics and Second Language Acquisition*. Routledge.

Bybee, J. L. (2010). *Language, usage and cognition*. Cambridge University Press.

Cadierno, T., & Eskildsen, S. W. (2015). *Usage-based perspectives on second language learning*. De Gruyter.

Candarli, D. (2021). A longitudinal study of multi-word constructions in L2 academic writing: The effects of frequency and dispersion. *Reading and Writing*, *34*(5), 1191–1223. https://doi.org/10.1007/s11145-020-10108-3

Che, W., Feng, Y., Qin, L., & Liu, T. (2021). N-LTP: A Open-source Neural Chinese Language Technology Platform with Pretrained Models. *arXiv:2009.11616 [Cs]*. http://arxiv.org/abs/2009.11616

Chen, H., & Xu, H. (2019). Quantitative linguistics approach to interlanguage development: A study based on the Guangwai-Lancaster Chinese Learner Corpus. *Lingua*, *230*, 102736. https://doi.org/10.1016/j.lingua.2019.102736

Clahsen, H., & Felser, C. (2006). Grammatical processing in language learners. *Applied Psycholinguistics*, *27*(1), 3–42. https://doi.org/10.1017/S0142716406060024

Croft, W., & Cruse, D. A. (2004). *Cognitive linguistics*. Cambridge University Press.

Crossley, S. A., Skalicky, S., Kyle, K., & Monteiro, K. (2019). Absolute frequency effects in second language lexical acquisition. *Studies in Second Language Acquisition*, *41*(4), 721–744.

DeKeyser, R. (2007). *Practice in a second language: Perspectives from applied linguistics and cognitive psychology* (1st ed.). Cambridge University Press. https://doi.org/10.1017/CBO9780511667275

Deshors, S. C., & Gries, S. T. (2022). Using corpora in research on second language psycholinguistics. In *The Routledge handbook of second language acquisition and psycholinguistics*. Routledge.

Eckerth, J., & Tavakoli, P. (2012). The effects of word exposure frequency and elaboration of word processing on incidental L2 vocabulary acquisition through reading. *Language Teaching Research*, *16*, 227–252. https://doi.org/10.1177/1362168811431377

Edmonds, A., & Gudmestad, A. (2023). Phraseological use and development during a stay abroad: Exploring sensitivity to frequency and cue contingency. *Language Learning*, *73*(2), 475–507. https://doi.org/10.1111/lang.12547

Elder, C., McNamara, T., Kim, H., Pill, J., & Sato, T. (2017). Interrogating the construct of communicative competence in language assessment contexts: What the non-language specialist can tell us. *Language & Communication*, *57*, 14–21. https://doi.org/10.1016/J.LANGCOM.2016.12.005

Ellis, N. C., & Ferreira–Junior, F. (2009). Construction learning as a function of frequency, frequency distribution, and function. *The Modern Language Journal*, *93*(3), 370–385. https://doi.org/10.1111/j.1540-4781.2009.00896.x

Ellis, N. C., Römer, U., O'Donnell, M. B., & Schleppegrell, M. J. (2016). *Usage-based approaches to language acquisition and processing: Cognitive and corpus investigations of construction grammar*. John Wiley & Sons.

Ellis, N. C., Simpson-Vlach, R., & Maynard, C. (2008). Formulaic language in native and second language speakers: Psycholinguistics, corpus linguistics, and TESOL. *TESOL Quarterly*, *42*(3), 375–396. https://doi.org/10.1002/j.1545-7249.2008.tb00137.x

Erickson, L. C., & Thiessen, E. D. (2015). Statistical learning of language: Theory, validity, and predictions of a statistical learning account of language acquisition. *Developmental Review, 37*, 66–108. https://doi.org/10.1016/j.dr.2015.05.002

Eskildsen, S. W. (2022). Usage-based SLA: From corpora to social interaction. In K. L. Geeslin (Ed.), *The Routledge handbook of second language acquisition and sociolinguistics*. Routledge.

Felser, C., Roberts, L., Marinis, T., & Gross, R. (2003). The processing of ambiguous sentences by first and second language learners of English. *Applied Psycholinguistics*, *24*, 453–489. https://doi.org/10.1017/S0142716403000237

Gablasova, D., Brezina, V., Mcenery, T., & Boyd, E. (2017). Epistemic stance in spoken L2 English: The effect of task and speaker style. *Applied Linguistics*, *38*(5), 613–637. https://doi.org/10.1093/applin/amv055

Gass, S. M., Behney, J., & Plonsky, L. (2020). *Second language acquisition: An introductory course*. Routledge.

Gass, S. M., & Mackey, A. (2002). Frequency effects and second language acquisition: A complex picture? *Studies in Second Language Acquisition*, *24*(2), 249–260.

Gelman, A., Hill, J., & Yajima, M. (2012). Why we (usually) don't have to worry about multiple comparisons. *Journal of Research on Educational Effectiveness*, *5*(2), 189–211. https://doi.org/10.1080/19345747.2011.618213

Goldberg, A. E. (1995). *Constructions: A construction grammar approach to argument structure*. University of Chicago Press.

Goldberg, A. E. (2015). Compositionality. In *The Routledge handbook of semantics*. Routledge.

Granena, G., & Long, M. H. (2013). Age of onset, length of residence, language aptitude, and ultimate L2 attainment in three linguistic domains. *Second Language Research*, *29*(3), 311–343. https://doi.org/10.1177/0267658312461497

Gries, S. T. (2010a). Useful statistics for corpus linguistics. *A Mosaic of Corpus Linguistics: Selected Approaches*, *66*, 269–291.

Gries, S. T., & Ellis, N. C. (2015). Statistical measures for usage-based linguistics. *Language Learning*, *65*(S1), 228–255. https://doi.org/10.1111/lang.12119

Gries, S. Th. (2010b). Dispersions and adjusted frequencies in corpora: Further explorations. In S. Th. Gries, S. Wulff, & M. Davies (Eds.), *Corpus-linguistic Applications*. Brill | Rodopi. https://doi.org/10.1163/9789042028012_014

Gries, S. Th. (2022a). On, or against?, (just) frequency. In H. C. Boas (Ed.), *Directions for pedagogical construction grammar* (pp. 47–72). De Gruyter. https://doi.org/10.1515/9783110746723-002

Gries, S. Th. (2022b). Toward more careful corpus statistics: Uncertainty estimates for frequencies, dispersions, association measures, and more. *Research Methods in Applied Linguistics*, *1*(1), 100002. https://doi.org/10.1016/j.rmal.2021.100002

Guo, J. (2023). *Usage-based Analysis of Schematic and Multifunctional Verb Argument Constructions in Three Contrasting Corpora: A Native Chinese Corpus, a CFL Textbook Corpus and a Chinese Learner Corpus.* Apollo - University of Cambridge Repository. https://doi.org/10.17863/CAM.104545

Hashimoto, B. J., & Egbert, J. (2019). More than frequency? Exploring predictors of word difficulty for second language learners. *Language Learning*, *69*(4), 839–872. https://doi.org/10.1111/lang.12353

Hilpert, M. (2014). Collostructional analysis. *Corpus Methods for Semantics: Quantitative Studies in Polysemy and Synonymy*, *43*, 391.

Hoey, M. (2012). *Lexical priming: A new theory of words and language.* Routledge.

Hopp, H. (2023). Sentence processing in a second language: Linguistic approaches. *The Routledge Handbook of Second Language Acquisition and Psycholinguistics*, 216–228.

Huang, S.-Z. (2006). Property theory, adjectives, and modification in Chinese. *Journal of East Asian Linguistics*, *15*(4), 343–369. https://doi.org/10.1007/s10831-006-9002-0

Hulstijn, J. H., Ellis, R., & Eskildsen, S. (2015). Orders and sequences in the acquisition of L2 morphosyntax, 40 years on: An introduction to the special issue. *Language Learning*, *65*. https://doi.org/10.1111/lang.12097

Ibbotson, P. (2013). The scope of usage-based theory. *Frontiers in Psychology*, *4*. https://doi.org/10.3389/fpsyg.2013.00255

Isbilen, E. S., Frost, R. L. A., Monaghan, P., & Christiansen, M. H. (2022). Statistically based chunking of nonadjacent dependencies. *Journal of Experimental Psychology: General*, *151*(11), 2623–2640. https://doi.org/10.1037/xge0001207

Johnson, R. A. (2009). *Statistics: Principles and methods* (6th edition). Wiley.

Kartal, G., & Sarigul, E. (2017). Frequency effects in second language acquisition: an annotated survey. *Journal of Education and Training Studies*, *5*(6), 1. https://doi.org/10.11114/jets.v5i6.2327

Kay, P., & Fillmore, C. J. (1999). Grammatical constructions and linguistic generalizations: The what's X doing Y? Construction. *Language*, *75*(1), 1–33. JSTOR. https://doi.org/10.2307/417472

Kyle, K. (2021). Natural language processing for learner corpus research. *International Journal of Learner Corpus Research*, *7*(1), 1–16. https://doi.org/10.1075/ijlcr.00019.int

Kyle, K., Crossley, S., & Verspoor, M. (2021). Measuring longitudinal writing development using indices of syntactic complexity and sophistication. *Studies in Second Language Acquisition*, *43*(4), 781–812.

Langacker, R. W. (1987). *Foundations of cognitive grammar: Theoretical prerequisites.* Stanford University Press.

Lenth, R. (2022). *emmeans: Estimated Marginal Means, aka Least-Squares Means_.* (Version R package version 1.8.0) [Computer software]. https://CRAN.R-project.org/package=emmeans

Levin, B., & Hovav, M. R. (2005). *Argument realization* (Vol. 10). Citeseer. https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=e960372b5813642963b1c66d527b1a05626bb7ed

Lim, J. D. O., Mark, G., Pérez-Paredes, P., & O'Keeffe, A. (2024). Exploring part of speech (pos) tag sequences in a large-scale learner corpus of L2 English: A developmental perspective. *Corpora*, *19*(1), 31–59. https://doi.org/10.3366/cor.2024.0297

Marinis, T., Roberts, L., Felser, C., & Clahsen, H. (2005). Gaps in second language sentence processing. *Studies in Second Language Acquisition*, *27*(1), 53–78.

Martinez, R., & Murphy, V. A. (2011). Effect of frequency and idiomaticity on second language reading comprehension. *TESOL Quarterly*, *45*(2), 267–290. https://doi.org/10.5054/tq.2011.247708

McElreath, R. (2020). *Statistical rethinking: A Bayesian course with examples in R and STAN* (2nd ed.). CRC.

McEnery, T., & Brezina, V. (2022). *Fundamental principles of corpus linguistics.* Cambridge University Press.

McManus, K. (2024). Introducing usage in SLA. In K. McManus (Ed.), *Usage in second language acquisition*. Routledge.

Meurers, D. (2015). Learner corpora and natural language processing. In S. Granger, G. Gilquin, & F. Meunier (Eds.), *The Cambridge handbook of learner corpus research* (1st ed., pp. 537–566). Cambridge University Press. https://doi.org/10.1017/CBO9781139649414.024

Monteiro, K. R., Crossley, S. A., & Kyle, K. (2020). In search of new benchmarks: Using L2 lexical frequency and contextual diversity indices to assess second language writing. *Applied Linguistics*, *41*(2), 280–300.

Murakami, A., & Ellis, N. C. (2022). Effects of availability, contingency, and formulaicity on the accuracy of English grammatical morphemes in second language writing. *Language Learning*. https://doi.org/10.1111/lang.12500

Nation, I. S. P. (2013). *Learning vocabulary in another language* (2nd ed.). Cambridge University Press. https://doi.org/10.1017/CBO9781139858656

Paquot, M., & Gries, S. T. (2020). *A practical handbook of corpus linguistics.* Springer Nature.

Pellicer-Sánchez, A. (2016). Incidental L2 vocabulary acquisition from and while reading: An eye-tracking study. *Studies in Second Language Acquisition*, *38*(1), 97–130. https://doi.org/10.1017/S0272263115000224

Pérez-Paredes, P. (2020). *Corpus linguistics for education: A guide for research*. Routledge.

Peters, E. (2018). The effect of out-of-class exposure to English language media on learners' vocabulary knowledge. *ITL - International Journal of Applied Linguistics*, *169*(1), 142–168. https://doi.org/10.1075/itl.00010.pet

Rebuschat, P., & Williams, J. N. (2012). *Statistical learning and language acquisition*. Walter de Gruyter.

Römer, U. (2019). A corpus perspective on the development of verb constructions in second language learners. *International Journal of Corpus Linguistics*, *24*(3), 268–290. https://doi.org/10.1075/ijcl.00013.roe

Römer, U., & Garner, J. (2019). The development of verb constructions in spoken learner English: Tracing effects of usage and proficiency. *International Journal of Learner Corpus Research*, *5*(2), 207–230. https://doi.org/10.1075/ijlcr.17015.rom

Sebastián-Gallés, N., Echeverría, S., & Bosch, L. (2005). The influence of initial exposure on lexical representation: Comparing early and simultaneous bilinguals. *Journal of Memory and Language*, *52*(2), 240–255.

Sinclair, J. (1991). *Corpus, concordance, collocation*. Oxford University Press.

Siyanova, A., & Schmitt, N. (2008). L2 learner production and processing of collocation: A multi-study perspective. *The Canadian Modern Language Review*, *64*(3), 429–458. https://doi.org/10.3138/cmlr.64.3.429

Sorace, A. (2006). Possible manifestations of shallow processing in advanced second language speakers. *Applied Psycholinguistics*, *27*, 88–91. https://doi.org/10.1017/S0142716406060164

Stefanowitsch, A., & Gries, S. Th. (2003). Collostructions: Investigating the interaction of words and constructions. *International Journal of Corpus Linguistics*, *8*(2), 209–243. https://doi.org/10.1075/ijcl.8.2.03ste

Stutterheim, C. V., Lambert, M., & Gerwien, J. (2021). Limitations on the role of frequency in L2 acquisition. *Language and Cognition*, *13*(2), 291–321. https://doi.org/10.1017/langcog.2021.5

Su, X. (2013). 现代汉语语义分类词典 *[A Thesaurus of Modern Chinese]*. Commercial Press.

Thompson, S. A., & Tao, H. (2010). Conversation, grammar, and fixedness: Adjectives in Mandarin revisited. *Chinese Language and Discourse. An International and Interdisciplinary Journal*, *1*(1), 3–30. https://doi.org/10.1075/cld.1.1.01tho

Tono, Y. (2004). Multiple comparisons of IL, L1 and TL corpora: The case of L2 acquisition of verb subcategorization patterns by Japanese learners of English. *Corpora and Language Learners*, *45*.

Trahey, M., & White, L. (1993). Positive evidence and preemption in the second language classroom. *Studies in Second Language Acquisition*, *15*(2), 181–204.

Uchihara, T., Webb, S., & Yanagisawa, A. (2019). The effects of repetition on incidental vocabulary learning: A meta-analysis of correlational studies. *Language Learning*. https://doi.org/10.1111/LANG.12343

Winke, P., Gass, S., & Sydorenko, T. (2013). Factors influencing the use of captions by foreign language learners: An eye-tracking study. *The Modern Language Journal*, *97*(1), 254–275. https://doi.org/10.1111/j.1540-4781.2013.01432.x

Winter, B., & Bürkner, P.-C. (2021). Poisson regression for linguists: A tutorial introduction to modelling count data with brms. *Language and Linguistics Compass*, *15*(11), e12439.

Winter, B., & Grice, M. (2021). Independence and generalizability in linguistics. *Linguistics*, *59*(5), 1251–1277. https://doi.org/10.1515/ling-2019-0049

Wolter, B., & Gyllstad, H. (2013). Frequency of input and L2 collocational processing: A comparison of congruent and incongruent collocations. *Studies in Second Language Acquisition*, *35*(3), 451–482.

Wulff, S. (2020). Usage-based approaches. In N. Tracy-Ventura & M. Paquot (Eds.), *The Routledge handbook of second language acquisition and corpora*. Routledge.

Xun, E., Rao, G., Xiao, X., & Zang, J. (2016). "大数据背景下 BCC 语料库存的研制" [The Construction of the BCC Corpus in the Age of Big Data]. 语料库语言学 *[Corpus Linguistics]*, *3*(1), 93–109.

Yan, J., & Liu, H. (2022). Semantic roles or syntactic functions: The effects of annotation scheme on the results of dependency measures. *Studia Linguistica*, *76*(2), 406–428. https://doi.org/10.1111/stul.12177

---