

1 Stefanie Wulff, Stefan Th. Gries and Nicholas Lester
2 **Optional *that* in complementation by**
3 **German and Spanish learners**
4

5
6 **1 Introduction**
7

8
9 This study examines the factors that govern the variable presence of the comple-
10 mentizer *that* in English object-, subject-, and adjectival complement construc-
11 tions as in (1) to (3):¹

- 12
13 (1) a. I thought that Nick likes candy.
14 b. I thought \emptyset Nick likes candy.
15
16 (2) a. The problem is that Nick doesn't like candy.
17 b. The problem is \emptyset Nick doesn't like candy.
18
19 (3) a. I'm glad that Stefan likes candy.
20 b. I'm glad \emptyset Stefan likes candy.
21
22

23 The conditions under which native speakers (NS) decide to realize or drop the
24 complementizer have been intensively studied (e.g., Jaeger 2010; Tagliamonte
25 and Smith 2005; Thompson and Mulac 1991; Torres Cacoullos and Walker
26 2009), while few studies have investigated this phenomenon in non-native
27 speakers (NNS) (e.g., Durham 2011; Wulff, Lester, and Martinez-Garcia 2014). In
28 the present study, we therefore address the following research questions:

- 29 – What factors govern *that*-variation in intermediate-level German and Spanish
30 L2 learners of English?
31 – How do these learners' preferences compare to those of native speakers?
32 More specifically, under what conditions, how much, and why do learners
33 deviate from native speaker behavior?
34

35
36 ¹ The complementizer is optional in other constructions as well, including appositions, relative
37 clauses of *it*-clefts, and with extraposed subjects; instances of these constructions, which are far
38 less frequent than the three constructions examined here, were not considered in this study.

39 **Stefanie Wulff**, University of Florida

40 **Stefan Th. Gries** and **Nicholas Lester**, University of California, Santa Barbara

<https://doi.org/10.1515/9783110572186-004>

1 The paper is structured as follows. Section 2 provides a compact overview of the
 2 factors suggested to impact *that*-variation; specifically, Section 2.1 discusses *that*-
 3 variation in L1 English whereas Section 2.2 briefly describes the equivalents of
 4 *that*-variation in L1 German and L1 Spanish, the native language backgrounds
 5 of the L2-learners investigated here. Section 3 gives a brief summary of previous
 6 studies on *that*-variation in learner populations. In Section 4, we describe our
 7 data sample in detail, explain how the data were annotated for the different
 8 variables included in the study, introduce the statistical method employed,
 9 MuPDAR, and explain how this method was applied to our data. Section 5
 10 summarizes the results, and Section 6 concludes by recapturing the main find-
 11 ings and their implications, in particular from the perspective of usage-based
 12 construction grammar.

13

14

15

2 Factors influencing *that*-variation

16

17

2.1 *That*-variation in native English

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

Over the last 25 years, *that*-variation has received a lot of attention. Space does not permit a detailed discussion of this body of research (see Wulff, Lester, and Martinez-Garcia 2014) so here we briefly summarize only those factors which have consistently emerged as relevant:

- *mode* (Biber 1999; Bryant 1962; Storms 1966): the complementizer is omitted more frequently in spoken than in written language; likewise, higher shares of zero-*that* are found in informal registers (both spoken and written).
- *structural complexity* (also referred to as *syntactic weight*; see Elsness 1984; Jaeger 2010; Kaltenböck 2006; Thompson and Mulac 1991; Torres Cacoullous and Walker 2009): syntactically light main and/or complement clause subjects as well as light complement clauses are correlated with zero-*that*, and these correlations are strongest with the structurally simple first person pronoun *I* in subject position in the matrix clause.
- *clause juncture* (Jaeger 2010; Kaltenböck 2006; Thompson and Mulac 1991; Torres Cacoullous and Walker 2009): chances of zero-*that* are highest when clause juncture is intact, i.e., when there is no intervening material anywhere. When material intervenes between the matrix clause subject and the verb, the matrix clause verb and the complementizer slot, or the complementizer slot and the ensuing complement clause, this raises the likelihood of *that* being realized. Some studies suggest that material preceding the matrix clause subject may also increase chances of *that* – while clause-initial material does not interrupt clause juncture, it adds to the overall complexity of the message.

- 1 – *properties of the matrix clause verb* (Dor 2005; Kaltenböck 2006; Rissanen
2 1991; Tagliamonte and Smith 2005): several studies point out that zero-*that*
3 is especially likely with (typically highly frequent) matrix clause verbs that
4 denote truth claim predicates (such as *think*, *know*, and *believe*). What is
5 more, Wulff, Lester, and Martinez-Garcia (2014) found that beyond their
6 absolute frequencies, some verbs are zero-favoring while others are *that*-
7 favoring, as can be expressed in the association strength between a given
8 verb and either construction, respectively.
- 9 – *surprisal* (Jaeger 2010; Levy 2008): matrix verb lemmas that are biased to
10 occur in the complement clause construction carry enough information about
11 the upcoming clause juncture to make the overt complementizer redundant.
12 This informational boost is quantified using an information-theoretic measure
13 known as surprisal, which Jaeger (2010) shows is positively correlated with
14 rates of *that*-mentioning.
- 15 – *individual variation*: just like in many other (psycho)linguistic phenomena,
16 there is individual variation among speakers.

17
18 In the next section, we provide a very brief overview of the equivalents of *that*-
19 variation in German and Spanish.

21 2.2 *That*-variation in native German and Spanish

22
23 Regarding complementizer optionality, German is slightly less permissive than
24 English: The complementizer *dass* is optional in subject and direct object com-
25 plements, but obligatory in adjectival complements. German also differs from
26 English in that the position of the verb in the complement clause is contingent
27 on whether the complementizer is realized or not: When the complementizer
28 is not realized, the verb follows the subject (which is the default word order
29 for main clauses in German); when the complementizer is realized, the verb
30 appears in clause-final position (which is the default word order for subordinate
31 clauses in German). Examples (4) to (6) provide German translations of (1) to (3)
32 respectively.

- 33
34 (4) a. *Ich dachte, dass Nick Suesses mag.*
35 I think.3SG.PST COMP Nick candy like.3SG.PRS
36 ‘I thought that Nick likes candy’
37
- 38 b. *Ich dachte, Ø Nick mag Suesses.*
39 I think.3SG.PST COMP Nick like.3SG.PRS candy
40 ‘I thought Nick likes candy’

- 1 (5) a. *Das Problem ist, dass Nick Suesses nicht mag.*
 2 the problem COP.3SG.PRS COMP Nick candy NEG like.3SG.PRS
 3 ‘The problem is that Nick doesn’t like candy’
 4
- 5 b. *Das Problem ist, Ø Nick mag*
 6 the problem COP.3SG.PRS COMP Nick like.3SG.PRS NEG
 7 *Suesses nicht.*
 8 candy NEG
 9 ‘The problem is Nick doesn’t like candy’
- 10 (6) a. *Ich bin froh, dass Stefan Suesses mag.*
 11 I COP.1SG.PRS glad COMP Stefan candy like.3SG.PRS
 12 ‘I’m glad that Stefan likes candy’
 13
- 14 b. **Ich bin froh, Ø Stefan mag Suesses.*
 15 I COP.1SG.PRS glad COMP Stefan like.3SG.PRS candy
 16 ‘I’m glad Stefan likes candy’
 17

18 Spanish, in turn, is even less permissive than German: the complementizer *que*
 19 is always obligatory. (7) to (9) are translations of (1) to (3), respectively.
 20

- 21 (7) a. *Pensé que a Nick le gustaban los dulces.*
 22 think.1SG.PST COMP to Nick CL.DAT. like.3.PL.IMP the candies
 23 ‘I thought that Nick likes candy’
 24
- 25 b. **Pensé Ø a Nick le gustaban los dulces.*
 26 think.1SG.PST COMP to Nick CL.DAT. like.3.PL.IMP the candies
 27 ‘I thought Nick likes candy’
- 28 (8) a. *El problema es que a Nick no le*
 29 the problem COP.3SG.PRS COMP to Nick NEG CL.DAT.
 30 *gustan los dulces.*
 31 like.3.PL.IMP the candies
 32 ‘The problem is that Nick doesn’t like candy’
 33
- 34 b. **El problema es Ø a Nick no le*
 35 the problem COP.3SG.PRS COMP to Nick NEG CL.DAT.
 36 *gustan los dulces.*
 37 like.3.PL.IMP the candies
 38 ‘The problem is Nick doesn’t like candy’
 39
 40

- 1 (9) a. *Me alegra que a Stefan le*
 2 CL.DAT. makes-happy.3SG.PRS COMP to Stefan CL.DAT.
 3 *gustan los dulces.*
 4 like.3.PL.PRS the candies
 5 ‘I’m glad that Stefan likes candy’
 6
- 7 b. **Me alegra Ø a Stefan le*
 8 CL.DAT. makes-happy.3SG.PRS COMP to Stefan CL.DAT.
 9 *gustan los dulces.*
 10 like.3.PL.PRS the candies
 11 ‘I’m glad Stefan likes candy’
 12

13 Given these contrasts between English, German, and Spanish, we can assume
 14 that native-like use of *that*-variation should be overall easier to attain for German
 15 learners of English than Spanish learners, who should be most reluctant to omit
 16 the complementizer. Previous research in fact supports this hypothesis (Wulff
 17 2016; Wulff, Lester and Martinez-Garcia 2014).
 18
 19

20 3 *That*-variation in L2 production

21
 22
 23 In contrast to the wealth of studies on native speakers, there are few studies to
 24 date that examine *that*-variation in L2 learners. One example is Durham (2011)
 25 on native speakers’ and French, German, and Italian ESL learners’ use of *that*-
 26 variation in emails. Durham reports that shares of zero-*that* hover around 35%
 27 overall; French and Italian learners are more likely to produce the comple-
 28 mentizer than the German learners and native English speakers. Furthermore,
 29 Durham confirms that, as in native speakers, combinations of the first person
 30 pronoun *I* as the matrix clause subject and verbs like *think* and *hope* trigger the
 31 highest shares of zero-*that*. The German and Italian learners display sensitivity
 32 also to clause juncture constraints while the French learners do not.

33 Wulff, Lester and Martinez-Garcia (2014) examine what comprises the written
 34 part of the data sample of the present study (i.e., native English speakers, German
 35 L2 learners, and Spanish L2 learners). They include all of the factors listed in
 36 Section 4.2.1 (except for mode, surprisal, and individual variation) in a multi-
 37 factorial regression analysis. Their findings suggest intermediate-advanced level
 38 German and Spanish learners are quite attuned to native-like choices: they
 39 appear to be sensitive to the same factors as native speakers, and the directions
 40 of the effects for these factors are identical. That said, compared to the native

1 speakers, both learner groups display a lower rate of zero-*that*. They also appear
 2 to be more impacted by processing-related factors such as structural complexity
 3 and clause juncture as opposed to lexical-semantic properties such as the choice
 4 of matrix clause verb.

5 Wulff (2016) expands Wulff, Lester, and Martinez-Garcia's (2014) study by
 6 adding spoken data to the sample. Her results are mainly in accord with the
 7 previous studies, and confirm that, like native speakers, second language
 8 learners (at least at an intermediate level of proficiency) are aware of the mode-
 9 dependent nature of *that*-variation.

10 In the present study, we are improving on Wulff's analysis in several impor-
 11 tant ways. First, the current analysis includes surprisal as a predictor. Second,
 12 the statistical analysis presented here is much more sophisticated than the
 13 binary logistic regression Wulff (2016) presents: firstly, we are using a two-step
 14 regression procedure that has been developed specifically for the analysis of
 15 differences between native and non-native language; secondly, the regressions
 16 we are using involve mixed-effects/multi-level models. This choice of model
 17 allows us to take complex hierarchical structures in the data into consideration,
 18 such as speaker- and verb-specific effects. We outline the specifics of this approach
 19 in Section 4.3.

22 4 Methods

24 4.1 Data

26 The data for this study were retrieved from different corpora. The NS data were
 27 obtained from the British component of the *International Corpus of English* (ICE-
 28 GB), a balanced, parsed, 1-million words corpus of British English, which com-
 29 prises 60% written and 40% spoken data. Using the ICE-CUP software packet
 30 that accompanies the corpus, all instances of the three complement construc-
 31 tions that are contained in the corpus were retrieved.

32 The written NNS data were obtained from the German and the Spanish sub-
 33 corpora of the second version of the *International Corpus of Learner English*
 34 (G-ICLE and SP-ICLE; see Granger et al. 2009). ICLE comprises 3.7 million words
 35 of EFL writing from learners from 16 different L1 backgrounds. The spoken
 36 learner data came from the German and Spanish sub-corpora of the LINDSEI
 37 corpus (see Gilquin, De Cock, and Granger 2010). LINDSEI is a 1-million-word
 38 corpus of informal interviews with high intermediate-advanced proficiency EFL
 39 learners.

1 Unlike the ICE-GB, neither ICLE nor LINDSEI are syntactically parsed, so in
 2 order to retrieve hits from these corpora, the following procedure was adopted:
 3 A list of all verb lemmas attested in the ICE-GB across the three constructions
 4 was created and used to retrieve all sentences with these verb lemmas in
 5 G-ICLE, SP-ICLE, and LINDSEI. The resulting candidate list was then manually
 6 checked for true hits.

7 Table 1 provides a breakdown of the final data sample of 9,445 hits by
 8 L1 background, construction (ADJ vs. OBJ vs. SUB complementation), mode
 9 (spoken vs. written), and whether the complementizer was absent or present.
 10 Two things stand out immediately when we look at the learner populations:
 11 both German and Spanish learners use complementation constructions far less
 12 frequently in speaking than in writing (this is especially true for adjectival and
 13 subject complementation), which reverses the trend we observe in the native
 14 speaker data. Secondly, adjectival complementation is very infrequent in the
 15 Spanish learner data.

16
 17 **Table 1:** Data sample of the present study

L1	Construction	Mode	<i>that</i> = absent	<i>that</i> = present	Total
English	ADJ	spoken	107	57	164
		written	41	35	76
	OBJ	spoken	2,446	1,235	3,681
		written	528	651	1,179
	SUB	spoken	85	296	381
		written	7	146	153
	Total		3,214	2,420	5,634
German	ADJ	spoken	2	4	6
		written	17	84	101
	OBJ	spoken	643	155	798
		written	224	853	1,077
	SUB	spoken	12	21	33
		written	9	213	222
	Total		907	1330	2,237
Spanish	ADJ	spoken	0	2	2
		written	0	3	3
	OBJ	spoken	437	173	610
		written	176	682	858
	SUB	spoken	4	35	39
		written	8	54	62
	Total		625	949	1,574
Total			4,746	4,699	9,445

4.2 Variables and operationalizations

4.2.1 Frequently-used predictors

The 9,445 hits retrieved from the corpora were coded for the factors listed below. In order to understand how each factor was operationalized, let us consider the (fictional) example sentence in (10).

(10) Seriously, I really hope very much that he likes this chocolate.

- **L1 background:** the native language of the speaker: English vs. German vs. Spanish;
- **Mode:** the sub-corpus from which an example came: spoken vs. written;
- **Complementizer:** complementizer presence: absent vs. present;
- **ComplementType:** the type of complement sentence: adjectival vs. object vs. subject;
- **LengthCIM:**² the length of any clause-initial material (before the matrix-clause subject) in number of characters;
- **LengthMatrixSubj:** the length of the matrix clause subject;
- **LengthComplementSubj:** the length of the complement clause subject;
- **LengthComplement:** the length of the complement clause;
- **LengthCCRemainder:** the length of any post-verbal material in the complement clause;
- **LengthMCSubjMCVerb:** the amount of material between the matrix clause subject and the matrix clause verb;
- **LengthMCVerbCC:** the amount of material between the matrix clause verb and the complement clause;
- **DeltaPWC/DeltaPCW:** the association of each verb attested in the data sample to *that* or *zero-that* was calculated and vice versa. The specific association measure employed here is a *Delta-P* association measure (using Stefan Th. Gries' *R-script coll.analysis 3.2*; Gries 2007), which involves two different scores: a *Delta-P_{WC}* value (*WC* stands for 'word-to-construction') quantifies how predictive the verb is of the absence or presence of *that*, and a *Delta-P_{CW}* value (*CW* stands for 'construction-to-word') indicates how predictive the absence or presence of *that* is for the verb in question (see Ellis 2006; Gries 2013). *Delta-P* values range between -1 when the first element strongly repels the second, via 0 (when there is no association), to 1 (when the first element strongly attracts the second).

² All length-related predictors were measured as the number of characters. While counting the number of syllables or words might seem more intuitive, for all intents and purposes, length counts in characters, words, phonemes, or syllables are so highly correlated (and, thus, come with no conceptual/interpretive disadvantages) that we opted for the ease of operationalizing length with automatically-countable character lengths.

1 Consider Table 2 for the annotation of (10):

2
3 **Table 2:** The annotation of example (10)

4 Complementizer: present	ComplementType: object
5 LengthCIM: 9 (“Seriously”)	LengthMatrixSubj: 1 (“I”)
6 LengthComplementSubj: 2 (“he”)	LengthComplement: 20 (“he likes this chocolate”)
7 LengthCCRemainder: 13 (“this chocolate”)	
8 LengthMCSbjMCVerb: 6 (“really”)	LengthMCVerbCC: 8 (“very much”)
9 Delta- P_{CW} for <i>hope</i> : 0.1148	Delta- P_{WC} for <i>hope</i> : 0.167

11
12 As previously mentioned, we also included a predictor measuring the surprisal
13 of the material spanning the clause juncture (i.e., the surprisal of moving from
14 *much* to *Nick* in (10)). Given the relative scarcity of such applications in SLA
15 research, we provide a more thorough discussion of this variable in Section
16 4.2.2. Finally, we added annotation to take into consideration speaker-specific
17 and lexically-specific effects: each example was annotated for the corpus and
18 the file it came from as well as for the verb form and the verb lemma of the
19 main clause.

20 21 22 **4.2.2 The information-theoretic notion of surprisal**

23
24 *That*-variation has been shown to be affected by various probabilistic relation-
25 ships between words (and larger units), both within and across the matrix and
26 complement clauses. Jaeger (2010) showed that one particularly important relation-
27 ship holds between the matrix verb lemma (uninflected stem, e.g., EAT for
28 *eat, eats, eating, ...*) and the syntactic juncture between the matrix and comple-
29 ment clause. When the verb lemma was highly informative about the presence of
30 an upcoming clause juncture, rates of *that* decreased. To measure the expecta-
31 tion of the clause juncture that is projected from the matrix verb lemma (in other
32 words, the *redundancy* of the complementizer), Jaeger used an information-
33 theoretic measure known as *surprisal* or *self-information*. Surprisal measures
34 how uncertain one would be about observing some event – how ‘surprising’
35 that event would be – given a known probability distribution of related events.
36 It is calculated by taking the negative binary log of the probability p of a given
37 event x belonging to probability distribution P , as in (11).

$$38$$

$$39 (11) S(x: x \in P) = -\log_2 p(x)$$

$$40$$

1 Because he was interested in the surprisal of the juncture given the matrix verb
 2 lemma, Jaeger (2010) substituted the conditional probability p (juncture | matrix
 3 verb lemma) for the simple probability p . The generalized form of this substitution
 4 tion, which we shall henceforth refer to as *conditional surprisal* S_c , is

5

$$6 \quad (12) \quad S_c(y|x: y, x \in P) = -\log_2 p(y/x)$$

7

8 In the present study, we replace Jaeger's (2010) conditional surprisal value with
 9 the bi-directional collostructional association measure ΔP , and so measure
 10 directly the preferences of each matrix verb for the presence or absence of
 11 the complementizer (as opposed to the presence or absence of a complement
 12 clause). However, the notion of conditional surprisal can be applied at a finer
 13 resolution to explore local negotiations of informational load at the clause juncture.
 14 For instance, as Jaeger points out, the relative (un)expectedness of the first
 15 word following the clause juncture (i.e., the complement clause onset) may
 16 influence *that*-mentioning, such that more surprising onsets correlate with
 17 greater shares of *that*. Jaeger proposes that, ideally, the surprisal of the onset
 18 should be conditioned on the joint probability of the matrix verb occurring in a
 19 complement clause construction, that is, $S_c(\text{onset} | \text{verb, complement construction})$.
 20 However, this measure misses the fact that different verbs are differently
 21 associated with rates of the *that*-mentioning apart from their likelihood of
 22 occurring within the complement-clause construction (consider the logically
 23 possible case of a verb that only occurs in the complement-clause construction,
 24 but prefers *that*). Moreover, Jaeger's proposal overlooks the possible fluctuations
 25 in informational load that can be attributed directly to the words standing at
 26 either edge of the clause juncture (the left edge may contain a word other than
 27 the matrix verb). The relationships between these words may incrementally or
 28 instantaneously overturn (or reinforce) the expectations triggered by the matrix
 29 verb. Finally, by taking his measurements at the level of the matrix verb lemma,
 30 Jaeger increases the statistical reliability of his estimates, but glosses over the
 31 possibility that the different inflected forms of a verb will correlate with different
 32 patterns of use.

33 Therefore, we include among our predictors an additional estimate of condi-
 34 tional surprisal: We take the surprisal of the first word of the complement clause
 35 onset conditioned on the last word of the matrix clause prior to the clause
 36 juncture, regardless of whether the complementizer separates the words or
 37 not. For example, the sequence from (10) *hope (that) he* would be measured as
 38 $S_c(\text{he}|\text{hope}) = -\log_2 p(\text{he}|\text{hope})$, which we operationalize based on data from the
 39 complete British National Corpus (World Edition). Thus, we measure how sur-
 40 prising the transition would be if no complementizer had been used, under the

1 assumption that more surprising local transitions will correlate with higher
 2 shares of *that*. Importantly, despite the criticisms mentioned above, we do not
 3 intend that our measure should be seen as an alternative to the one employed
 4 by Jaeger (2010). Rather, we propose that our measure be seen to complement
 5 his at a finer granularity.

6 7 8 **4.3 Statistical evaluation: MuPDAR**

9 In order to tease apart how and why the NNS differ from the NS choices of *that*-
 10 complementation, we are using an approach called MuPDAR (Multifactorial
 11 Prediction and Deviation Analysis using Regressions), which was recently de-
 12 veloped in Gries and Deshors (2014) and Gries and Adelman (2014). MuPDAR
 13 involves the following three steps:

- 14 – fit a regression R_1 that models the choices of speakers of the target language
 15 (here, English as operationalized by the ICE-GB) with regard to the phenom-
 16 enon in question;
- 17 – apply the results of R_1 to the other speakers in the data (here, German and
 18 Spanish learners of English) to predict for each of their data points what the
 19 native speakers of the target language would have done in their situation;
- 20 – fit a regression R_2 that explores how the non-native speakers' choices differ
 21 from those of the speakers of the target/reference variety.

22
23 Crucially, in this study, both R_1 and R_2 are mixed-effects models that take into
 24 consideration the potential variability that is shared by all examples retrieved
 25 from one file and by all examples sharing the same verb (lemma), as will be
 26 detailed below; note that one can use any kind of classifier, not just regression.

27 After preparation of the data (logging several variables and factorizing
 28 others, see below), for R_1 , we began with a regression model that predicted
 29 *that*-complementation patterns of the NNS on the basis of the following pre-
 30 dictors, to which interactions were added as required by likelihood ratio tests:
 31 ComplementType, Mode, LengthCIM (factorized into three different levels given
 32 the highly skewed distribution of the data), LengthMatrixSubj (factorized into
 33 two levels), LengthMCSubjMCVerb (factorized into two levels), LengthMCVerbCC
 34 (factorized into two levels), both Delta-*P* values, and (the logged values of)
 35 LengthComplementSubj, LengthComplement, and LengthCCRemainder.³

36
37
38 ³ While factorizing numeric predictors is typically not recommended given the loss of informa-
 39 tion it incurs, we nonetheless opted for it here because initial exploratory analyses indicated
 40 potentially problematic distributional characteristics for several numeric predictors. For instance,

1 We then applied the final version of R_1 to the NNS data and added four
 2 columns to them: a column PredictionsNum (the predicted probabilities of a
 3 NS using *that* in the situation the NNS is in), PredictionsCat (the dichotomized
 4 decision following from PredictionsNum whether a NS would use *that* or not),
 5 Correct (whether the NNS made the nativelike choice or not), and, most impor-
 6 tantly at present, a column called Deviation. Deviation contains a 0 if the NNS
 7 made the nativelike choice, and it contains $0.5 - \text{PredictionsNum}$ if the NNS did
 8 not make the nativelike choice. That means, Deviation is >0 when the NNS
 9 used *that* while the NS wouldn't have, and Deviation is <0 when the NNS did
 10 not use *that* while the NS would have.

11 Finally, we developed a regression model R_2 that tries to predict Deviation,
 12 i.e. how nativelike the NNS choices were on the basis of the same predictors as
 13 in R_1 , but also adding L1 as a predictor that could interact with all others. This
 14 last predictor, through interactions, allows us to determine which factors have
 15 L1-specific effects. We began with a model involving only main effects, then
 16 added interactions of those with L1, then interactions among all predictors
 17 (using LR-tests), testing for collinearity at each step and not admitting predictors
 18 that would raise variance inflation factors (VIFs) to ≥ 5.1 . The final model of R_2 we
 19 adopted includes one predictor that was only marginally significant but interest-
 20 ing and was then explored and visualized, as outlined in the next section.

21

22

23 5 Results

24

25 5.1 Results of R_1 on the NS data

26

27 The result of the model selection process for R_1 were encouraging: R_1 featured a
 28 variety of highly significant predictors and arrived at a very good classification
 29 accuracy: 85.7% of the native speakers' *that* choices were classified correctly,
 30 which, according to exact binomial tests, is highly significantly better than
 31 either making the more frequent choice all the time (baseline₁: 68.5%) or making
 32 random choices proportional to the complementation frequencies (baseline₂:
 33 56.8%); both p 's $< 10^{-10}$. The C -value for this regression model is 0.91, thus

34

35

36 when $<10\%$ of all data points of LengthMCSbjMCVerb cover character lengths from 2 to 121, then
 37 estimating a regression slope for such a large but sparsely populated range of values is not going
 38 to yield reliable results, and a binary factorization of this predictor does not adversely affect the
 39 degrees of freedom. Also, note that factorization is a purely methodological choice – it does not
 40 reflect particular assumptions of ours regarding the cognitive mechanisms that go into selecting
 (to omit) a complementizer.

1 exceeding the typical threshold of 0.8, and the marginal and conditional R^2 are
 2 a reassuring 0.48 and 0.59. As for the random-effects structure of the model,
 3 we accounted for varying baselines of speakers to use/omit *that* (by including
 4 varying intercepts for files in the model) as well as varying preferences of verbs
 5 to use/omit *that* (by including varying intercepts for verb forms nested into
 6 lemmas in the model).

9 5.2 Applying R_1 to the NNS data

10 The application of the above regression model to the NNS data also yielded en-
 11 couraging results: the NS regression model predicted 75.2% of the NNS choices
 12 correctly, which again highly significantly (both p 's $< 10^{-100}$) exceeds both base-
 13 lines (at 0.5, because the NNS chose to realize *that* nearly half of the time); the
 14 C-value for this prediction was 0.86.

17 5.3 Results of R_2 on the NNS data

19 Computing R_2 , the model exploring to what degree NNS made nativelike choices,
 20 required a few tweaks: because of their high intercorrelations, the two Delta- P
 21 values as well as LengthCCRemainder and ComplementLength were each com-
 22 bined into a single variable (using principal component scores); the principal
 23 component for the Delta- P s, however, did not survive the model selection
 24 process. As above, we included a simple random-effects structure for files and
 25 verbs (forms nested into lemmas). R_2 returned a variety of significant predictors,
 26

28 **Table 3:** Summary results of R_2

30 Fixed effects predictor	Likelihood ratio test	p
31 LengthCIM	40.103 ($df = 2$)	< 0.0001
32 Surprisal	10.434 ($df = 1$)	0.0012
33 ComplementType : LengthComplementSubject	23.902 ($df = 2$)	< 0.0001
34 Mode : LengthComplementSubject	18.792 ($df = 1$)	< 0.0001
35 Mode : LengthMatrixSubj	19.7 ($df = 2$)	< 0.0001
36 ComplementLength/LengthCCRemainder : LengthMatrixSubj	7.531 ($df = 2$)	0.0232
37 L1 : LengthMCSbjMcVerb	8.282	0.004
38 L1 : LengthComplementSubject	2.896	0.0089 ms

1 both main effects and interactions (some pointing to L1-specific effects of the
 2 learners, some applying to both learner groups). The overall model R^2 -values
 3 are less high than those of R_1 : marginal and conditional R^2 are 0.13 and 0.3
 4 respectively. Table 3 gives a brief overview of the highest-level predictors in the
 5 final model of R_2 .

6 For reasons of space, we can unfortunately not discuss all effects in much
 7 detail; here, we will leave out the predictors involving the matrix subject. In
 8 our discussion, we will first turn to the main effects (Section 5.3.1), then we will
 9 turn to interactions, first those that apply to both learner groups (Section 5.3.2),
 10 then the ones that reveal differences between the German and Spanish learners
 11 (Section 5.3.3).

14 5.3.1 Main effects in R_2

16 Figure 1 shows the main effect of LengthCIM on Deviation: The more material
 17 precedes the main clause, the more the NNS make nativelike choices. What are
 18 the NS choices? The more material precedes the main clause, the more the NS
 19 use *that*, from 29.5% (for none) over 43.6% (for some) to 59.4% (for much). Our
 20 results show that the NNS exhibit the same tendency, but with higher propor-
 21 tions of *that*-use throughout: 44.6% over 67.5% to 77.4%. One possible explana-
 22 tion for this pattern is that, as the amount of material before the main clauses
 23 grows, both NS and NNS benefit more from inserting *that* as a structural marker
 24 between main clause and complement clause.

25 Figure 2 shows that, as the first word of the complement clause becomes
 26 more surprising given the last word of the main clause, NNS make significantly
 27 more nativelike choices. Both NS and NNS increase their complementizer use
 28 with higher rates of surprisal, and as before, the NNS just do this with a higher
 29 overall baseline of *that*-use. This difference reflects the fact that even what is
 30 expected by NS remains rather unexpected to NNS, a likely consequence of their
 31 lesser experience with naturalistic English use. Nevertheless, under conditions
 32 of high uncertainty, both groups appear to use *that* to smooth spikes in informa-
 33 tional load (as reported for NS by Jaeger 2010).

34 In sum, both overall main effects are compatible with the interpretation that
 35 NS and NNS are subject to similar processing pressures and react to them in
 36 similar ways even though NNS have a much higher baseline of *that*-use.

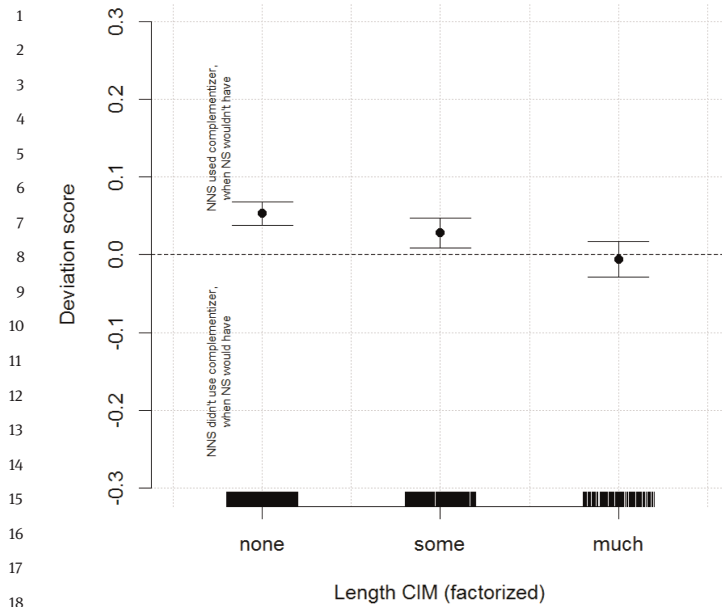


Figure 1: The effect of LengthCIM in R_2

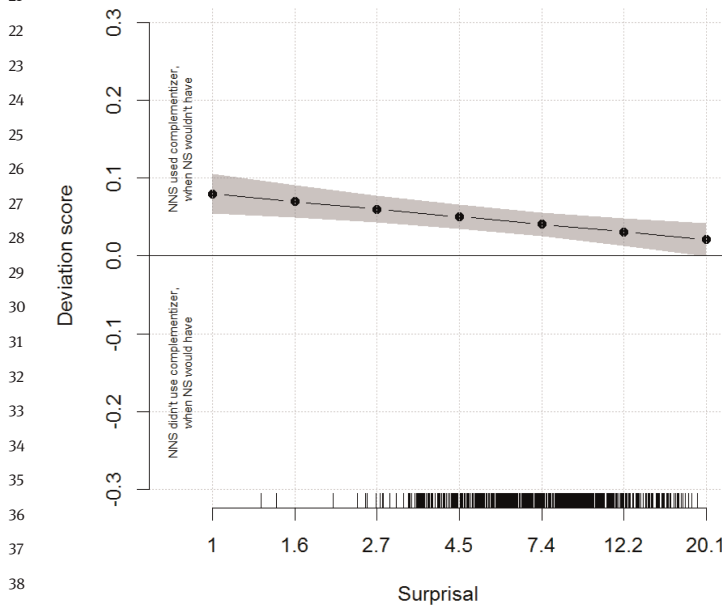


Figure 2: The effect of surprisal in R_2

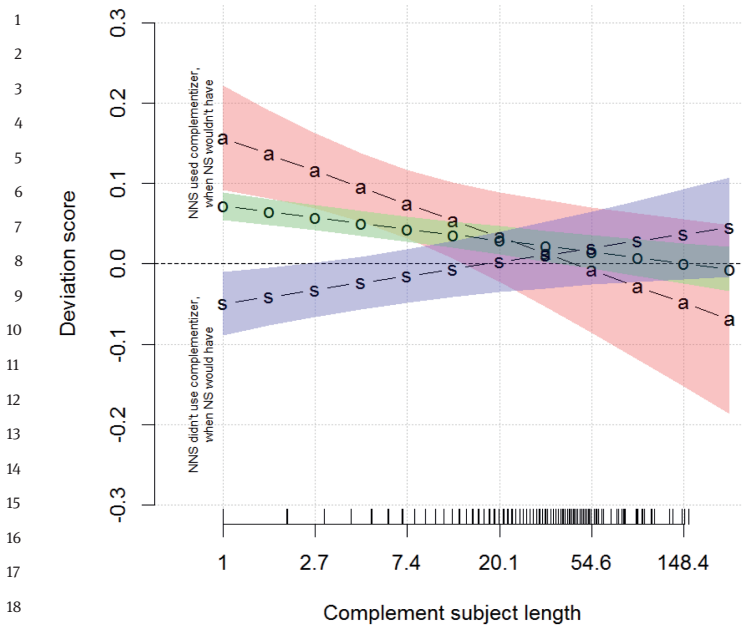


Figure 3: The effect of ComplementType : LengthComplementSubj in R_2

5.3.2 Interactions in R_2 that do not involve L1

Figure 3 shows the interaction ComplementType : LengthComplementSubj; the former predictor is represented by three regression lines with the initial letters of the complement types, the latter is represented on the x-axis. While the sample size in particular for ComplementType: Adjective is very small, as reflected in the wider confidence band, the corresponding effect in the NS data is that, with increasing length of the subject of the complement clause, speakers use *that* more. The NNS exhibit a similar trend: As the length of the subject of the complement clause increases, they also use *that* more, just like the NS. However, when the subjects of the complement clauses are short, the NNS overuse *that* in adjectival and object complement clauses and are fairly close to NS all the time in subject complement clauses. It is very plausible that this is due to transfer: In Spanish, the complementizer is obligatory in object and adjectival complement clauses, and in German, it is obligatory in adjectival complement clauses. The fact that both NNS L1s require the complementizer in at least one complement construction suggests that functionally specific transfer could be responsible for the overuse of *that* by our sample of NNS.

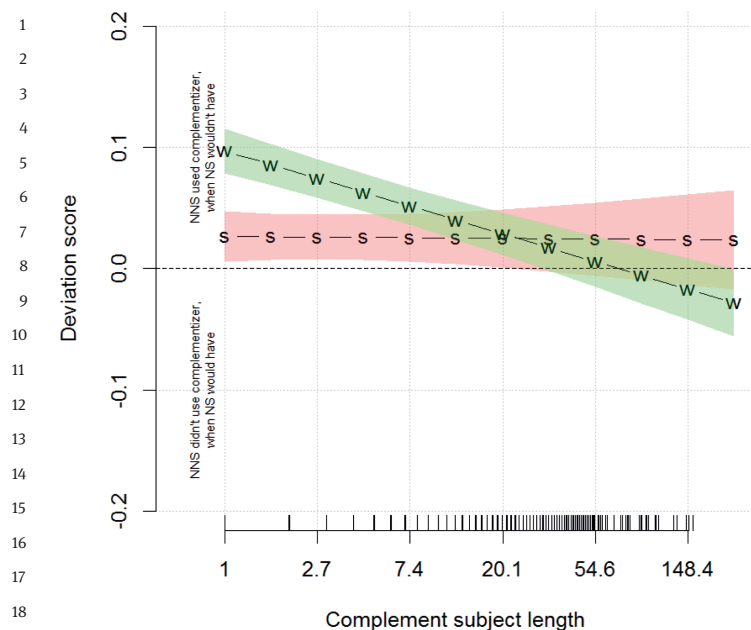


Figure 4: The effect of Mode : LengthComplementSubj in R_2

Figure 4 reflects a clear-cut effect. NS use *that* more in writing and less in speaking while the NNS are fairly close to the NS in speaking but still overuse the complementizer regardless of the length of the complement subject. In writing, on the other hand, the NNS are more nativelike with longer subjects, but overuse *that* with short subjects (in particular *I*).

Both effects show that the length of the complement clause subject is important for all speaker groups and that the learners ‘get’ the overall preference; however, due to transfer from complementizer use in their L1s and exaggerating the difference between modes, intermediate learners still need to fine-tune their preferences.

5.3.3 Interactions in R_2 that involve L1

Let us finally turn to two interactions that reveal differences between German and Spanish learners. Figure 5 shows how the two learner groups (represented with separate regression lines) react differently to the length of the subject of the complement clause. As discussed above, all speakers – NS and NNS – are

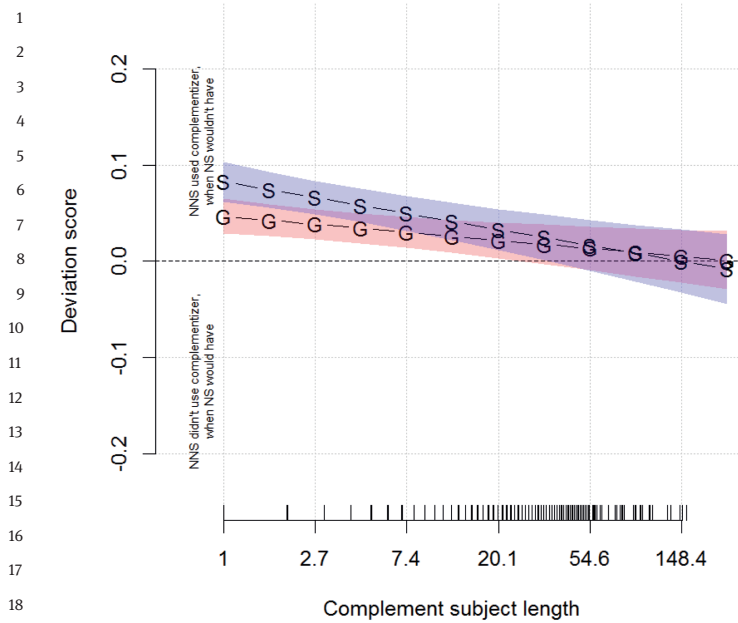


Figure 5: The effect of L1 : LengthComplementSubj in R_2

more likely to use *that* with longer complement clause subjects. However, the Germans are marginally significantly more similar to the NS with short complement subjects than the Spanish learners, who with short subject overuse *that* more than the Germans.

Finally, Figure 6 shows that, if there is material intervening between the subject and the verb of the main clause, then both German and Spanish speakers behave nativelike and use *that*, but when there is none, then both learner groups overuse *that*, and the Spanish speakers particularly much.

In sum, the German learners produce more nativelike rates of *that*-mentioning than the Spanish learners when it comes to the length effects studied in this section.

Space only permits a brief comment regarding the random-effects structure of the final model of R_2 . The largest amount of the variance of the random effects by far was accounted for by the file names, i.e. our proxy for different speakers, namely 12.5%. The second most useful random effect was the verb forms (nested into the verb lemmas), which accounted for an additional 3.5%; verb lemmas contributed an additional 3.1%. While these numbers may not seem high, they point to the need for including such effects for more accurate results than

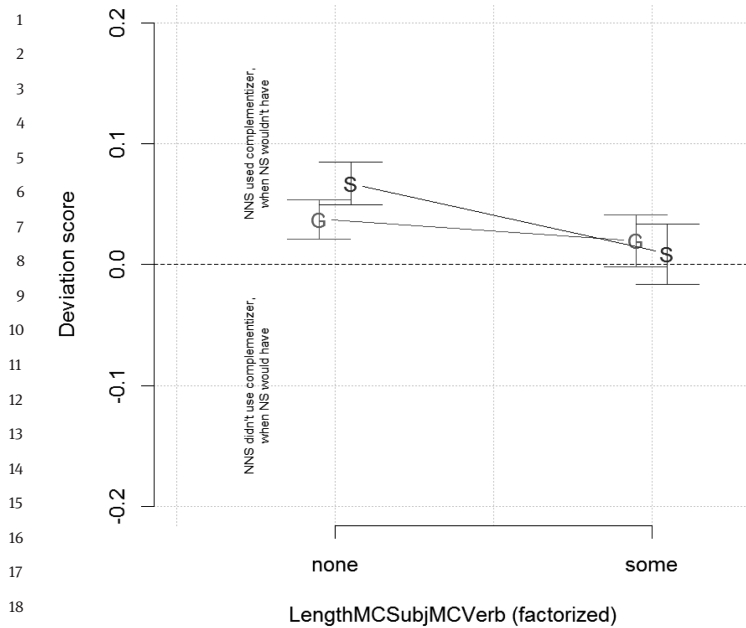


Figure 6: The effect of L1 : LengthMCSbjMCVerb in R_2

are usually provided in SLA research, and it needs to be borne in mind that our random-effects structure was restricted to varying intercepts only (given data sparsity) – more complex structures might well explain (much) more variability.

6 Discussion

The results of the MuPDAR analysis suggest that the intermediate-advanced German and Spanish learners are quite well aligned with NS norms overall. Minor (yet significant) differences were identified in the second regression: both learner groups employ comparatively higher shares of *that* as the processing demands increase, be it in the form of more material occurring at the onset of the clause or with longer complement subjects. More pronounced differences between NS and learners become visible when we consider construction-specific uses of *that* – learners overuse the complementizer in adjectival and object constructions – and register-specific uses of *that*: both learner groups overuse the complementizer especially in writing when the main clause subject is *I*. Finally, a few L1-specific differences emerge: the Spanish learners overuse the

1 complementizer more frequently than their German peers do in contexts with
2 short complement clause subjects and when clause juncture is interrupted.

3 These findings suggest that the intermediate-advanced learners examined
4 here rely on the same basic mechanisms governing *that*-variation as native
5 speakers, but at the same time display a comparatively more conservative
6 behavior than the native speakers: learners produce the complementizerless
7 utterances only in what we may call “ideal contexts” associated with low shares
8 of *that* also in NS use, namely in speaking, with short subject and complement
9 clause subjects, and with little or no increased processing costs imposed by
10 optional additional and/or intervening material. When the context is less than
11 ideal, the learners – and the Spanish learners more so than their German peers,
12 arguably reflecting transfer from the L1 – resort to the “safe” strategy of realizing
13 the complementizer as this choice is never, strictly speaking, ungrammatical, if
14 only, at times, non-idiomatic.

15 Generally speaking, the learner behavior is not fundamentally different from
16 NS behavior; rather, the thresholds for producing the complementizer are signif-
17 icantly lower compared to NS speakers, and they are reactive to the factors
18 mentioned above. This stands in accord with usage-based models of L2 learning
19 such as N. Ellis’ Associative-Cognitive CREED model (Ellis and Wulff 2015) or
20 Goldberg’s (2006) usage-based Construction Grammar, to name but two exam-
21 ples; these models share the assumption that L2 learning is best characterized
22 as the gradual approximation towards native-like representations. As one reviewer
23 pointed out, questions regarding which specific mechanisms underlie the factors
24 included here – cognitive load, learning as a result of usage, transfer effects,
25 and/or instructional effects –, and how exactly each these mechanisms operate
26 in the individual learner – even something as seemingly straightforward as
27 cognitive load can be manifested on different levels of linguistic analysis and
28 can interact with general intelligence, working memory, age, etc. – are beyond
29 the scope of the present analysis, and possibly beyond a purely corpus-based
30 approach. In the following, we can only speculate about the relationship
31 between these factors and the cognitive mechanisms they potentially tap into.

32 Firstly, it is with regard to processing-related factors such as clause com-
33 plexity and juncture that we see learners in need to further improve their align-
34 ments to the target norm. This reminds us of psycholinguistic accounts such
35 as that of Kroll and her colleagues, who argue in favor of a tight link between
36 bilingualism and cognitive cost: according to Kroll and Dussias (2013), speaking
37 a second language entails a higher cognitive load because the speaker const-
38 antly has to juggle between the two (or more) languages (Kroll and Dussias
39 2013). From that perspective, it makes sense that our learners display lower
40 tolerance thresholds for factors that themselves are directly related to cognitive

1 cost, such as complexity or clause juncture: compared to native speakers, the
2 learners have fewer cognitive resources to allot in the first place. As a result,
3 they produce the complementizer more frequently.

4 In addition, we found that NS and NNS both responded in the expected
5 fashion to spikes in uncertainty (based on Jaeger 2010) as captured by the con-
6 ditional surprisal of the first word of the complement clause given the last word
7 of the matrix clause. Both groups were more likely to produce *that* at high-
8 uncertainty transitions. However, NNS also tended to overproduce *that* at lower
9 surprisal junctures, suggesting again a conservative strategy. This effect, like
10 that discussed above, is amenable to explanation in terms of cognitive cost,
11 with NNS experiencing greater difficulty with transitions that are otherwise
12 unproblematic for native speakers, but converging on native performance when
13 the transitions reach a certain threshold of uncertainty.

14 As far as the implications of the present study for language teaching are
15 concerned, one may conclude that overall, *that*-variation does not constitute an
16 insurmountable challenge to learners: in spite of the fact that proper com-
17 plementizer use is hardly if ever a topic of explicit classroom instruction, the
18 intermediate-advanced learners investigated here seem to be well on their
19 way to nearly native-like behavior. *That*-variation may be taken as a powerful
20 example of how much learners can pick up by implicitly scrutinizing the dis-
21 tributional patterns of their input even though the random effects also showcase
22 considerable individual variation. That said, the results, of course, point to room
23 for improvement. For one, instruction could focus more on complementizer
24 variability by comparing the L1 with the L2; especially the Spanish learners
25 may benefit from their attention being directed at the optionality of *that* in
26 adjectival and object complements in particular. Similarly, increasing awareness
27 for mode-dependent differences may be useful for both learner groups examined
28 here.

31 References

- 32
33 Biber, Douglas. 1999. A register perspective on grammar and discourse: variability in the form
34 and use of English complement clauses. *Discourse Studies* 1. 131–50.
35 Bryant, Margaret M. 1962. *Current American usage*. New York: Funk & Wagnalls.
36 Dor, Daniel. 2005. Toward a semantic account of *that*-deletion in English. *Linguistics* 43(2).
37 345–382.
38 Durham, Mercedes. 2011. I think (that) something's missing: Complementizer deletion in non-
39 native emails. *Studies in Second Language Learning and Teaching* 1(3). 421–445.
40 Ellis, Nick C. 2006. Language acquisition as rational contingency learning. *Applied Linguistics*
27. 1–24.

- 1 Ellis, Nick C. & Stefanie Wulff. 2015. Usage-based approaches in second language acquisition.
 2 In Bill VanPatten & Jessica Williams (eds.), *Theories in second language acquisition: An*
 3 *introduction*, 75–93. London & New York: Routledge.
- 4 Elness, Johan. 1984. *That* or zero? A look at the choice of objective clause connective in a
 5 corpus of American English. *English Studies* 65. 519–533.
- 6 Gilquin, Gaëtanelle, Sylvie de Cock & Sylviane Granger. 2010. *Louvain International Database*
 7 *of Spoken English Interlanguage*. Louvain-la-Neuve: Presses universitaires de Louvain.
- 8 Goldberg, Adele E. 2006. *Constructions at work: The nature of generalization in language*.
 9 Oxford: Oxford University Press.
- 10 Granger, Sylviane, Estelle Dagneaux, Fanny Meunier & Magali Paquot. 2009. *International*
 11 *Corpus of Learner English v2*. Louvain-la-Neuve: Presses Universitaires de Louvain.
- 12 Gries, Stefan Th. 2007. *Coll.analysis 3.2*. A program for R for Windows.
- 13 Gries, Stefan Th. 2013. 50-something years of work on collocations: What is or should be next ...
 14 *International Journal of Corpus Linguistics* 18(1). 137–165.
- 15 Gries, Stefan Th. & Allison S. Adelman. 2014. Subject realization in Japanese conversation by
 16 native and non-native speakers: Exemplifying a new paradigm for learner corpus research.
 17 In Jesús Romero-Trillo (eds.), *Yearbook of corpus linguistics and pragmatics 2014: New*
 18 *empirical and theoretical paradigms*, 35–54. Cham: Springer.
- 19 Gries, Stefan Th. & Sandra C. Deshors. 2014. Using regressions to explore deviations between
 20 corpus data and a standard/target: Two suggestions. *Corpora* 9(1). 109–136.
- 21 Jaeger, T. Florian. 2010. Redundancy and reduction: Speakers manage syntactic information
 22 density. *Cognitive Psychology* 61. 23–62.
- 23 Kaltenböck, Gunther. 2006. ‘... *That* is the question’: Complementizer omission in extraposed
 24 *that*-clauses. *English Language and Linguistics* 10(2). 371–96.
- 25 Kroll, Judith F. & Paola E. Dussias. 2013. The comprehension of words and sentences in two
 26 languages. In Tej K. Bhatia & William C. Ritchie (eds.), *The handbook of bilingualism and*
 27 *multilingualism*, 216–243. Malden, MA: Wiley-Blackwell Publishers.
- 28 Levy, Roger. 2008. Expectation-based syntactic comprehension. *Cognition* 106(3). 1126–1177.
- 29 Rissanen, Matti. 1991. On the history of *that/zero* as object clause links in English. In Karin
 30 Aijmer & Bengt Altenberg (eds.), *English corpus linguistics*, 272–289. London: Longman.
- 31 Storms, G. 1966. *That*-clauses in Modern English. *English Studies* 47. 249–70.
- 32 Tagliamonte, Sali A. & Jennifer Smith. 2005. *No momentary fancy!* The zero ‘complementizer’ in
 33 English dialects. *English Language and Linguistics* 9(2). 289–309.
- 34 Thompson, Sandra A. & Anthony Mulac. 1991. The discourse conditions for the use of the
 35 complementizer *that* in conversational English. *Journal of Pragmatics* 15. 237–51.
- 36 Torres Cacoullos, Rena & James A. Walker. 2009. On the persistence of grammar in discourse
 37 formulas: A variationist study of *that*. *Linguistics* 47. 1–43.
- 38 Wulff, Stefanie. 2016. A friendly conspiracy of input, L1, and processing demands: *that*-
 39 variation in German and Spanish learner language. In Andrea Tyler, Lourdes Ortega, Hae
 40 In Park, & Mariko Uno (eds.), *The usage-based study of language learning and multi-*
lingualism. Georgetown: Georgetown University Press.
- Wulff, Stefanie, Nicholas A. Lester & Maria M. Martinez-Garcia. 2014. *That*-variation in German
 and Spanish L2 English. *Language and Cognition* 6. 271–299.