

# Lexical access and lexical diversity in first language attrition\*

MONIKA S. SCHMID  
*University of Essex/University of Groningen*  
SCOTT JARVIS  
*Ohio University*

(Received: August 7, 2012; final revision received: November 14, 2013; accepted: November 19, 2013; first published online 16 January 2014)

*This paper presents an investigation of lexical first language (L1) attrition, asking how a decrease in lexical accessibility manifests itself in long-term residents in a second language (L2) environment. We question the measures typically used in attrition studies (formal tasks and type–token ratios) and argue for an in-depth analysis of free spoken data, including factors such as lexical frequency and distributional measures. The study is based on controlled, elicited and free data from two populations of attriters of L1 German (L2 Dutch and English) and a control population (n = 53 in each group). Group comparisons and a Discriminant Analysis show that lexical diversity, sophistication and the distribution of items across the text in free speech are better predictors of group membership than formal tasks or elicited narratives. Extralinguistic factors, such as frequency of exposure and use or length of residence, have no predictive power for our results.*

Keywords: language attrition, lexical attrition, mental lexicon, methodology

When (monolingual) native speakers leave their country of origin, take up life in a different linguistic environment, become bilingual and consequently have less input in and make less use of their first language (L1), they often find this language changes. This change can manifest itself in lexical access difficulties, disfluency phenomena, cross-linguistic interference, increased optionality in grammatical features and a foreign accent (among other things). The linguistic development experienced in such situations is commonly referred to as L1 ATTRITION. It has often been described as a selective process (e.g. Sorace, 2005; Tsimpli, 2007) which affects different components of linguistic knowledge in a different order, at a different rate and to different extents (Paradis, 2007), and which is sensitive to external factors such as age of emigration, length of residence and amount of L1/L2 exposure.

It has become almost axiomatic in language attrition research to assume that lexical-semantic knowledge is the most vulnerable part of the linguistic repertoire, deteriorating first, fastest and most dramatically as compared to, for example, grammar or phonetics (Hulsen, 2000; Köpke & Nespoulous, 2001; Köpke & Schmid, 2004; Montrul, 2008; Opitz, 2011, to name but a few). It is also commonly believed that lexical deterioration will be more or less linearly related to frequency of L1 use (e.g. Paradis, 2007). However, the empirical record

regarding L1 lexical attrition, its relationship with other areas of attrition and the factors that drive it, is quite weak. We argue that this is due to a paucity of studies that have approached L1 lexical attrition with a sufficiently stringent methodological design and sufficiently solid theoretical underpinnings. In particular, there are few studies comparing data elicited by means of multiple tasks or contrasting different theoretical predictions and assessing the impact of environmental variables, such as language use.

The present study focuses on lexical accessibility and explores the ways in which the twofold phenomenon of less exposure to the L1 and more exposure to an L2 may result in reduced L1 lexical accessibility. We also enquire whether the detectability of L1 lexical reduction varies across different tasks. In order to obtain better insights into this phenomenon, the present study combines and compares different tasks and different measures and assesses their explanatory potential.

## Language attrition in the lexicon

The term ‘language attrition’<sup>1</sup> refers to changes in a native language that has either fallen into disuse or is used alongside an environmental one. In accordance with this definition, attrition is a process that is driven

\* The authors are particularly grateful for the kind help of Franck Bodmer Mory of the *Institut für deutsche Sprache* at Mannheim for assessing the frequency of the lemmas from our study against the COSMAS II Corpus.

<sup>1</sup> In keeping with common practice, we reserve the term ‘language attrition’ for the attrition of a native language, while the attrition of later learned languages is referred to as second or foreign language attrition.

Address for correspondence:

Monika S. Schmid, Department of Language and Linguistics, University of Essex, Wivenhoe Park, Colchester CO4 3SQ, United Kingdom  
mschmid@essex.ac.uk

by two factors: (a) the presence, development and regular use of a second linguistic system, leading to crosslinguistic interference (CLI), competition and other effects associated with bilingualism, and (b) a decreased use of the attriting language, potentially leading to access problems (Schmid & Köpke, 2007).

A number of quantitative investigations have investigated the extent to which *lexical access* is compromised in L1 attrition due to long-term underuse, that is, whether attriters take longer to retrieve items from the mental lexicon or have a more restricted set of items at their disposal than unattrited monolinguals. The earliest such studies mainly relied on verbal fluency (VFT, e.g. Waas, 1996; Yağmur, 1997) or picture naming tasks (PNT, e.g. Ammerlaan, 1996; Hulsen, 2000). More recently, it has been argued that such experimental approaches should be combined with analyses of lexical access and lexical diversity in (relatively) unguided free speech (Schmid, 2004) as well as with investigations of disfluency phenomena (Schmid & Beers Fägersten, 2010). Such comparisons not only allow an assessment of lexical diversity across the full range of a speaker's repertoire (while controlled tasks are necessarily limited by the stimuli used); they may also show whether attrition effects are less (or more) pronounced in controlled tasks that allow the participant to fully focus on lexical retrieval than in naturalistic language production situations where information from all linguistic levels must be rapidly integrated in real time. If it could be shown that attrition effects are more pronounced in the latter type of data, this might indicate that attrition itself is largely the outcome of the increased cognitive load involved in managing two linguistic systems at the same time, as opposed to lexical access problems arising from an increase in activation thresholds.

### **The role of L1 use**

The assumption that (lexical) attrition is closely related to frequency of L1 use bears close resemblance to Michel Paradis' Activation Threshold Hypothesis (Paradis 1993, 2004, 2007; henceforth ATH). This hypothesis takes as its point of departure the fact that accessing items stored in memory requires a certain amount of neuronal excitation. The level of energy necessary to retrieve an item is determined by the frequency and recency with which it has previously been called upon, so that less frequent items and items that have not been used for a long time become harder to access. Attrition is therefore hypothesized to predominantly affect less frequent linguistic/lexical items, and to be more pronounced for speakers who do not use their L1 on a regular basis (Andersen, 1982; Paradis, 2007).

Furthermore, the activation threshold is determined not solely by activation but also by inhibition. Each

time a speaker selects a target item from the lexicon, its competitors (semantic/phonological neighbors, cognates, translation equivalents, etc.) must be inhibited. This inhibition process raises the activation threshold, making the inhibited items harder to access subsequently (Green, 1986, 1998). The phenomena we can witness in the attritional process are thus not only dependent on the underuse of the attriting system (raising of the AT due to non-activation), but also on the presence, use and development of the environmental language (raising of the AT due to inhibition). Bilinguals routinely have to inhibit the language that is not chosen for activation (e.g. Bialystok, 2005), and this may lead to phenomena such as word-finding difficulties. It is therefore important to take into consideration not only how much each language is used, but also its contexts of use (formal/informal, monolingual/bilingual interlocutor), which affect the degree to which the other language is inhibited in these situations (Schmid, 2007).

In this connection it is interesting to note that there is only one language use factor which has emerged as a significant predictor of attrition across a range of skills and linguistic levels in a number of studies (e.g. de Leeuw, Schmid & Mennen, 2010; Schmid, 2007; Schmid & Dusseldorp, 2010), namely the use of the L1 for professional purposes: the more migrants use the L1 at work, the less L1 attrition they exhibit. The differential impact of informal vs. formal (i.e., professional) L1 use can be explained in terms of inhibition: in the daily life of a migrant, most interlocutors with whom the L1 is spoken informally (friends, family members) are also bilingual, so that code-switches and code-mixing do not need to be inhibited. In a professional context, however, it is usually not considered appropriate to mix the two languages, so that speakers who use their L1 at work will probably have more practice at inhibiting the L2 and resisting any 'intrusions' (lexical or otherwise) from this language system. This increased practice at inhibiting the environmental linguistic system may explain why such speakers can perform better in their L1 than speakers who do not have occasion to use it in these types of formal situations (Schmid, 2007).

A further intriguing question in the context of the ATH is to what extent lexical attrition phenomena differ between contact languages of varying typological distance. If attrition is predominantly determined by frequency/recency of activation, speakers of two languages that share a large part of their lexicon should have an advantage over migrants who are bilingual in typologically more distant languages, since the activation of cognates can be assumed to spread across languages (Berthele, 2011; Dijkstra, 2005; Jarvis, 2009). If, on the other hand, inhibition is the driving force of lexical access difficulties, speakers of more distant languages can be expected to encounter fewer difficulties, as there

will be less competition between translation equivalents which share no surface similarity. This implies that attrition studies would benefit from making comparisons of populations with either different L1s or different L2s instead of focusing on a single language combination (as has almost invariably been the case so far).

### Measuring lexical diversity

It was pointed out above that it is desirable for investigations of lexical attrition to combine a range of tasks probing lexical accessibility. Controlled tasks, such as measuring a person's response times in naming a picture stimulus, allow participants to focus all their attention on the lexical retrieval process. In free speech, on the other hand, language processing takes place across many linguistic levels simultaneously, and there may be trade-off effects, for example between morphosyntactic complexity/accuracy and lexical diversity, or between the use of less frequent linguistic items and fluency. This suggests that attriters, who may be experiencing difficulties with other parts of the linguistic repertoire, might show larger differences from controls in free speech than in controlled tasks.

Investigations of (elicited) free speech have to tackle the problem of how to measure lexical diversity in such data. This question has received considerable attention recently (for a comprehensive overview, see Jarvis, 2012, 2013a, 2013b). Traditionally, it has been addressed with type-token ratio (TTR) based measures, which relate the total number of words in a text to the total number of different lemmas (e.g. Jarvis, 2002; Schmid, 2011). Simple TTRs have been shown to be a problematic measure of lexical diversity in free speech because they vary as a function of text length: the rate of word repetition inevitably increases as the text grows longer. McCarthy and Jarvis (2007) demonstrate that this problem is persistent even for measures of lexical diversity which have been devised to overcome the impact of length, such as Guiraud's index, Yule's *K* and VOCD (McCarthy & Jarvis, 2007; Jarvis & Daller, 2013). While the effects of length on probability-based measures such as Yule's *K* and VOCD are relatively subtle, the principles of probability render them sensitive to a second factor: evenness, which is a matter of how evenly the tokens in a sample are distributed across types. Jarvis (2012, 2013a, 2013b) argues that length and evenness are inherent properties of diversity. This consideration is particularly relevant for situations such as language attrition: where lexical access is compromised, activation of specific items might be contextually determined and spread across a narrower range of the semantic field, leading to a more uneven distribution of the tokens across the entire text than would be the case for monolinguals.

The first problem (text length or volume) can be solved through the use of a measure of textual lexical diversity (MTLD), which has been applied in recent studies to remove the effects of text length (McCarthy & Jarvis, 2010, 2013). It is calculated as the average number of running words in a text that remain above a certain TTR (usually .72). Evenness can be assessed through measures used in the field of ecology in studies dealing with biodiversity. Although numerous indices of evenness exist (Smith & Wilson, 1996), the one we have adopted for the present study is that based on Shannon's entropy index (Shannon, 1948), which has been described by Chao and Jost (in press) as the most appropriate general-purpose measure of diversity because it avoids giving too much weight to either rare or abundant species. Evenness is calculated as a ratio of the observed diversity of a sample (i.e., Shannon's index) divided by the maximum possible diversity of the sample that would occur if all types (or species) in the data were equally abundant (Pielou, 1969).

So far, no measures that attempt to correct the problems of volume and evenness have been applied in language attrition studies, and lexical diversity in free speech has invariably been described in terms of measures such as TTR, VOCD or Guiraud. These are also the measures included in the test battery proposed by Schmid (2004, 2011, see also [www.let.rug.nl/languageattrition](http://www.let.rug.nl/languageattrition)), which combines controlled and free speech tasks and has recently been applied in a number of investigations of the attrition of a variety of languages.

As can be seen in Table 1, the results from these studies appear to be somewhat inconsistent: with one exception, all the investigations did find a statistically significant increase in (at least some) disfluency markers (e.g. empty pauses, filled pauses, repetitions). Lexical diversity (VOCD) in free speech and performance on the verbal fluency task, on the other hand, did not differ consistently across populations and studies.

A closer look at Table 1 reveals a relationship between findings and sample size for these measures: in investigations with relatively small population sizes (25 or lower), statistical significance is not reached (and, with the exception of Dostert, 2009, whose attriters outperform the controls on a number of measures, all the authors note that the descriptive statistics show better performance of the controls on most tasks). The lack of significant findings in the investigations by Cherciov (2011), Opitz (2011) and Varga (2012) may therefore be a Type II error due to the limited sample size. All of the larger studies show consistent differences between attriters and controls in free speech and on the VFT. The PNT used by Yılmaz and Schmid (2012), on the other hand, does not reveal a significant difference between attriters and controls, despite the comparatively large sample of 54 speakers in each population. This suggests that there may be

Table 1. *Populations, tasks and findings of recent investigations of lexical attrition*

	Cherciov 2011	Dostert 2009	Keijzer 2007	Opitz 2011	Schmid & Dusseldorp 2010	Varga 2012	Yilmaz & Schmid, 2012
L1	Romanian	English	Dutch	German	German	Hungarian	Turkish
L2	English	German	English	English	English/ Dutch	Danish	Dutch
n attriters	20	25	45	13	106	20	54
n controls	15	20	45	17	53	20	54
Formal task	VF	VF	VF	VF	VF	VF	PNT
Free speech elicitation task	CC	CC, picture description	CC	CC	CC	CC	interview
Significant difference attriters/ controls on controlled task	no	no	yes	no	yes	yes	no
Significant difference attriters/ controls on lexical diversity in free speech (VOCD)	no	no	yes	no	yes	no	yes
Significant difference attriters/ controls on fluency in free speech	yes	yes	not assessed	no, except for repetitions	yes	no	yes

population-specific constraints on lexical attrition, since the speakers investigated in this particular study (Turks in the Netherlands) are drawn from the numerically largest group of non-Western immigrants in the host country and may therefore have a larger L1 network providing support for language maintenance than is the case in the other investigations. Similarly, the fact that Dostert did not find any attrition of L1 English may indicate that the global importance and presence of this language exerts a protective effect even in the absence of dense personal networks where the L1 is spoken.

### Summary and research questions

Although lexical attrition has long been assumed to be one of the earliest and most noticeable symptoms of attrition in general, relatively few studies so far have attempted to gain a comprehensive view of exactly what this process entails. Especially with respect to the predictions concerning lexical reduction made by Andersen (1982) and echoed in the ATH, there is to date little solid evidence of whether and to what extent these hold true. To the extent that these predictions have been tested, this has generally been done exclusively through the use of controlled tasks. However, controlled tasks are often limited in their potential to probe the reduction of lexical accessibility, as they rely on certain characteristics of the stimuli (e.g. imageability in picture naming, limitations to a certain lexical field in semantic verbal fluency tasks). Previous studies that

have examined lexical attrition in free speech have also suffered from a second shortcoming in that they have so far relied on measures of type–token relationships, which do not allow insights into lexical frequency/sophistication or into characteristics of the distribution of items across the text (evenness).

A second issue that is of theoretical interest is the role of activation vs. inhibition in lexical attrition: if attriters do have lexical access problems, are these linked to the fact that L1 items are more difficult to access (because they are used infrequently) or to problems involving the inhibition of closely related L2 items? This question may to some extent be resolved through a comparison of two groups of attriters who speak the same L1 but a different L2.

The present study address both of these issues and also explores the impact of personal background and language use factors on lexical attrition. It is guided by the following research questions:

RQ1: Lexical access: Is lexical access in the L1 compromised in long-term migrants (henceforth: attriters)? To what extent are such access problems differentially detectable by formal tasks vs. in spontaneous language use?

RQ1a: Are attriters outperformed by predominantly monolingual speakers in their country of origin (henceforth: controls) on a controlled lexical access task?

Table 2. *Participant characteristics*

		Control group	Attriters in Canada	Attriters in the Netherlands
Age	Mean	60.89	63.23	63.36
	Maximum	91	88	85
	Minimum	39	37	37
Age at emigration	Mean	.	26.13	29.08
	Maximum	.	47	51
	Minimum	.	14	17
Length of residence	Mean	.	37.09	34.28
	Maximum	.	54	58
	Minimum	.	9	14

RQ1b: Do attriters use a less diverse lexical repertoire in free speech, as assessed by a number of diversity measures?

RQ1c: Do attriters show an overall preference for more “common, high-frequency” lexical items and a tendency to underuse “less-common, low-frequency” items (as predicted by Andersen, 1982)?

RQ1d: Can attrition effects across items of varying lexical frequency be predominantly ascribed to problems of activation due to non-use, or to problems of inhibition related to the increased use of L2?

RQ2: Measuring attrition: What measures of lexical diversity are best suited to the detection of lexical attrition in free speech?

RQ2a: Are type–token-based measurements good indicators of lexical attrition?

RQ2b: What measurements are suitable for detecting differential attrition effects across items of varying lexical frequency?

RQ2c: What measures are suitable for detecting attrition effects linked to the distribution, as opposed to the selection, of lexical items in the wider discourse?

RQ3: Extralinguistic variables: Which sociolinguistic and personal background variables have an impact on lexical attrition?

RQ3a: How does L1 use in a variety of settings (formal vs. informal, interactional vs. exposure etc.) affect lexical attrition?

RQ3b: How do personal background factors, such as age at emigration, length of residence and level of education affect lexical attrition?

RQ3c: To what extent does the degree of similarity between L1 and L2 lexicons affect lexical attrition?

## The study

### Participants

The present investigation is based on two verbal fluency tasks as well as on free and elicited speech collected from 159 native speakers of German, 53 of whom emigrated to Anglophone Canada (Vancouver area) and 53 to the Netherlands at least nine years before the data collection at age 14 or older. A third group of 53 had lived in Germany all their lives and never routinely used a language other than German. The three groups were matched for gender, age and education. Participant characteristics are summarized in Table 2.

A sociolinguistic questionnaire comprising 78 items was used to assess the bilingual speakers’ history, linguistic habits and attitudes towards the L1 and L2 (for details see [www.let.rug.nl/languageattrition/SQ](http://www.let.rug.nl/languageattrition/SQ) and Schmid & Dusseldorp, 2010), and also formed the basis for free speech production (see below). All variables were coded on a scale between 0 and 1, where 0 means no use of the L1 or a strong preference for the L2, while 1 indicates very high use or strong preference for the L1. In order to reduce these variables to a realistic number of predictors for statistical analysis, the following compound factors were calculated, based on the procedures suggested by Schmid and Dusseldorp (2010):

*Total L1 use:* This variable was an average of a total of 17 questions pertaining to the use of German in informal settings (Cronbach  $\alpha = .919$ ). It comprised the following subfactors:

- overall frequency of use of L1 (one question)
- frequency of use of L1 within the family (one question)
- language used most frequently on a daily basis (one question)
- frequency of use of L1 with the partner (four questions)

- frequency of use of L1 with children (where applicable, four questions)
- frequency of use of L1 with friends (four questions)
- frequency of visits to Germany (one question)
- frequency of contacts (telephone, letters, email, etc.) with Germany (one question).

*Affiliation with L1 language and culture:* This variable was based on a total of six questions about linguistic and cultural affiliations and preferences (Cronbach  $\alpha = .706$ ). It comprised the following factors:

- importance of maintaining the L1
- preferred culture
- preferred language
- language and culture with which there are the strongest emotional ties
- language and culture with which speaker identifies most
- if all external constraints were removed, would speaker like to move back to Germany?

*L1 at work:* Previous investigations of L1 attrition have established that the use of the L1 for professional purposes often emerges as the most important external predictor for attrition across a range of skills (see above). This factor, based on one question (how frequently do you use German for professional purposes?) is therefore also included here.

The two attriting populations behave very similarly with respect to the amount they use their L1 in their daily lives (none of the predictors differs significantly for the two populations; see Table 1 in the associated Supplementary Materials Online). Of course the geographical proximity of the Netherlands to Germany may impact to some extent on the opportunity that individual speakers have to use their L1 (e.g. in terms of how often they can visit their home country). However, in their close social networks, daily contacts and professional lives, both groups appear to use the L1 with approximately equal frequency, and they are also similar in their cultural and linguistic preferences.

While the two bilingual populations are similar to each other on these factors, there is considerable individual variability, some individuals in both groups having frequent contact and high affiliation with German while others use it only sporadically and have low affiliation.

### **Verbal fluency tasks**

The investigation used two semantic verbal fluency (VF) tasks. In each task, the participant was given 60 seconds to produce as many items from a given semantic category as possible. The two categories used were a) fruit and vegetables and b) animals. Overall production on the two tasks was averaged for each participant (VFTot). Furthermore, each task was divided into six segments of

equal length (ten seconds), and the average productivity of each speaker in each segment was calculated.

### **Speech samples**

Two speech samples were collected from each speaker. The first was a conversation about the individual's history and biography, which typically lasted between 30 and 90 minutes. This interview allowed the participants to make use of the full range of their language skills on a variety of topics that most migrants frequently talk about, and thus imposed few contextual constraints. The second speech sample was a monologue-narrative in which participants watched and subsequently narrated a ten-minute film sequence from the silent Charlie Chaplin movie "Modern Times" (Perdue, 1993; for the full elicitation procedure see Schmid, 2011). This task was assumed to be somewhat more challenging, as it required the participant to describe a specific set of events and items, while also imposing the extra cognitive task of remembering the narrative line.

Due to equipment failure and time constraints, there were a few participants from whom one of the two speech samples could not be elicited. The sociolinguistic interviews ( $n = 153$ ) comprised a total of 378,740 tokens (not counting filled pauses, false starts, repetitions and self-corrections). The distribution of these data across the three groups is given in Supplementary Table 2. The Charlie Chaplin film retellings ( $n = 155$ ) comprised a total of 110,270 tokens; for full details see Supplementary Table 3.

### **Lexical diversity analyses**

Our approach to lexical diversity in the present study involves the use of type-token-based measures, such as VOCD, as well as measures designed to solve the problematic areas discussed above. The full list of the diversity-related measures and indices used in this study is:

- Types (the number of different words in each sample)
- Tokens (the total number of words in each sample)
- VOCD (a probability-based measure of diversity derived through random sampling)
- MTLD (the mean number of running words that remains above a TTR threshold of .72)
- Shannon's index (a probability-based measure of diversity that takes both types and evenness into account)
- Evenness (how evenly tokens are spread across types, derived from Shannon's index)
- Effective types (the number of types in the sample adjusted downward in accordance with the degree of unevenness in the sample; see Chao & Jost, in press)

Table 3. Frequency bands of lemmatized lexical items in the two corpora

	present corpus				COSMAS II corpus			
	INT		CC		INT		CC	
	# of tokens	# of types	# of tokens	# of types	# of tokens	# of types	# of tokens	# of types
Frequency band 1 (20% most frequent items)	18,560	10	6,298	10	133,966,582	28	77,562,547	20
Frequency band 2	18,243	31	6,238	29	133,551,407	131	78,407,451	60
Frequency band 3	18,238	116	6,259	74	133,778,194	328	78,781,247	131
Frequency band 4	18,242	612	6,293	262	133,872,371	762	78,860,066	293
Frequency band 5 (20% least frequent items)	17,959	7,946	6,354	2,707	133,676,872	7466	80,801,470	2574

- Rarity (the mean frequency rank of the types in the sample; based on the COSMAS II corpus, which we describe in more detail below)
- Dispersion (the mean number of words between tokens of the same type).

We also investigated the development of lexical sophistication (that is, the use of words with varying frequencies) in the attritional process. For both sets of files (interview and film retelling), our analysis of lexical diversity was based on the lemmatized version of the files created by means of the procedure for morphological tagging and disambiguation described above. The VOCD measure was calculated for this lemmatized corpus with the help of the CLAN program.

Lemmatization was checked manually in the lists that had been extracted from each file. From this list, all lexical items (nouns, full verbs and adjectives) were retained. This yielded a total of 8,715 lemmatized lexical types (91,242 tokens) in the interview files and 3,082 lemmatized lexical types (31,442 tokens) in the film retelling files.

Since one of the predictions for lexical attrition is that attriters will come to prefer more frequent words and underuse less frequent ones, overall frequency was assessed in two ways. First, it was established how frequently each lemma was used within its relevant corpus (the interview or film retelling files). Second, the total frequency of use for each lemma was also measured in the COSMAS II corpus (Corpus Search Management and Analysis, maintained by the *Institut für deutsche Sprache* at Mannheim and based on ca. 5.4 billion word forms)<sup>2</sup>.

<sup>2</sup> The authors would like to express their gratitude to the *Institut für deutsche Sprache* for making these data available, and in particular to Franck Bodmer Mory, who very kindly conducted the analyses for us. The COSMAS II corpus contains texts from Germany, Austria and Switzerland. The frequency with which lexical items were used in our own corpora (interview and film retelling) correlated significantly ( $p < .001$ ) with the COSMAS II frequencies, and these correlations were highest when only the data from Germany were considered (interview:  $r = .239$ , film retelling:  $r = .543$ ; for the data from Austria

Lexical sophistication was then assessed by dividing all types into five frequency bands on the basis of both our own corpora and the COSMAS II corpus. Each of these bands contained (as closely as possible) the same number of tokens (as suggested by Schmid, Verspoor & MacWhinney, 2011). An overview of these frequency bands is given in Table 3.

It is notable that the most frequent band in both corpora contains a number of topic-specific items. For the interview, these were words related to the migration experience, such as *Jahr* “year”, *sprechen* “talk”, *Deutsch* “German”, *Englisch* “English”; in the film retelling, among the most frequent words were nouns relating to important protagonists or items from the film, such as *Polizist* “policeman” (the ten-minute film sequence features six different police officers), *Frau* “woman”, *Mädchen* “girl”, *Brot* “bread” and *Haus* “house”.

For each speaker it was then calculated what proportion of the lemmatized lexical items that s/he had used fell into each of these frequency bands. Furthermore, it was determined what proportion of the items each speaker had used were among the 50 most frequent in the relevant corpus (as suggested by Paul Meara, p.c.).

### Cognates

In order to determine to what degree crosslinguistic similarities might have affected the attritional process, an assessment was made of how many of the lexical items used in the corpus shared a similar form across the contact languages (German–English for the Canadian group, German–Dutch for the Netherlands group). The two combined corpora (interview and film retelling) contained a total of 10,481 lemmatized lexical types, which were classified for similarities in lexical form<sup>3</sup>. Of these, 4,445

$r = .169/.511$ ; for the data from Switzerland  $r = .196/.508$ ; for all data combined  $r = .207/.537$ ). Since all participants originated in Germany, we based our further analyses on these data only.

<sup>3</sup> The classification was performed by a student assistant who is an early Dutch–German bilingual and a third-year student of English.

Table 4. *VF tasks – descriptive and inferential statistics*

	Group means			Anova					Post hoc	
	Control group	Attriters in	Attriters in the	Levene's F	Group Sig.	Planned contrasts			CG vs. CA	CG vs. NL
		Canada	Netherlands			F	p	$\eta^2$		
Total	25.76	21.77	21.58	.058	.944	13.776	.000	.151	<.001	<.001
VF1 (first ten seconds)	7.85	7.07	6.59	1.257	.287	9.014	.000	.104	<.05	<.001
VF2	4.99	4.24	4.10	.731	.483	6.591	.002	.078	<.01	<.01
VF3	4.03	3.25	3.14	.389	.678	6.947	.001	.082	<.05	<.01
VF4	3.37	2.88	2.78	.898	.409	2.805	.064	.035	.070	<.05
VF5	2.85	2.16	2.55	.091	.913	5.281	.006	.064	<.01	.210
VF6 (last ten seconds)	2.68	2.16	2.42	.094	.910	2.497	.086	.031	<.05	.230

(42.41%) were German–Dutch (but not German–English) cognates, 3,080 (29.39%) were German–Dutch–English cognates, 98 (0.94%) were German–English (but not German–Dutch) cognates, and in 2,858 (27.27%) cases there were no similarities between German and either of the contact languages. This means that, of the lexical items (types) used in the present corpus, a total of 7,525 (71.80%) were similar in form in Dutch and German, and 3,178 (30.32%) were similar in English and German. It was then assessed for all bilingual speakers what proportion of the speech they had produced consisted of items with cognates in their own L2.

## Results

### *Verbal fluency*

The total average score achieved in the two VF tasks and the average productivity of all participants in each ten-second segment are shown in Table 4. The descriptive statistics show that the German controls outperform the two bilingual populations in each segment of the task as well as in their total score. These differences were assessed by means of a multiple analysis of variance (MANOVA) and shown to be highly significant (Wilks' Lambda:  $F = 3.135$ ,  $p < .001$ ,  $\eta^2 = .111$ ). Tests of between-subject effects for the individual dependent variables show that effect size has a tendency to decrease as the task proceeds (see Table 4).

A further interesting effect is illustrated in Figure 1: all groups begin at a fairly high level of productivity, which drops sharply in the next segments and then levels off. The differences between attriters and controls seem particularly pronounced in segments 2 and 3. Towards the end of the task, the attriters in the Netherlands appear

to “catch up” with the controls, whereas the Germans in Canada (who had started at a higher level of productivity than the participants with Dutch L2) show a further drop.

### *Free speech*

#### *Lexical diversity*

Table 5 summarizes the descriptive group results for lexical diversity and sophistication based on the present corpus.

Multivariate analyses of variance by group were conducted for the interview and film retelling data; these analyses are summarized in Table 6. Since Jarvis (2013a, b) has pointed out that most measures of lexical diversity are affected to some extent by volume (i.e. text length), the number of tokens produced in each sample was included as a covariate. With all variables represented in Table 5 included, Box's M was highly significant, indicating a violation of the homogeneity of the covariance matrix. Removing Shannon's index from the analyses reduced Box's M to acceptable levels for both datasets (interview: Box's M 184.427,  $p = .026$ ; film retelling: Box's M = 174.482,  $p = .071$ ).<sup>4</sup> Group differences were significant for both datasets at  $p < .01$  (Wilks' Lambda  $F = 11.439$ , partial  $\eta^2 = .475$  for the interview data and  $F = 2.082$ , partial  $\eta^2 = .140$  for the film retelling data). The covariate, number of tokens, was highly significant at  $p < .001$  for both datasets, with a very high effect size (interview:  $F = 204.138$ , partial  $\eta^2 = .942$ ; film retelling:  $F = 208.357$ , partial  $\eta^2 = .942$ ).

The analysis of the individual dependent variables (see Table 6) shows that volume consistently affects the first set of variables designed to assess lexical diversity and distribution, with the exception of the MTLTD measure.

She rated those items as cognate which were predictable on the basis of regular phonological correspondences.

<sup>4</sup> MANOVA is robust in relation to homogeneity assumptions except when sample sizes are unequal and Box's M shows a significance of  $p < .001$  (see Tabachnick & Fidell, 1996, p. 382).



Table 5. *Lexical diversity and sophistication, based on present corpus (group means)*

	INT			CC		
	CG	CA	NL	CG	CA	NL
VOCD	94.95	98.91	95.34	77.69	72.65	71.84
MTLD	52.42	49.61	50.17	45.62	41.96	43.23
Shannon	5.01	5.23	5.13	4.62	4.60	4.58
Evenness	.84	.81	.82	.86	.84	.81
Rarity	1736.81	1803.60	1691.09	780.70	772.19	758.72
Dispersion	105.24	163.77	139.27	60.68	60.20	60.89
% Frequency band 1 (20% most frequent items)	17.53	20.66	21.83	19.18	20.62	20.36
% Frequency band 2	19.16	19.93	19.96	18.03	20.61	21.04
% Frequency band 3	18.11	20.63	19.76	20.35	19.35	20.00
% Frequency band 4	18.87	20.15	18.98	20.23	19.74	20.07
% Frequency band 5 (20% least frequent items)	21.62	18.63	16.73	22.22	19.69	18.53
% 50 most frequent items	39.33	43.95	45.10	41.50	45.65	46.41

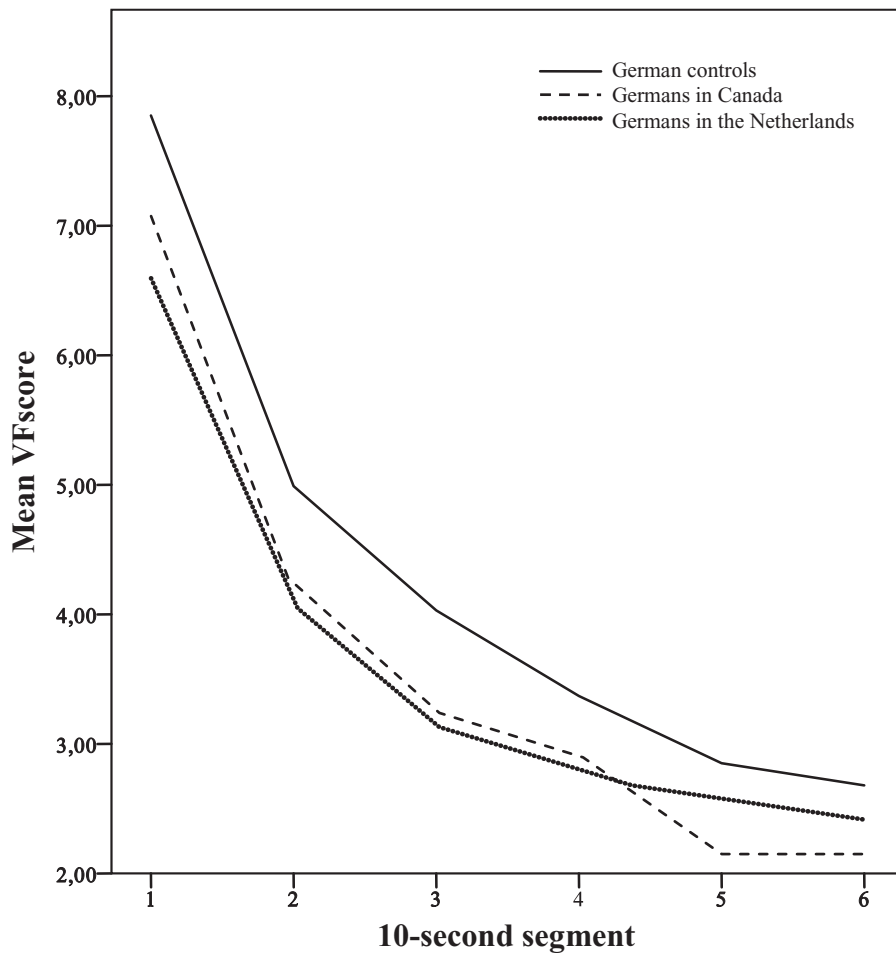


Figure 1. Average productivity in each ten-second segment of the two verbal fluency tasks across populations

Table 6. ANOVAs and planned contrasts for the lexical diversity and sophistication measures

	Levene's		Tokens			Group			CA vs. CG	NL vs. CG
	F	p	F	p	partial $\eta^2$	F	p	partial $\eta^2$		
INT										
VOCD	.367	.694	1.478	.226	.010	2.176	.117	.028	.05	.63
MTLD	.171	.843	.835	.362	.006	1.143	.322	.015	.19	.18
Evenness	5.920	.003	398.753	<.001	.728	7.214	<.01	.088	.06	<.001
Dispersion	2.342	.100	1146.428	<.001	.885	14.803	<.001	.166	<.001	<.001
Rarity	.370	.692	41.439	<.001	.218	3.876	<.05	.049	<.05	<.01
Freq.band 1	.144	.866	4.118	<.05	.027	21.632	<.001	.225	<.001	<.001
Freq.band 2	.307	.736	1.458	.229	.010	1.719	.183	.023	.09	.12
Freq.band 3	.806	.448	3.552	.061	.023	13.819	<.001	.156	<.001	<.001
Freq.band 4	.283	.754	.888	.347	.006	2.972	.054	.038	<.05	.99
Freq.band 5	.909	.405	3.847	.052	.025	14.615	<.001	.164	<.001	<.001
50 most freq.	.210	.811	5.484	<.05	.036	19.276	<.001	.206	<.001	<.001
CC										
VOCD	.148	.862	8.264	<.01	.052	1.975	.142	.025	.116	.070
MTLD	.143	.867	.894	.346	.006	2.453	.089	.031	<.05	.156
Evenness	.752	.473	129.943	<.001	.463	1.822	.165	.024	.068	.646
Dispersion	1.600	.205	990.405	<.001	.868	1.152	.319	.015	.159	.243
Rarity	.894	.411	38.017	<.001	.201	1.193	.306	.016	.484	.125
Freq.band 1	1.480	.231	2.187	.141	.014	1.481	.231	.019	.131	.148
Freq.band 2	1.632	.199	20.291	<.001	.118	12.029	<.001	.137	<.001	<.001
Freq.band 3	3.000	.053	.743	.390	.005	1.176	.311	.015	.159	.855
Freq.band 4	.698	.499	4.091	.045	.026	.458	.633	.006	.435	.388
Freq.band 5	.199	.820	5.655	.019	.036	6.549	<.01	.080	<.01	<.001
50 most freq.	2.922	.057	15.825	<.001	.095	11.420	<.001	.131	<.001	<.001

VOCD is only affected in the film retelling but not in the interview data; all other measures (evenness, dispersion and rarity) appear highly sensitive to volume. Where the frequency analyses are concerned, there are hardly any volume effects (with the exception of the highest frequency band in the interview data and the second highest in the film retelling).

Turning to the group comparison, there appear to be few differences across populations in the two overall measures of lexical diversity, VOCD and MTLD. The only significant finding is that MTLD is somewhat lower for the Canadians in the film retelling data than for the other populations. Evenness, dispersion and rarity are affected consistently in the interview data (with a marginal significance for evenness among the Canadian group and (highly) significant levels on all other measures) but not at all in the film retelling.

In the interview, frequency bands 1, 3 and 5 as well as the 50 most frequent items differ between populations, while for the film retelling data, only frequency bands 2 and 5 as well as the 50 most frequent items are significantly

different. Effect sizes for these group differences are, however, quite weak, ranging from  $\eta^2 = .080$  (frequency band 5 in film retelling) to  $\eta^2 = .162$  (frequency band 5 in interview), indicating that the group differences, while consistent, are hardly dramatic. Supplementary Figures 1a/b show that for the lexical frequency bands the group differences are indeed the effect of the attriters differing from the controls in the predicted direction (i.e. an overuse of the more frequent items and an underuse of the least frequent ones).

#### *COSMAS II frequencies*

Table 7 summarizes the results for the frequencies of items used by each speaker, based on the COSMAS II corpus. In these data, there do not seem to be straightforward tendencies among the attriters to overuse the more frequent or underuse the less frequent items, as was the case for the frequencies based only on the corpora at hand in the analyses above. The distribution across frequency bands is graphically represented in

Table 7. *Lexical diversity and sophistication, based on COSMAS II corpus*

		Group means			Levene's		Group			Post hoc	
		CG	CA	NL	F	p	F	p	partial $\eta^2$	CA vs. CG	NL vs. CG
INT	% FB 1	15.01	14.42	16.44	4.292	.015	5.751	<.01	.072		
	% FB 2	11.33	12.79	11.75	1.320	.270	8.540	<.001	.103	<.001	0.16
	% FB 3	11.14	12.48	11.90	.966	.383	3.113	<.05	.040	<.05	0.11
	% FB 4	13.92	14.89	13.55	.233	.793	6.467	<.01	.080	<.05	0.58
	% FB 5	43.89	45.41	43.64	2.446	.090	4.249	<.05	.054	<.05	0.98
CC	% FB 1	15.06	15.00	16.15							
	% FB 2	12.13	11.13	11.02							
	% FB 3	12.36	12.98	14.04							
	% FB 4	11.55	11.55	10.84							
	% FB 5	48.67	49.14	47.85							

FB1 = Frequency band 1 (20% most frequent items), FB2 = Frequency band 2, FB3 = Frequency band 3, FB4 = Frequency band 4, FB5 = Frequency band 5 (20% least frequent items)

Supplementary Figures 2a (interview data) and 2b (film retelling data).

Multivariate analyses of variance by group were conducted for the interview and film retelling data based on the COSMAS II frequencies summarized in Table 7. As in the analyses presented above, number of tokens was included as a covariate. The homogeneity of the covariance matrix was assessed by means of Box's M, which was not significant for either dataset (interview: Box's M 27.549,  $p = .038$ ; film retelling: Box's M 52.926,  $p = .188$ ). Group differences were significant for the interview data (Wilks' Lambda  $F = 15.035$ ;  $p < .001$ ; partial  $\eta^2 = .341$ ), but not for the film retelling data (Wilks' Lambda = 1.335;  $p = .086$ ). The covariate, number of tokens, was marginally significant for the interview data (Wilks' Lambda  $F = 2.559$ ,  $p < .05$ , partial  $\eta^2 = .081$ ) but not for the film retelling data (Wilks' Lambda  $F = 1.335$ ,  $p = .245$ ).

The analysis of the individual dependent variables in the interview data shows significant group differences for all frequency bands. Planned comparisons show that the attriters in Canada differ from the controls on frequency bands 2–5, showing a higher use of the high-frequency items in band 2 than the controls, but also somewhat higher proportions of the low-frequency items in bands 4 and 5. The attriters in the Netherlands are only different from the controls on the extremely high-frequency band 1, which they use more frequently. The full tests are summarized in Table 7.

### Cognates

For all speech samples an assessment was made of the proportion of the lemmatized lexical items used that consisted of German–English (GE–EN) and of German–Dutch (GE–NL) cognates. The use of GE–EN cognates

by the controls and the Germans in Canada looks very similar: the controls use 48.97% in the interview and 44.93% in the film retelling, while the Canadian attriters use 48.01 and 46.30%, respectively. For the GE–NL cognates, a slightly higher use of cognates by the Dutch–German bilinguals can be observed; they account for 83.25% in the interview and 82.19% in the film retelling for the controls, and for 86.31% in the interview and 84.66% in the film retelling for the attriters in the Netherlands.

A multivariate analysis of variance by group was conducted for the use of GE–EN and GE–NL cognates in interview and film retelling. The homogeneity of the covariance matrix was assessed by means of Box's M, which was not significant (Box's M = 25.460,  $p = .223$ ). Group differences were significant (Wilks' Lambda  $F = 4.602$ ;  $p < .001$ ; partial  $\eta^2 = .142$ ).

Planned contrasts established no overuse of GE–EN cognates by the Germans in Canada, but significant overuse ( $p < .01$ ) of GE–NL cognates by the Germans in the Netherlands. (The fact that the use of GE–EN cognates in the film retelling data was marginally significant for this group as well might be due to the fact that there is a substantial overlap between GE–EN and GE–NL cognates; see above.) Effect sizes were small, with a maximum  $\eta^2$  of .172. Table 8 gives the details of these analyses, while Supplementary Figures 3a/b show the group differences in the use of cognates.

### Summary of group comparisons

The analyses presented so far have looked at lexical productivity, lexical access and lexical diversity in two attriting populations through a range of tasks and measures. Where the controlled VF task was

Table 8. *Proportion of cognates: ANOVAs and planned contrasts*

	Levene's		Group			CA vs. CG	NL vs. CG
	F	p	F	p	partial $\eta^2$		
German–English cognates in INT	3.741	.026	.439	.646	.006	.765	.557
German–English cognates in CC	.175	.840	3.326	.039	.044	.149	<.05
German–Dutch cognates in INT	2.292	.105	15.200	<.001	.172	.050	<.001
German–Dutch cognates in CC	.562	.571	5.279	.01	.067	.183	<.01

Table 9. *Overview of results*

	Germans in Canada		Germans in the Netherlands	
	INT	CC	INT	CC
VFTot	***		***	
VF1	*		***	
VF2	**		**	
VF3	*		**	
VF4			*	
VF5	**			
VF6	*			
VOCD				
MTLD		*		
Evenness			***	
Dispersion	***		***	
Rarity	*		**	
Freq.band 1	***		***	
Freq.band 2		***		***
Freq.band 3	***		***	
Freq.band 4	*			
Freq.band 5	***	**	***	***
50 most freq.	***	***	***	***
Cognates			***	**

\*: p < .05; \*\*: p < .01; \*\*\*: p < .001

concerned, the analyses found robust differences between populations, in particular for the early segments of the task, indicating that the attriters were less productive in naming items belonging to the two categories and took longer to “get into” the task. For the two free speech samples, it is interesting to see that the measures involving ratios of types and tokens (VOCD and MTLD) show no substantial difference between populations. Evenness, dispersion and rarity do differ, but only for the interview, not for the film retelling data. These differences may thus to some extent be related to text length – even though this was included as a covariate in the analyses – or to the fact that the film retelling refers to a predetermined sequence of events and items and is thus more constrained. In both

corpora, the attriters overuse the most and underuse the least frequent vocabulary when frequency is based on the present corpora alone, but not when it is based on a large corpus of German data of all types. Cognates are overused by L2 speakers of a closely related language (Dutch) that has a similar lexicon to the L1, but not by L2 speakers of English, whose lexicon differs more substantially from German despite the historical and typological relationship between these two languages. The differences between the L2 Dutch and L2 English groups are summarized in Table 9.

These findings suggest a complex pattern of how the lexicon is affected in the attritional process. The following section will explore how and to what extent the different

measures investigated here can best be combined in order to predict whether a particular speaker is an attriter or not.

### **Discriminant analysis**

Linear discriminant analysis (DA) is similar in many respects to MANOVA, in other respects to regression analysis, and in yet others to factor analysis. As described by Huberty and Olejnik (2006), DA and MANOVA rely on similar statistical techniques but are in effect opposites: the grouping variable in a MANOVA becomes the dependent variable in a DA. The purpose of DA is to use outcome variables as predictors of the group membership of individual cases by constructing a model that best predicts the dependent variable. In this sense, DA is similar to multiple logistic regression, which also involves the use of a categorical (or nominal) dependent variable. Finally, DA is also similar to factor analysis in that both types of analysis identify clusters of variables along multiple dimensions.

The DA in the present study was run with SPSS 19.0. The participants' group membership was used as the dependent variable. The independent variables included the following measures, totaling 54 items:

- VF task: the total average number of elements named in both tasks per participant, and the total average number of elements named in each of the six ten-second segments in both tasks (total of seven items)
- interview corpus: VOCD, MTL, effective types, Shannon, rarity, evenness, dispersion, frequency band 1–5 based on present corpus, proportion of 50 most frequent items based on present corpus, average frequency of lexical items based on present corpus, frequency band 1–5 based on COSMAS II corpus, proportion of 50 most frequent items based on COSMAS II corpus, average frequency of lexical items based on COSMAS II corpus (21 items)
- film retelling corpus: VOCD, MTL, effective types, Shannon, rarity, evenness, dispersion, frequency band 1–5 based on present corpus, proportion of 50 most frequent items based on present corpus, average frequency of lexical

items based on present corpus, frequency band 1–5 based on COSMAS II corpus, proportion of 50 most frequent items based on COSMAS II corpus, average frequency of lexical items based on COSMAS II corpus (21 items)

Fluency in film retelling: The proportion of disfluency markers (filled pauses (FP), empty pauses (EP), repetitions (REP) and self-corrections or retractions (RETR)) were described and analyzed for the film retelling data by Schmid and Beers Fägersten (2010). These measures were included in the present analysis, as well as the number of words spoken per minute (excluding disfluency phenomena) in these retellings (five items).

The results from the cognate analysis were not included in this analysis, since they are presumed to affect the three populations differentially (it is of no consequence to an English L2 speaker whether a particular German word has a cognate equivalent in Dutch, or vice versa, nor does cognate status matter for the monolinguals).

The DA method was set to stepwise, which meant that it chose only one variable at a time in accordance with how much that variable contributed to the strength of the model. The criteria used for variable entry and removal were the default Wilks' Lambda F values of 3.84 for entry and 2.71 for removal. This ensured that a variable would be added to the model only if it contributed significantly to the strength of the model, and that if it no longer made a significant contribution to the developing model it would be removed as new variables were added. Results were cross-validated using leave-one-out cross-validation, a procedure that iteratively builds a model using all but one of the cases, and then tests the model's predictive accuracy blindly (without access to its group membership) on the case that was left out during that iteration. The number of iterations during the leave-one-out cross-validation phase is equal to the number of cases, which means that each case is used once as the test case. The results of the cross-validation phase show the number of cases whose group membership was predicted correctly.

Due to missing variables for some of the participants, the DA omitted 13 participants from the analysis. The stepwise procedure selected 11 predictor variables (these are listed in Supplementary Table 4). The overall cross-validated classification accuracy of the 11-variable

Table 10. *Cross-validated classification results*

	Group	Predicted Group Membership			Total
		CG	CA	NL	
Count	CG	49	0	3	52
	CA	0	42	8	50
	NL	2	9	34	45
%	CG	94.2	0	5.8	100.0
	CA	0	84.0	16.0	100.0
	NL	4.4	20.0	75.6	100.0

model is 85.0%, meaning that the group memberships of 85.0% of the participants (125 out of 147) were predicted correctly based on a model consisting of just the 11 variables selected (see Table 10 for the numbers and percentages of participants in each group whose membership was predicted accurately during the cross-validation phase). This is significantly above the chance level of 33% ( $df = 4, n = 147, X^2 = 181.885, p < .001$ ). As shown in the classification matrix in Table 10, the model is best at identifying the German speakers still living in Germany (94.2% classification accuracy), followed by the attriters living in Canada (84.0%). It is least effective at identifying the attriters living in the Netherlands (75.6%), some of whom are misclassified as attriters living in Canada (20.0%) and two as German speakers still living in Germany (4.4%).

As already mentioned, DA explores relationships between variables along different dimensions – referred to as functions – and the number of functions is always one less than the number of groups. The results of the DA showed that Function 1 accounts for 90.1% of the variance in the data, and that it is most characterized by qualities related to lexical diversity in the interview task. Of the 11 variables included in the classification model, four are primarily affiliated with Function 1. These include (a) interview evenness, (b) interview effective types, (c) interview percentage of items in frequency band 3, and (d) interview MTLTD. The seven other variables are more closely aligned with Function 2, which accounts for 9.9% of the variance, and is characterized primarily by qualities related to lexical sophistication in the interview task.

DA assigns specific weights (or canonical discriminant function coefficients) to each variable, and uses these to calculate overall function scores for each participant. The group means for the three groups in relation to both functions are plotted in Supplementary Figure 4, which shows the overall separation between the groups in relation to the 11 predictor variables in the model. These results suggest that the participants' patterns of lexical diversity in the interview task (i.e., Function 1)

separate the German speakers living in Germany from the attriters, whereas the participants' patterns of lexical sophistication in the same task (Function 2) separate the attriters in the Netherlands from the other two groups.<sup>5</sup>

It is interesting that all but one of the predictive variables in the DA model relate to lexical qualities in the interview task. Only one predictor is associated with the film retelling task, and this variable is more a matter of fluency (i.e., words per minute) than of word choice. The variables pertaining to the VF task, one of the instruments most often used in attrition studies (and on which the populations differed significantly in the individual analyses presented above), do not appear in the final model at all.

### *The impact of extralinguistic variables*

The analyses presented above reveal a number of differences between the attriting and the control populations. The discriminant analysis provided a model which included those factors that collectively have the most predictive power in assigning group membership to each individual participant. For each speaker, a new variable was calculated, representing his or her cumulative score on the two functions that emerged from the DA.

In order to assess to what extent external factors such as the biological age of the speaker, the length of residence in the country of migration, as well as frequency of use and attitudes towards the L1 and L2 might have had an impact on these scores (as was asked in RQ3), multiple linear regressions were conducted on the data from the bilingual speakers. The outcome variables for these analyses were those measures for each of the three tasks for which effect size had been strongest in the MANOVAs:

<sup>5</sup> This interpretation is supported by the participants' aggregated function scores rather than by a comparison of the group means for individual variables.

- the total number of items on the VFT
- dispersion in the interview
- the proportion of items in frequency band 1 in the interview
- the proportion of items in frequency band 2 in the film retelling.

In addition, the two functions that were calculated by the DA were used as outcome variables. The predictors entered into the models were the extralinguistic factors discussed above, namely:

- age at emigration
- length of residence
- overall L1 use
- affiliation with the L1 language and culture
- L1 use for professional purposes.

One further variable was included in these models, based on the prediction made by Segalowitz (1991) that a higher L2 proficiency may lead to reduced automaticity in the L1. The original test battery contained a measure of overall L2 proficiency, namely a C-Test. The C-Test is a version of the cloze test which is frequently used in research on both L2 acquisition and L1 attrition. It consists of short texts in which parts of words are deleted according to a predetermined schema (see Schmid, 2011, for an in-depth discussion): each correctly completed word is awarded one point. Our C-Test used five short texts with a maximum possible score of 100. Although the two attriting populations completed two different versions of this test (English for the Germans in Canada and Dutch for the Germans in the Netherlands), performance on this task was exactly the same: both groups achieved a mean score of 73.23 (English group: st.dev 16.31, range 18–96; Dutch group: st.dev. 15.57, range 19–99).

Of the six models arrived at by these linear regression analyses, only VFT emerged as significantly impacted by the predictors (for the full details of the analyses, see Supplementary Table 5), with a weak impact of L1 use for professional purposes and L2 proficiency (speakers who used German in their workplace and who achieved higher scores on the L2 C-Test were somewhat more productive on the VFT). Age at testing was removed from all analyses as having no predictive power overall, and for all the other analyses, none of the predictors reached significance. This implies that the variance present on those measures that can best discriminate the attriting populations from the monolingual controls is not affected by factors pertaining to age of emigration, length of residence, amount of L1 use, attitudes, or proficiency in L2. This finding may be surprising and appear counterintuitive, but it does underscore the results from a number of earlier studies, all of which failed to establish a link between such extralinguistic factors and language attrition (e.g. Schmid

& Dusseldorp, 2010; Varga, 2012; Yilmaz & Schmid, 2012).

## Discussion

The aim of the present study was to provide an analysis of first language attrition in the domain of the lexicon which would go beyond what previous studies have achieved in both breadth and depth. In order to achieve this, a set of data was analyzed that comprised both controlled, formal tasks probing lexical access (verbal fluency tasks) and spontaneous speech elicited by means of a film retelling and a semi-structured interview. Our first set of research questions addressed the possible loss of lexical accessibility in a situation where the speaker has little opportunity to use her L1 and relies mainly on the L2 in everyday interactions. In particular, we were interested in the impact of activation vs. inhibition on the attritional process.

Next, we aimed to establish what types of measures might be most suited to detect L1 attrition in free speech. We augmented type–token-based measures with an analysis of measures that take into account the rarity and distribution of items across the discourse, as well as assessing the proportion of items of varying lexical frequency. We also assumed that in free speech, the distribution of items might be more uneven for attriters than for controls, due to the limited accessibility of the entire range and a more local spread of activation across semantic fields.

A last set of research questions related to the impact of individual external background factors such as age at emigration and L1 exposure, as well as the degree of lexical difference between languages (English vs. Dutch in comparison to German), on the degree of attrition.

### *Verbal fluency task*

The statistical comparisons showed that the attriters were indeed less productive than the controls on a verbal fluency task, indicating that they were able to access fewer items within a given lexical category and specified time limit. With an overall effect size of  $\eta^2 = .151$ , however, these findings indicate only a relatively weak difference in accessibility. Whether they are due to the fact that bilinguals have a much larger repertoire and also have to suppress L2 items in order to be able to produce items in the L1, or whether they can actually be considered effects of an attritional process that has made lexical access more difficult and effortful, is difficult to establish based on the present data. The fact that bilinguals become slower at naming objects in their first language very shortly after the onset of bilingualism (so probably not because of attrition) has been demonstrated by, for example, Mägiste (1979) and Linck, Kroll and Sundermann (2009). Future

work might address this question by investigating fluent bilinguals in the country of origin.

It is interesting to observe, however, that the attriters in Canada start out closer to the native norm than the attriters in the Netherlands. During the last third of the task, the effect reverses and it is the Dutch L2 speakers who catch up with the native norm, while the Canadians drop back. It is possible that this crossover effect is due to the differential amounts of competition and consequently inhibition: for the L2 English speakers, competition from the L2 is less strong than for the attriters in the Netherlands, whose L2 shares more cognates with their L1 (as was shown above, less than one third of English items but more than two thirds of Dutch items have cognates in German). Presumably, speakers begin by naming those items that are most accessible in their minds, so that in the early stages of the task, what slows bilinguals down will be predominantly the added cognitive effort of inhibiting competitors. This will disadvantage speakers whose languages share a large part of their vocabulary. Later on, when the speaker has to “dig deeper”, it may become easier to activate less easily accessible items if they are supported by similar forms in both languages.

The verbal fluency task is a very popular task in language attrition studies, since it is, for example, easy to construct, administer and score, does not rely on specialized equipment (unlike picture naming and other reaction time tasks) and needs no stimuli. However, despite the robust group differences detected by the group comparisons, none of the variables pertaining to this task were selected in the DA. In other words, the individual results from this task had very limited predictive power when it came to classifying an individual participant as an attriter or a control speaker. This means that the use of the verbal fluency task to detect lexical attrition may need to be reassessed. As was noted above, a range of studies have administered this task alongside other elicitation techniques, such as narratives and interviews. A metastudy of these investigations, using the methods described here, may be a valuable step forward in this context. In summary, the answer to RQ1a is that attriters are indeed outperformed by controls on the VFT, but that effect sizes are weak and performance on this task does not contribute much to the profile of an attriter vs. a control.

### **Free speech**

Where free speech was concerned, we first observed that type–token frequency based measures failed to distinguish between attriters and controls. There were some differences related to the distribution of items across the interview, indicating that for some of the attriters there were longer intervals between items of the same type than in the data from the controls. This may be related

to the nature of the interview data, where attriters may have more often gone off on “tangents” related to the emigration experience and also produced longer speech samples overall. No such distributional effects were found in the film retelling data, probably because the sequence of events was predetermined by the stimulus, so lexical selection was more constrained. In other words, in answer to RQ1b, it is likely that such distributional differences between populations are the outcome of the type (and length) of data produced and not an indication of changes in lexical access due to language attrition.

Word frequency was assessed on the basis of two corpora: one that consisted only of the data collected in the present study (that is, only texts of a very similar nature) and one that measured lexical frequency in a large corpus of texts comprising 5.4 billion words from a wide range of written and spoken sources (the COSMAS II corpus collected by the *Institut für Deutsche Sprache*). In order to assess the use of high-, medium- and low-frequency items, we divided all lexical lemmas into five frequency bands, each representing 20% of all tokens. The group comparisons showed an overuse of the high-frequency and an underuse of the low-frequency items among the attriters. It was shown that for the interview data, overuse affected the items of the very highest frequency, while in the film retelling, it was the second highest frequency band that was overused. Again, this is probably due to the nature of the stimuli: as was pointed out above, the highest frequency band in the film retelling data contained a large number of items that figured prominently in the film sequence and were therefore necessary for all speakers.

We also found an overuse in both attriting populations of the mid-frequency band 3 in the interview data. This is probably related to the fact that very low-frequency items are dispreferred by these speakers, who then come to rely more on the mid-range. However, it should be pointed out again that effect sizes here were very modest indeed. It is certainly not the case that attriters stop using low-frequency items overall: every single speaker used at least a certain number of words from this category, as is illustrated in the histograms in Supplementary Figures 5a/b.

The interpretation of the frequencies that were derived from the COSMAS II corpus was far less straightforward. This indicates that it may be difficult to rely on overall lexical frequencies, and preferable for reference purposes to use a corpus that is very similar in type to the data under investigation. This finding is helpful for researchers working on less well documented languages, for which large corpora may not be available, as it suggests that the corpus at hand may yield the most reliable results (provided that it is large enough). In answer to RQ1c, the analyses presented here have established that attriters do indeed develop a preference for somewhat more common and more high-frequency lexical items within the context



specified by a certain task, but that this preference does not appear to be reflected when lexical frequency is based on overall language use from a wide range of text types and tasks.

RQ1d focused on the important but complex issue of whether attrition is determined more by a reduction in lexical accessibility, due to infrequent use of the L1, or to problems related to inhibiting the more frequently used L2. It was predicted that speakers of an L2 that was very closely related to the L1 would have fewer problems related to activation due to underuse, since activation levels degrade over time, but would be maintained because of the use of similar items in the L2. Such speakers would, on the other hand, suffer more from problems related to inhibition than bilinguals with a less similar L2, since the closeness of items in L1 and L2 would make inhibiting the competing items more costly. Two findings stand out in this respect:

1. The attriters in the Netherlands differ significantly from the controls with respect to evenness in the interview data, but there is no difference between the Canadians and the controls. This effect may be linked to some extent to the similarity between the two languages, which requires the Dutch L2 speakers to inhibit their L1 more strongly in their daily lives when interacting in their L2. Potentially, this inhibition effort makes the L1 less available, so that those items whose activation is boosted by semantic similarity to the topic at hand are preferred.
2. The Dutch attriters overuse cognates while the Canadian attriters do not. This finding may again relate to the proportion of the lexicon that is shared between a speaker's L1 and L2, and the assumption proposed by, for example, de Groot, Dannenburg and van Hell (1994) that cognates share the same semantic representation (recall that more than two thirds of the lexical items used in the present corpus were cognate between Dutch and German, but less than a third of the English–German pairs). L2 speakers of English have more reason to distrust forms that look the same, and may therefore opt for a different word if they have the choice.

The lack of clearer results in this respect may to some extent be due to the fact that, although English is less closely related to German, it still belongs to the same language family. Further studies invoking more distant languages may be able to shed more light on this issue.

### *Diagnosing attrition*

We then attempted to combine the measures that were derived from the data collected for the present study in order to determine which combination of factors would

produce the most reliable model that would discriminate between attriters and non-attriters. We entered the measures from the verbal fluency task and those for lexical diversity and sophistication, plus a number of fluency measures from the film retelling task into the analysis, which yielded a model comprising eleven factors. What was most immediately striking was that both the formal task and the film retelling task did not contribute to the final model (with the exception of the speech rate from the film retelling task).

The model revealed that only four of the 54 measures included in the analysis together accounted for more than 90% of the overall variance. These measures were mainly to do with the diversity and distribution of types across the interview data, and also included the use of items in the mid-range frequency band (which, as mentioned above, was overused by the attriters). A second set of measures, accounting for a further 9% of the variance, distinguished the Dutch attriters from the other two groups and mainly contained measures relating to word frequency in the interview data, in particular to the high- and low-frequency lexical items.

Where the use of different tasks and various measures of lexical diversity in language attrition research are concerned (RQs 1 and 2), the findings from the DA suggest that the predominantly formal tasks and even elicited data that have provided the data for most previous investigations of attrition are less powerful and yield less valid comparisons of attriters and controls than true spontaneous speech. Similarly, the relatively easily derived measures that are based on type–token frequencies appear to give less reliable insights into the attritional process than measures that assess the frequency of lemmas, either in the corpus at hand or (where available) in larger corpora of the language. It also appears important not only to look at the presence/absence of items in the data but also at their distribution (evenness) across the entire discourse, which also may have changed in the attritional process.

### *The impact of extralinguistic variables*

While the group comparisons showed robustly significant differences between attriters and controls and the model yielded by the DA was quite reliable in assigning individuals to the different populations (the classification was accurate in 85% of the cases), the analyses did leave us with the question of why some of the attriters in the individual populations showed stronger attrition effects than others. To gain insight into the questions formulated under RQ3, we investigated the impact of external factors such as the length of time spent in an L2 environment, the use of the L1 in different situations and settings, and attitudes towards the native language and culture, as well as L2 proficiency on these functions. To do this, we

conducted linear regression analyses, with these factors as predictors, on the outcome variables that had most strongly differentiated attriters and controls (the number of items produced on the VFT, the use of high-frequency items in interview and film retelling, dispersion in the interview and the two functions that had been calculated by the DA). The only significant impact revealed by these analyses, however, concerned the formal task, in which speakers who used the L1 in their workplace and who were more proficient in the L2 were shown to be more productive. Other than that, these analyses produced no significant results and left us no nearer to answering the puzzling question of why some individuals may be more susceptible to attrition than others. Based on the data investigated here, length of residence, frequency of L1 use, attitudes and L2 proficiency (as addressed by RQs 3a and 3b) do not appear to have an impact on individual variation in lexical access.

### Conclusion

This study has presented an in-depth examination of lexical L1 attrition in relation to the likelihood that attriters will experience decreased lexical accessibility as a consequence of their reduced use of the L1, and as a result of the limited number of contexts in which they use it. We investigated lexical diversity and distribution, lexical sophistication and verbal fluency, and looked for any generalizable and predictable differences between Germans who have remained in Germany, Germans living in Canada and Germans living in the Netherlands.

Lexical diversity in this study was examined as a matter of multidimensional compositional complexity that extends beyond the relationship between types and tokens to other dimensions of word choice, including the degree to which tokens are distributed evenly across different types, the distances to which repetitions of the same type are dispersed, and the relative rarity of the words that are used in a sample of speech. The last of these also overlaps with the construct of lexical sophistication, which we examined carefully in relation to the percentage of words in each participant's speech samples that can be found in each of the lexical frequency bands in the present data as well as in the COSMAS II native German corpus.

Overall, we found that the attriters do indeed differ significantly from each other and from the German controls on a number of measures relating to lexical diversity and distribution, lexical sophistication and verbal fluency. In most cases, though, the differences in lexical diversity are quite subtle and cannot be detected with measures that are based only on type and token frequencies. The significance levels and effect sizes also differ by task, such that significant effects for lexical diversity are found only in the interview. Importantly,

this is the task where the largest differences were found between groups in terms of number of words produced, so there is a strong possibility that the lexical diversity differences between groups in this task were to some extent affected by sample size. Although we used sample size as a covariate, we believe that future studies will benefit by using these measures of lexical diversity applied to equal-sized subsamples in order to further reduce the potential effects of sample size.

Regarding rarity and lexical sophistication, we recognize that the frequency rank of a word in a corpus is not a reliable measure of its perceived sophistication or of whether it has been used appropriately in its context. The findings reported here do nevertheless reveal that the three groups of participants did not use all of the same words, and that there are significant and predictable differences between the words they used and the patterns with which they used them. Future work that explores these qualitative differences could substantially further our understanding of the nature of lexical L1 attrition and the mechanisms, such as lexical accessibility, through which it occurs.

Lastly, the finding that language use in the interview was the most strongly predictive task in classifying a speaker as an attriter or a non-attriter (all but one of the eleven factors retained in the Discriminant Analysis originated from this task) underscores the fact that attrition affects the skill that is most characteristic of what native speakers know how to do: use language in free speech. This finding is important in light of the fact that most attrition studies, in particular those assessing lexical attrition, typically use controlled tasks. To what extent such tasks actually assess attrition, or whether they may be detecting a weakening of metalinguistic skills or knowledge instead, remains to be seen.

### References

- Ammerlaan, T. (1996). You get a bit wobbly. Exploring bilingual lexical retrieval processes in the context of first language attrition. Ph.D. dissertation, Katholieke Universiteit Nijmegen.
- Andersen, R. W. (1982). Determining the linguistic attributes of language attrition. In R. D. Lambert & B. F. Freed (eds.), *The loss of language skills*, pp. 83–118. Rowley, MA: Newbury House.
- Berthele, R. (2011). On abduction in receptive multilingualism. Evidence from cognate guessing tasks. *Applied Linguistics Review*, 2, 191–220.
- Bialystok, E. (2005). Consequences of bilingualism for cognitive development. In J. F. Kroll & A. M. B. de Groot (eds.), *Handbook of bilingualism*, pp. 417–432. Oxford: Oxford University Press.
- Chao, A., & Jost, L. (in press). *Diversity analysis*. London: Chapman & Hall/CRC Press.

- Cherciov, M. (2011). Between attrition and acquisition: The case of Romanian in Immigrant contexts. Ph.D. dissertation, University of Toronto.
- de Groot, A. M. B., Dannenburg, L., & van Hell, J. G. (1994). Forward and backward word translation by bilinguals. *Journal of Memory and Language*, 33, 600–629.
- de Leeuw, E., Schmid, M. S., & Mennen, I. (2010). Perception of foreign accent in native speech. *Bilingualism: Language and Cognition*, 13, 33–40.
- Dijkstra, T. (2005). Bilingual visual word recognition and lexical access. In J. F. Kroll & A. M. B. de Groot (eds.), *Handbook of bilingualism: Psycholinguistic approaches*, pp. 179–201. Oxford: Oxford University Press.
- Dostert, S. (2009). Multilingualism, L1 attrition and the concept of 'native speaker'. Ph.D. dissertation, Heinrich-Heine Universität Düsseldorf.
- Green, D. W. (1986). Control, activation and resource: A framework and a model for the control of speech in bilinguals. *Brain and Language*, 27, 210–223.
- Green, D. W. (1998). Mental control of the bilingual lexico-semantic system. *Bilingualism: Language and Cognition*, 1, 67–81.
- Huberty, C. J., & Olejnik, S. (2006). *Applied MANOVA and Discriminant Analysis*. Hoboken, NJ: Wiley.
- Hulsen, M. (2000). Language loss and language processing: Three generations of Dutch migrants in New Zealand. Ph.D. dissertation. University of Nijmegen.
- Jarvis, S. (2002). Short texts, best-fitting curves and new measures of lexical diversity. *Language Testing*, 19, 57–84.
- Jarvis, S. (2009). Lexical transfer. In A. Pavlenko (ed.), *The bilingual mental lexicon*, pp. 99–124. Bristol, UK: Multilingual Matters.
- Jarvis, S. (2012). Lexical challenges in the intersection of applied linguistics and ANLP. In C. Boonthum-Denecke, P. M. McCarthy & T. Lamkin (eds.), *Cross-disciplinary advances in applied natural language processing: Issues and approaches*, pp. 50–72. Hershey, PA: IGI Global.
- Jarvis, S. (2013a). Capturing the diversity in lexical diversity. *Language Learning*, 63 (supplement 1), 87–106.
- Jarvis, S. (2013b). Defining and measuring lexical diversity. In S. Jarvis & M. Daller (eds.), pp. 13–44.
- Jarvis, S., & Daller, M. (eds.) (2013). *Vocabulary knowledge: Human ratings and automated measures*. Amsterdam: Benjamins.
- Keijzer, M. (2007). First language attrition: A crosslinguistic investigation of Jakobson's regression hypothesis. Ph.D. dissertation, Vrije Universiteit Amsterdam.
- Köpke, B., & Nespoulous, J.-L. (2001). First language attrition in production skills and metalinguistic abilities in German–English and German–French bilinguals. In T. Ammerlaan, M. Hulsen, H. Strating & K. Yağmur (eds.), *Sociolinguistic and psycholinguistic perspectives on maintenance and loss of minority languages*, pp. 221–234. Münster: Waxmann.
- Köpke, B., & Schmid, M. S. (2004). First language attrition: The next phase. In M. S. Schmid, B. Köpke, M. Keijzer & L. Weilemar (eds.), pp. 1–43.
- Köpke, B., Schmid, M., Keijzer, M., & Dostert, S. (eds.) (2007). *Language attrition: Theoretical perspectives*. Amsterdam: Benjamins.
- Linck, J. A., Kroll, J. F., & Sunderman, G. (2009). Losing access to the native language while immersed in a second language: Evidence for the role of inhibition in second language learning. *Psychological Science*, 20, 1507–1515.
- Mägiste, E. (1979). The competing language systems of the multilingual: A developmental study of decoding and encoding processes. *Journal of Verbal Learning and Verbal Behavior*, 18, 79–89.
- McCarthy, P. M., & Jarvis, S. (2007). VOCD: A theoretical and empirical evaluation. *Language Testing*, 24, 459–488.
- McCarthy, P. M., & Jarvis, S. (2013). From intrinsic to extrinsic issues of lexical diversity assessment: An ecological validation study. In S. Jarvis & M. Daller (eds.), pp. 45–78.
- Montrul, S. (2008). *Incomplete acquisition in bilingualism. Re-examining the age factor*. Amsterdam/Philadelphia: Benjamins.
- Opitz, C. (2011). First language attrition and second language acquisition in a second-language environment. Ph.D. dissertation, Trinity College Dublin.
- Paradis, M. (1993). Linguistic, psycholinguistic, and neurolinguistic aspects of interference in bilingual speakers: The activation threshold hypothesis. *International Journal of Psycholinguistics*, 9, 133–145.
- Paradis, M. (2004). *A neurolinguistic theory of bilingualism*. Amsterdam: Benjamins.
- Paradis, M. (2007). L1 attrition features predicted by a neurolinguistic theory of bilingualism. In B. Köpke, M. S. Schmid, M. Keijzer & S. Dostert (eds.), pp. 121–134.
- Perdue, C. (1993). *Adult language acquisition: Cross-linguistic perspectives*. Cambridge: Cambridge University Press.
- Pielou, E. C. (1969). *An introduction to mathematical ecology*. Hoboken, NJ: Wiley.
- Schmid, M. S. (2004). First language attrition: The methodology revised. *International Journal of Bilingualism*, 8, 239–255.
- Schmid, M. S. (2007). The role of L1 use for L1 attrition. In M. S. Schmid, B. Köpke, M. Keijzer & L. Weilemar (eds.), pp. 135–153.
- Schmid, M. S. (2011). *Language attrition*. Cambridge: Cambridge University Press.
- Schmid, M. S., & Beers Fägersten, K. (2010). Fluency and language attrition. *Language Learning* 60, 753–791.
- Schmid, M. S., & Dusseldorp, E. (2010). Quantitative analyses in a multivariate study of language attrition. *Second Language Research* 26, 125–160.
- Schmid, M. S., & Köpke, B. (2007). Bilingualism and attrition. In B. Köpke, M. S. Schmid, M. Keijzer & S. Dostert (eds.), pp. 1–8.
- Schmid, M. S., Köpke, B., Keijzer, M., & Weilemar, L. (2004). First language attrition: Interdisciplinary perspectives on methodological issues. Amsterdam: Benjamins.
- Schmid, M. S., Verspoor, M., & MacWhinney, B. (2011). Coding and extracting data. In M. Verspoor, W. Lowie & K. de Bot (eds.), *A dynamic approach to second language development: Methods and techniques*, pp. 39–54. Amsterdam/Philadelphia: Benjamins.
- Segalowitz, N. (1991). Does advanced skill in a second language reduce automaticity in the first language? *Language Learning*, 41, 59–83.

- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27, 379–423, 623–656.
- Smith, B., & Wilson, J. B. (1996). A consumer's guide to evenness indices. *Oikos*, 76, 70–82.
- Sorace, A. (2005). Selective optionality in language development. In L. Cornips & K. P. Corrigan (eds.), *Syntax and variation. Reconciling the biological and the social*, pp. 55–80. Amsterdam: Benjamins.
- Tabachnick, B. G., & Fidell, L. S. (1996). *Using multivariate statistics* (3rd edn). New York: Harper Collins.
- Tsimpli, I. (2007). First language attrition from a minimalist perspective: Interface vulnerability and processing effects. In B. Köpcke, M. S. Schmid, M. Keijzer & S. Dostert (eds.), pp. 83–98.
- Varga, Z. (2012). First language attrition and maintenance among Hungarian speakers in Denmark. Ph.D. dissertation, Århus University.
- Waas, M. (1996). *Language attrition downunder*. Frankfurt: Peter Lang.
- Yağmur, K. (1997). *First language attrition among Turkish speakers in Sydney*. Tilburg: Tilburg University Press.
- Yılmaz, G., & Schmid, M. S. (2012). L1 lexical accessibility among Turkish–Dutch bilinguals. *The Mental Lexicon*, 7, 249–274.