

## Comparing MOSAIC and the Variational Learning Model of the Optional Infinitive Stage in Early Child Language

Daniel Freudenthal ([D.Freudenthal@Liv.Ac.Uk](mailto:D.Freudenthal@Liv.Ac.Uk))

Julian M. Pine ([Julian.Pine@Liv.Ac.UK](mailto:Julian.Pine@Liv.Ac.UK))

School of Psychology, University of Liverpool

Fernand Gobet ([Fernand.Gobet@Brunel.Ac.Uk](mailto:Fernand.Gobet@Brunel.Ac.Uk))

School of Social Sciences, Brunel University

### Abstract

This paper compares MOSAIC and the Variational Learning Model (VLM) in terms of their ability to explain the level of finiteness marking in early child Dutch, English, Spanish, German and French. It is shown that both models are successful in explaining cross-linguistic variation in rates of Optional Infinitive (OI) errors, although both models underestimate the error rate in English. A second set of analyses shows strong lexical effects in the pattern of errors across all five languages studied. This finding is problematic for the Variational Learning Model and provides strong support for the notion that OI errors are incomplete compound finites as instantiated in MOSAIC.

### Introduction

In many languages, children go through a stage in which they produce non-finite forms in contexts where a finite form is obligatory in the adult language. Thus, English speaking children produce utterances such as *he go* instead of the correct *he goes*, while Dutch speaking children produce utterances such as *Papa lopen* (Daddy walk) instead of *Papa loopt* (Daddy walks). Since these errors occur at a stage when the child is also producing correctly marked finite verb forms, they have come to be known as Optional Infinitive (OI) errors (Wexler, 1994).

Early theories of the OI stage attempted to explain the finding that children typically produce OI errors in obligatory subject languages such as English and Dutch, but not in optional subject languages such as Spanish and Italian. What these theories failed to explain, however, is that there is considerable cross-linguistic variation in the rates at which children produce OI errors even in obligatory subject languages. For example, Phillips (1995) reports data from children learning 9 different languages and concludes that rates of OI errors vary on a continuum from high in English and Swedish, through medium in French, Dutch and German to low (but by no means zero) in Spanish, Italian, Catalan and Hebrew. These data present a problem for theories that assume a qualitative distinction between OI and non-OI languages as they suggest that the OI stage is best viewed as a graded quantitative phenomenon.

Two recent theories that attempt to explain quantitative variation in OI errors are the Variational Learning Model (VLM; Legate & Yang, 2007), and MOSAIC (Freudenthal, et al. 2005, 2007). While both theories attempt to explain cross-linguistic variation as a function of the distributional statistics of the input, the underlying assumptions of the

theories are very different. The VLM is a generativist theory that states that the grammar that children employ to parse the input varies probabilistically as a function of the amount of evidence in the input that supports the hypothesis that they are learning a language with obligatory tense marking. MOSAIC is a constructivist theory that views OI errors as incomplete compound finites (modal/auxiliary plus infinitive constructions). MOSAIC explains cross-linguistic differences in rates of OI errors in terms of differences in the rates at which compound finites occur in the input, as well as the way in which compound finites pattern across languages.

This paper aims to compare the VLM and MOSAIC by determining the success of both theories in explaining levels of OI errors across five different languages. Since both theories predict that these levels will be correlated with measures of the distributional statistics of the input, there is an obvious possibility that the theories will be equally successful. A further test of the theories is therefore conducted by looking at an area where the theories clearly make different predictions: lexical effects on the pattern of OI errors. According to the VLM, OI errors are caused by the fact that the child (probabilistically) entertains the possibility that she is learning a language that does not mark Tense. Since this is an attribute of the grammar, the VLM predicts that rates of OI errors should be roughly equal across the verbs that the child uses. According to MOSAIC, OI errors are incomplete compound finites that have been learned from the input. MOSAIC therefore predicts that children produce OI errors at differential rates for individual verbs: OI errors will be frequent for verbs that frequently occur in compound finites, and infrequent for verbs that rarely occur in compound finites.

### The Variational Learning Model

Legate and Yang's (2007) Variational Learning Model (VLM) is designed to account for the quantitative, cross-linguistic pattern of rates of OI errors from a generativist perspective. According to the VLM, children probabilistically entertain the hypothesis that they are learning a language that requires tense marking. The probability associated with this hypothesis varies dynamically in response to the linguistic input. Thus, if the child attempts to parse an utterance that is inconsistent with the grammar (or parameter setting) that the child entertains, this grammar will be punished, and its associated

probability lowered. According to the VLM, children initially entertain the possibility that they are learning a language that does not have tense marking (such as Mandarin Chinese). For children learning tense-marking languages, this means that they will initially produce incorrect utterances that lack inflection (i.e. OI errors). As children are subjected to input that is inconsistent with this hypothesis, they will gradually abandon this hypothesis in favour of the correct hypothesis that they are learning a tense-marking language. As a result, the number of untensed utterances that they produce will gradually decrease. An important feature of the VLM is that the speed with which a hypothesis is abandoned is dependent on the amount of evidence against that hypothesis that the child encounters in the linguistic input. Thus, in languages with rich verb morphology (such as Spanish) children will quickly converge on the hypothesis that they are learning a tense-marking language. In languages with impoverished verb morphology (such as English), this will take longer, and children are expected to produce more OI errors.

Legate and Yang (2007) test their model by analysing input from three different languages: English, French and Spanish. They determine the amount of evidence for a tense-marking language (or a [+Tense] grammar) by analysing utterances from corpora of child directed speech. Utterances that contain verb forms that are overtly tensed or display tense-dependent morphology are counted as rewarding a [+Tense] grammar, while (verbal) utterances that do not display tense or tense-dependent morphology punish a [+Tense] grammar. Applying their analysis to English, French and Spanish, Legate and Yang find that the proportions of verbal utterances that reward the [+Tense] grammar are .53, .70 and .80 respectively. These numbers correspond well with rates of OI errors for these languages that are reported in the literature; English is generally considered a language with high rates of OI errors, French is an intermediate language, while OI errors in Spanish are rare. It could, however, be argued that these languages do not constitute a particularly strong test of the VLM, as they are relatively far apart on the continuum of OI languages. One of the aims of this paper is therefore to provide a stronger test of the VLM by including two additional intermediate OI languages: Dutch and German.

## MOSAIC

MOSAIC (Model Of Syntax Acquisition In Children) is a computational model that has already been applied to the simulation of the OI phenomenon in English, Spanish, Dutch and German (Freudenthal et al., 2005, 2007). MOSAIC learns to produce output corpora of increasing Mean Length of Utterance (MLU) from input consisting of orthographically transcribed child directed speech. Importantly, MOSAIC does this without employing any abstract linguistic knowledge. That is to say, MOSAIC is a simple distributional analyser that is dependent on cross-linguistic differences in the distributional statistics of its input to simulate cross-linguistic differences in the

characteristics of early child speech. The version of MOSAIC used for the simulations reported in this paper combines MOSAIC's utterance-final bias (as described in Freudenthal et al. 2007) with a smaller utterance-initial bias (A preliminary version of this model is described in Freudenthal et al. (2005)). It employs right edge learning, in which the representation of an utterance is slowly built up by starting at the end (right edge) of the utterance and slowly working its way to the beginning of the utterance. Right edge learning allows MOSAIC to produce (incomplete) utterance-final phrases, which resemble certain types of child error, including subject omission errors (e.g. *go to the shops*). MOSAIC's right edge learning is complemented with a (slower) left edge learning mechanism. This mechanism, which builds up representations from the beginning of the utterance, allows for the association of short utterance-initial phrases with (longer) utterance-final phrases, resulting in output with missing sentence-internal elements. MOSAIC simulates OI errors by producing compound constructions (modal/auxiliary plus infinitive constructions) with missing modals/auxiliaries (e.g. *he (can) go to the shops*). MOSAIC has been shown to successfully simulate differential rates of OI errors in English, Dutch, German and Spanish. The differential rates of OI errors in these languages are explained by the number of compound constructions in the input and the way in which these compounds pattern. English, which has many compound constructions, has high rates of OI errors. Dutch and German do not have very high rates of compound constructions, but the fact that these languages are SOV/V2 languages means that non-finite forms occur in utterance-final position. Since MOSAIC has an utterance-final bias, early output from MOSAIC contains many nonfinite verb forms. Spanish has rates of compounds that are comparable to Dutch and German, but since object complements occur after verbs in Spanish, the number of utterance-final non-finites is not very high. Freudenthal et al. (2007) show that, cross-linguistically, the rate of utterance-final non-finites is closely related to the proportion of OI errors that children produce. However, certain features of the coding scheme used by Freudenthal et al. (in particular the restriction of English to utterances containing a 3<sup>rd</sup> singular subject) mean that a direct comparison of English with the other languages is not possible. The present paper introduces a new coding scheme that overcomes this problem. The present paper also introduces an additional language (French) in order to compare MOSAIC and the VLM across 5 languages that take up different positions on the OI continuum.

## Cross-linguistic rates of Children's OI errors

Cross-linguistic rates of OI errors in child speech were determined by analysing corpora of child speech at an MLU of approximately 2.0 and determining the proportion of utterances containing a main verb that did not contain any inflected verb forms (i.e. that only contained non-finite verb forms). The following languages were studied: English,

Dutch, German, French and Spanish. The children analysed for the respective languages were drawn from the following corpora: Manchester (Theakston et al., 2001), Groningen (Bol, 1995), Leo (Behrens, 2006), Lyon (Demuth & Tremblay, 2008), and Nottingham (Aguado-Orea, 2004). Apart from the Leo corpus, all these corpora are available through the CHILDES data-base (MacWhinney, 2000). With the exception of English, distinguishing between finite and non-finite utterances is relatively straightforward as non-finite forms are easily distinguished from finite forms<sup>1</sup>. In English, however, such a procedure is hampered by the fact that all present tense forms except for the 3<sup>rd</sup> singular match the infinitive. Thus, in contrast to the other languages studied, English OI errors can only be identified in 3<sup>rd</sup> singular contexts. In previous work with MOSAIC this difficulty was solved by restricting the analysis of English to utterances containing a 3<sup>rd</sup> singular subject. However, the present study specifically aims to compare OI rates across languages. Since the restriction of English utterances to utterances containing a 3<sup>rd</sup> singular subject precludes such a comparison, the child data for English were hand coded and the context in the transcript was used to restrict the analysis to all utterances (including subjectless utterances) that were produced in a 3<sup>rd</sup> singular context. For the other languages, where non-finite and finite forms are easily distinguished, an automated analysis was performed. Table 1 shows the results of this analysis. As can be seen in Table 1, there is clear cross-linguistic variation in the rates of OI errors. OI errors are rare in Spanish and extremely frequent in English and Dutch, with intermediate levels of error in German and French.

In order to determine whether the data from these children are representative for their languages, one additional child from each corpus was analysed. The proportion of OI errors for these children were: Anne (English): .87 (109); Peter (Dutch) .74 (290); Rah (German): .58 (178), Anais (French): .41 (203), Lucia (Spanish): .05 (62). These numbers suggest that the data presented in Table 1 are reasonably representative of children learning these languages, and are therefore suitable for a comparison of the VLM and MOSAIC.

Table 1: Cross-linguistic rates of OI errors at an MLU of approximately 2.0.

	MLU	Prop. OI errors (Total number of utterances)
Becky (English)	2.17	.97 (98)
Matthijs (Dutch)	2.06	.77 (347)
Leo (German)	2.08	.58 (3967)
Tim (French)	1.96	.32 (250)
Juan (Spanish)	2.17	.20 (305)

<sup>1</sup> This is not strictly true for Dutch and German where plural present tense forms match the infinitive. However, since plural forms are relatively infrequent, this complication is unlikely to greatly affect the estimate of error rates.

## Cross-linguistic predictions of the VLM

The predictions of the VLM were tested by analysing the child-directed speech from the corpora analysed in the previous section, and determining the proportion of utterances that provide evidence that the language that the child is learning is a tense-marking language (i.e. a language with a [+Tense] grammar). Whether or not an utterance provides evidence for a [+Tense] grammar depends on whether or not the utterance contains verb forms in tensed position that are either overtly marked for tense or display tense-dependent morphology. Since languages differ in terms of the number of verb forms that are overtly marked for tense, this analysis differs for different languages. The analysis performed on English, French and Spanish matches that of Legate & Yang (2007), to whom the reader is referred for further details. Due to space restrictions, only the analysis of English is briefly outlined here. In English, all 3<sup>rd</sup> singular present tense forms carry the -s morpheme and are thus counted as rewarding the [+Tense] grammar. The other present tense forms do not display tense-dependent morphology and are thus counted as punishing the [+Tense] grammar, as are present tense modals. Past tense modals, as well as (inflected forms of) copula and auxiliary *be* and *have* are marked for tense and thus reward the [+Tense] grammar. Thus, utterances like *He can go*, and *They go* do not display tense-dependent morphology and punish the [+Tense] grammar, while utterances like *He is going* and *He goes* reward the [+Tense] grammar. Regular past tense forms are slightly ambiguous as they match the participle. In line with Legate and Yang's analysis such forms were counted as tensed when they occurred in tensed position (i.e. in the absence of an auxiliary as in *He walked*), and as untensed when they did not occur in tensed position (i.e. in the presence of an auxiliary as in *They have walked*). The latter utterance was thus counted as punishing the [+Tense] grammar, unless the auxiliary itself was tensed (*He has walked*).

Legate and Yang do not perform an analysis of Dutch and German, and one therefore needs to assess Dutch and German verb morphology in order to determine what utterances reward and punish the [+Tense] grammar. Since Dutch and German verb morphology are almost identical, the analysis will be illustrated using an example of a (regular) Dutch verb: *werken* (work).

Table 2: Conjugation of a regular Dutch verb.

	Infinitive: werken		Participle: gewerkt
	Present Tense	Past Tense	
1 <sup>st</sup> singular	Ik werk	Ik werkte	
2 <sup>nd</sup> singular	Jij werkt	Jij werkte	
3 <sup>rd</sup> singular	Hij werkt	Hij werkte	
1 <sup>st</sup> plural	Wij werken	Wij werkten	
2 <sup>nd</sup> plural	Jullie werken	Jullie werkten	
3 <sup>rd</sup> plural	Zij werken	Zij werkten	

As can be seen in Table 2, plural present tense forms match the infinitive and were thus counted as punishing the

[+Tense] grammar. Past tense forms and the 2<sup>nd</sup> and 3<sup>rd</sup> present tense singular display Tense (or Tense-dependent morphology) and were counted as rewarding the [+Tense] grammar. The 1st singular present tense form, while distinct from the infinitive, matches the verb stem. In line with Legate and Yang's analysis of English, French and Spanish, such forms were counted as punishing the [+Tense] grammar. Dutch and German modals differ from English modals in the sense that they inflect as main verbs (and can be used as main verbs). Thus, inflected modals (past tense and singular present tense, except when matching the stem) were counted as rewarding the [+Tense] grammar.

The inflection of German verbs is almost identical to the inflection of Dutch verbs, with the following exceptions. The 1st singular present tense has an *-e* suffix. Since this suffix consists of a single schwa (which is often not pronounced), this form was counted as a [+Tense] punishing stem. The 2<sup>nd</sup> singular present tense form has an *-st* (rather than *-t*) suffix. Like the Dutch 2<sup>nd</sup> singular, this form was counted as rewarding the [+Tense] grammar. The German 2<sup>nd</sup> plural present Tense, has a *-et*, rather than *-en* suffix. Since this form does not match the infinitive, it was counted as rewarding the [+Tense] grammar. One final difference between Dutch and German is that in Dutch, but not in German, the 2nd singular present tense *-t* suffix is dropped in questions (which are formed through main verb inversion). Thus, the declarative *Jij werkt* (You work) is marked for Tense and rewards the [+Tense] grammar, but the interrogative *Werk jij?* (Work you?) is not. This feature of Dutch results in the number of stems in Dutch input being higher than in German.

The analysis of the amount of evidence for the [+Tense] marking grammar was performed in an identical manner for all five languages. Based on the analysis of the verb morphology, a list of verb forms that were and were not marked for Tense was drawn up. Next, an (automated) search was performed on all parental utterances in the corpus. All utterances that contained (at least one) verb form that was marked for Tense were counted as rewarding the [+Tense] grammar, while utterances that only contained verb forms that were not marked for Tense were counted as punishing the [+Tense] grammar. The proportion of utterances with Tensed verbs then constituted the amount of evidence for the [+Tense] grammar in that language. Table 3 displays the results of this analysis. The numbers derived for English, French and Spanish closely match those reported by Legate & Yang (.53, .70 and .80 respectively)

As can be seen in Table 3, the VLM correctly predicts the ordering of Spanish, French, German and Dutch. However, it fails to predict the high levels of OI errors in English, which actually provides more evidence for the [+Tense] grammar than Dutch. The VLM therefore predicts that Dutch children should make more OI errors than English children, when the data show the opposite pattern. One possible way of resolving this discrepancy would be to argue that the present analyses underestimate the amount of evidence for the [+Tense] grammar in Dutch and German

by treating stems (that do not match the infinitive) as punishing the [+Tense] grammar when they should be treated as rewarding the [+Tense] grammar. In fact, however, treating stems as if they rewarded the [+Tense] grammar results in a worse rather than a better fit to the child data. Thus, when Dutch and German stems are treated as tensed forms, the amount of evidence in Dutch and German increases to .85 and .83 respectively, levels that are even higher than the estimates for French and Spanish. This would result in the VLM predicting low levels of OI errors in Dutch and German. What's more, there is no longer a difference between German and Dutch in the amount of evidence for a [+Tense] grammar. This means that, if the VLM treated stems as rewarding the [+Tense] grammar, the model would fail to predict the differential rates of OI errors that exist between these languages.

Table 3: Cross-linguistic proportion of tensed utterances in child-directed speech.

	Number of clauses	Proportion Tensed
Becky (English)	16138	.54
Matthijs (Dutch)	8176	.49
Leo (German)	18413	.62
Tim (French)	14169	.67
Juan (Spanish)	19044	.81

### Cross-linguistic predictions of MOSAIC

MOSAIC's predictions about cross-linguistic variation in rates of OI errors were derived by training MOSAIC on the child-directed speech of the children analysed in Table 1. The models were trained until their output reached an MLU of approximately 2.0. The models' output was then analysed in the same manner as the child data. Utterances were divided into utterances containing only non-finite verb forms (OI errors), and utterance that contained (at least one) finite verb form (or utterances containing both finite and non-finite verb forms). Next, the proportion of OI errors was determined. As with the child analysis, such an analysis is problematic in English where an OI error can only be reliably diagnosed in a 3rd singular context. The solution to this problem was similar to the one employed in the child analysis (where the corpus was hand coded for 3<sup>rd</sup> singular contexts). The entire input corpus (~ 25,000 utterances) was hand coded to determine which verbs were uttered in a 3<sup>rd</sup> singular context. All such verbs were tagged for having occurred in a 3<sup>rd</sup> singular context. This made it possible to differentiate utterances in MOSAIC's output that had been learned from a 3<sup>rd</sup> singular context (e.g. *(he can) go-3<sup>RD</sup>*) from those that had been learned from a non-3<sup>rd</sup> singular context (e.g. *you can go*). Table 4 presents the rates of OI errors in MOSAIC's output. As can be seen in Table 4, MOSAIC performs slightly better than the VLM in that it correctly predicts the ordering of OI errors across the five languages: Spanish < French < German < Dutch < English. It is also apparent, however, that MOSAIC underestimates the levels of OI errors in English to a far greater extent than

it does for the other languages. Thus, like the VLM, MOSAIC has difficulty explaining the high levels of OI errors found in English children. The fact that English is a special case is underscored by an analysis of the proportion of utterance-final non-finites in the input (.89 for English, .87 for Dutch, .69 for German, .40 for French and .21 for Spanish). Rates of OI errors in English children are higher than the proportion of utterance-final non-finites, whereas they are lower for all other languages. A possible explanation for this discrepancy will be discussed in the conclusions.

Table 4: Proportion of OI errors in MOSAIC models

	MLU	Prop. (N)
Becky (English)	2.07	.72 (119)
Matthijs (Dutch)	1.95	.65 (561)
Leo (German)	1.96	.49 (1508)
Tim (French)	1.95	.32 (510)
Juan (Spanish)	2.08	.15 (1514)

### Distinguishing between MOSAIC and VLM

While MOSAIC and the VLM make surprisingly similar predictions regarding the levels of OI errors in the different languages, there is an area where the two theories clearly make different predictions. Since the VLM operates at the level of the grammar, it predicts that children's rates of OI errors will be roughly equal across different verbs. That is, since the child probabilistically represents the need for verbs to be tensed in the language it is learning (and this probability is determined by the amount of evidence for the [+Tense] grammar in the input), it should apply Tense with equal probability regardless of the identity of the verb it produces. According to MOSAIC, the nature of the child's output is more directly determined by the input it hears. OI errors are compound finites with missing modals or auxiliaries. MOSAIC therefore predicts that the verbs produced as OI errors will be those that occur in compound finites in the input, whereas verbs that tend to occur in finite form in the input will rarely be used in non-finite form.

In order to test for such lexical effects, finite and non-finite uses of verbs by the children analysed earlier were tabulated on a verb-by-verb basis. That is, across the languages we assessed for each individual (main) verb how often it was used as a correct finite form and an incorrect infinitive form. Next, the parental speech directed at these children was analysed to determine how often the target verbs used by the children were used as a finite form (in a simple finite utterance) or as an infinitive in a compound construction. A high correlation (across verbs) between these rates of finite and infinitive use constitutes strong evidence that children's OI errors are not finite constructions with missing inflection, but rather compound constructions with missing modals or auxiliaries. Such a finding would be problematic for the VLM. Where possible (in English, Dutch and German) the analysis was carried out on a sample of child speech where finite and non-finite verbs occurred at roughly equal rates. For Spanish and

French this was not possible (as the children produced few OI errors), and the sample analysed in Table 1 (at an MLU of 2.0) was used. Using a sample with few OI errors inevitably means that the likelihood of finding a significant correlation is reduced, as the variation in the sample is decreased relative to a sample where finite and infinitive uses occur at roughly the same rates.

As in the earlier analyses, the analysis of English was restricted to utterances produced in 3<sup>rd</sup> singular contexts. In addition, in the German and Dutch analyses, plural forms were counted as finites, rather than as the infinitives they match. The correlations between children's infinitive uses and parental compound uses of individual verbs are displayed in Table 5. Two correlations are reported for each language. The first correlation was computed on all verbs produced by the child. This set includes verbs that were used only once. Such verbs are likely to contribute to sampling error as their estimate of finite usage is likely to be somewhat unreliable. A second correlation was therefore computed on a restricted set of verbs: verbs that were used at least three times by the child. This restricted set is likely to provide a better estimate of the child's finite and infinitive usage.

Table 5: Correlations between children's level of finiteness marking and parental use of compound constructions. (Number of contributing verbs in parentheses). + =  $p < .10$ , \* =  $p < .05$ , \*\* =  $p < .01$

	Full set	Restricted set
English	.35* (43)	.55* (15)
Dutch	.71** (102)	.83** (59)
German	.48** (143)	.68** (69)
French	.45** (75)	.57** (37)
Spanish	.40** (69)	.29+ (43)

As can be seen in Table 5, there are clear lexical effects in all of the five languages. Thus, nine out of ten correlations are significant at  $p < .05$ , and the remaining correlation (Spanish) is marginally significant ( $p = .056$ , two-tailed). These results are difficult to explain in terms of the VLM, and provide strong evidence for the claim that OI errors are learned from compound finites in the input.

### Conclusions

This paper set out to test two accounts (MOSAIC and the VLM) of cross-linguistic differences in the rate at which children produce OI errors. The analyses reported suggest that both MOSAIC and the VLM fit the cross-linguistic data surprisingly well. Both models correctly predict the relative order of rates of OI errors in Spanish, French, German and Dutch. However, both models also struggle with the high rates of OI errors in English. A further analysis showed clear evidence of lexical effects in the data, which are predicted by MOSAIC, but are problematic for the VLM.

The reason that the VLM fails to predict the high levels of OI errors in English appears to be that it operates at too high a level of abstraction. Thus, while the evidence for tense

marking on English *lexical verbs* is indeed very low (lower than it is for Dutch), English input actually provides a substantial amount of evidence for tense marking in the form of tensed copulas and auxiliaries. According to the VLM, this evidence ought to drive down the level of OI errors in English children's speech, but it does not appear to do so. The implication is that for a model like the VLM to explain the data on English, it would need to be made preferentially sensitive to the very low level of tense marking on lexical verbs.

The reason that MOSAIC fails to predict the high levels of OI error in English is that the proportion of utterance-final verbs that are non-finite in English is simply not high enough to explain the almost exclusive production of OI errors during the early stages. One reason why English might be a special case in this respect is that, in English, the infinitive is indistinguishable from the bare stem. Since the only present tense form that is not a bare stem in English is the 3<sup>rd</sup> singular, a much higher proportion of lexical verb forms in the input are either infinitives or forms that are indistinguishable from the infinitive. This fact is likely to slow down the process of paradigm building in English and result in default effects where the child produces a bare stem/infinitive in the absence of knowledge of the relevant 3<sup>rd</sup> singular or past tense form. Since MOSAIC is insensitive to the morphological structure of the verbs that it encodes, it is clearly unable to simulate this kind of default effect.

Of course, if it is necessary to supplement MOSAIC's account of OI errors with some kind of paradigm-building account in order to explain the English data, one might wonder whether it is possible to simply replace MOSAIC with a paradigm-building account (e.g. MacWhinney, 1978). We would argue that there are at least three reasons for seeing MOSAIC and paradigm building as complementary rather than competing accounts. First, although a paradigm-building account provides a very natural way of explaining OI errors in English, it fares much less well as an account of OI errors in other languages. This is because, in languages other than English, infinitives are typically not the most frequent form in the language, and tend to carry infinitival morphology that distinguishes them from other more frequent forms. It is therefore difficult to see why children learning these languages would default to the infinitive rather than to the bare stem or to some other more frequent inflected form.

Second, a paradigm-building account of OI errors would seem to predict that defaulting to the infinitive would reflect some more general confusion on the part of the child between finite and non-finite forms. In fact, however, a key feature of the OI stage is that, although children produce infinitives in tensed position, their use of these forms is highly sensitive to differences in the distributional properties of finite and non-finite forms in the input. For example, in Dutch and German, children correctly place finite forms before their complements and the infinitives in OI errors after their complements (Wexler, 1994).

Finally, a paradigm-building account offers no obvious explanation of the lexical effects found in the present study. These effects, which can be seen in all 5 of the languages under investigation, suggest that those verbs that occur as OI errors in children's speech also tend to occur as infinitives in adult compound finites. They thus provide strong support for the idea, instantiated in MOSAIC, that OI errors are learned from compound structures in the input.

## Acknowledgments

This research was funded by the ESRC under grant number RES000230211.

## References

- Aguado-Orea, J. (2004). *The acquisition of morpho-syntax in Spanish: Implications for current theories of development*. Unpublished doctoral dissertation, University of Nottingham, United Kingdom.
- Behrens, H. (2006). The input-output relationship in first language acquisition. *Language and Cognitive Processes*, 21, 2-24.
- Bol, G. W. (1995). Implicational scaling in child language acquisition: the order of production of Dutch verb constructions. In M. Verrips & F. Wijnen, (eds.), *Papers from the Dutch-German Colloquium on Language Acquisition*, Amsterdam Series in Child Language Development, 3, Amsterdam: Institute for General Linguistics.
- Freudenthal, D., Pine, J. M. & Gobet, F. (2005). Simulating the cross-linguistic development of optional infinitive Errors in MOSAIC. In B. G. Bara, L. Barsalou & M. Buchiarelli (Eds.), *Proceedings of the 27th Annual Meeting of the Cognitive Science Society* (pp. 702-707). Mahwah, NJ: Erlbaum.
- Freudenthal, D., Pine, J. M., Aguado-Orea, J. & Gobet, F. (2007). Modelling the developmental patterning of finiteness marking in English, Dutch, German and Spanish using MOSAIC. *Cognitive Science*, 31, 311-341.
- Legate, J. A. & Yang, C. (2007). Morphosyntactic learning and the development of tense. *Language Acquisition*, 14, 315-344.
- MacWhinney, B. (2000). *The CHILDES project: Tools for analysing talk (3<sup>rd</sup> Edition)*. Mahwah, NJ: Erlbaum.
- MacWhinney, B. (1978). The acquisition of morphophonology. *Monographs of the Society for Research in Child Development*, 43, 1-123.
- Phillips, C. (1995). Syntax at age two: Cross-linguistic differences. In C. Schütze, J. Ganger & K. Broihier (eds.), *Papers on Language Processing and Acquisition. MIT Working Papers in Linguistics*.
- Theakston, A. L., Lieven, E. V. M., Pine, J. M. & Rowland, C. F. (2001). The role of performance limitations in the acquisition of verb-argument structure: an alternative account. *Journal of Child Language*, 28, 127-152.
- Wexler, K. (1994). Optional infinitives, head movement and the economy of derivation in child grammar. In N. Hornstein & D. Lightfoot (Eds.), *Verb Movement* (pp. 305-365). Cambridge: Cambridge University Press.