

Development of Verb Morphology: From Item-Specificity to Proficient Use

Jekaterina Mažara (jekaterina.mazara@uzh.ch)

University of Zurich, Zurich, Switzerland

Sabine Stoll (sabine.stoll@uzh.ch)

University of Zurich, Zurich, Switzerland

Abstract

The initial phase of linguistic production by children is characterized by rote-learned, lexically restricted forms and constructions. Only during later phases of language acquisition do they develop flexibility across a paradigm and mix lexical and grammatical material more freely. In the development of verb morphology, a correlation between the use of tense and aspect has been observed in many languages. It has been suggested that this leads to an intermediary state of paradigm categorization based on temporal categories. So far the flexibility of individual verbs occurring in different tense-aspect combinations has not been examined in detail. Here we evaluate the flexibility of verb use in a large longitudinal corpus of 4 Russian children. We compute the Shannon entropy of verb stems distributed over individual grammatical forms. Results show that children do not pass through a stage of paradigm categorization based on aspecto-temporal categories. After a brief item-specific phase of rote learned forms, they quickly become flexible users of verbs in both aspects.

Keywords: language acquisition; corpus study; item-specificity; verb morphology; aspect; Russian

Introduction

Usage-based approaches to language acquisition propose an early phase during which children use a small number of lexically specific constructions which are presumably rote-learned (Lieven, Pine, & Baldwin, 1997; Pine & Lieven, 1997; Tomasello, 2000, 2003). During this short phase of lexical specificity, flexibility of word form use is very low, but soon after using the first rote-learned constructions, children start to produce new forms and apply them to new contexts. So far, relatively little is known about this generalization process from lexically specific constructions to full productivity.

In this study, we focus on the acquisition of Russian verb morphology and the role of aspect. Grammatical aspect is the expression of the viewpoint on the temporal structure of an event. *Perfective aspect* describes an external and temporally bounded view of a completed event, while *imperfective aspect* focuses on the internal stages or temporal extension of an event (Comrie, 1976).

Languages differ vastly in how (and if) they mark grammatical aspect but independent of the realizations, aspect has been found to play a pivotal role in the acquisition of the verbal system in relation with tense (Shirai & Anderson, 1995; Shirai, Slobin, & Weist, 1998). Correlations between verbs with a defined end-point (telic verbs) and perfective past marking as well as verbs without a defined end-point (atelic) and non-past imperfective marking have been

found in early acquisition of a number of different languages (cf. Bloom, Lifter, and Hafitz (1980); Harner (1981); Shirai and Anderson (1995); Clark (1996); Johnson and Fey (2006) for English, Bronckart and Sinclair (1973) for French, Antinucci and Miller (1976) for Italian, Li and Bowerman (1998); Shirai and Anderson (1995); Shirai et al. (1998); Li and Shirai (2000) for Japanese; Stoll (1998, 2005); Stoll and Gries (2009); Gagarina (2000); Bar-Shalom (2002) for Russian; Li (1990); Li and Shirai (2000) for Mandarin; Aksu-Koç (1998) for Turkish; Stephany (1985) for Greek; Weist, Wysocka, Witkowska-Stadnik, Buczowska, and Konieczna (1984); Weist and Konieczna (1985) for Polish; as well as self-organizing feature map models (cf. Li (2000); Li and Shirai (2000)).

It has been suggested that due to the presence of this correlation, after the lexically-specific phase, the development of productivity passes through an intermediary stage, during which children are more productive in their use of verbal morphology with the appropriate prototypes of a category (also known as the *Aspect Hypothesis* see Shirai and Anderson (1995)). These correlations are also present in the speech of adults, albeit to a lesser degree. However, to date, only a few studies have systematically compared these correlations in child and child-surrounding speech. For Russian children, Stoll and Gries (2009) have found a gradual decrease of this association in children over the course of development.

The goal of this study is to examine the development of flexibility of verb form use in Russian children. We test whether there is indeed a transition phase based on the tense-aspect correlation during which children are more productive within sub-categories of the verb paradigm before becoming fully productive verb users.

We first establish phases in production based on verb form inventory size. We then compare both type and token distributions in children's use during these phases to that of adults. We show that in token use, both adults and children display distributional bias of tense-aspect correlations. The bias is stronger in children in the first phase of production and approaches adult levels in the second phase. We evaluate the flexibility of use over time by measuring the entropy of lemmas used with individual grammatical forms. We show that, as item-specificity decreases, a great variety of forms is introduced early on and quickly generalized so that both past and non-past marking is used with verbs of both aspects.

Verb morphology in Russian

Russian has relatively complex verbal morphology centering on a semantically and morphologically complex category of grammatical aspect which interacts with tense. Grammatical aspect in Russian is characterized by a perfective/imperfective distinction and each verb is either perfective or imperfective. In contrast to English which has one single aspectual marker (*-ing*), Russian has many different markers for the perfective aspect (mainly prefixes and one suffix) and one suffix (with various allomorphs) for the imperfective aspect or zero marking.

On the functional level, several temporal and contextual features influence the use of the two aspects. Russian imperfective verbs are used when the duration of an action is relevant (e.g. *ona čitaet ves' den'* 'she reads all day') and if the action is presented as a completed event (e.g. *ona včera čitala ves' den'*, 'she spent all day reading yesterday'). Perfective verbs are used when the focus of the utterance is a boundary of the action; this can be either the beginning of an action, the end/result or both (e.g. *ona dočitala knigu*, 'she finished reading the book'). Morphologically, perfectives are typically derived from imperfectives by prefixation. To complicate things, however, the meaning cannot be derived via simple rules (Timberlake, 2004) and always involves some degree of rote-learning. There is no one-to-one relationship between prefixes and the resulting meaning change in the verb they are attached to. Further, most verbs can combine with multiple prefixes, while others are restricted in their combinability.

Verbs of both aspects express other verb categories (person, number, tense, voice, and mood) with the same morphemes. There are, however, some differences in meaning. Non-past morphology denotes present tense when it appears with imperfectives, but expresses the future in combination with perfectives. To express imperfective future, an analytic form is used (consisting of a finite 'to be' auxiliary and the infinitive of the main verb). In this paper, we focus on the acquisition of synthetic morphology and, therefore, exclude the analytic future. Past morphology can be used with both aspects equally.

The broad generalization found in the works cited in Shirai et al. (1998) states that children begin their acquisition of verb forms by using past morphology with achievement verbs and progressive morphology with activity verbs and only later extend it to the other group. Since lexical aspect is not annotated in the corpus we use, we focus on correlations between grammatical aspect and tense. However, this still allows us to assess this hypothesis, since achievements are necessarily perfectives and activities are necessarily imperfectives in Russian. We will, therefore, focus on whether Russian children display correlations between perfective aspect and past tense (e.g. *On doel sup*, 'he ate the soup' (meaning: he finished the bowl)) and imperfective aspect and non-past marking (*On smotrit televizor*, 'he is watching TV').

Methods

Data

The data is extracted from an audio-visual longitudinal corpus of Russian language acquisition (Stoll & Meyer, 2008) comprising data of six monolingual children living in St.Petersburg, Russia. All recordings were done in naturalistic settings at the home of the children and include the focal child and a varying number of surrounding speakers including siblings (excluded here) and adults. The children were recorded for one hour each week. We focus on 4 children, whose recordings started before the age of three. The entire corpus is transcribed and words are annotated for part of speech and morphology. Table 1 summarizes the number of utterances, words, and verbs uttered by each focal child as well as the age range of recording.

Table 1: Age spans of the focal children and number of words produced by the children and surrounding adults

Focal Child	Age span	Number of recordings	N(tokens)			
			Child		Adults	
			words	verbs	words	verbs
1	1;8.10 - 4;8.21	130	241,948	38,843	301,418	60,987
2	1;4.23 - 4;1.24	109	57,929	5,411	354,034	65,173
3	1;3.24 - 4;9.29	123	74,926	10,733	423,078	84,659
5	1;11.28 - 4;3.12	67	97,397	16,585	223,289	43,149

Finding phases in acquisition

First, we establish whether there are phases in verb form acquisition. The phases were derived directly from the target children's verb form production. We computed the additive growth in full verb forms (stem+grammatical markers) over time. The growth curves show a slow rate of increase in the earlier sessions followed by a sudden increase in the rate of newly observed forms¹. To estimate the age at which this change in rate of acquisition occurs, we conducted a segmented regression on the growth curve of each child. The break points at which the regression created a new segment are summarized in Table 2. We use these points as the estimated end of the first phase of production for the next analysis.

Table 2: Break points in growth curve as identified by segmented regression.

Child	Break-point
Child 1	2;2
Child 2	3;3
Child 3	2;3
Child 5	before recordings started

¹This was the case for all but Child 5 who already had highly developed speech at the onset of the recordings. Child 5, therefore, did not exhibit this change in rate of newly observed forms.

Entropy of verb form use

To assess the development of flexibility of form use in the observed production, we used Shannon Entropy (Shannon, 1948), the rate at which a process produces information by characterizing the balance of frequency distributions over a set of elements. If the probability of produced elements is distributed equally among them, the output is less predictable. Early child language is usually characterized by the repeated use of a few forms, while other forms might appear only once. This would result in a highly predictable output and low entropy. The formula for Shannon entropy is given in Eq. 1

$$H(X) = - \sum_{i=1}^N p(x_i) \log p(x_i) \quad (1)$$

where N is the number of distinct forms and $p(x)$ is the probability of occurrence of a specific form ².

Analysis 1: Distribution of forms in the first and second phase To gain a better understanding of verb form production during the first and the second phase of development, we extracted the verb lemmas (lexical elements) used before the break point in development. To obtain a sample comparable in size and lexical coverage, we extracted the same lemmas from the adults' production during this phase and sampled the same number of tokens as produced by the focus child. Finally, we conducted the same procedure for both focus child and surrounding adults for the second phase. To gain a first insight into the form use and assess the level of item-specificity, we visualised the data in mosaic plots showing both type and token use of children and adults in both phases. To characterize the difference between the distributions, we computed the Jensen-Shannon divergence (Lin, 1991) between each child and their surrounding adults, and between the child's own first and second phases. Jensen-Shannon divergence (JSD) measures the distance between two probability distributions over the same elements (i.e. verb lemmas in this case). The formula for JSD for two distributions P and Q with equal weight (0.5) is given in Eq. 2.

$$JSD(P||Q) = \frac{1}{2}D(P||M) + \frac{1}{2}D(Q||M) \quad (2)$$

M represents the average distribution $M = \frac{1}{2}(P + Q)$; and D stands for Kullback-Leibler divergence (KLD, sometimes also called *relative entropy*), given in 3.

$$KLD(P||Q) = - \sum_{x \in X} P(x) \log \frac{Q(x)}{P(x)} \quad (3)$$

JSD is based on KLD, but it is symmetric and its value is always finite and non-negative. When $P = Q$, $JSD = 0$ (i.e. the two distributions are equal). To evaluate the development

²For the computation of entropy used over time, where we looked at each session individually, we did not treat the system as a complete survey of the forms that have been acquired by that point in time.

of forms use across verbs from the two aspects and different grammatical categories, we computed the distribution of grammatical markers over the lexical elements extracted from the first and second phase of each child and the surrounding adults. First, we compare the probability distributions of the child and the adults during phase 1 and phase 2; then, we also compare the distribution of the child in phase 1 and the same child during phase 2. We do this both for types and tokens of verb forms.

Analysis 2: Flexibility of form use over time While JSD is useful when we can compare the probability distributions for a set of identical items, it is impossible to assess the week-to-week development in this way, since we have no way of controlling the context and lexical content of individual recording session. Cutting the production down to forms that appear in both adults' and children's production would also result in a severe underestimation of the development and a distortion of the actual production. Therefore, we compute the entropy of all elements occurring in an individual recording session. To assess whether certain tense-aspect combinations indeed aid in acquisition, we divide the data into past and non-past marked verbs and compute the entropy of perfective and imperfective verb lemmas used with past and non-past marking. As the children develop away from the item-specific phase, we expect their use of individual grammatical markers to become more flexible, i.e. they learn to combine a variety of verb lemmas with individual forms.

To estimate the time at which children start approaching adult levels of flexibility in their verb form use, we use the entropy computations of adults as a comparison within each session. The children's entropy is divided by the corresponding surrounding adults' entropy within each session. A value below 1 signifies that the child is below the adult level of entropy, values above 1 mean that the child's verb production has a higher entropy than that of surrounding adults. To control for contextual influence and other effects that might lead to particularly high or low entropies, we bootstrapped the data in each recording session for 100 iterations.³

Since the corpus consists of naturalistic data, it is difficult to normalize the production for comparative reasons. Sampling a fixed number of tokens from children and adults in each session would distort the data in a number of ways: i) if a fixed number of tokens is sampled across the recording span (e.g. 500 tokens from children and adults), the children's initial production is inflated, while adults and children's later production are underestimated; ii) if the number of tokens is determined by the number produced by the target child in each recording, this – again – severely underestimates the adults' production in the early recordings and does not represent a realistic measure for comparison. Same goes for a restriction of lexical elements used for the computation of entropy, since the fact that children's vocabulary size is growing is also an important factor and should not be ignored. This is

³The relatively low count of bootstraps was chosen for reasons of graphy clarity.

especially important for Russian, where aspect is encoded as part of the lemma.

To evaluate whether the age at which significant changes in the entropy of lemma use with individual forms happen, we fitted a generalized additive model to the data and estimated the change points of the regression to find the age at which diversification starts and when it levels off.

Results

Analysis 1: Distribution of forms in the first and second phase

Looking at the sample of matched verb lemmas and number of tokens in the two phases of each child and their surrounding adults, we see that the type distribution is slightly more diversified than the distribution of tokens. While there are tendencies to use more non-past forms with imperfective verbs and more past forms with perfectives, even during the earliest phase this tendency is not absolute and both past and non-past forms appear with verbs of both aspects early on. While types are distributed fairly equally, the token distribution is less even during both phases. This holds for both adults and children.

Table 3 shows the JSD computed for each child’s early production compared to that of surrounding adults and the child’s own production during the later phase. In the case of child 5, it was not possible to establish an early phase similar to that of the other children. Additionally, Child 5’s earliest recorded production is so varied that it was impossible to obtain a sample of the same lexical items within the same time window from the surrounding adults. Therefore, the results shown for Child 5 represent a comparison of Child 5’s production during the first 5 recordings sessions compared to a lexically and size-matched sample from his surrounding adults across the entire corpus.

Table 3: Jensen-Shannon divergence per child.

	Phase1		Phase2		Child	
	Child-to-Adults		Child-to-Adults		Phase1-to-Phase2	
	Types	Tokens	Types	Tokens	Types	Tokens
Child 1	0.124	0.501	0.020	0.056	0.122	0.421
Child 2	0.143	0.505	0.032	0.135	0.115	0.508
Child 3	0.121	0.647	0.074	0.107	0.112	0.609
Child 5	n/a	n/a	0.097	0.327	n/a	n/a

In all samples, the difference between the distributions of tokens is more pronounced than that of types, and shows less of a decrease between the two phases. However, the difference between each child’s first phase and second phase sample is comparable to the difference between the child’s production and that of adults in phase 1. This suggests that their development approaches a stage where their use of verb forms in spontaneous home interactions is very similar to that of the adults.

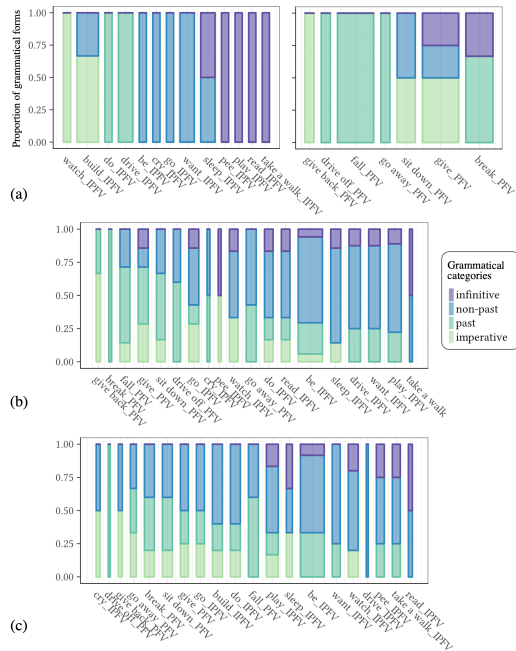


Figure 1: (a) Distribution of full form verb types in the production of Child 1 during phase 1; (b) Distribution of types in a sample of same lemmas and same number of tokens in Child 1’s production during phase 2; (c) Distribution of types in a sample of same lemmas and same number of tokens in adults’ production during phase 1.

To gain insight into the actual combinations of lemmas and forms used in each phase, Figures 1 a–c and 2 a–c exemplify the visualization of the type and token use within the sample of Child 1 and surrounding adults (we are not able to show the corresponding visualizations of the other children for space reasons). The thickness of the bars corresponds to the distribution of forms across the verb lemmas, while the colors stand for grammatical categories to which the forms belong. For types, only the plot for the child’s first phase distribution was split by aspect, because phase 2 did not show the difference as strongly. For the token use, however, all plots are split by aspect, since the token distribution of both children and adults shows more differences between the use of grammatical markers with verbs in the two aspects.

Analysis 2: Flexibility of form use over time

Entropy ratios (child/adults) of the use of lemmas with individual grammatical markers from the sub-sets of non-past and past morphology and the segmented regression reveal that difference in the onset of diversification is not large. For past morphology, perfective lemmas show an earlier increase of entropy, but imperfective lemmas follow suit only a few weeks later and vice versa. The onset of use starts with imperfective+non-past and perfective+past for all children except Child 5, whose production is already diversified at the start of recordings. Only Child 3 shows a lag of more

through an intermediary phase during which the generalization first occurs within subdivisions of the verb paradigm (for perfective verbs with past marking, for imperfective verbs with non-past), the generalization starts early across the entire paradigm. Soon after item-specificity starts decreasing, children begin applying forms of a grammatical category to verbs of both aspects. This is strengthened by the observation that verb use in the first phase of production shows a stronger distributional bias in the distribution of tokens than in that of types. Coupled with the observation that the same holds for adult production — albeit in a weaker form — this finding suggests that the patterns of aspect-tense combinations found in the literature might be a mirroring of adult distributional patterns. Supporting this view is the fact that hardly any of these studies took the diversity of forms into account and thus have mostly confirmed the Aspect Hypothesis for the preferred use of forms, while making a less firm statement about availability of different forms at any stage of development. Given that distributional bias also factors into adult speech, it is important not to overstate the effect of preferred aspect-tense combinations on learnability of forms in the paradigm. Since children are able to pick up on distributional cues, their initial use of forms might simply be a reflection of the distributions found in adults as well as personal needs (cf. Figure 2b and the large proportion of the imperative form of *give*). A similar observation was already made by one of the authors of the Aspect Hypothesis Shirai (1998), who found that Japanese children do not follow the predictions of the Aspect Hypothesis and, therefore, suggested that multiple factors should be taken into account when examining early acquisition of tense-aspect morphology.

By looking at the use of different lemmas with the individual grammatical forms and thus measuring how flexibly a form is used, we were able to show that the development of form use might be more advanced than indicated by preferential use of certain tokens which skew the distributions. Going forward, it is important to disentangle the issue of tense-aspect marking further and take into account the differences between token and type distributions as well as further factors, such as lexical development and underlying distributions of grammatical markers in individual languages.

Acknowledgments

This work was supported by the European Research Council (ERC Consolidator Grant, ACQDIV 615988, to S. Stoll).

References

Aksu-Koç, A. A. (1998). The role of input vs. universal predispositions in the emergence of tense-aspect morphology: evidence from Turkish. *First Language*, 18, 255-280.

Antinucci, F., & Miller, R. (1976). How children talk about what happened. *Journal of Child Language*, 3, 167-189.

Bar-Shalom, E. (2002). Tense and aspect in early child Russian. *Language Acquisition: A Journal of Developmental Linguistics*, 10(4), 321 - 337.

Bloom, L., Lifter, K., & Hafitz, J. (1980). Semantics of verbs and the development of verb inflection in child language. *Language*, 56, 386-412.

Bronckart, J.-P., & Sinclair, H. (1973). Time, tense and aspect. *Cognition: International Journal of Cognitive Psychology*, 2, 107-130.

Clark, E. V. (1996). Early verbs, event-types, and inflections. *Children's Language*, 9, 61-73.

Comrie, B. (1976). *Aspect*. London: Cambridge University Press.

Gagarina, N. (2000). The acquisition of aspectuality by Russian children: the early stages. *ZAS-Papers in Linguistics*, 15, 232-246.

Harner, L. (1981). Children talk about the time and aspect of actions. *Child Development*, 52, 498-506.

Johnson, B. W., & Fey, M. E. (2006). Interaction of lexical and grammatical aspect in toddlers' language. *Journal of Child Language*, 33(02), 419-435.

Li, P. (1990). *Aspect and aktionsart in child Mandarin* (Doctoral Thesis).

Li, P. (2000). The acquisition of lexical and grammatical aspect in a self-organizing feature-map model. In *Proceedings of the 22nd annual meeting of the cognitive science society*. Mahwah, NJ: Lawrence Erlbaum.

Li, P., & Bowerman, M. (1998). The acquisition of lexical and grammatical aspect in Chinese. *Journal of Child Language*, 54, 311-350.

Li, P., & Shirai, Y. (2000). *The acquisition of lexical and grammatical aspect*. Berlin/New York: Mouton de Gruyter.

Lieven, E. V., Pine, J. M., & Baldwin, G. (1997). Lexically-based learning and early grammatical development. *Journal of Child Language*, 24, 187-219.

Lin, J. (1991). Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory*, 37(1), 145-151.

Pine, J. M., & Lieven, E. V. (1997). Slot and frame patterns and the development of the determiner category. *Applied Psycholinguistics*, 18, 123-138.

R Core Team. (2015). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org>

Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27(3), 379-423.

Shirai, Y. (1998). The emergence of tense-aspect morphology in Japanese: universal predisposition? *First Language*, 18(54), 281-309.

Shirai, Y., & Anderson, R. W. (1995). The acquisition of tense-aspect morphology: a prototype account. *Language*, 71, 743-762.

Shirai, Y., Slobin, D. I., & Weist, R. M. (1998). *The acquisition of tense-aspect morphology* (Vol. 18). Alpha Academic.

Stephany, U. (1985). *Aspekt, Tempus und Modalität: Zur Entwicklung der Verbalgrammatik in der neugriechischen*

- Kindersprache. [Aspect, tense, and modality: The development of grammar in young greek children].* Tübingen, Germany: Gunther Narr.
- Stoll, S. (1998). The role of aktionsart in the acquisition of russian aspect. *First Language, 18*, 351-378.
- Stoll, S. (2005). Beginning and end in the acquisition of the russian perfective aspect. *Journal of Child Language, 32*, 805-825.
- Stoll, S., & Gries, S. (2009). How to measure development in corpora? an association-strength approach to characterizing development in corpora. *Journal of Child Language, 36*, 1075-1090.
- Stoll, S., & Meyer, R. (2008). *Audio-visional longitudinal corpus on the acquisition of Russian by 5 children.*
- Timberlake, A. (2004). *A reference grammar of Russian.* Cambridge University Press.
- Tomasello, M. (2000). First steps toward a usage-based theory of language acquisition. *Cognitive Linguistics, 11*(1/2), 61-82.
- Tomasello, M. (2003). *Constructing a language: a usage-based theory of language acquisition.* Harvard, MA: Harvard University Press.
- Weist, R. M., & Konieczna, E. (1985). Affix processing strategies and linguistic systems. *Journal of Child Language, 12*, 27-35.
- Weist, R. M., Wysocka, H., Witkowska-Stadnik, K., Buczowska, E., & Konieczna, E. (1984). The defective tense hypothesis: on the emergence of tense and aspect in child polish. *Journal of Child Language, 11*, 347-374.