# Developing the Global Scale of English Vocabulary for Young Learners (6 to 11)

Dr Veronica Benigno
Prof. John de Jong

—

# Contents

# Purpose of this document

This document reports on an ongoing project to develop the Global Scale of English (GSE) Vocabulary for Young Learners, a graded lexical framework for EFL learners aged 6 to 11. The framework is aligned to the GSE and the Common European Framework of Reference for Languages (CEFR; Council of Europe, 2001). This project follows up on a previous study carried out to identify vocabulary targets of adult learners of English (Benigno & De Jong, 2017, available at https://prodengcom.s3.amazonaws.com/GSE-Vocab.pdf and, together with the *GSE Learning Objectives for Young Learners*, aims to help teachers identify pupils' learning targets at increasing levels of proficiency.

# Executive summary

**The GSE Vocabulary for Young Learners identifies vocabulary targets of EFL learners aged 6 to 11, an age group which roughly corresponds to primary education in many countries in the world. It is our intention to carry out a separate project to identify lexical needs of older children.**

The GSE Vocabulary for Young Learners is a database consisting of about 3,000 lexical units. Each lexical unit is a distinct word meaning - in line with research evidence demonstrating that vocabulary learning gradually develops from basic to complex and that, therefore, not all meanings of a word are learned at once (Wolter, 2009). Take, for example, two different meanings of *can*: our analysis suggests that the verb *can* (to express ability) should be learned before the noun *can* (as in a can of coke). Similarly, the word *bat* features in two distinct meanings. We suggest that the meaning of "a small animal that flies at night" is learned before the meaning of "a long wooden stick used to hit the ball in some games". For each entry in the database, the following information is provided: word definition, topic and/or subtopic (and sometimes sub-subtopic), grammatical category, GSE value/CEFR level, and an example sentence. For example, in the topic "Family and self", subtopic "Family members", entries such as *brother* and *sister* can be found alongside entries like *baby brother/sister*, *big brother/sister*, and an *only child*. In the topic "School", subtopic "In the classroom", sub-subtopic "Classroom instructions", are entries such as  *look* (verb), *choose*, *partner* as well as many useful word combinations often used to interact in the classroom, e.g. *what is it*, *sit down*, etc.

—

# Rationale

**The development of the GSE Vocabulary for Young Learners is informed by research on second language acquisition and vocabulary learning – in order to ensure that the specific needs of this particular target group are met. There is evidence that children learn vocabulary by chunks, preferably associating words with concrete objects, and mainly indirectly, e.g. by exposure to oral language (e.g. listening to adults) and written language (e.g. reading activities) (Cameron, 2001; Wray, 2002). Research on vocabulary size of English learners generally suggests that the 2,000/3,000 vocabulary level is a crucial learning goal for low-level EFL learners (Dale and Chall, 1948; Nation, 2001; Schmitt and Schmitt, 2014; Staehr, 2008). Although there are only a few studies on the vocabulary size of EFL young learners, there seems to be a general agreement that children are able to learn 500 words per year in good learning conditions (Nation, 1990 cited by Cameron, 2001, p.75). As an indicative cut-off, we chose 3,000 vocabulary units as target when we selected how many words primary EFL learners should learn.**

The same key theoretical principles applied in the development of the GSE Vocabulary for Adults were applied in developing the GSE Vocabulary for Young Learners. The framework was created using a mixed methodology which combines corpus frequency data and teacher judgments to identify level-appropriate vocabulary. Whilst frequency is a valid criterion to establish what vocabulary to teach first, frequency data can be biased by the nature of the corpus and therefore be misleading. Equally, it is not currently possible to accurately extract the corpus frequency of different meanings of the same word automatically. In order to compensate for the limitations deriving from a purely frequency-based approach, we asked a group of primary teachers to evaluate the communicative usefulness of each word meaning. We therefore obtained a weighted measure which combines both quantitative (frequency) and qualitative (teacher judgments) criteria to rank lexical items on the GSE and the CEFR. The reader should refer to Benigno & De Jong (2017), available at https://prodengcom. s3.amazonaws.com/GSE-Vocab.pdf for further details on the theoretical framework which informed the methodological choices outlined below.

# Methodology and data

The methodology consisted of four main steps.

### Step 1. Creation of a frequency list from English L1 corpora and materials for children

The first step consisted in the creation of a lexical database relevant to the needs of EFL young learners in the age range 6 to 11.  A lexical database of 3,081 word meanings was compiled by drawing words from two different sources, corpora and ELT materials and wordlists created for children. The purpose of drawing items from the above resources was to add low-frequency words (not retrieved by corpus analysis) which are representative of the language of young learners in this age range. The combination between the corpus-based frequency list and the list of items found in the ELT materials resulted in a list of 5k lemmas which was further filtered/cleaned, resulting in a final list of  3,081 word meanings, which corresponds to the target size of about 3k items to cover relevant vocabulary for primary EFL learners.

The chosen corpora were:

- the British National Corpus (http://www.natcorp.ox.ac.uk/corpus/): Spoken; ~25k tokens; BBC radio broadcasts, school lessons, & spontaneous conversations; age 7 to 11

- the Child Language Exchange Data System –CHILDES- (MacWhinney, 2000): Spoken; ~ 160k tokens; age 3-7; Interview, Free/guided conversation and spontaneous speech

- SUBTLEX-UK (van Heuven et al., 2014), which gathers word frequencies based on subtitles of BBC broadcasts: Spoken; ~ 13m tokens; CBBC broadcasts; text types – comedy, entertainment, drama, animation, news, and factual; age range 6-12

The chosen ELT materials were young learner coursebooks/flashcard lists, the Seward 4K Teaching List (aka Zeno, Ivenz, Millard, and Duvvuri's (1995) Educator's Word Frequency Guide) and the Dale-Chall list (Dale and Chall, 1948) of 3,000 simple words (https://readable.io/blog/the-dale-chall-word-list/).

**Step 2: Semantic annotation of word meanings to assign them to topics and subtopics**

The second methodological step consisted in manually annotating each entry in the database to add the information about which topic, and sometimes subtopic, each word meaning belongs to. In this phase we annotated all word meanings by topics which are relevant to the young learners' context of usage. The starting point for the list of topics is the categorization found in Nikolov (2016), which was compared with our Adult GSE Vocabulary topic tree (accessible via the GSE Teacher Toolkit at https://www.english.com/gse/teacher-toolkit/user/lo) and also with topic names used in young learner course books. The GSE Young Learner Learning Objectives were also consulted to construct the list of topics. The GSE Young Learner Learning Objectives are skills-based, not vocabulary-based. Nevertheless we were able to extract some useful vocabulary areas. The list of topics contains 27 topic headings and a number of subtopics and sub-subtopics.

**Step 3. Rating exercise on word usefulness by primary teachers**

In this phase, teacher ratings were collected to obtain an insight into the usefulness of vocabulary. The term "usefulness" refers to the importance of a word or phrase to a young learner in terms of representativeness of his/her experiential domain, with a focus on "common" communicative acts and topics which are more likely to come up in an EFL/young learner context. For this reason, during the teacher training we asked raters to give particular importance to the vocabulary used in the classroom. The concept of usefulness only partially relates to the concept of ease/difficulty. Some words are inherently difficult, for example because they are long, but they are still learnt very early because they are frequent/useful.

The rater group was made up of 18 EFL primary teachers (10 with English L1 and 8 with a C level of proficiency on the CEFR) teaching in 11 different countries worldwide: Switzerland, Portugal, Greece, Spain, South Korea, China, Hungary, Croatia, Brazil, Russia, and Germany. Each word meaning was rated by 12 teachers (6 having English as their L1 and 6 having a different L1 than English) in an overlapping design. Teachers were asked to evaluate the degree of usefulness of vocabulary items for communication in the context of EFL at primary level. For each word, they were presented with the word definition and the part of speech. In this way they rated words by word meaning, not by word form. They were asked to rate word meanings on a scale of usefulness ranging from 1 to 5 (see table 1 below). They were also given the possibility to indicate whether they didn't know a word or whether a particular word

was not relevant in the context of primary EFL – using a 99. Before starting rating, the teachers took part in an online training/standardization session to learn about the project, the rating scale and the task itself. Following the training, they received a briefing document.

| 1 Essential | "Essential" items are the word meanings/phrases that children would want to acquire first. They are essential for basic communication. |
|---|---|
| 2 Important | "Important" items are word meanings/phrases that become necessary at a next stage; they are still very common. They are perhaps a little more detailed or a little more specific in their meaning. |
| 3 Useful | "Useful" items are word meanings/phrases that expand the children's vocabulary enabling more detailed and specific language use. |
| 4 Nice to have | "Nice to have" word meanings/phrases are for children to express themselves accurately and precisely. |
| 5 Extra | The "extra" category is for word meanings/phrases that some children will use occasionally, but they are not needed for everyday communication. |

Table 1. The rating scale.

## Step 4. Statistical analysis of the frequency data and the teacher ratings and data modelling

The corpus frequency data and the ratings provided by the teachers were then analysed and combined to assign a CEFR level/GSE value to each word meaning. A total of 36,972 teacher ratings were collected and analysed by means of descriptive statistics in order to identify deviant ratings in the data set. Removing misfitting ratings improved the inter-rater correlation from 0.70 to 0.81.Frequency values were rescaled to the same scale as the ratings, i.e. from 1 to 5. This action was performed to make the comparison between the two measures more transparent. The reliability (certainty) of the ratings was calculated and adjusted depending on the number of ratings available for each entry. Next, a combined algorithm which used both frequency and rating information was produced to rank vocabulary from the most to the least useful word meaning. Finally, modelling and regression analyses were carried out to transform the ranking into GSE values, based on vocabulary size estimates available in the literature. The finalized values are published in the GSE Teacher Toolkit at English.com.

# References

Benigno, V. & De Jong, J. (2017). Developing the GSE Vocabulary. Global Scale of English Research Series.
Available at https://prodengcom.s3.amazonaws.com/GSE-Vocab.pdf

Cameron, L. (2001). Teaching languages to young learners. Cambridge University Press, Cambridge

Council of Europe (2001). The Common European Framework of Reference for Languages: learning, teaching, assessment. Cambridge University Press, Cambridge, UK

Dale, E. & Chall, J. (1948). A Formula for Predicting Readability. Educational Research Bulletin. 27, 11–20

MacWhinney, B. (2000). The CHILDES project: Tools for analyzing talk. Lawrence Erlbaum Associates

Nation, I., S., P. (2001). Learning vocabulary in another language. Cambridge University Press, Cambridge

Nikolov, M. (2016). A framework for young EFL learners' diagnostic assessment: Can do statements and task types. In M. Nikolov (Ed.), Assessing young learners of English: Global and local perspectives. Springer, New York

Schmitt, N., & Schmitt, D. (2014). A reassessment of frequency and vocabulary size in L2 vocabulary teaching. In Language Teaching, 47, 484-503

Stæhr, L., S. (2008). Vocabulary size and the skills of listening, reading and writing. In Language Learning Journal, 36(2), 139-152

Van Heuven, W., J., Mandera, P., Keuleers, E., & Brysbaert, M. (2014). SUBTLEX-UK: a new and improved word frequency database for British English. In Quarterly Journal of Experimental Psychology, 67, 1176-1190

Wolter, B. (2009). Meaning-last vocabulary acquisition and collocational Productivity. In Fitzpatrick T. & Barfield A. (Ed.), Lexical Processing in Second Language Learners: Papers and Perspectives in Honour of Paul Meara (Second Language Acquisition), Multilingual Matters, Bristol, 128-140

Wray, A. (2002). Formulaic language and the lexicon. Cambridge University Press, Cambridge, UK

Zeno, S. M., Ivenz, S. H., Millard, R. T., & Duvvuri, R. (1995). The educator's word frequency guide. Brewster, Touchstone Applied Science Associates, NY

Be yourself
in English.

Pearson