

UC Merced

Proceedings of the Annual Meeting of the Cognitive Science Society

Title

Levels of Representation in Language Development

Permalink

<https://escholarship.org/uc/item/6d386703>

Journal

Proceedings of the Annual Meeting of the Cognitive Science Society, 27(27)

ISSN

1069-7977

Authors

Chater, Nick
Hulme, Charles
Monaghan, Padraic

Publication Date

2005

Peer reviewed

Levels of Representation in Language Development

Padraic Monaghan (pjm21@york.ac.uk)

Department of Psychology, University of York
York, YO10 5DD, UK

Nick Chater (N.Chater@warwick.ac.uk)

Department of Psychology, University of Warwick
Coventry, CV4 7AL, UK

Charles Hulme (ch1@york.ac.uk)

Department of Psychology, University of York
York, YO10 5DD, UK

Abstract

In this paper we provide a computational exploration of changes in levels of representation in children's language development. Children demonstrate a progression from awareness of larger units in their language, such as words or syllables, and only later do they develop the ability to manipulate phoneme level representations. We employed a minimum description length approach to encode a corpus of child-directed speech. The analysis indicated that, for early stages of language learning, a word level representation was most efficient, whereas after extended exposure, a phoneme level representation was more efficient. The analysis also accurately predicted which phonemes children would be able to distinguish from syllables at different stages of development.

Introduction

Language contains structure at several different levels of representation, from such large units as discourses, to subphonemic features at the other end of the scale. The child learning her first language will become adept at processing language at each of these structural levels, but how does this knowledge of the units and the relations between them develop?

The child's awareness of structure at a particular level can be assessed by requesting the child to tap out the number of units present in the language stimulus. So, for the word "penguin", the child with awareness of the syllable level of representation would be able to tap twice, to indicate two syllables. For the same stimulus, the child must tap six times when the task is to mark the number of phonemes in the word (/pɛŋ-wɪn/). Alternative methods for determining a child's awareness of structure at the phoneme level are tasks where the child must delete a phoneme from the word (e.g., say "scow" without the "s"), or produce the first sound in the word (so respond /s/ to "scow").

Longitudinal studies suggest that children develop from awareness of larger units to smaller units (Ziegler & Goswami, 2005). Anthony, Lonigan, Driscoll, Phillips, and Burgess (2003) found that, in a study of 1000 children between the ages 2 and 6, there was a developmental progression from word-level representations, to the syllable level, then the onset/rime level, and finally to the phoneme level. Similarly, Liberman, Shankweiler, Fischer, and Carter

(1974) found that 4 year old children could tap out the number of syllables but could not tap out the number of phonemes within a word, but 70% of 6 year olds could manage the phoneme tapping task. In terms of segmenting continuous speech into units, Jusczyk, Cutler, and Redanz (1993) proposed that children learn to segment language into successively smaller units. So, at one stage in language development, children learn to segment utterances into words, and at later stages, they learn to segment into sub-word units. Brent and Cartwright (1996) provided a computational model of segmentation that showed a similar progression from larger units to smaller, word-like units of representation.

Awareness of structure at the phoneme level has attracted considerable interest because this ability is a strong predictor of reading success (Hulme, 2002; Muter, Hulme, Snowling & Stevenson, 2004.). Indeed, it has sometimes been suggested that phoneme awareness may be a consequence of learning to read (Bruce, 1964; Ehri & Wilce, 1980; Goswami & Bryant, 1990; Seymour & Evans, 1994; Treiman, 1983; Tunmer & Nesdale, 1985). In adults, too, phoneme awareness was much poorer in illiterates than a literate control group (Morais, Cary, Alegria, & Bertelson, 1979).

Yet, there is evidence to suggest that phoneme awareness, though a reduced ability, is present in children and adults who have not yet learned to read. The illiterate adults in Morais et al.'s (1979) study were poorer but could still perform the task to a certain degree. Hulme, Caravolas, Málková, and Brigstocke (in press) showed that the child's awareness of a particular phoneme did not depend on knowing the letter corresponding to that sound, as the child could isolate /p/ from the spoken nonword "pag", without being able to identify and name the letter "p". Also, some of the children in the study could not name any letters and yet could isolate some of the phonemes from spoken words. Durgunoglu and Oney (1999) found that pre-literate Turkish children could tap out the phonemes in a word to a high degree of accuracy. So, it seems that, while learning to read may speed up and motivate phoneme awareness, some ability to manipulate phonemes is present in pre-literate children and adults (Hulme et al., in press).

So, why is there this progression from large to small units in the child's development? Turkish children develop

phoneme awareness relatively earlier than children learning some other languages, and Durgunoglu and Oney (2002) suggest this is because morphological properties of the language are realized at the phoneme level, such as vowel harmony for pluralisation, thus requiring a phoneme level representation. Similarly, Caravolas and Bruck (1993) discovered that Czech children found isolating phonemes in consonant clusters easier than English children. They suggested this was due to Czech children being more aware of phonemes within consonant clusters due to the greater number of consonant clusters in Czech.

Ziegler and Goswami (2005) performed a corpus study of English, German, French, and Dutch, all languages where the onset/rime distinction can be made by children. They found that there were more neighbours for syllables (“cow”) when the onset and rime were distinguished (“c”, “ow”) than when neighbours for the onset-vowel and coda were determined (“co”, “w”). They suggested that this neighbourhood property would encourage representation at the onset/rime division, in contrast to an onset-vowel/coda division.

An alternative view is that as the child’s vocabulary grows, the distinctions the child must make between words in the lexicon are more fine-grained as the neighbourhood of a particular word becomes more dense. So, the child must learn to represent subsyllabic properties of the word in order to identify it correctly (Charles-Luce & Luce, 1995; Walley, 1993).

In this paper we take a rational analysis approach to the development of levels of language representation, and provide a complementary explanation for the development of levels of language processing. We assume that an efficient representation of the language is one that is most likely to be used by the child. We employ a minimum description length (MDL) approach to determine the level of representation that is most efficient for encoding the language after different quantities of exposure to the language.

We focus on levels of language representation determinable in a corpus of child-directed speech, with the utterance as the largest unit, then the word, the syllable, and then sub-syllabic levels of onset/rhyme, and the phoneme as the smallest unit we consider here. We provide three tests of the model. First, we explore which level of language representation is most efficient for different exposure to the language. Second, we report the efficiency of encoding *new* syllables after different exposure to the language. Third, we test the model’s predictions for the ease with which certain phonemes can be isolated from syllables.

We next describe the MDL method as a means for assessing efficient encoding at different levels of language representation.

Minimum Description Length

We selected MDL as a method for reflecting the efficiency of encoding language at different levels of granularity. The MDL approach contends that the best theory to explain a set

of data is the one that minimizes the sum of (a) the length in bits of the description of the theory; and (b) the length in bits of the data when encoded with the theory (Rissanen, 1982). In the case of encoding a language, the MDL approach measures the amount of information in bits required to describe the language, given a particular level of representation. For an inefficient level of representation for the language, the amount of information required to describe the language will be high, whereas for a more optimal level of representation for the language, the number of bits in the description will be lower.

As mentioned above there are two parts to the MDL encoding. First, the theory of the language is implemented as a list of the units used in the encoding. So, for a phoneme-level encoding of the language corpus, all the phoneme types that occur in the corpus are listed. For a word-level encoding, the list comprises all the words in the corpus. The cost in bits of describing each element is related to the frequency of each element in the list. This provides a lower-bound for the cost of encoding the list of units, because the learner does not know the frequencies in advance and hence cannot optimize the code to those frequencies. Items in the list that occur very frequently can be encoded more efficiently than items that occur infrequently. For this measure, we used Shannon’s noiseless coding theorem, where $b_i = \log_2(1/p_i)$, where b_i is cost in bits for describing unit i , and p_i is the probability of occurrence of unit i . If the unit has high probability of occurrence, the length in bits is short, but length is long if the probability of a unit occurring is very low. Rare events have low expectancy and so encoding of them cannot be as efficient as common events.

The second part of the MDL encoding is the description of the corpus. The cost in bits of this encoding is the number of units in the corpus at the chosen level of representation. So, for a phoneme-level description, it is the number of phonemes in the corpus and for a word-level description, it is the number of words in the corpus.

Encoding the language using large units, such as utterances, results in a large cost for describing the list, as each distinct utterance has to be encoded separately, and the probability of occurrence of each utterance is low as there are many unique utterances in a language corpus. However, the cost for encoding the corpus is then very low, as the unit size is large. Encoding using smaller units, such as phonemes, results in a short list of high-frequency units, but a larger cost in encoding the corpus as each utterance contains several phonemes. An efficient level of description therefore entails finding the optimal trade-off between encoding both the list and the corpus.

MDL approaches have been widely used in language processing research. Ellison (1992), for example, illustrated how a MDL approach could result in the discovery of phonotactic constraints in Turkish, Goldsmith (2001) provided a MDL algorithm for discovering morphology, and Grünwald (1994) showed how aspects of grammar could be successfully inferred using a MDL algorithm. We assume

that the child attempts to derive an encoding of the language that is as efficient as possible, and will have available a level of representation that minimizes the code length of the description.

We made three predictions about the MDL approach. First, we hypothesized that early in development, i.e., for a small language corpus, larger units will be more efficient for encoding the language than smaller units. This is because the length of encoding the list of units will be small for both utterance-level encoding as well as a phoneme-level encoding as the diversity of utterances will not be very great. However, the cost of encoding the language will be smaller for the utterance-level as the number of units in the corpus will be much smaller. Later in development, i.e., for a large language corpus, representing the language in terms of smaller units may be more efficient. This is because the size of the list encoding will be substantially greater for larger unit encoding because there will be many different utterance types each with low probability of occurrence, sufficient to counteract the smaller description length of the corpus when using larger units. This hypothesis was tested in Study 1. The second hypothesis was that encoding *new* information at the syllable level would be most efficient after small language exposure, whereas encoding at the phoneme level would be more efficient after more exposure. We tested this in Study 2. Finally, we hypothesized that phoneme awareness would be easier for certain phonemes and harder for others, and that the model's predictions should match a developmental profile. We tested the extent to which the model could reflect this detailed data in Study 3.

Study 1: Encoding the whole language

Method

Corpus preparation All the sentences spoken by adults in the English component of the CHILDES (MacWhinney, 2000) database were selected for the analysis. Sentences were cross-classified with CELEX (Baayen, Pipenbrock, & Gulikers, 1995) to gain phonological transcriptions for each word. Several words were not found in the CELEX database, and utterances containing these words were omitted. This resulted in approximately 14.8 million phonemes, in 5.6 million words, in 1.2 million utterances. For the following analyses we used only the first 1 million utterances from the corpus.

Corpus analysis We encoded the corpus at 5 levels of representation: (1) the utterance level, where each utterance was coded as one unit; (2) the word-level, where words were separately coded; (3) the syllable-level, where syllables were coded as separate units; (4) the onset/rime-level, where the onset of each syllable was encoded separately from the rime; and (5) the phoneme-level.

In order to test at different stages of language exposure, we selected the first n utterances from the corpus and

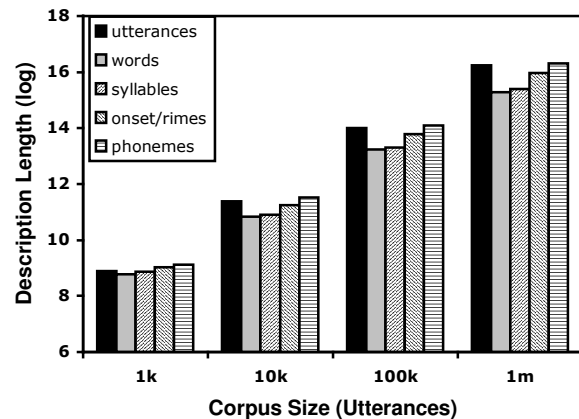


Figure 1. MDL of unigram encoding at 5 different levels of granularity, for increasing corpus size.

assessed the description length of this fragment of the corpus for each of the 5 levels of representation.

Results and discussion

The results for corpus sizes 1000, 10000, 100000, and 1 million utterances are shown in Figure 1. For the smallest corpus, describing the language at the word level is most efficient, with decreasing efficiency as unit size becomes smaller. The word level is also more efficient than the utterance level, which had approximately the same efficiency as the syllable level encoding. As corpus size increases, the word level representation remains the most efficient encoding, and the same ordering of efficiency is found for decreasing unit sizes: syllables are more efficient than onset/rime which is more efficient than phoneme level representations. However, the largest unit size considered – that of the utterance – becomes less efficient as the corpus increases, until at 1 million utterances, it is less efficient than both the syllable and the onset/rime encodings.

The results did not support our hypotheses about levels of encoding and length of exposure to the language. We did find that large units were more efficient for smaller corpora, but, with the exception of the utterance level, this pattern did not alter as corpus size increased. The benefit of describing a small list of distinct phonemes did not override the extra length required to describe the corpus at the phoneme level.

Yet, describing the whole language may not be the most important task facing the child learning the language. Instead, efficiently encoding novel language information may be a processing priority. The next Study explores which level of encoding is most efficient for adding new words to the language.

Study 2: Encoding new language information

Method

Corpus preparation We used the same corpus as for Study 1.

Materials We took 20 single-syllable nonwords used in Study 2 of Hulme et al. (in press), used to test phoneme awareness in children. There were 10 sets of 2 syllables beginning with each of the phonemes /p/, /b/, /d/, /f/, /l/, /m/, /n/, /t/, /v/, and /z/.

Corpus analysis We encoded the language at each of the five levels of representation as in Study 1, for varying corpus sizes. After encoding the language corpus, we measured the efficiency of encoding each of the 20 nonwords given the encoding at each level of representation.

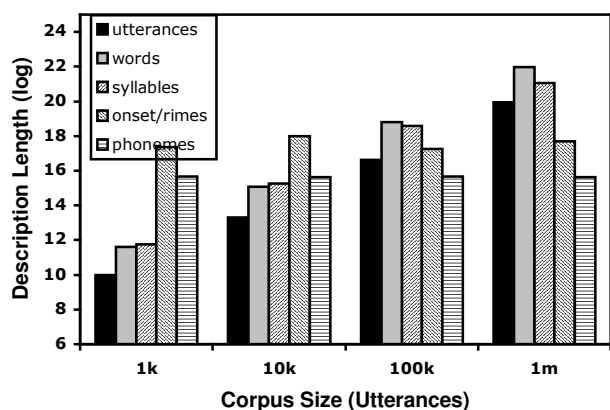


Figure 2. Cost of encoding new words (nonwords) after different exposure to the language at 5 different levels of granularity.

Results and discussion

The results for 1000, 10000, 100000, and 1 million utterances taken from the corpus are shown in Figure 2. For the smallest corpus size, encoding the new information at the utterance level is most efficient, followed by the word level, then the syllable level, and then the phoneme level. Surprisingly, the onset/rime encoding is least efficient for this corpus size. The results support the hypothesis that, when exposure to the language is limited, encoding at a large unit size is most efficient.

As the corpus size increases, the efficiency of encoding new information at larger unit sizes becomes less efficient. After 10000 utterances, the word and syllable levels are equivalent to the phoneme level representation for efficiency, though the utterance level is most efficient. For the 100000 utterance corpus, the phoneme level becomes the most efficient, with the onset/rime level also more efficient than the word and syllable level representations. For the largest corpus, the phoneme level representation is the most efficient encoding.

The results fit the hypothesis extremely closely. For a small corpus, the small list of the large unit encoding meant that a new word added to the language presents only a small cost for encoding the new information at the same level. However, for a large corpus, the novel stimulus has very

low expectancy in the utterance, word and syllable level representations, as it is an entirely novel item. For the onset/rime and phoneme level representations, however, the encoding can exploit the prior occurrence of tokens of elements of the word in the corpus. So, for the word “pag”, /p/ and /æg/ have occurred as onsets and rimes in the corpus many times, so do not need to be added to the description of the onset/rime encoding. Similarly, /p/, /æ/, and /g/ occur in the phoneme level description so require no additional description length in the theory.

The MDL approach to reflecting efficiency of language encoding at different levels of representation raises specific hypotheses about which phonemes children should find easiest and hardest to extract from syllables. Phonemes that are more frequent, and syllables that are less frequent, should provide the best chance for a phoneme level representation to be optimal. Phonemes that are less frequent, and syllables that are highly frequent will reduce the effectiveness of smaller unit size encoding. The next Study tests this prediction.

Study 3: Phoneme awareness for different phonemes

Method

Corpus analysis We used the same corpus as for Study 1.

Materials We used the 20 single-syllable nonwords from Hulme et al. (in press).

Procedure We reanalyzed the child data from Hulme et al. to determine the order in which children could isolate particular phonemes from syllables. So, if the child could isolate only one phoneme then this was recorded as the first phoneme learned. If the child could isolate two phonemes then these were recorded as the first two phonemes. Across the 77 children from the study, we determined whether there was a consistent order of isolating phonemes.

For the model data, we rated which syllables were most efficiently encoded at the phoneme level, and produced a rank to predict the order in which phonemes would be isolated, assuming that the phonemes learned earliest would be those that were most efficiently encoded at the phoneme level.

Results and discussion

We generated an order of acquisition of phoneme awareness for the phonemes. This was achieved by assessing for each child the phonemes that they managed to produce and those that they did not. We found that this was consistent for every child, and an order of acquisition was determined:

/b/ /d/ < /n/ < /l/ /t/ < /f/ /p/ < /m/ /v/ < /z/

where “x < y” indicates that for every child phoneme awareness for phoneme x is acquired before phoneme y.

When two phonemes occur together, as for /b/ and /d/ this means that the data did not determine for which of these phonemes was phoneme awareness prior: all children either isolated both /b/ and /d/, or neither phoneme.

For the MDL data the likelihood of phoneme awareness for a particular phoneme was taken to be the advantage for an encoding at the onset/rime level or the phoneme level over the syllable level of analysis. The predicted order from the onset/rime level encoding was:

/l/ < /n/ < /d/ < /b/ < /m/ < /f/ < /p/ < /v/ < /z/ < /t/

For the phoneme level encoding the predicted order was:

/n/ < /z/ < /d/ < /b/ < /l/ < /t/ < /m/ < /f/ < /v/ < /p/

The correlation between the order predicted by the onset/rime level encoding and the children's data was $rho = .417$, $p = .230$. For the phoneme level encoding the correlation was $rho = .772$, $p = .021$. For the syllable level, each phoneme had equal cost of encoding as it comprised part of an undivided syllable each of which had equal probability of occurrence, and so equal code length. This was also true for the word and utterance level representations.

We also tested whether order of awareness of individual phonemes was related purely to the frequency of phonemes. This differed from the MDL phoneme level approach in that, in the latter analysis, all phonemes in the nonword are considered. For the phoneme frequency count, the predicted order did not result in a significant correlation with the children's order, $rho = .460$, $p = .181$. The MDL phoneme level encoding therefore provided the best fit to the children's data, accounting for 51% of the variance in response.

General Discussion

Children develop awareness of the structure of their language starting with larger units, such as words or syllables, and gradually define smaller unit representations, such as onsets/rimes and finally phonemes. There is evidence that some phoneme awareness may occur prior to literacy, and hence the child's ability to manipulate phonemes within words is not entirely due to learning to read. This paper attempted to address the issue of why the development occurs from large to small representational units. We have shown that a model that takes into account the efficiency of encoding at different levels of representation mimics this development for the learning of new material. This developmental pattern in the MDL analyses was due to the length of exposure to the language. For small corpora, larger chunks are more efficient, whereas for large corpora, finer-grained distinctions begin to pay off, despite the cost of counting several elements within a single syllable.

The MDL analyses are not inconsistent with other theories about why phoneme awareness develops later in children, and why it may occur at different times for different languages (Caravolas & Bruck, 1993; Ziegler & Goswami, 2005). Our analyses provide a mechanism for

determining the efficiency of encoding at different levels, and consequently provide predictions for the ease of processing of different speech stimuli for children.

In particular, the MDL analyses generated the prediction that certain phonemes would be more easily isolable from syllables than others, and that this depended not only on the frequency of the phoneme to be isolated but also the frequency of the other phonemes within the syllable. The reanalysis of 77 children's data indicated that there was a clear pattern of development for which phonemes could be isolated, and that this pattern was best accounted for by the MDL analysis with encoding at the phoneme level.

There were points where our hypotheses were not supported. Our first hypothesis was that encoding the whole language that the child experiences is most efficiently accomplished by small unit representations after extended exposure. This was not found to be the case in the results of Study 1, where the phoneme level representation was the least efficient, due to the length of the entire corpus when encoded as a list of phoneme tokens. Encoding this list at a word level provided the most efficient and effective, which is consonant with Zipf's assessments of natural language (Zipf, 1935).

The current analyses treat each unit size as equal in terms of its cost for encoding. A more realistic assumption is that larger units require more bits for encoding. An analysis that takes this differential cost related to unit size may produce results that more closely map onto the child's development. Another improvement to the MDL analyses we have presented would be to take account of multiple, simultaneous levels of representation of language that the child presumably possesses. Discovering structure at the phoneme level does not preclude using structure at other levels of description, and combining levels is likely to provide the most efficient encoding. Certain words are therefore better encoded as words, but others, particularly novel words as seen in Study 2, are better encoded at a subsyllabic level. Gobet, Lane, Croker, Cheng, Jones, Oliver, and Pine (2001) have undertaken modeling using such principles of multiple representational levels in language acquisition.

In addition, we have ignored the possibility that each level of language representation is equally available as a possible level of encoding of the language. Yet, variation in the acoustic-phonetic realization of individual phonemes may make an important contribution to the availability of this level (McClelland, 2004), and contribute to the development of phoneme-awareness skills. The results of the analyses are an indication of how much specific child development data can be modeled without referring to such additional factors.

We have provided a first attempt at a computational account of development of awareness of structure at different levels of granularity of representation. The progression from large units to small units in the child's awareness of their language is mirrored in the MDL approach. Critical in the account of this progression is that

the levels of representation available to the child are a consequence of whether the level provides an efficient level of encoding of the language. The MDL analyses we performed were on up to 1 million words, which is an underestimate of the quantity of language a child may hear – the MDL approach operates in a noise-free environment and is an optimal learner which is likely not the case for the child. It is the general pattern of development, and not precisely when it occurs, that is important. However, the model we have presented predicts that limited exposure to speech will result in poorer phoneme awareness skills – a potential contributory factor to phonological ability, and hence reading ability, in children.

References

- Anthony, J.L., Lonigan, C.J., Driscoll, K., Phillips, B.M., & Burgess, S.R. (2003). Phonological sensitivity: A quasi-parallel progression of word structure units and cognitive operations. *Reading Research Quarterly*, 38, 470-487.
- Baayen, R.H., Pipenbrock, R. & Gulikers, L. (1995). *The CELEX Lexical Database* (CD-ROM). Linguistic Data Consortium, University of Pennsylvania, Philadelphia, PA.
- Brent, M.R. & Cartwright, T.A. (1996). Distributional regularity and phonotactic constraints are useful for segmentation. *Cognition*, 61, 93-125.
- Bruce, D.J. (1964). The analysis of word sounds. *British Journal of Educational Psychology*, 34, 158-170.
- Caravolas, M. & Bruck, M. (1993). The effect of oral and written language input on children's phonological awareness: A cross-linguistic study. *Journal of Experimental Child Psychology*, 55, 1-30.
- Charles-Luce, J. & Luce, P.A. (1995). An examination of similarity neighbourhoods in young children's receptive vocabularies. *Journal of Child Language*, 22, 727-735.
- Durgunoglu, A.Y. & Oney, B. (1999). A cross-linguistic comparison of phonological awareness and word recognition. *Reading and Writing*, 11, 281-299.
- Durgunoglu, A.Y. & Oney, B. (2002). Phonological awareness in literacy development: It's not only for children. *Scientific Studies of Reading*, 6, 245-266.
- Ehri, L.C. & Wilce, L.S. (1980). The influence of orthograph on readers' conceptualization of the phonemic structure of words. *Applied Psycholinguistics*, 1, 371-385.
- Ellison, T.M. (1992). *The Machine Learning of Phonological Structure*. Unpublished PhD thesis, University of Western Australia.
- Gobet, F., Lane, P.C.R., Croker, S., Cheng, P.C-H., Jones, G., Oliver, I., & Pine, J.M. (2001). Chunking mechanisms in human learning. *Trends in Cognitive Sciences*, 5, 236-243.
- Goldsmith, J. (2001). Unsupervised learning of the morphology of a natural language. *Computational Linguistics*, 27, 153-189.
- Goswami, U. and Bryant, P. (1990). *Phonological Skills and Learning to Read*. Hillsdale, NJ: Lawrence Erlbaum.
- Grünwald, P. (1994). A minimum description length approach to grammar inference. In G. Scheler, S. Wermter, & E. Riloff (Eds.), *Connectionist, Statistical and Symbolic Approaches to Learning for Natural Language*, pp. 203-216. Berlin: Springer Verlag.
- Hulme, C. (2002). Phonemes, rimes and the mechanisms of early reading development. *Journal of Experimental Child Psychology*, 82, 58-64.
- Hulme, C., Caravolas, M., Málková, G., & Brigstocke, S. (in press). Phoneme isolation ability is not simply a consequence of letter-sound knowledge. *Cognition*.
- Jusczyk, P.W., Cutler, A., & Redanz, N. (1993). Preference for the predominant stress patterns of English words. *Child Development*, 64, 675-687.
- Lieberman, I., Shankweiler, D., Fischer, F.W. and Carter, B. (1974). Explicit syllable and phoneme segmentation in the young child. *Journal of Experimental Child Psychology*, 18, 201-212.
- MacWhinney, B. (2000). *The CHILDES Project: Tools for Analyzing Talk*, Third Edition. Mahwah, NJ: Lawrence Erlbaum Associates.
- McClelland, J.L. (2004). Phonology without phonemes. *Proceedings of the Twenty-Sixth Annual Conference of the Cognitive Science Society*, p.29. Mahwah, NJ: Erlbaum.
- Morais, J., Cary, L., Alegria, J., & Bertelson, P. (1979). Does awareness of speech as a sequence of phones arise spontaneously? *Cognition*, 7, 323-331.
- Muter, V., Hulme, C., Snowling, M.J., & Stevenson, J. (2004). Phonemes, rimes and language skills as foundations of early reading development: Evidence from a longitudinal study. *Developmental Psychology*, 40, 665-681.
- Rissanen, J. (1982). A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11, 416-431.
- Seymour, P.H.K. & Evans, H. (1994) Sources of constraint and individual variations in normal and impaired spelling. In: G.D.A. Brown and N.C. Ellis (eds.), *Handbook of Spelling: Theory, Process and Intervention*, pp. 129-154. Chichester: John Wiley and Sons.
- Shannon, C. E., & Weaver, W. (1949). *The Mathematical Theory of Communication*. Urbana IL: University of Illinois Press.
- Treiman, R. (1983). The structure of spoken syllables: Evidence from novel word games. *Cognition*, 15, 49-74.
- Tunmer, W.E. & Nesdale, A.R. (1985). Phonemic segmentation skill and beginning reading. *Journal of Educational Psychology*, 77, 417-427.
- Walley, A. (1993). The role of vocabulary development in children's spoken word recognition and segmentation abilities. *Developmental Review*, 13, 286-350.
- Ziegler, J.C. & Goswami, U.C. (2005). Reading acquisition, developmental dyslexia and skilled reading across languages: A psycholinguistic grain size theory. *Psychological Bulletin*, in press.
- Zipf, G.K. (1935). *Psycho-Biology of Languages*. Cambridge, MA: MIT Press.