

Lexicographic Profiling: an Aid to Consistency in Dictionary Entry Design

Sue Atkins

Lexicography MasterClass Ltd.

River House

Church Lane

Lewes

East Sussex BN7 2JA

UK

Valerie Grundy

Lexicographic Consultant

Bellevue

07440 Boffres

France

Abstract

The first phase of any new dictionary project includes the detailed design of the dictionary entries. This involves writing sample entries. Deciding which lemmas to write entries for tends to be random, for lack of theory-based criteria for selection. This paper describes a database of the lexicographic properties of English lemmas, reflecting the work of theorists such as Apresjan, Cruse, Fillmore, Lakoff, Levin and Mel'Ëuk. It was used successfully in the selection process during the planning phase of a major new English bilingual dictionary, ensuring that all major policy decisions on entry structure were made early in the project, and resulting in a DTD well able to cope with 50,000 entries and a Style Guide already comprehensive at the start of the main project. The method described here ("lexicographic profiling") is applicable to any language, although of course the actual properties of lemmas are to some extent language specific.

1 The background

The paper describes a systematic approach to the microstructure design for a totally new corpus-based bilingual dictionary, tried and tested in Phase 1 of the *New English Irish Dictionary* (NEID) project.¹ The following packages were completed on time and within the rather modest budget in the course of the one-year first phase of this project:

¹ This project was launched in the summer of 2003. The contract for Phase 1 of the project was awarded to the Lexicography MasterClass Ltd. The project was directed by Sue Atkins, Adam Kilgarriff, and Michael Rundell, with Valerie Grundy as Managing Editor.

1. a design for a customised lexicographical corpus;
2. a 225 million word corpus of English (Hiberno-, British and American), and a 30 million word corpus of Irish, all linguistically annotated to a high level;
3. software: a corpus query system with the corpora loaded;
4. a user profile, a set of headword selection principles and a headword list;
5. a list of linguistic labels for marking register, style, domain, variety etc.;
6. 100 sample bilingual entries (English-Irish) covering the full range of entry types;
7. 50 'template' (model) entries, for certain lexical sets;
8. a detailed description of the proposed entry structures needed for the dictionary;
9. a document type definition (DTD);²
10. software: a dictionary production system with the DTD loaded;
11. a comprehensive style guide (lexicographers' manual);
12. a functioning web-based 'reading-and-marking' programme;
13. software: project-management tools (scheduling, textflow, budgets etc.);
14. a Linguistic Advisory Board consisting of 30 leading linguists, to advise on dictionary policy and comment on dictionary sample text;
15. a business plan for Phase 2 (dictionary production), with detailed schedules and budgets.

The work on lexicographic profiling described in this paper contributed substantially to tasks 6-11 inclusive of the above list.

2 How we set about it

For the dictionary writing software to serve the project correctly, the major policy decisions on entry structure, content and layout must be taken before the DTD is finalized. Similarly, the Style Guide must be comprehensively drafted, although of course this document is inevitably fine-tuned during the first year or so of a dictionary project. Our aim therefore in selecting the headwords for the 100 sample entries which would inform the DTD and Style Guide was to choose words which would bring out, as far as possible, the whole spectrum of problems, situations, and combinations likely to be encountered by the lexicographers and translators working on Phase 2 of the project.

2.1 Identifying lexicographic problems

In the course of our work as senior dictionary editors we had already compiled a database of words that had proved problematic in other dictionaries, together with the specific problems they raised. Lemmas were grouped according to part of speech, and within these groups there were a number of sub-categories relevant to that part of speech. Each of the problems in our original list had been analysed in detail, resulting in a relational database in which were recorded, for each lemma, the lexicographic properties that would require policy deci-

² The commercial dictionary production software, supplied by IDM Paris (www.idm.fr), required customization for the project: in particular of course the DTD had to be written, defining the structure of the various types of entry.

sions for the dictionary macro- and microstructure. Table 1 shows some examples of the kinds of issue we recorded at this point.

headword	issues
box	homograph issues: whether or not to make different entries with homograph numbers for the different core meanings of the lemma ('container', 'sport', 'botanical' etc.)
box	whether to organize the senses according to part of speech or according to meaning; e.g. should <i>box</i> the noun, meaning 'container', be located close to <i>box</i> the verb meaning 'put into a container' and separate from the other noun <i>box</i> (the plant), or should all the noun senses be grouped, then all the verb senses?
speck	an instance of a type of word that always raises problems in bilingual dictionaries: when used as itemizers, ³ such words often give rise to a plethora of translations according to the noun itemized (<i>speck of dust, speck of mud, specks of yellow</i> etc.)
echo	organization of multiple parts of speech (noun, modifier, verb), including representatives of main subtypes of verb (intransitive, transitive, report verb) and wide range of directional PPs.
echo	participation in multiword expressions (<i>cheer to the echo, echo chamber, echo back</i> etc.): where these should be shown (as headwords, nested headwords, within the entry in a special section etc.)
little	how to deal with not just the adjective but parts of speech that are notoriously awkward to pin down in a bilingual entry (broad negative quantifier/determiner)
might	as a modal verb, frequently raises problems of target-language equivalence: where no direct equivalent exists, particular lexicographical strategies need to be developed (specially adapted entry layout, discursive notes etc.)

Table 1. Some of the issues noted in our survey

2.2 Classifying the problems

The structured data constituted a database of *lexicographic properties of lemmas* (henceforth the lemma properties database, or LPD), i.e. those properties that have an immediate impact on lexicographic decisions. We included in this database concepts from linguists such as Apresjan, Cruse, Fillmore, Lakoff, Levin and Mel'Ëuk, whose work in lexical semantics gave us a clearer view of what we had to do. An example is the use we made of Apresjan's seminal study of regular polysemy (see 4.2.2).

This allowed us to produce lexicographic profiles of words which we were considering as candidates for the 100 sample entries, and to select the final 100 headwords on a systematic basis, ensuring that most of the problems the dictionary team would come across in Phase 2 of the project had been tackled and solved in this preliminary stage.

3 The problems: principal microstructure decisions

The main microstructure issues on which decisions had to be taken are summarised in

³ This is Mel'Ëuk's 'Sing' lexical function; see Mel'Ëuk *et al* (1984/1988/1992)

this section: although our particular work was establishing the infrastructure for an English-Irish dictionary with a detailed user profile (and therefore focussing only on English), the decisions discussed here are relevant to every bilingual dictionary. We provided the publisher with seven batches of sample entries. Each batch focussed on one of the key issues listed below, and consisted of 12–15 entries and a discussion document offering alternative ways of handling the data and proposing our preferred method.⁴

3.1 High-level structural division of lexical entries

This is the first decision regarding the overall microstructure organization: should the top-level organizing principle for grouping the various *lexical units*⁵ of a lemma be (1) meaning or (2) part of speech? This is a decision which must be made very early in the planning phase because of its implications for the DTD.

3.2 Treatment of multi-word expressions and derived forms

How should multiword expressions (MWEs: essentially, compounds, phrasal verbs, idioms) and derived forms be dealt with within the microstructure and what principles could be established for deciding which types of MWE (if any) should be given headword status. An early decision is needed on this point too, in order to establish a basic DTD.

3.3 Treatment of linguistic labels in dictionary and database

The actual types of label (domain, register, style, variety etc.) and the specific labels to be made available to the lexicographers were established separately. This batch of sample entries focussed on the circumstances in which a label should be used – in the initial source-language analysis, in the translation material, and in the final version of the dictionary entry – and those where some other type of indicating material would be more appropriate. The scope of the label, both in SL and in TL text, was also discussed.

3.4 Treatment of Irish regional variants in dictionary and database

The issue here was to determine which regional vocabulary and syntactic variants should be included in the target-language material, how these should be labelled (if at all), and – in the absence of an accepted ‘standard’ Irish – whether there should be any prioritizing distinctions made among the various regional equivalents.

3.5 Function word entries: content, style and layout

The main decisions here related to the words to be considered as ‘function words’ (in particular, the grey area of English adverbial and prepositional particles), and how the *functional* equivalents in the target language should be set out and exemplified. The methodology de-

⁴ Final decisions on all of these points were the responsibility of the publishing institution, Foras na Gaeilge.

⁵ This term is used as defined in Cruse (1986): essentially, a lemma in one of its meanings.

vised for dealing lexicographically with *semantic* equivalence is not very helpful in such cases.

3.6 Miscellaneous style issues

This batch of entries dealt with problems not addressed elsewhere. Each of them impacted less critically on DTD design than those discussed in 3.1 – 3.5, but nonetheless had to be accounted for in the Style Guide and DTD. They included dealing with abbreviations and acronyms, productive prefixes and suffixes, and contracted forms, as well as showing verbs with obligatory adjuncts, and resolving the various problems raised by nouns functioning as itemizers.

3.7 Issues relating to Irish equivalence

The issues raised in this batch of entries are those which cause problems in every bilingual dictionary: source-language (SL) items with no direct target-language (TL) translation; those with partial semantic congruence, where either the SL item or the TL item is more specific than the other; instances where the denotational meaning of SL and TL items matched but they diverged along the axes of register, style, or pragmatic force etc.; the approach to finding equivalents of culture-bound encyclopedic items such as *Downing Street*; and in what circumstances a simple TL paraphrase of the SL item should be acceptable.

4 Analysing the problems: the lemma properties database

The database is designed to be stored in Microsoft Access or a similar database package and to be queried on the basis of property type, actual property, or instantiating lemma. Still incomplete, it contains 27 classes of lexicographic properties, which account for 581 actual properties.

4.1 Classes of lexicographic properties

While some of these properties relate to the lemma (headword) itself, the great majority relate to specific parts of speech and therefore to the lexical unit (LU) rather than the lemma. The principal properties recorded about the lemma itself were:

- the various wordclasses of its LUs;
- its lexical form (simple word, prefix, abbreviation, hyphenated, multiword etc.);
- its corpus profile (corpus frequency, preferred text type, regional variety etc.);
- the labels needed for its lexicographic description (domain, register etc.);
- miscellaneous properties (monosemous, polysemous, with non-homophone homograph etc.);
- its participation in multiword expressions;
- its participation in regular polysemy types (see 4.2.2);
- its participation in a lexical set for which a model ('template') entry was envisaged.

The remainder of the properties in the database were recorded individually for nouns, verbs, adjectives, adverbs and function words and focussed on:

- noun, verb etc. types (common/proper nouns; lexical, auxiliary, copular verbs and subtypes of these; gradable, ungradable etc. adjectives and adverbs);
- semantic classes (human, artifact etc. nouns; motion, sound etc. verbs; colour, sound, existential etc. adjectives; manner, degree, time etc. adverbs);
- morphological properties, both inflectional and derivational;
- syntactic properties (nouns that modify other nouns, nouns postmodified by an adverbial phrase, as in *the journey home*; the wide range of syntactic behaviour, especially complementation, of verbs, adjectives and adverbs);
- participation in Levin-type alternations (for verbs and nouns);
- some other miscellaneous properties for each part of speech.

4.2 Some examples of lexicographic properties

Examples of two of these classes follow: *lexical form of the lemma (headword)* and *regular polysemy participant*.

4.2.1 Properties relating to the lexical form of the headword

The lexical form – or variations on this – of the lemma raises issues of content and layout which impact on both macrostructure and microstructure. Some examples are shown in Table 2.

property	e.g.	at headword	issue
simple word	<i>book, give, France</i>		(standard headword)
abbreviation (initial letters, pronounced as letters)	<i>the EC, the BBC, WMD</i>	European, British, weapon	Where will the full information be given? Will there be cross-references? If the target language has also an abbreviated form, will the full form of the target language be included and if so where? (as for abbreviation above)
acronym (initial letters, pronounced as a word)	<i>Unesco, Nato</i>	united, north	
prefix (productive)	<i>anti-war, ex-wife</i>	anti-, ex- ...	How to handle semantically rich prefixes which may be attached to a wide range of words.
suffix (productive)	<i>talkfest, additive-free, work-shy</i>	-fest, -free, -shy	How to handle similarly productive suffixes.
multiword lemmas	<i>all right, in spite of</i>	all, right, spite ...	Should there be multiword headwords? If so, define criteria for headword status. Decide how (if at all) to show links with other headwords.
hyphenated word	<i>broad-leafed, flop-eared</i>	broad, leaf, flop, ear	How to handle compounds (hyphenated or non-hyphenated) – as headwords? Secondary headwords? Defined? Not defined? etc.
has variant spelling	<i>aluminium /aluminum</i>	aluminium	Should both versions have headword status or should one variant be given within the entry for the other? If two headwords, where should the information be given?

Table 2. Issues raised by the lexical form of the lemma

4.2.2 Properties relating to regular polysemy relationships

This class of properties is very large (163 items): it was compiled initially on the basis of Apresjan's work (see Apresjan 1973), and expanded as other instances of the phenomenon surfaced in our analysis. This property does not generate lemma-specific problems. However, the identification of types of regular polysemy allows the Style Guide to include instructions on how to deal with headwords which participate in each particular type of relationship, making the overall approach to this aspect of the lexicography much more consistent. Table 3 contains a very short extract showing some types of regular polysemy.

property	instantiating lemmas	e.g.
container n->contents n	<i>box, tin, can, pocket, bag, phial, case</i>	<ul style="list-style-type: none"> • <i>six boxes on the floor</i> • <i>he ate the whole box himself</i>
container n->amount n	<i>box, glass, phial, case</i>	<ul style="list-style-type: none"> • <i>six sherry glasses</i> • <i>add half a glass of wine</i>
container n->to place in container vt	<i>pocket, box, can, tin</i>	<ul style="list-style-type: none"> • <i>she had some money in her pocket</i> • <i>he pocketed the change</i>
dance n->dance-music n	<i>tango, quickstep, reel</i>	<ul style="list-style-type: none"> • <i>I'm learning the tango</i> • <i>they played several tangos</i>
dance n->do that dance vt	<i>tango, quickstep, foxtrot ..</i>	<ul style="list-style-type: none"> • <i>I'm learning the tango</i> • <i>they tangoed round the room</i>
fruit of plant n->plant n	<i>raspberry, apple, orange, strawberry</i>	<ul style="list-style-type: none"> • <i>a bowl of raspberries</i> • <i>a field of raspberries</i>
tree n->its wood n	<i>oak, cedar, elm, pine, mahogany ..</i>	<ul style="list-style-type: none"> • <i>three old pines behind the house</i> • <i>table made of pine</i>
tree / wood n->of the colour of that wood n mod	<i>mahogany, walnut, oak, yew</i>	<ul style="list-style-type: none"> • <i>made of mahogany</i> • <i>the rich mahogany of her hair</i>
animal n->its fur or skin n	<i>mink, squirrel, fox, leopard, crocodile</i>	<ul style="list-style-type: none"> • <i>a lake full of crocodiles</i> • <i>a crocodile handbag</i>

Table 3. Some types of regular polysemy recorded in the database

5 Solving the problems: lexicographic profiles

A list was drawn up of over 400 lemmas known to be problematic, each lemma was checked against the appropriate property classes and its lexicographic profile was extracted. Our aim was to collect words as disparate as possible with regard to all aspects of inherent properties and corpus use, thus ensuring policy decisions on as many as possible points of macrostructure and microstructure. One example suffices, that of *echo*.

5.1 Case study: lexicographic profile of *echo*

Table 4 contains the lexicographic profile of the lemma *echo*, drawn from our database. The order in which the properties are listed has no significance. These profiles help to ensure that the words chosen as headwords of sample entries cover as many issues as possible. The implications for dictionary entry structure of most of these properties are for the most part transparent, but an explanatory note has been included where there could be some confusion.

	property	of lemma <i>echo</i>	note on policy decision
1	lexical form: simple word		
2	corpus profile ⁶ : raw frequency: lemma	30.65 per million words	Frequency and ranking of lemmas and LUs can affect the depth of treatment, particularly in learners' dictionaries.
3	corpus profile: frequency rank: noun	5089	<i>echo</i> is the 5089 th most frequent noun in the corpus (see note above)
4	corpus profile: frequency rank: verb	3581	<i>echo</i> is the 3581 st most frequent verb in the corpus (see note above)
5	belongs to a 'word family'	<i>echoing</i> adjective (<i>echoing footsteps</i>)	Morphologically related headwords can influence treatment of a lemma.
6	polysemous lemma	several dictionary senses e.g. 1) <i>his voice echoed</i> 2) <i>I echo these sentiments</i>	Relates to sense ordering, also hierarchical or non-hierarchical sense numbering etc.
7	p-o-s = noun	<i>Can you hear the echo?</i>	
8	p-o-s = verb	<i>His voice echoed down the hall.</i>	
9	noun subtype: count noun	... <i>sending out sounds and listening to the echoes</i>	Subtypes of noun relate to decisions on p-o-s markers in dictionary. e.g. <i>ne, ni</i> ...
10	noun subtype: uncount noun	... <i>using string sections and lots of echo to make records that appealed to older listeners.</i>	See note at 9.
11	noun syntax: (head of NP) subject of VP	<i>The unexpected echo startled us all.</i>	Making sure all syntactic possibilities are covered: some nouns rarely or never occur in this position, and this must be noted where relevant.
12	noun syntax: (head of NP) object of VP	<i>We heard a loud echo.</i>	See note above at 11.
13	noun: semantic frame ⁷ : sound	<i>We heard a loud echo.</i>	This is a concept from frame semantics; see Atkins <i>et al</i> (2003) and Fillmore <i>et al</i> (2003). Noting the semantics helps to ensure that the sample headwords are as disparate as possible, semantically.
14	noun: semantic frame: becoming aware	<i>It was an echo of former happiness.</i>	See note at 13.
15	noun: morphology: plural in -es	<i>echoes</i>	How to show source language plurals in the dictionary.
16	noun: morphology: verb-derived zero nominal	1) <i>it will echo</i> 2) <i>hear an echo</i>	See Levin (1993). Check if relationship holds across semantically distinct LUs. Policy on how to handle this if so and if not.

⁶ Statistics taken from the British National Corpus (www.natcorp.ox.ac.uk/).

⁷ "Semantic frames are schematic representations of situation types (eating, spying, removing, classifying, etc.) together with lists of the kinds of participants, props, and other conceptual roles that are seen as components of such situations." FAQs on FrameNet website, www.icsi.berkeley.edu/~framenet/. See also Fontenelle (2003).

	property	of lemma <i>echo</i>	note on policy decision
17	noun is element in MWE	<i>echo chamber</i> <i>to cheer to the echo</i>	
18	in MWE type: idiom: morpho-syntactic flexible	<i>cheered/cheering to the echo</i>	How to show flexibility (or not) of MWEs.
19	in MWE type: compound noun	<i>echo chamber</i> , plural <i>echo chambers</i>	How to handle plural of compound nouns: here the plural marked on 2 nd component, but some compound nouns (e.g. <i>court martial</i>) mark plural on 1 st component.
20	noun: metaphorical extension	... <i>uncanny echoes of past events</i> .	How to handle: productive use of metaphor in language, yielding uses which are not 'set' enough to be treated as dictionary senses. See Lakoff & Johnson (1980).
21	verb subtype: intransitive	<i>The shout echoed across the valley.</i>	Subtypes of verb relate to decisions on p-u-s markers in dictionary: e.g. vt, vi, vii, ... How to deal with transitivity/ intransitivity in entry layout.
22	verb subtype: transitive	<i>He echoes this view.</i>	See note above at 21.
23	verb subtype: reporting verb	"Because of the play?" he <i>echoed</i> mockingly.	See note above at 21.
24	verb semantic frame: sound	<i>The sound echoed round the room.</i>	See note at 13.
25	verb: semantic frame: communication	<i>I am happy to echo those sentiments.</i>	See note at 13.
26	verb: morphology: weak vb -es, -ing, -ed	<i>echoes, echoed, echoing</i>	How to show verb inflections, which type to include, and which to consider default.
27	verb: metaphorical extension	<i>We painted the beams green to echo the colour of the furniture.</i>	How to handle: metaphorical meaning extensions which are not 'set' enough to be treated as dictionary senses. See Lakoff & Johnson (1980).
28	verb: in alternation: swarm: X verbs in Y → Y verbs with X	1) <i>their voices echoed in the hall</i> 2) <i>the hall echoed with their voices</i>	See Levin (1993).
29	verb is element in MWE type: phrasal verb - V+ADV intransitive	<i>The sound ricocheted off the walls and <u>echoed back</u>.</i>	How to handle intransitive verb + particle phrasal units.
30	verb is element in MWE type: phrasal verb - V+ADV transitive	<i>The stage but <u>echoes back</u> the public voice.</i>	How to handle transitive verb + particle phrasal units.
31	label type; domain: telecommunications	<i>There's an <u>echo</u> on the line.</i>	A labelling policy is essential, together with a list of agreed labels and how they will be used.

Table 4. Lexicographic profile of *echo*

6 Applying lexicographic profiles

The lexicographic profiles of the 100 headwords selected for the sample entries generated principled decisions on the type of information to be included in the dictionary entry, how that information was to be shown, how different types of entry were to be structured and how the information types fitted together to make up that structure. From there, we were able to go on and establish the DTD and write the Style Guide.

6.1 The DTD

The DTD defines the structure of the different types of entry in the dictionary, and the use of the lemma properties database in the selection of sample entries enabled us to write the DTD with a high level of confidence that no major issue had been overlooked.

Every aspect of the lexicography of each headword was thoroughly studied: for instance, our analysis of different multiword expressions led us to conclude that – in this one-volume print dictionary – we needed to retain a good deal of flexibility with regard to the status of compounds. Thus, in the DTD, a Compound Container can contain essentially the same structure as a full lexical Dictionary Entry, allowing lexicographers at the final stage of editing⁸ to decide (on the basis of the Style Guide) whether a specific compound should be given headword status or shown within the body of an entry, as this DTD extract⁹ shows:

```
<!ELEMENT DEnt HwdGp, (SenBlk | FwkMWEBLk | PhrVBLk | UsgNoteCnt | XRefCnt | FwkSenCnt | RegVarNoteCnt | COMMENT)* >  
<!ELEMENT CpdCnt (CpdGp, (SenBlk | FwkSenCnt | UsgNoteCnt | XRefCnt | COMMENT)* ) >
```

6.2 The Style Guide

On the basis of a comprehensive sampling of headwords participating in multiword constructions, the Style Guide contained quite explicit guidelines on how to handle these, as is shown in the following extract from the manual used by editors compiling the bilingual entry:

14.8 Compounds - nouns

In the framework, all compounds associated with a headword are contained in the CPDBlk within the MWEBLk. You should consider each carefully. All those that merit full compound-headword status according to the principles discussed above should be promoted to full headword entries and treated separately as such (CpdEnt). Any that do not should be subsumed in the modifier sense within the main entry.

⁸ The three stages of editing were (1) initial source-language analysis, recorded in a database ("the framework"); (2) target-language translations added to database; (3) extraction from this material of the polished bilingual entry.

⁹ The abbreviations CPDBlk etc. are labels used to define elements of entry structure in the DTD and to signpost them in the dictionary editing software.

7 Conclusion

The lemma properties database provided a systematic approach to policy decision-making in a new corpus-based dictionary project. By applying the ideas of theoretical linguists, and using our many years of lexicographic experience, we ensured that most difficulties which would be encountered during the course of dictionary editing had been carefully considered and instructions on how to handle them encoded into the style guide, and that the software was running on a comprehensive and flexible DTD.

References

- Apresjan, J. D. (1973), 'Regular Polysemy', *Linguistics* 124, pp. 5-39.
- Atkins, B. T. S., C. Fillmore, J., Johnson, C. R. (2003), 'Lexicographic relevance: selecting information from corpus evidence', in *International Journal of Lexicography*, Oxford, OUP: 16:3, pp. 251-280.
- Cruse, D. A. (1986), *Lexical Semantics*, Cambridge University Press, Cambridge, UK.
- Fillmore, C. J., Johnson, C.R., Petruck, M. (2003), 'Background to FrameNet', in *International Journal of Lexicography*, Oxford, OUP: 16:3, pp. 235-250.
- Fontenelle, Th. (2003) (guest editor), *International Journal of Lexicography*, Oxford, OUP: 16:3 (issue devoted to FrameNet Project).
- Lakoff, G., Johnson, M. (1980), *Metaphors We Live By* University of Chicago Press, Chicago, USA.
- Levin, B. (1993), *English Verb Classes and Alternations: A Preliminary Investigation*, University of Chicago Press, Chicago, USA.
- Mel'žuk, I., Arbatchewsky, N., Dagenais, L., Elnitsky, L., Iordanskaja, L., Lefebvre, M. N., Mantha, S., Lessard, A. (1984/1988/1992), *Dictionnaire Explicatif et Combinatoire du Français Contemporain: Recherches Lexico-Sémantiques I, II, III*, Montréal, Les Presses de l'Université de Montréal.