

UNNC-CAWSE

University of Nottingham Ningbo China

Corpus of Academic Written and Spoken Corpus



TRANSCRIPTION CONVENTIONS

For the Spoken Subcorpus [v.1]

December 2017

Yu-Hua Chen & Qianqian Zhou

UNNC

Ningbo, China

Transcription Conventions (with Preliminary Annotation) for Spoken Data v.1

Overview

1. Interactional	Speaker Turns, Latching, Overlaps
2. Verbal	Acronyms, Capitalisation, Code Switching, Contractions, Fillers, Lengthening, Lexicalised Reduced Forms, Numbers & Dates, Orthography & Hyphenation, Punctuation, Repetition, Unintelligible Speech, Truncation
3. Vocal (non-verbal)	Exhalation/Inhalation, Laughing, Pauses
4. Non-vocal	Non-vocal Communicative (NVC) Events
5. Others	Anonymisation, Time Stamps, Deviations

Type (subtype in alphabetical order)	Definition	Example
1. Interactional		
1.1 Speaker Turns	Speaker turns are indicated by individual speaker IDs followed by a colon ':' and utterances. Each speaker is labeled with their unique ID in the format of t/s/r plus four digits of Arabic numbers: 't' for teacher, 's' for student, 'r' for researcher and 'x' for unknown speaker.	t0001: <i>right okay thank you that's the end of the test</i> s0001: <i>okay</i> t0001: <i>you're free to go</i>
1.2 Latching	The lack of a pause between different speakers is marked by '=' in the utterances: one at the end of the first speaker and the other at the beginning of another where latching occurs.	s0001: <i>so er: I think from you two=</i> s0002: <i>=OK</i>
1.3 Overlaps	The tagset and are used for overlaps, and two pairs are used where the overlap occurs: one with the first speaker, and the other with the second speaker.	s0005: <i>oh oh oh</i> s0002: <i>you know this one is much sweeter than this one</i>
2. Verbal		
2.1 Acronyms	2.1.1 If an acronym is pronounced as a sequence of letters, it is transcribed as a sequence of capital letters separated by spaces.	<i>I'll now ask you some general questions about U N N C life OK</i>
	2.1.2 If an acronym is uttered as a word, it is transcribed as a sequence of capital letters without any spaces.	<i>I'm doing a TESOL</i>
2.2 Capitalisation	Capital letters are NOT used at the beginning of sentences. They are only retained as required in spelling conventions such as	<i>so er: I think from you two</i>

Transcription Conventions (with Preliminary Annotation) for Spoken Data v.1

	proper nouns (e.g. 'New York', 'Ningbo'), first personal pronoun 'I', or 'Mr', 'Mrs', 'Dr', etc.	
2.3 Code Switching	Utterances in Chinese are marked up by the tagset <cs n="zh"> and </cs>. The language code is specified after the attribute "n", and in this case the code for Chinese is "zh". Translations into English are provided wherever possible in curly brackets {} after the utterance of code switching for those who may not understand Chinese. If it is neither English nor Chinese, the language code can be found in the Library of Congress (US): https://www.loc.gov/standards/iso639-2/php/code_list.php	<cs n="zh">那个地壳是不是{is that earth crust}</cs> or maybe on some fa- some face <cs n="zh">表情怎么说啊{how to say 'facial expression'}</cs>
2.4 Contractions	Standard spelling of all contractions are retained.	I'm she's you'd we've
2.5 Fillers/Filled Pauses	All filled pauses are standardized in orthography and marked as one of the following: <i>ah, en, er, erm, huh, mhm, mm, oh, uh, or uhu.</i> No other fillers are used.	s0001: mm: when I get up early
2.6 Lengthening	Lengthened sounds are represented by the symbol of a colon ':'. s0008: (2.3) hh (1.4) just er: ok animals human s0001: oh: oh:	
2.7 Lexicalised Reduced Forms	When a lexicalised reduced form is uttered, the original shortened form (e.g. <i>gonna</i>) is rendered as opposed to a full standard form.	cos kinda gonna gotta wanna
2.8 Numbers & Dates	Numbers and dates are written out in full words in the same way of how they are uttered.	nineteen ninety nine (rather than 1999)
2.9 Orthography & Hyphenation	2.9.1 Both hyphenation and spelling follows the rules of British English, and the online Oxford Dictionary is consulted: (https://en.oxforddictionaries.com/). For example, for the suffix of <i>-ise/-ize</i> , only the <i>-ise</i> variant is used rather than the <i>-ize</i> variant although both morphemes are used in British English. Other common instances include suffixes such as <i>-our</i> (e.g. <i>colour</i>), <i>-re</i> (e.g. <i>centre</i>).	
	2.9.2 In the case of more than one spelling variations available for a word in the Oxford Dictionary, the first variation is chosen unless it involves British/American spelling. For	OK all right we'll now move on to part two

Transcription Conventions (with Preliminary Annotation) for Spoken Data v.1

	example, the spellings of <i>OK</i> and <i>all right</i> are standardised (as opposed to, for example, <i>okay</i> or <i>alright</i>).	
2.10 Punctuation	No punctuation is used for sentence or clause boundaries.	<i>okay right now I'm gonna move on to part two within part two I want you to give a short speech on a topic that I will give you</i>
2.11 Repetition	All repetitions of utterances are transcribed.	<i>sometimes er: er the first thing I I will consider is the quality of the er: product</i>
2.12 Unintelligible	Unintelligible speech is represented by <ut>x</ut> regardless of word length. The number of 'x' approximates the number of words heard. For example, <ut>xxx</ut> indicates three words.	<i>can you can you discuss the article</i> <ut>xx</ut>
2.13 Word Fragments/ Truncation	A hyphen is used to mark word fragments (such as self-interruption or self-correction), which indicates a truncation at the beginning or the end of a word.	<i>=wha- what you mean biology</i>
3. Vocal (non-verbal)		
3.1 Exhalation/ Inhalation	Noticeable breathing is marked by 'hh'.	<i>I think (1.3) it depends hh er: they can the school test can test something</i>
3.2 Laughing	3.2.1 If the current speaker is laughing, use <laughing> in the utterances where the laughter occurs.	<i>so confused <laughing></i>
	3.2.2 If more than one person is laughing, speaker codes are added in the tag to specify who is laughing.	<i><s0001 and t0001 laughing></i>
	3.2.3 If everybody is laughing, use <all laughing>.	<i><all laughing></i>
3.3 Pauses	A pause equivalent to or longer than one second is represented by the duration of time within parentheses (), correct to one decimal place.	<i>hh (1.4) just er: ok animals human</i>
4. Non-vocal		
4.1 Non-vocal Communicative (NVC) Events	Three frequent NVC events are noted: <writing>, <nodding> and <shaking head>. No other NVC event is indicated.	<i>s0005: <nodding></i>
5. Others		
5.1 Anonymisation	Anonymisation: The tagset <anm>x</anm> is used to anonymise the data for the sake of ethics, e.g. names of students, staff, family, friends, or any other sensitive information. The number of 'x' approximates the number of words heard. If it is unclear whether it is sensitive information, the rule of thumb for the transcribers is to simply anonymise the	<i>t0001: great what is the South African student's name</i> <i>s0001: er <anm>xx</anm></i>

Transcription Conventions (with Preliminary Annotation) for Spoken Data v.1

	<p>names.</p> <p>However, since the data is collected from UNNC, UNNC and its associated proper nouns such as building names do not need to be anonymised.</p>	
5.2 Time Stamps	5.2.1 For classroom video recordings, time stamps for every ten minutes are added.	<10 mins> <20 mins>
	5.2.2 For audio recordings of speaking assessment, preparation times are edited out, but the duration of time is indicated in the tag <prep time XX min(s)>.	<prep time 1 min> <prep time 2 mins>
5.3. Deviation	<ul style="list-style-type: none"> Two types of deviations are covered here: lexical and pronunciation (including a slip of the tongue which may otherwise be considered transcribing errors). 1) Lexical errors in relation to forms (i.e. formal misselection and misformation) and semantics (i.e. collocation/word choice) and 2) pronunciation errors which may hinder comprehension are preliminarily tagged and will be illustrated below. Note that this is not intended to be used for grammatical errors (e.g. subject-verb agreement). 	
	<p>5.3.1 If 1) there is a clear deviation (in terms of pronunciation or lexis) and 2) both the intended and deviation words are clear in the context, the tagset <dvp/dvl/dv>DEVIATION{CORRECTION}</dvp/dvl/dv> is used. The corrected word is provided in the curly brackets {} following the deviation, unless the speaker corrects it themselves immediately.</p> <p>The attribute 'dvp' is used for pronunciation deviation while 'dvl' refers to lexical deviation. When it is unclear which category the deviation falls into or the deviation may involve both pronunciation and lexical levels, then 'dv' is used.</p>	<p>Pronunciation deviation (with an issue of intelligibility) <dvp>: <i>it said that our <dvp>hurt {heart}</dvp> er have four chambers</i></p> <p>Lexical deviation <dvl> (1) Slip of the tongue: <i>the description of the movement of the water between ocean er</i> <dvl>background{underground}</dvl> <i>water and er atmosphere</i></p> <p>(2) Part-of-speech (POS) <i>er: some music bands will er: maybe</i> <dvl>creative{create}</dvl> <i>some music</i></p> <p>(3) Derivation <i>I stay in the library</i> <dvl>alonely{alone}</dvl></p> <p>(4) Inflection <i>if you have so many</i> <dvl>childrens{children}</dvl><i>in your in family</i></p> <p>(5) Collocation/ word choice: <i>I think sometimes we buy a</i> <dvl>purchase{product}</dvl> <i>not only for its good quality</i></p>

Transcription Conventions (with Preliminary Annotation) for Spoken Data v.1

		<p>Mixed/uncertain deviation <dv>: <i>maybe it's very bad for the <dv>economic{economy}</dv> er: to the country</i></p>
	<p>5.3.2 If the deviation of utterance is unclear/un-transcribable or the transcriber is uncertain about the intended correct word, 'x' is used. The number of 'x' approximates the number of words heard, e.g. <dv>xx{CORRECTION}</dv> or <dv>DEVIATION{x}</dv>.</p>	<p>Pronunciation deviation <dvp>: <i><dvp>x{osmosis}</dvp> is a kinds of maybe a material and which only let's er water to go through it</i></p> <p>Lexical deviation <dvl>: <i>I think art is a <dvl>virus{x}</dvl> it can include traditional art and also the art er graffiti</i></p>