

Available Resources and type of annotation

LDC corpora

<u>LDC96S36</u>	Boston University Radio Speech Corpus: <ul style="list-style-type: none"> ➤ orthographic transcription, phonetic alignments (TIMIT phonetic labeling system; Arpabet), part-of-speech tags and prosodic markers
<u>LDC2002S28</u>	Emotional Prosody Speech and Transcripts: <ul style="list-style-type: none"> ➤ emotionally transcribed
<u>LDC2003S06</u>	Santa Barbara Corpus of Spoken American English Part-II: <ul style="list-style-type: none"> ➤ some limited prosodic annotation, lots of data on conversation type, speaker background, etc.

In addition we have tons of LDC corpora that are transcribed + speech or only speech but no additional annotation.

Other corpora

- Part of Switchboard (TIMIT phonetic labeling system; Arpabet)
 - ~9k sentences (60k words)
- Similar Chinese corpus (but rather phonemic transcription)
- TIMIT
 - 6300 sentences from 8 English dialects
 - phonetically transcribed (TIMIT phonetic labeling system; Arpabet)
- VerbMobil I and II (German, English, Japanese)
 - We don't have the sound files, but a lot annotated files
 - Superimposed Speech - SUP
 - Phonetic Segmentation PhonDat - PHO
 - Phonetic Segmentation (SAM-PA phonetic system) - SAP
 - Automatic Segmentation (SAM-PA phonetic system) - MAU
 - Word Segmentation - WOR
 - Dialogact Segmentation - DAS
 - Prosodic Segmentation - PRB
 - Symbolic prosodic Segmentation - PRS
 - Signal-based Prosodic accents labeling - LBP
 - Signal-based Prosodic boundaries labeling - LBG
 - Syntactic-prosodic labeling - PRO
 - Syntactic trees - SYN,FUN,LEX
 - Parts of Speech
 - Phonetic Segmentation
 - Segmentation in turns/sentences/chunks/etc
 - SmartKom Transliteration, Gesture Labeling, User State Labeling holistic, User State Labeling by mimic expression, User State Labeling Occlusions, Meta Linguistic Features
 - Translation - TLN

- The London-Lund Corpus of Spoken English (part of the ICAME)
 - consists of 100 texts which are ToBI-labeled. ~ 500,000 words.
 - Phonetically annotated:

	BOOK	SLIP	TAPE	
GRAPHICS	full-time	full-time	045	Hyphen ¹⁷
	that's	that's	096	Apostrophe
	à	à	097 052	eg <i>à la maison</i>
	ä	ä	097 053	eg <i>Händel</i>
	é	é	101 049	eg <i>donné</i>
	è	è	101 050	eg <i>très</i>
	ê	ê	101 051	eg <i>être</i>
	ë	ë	101 052	eg <i>Citroën</i>
	ö	ö	111 049	eg <i>Höllgarten</i>
	ü	ü	117 049	eg <i>Dürer</i>
ç	ç	099 053	eg <i>français</i>	
PHONETICS ¹⁸	[ə]	[ə]	064	
	[tʃ]	[tʃ]	116 115 104	
	[dʒ]	[dʒ]	100 122 104	
	[θ]	[θ]	116 104	
	[ð]	[ð]	100 104	
	[ʃ]	[ʃ]	115 104	
	[ʒ]	[ʒ]	122 104	
	[ŋ]	[ŋ]	110 103	
	[ɑ]	[i], [ɪ]	126	Voiceless clicks
	[œ]	[oe]	111 101	eg in Fr. <i>peur</i>
	[ʔ]	[ʔ]	063	Glottal stop
	[ã]	[ã]	097 049	Nasal

- 9 urban dialects from Britain are collected in the IViE-corpus
 - Contains about 36 hours of recordings.
 - Very small subset is annotated for prosody and prominence.
 - F0-tier
- RNC – Corpus of German radio news
 - ~ 160 news stories
 - orthographically transliterated, words tier, morphosyntactically annotated
 - automatically word aligned
 - manually prosodically labeled, full ToBI labelling
 - phone transcription (also available as syllable-based transcription):
 - 0.570000 122 <P>
 - 0.700000 122 d
 - 0.840000 122 e:
 - 0.870000 122 R

TIMIT

Word label (. wrd):

```
7470 11362 she
11362 16000 had
15420 17503 your
17503 23360 dark
23360 28360 suit
28360 30960 in
30960 36971 greasy
```

Phonetic label (. phn):

(Note: beginning and ending silence regions are marked with h#)

```
0 7470 h#
7470 9840 sh
9840 11362 iy
11362 12908 hv
12908 14760 ae
14760 15420 dcl
15420 16000 jh
16000 17503 axr
```

Switchboard

