

A Multi-criteria Text Selection Approach for Building a Speech Corpus

Chiragkumar Patel^(✉) and Sunil Kumar Kopparapu

TCS Innovation Labs - Mumbai, Thane (West) 400601, Maharashtra, India
{patel.chiragkumar,sunilkumar.kopparapu}@tcs.com
<http://www.tcs.com>

Abstract. Speech corpus is an important and primary requirement for several speech tasks. Building a speech corpora is a lengthy, time consuming and expensive process, it typically involves collection of a large set of textual utterances and then selective distribution of these text utterances among a set of speakers, called speaker sheets. These speaker sheets are articulated by speakers to generate the speech corpora. Depending on the task at hand the speech corpora needs to satisfy certain criteria; For example, a phonetically balanced speech corpora is essential for building an automatic speech recognition (ASR) engine, while for a text dependent speaker recognition engine there is a need for several spoken repetition of the same text by several speakers. In this paper, we formulate a method that enables creation of speaker sheets from a predetermined set of text utterances such that the speech corpora satisfies the desired requirement.

Keywords: Speech corpora · Speaker sheet generation · Optimization

1 Introduction

A speech corpus is a collection of speech audio files and their text transcripts. Speech corpora find use in building speech based solutions; the most common use being to build acoustic models for automatic speech recognition (ASR) purpose. The traditional approach to build a speech corpus (example SPEECON [1]) is to construct text speaker sheets which satisfy some desired criteria; recruited speaker in turn speak the text utterance to generate the speech audio data. The process of creating speaker sheets, generally picked up from a repository of textual utterances, satisfying a certain criteria is one of the important steps in building a speech corpora. In this paper, we address this problem of building speaker sheets so that the speech corpus developed satisfies multiple criteria.

Specifically, the problem that we are addressing can be stated as, given a set $\mathcal{U} = \{U_1, U_2, U_3, \dots, U_N\}$ of N utterances in a language \mathcal{L} having K phonemes denoted by $P = \{P_1, P_2, P_3, \dots, P_K\}$, create m sets $S_1, S_2, S_3, \dots, S_m$, each having p utterances, such that $S_i = \{S_{i1}, S_{i2}, \dots, S_{ip}\}$ and $S_{ij} \in \mathcal{U}$ and $S_i \subset \mathcal{U}$. Note that $\mathcal{S} = \bigcup_{i=1}^m S_i$ is the generated speech corpus. Both \mathcal{S} and $\{S_i\}_{i=1}^m$ need

to satisfy some criteria jointly or individually depending on the requirement. For example, the criterion could be that all the K phonemes occur in $\bigcup_{i=1}^m S_i$ the same number of times (phonetically balanced) as is required to build a speech corpus for building an ASR or $\{S_i\}_{i=1}^m$ be such that it can be used for text dependent speaker identification, namely, $S_1 = S_2 = \dots = S_m$.

In this paper, we propose a novel optimization approach which allows us to construct \mathcal{S} given \mathcal{U} . The approach is based on construction of multiple cost functions which when minimized generates a \mathcal{S} with desired requirement. The rest of the paper is organized as follows, in Section 2 we review related literature; we identify criteria that is required to build speech corpora for a particular kind of speech analysis in Section 3 and discuss our approach in detail. Experimental validation of the proposed approach is discussed in Section 4 and we conclude in Section 5.

2 Literature Review

Speech corpora development is by and far restricted to that of building a phonetically balanced corpora for ASR applications (for example, [2–6]). The general rule of thumb is that the more distributed the available training textual data, the better the utility of the data to enable building automatic speech recognition (ASR) systems. For example, in [7] a method for selecting training data from text databases is discussed for the task of syllabification. A proposal to choose data uniformly according to the distribution of some target speech unit (phoneme, word or character etc.) is discussed in [8]. They show that it is possible to select a highly informative subset of data that produces recognition performance comparable to a system that makes use of a much larger amount of data. Their experiments negate the common belief that there is no data like more data.

Optimal selection of speech data for ASR systems is proposed in [9]. They propose a method for selecting a limited set of maximally information rich speech data from a larger speech database for ASR training. It uses principal component analysis (PCA) to map the variance of speech database into a low-dimensional space, followed by clustering and a selection technique. A rapid method for optimal text selection is discussed in [10] and propose an implementation of a faster version of an iterative greedy algorithm. Using diphone as the basic unit their selection criteria is to maximize the diphone coverage. In [11], with the aim of developing a Bengali speech corpus for a phone recognizer, they use an optimal text selection technique. They maximize the less frequent phones and minimize more frequent phones with minimum text. As can be observed, the criteria for building a speech corpora is majorly defined by the phonetic balance to automatic speech recognition (ASR). In this paper, we propose an approach which enables creation of a speech corpora by generating speaker sheets which can be used for different speech application, including ASR.

3 Proposed Approach

As mentioned earlier, assume that we are given a set $\mathcal{U} = \{U_1, U_2, U_3, \dots, U_N\}$ of N utterances in a language \mathcal{L} having K phonemes denoted by $P = \{P_1, P_2, P_3, \dots, P_K\}$. Let α_{ij} denote the total number of occurrences of the phoneme P_j in the utterance U_i . Observe that

$$\#P_j(\mathcal{U}) = \sum_{i=1}^N \alpha_{ij} \quad (1)$$

denotes the total number of phoneme P_j in the set \mathcal{U} . Note that $\sum_{i=1}^N \#P_j(U_i) = \#P_j(\mathcal{U})$.

Say we are required to create m sets $S_1, S_2, S_3, \dots, S_m$ (speaker sheets) such that each speaker sheet S_i contains p utterances, namely, $S_i = \{S_{i1}, S_{i2}, \dots, S_{ip}\}$.

Additionally, both $\bigcup_{i=1}^m S_i$ ($\stackrel{def}{=} \mathcal{S}$) and $\{S_i\}_{i=1}^m$ satisfy the criteria for a speech recognition application; then \mathcal{S} should be phonetically balanced, namely,

$$\#P_1(\mathcal{S}) = \#P_2(\mathcal{S}) = \#P_3(\mathcal{S}) = \dots = \#P_K(\mathcal{S})$$

which implies that all the K phones in the corpora \mathcal{S} occur equal number of times. One of the known methods adopted is to construct

$$f_i = \sum_{j=1}^K \overbrace{\frac{1}{\#P_j(\mathcal{U})}}^{w_j} \#P_j(U_i) \quad (2)$$

for each utterance $i = 1, 2, 3, \dots, N$. Note that w_j is inversely proportional to $\#P_j(\mathcal{U})$ implying that if a phoneme j occurs more frequently in \mathcal{U} compared to a phoneme l , then $w_l > w_j$. Subsequently, an utterance with higher number of rare phonemes will result in a higher value of f_i score. It is immediately clear that the utterance with higher f_i score must occur more number of times in \mathcal{S} so as to enable phonetic balance of \mathcal{S} .

One of the approaches to build a phonetically balanced speaker sheet set \mathcal{S} is to first sort the N utterances ($\in \mathcal{U}$) in the descending order of their f_i scores and select a value k (where $1 < k < N$) and partition the sorted N utterances into two sets; the top k utterances (\mathcal{U}_t) and the bottom $(N - k)$ utterances (\mathcal{U}_b). Note that $\mathcal{U} = \mathcal{U}_t \cup \mathcal{U}_b$; note that the set \mathcal{U}_t will have most of the rare phonemes.

If every speaker sheet S_i contains p utterances, a percentage $\gamma_p\% = \left(\frac{\gamma}{p}\right) \times 100$ of utterances can be chosen from the set \mathcal{U}_t and the rest, namely, $(100 - \gamma_p\%)$ can be selected from the utterances set \mathcal{U}_b . A good choice of \mathcal{U}_t and $\gamma_p\%$ will ensure that \mathcal{S} , represented by $\mathcal{S}(\mathcal{U}_t, \gamma_p\%)$ has the desired property (say, phonetically balanced). We now formulate the desired criteria that \mathcal{S} needs to satisfy,

C_0 A measure of phonetically balanced corpus would be to compute

$\mathcal{P} = \{\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%))\}_{k=1}^K$ and find

$$C_0(\mathcal{U}_t, \gamma_p\%) = \frac{1}{K} \sum_{k=1}^K (\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)) - \bar{\mathcal{P}})^2 \quad (3)$$

where $\bar{\mathcal{P}} = \frac{1}{K} \sum_{k=1}^K (\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)))$ is the mean. Note that $C_0(\mathcal{U}_t, \gamma_p\%)$ is the variance of \mathcal{P} . The configuration $(\mathcal{U}_t, \gamma_p\%)$ for which (3) is minimum is desired and gives the best phonetically balanced \mathcal{S} .

However, the phonetically balanced corpus is not the only desired criteria on $\{S_i\}_{i=1}^m$ or \mathcal{S} . We now elaborate on criteria which can allow for $\{S_i\}_{i=1}^m$ or \mathcal{S} to have certain requirements imposed on them.

- C_1 The minimum occurrence of every phoneme in \mathcal{S} should be maximized, namely,

$$C_1(\mathcal{U}_t, \gamma_p\%) = \min_k \{\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%))\} \quad (4)$$

Subsequently maximizing (4) ensures that even the phoneme that occur the least number of times in $\mathcal{S}(\mathcal{U}_t, \gamma_p\%)$ is maximized.

- C_2 Let $\#U_n(\mathcal{S})$ denote the count of utterance U_n in the corpus \mathcal{S} . A measure of equal distribution of utterances in the corpus would be

$$C_2(\mathcal{U}_t, \gamma_p\%) = \frac{1}{N} \sum_{n=1}^N (\#U_n(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)) - \bar{U})^2 \quad (5)$$

where $\bar{U} = \frac{1}{N} \sum_{n=1}^N (\#U_n(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)))$ is the mean. Note that $C_2(\mathcal{U}_t, \gamma_p\%)$ captures the distribution of the utterances in $\mathcal{S}(\mathcal{U}_t, \gamma_p\%)$. The configuration $(\mathcal{U}_t, \gamma_p\%)$ for which (5) is minimum is desired so that all utterances occur uniformly in \mathcal{S} .

- C_3 Common utterances between any two speaker sets, namely,

$$C_3(\mathcal{U}_t, \gamma_p\%) = \sum_{i,j=i+1}^m |S_i(\mathcal{U}_t, \gamma_p\%) \cap S_j(\mathcal{U}_t, \gamma_p\%)| \quad (6)$$

where $S_i(\mathcal{U}_t, \gamma_p\%) \cap S_j(\mathcal{U}_t, \gamma_p\%)$ captures the utterances that are common to both S_i and S_j and $|S_i \cap S_j|$ gives the count of common utterances. The configuration for which (6) is minimum is desired so that there is a rich utterance variability in corpus \mathcal{S} .

We hypothesize that the combination of these criteria jointly (7) produces the best possible dataset for a given speech application rather than the dataset which is based on individual criteria. Namely,

$$(\mathcal{U}_t^*, \gamma_p\%^*) = \arg \min_{(\mathcal{U}_t, \gamma_p\%)} \left\{ w_1 C_0(\mathcal{U}_t, \gamma_p\%) + w_2 \left(\frac{1}{C_1(\mathcal{U}_t, \gamma_p\%)} \right) + \right. \\ \left. w_3 C_2(\mathcal{U}_t, \gamma_p\%) + w_4 C_3(\mathcal{U}_t, \gamma_p\%) \right\} \quad (7)$$

where w_i are the weights and $\sum_{i=1}^4 w_i = 1$. Algorithm (1) describes this in more detail.

Note that in literature C_0 is the only criteria that is used to build a phonetically balanced speech corpus. The main contribution of this paper is to identify criteria that make the speech corpus usable. For example, $w_1 = 1$ and $w_2 = w_3 = w_4 = 0$ would reduce to what is done in the literature.

Algorithm 1. Multi Criteria approach for constructing \mathcal{S} .

```

for  $(\mathcal{U}_t, \gamma_p\%)$  do
  Generate  $\mathcal{S}(\mathcal{U}_t, \gamma_p\%)$ 
  for  $k = 1, 2, \dots K$  do
    Compute  $\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%))$ 
  end for
  for  $n = 1, 2, \dots N$  do
    Compute  $\#U_n(\mathcal{S}(\mathcal{U}_t, \gamma_p\%))$ 
  end for

  Find  $\bar{P} = \frac{1}{K} \sum_{k=1}^K (\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)))$ ;
   $\bar{U} = \frac{1}{N} \sum_{n=1}^N (\#U_n(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)))$ 

   $C_0(\mathcal{U}_t, \gamma_p\%) = \frac{1}{K} \sum_{k=1}^K (\#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)) - \bar{P})^2$ 

   $C_1(\mathcal{U}_t, \gamma_p\%) = \min_k \{ \#P_k(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)) \}$ 

   $C_2(\mathcal{U}_t, \gamma_p\%) = \frac{1}{N} \sum_{n=1}^N (\#U_n(\mathcal{S}(\mathcal{U}_t, \gamma_p\%)) - \bar{U})^2$ 

   $C_3(\mathcal{U}_t, \gamma_p\%) = \sum_{i,j=i+1} |S_i(\mathcal{U}_t, \gamma_p\%) \cap S_j(\mathcal{U}_t, \gamma_p\%)|$ 
end for
Normalize using (8)
 $C_0(\mathcal{U}_t, \gamma_p\%), C_1(\mathcal{U}_t, \gamma_p\%), C_2(\mathcal{U}_t, \gamma_p\%), C_3(\mathcal{U}_t, \gamma_p\%)$ 
to produce  $nC_0(\mathcal{U}_t, \gamma_p\%), nC_1(\mathcal{U}_t, \gamma_p\%), nC_2(\mathcal{U}_t, \gamma_p\%), nC_3(\mathcal{U}_t, \gamma_p\%)$ 
 $(\mathcal{U}_t^*, \gamma_p\%^*) = \arg \min_{(\mathcal{U}_t, \gamma_p\%)} \left\{ \begin{array}{l} w_1 nC_0(\mathcal{U}_t, \gamma_p\%) + w_2 \left( \frac{1}{nC_1(\mathcal{U}_t, \gamma_p\%)} \right) + \\ w_3 nC_2(\mathcal{U}_t, \gamma_p\%) + w_4 nC_3(\mathcal{U}_t, \gamma_p\%) \end{array} \right\}$ 

```

4 Experimental Results

For the purpose of analysis we collected $N = 1493$ unique English utterances, namely, $\mathcal{U} = \{U_1, U_2, \dots, U_{1493}\}$ [12]. The number of phonemes is $K = 39$. The distribution of the phonemes in \mathcal{U} is shown in Figure 1. It can be observed that the phoneme ‘AH’ occurs the most number of times (11.55%) in \mathcal{U} while the ‘OY’ occurs the least number of times (0.05%). All our experimental results are based on this set of utterances.

Using (2) we arranged all the 1493 utterances in the descending order of their f_i score. The sorted utterances were partitioned into two sets. First set $(\mathcal{U}_t(k))$ contained the first $k = 1, 2, \dots 1493$ utterances while the second set $(\mathcal{U}_b(k))$ contained $(1493 - k)$ utterances. The task was to build $m = 500$ sets with each S_i containing $p = 10$ utterances, such that $|\mathcal{S}| = 5000$. Each S_i gets $\gamma_p\% = 10, 20, \dots, 90$ utterances from $\mathcal{U}_t(k)$ while the remaining $(100 - \gamma_p\%)$ utterances came from $\mathcal{U}_b(k)$. In all we constructed $1493 \times 9 (= 13437)$ different sets of speaker sheets, namely, $\mathcal{S}(\mathcal{U}_t, \gamma_p\%)$ for $\mathcal{U}_t = 1, 2, \dots 1493$, $\gamma_p\% = 10, 20, \dots, 90$.

The first set of experiments were based on Algorithm 1 with $w_1 = 1$ and $w_2, w_3, w_4 = 0$ which is the generic approach adopted to build a phonetically balanced corpus in literature. For each of these speaker sheet sets we computed

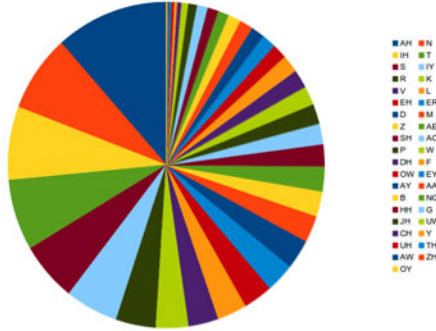


Fig. 1. Distribution of phoneme in \mathcal{U} .

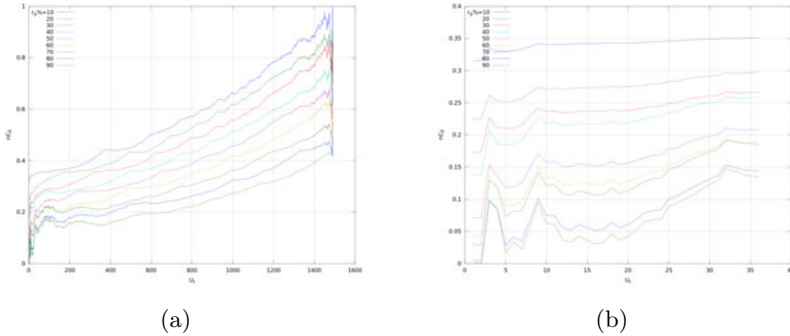


Fig. 2. (a) $C_0(\mathcal{U}_t, \gamma_p\%)$ and (b) for the first $U_t = 36$. The minimum occurs for $U_t = 2$ and $\gamma_p\% = 80$.

Table 1. ($\mathcal{U}_t, \gamma_p\%$) determined for different criteria.

w_1	w_2	w_3	w_4	$(\mathcal{U}_t^*, \gamma_p^{\text{OT}*})$
1	0	0	0	(2, 80)
0	1	0	0	(126, 90)
0	0	1	0	(49, 10)
0	0	0	1	(125, 10)
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	(367, 90)

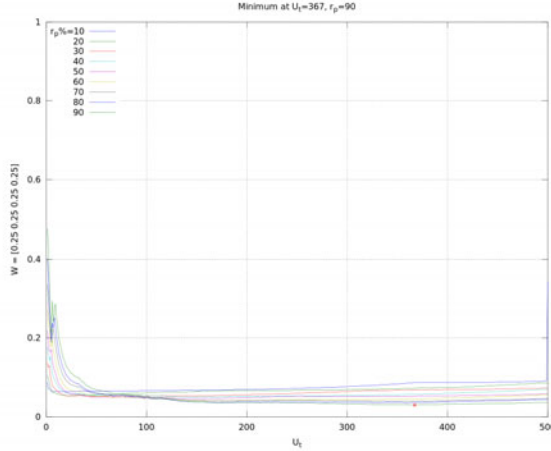
 $C_0(\mathcal{U}_t, \gamma_p \%)$ (3) and normalized it

$$nC_0 = \frac{(C_0 - \min(C_0))}{(\max(C_0) - \min(C_0))} \quad (8)$$

Figure 2 shows the plot of nC_0 for different values of $(\mathcal{U}_t, \gamma_p\%)$. The speaker sheet set (among the 13437 speaker sheet sets) with the least $C_0(\mathcal{U}_t, \gamma_p\%)$ is the set that is phonetically best balanced. As can be seen $\mathcal{U}_t(k=2), \gamma_p\% = 80$ (see Figure 2(b)) produces the best phonetically balanced dataset suitable for ASR

Table 2. Criteria cost for different speaker sheet set.

$(\mathcal{U}_t, \gamma_p\%)$	nC_0	$1/nC_1$	nC_2	nC_3
(2, 80)	0	0.00042	0.556	0.888
(126, 90)	0.104	0	0.00714	0.0645
(49, 10)	0.260	0.00002	0	0.00098
(125, 10)	0.268	0.00007	0.00047	0
(367, 90)	0.097	9.0363e-06	0.0077	.017312

**Fig. 3.** $(\mathcal{U}_t^*, \gamma_p\%^*) = (367, 90)$ for $w_{1,2,3,4} = 1/4$.

type applications. Clearly, one can observe that using only the C_0 criteria does not produce the best data set (even though it is best in the sense of phonetically being balanced) because the majority of the utterances, namely $\gamma_p\% = 80$ of the dataset consists of just $\mathcal{U}_t = 2$ utterances. This motivates the need for other criteria to construct a speech corpora.

Table 1 gives the speaker sheet set (denoted by $(\mathcal{U}_t, \gamma_p\%)$) that best produces \mathcal{S} if we consider different combination of the proposed criteria. Clearly the choice of speaker sheet sets depends on the emphasis given to a criteria.

When all the four criteria are given equal weightage, speaker sheet set denoted by $(\mathcal{U}_t^*, \gamma_p\%^*) = (367, 90)$ is the best (see last row in Table 1). It is clear from Table 2 that $(\mathcal{U}_t^*, \gamma_p\%^*) = (367, 90)$ is the best in terms of individual criteria ($nC_0, 1/nC_1, nC_2, nC_3$) being minimum together.

5 Conclusion

In this paper, we proposed a multi criteria approach to generate speaker sheets which satisfy the desired requirements on the speech corpora. We believe the formulation can be used to generate speaker sheets which will assist in building

a speech corpora that might be required for different speech applications and research. For example, a researcher who is doing in-depth analysis on phones may want to maximize the occurrence of the phone that occurs the least number of times (using C_2). We believe that satisfying all the proposed criteria jointly will produce the best speech corpora in terms of its being useful for different speech research and development. The main contribution of this paper is (a) identification of several criteria which need to be satisfied to generate a speech corpora, (b) formulation of a multi-criteria approach by combining the criteria and (c) experimental validation of the proposed approach for speaker sheet generation.

References

1. SPEECON, Speech-driven interfaces for consumer devices (2014). <http://www.speechdat.org/speecon/index.html>
2. Abushariah, M.A., Ainon, R.N., Zainuddin, R., Elshafei, M., Khalifa, O.O.: Phonetically rich and balanced text and speech corpora for Arabic language. *Lang. Resour. Eval.* **46**(4), 601–634 (2012)
3. Pineda, L.A., Pineda, L.V., Cuétara, J., Castellanos, H., López, I.: DIMEx100: a new phonetic and speech corpus for Mexican Spanish. In: Lemaitre, C., Reyes, C.A., González, J.A. (eds.) *IBERAMIA 2004*. LNCS (LNAI), vol. 3315, pp. 974–983. Springer, Heidelberg (2004)
4. Uraga, E., Gamboa, C.: VOXMEX speech database: design of a phonetically balanced corpus. In: *Proceedings of the Fourth International Conference on Language Resources and Evaluation. LREC 2004*, Lisbon, Portugal, May 26–28. European Language Resources Association (2004)
5. Asinovsky, A., Bogdanova, N., Rusakova, M., Ryko, A., Stepanova, S., Sherstinova, T.: The ORD speech corpus of Russian everyday communication “One Speaker’s Day”: creation principles and annotation. In: Matoušek, V., Mautner, P. (eds.) *TSD 2009*. LNCS, vol. 5729, pp. 250–257. Springer, Heidelberg (2009)
6. van Heerden, C., Davel, M.H., Barnard, E.: The semi-automated creation of stratified speech corpora (2013). <http://www.nwu.ac.za/sites/www.nwu.ac.za/files/files/v-must/Publications/prasa2013-17.pdf>
7. Tian, J., Nurminen, J., Kiss, I.: Optimal subset selection from text databases. In: *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, (ICASSP 2005)*, vol. 1, pp. 305–308, March 2005
8. Wu, Y., Zhang, R., Rudnicky, A.: Data selection for speech recognition. In: *IEEE Workshop on Automatic Speech Recognition Understanding, ASRU*, pp. 562–565, December 2007. <http://www.cs.cmu.edu/~yiwu/paper/asru07.pdf>
9. Nagroski, A., Boves, L., Steeneken, H.: Optimal selection of speech data for automatic speech recognition systems. In: *ICSLP*, pp. 2473–2476 (2002)
10. Chitturi, R., Mariam, S.H., Kumar, R.: Rapid methods for optimal text selection. In: *Recent Advances in Natural Language Processing*, September 2005
11. Mandal, S., Das, B., Mitra, P., Basu, A.: Developing Bengali speech corpus for phone recognizer using optimum text selection technique. In: *2011 International Conference on Asian Language Processing (IALP)*, pp. 268–271, November 2011
12. Awaz, Y.P.: Data: Speaker sheet generation for building speech corpora (2015). <https://sites.google.com/site/awazyp/data/speaker>

Text, Speech, and Dialogue

18th International Conference, TSD 2015, Pilsen, Czech
Republic, September 14-17, 2015, Proceedings

Král, P.; Matoušek, V. (Eds.)

2015, XVIII, 612 p. 122 illus. in color., Softcover

ISBN: 978-3-319-24032-9