

# Statistical Testing on ASR Performance via Blockwise Bootstrap

Zhe Liu, Fuchun Peng

Facebook AI, Menlo Park, CA, USA

zheliu@fb.com, fuchunpeng@fb.com

## Abstract

A common question being raised in automatic speech recognition (ASR) evaluations is how reliable is an observed word error rate (WER) improvement comparing two ASR systems, where statistical hypothesis testing and confidence interval (CI) can be utilized to tell whether this improvement is real or only due to random chance. The bootstrap resampling method has been popular for such significance analysis which is intuitive and easy to use. However, this method fails in dealing with dependent data, which is prevalent in speech world - for example, ASR performance on utterances from the same speaker could be correlated. In this paper we present blockwise bootstrap approach - by dividing evaluation utterances into nonoverlapping blocks, this method resamples these blocks instead of original data. We show that the resulting variance estimator of absolute WER difference between two ASR systems is consistent under mild conditions. We also demonstrate the validity of blockwise bootstrap method on both synthetic and real-world speech data. **Index Terms:** automatic speech recognition, word error rate, statistical hypothesis testing, confidence interval, resampling, blockwise bootstrap

## 1. Introduction

The most widely used metric for measuring the performance of an automatic speech recognition (ASR) system is the word error rate (WER), which is derived from the Levenshtein distance [1] working at the word level:

$$WER = \frac{\sum_{s=1}^n e_s}{\sum_{s=1}^n m_s}, \quad (1)$$

where  $m_s$  is the number of words in the  $s$ th sentence (i.e. reference text of audio) of the evaluation dataset, and  $e_s$  represents the sum of insertion, deletion and substitution errors computed from the dynamic string alignment of the recognized word sequence with the reference word sequence. The WER may also be referred to as the length normalized edit distance [2].

A practical question being commonly raised in ASR system evaluations is that how reliable is an observed improvement of ASR system  $B$  comparing to ASR system  $A$ . For example, if we obtained an absolute 0.2% WER reduction, how could we tell if this improvement is real and not due to random chance.

Here is where the statistical hypothesis testing comes into play. The use of statistical testing in ASR evaluations has been previously explored [3, 4, 5, 6, 7, 8]. In particular, the work of [6] presents a *bootstrap* method for significance analysis on ASR evaluations which makes no distributional approximations and the results are immediately interpretable in terms of WER.

To be more specific, suppose we have a sequence of independent and identically distributed (i.i.d.) random variable  $\{Z_s\}_{s=1}^n$  and we are interested in estimating the variance of some statistic  $T(Z_1, \dots, Z_n)$ . The bootstrap method [9, 10] resamples data from the empirical distribution of  $\{Z_s\}_{s=1}^n$ , and

then recalculates the statistic  $T$  on each of these “bootstrap” samples. Then the variance of  $T(Z_1, \dots, Z_n)$  can be estimated from the sample variance of these computed statistics.

For the ASR systems comparison problem that we raised previously, authors in [6] proposed the idea of using bootstrap approach above to resample (with replacement) the utterances in the evaluation dataset for each replicate, and then estimate the probability that the absolute WER difference of ASR system  $B$  versus system  $A$

$$\Delta W := WER^B - WER^A = \frac{\sum_{s=1}^n (e_s^B - e_s^A)}{\sum_{s=1}^n m_s} \quad (2)$$

is positive, where ASR systems  $A$  and  $B$  have word error counts  $e_s^A$  and  $e_s^B$  on the  $s$ th sentence, respectively. Notice that we should calculate the difference in the number of errors of the two systems on identical bootstrap samples.

However, one of the key issue confronting bootstrap resampling approximations is how to deal with dependent data [11]. This is particularly the case for speech data since the speech recognition errors could be highly correlated across different utterances if they are 1) from the same speaker, 2) similar in the uttered sentence (e.g. in the same domain or topic), or both. When the dependent structure across  $\{Z_s\}_{s=1}^n$  is nontrivial, the true sampling distribution of  $T(Z_1, \dots, Z_n)$  would depend on the joint distribution of  $\{Z_s\}_{s=1}^n$  and thus the bootstrap samples should preserve such dependent structure as well. Unfortunately, the reshuffled samples obtained from the ordinary bootstrap method break such dependence, and thus lead to wrong variance estimations of the statistic. In particular, ASR errors on dependent utterances could always be positively correlated, which typically makes confidence intervals computed by bootstrap much narrower than what they should be. This might lead to over-optimistic conclusions due to false-positive discovery.

In this paper, we present the *blockwise bootstrap* approach for statistical testing of ASR performance. By dividing evaluation utterances into nonoverlapping blocks, we resamples these blocks instead of original data points (i.e. word errors) of utterances. The idea of blockwise bootstrap was initially developed in the work of [12, 13] for dealing with dependent time series data. Since then, there has been a line of research in statistical literature on various block construction methods and theoretical comparisons of them [14, 15, 16, 17, 18, 19]. To the best of our knowledge, our work is the first to introduce the blockwise bootstrap approach to address dependent speech data in ASR performance evaluations and illustrate how it helps calculate valid confidence intervals for absolute WER difference.

The rest of the paper is organized as follows. Section 2 introduces the use of blockwise bootstrap on ASR evaluation problem. Section 3 shows the statistical property that the blockwise variance estimator is consistent under mild conditions. Section 4 and Section 5 demonstrate the validity of blockwise bootstrap method on simulated synthetic data and real-world speech data. We conclude in Section 6.

## 2. Methods

In this section, we describe the blockwise bootstrap method to compute the confidence interval of the absolute WER difference  $\Delta W$  as in the formula (2).

Given two ASR systems  $A$  and  $B$  in the comparison, assume we have their evaluation results on  $n$  utterances as follows:

$$(m_1, e_1^A, e_1^B), (m_2, e_2^A, e_2^B), \dots, (m_n, e_n^A, e_n^B) \quad (3)$$

where for any  $s = 1, \dots, n$ , let  $m_s$  be the number of words in the reference of  $s$ th utterance,  $e_s^A$  and  $e_s^B$  represent the numbers of word errors in ASR systems  $A$  and  $B$ , respectively. The statistic that we are interested in is the absolute WER difference  $\Delta W$  comparing system  $B$  versus system  $A$ .

Suppose the evaluation results data above can be partitioned into  $K$  nonoverlapping blocks (e.g. by speakers) such that word error counts within each block are correlated while the dependency of word error counts from different blocks are negligible

$$\{(m_s, e_s^A, e_s^B)\}_{s \in S_k}, k = 1, \dots, K \quad (4)$$

where  $\cup_k S_k = \{1, 2, \dots, n\}$ , and  $S_i \cap S_j = \emptyset$  for any  $i, j$ . Then the blockwise bootstrap method works as follows.

For any  $b = 1, \dots, B$  where  $B$  is a large number, we randomly sample (with replacement)  $K$  elements  $\{S_{k'}^{(b)}\}_{k'=1, \dots, K}$  from the set  $\{S_k\}_{k=1, \dots, K}$  to generate a bootstrap sample of evaluation results data

$$\{(m_s, e_s^A, e_s^B)\}_{s \in S_{k'}^{(b)}}, k' = 1, \dots, K \quad (5)$$

where each  $S_{k'}^{(b)} \in \{S_k\}_{k=1, \dots, K}$ . Then for this bootstrap replicate sample, the statistic is computed as

$$\Delta W^{(b)} = \frac{\sum_{k'=1}^K \sum_{s \in S_{k'}^{(b)}} (e_s^B - e_s^A)}{\sum_{k'=1}^K \sum_{s \in S_{k'}^{(b)}} m_s}. \quad (6)$$

Once we have all  $\{\Delta W^{(b)}\}_{b=1, \dots, B}$ , then the 95% confidence interval for  $\Delta W$  can be determined by the empirical percentiles at 2.5% and 97.5% of the bootstrap sample statistics

$$(\Delta W_{2.5\%}^{blockBoot}, \Delta W_{97.5\%}^{blockBoot}). \quad (7)$$

Alternatively, the uncertainty of  $\Delta W$  can be quantified by its standard error, which can be approximated by the sample standard deviation of the  $B$  bootstrap sample statistics

$$se^{blockBoot}(\Delta W) = \sqrt{\frac{\sum_{b=1}^B (\Delta W^{(b)} - m^{blockBoot}(\Delta W))^2}{B-1}} \quad (8)$$

where  $m^{blockBoot}(\Delta W) = \frac{1}{B} \sum_{b=1}^B \Delta W^{(b)}$ . Then based on the Gaussian approximation, the 95% confidence interval for  $\Delta W$  can be obtained by

$$(m^{blockBoot}(\Delta W) \pm 1.96 \cdot se^{blockBoot}(\Delta W)). \quad (9)$$

Generally speaking, the percentile confidence intervals (7) always give similar results with Gaussian approximation confidence intervals (9) when  $B$  is large, unless the corresponding bootstrap sample statistics are highly skewed, in which case the former ones are preferred.

Note that in this paper we mainly focus on the confidence intervals of absolute WER difference. Similarly, the blockwise bootstrap method can also be utilized to compute the confidence intervals for the WER itself as well as relative WER difference between two ASR systems.

## 3. Theoretical Properties

We work out some statistical theories to show that the blockwise variance estimator of  $\Delta W$  is consistent under mild conditions. That is, as the number of evaluation data points increases indefinitely, the resulting sequence of variance estimates converges to the truth variance [20].

For simplicity, we assume all utterances in the evaluation dataset have the same number of words, that is,  $m_i = m$  for all  $i = 1, \dots, n$ . Let's denote  $Z_i := (e_i^B - e_i^A)/m$ . Then we have the statistic of interest written as

$$\Delta W_n = \frac{1}{n} \sum_{i=1}^n Z_i \quad (10)$$

where the subscript  $n$  in  $\Delta W_n$  indicates the number of samples (i.e. utterances) corresponding to the quantity. Further, by dividing the  $n$  utterances into nonoverlapping blocks, suppose each block has same number of utterances, denoted as  $d_n$ . Let  $K_n = \lfloor n/d_n \rfloor$  as the number of blocks. Note that the assumptions of having same number of words in each utterance as well as same number of utterances in each block are for the sake of simplicity, the results below still hold if the equality assumptions here are relaxed to be in the same order of  $n$ .

Without the loss of generality, assume that the blocks are consecutive and thus the  $k$ th block consists of the sequence  $Z_{(k-1)d_n+1}, \dots, Z_{kd_n}$  where  $k = 1, \dots, K_n$ . We further let

$$\Delta W_{d_n}^{(k-1)d_n+1} := \frac{1}{d_n} \sum_{i=(k-1)d_n+1}^{kd_n} Z_i \quad (11)$$

where the subscript  $d_n$  in  $\Delta W_{d_n}^{(k-1)d_n+1}$  for the  $k$ th block represents the block size and the superscript  $(k-1)d_n+1$  indicates the index of starting variable in the block. Notice that  $\frac{1}{K_n} \sum_{k=1}^{K_n} \Delta W_{d_n}^{(k-1)d_n+1} = \Delta W_n$ .

Consider the blockwise variance estimator (i.e. standardized sample variance of blockwise estimators of the statistic)

$$\hat{\sigma}_n^2 := \frac{d_n}{K_n} \sum_{k=1}^{K_n} \left( \Delta W_{d_n}^{(k-1)d_n+1} - \Delta W_n \right)^2. \quad (12)$$

The following theorem establishes its  $L_2$ -consistency result.

**Theorem 1.** Assume the asymptotic variance of  $\Delta W_n$  is

$$\lim_{n \rightarrow \infty} n \mathbf{E}(\Delta W_n - \mathbf{E}(\Delta W_n))^2 = \sigma^2 \in (0, \infty) \quad (13)$$

and  $\mu = \mathbf{E}(Z_i)$  for any  $i = 1, \dots, n$ . Let  $d_n$  be s.t.  $d_n \rightarrow \infty$  and  $K_n \rightarrow \infty$  as  $n \rightarrow \infty$ . If  $n^2 \mathbf{E}(\Delta W_n - \mu)^4$  is uniformly bounded, and for any  $1 \leq k < k' \leq K_n$ , the sequence of  $Z_{(k-1)d_n+1}, \dots, Z_{kd_n}$  and the sequence of  $Z_{(k'-1)d_n+1}, \dots, Z_{k'd_n}$  are uncorrelated, then

$$\hat{\sigma}_n^2 \rightarrow_{L_2} \sigma^2 \text{ as } n \rightarrow \infty. \quad (14)$$

The proof of the theorem is deferred to Appendix (accompanied by additional files). According to Theorem 1, we require both the number of blocks  $K_n$  and number of utterances in each block  $d_n$  go to infinity as the number of utterances  $n$  grows to infinity. This is reasonable for speech evaluation data collection if the blocks are partitioned by different speakers or topics.

On the other hand, the uncorrelated assumption between different blocks seems strong, because in practice it's possible that word errors over different blocks are also weakly correlated. We relax this assumption in the corollary below.

**Corollary 1.1.** *Theorem 1 still holds if the assumption of uncorrelated blockwise variables is relaxed as follows: for any  $\epsilon > 0$ , if  $n$  is large enough, for any  $1 \leq k < k' \leq K_n$  and  $(k-1)d_n+1 \leq i < j \leq kd_n$ ,  $(k'-1)d_n+1 \leq i' < j' \leq k'd_n$  assume*

$$\mathbf{E}(|(Z_i - \mu)(Z_j - \mu)(Z_{i'} - \mu)(Z_{j'} - \mu)|) \leq \epsilon. \quad (15)$$

The proof of the corollary is also deferred to Appendix (accompanied by additional files).

## 4. Simulation Experiments

In this section, we conduct simulation experiments to show that the blockwise bootstrap approach is capable to generate valid confidence intervals for absolute WER differences between two ASR systems and is superior to the ordinary bootstrap when the utterances in the evaluation dataset are dependent.

### 4.1. Setup

In this simulation experiments, we generate synthetic data, i.e. counts of ASR errors, to measure the performance of ordinary bootstrap and blockwise bootstrap methods. We assume the total number of utterances in the evaluation set is  $n = 3,000$  and number of words in each utterance is equally  $m = 100$ . For the two ASR systems  $A$  and  $B$  in the comparison, suppose the ‘‘ground-truth’’ WERs are given by  $WER^A = 10.0\%$  and  $WER^B = 9.5\%$  respectively. Thus the absolute WER difference  $\Delta W$  between them is  $-0.5\%$ .

Under the scenario that the numbers of errors from different utterances are independent with each other, we generate the number of errors for each utterance from the binomial distribution  $Binom(m = 100, p)$  where  $p = WER^A$  or  $WER^B$  depending on which ASR system was used. On the other hand, when the ASR errors are dependent across different utterances, we need to make additional correlation structure assumption while keeping the marginal distribution of error count on each utterance to be binomial distributed.

Here we assume the numbers of errors across different utterances are block-correlated, that is, for any two utterances, their ASR errors are correlated if they belong to the same block while their errors are independent if they belong to different blocks. Without the loss of generality, suppose the blocks are consecutive and the size of block (i.e. number of utterances in each block) is denoted as  $d$ . Follow the steps below to generate ASR errors for each block and each ASR system:

1. Generate a sample  $(v_1, \dots, v_d)$  from multivariate Gaussian distribution  $N(0, \Sigma_d)$ , where  $\Sigma_d$  is an  $d$ -by- $d$  covariance matrix with  $\sigma_{ij} = 1$  if  $i = j$  and  $\sigma_{ij} = \rho$  if  $i \neq j$ . Here  $\sigma_{ij}$  is the  $(i, j)$ -th element of  $\Sigma_d$ ;
2. Turn  $(v_1, \dots, v_d)$  into correlated uniforms  $(u_1, \dots, u_d)$  where  $u_s = \Phi(v_s)$  for  $s = 1, \dots, d$  and  $\Phi(\cdot)$  is the Gaussian cumulative distribution function (CDF);
3. Generate correlated Binomial samples  $(e_1, \dots, e_d)$  by inverting the Binomial CDF:  $e_s = Q_{binom}(u_s, m, p)$  for  $s = 1, \dots, d$ , where  $Q_{binom}(\cdot, m, p)$  is the inverse of the Binomial CDF.

The counts of word errors for different blocks and different ASR systems are generated independently. For both ordinary bootstrap and blockwise bootstrap methods in the comparison, we set the resampling size  $B = 1,000$ . The block size  $d$  and correlation parameter  $\rho$  are varied in our experiment.

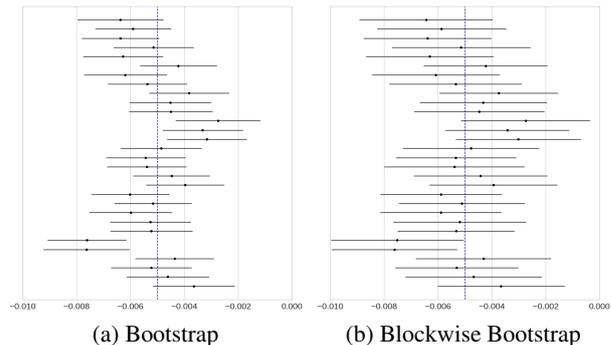


Figure 1: *Visualisation of confidence intervals computed on the first 30 simulated data ( $d = 5, \rho = 0.4$ ), from both ordinary bootstrap and blockwise bootstrap. The vertical blue line represents the true absolute WER difference  $\Delta W = -0.5\%$ .*

While in this and next section we mainly focus on the comparison against ordinary bootstrap method, there also exists other parametric approaches (e.g. in [8]) without the need of Monte Carlo resampling. They typically give similar results with bootstrap but are less preferred than bootstrap especially when these distributional assumptions are violated in practice.

### 4.2. Results

Strictly speaking, a 95% confidence interval means that if we were able to have 100 different datasets from the same distribution of the original data and compute a 95% confidence interval based on each of these datasets, then approximately 95 of these 100 confidence intervals will contain the true value of the statistic of interest [21, 22, 23]. In our experiments, for each setup of various block size and correlation parameter, we replicate the simulation for 1,000 times. Therefore if confidence intervals were computed correctly, then approximately 950 of these confidence intervals would contain the true absolute WER difference  $\Delta W = -0.5\%$ .

Seen from Table 1, the blockwise bootstrap method always gives valid confidence intervals since the percentage of confidence intervals that contain the true  $\Delta W$  is very close to 95%, regardless of block size ( $d = 5$  or  $30$ ) or correlation parameter ( $\rho = 0, 0.05, 0.1, 0.2$  or  $0.4$ ). As the block size or correlation increases, the width of confidence intervals becomes larger. On the other hand, the ordinary bootstrap method fails to generate correct confidence intervals when the data is dependent ( $\rho > 0$ ), since its percentage of confidence intervals that contain the true  $\Delta W$  is much lower than 95%.

Figure 1 plots the confidence intervals computed on the first 30 simulated data for  $d = 5$  and  $\rho = 0.4$ , where we can see that the confidence intervals from the blockwise bootstrap method are wider and capture more true values of  $\Delta W = -0.5\%$ .

It is worth noting that when the dependency across different utterances is relatively strong, the width of confidence intervals generated by ordinary bootstrap is less than the half of the valid ones generated by blockwise bootstrap, in which case ordinary bootstrap might lead to over-optimistic discovery.

## 5. Real Data Experiments

In this section, we also apply the blockwise bootstrap approach on two real-world speech datasets and demonstrate how it helps compute the confidence intervals of absolute WER difference

Table 1: Comparison results of ordinary bootstrap and blockwise bootstrap methods on simulated data with various block size ( $d$ ) and correlation parameter ( $\rho$ ), where the average width of confidence intervals and the percentage of times that confidence intervals contain the true  $\Delta W$  are shown.

Block Size ( $d$ )	Correlation ( $\rho$ )	Bootstrap		Blockwise Bootstrap	
		Width	% contains the truth	Width	% contains the truth
$d = 5$	$\rho = 0$	0.0030	94.1%	0.0030	94.7%
	$\rho = 0.05$	0.0030	92.7%	0.0033	95.2%
	$\rho = 0.1$	0.0030	90.1%	0.0035	94.3%
	$\rho = 0.2$	0.0030	86.2%	0.0040	94.9%
	$\rho = 0.4$	0.0030	76.9%	0.0048	94.0%
$d = 30$	$\rho = 0$	0.0030	94.1%	0.0030	94.7%
	$\rho = 0.05$	0.0030	78.1%	0.0046	95.2%
	$\rho = 0.1$	0.0030	69.2%	0.0058	94.9%
	$\rho = 0.2$	0.0030	54.4%	0.0077	94.7%
	$\rho = 0.4$	0.0030	41.2%	0.0105	95.9%

between two ASR systems. We consider the following two ASR evaluation datasets in this analysis

- *Conversational Speech* data. This dataset was collected through crowd-sourcing from a data supplier for ASR, and the data was properly anonymized. It consists of 235 conversions with more than 20 topics that are common in daily life, including family, travel, foods, etc;
- *Augmented Multi-Party Interaction (AMI) Meeting* data. The dataset [24, 25] includes scenario meetings (with roles assigned for participants) and non-scenario meetings (where participants were free to choose topics). For scenario meetings, each session is divided into 4 one-hour meetings. Each meeting has 4 participants.

Since our purpose is to evaluate ASR performance, we only use the “dev” and “eval” splits of the entire datasets above.

We use in-house developed conversation ASR system in this investigation: a baseline model (denoted as system  $A$ ) and an improved model (denoted as system  $B$ ), and we are interested in computing a 95% confidence interval of the absolute WER difference between the two ASRs. If the upper bound of such confidence interval is negative, then we can tell that this improvement is real and not due to random chance.

To apply the blockwise bootstrap method, we need to define the correlated block structures among the utterances in the evaluation data (merged “dev” and “eval” splits). For Conversation data, it’s natural to treat each conversation as a single separated block since the same topics were being discussed. For AMI Meeting dataset, we treat the utterances from each speaker in each (either scenario or non-scenario) meeting as a block. By doing that, for any two utterances, we assume their ASR errors are correlated if they belong to the same block while the errors have very weak correlations if they belong to different blocks. Table 2 shows details of the two evaluation datasets in terms of number of utterances, number of total words, and number of correlated blocks.

We apply both ordinary bootstrap and blockwise bootstrap methods on the two evaluation datasets. Results are shown in Table 3. Again, we observe that the confidence intervals computed from blockwise bootstrap are much wider than the ones generated from ordinary bootstrap: around 1.5 times wider on Conversation data and 2 times wider on AMI Meeting data. Also, we can see that confidence intervals computed from the empirical percentiles at 2.5% and 97.5% of bootstrap samples

Table 2: Summary of the Conversation Speech and AMI Meeting datasets in the experiments of real data analysis.

Feature	Evaluation Dataset	
	Conversation	AMI Meeting
Number of Utterances	13,987	25,741
Number of Words	160,338	189,590
Number of Correlated Blocks	235	135

Table 3: Results of bootstrap and blockwise bootstrap methods on real-world Conversation and AMI meeting datasets.

Method	Metric	Evaluation Dataset	
		Conversation	AMI Meeting
<b>Bootstrap</b>	$\Delta W$	-1.47%	-1.80%
	$se^{boot}(\Delta W)$	0.074%	0.067%
	Percentile CI	(-1.61%, -1.32%)	(-1.94%, -1.67%)
	Gaussian Approx. CI	(-1.62%, -1.33%)	(-1.94%, -1.67%)
<b>Blockwise Bootstrap</b>	$\Delta W$	-1.47%	-1.80%
	$se^{blockBoot}(\Delta W)$	0.116%	0.153%
	Percentile CI	(-1.69%, -1.24%)	(-2.09%, -1.51%)
	Gaussian Approx. CI	(-1.69%, -1.24%)	(-2.10%, -1.50%)

are almost the same with the ones computed from Gaussian approximation.

Figure 2 displays the histograms of absolute WER difference computed from the bootstrap samples, where we can see again that the data distribution from the blockwise bootstrap method is more spread out. Thus confidence intervals computed by ordinary bootstrap underestimate the standard errors and are much narrower than what they should be, which could lead to over-optimistic conclusions due to false-positive discovery.

## 6. Conclusion

In this paper, we present blockwise bootstrap approach for statistical testing of ASR performance - by dividing the evaluation utterances into nonoverlapping blocks, this method resamples these blocks instead of original data. We show that the resulting variance estimator of the absolute WER difference is consistent under mild conditions. We also illustrate the validity of blockwise bootstrap method on synthetic and real-world speech data.

Future work might include how to infer the correlated block structures from data, for example, estimating a blockwise sparse correlation matrix across evaluation utterances based on the embeddings of speakers and text sentences.

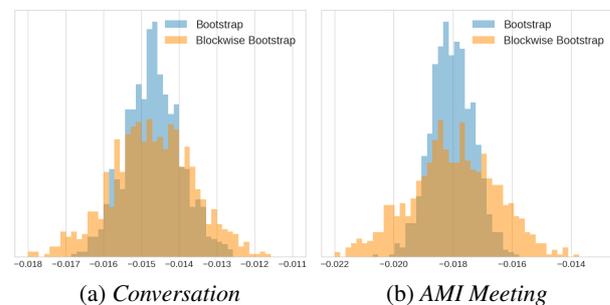


Figure 2: Histograms of absolute WER difference calculated from the bootstrap samples of both bootstrap and blockwise bootstrap methods.

## 7. References

- [1] G. Navarro, "A guided tour to approximate string matching," *ACM Computing Surveys (CSUR)*, vol. 33, no. 1, pp. 31–88, 2001.
- [2] S. Nießen, F. J. Och, G. Leusch, H. Ney *et al.*, "An evaluation tool for machine translation: Fast evaluation for MT research." in *LREC*, 2000.
- [3] L. Gillick and S. J. Cox, "Some statistical issues in the comparison of speech recognition algorithms," in *International Conference on Acoustics, Speech, and Signal Processing*, 1989, pp. 532–535.
- [4] D. S. Pallett, J. G. Fiscus, W. M. Fisher, and J. S. Garofolo, "Benchmark tests for the DARPA spoken language program," in *Proceedings of the Workshop on Human Language Technology*, 1993, pp. 7–18.
- [5] H. Strik, C. Cucchiari, and J. M. Kessens, "Comparing the recognition performance of CSRs: in search of an adequate metric and statistical significance test," in *Proceeding of International Conference on Spoken Language Processing*, 2000, pp. 740–743.
- [6] M. Bisani and H. Ney, "Bootstrap estimates for confidence intervals in ASR performance evaluation," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, 2004, pp. I–409.
- [7] D. Vilar, G. Leusch, H. Ney, and R. E. Banchs, "Human evaluation of machine translation through binary system comparisons," in *Proceedings of the Second Workshop on Statistical Machine Translation*. Association for Computational Linguistics, 2007, pp. 96–103.
- [8] J. M. Vilar, "Efficient computation of confidence intervals for word error rates," in *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2008, pp. 5101–5104.
- [9] B. Efron and R. J. Tibshirani, *An introduction to the bootstrap*. CRC Press, 1994.
- [10] B. Efron, "Second thoughts on the bootstrap," *Statistical Science*, vol. 18, no. 2, pp. 135–140, 2003.
- [11] J.-P. Kreiss and E. Paparoditis, "Bootstrap methods for dependent data: A review," *Journal of the Korean Statistical Society*, vol. 40, no. 4, pp. 357–378, 2011.
- [12] P. Hall, "Resampling a coverage pattern," *Stochastic Processes and their Applications*, vol. 20, no. 2, pp. 231–246, 1985.
- [13] E. Carlstein, "The use of subseries values for estimating the variance of a general statistic from a stationary sequence," *The Annals of Statistics*, vol. 14, no. 3, pp. 1171–1179, 1986.
- [14] R. Y. Liu and K. Singh, "Moving blocks jackknife and bootstrap capture weak dependence," *Exploring the Limits of Bootstrap*, vol. 225, p. 248, 1992.
- [15] D. N. Politis and J. P. Romano, "A circular block-resampling procedure for stationary data," *Exploring the Limits of Bootstrap*, vol. 2635270, 1992.
- [16] ———, "The stationary bootstrap," *Journal of the American Statistical Association*, vol. 89, no. 428, pp. 1303–1313, 1994.
- [17] G. Morvai, S. Yakowitz, and L. Györfi, "Nonparametric inference for ergodic, stationary time series," *The Annals of Statistics*, vol. 24, no. 1, pp. 370–379, 1996.
- [18] S. N. Lahiri, "Theoretical comparisons of block bootstrap methods," *Annals of Statistics*, pp. 386–404, 1999.
- [19] D. J. Nordman, "A note on the stationary bootstraps variance," *The Annals of Statistics*, vol. 37, no. 1, pp. 359–370, 2009.
- [20] T. Amemiya, *Advanced econometrics*. Harvard university press, 1985.
- [21] J. Neyman, "Xoutline of a theory of statistical estimation based on the classical theory of probability," *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences*, vol. 236, no. 767, pp. 333–380, 1937.
- [22] A. Stuart and M. G. Kendall, *The advanced theory of statistics*. Griffin, 1963.
- [23] D. R. Cox and D. V. Hinkley, *Theoretical statistics*. CRC Press, 1979.
- [24] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, and M. Kronenthal, "The AMI meeting corpus: A pre-announcement," in *International Workshop on Machine Learning for Multimodal Interaction*. Springer, 2005, pp. 28–39.
- [25] O. Aran, H. Hung, and D. Gatica-Perez, "A multimodal corpus for studying dominance in small group conversations," *Multimodal Corpora: Advances in Capturing, Coding and Analyzing Multimodality*, vol. 22, 2010.

## 8. Appendix

### 8.1. Proof of Theorem 1

*Proof.* Let's denote

$$t_{d_n}^{(k-1)d_n+1} := \sqrt{d_n}(\Delta W_{d_n}^{(k-1)d_n+1} - \mu). \quad (16)$$

Then the variance estimator can be written as

$$\hat{\sigma}_n^2 = \frac{1}{K_n} \sum_{k=1}^{K_n} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 - (t_n)^2 \quad (17)$$

where

$$t_n = \sum_{k=1}^{K_n} t_{d_n}^{(k-1)d_n+1} / K_n. \quad (18)$$

Note that  $\mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} \right) = \mathbf{E}(t_n) = 0$ .

We will first show that the first term of the right hand side of (17) converges to  $\sigma^2$  in  $L_2$ . Notice that

$$\mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 = d_n \mathbf{E} \left( \Delta W_{d_n}^{(k-1)d_n+1} - \mu \right)^2 \quad (19)$$

$$\rightarrow \sigma^2 \quad (20)$$

as  $n \rightarrow \infty$  and (thus)  $d_n \rightarrow \infty$ . Then it suffices to show

$$\mathbf{Var} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 \right) \rightarrow 0. \quad (21)$$

This is true since

$$\mathbf{Var} \left( \sum_{k=1}^{K_n} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 \right) \leq \sum_{k=1}^{K_n} \mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} \right)^4 \quad (22)$$

$$\leq K_n C \quad (23)$$

when  $n$  and  $d_n$  are sufficiently large. Here the first less than or equal to sign follows from the assumption that for any two blocks  $k < k'$ ,  $\Delta W_{d_n}^{(k-1)d_n+1}$  and  $\Delta W_{d_n}^{(k'-1)d_n+1}$  are uncorrelated, and the second less than or equal to sign follows from the assumption that  $n^2 \mathbf{E}(\Delta W_n - \mu)^4$  is uniformly bounded.

Now we only need to show that the second term of the right hand side of (17) converges to 0 in  $L_2$ , or equivalently,  $t_n$  converges to 0 in  $L_4$ .

Note that

$$t_n = \sqrt{d_n}(\Delta W_n - \mu) \quad (24)$$

and thus

$$\mathbf{E}(t_n^4) = d_n^2 \mathbf{E}(\Delta W_n - \mu)^4 \rightarrow 0 \quad (25)$$

as  $n \rightarrow \infty$  since  $n^2 \mathbf{E}(\Delta W_n - \mu)^4$  is bounded.  $\square$

### 8.2. Proof of Corollary 1.1

*Proof.* It suffices to show

$$\mathbf{Var} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 \right) \rightarrow 0 \quad (26)$$

under the relaxed assumption. Consider for any  $k < k'$

$$\mathbf{Cov} \left( \left( t_{d_n}^{(k-1)d_n+1} \right)^2, \left( t_{d_n}^{(k'-1)d_n+1} \right)^2 \right) \quad (27)$$

$$\leq \mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} t_{d_n}^{(k'-1)d_n+1} \right)^2 \quad (28)$$

$$= d_n^2 \mathbf{E} \left( \Delta W_{d_n}^{(k-1)d_n+1} - \mu \right)^2 \left( \Delta W_{d_n}^{(k'-1)d_n+1} - \mu \right)^2 \quad (29)$$

$$= \frac{1}{d_n^2} \mathbf{E} \left( \sum_{i=(k-1)d_n+1}^{kd_n} (Z_i - \mu) \right)^2 \left( \sum_{i'=(k'-1)d_n+1}^{k'd_n} (Z_{i'} - \mu) \right)^2. \quad (30)$$

Then under the assumption that  $\mathbf{E}(|(Z_i - \mu)(Z_j - \mu)(Z_{i'} - \mu)(Z_{j'} - \mu)|) \leq \epsilon$  if  $n$  is sufficiently large for any  $i, j$  and  $i', j'$ , we have

$$\mathbf{Var} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \left( t_{d_n}^{(k-1)d_n+1} \right)^2 \right) \quad (31)$$

$$\leq \frac{1}{K_n^2} \sum_{k=1}^{K_n} \mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} \right)^4 \quad (32)$$

$$+ \frac{2}{K_n^2} \sum_{1 \leq k < k' \leq K_n} \mathbf{E} \left( t_{d_n}^{(k-1)d_n+1} t_{d_n}^{(k'-1)d_n+1} \right)^2 \quad (33)$$

$$\leq \frac{1}{K_n^2} (K_n C + K_n^2 \epsilon) = \frac{C}{K_n} + \epsilon. \quad (34)$$

which converges to 0 as  $n \rightarrow \infty$ .  $\square$