THE 1998 HTK SYSTEM FOR TRANSCRIPTION OF CONVERSATIONAL TELEPHONE SPEECH

T. Hain, P.C. Woodland, T.R. Niesler and E.W.D. Whittaker

Cambridge University Engineering Department, Trumpington Street, Cambridge CB2 1PZ, UK. Email: {th223,trn,ewdw2,pcw}@eng.cam.ac.uk

ABSTRACT

This paper describes the 1998 HTK large vocabulary speech recognition system for conversational telephone speech as used in the NIST 1998 Hub5E evaluation. Front-end and language modelling experiments conducted using various training and test sets from both the Switchboard and Callhome English corpora are presented. Our complete system includes reduced bandwidth analysis, sidebased cepstral feature normalisation, vocal tract length normalisation (VTLN), triphone and quinphone hidden Markov models (HMMs) built using speaker adaptive training (SAT), maximum likelihood linear regression (MLLR) speaker adaptation and a confidence score based system combination. A detailed description of the complete system together with experimental results for each stage of our multi-pass decoding scheme is presented. The word error rate obtained is almost 20% better than our 1997 system on the development set.

1. INTRODUCTION

Transcription of conversational telephone speech is a complex task, which has to deal with many severe degradations in speech quality. These degradations continue to lead to word error rates in the range of 30 to 50 %, which are almost twice as high as for other difficult tasks like Broadcast News Transcription [7]. The difficulties result from a limited bandwidth, distorted audio channels, cross-talk and other acoustic interference, as well as highly variable speaking rates and conversational styles in which grammatical rules are less important.

Current experiments for conversational telephone speech are usually conducted on three corpora distributed by the Linguistic Data Consortium (LDC): Switchboard-I (Swbd-I), Switchboard-II (Swbd-II) and Callhome English (CHE). Both Switchboard corpora consist of telephone conversations within the USA between strangers. For Swbd-I speakers are given a topic, whereas for Swbd-II the topic is merely suggested. CHE data consists of calls to friends or relatives abroad. This leads not only to severe acoustic channel distortions caused by long distance telephone connections, but also to a higher number of non-English (and hence unknown) words. Furthermore, multiple speakers per conversation side are not uncommon. These factors usually lead to 10% difference in word error rate between Switchboard and Callhome recognition tests.

The Swbd-I, Swbd-II and CHE corpora are the subject of the yearly Hub5 evaluation conducted by the National Institute for Standards and Technology (NIST). In the following sections we

describe the system we prepared for participation in the 1998 Hub5E evaluation, and present its final performance.

The remainder of this paper is organised as follows: First we give a brief overview over our system and the development objectives. Then we describe our front-end experiments, including analysis bandwidth, cepstral feature normalisation and vocal tract length normalisation (VTLN). Subsequently we present the speaker adaptation tests using the new front-end, followed by details about language models used. The final section gives the overall system performance.

2. SYSTEM OVERVIEW

The basis for developing our 1998 system was formed by our 1997 conversational telephone speech transcription system [6]. This system employed gender independent decision-tree state-clustered triphone models, a 3-gram language model trained on 2 million words (MW) from Swbd-I and CHE, and a 22K word dictionary based on the LIMSI 1993 WSJ pronunciation dictionary. Standard techniques for telephone speech were employed, with the only major refinement being the introduction of VTLN. All available bandwidth were used with per segment cepstral mean normalisation.

During the development of our 1997 system we found the vocal tract length normalisation (VTLN) process to give unreliable results in terms of word error rate (WER) across speakers and test sets. In particular, the performance gain when using VTLN was considerably lower than expected especially for the 1997 Hub5E evaluation set. Furthermore the front-end processing did not account for the existence of very short speech segments, nor the special characteristics of the telephone channel. These issues have been addressed in a series of experiments, and subsequent improvements have been implemented and tested with our current system.

Our 1998 system uses an eight-pass decoding strategy with multiple gender independent and gender dependent state-clustered triphone and quinphone HMM model sets, and multiple stages of speaker adaptation.

Each frame of input speech is represented by a 39 dimensional feature vector that consists of 13 (including c_0) MF-PLP cepstral parameters and their first and second differentials. The results from experiments described in section 3.1 suggested the use of reduced bandwidth analysis and cepstral mean and variance normalisation per conversation side.

Three different types of HMM model sets were used. First, a gender independent state clustered triphone model set was built and trained using a subset of Swbd-I containing 65 hours of speech

(WS96train). The resulting system contained 6039 speech states with 12 Gaussian mixture components per state. The final model set (M1) was obtained from this by further reestimation and mixture splitting steps using a training set (h5train98) consisting of 163 hours of speech from Swbd-I and 17 hours from CHE. It was found that 16 Gaussians per speech state was optimal. The M1 model set was used in the first decoding pass to obtain the transcripts for gender detection and VTLN warp factor computation.

Secondly, a gender independent triphone model set that uses VTLN warped training data was obtained in a similar fashion to M1. Gender dependent versions were then derived by a single gender dependent reestimation step. This model-set pair, subsequently referred to as M2, was used in the second and the third recognition passes.

Finally, decoding passes 4-7 used a gender dependent pair of quinphone HMM models (M3) trained on the VTLN-warped h5train98 set. A further speaker adaptive training (SAT) [5] iteration has been used. The resulting model set contained 8763 speech states, each characterised by a 16 component mixture Gaussian.

In passes 3-7, maximum likelihood linear regression (MLLR) [2] was employed for updating both means and variances for each conversation side. Whereas one global MLLR transform was used in passes 3 and 4, the following stages used a maximum of 2, 4 and 8 transforms per side respectively.

The final stage combined two different system outputs according to a computed confidence score for each word. The confidence scores were generated using an N-best homogeneity measure found using the 1000-best hypotheses from the lattices generated at the appropriate stage. A decision tree pruned using 10-fold cross-validation was used to convert the N-best homogeneity scores to confidence probabilities. This decision tree was trained on the development data also using 10-fold cross validation. System output from the best triphone system (pass 3) and the best quinphone system (pass 7) were combined using ROVER [1].

3. FRONT-END EXPERIMENTS

For fast turnaround on front-end experiments, a small subset of the Swbd-I corpus was chosen for training. This subset (referred to as MiniTrain) covers 398 sides containing 17.8 hours of speech and is approximately gender balanced. For testing a gender balanced half-hour set (MTtest) containing Swbd-I data was chosen. All front-end experiments have been conducted using a 2MW Switchboard trigram language model. The following sections detail experiments that investigate analysis bandwidth, cepstral normalisation, and VTLN.

3.1. Coding Bandwidth and Cepstral Normalisation

Due to the special characteristics of telephone channels, the lower and upper frequency regions are either distorted or blocked by filtering operations. We compared systems using Mel-scale Filterbanks within the full 4kHz range and a reduced range between 125-3800 Hz.

Speech recognition systems designed for read speech or even Broadcast News data usually apply a per segment cepstral mean normalisation scheme to reduce the effects of constant channel characteristics. However, for telephone conversations the average utterance duration is less than 3 seconds, thus providing poor estimates for the segment means. To overcome this, the mean was calculated over a complete conversation side.

Results in Table 1 show the performance of HMM model sets trained on MiniTrain for different bandwidths using both of the cepstral mean subtraction strategies. Surprisingly the reduced bandwidth system performs worse on the MTtest set. Nevertheless both coding strategies show a gain of about 1%.

	0-4000Hz	125-3800Hz
Seg-CMN GI	46.58	47.33
Side-CMN GI	45.67	46.17

Table 1: % Word Error Rates (WER) for full and reduced bandwidth coding using models trained on the MiniTrain set and tested on MTtest

We also tested the performance of variance normalisation in conjunction with side-based mean normalisation using several different techniques. A standard segment based scheme was compared with side-based variance normalisation and normalisation using a time constant decay. Each feature vector component was normalised to obtain a target variance, which was chosen to be the overall test data variance. Linear side-based variance normalisation produced the best results.

	0-4000Hz	125-3800Hz
Side-CMN, Side-CVN, GI	44.82	44.35
Side-CMN, Side-CVN, GD	44.33	43.00

Table 2: % WER on MTtest using different bandwidth and gender independent (GI) and gender dependent (GD) MiniTrain model sets.

Table 2 shows the effect of side based variance normalisation using both full and reduced bandwidth coding on gender dependent and gender independent models. The performance gain is high especially for the reduced bandwidth case. Reduced bandwidth coding outperforms full bandwidth analysis further using gender dependent HMMs.

3.2. Maximum Likelihood Vocal Tract Length Normalisation

Maximum likelihood vocal tract length normalisation implements a per speaker linear frequency scaling of the speech spectrum. The scale factor is obtained using a search procedure and is then applied in speaker specific feature stream computation. The scaling can also be implemented by scaling the Mel filterbank centre frequencies with the inverse warping factor. Smoothing of the upper frequency filterbank contents is required when using scale factors larger than one. Instead of achieving this by mirroring the contents of the upper frequency contents [5], our new approach introduces a piecewise linear warping function with lower and upper cut-off frequencies (see Figure 1). The upper threshold improved the stability of our implementation (Table 3), while the lower cut-off frequency affected performance only slightly. Warp factors were found by conducting a parabolic search over data likelihoods versus warp factors. Given a previously obtained word level transcript, the average per-frame log-likelihood given some HMM model set was computed by feature recomputation and alignment with the transcripts. The next warping factor was

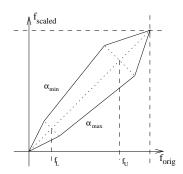


Figure 1: Piecewise linear VTLN frequency scaling function for warp factors α . f_L and f_U denote lower and upper threshold frequencies.

then selected and the procedure repeated. Since the per-frame loglikelihood tends to be a parabolic function of the warp factor, a suitable search method was chosen to allow rapid estimation of the warp factors.

	test	train & test
old	43.21	42.66
new	42.52	41.56

Table 3: % WER comparison for different VTLN implementations on MTtest using full bandwidth coding. Test denotes test-set VTLN only, train & test denotes single iteration VTLN models

Warp factors were computed using standard HMM models trained on a particular dataset for the use of VTLN in training. New models were generated by single pass retraining with the appropriately warped training data. Since multiple iterations of VTLN training are necessary to allow the warp factor distribution to converge, the models generated in one iteration serve as a warp factor estimator as well as the base for single pass retraining in the next iteration.

VTLN Experiments for both coding bandwidths are shown in Table 4. In the full bandwidth case the warp factor distribution settled after two iterations, whereas four iterations were necessary for the reduced bandwidth models. Even though full-band coding seems to perform well for GI models, the gain using GD models for reduced-band coding is 0.4% greater. Since the relative improvement is small, this result was cross-checked with another test set.

		0-4000Hz	125-3800Hz
GI	test	43.21	43.18
GD	test	41.92	41.33
GI	train/test	41.45	41.61
GD	train/test	40.75	40.20

Table 4: % WER for systems trained on MiniTrain and tested on MTtest. VTLN warping in training and test both for GI and GD models for full and reduced bandwidth coding.

4. SPEAKER ADAPTATION

The more robust implementation of VTLN together with gender dependent modelling reduced the improvement in word error rate achieved with MLLR speaker adaptation significantly. Our 1997 system achieved a gain of 4.8% on the 1997 Hub5E evaluation set using only global mean and variance speech transforms, and a further 1.6% by subsequent iterations with larger numbers of speech transforms. In comparison, only 2.5% improvement has been obtained using our 1998 front-end. The contribution of variance adaptation was only 0.2%, and the smallness of this figure may be attributed to the per-speaker variance normalisation. Subsequent MLLR iterations gave approximately similar improvements as for our 1997 system.

5. LANGUAGE MODELLING

Approximately 3 million words of Switchboard and Callhome English transcriptions were available for language model training (the h5trainLM set). From this, a 27k word recognition vocabulary containing only English words was determined. Furthermore, backoff bigram (bgH5), trigram (tgH5) and 4-gram (fgH5) models were trained from h5trainLM. To evaluate the effect of the increase in training data, a trigram tgH5_97 was built using the approximately 2 million words of Switchboard transcriptions used in our 1997 system.

Using the 27k wordlist, bigram (bgBN), trigram (tgBN) and 4-gram (fgBN) models were trained from Broadcast News data ranging in epoch from January 1992 to December 1997.

Corresponding H5 and Broadcast News models were merged by linear interpolation into a single resultant language model file, allowing them to be used directly in the recognition search. Thus bgH5 was merged with bgBN to form bgint98, tgH5 with tgBN to form tgint98, and fgH5 with fgBN to yield fgint98.

Finally, a class-based trigram language model (cat98) was produced using 350 automatically generated word classes based on word bigram statistics [3]. Bigrams and trigrams were only added if they improve the training set leave-one-out perplexity [4]. Both the categories as well as the trigram category model were built using only h4trainLM. An optimal interpolation (in terms of perplexity on the complete 1997 Hub5E evaluation set) was produced between fgH5, fgBN and cat98 with respective weights of 0.42, 0.28 and 0.30, and will be referred to as fgintcat98.

Table 5 displays the performance of these language models. Note that the 1997 NIST scoring conventions were used in WER calculation. The WER results for tgint98, fgint98, and fgintcat98 were obtained by rescoring lattices produced with bgint98.

LM	PP	WER
tgH5_97	98.3	-
tgH5	94.1	-
cat98	101.8	-
bgint98	101.7	45.8
tgint98	82.0	42.7
fgint98	79.2	42.3
fgintcat98	76.4	41.5

Table 5: Perplexity (PP) on eval97 and WER on eval97sub for various language models.

6. SYSTEM RESULTS ANALYSIS

Table 6 shows the performance of the individual stages on a subset of the 1997 Hub5E evaluation set (eval97sub) and the full 1998 Hub5E evaluation set (eval98). The eval97sub set was used for system development and consisted of 20 conversation sides from Swbd-II and CHE. This set was selected to give approximately the same performance as the full 1997 evaluation set. The eval98 set is gender balanced on Swbd-II data, but only contains 6 male speakers from CHE.

PASSES	Total	Swbd-II	CHE
P1	51.1	43.6	58.7
P2	44.6	36.5	52.8
P3	39.5	31.1	48.0
P4	38.1	29.9	46.4
P5	37.5	29.0	46.0
P6	37.3	29.1	45.6
P7	37.1	28.7	45.5
P8	36.6	28.5	44.7
(a)			

(a)			
PASSES	Total	Swbd-II	CHE
P1	49.3	47.0	51.6
P2	45.6	42.9	48.2
P3	42.6	39.9	45.3
P4	40.9	38.3	43.4
P5	40.5	37.9	43.2
P6	40.4	37.7	43.0
P7	40.3	37.7	42.8
P8	39.5	36.7	42.2

Table 6: % WER for the eval97sub set (a) and the eval98 (b) set for each decoding pass P1-P8. Word error rates are computed using the 1998 Hub5E scoring rules.

In the first pass (P1) a word level transcript has been obtained using M1 models and the tgint98 language model. In the second pass (P2) all VTLN warping factors for all sides were computed using gender dependent warp estimation models and the output from the first pass. Secondly the likelihood for the best warp factor for both genders were compared and gender selected according to the more likely model set. Whereas on eval97sub this gave no gender detection errors, this was not the case on eval98, where 3 sides out of 80 were misclassified.

Afterwards M2 models and the tgint98 language model were used to produce better MLLR supervision for the next stage. On the eval97sub gender dependent modelling plus VTLN brought a 6.5% gain in WER compared to only 3.7% on eval98.

In the third pass (P3) M2 models, MLLR speaker adaptation and the interpolated bigram model bgint98 were used to produce lattices, which were expanded using tgint98 and fgintcat98 language models. The gain of MLLR plus 4-gram language modelling was 5.1% on eval97sub compared to only 3% on eval98. The use of the M3 quinphone models with further MLLR passes brought 2.4% on eval97sub and 2.3% on eval98. The final system combination using ROVER [1] performed approximately equally on both sets with 0.8% in eval98 and only 0.6% on eval97sub.

The 1997 HTK system [6] on eval97sub had an error rate of

47.7%. The final word error rate (P8 in Table 6a) on the same data using the 1997 NIST scoring procedure was 38.5%, or a relative reduction in error rate of nearly 20%.

7. CONCLUSIONS

The 1998 conversational telephone speech transcription system has been described and shown significant gain in performance based on improved acoustic and language modelling concepts.

The improvements include reduced bandwidth analysis, sidebased cepstral feature normalisation, improved VTLN and SAT trained quinphone models. For language modelling 4-grams and 3-fold interpolation including a class-based model was used.

8. ACKNOWLEDGEMENTS

We would like to thank BBN for providing the MiniTrain training and MTtest test set definitions. This work was in part supported by GCHO.

9. REFERENCES

- Fiscus, J.G. (1997) A Post-Processing System to Yield Reduced Word Error Rates: Recogniser Output Voting Error Reduction (ROVER). *Proc. IEEE ASRU Workshop*, pp. 347-352, Santa Barbara.
- [2] Gales M.J.F. & Woodland P.C. (1996). Mean and Variance Adaptation Within the MLLR Framework. *Computer Speech* & *Language*, Vol. 10, pp. 249-264.
- [3] Kneser R., & Ney H. (1993). Improved Clustering Techniques for Class-Based Statistical Language Modelling. *Proc. Eurospeech'93*, pp. 973-976, Berlin.
- [4] Niesler T.R., Whittaker E.W.D & Woodland P.C. (1998). Comparison of Part-Of-Speech and Automatically Derived Category-Based Language Models for Speech Recognition. *Proc. ICASSP'98*, pp. 177-180, Seattle.
- [5] Pye D. & Woodland P.C. (1997). Experiments in Speaker Normalisation and Adaptation for Large Vocabulary Speech Recognition. *Proc. ICASSP'97*, pp. 1047-1050, Munich.
- [6] Woodland P.C., S. Kapadia, H.J. Nock & S.J. Young (1997). The HTK System for the March 1997 Hub 5 Evaluations. Presented at the *DARPA Hub5E Conversational Speech Recognition Workshop*, May 13-15,1997, Baltimore, Maryland.
- [7] Woodland P.C., Hain T., Johnson S.E., Niesler T.R., Tuerk A., Whittaker E.W.D. & Young S.J. (1998) The 1997 HTK Broadcast News Transcription System. *Proc. DARPA Broadcast News Transcription and Understanding Workshop*, pp. 41-48, Lansdowne, Virginia