

## **08/15/04 — 03/31/05: RESEARCH AND EDUCATIONAL ACTIVITIES**

The overall goal of this ITR was to create a strong synergy between speech recognition (ASR) and natural language processing (NLP). At the time this project began, integration of ASR and NLP was in its infancy, particularly for conversational speech applications. Over the duration of this project, two significant things happened. First, through the parallel efforts of DoD-funded research, community-wide focus on conversational speech was achieved. Progress was impressive as error rates on tasks such as Switchboard and Call Home English decreased from 50% to 10%. ASR technology was now producing transcripts that were useful to NLP systems, and could support information retrieval applications involving important quantities such as named entities.

Second, NLP research began to focus on the problem of parsing speech recognition output, which lacks punctuation and formatting that was previously considered crucial to high performance parsing. This latter issue was the main focus of this ITR, and to some extent served as a beacon for the community. We produced resources that were extremely valuable, such as the extensions to the Penn Treebank that were released in 2003 (reconciliation of the ISIP Switchboard segmentations and transcriptions with the Penn Treebank segmentations and transcriptions). We introduced the mainstream community to advanced statistical modeling techniques such as Support Vector Machines and enhanced these for NLP applications.

Further, in line with the primary goal of the ITR program, this project created close collaborations between groups who did not previously work together. The PIs collaborated on a number of new initiatives as offshoots of this project, including applications in parsing, information retrieval, and homeland security. A subset of the PIs participated in conversational speech evaluations and workshops (e.g., DARPA EARS). Hence, we can conclude that this project created new synergies and new research directions that will continue beyond the timeframe of this project.

In this final report, we briefly describe some of the significant findings of our research below.

### **A. Laboratory for Linguistic Information Processing, Brown University**

Learning general functional dependencies, i.e. functions between arbitrary input and output spaces, is one of the main goals in supervised machine learning. Recent progress has to a large extent focused on designing flexible and powerful input representations, for instance by using kernel-based methods such as Support Vector Machines. We have addressed the complementary issue of problems involving complex outputs such as multiple dependent output variables and structured output spaces. In the context of this project we have mainly dealt with the problem of label sequence learning, a class of problems where dependencies between labels take the form of nearest neighbor dependencies along a chain or sequence of labels. The latter is a natural generalization of categorization or multiclass-classification that has many applications in the context of natural language processing and information extraction. Special cases include part-of-speech tagging, named entity recognition, and speech-accent prediction. More specifically, we have developed and empirically investigated several extensions of state-of-the-art categorization algorithms such as AdaBoost, Support Vector Machines, and Gaussian Process classification. We have designed and implemented several scalable learning algorithms that combine standard optimization techniques employed in the context of the above mentioned methods with dynamic programming techniques that account for the nearest neighbor dependencies. Experimental evaluations on a wide variety of tasks have shown the competitiveness of these methods compared to existing techniques like Hidden Markov Models and Conditional Random Fields.

A second line of research conducted in the context of the present ITR has dealt with ways to systematically exploit class hierarchies and taxonomies. The main question we have investigated is

whether or not a priori knowledge about the relationships between classes helps in improving classification accuracy, in particular in cases with many classes and few training examples. This is highly relevant for applications like word sense disambiguation and text categorization, where the number of classes can easily be in the tens of thousands. To that extend we have focused on a hierarchical version of the well-known perceptron learning algorithm as well as an extension of multiclass Support Vector Machines. We have shown that this approach can be effective in situations with sparse training data.

## B. Center for Language and Speech Processing, Johns Hopkins University

The Structured Language Model (SLM) aims at making a prediction of the next word in a given word string by making a syntactical analysis of the preceding words. However, it faces the data sparseness problem because of the large dimensionality and diversity of the information available in the syntactic parses. A neural network model is better suited to tackle the data sparseness problem and its use has been shown to give significant improvements in perplexity and word error rate over the baseline SLM (Emami et al, 2003).

In this work we have investigated a new method of training the neural net based SLM. Our model makes use of a neural network for that component of the SLM that is responsible for predicting the next word given the previous words and their partial syntactic structure. We have investigated both a mismatched and a matched training scenario. In matched training, the neural network is trained on partial parses similar to those that are likely to be encountered during evaluation. On the other hand in the mismatched scenario, faster training time is achieved but at the cost of mismatch between training and evaluation and hence, possible degradation in performance.

The Structured Language Model works by assigning a probability  $P(W,T)$  to every sentence  $W$  and every possible binary parse  $T$  of  $W$ . The joint probability  $P(W,T)$  of a word sequence  $W$  and a complete parse  $T$  is broken into:

$$P(W,T) = \prod_{k=1}^{n+1} P(W_k | W_{k-1}T_{k-1}) \cdot P(t_k | W_{k-1}T_{k-1}, W_k) \cdot \prod_{i=1}^{N_k} P(p_i^k | W_{k-1}T_{k-1}, w_k, t_k, p_1^k \cdots p_{i-1}^k)$$

where  $W_{k-1}T_{k-1}$  is the word-parse (k-1)-prefix,  $t_k$  is the tag assigned to  $w_k$  by the TAGGER,  $N_k - 1$  is the number of operations the CONSTRUCTOR executes at sentence position k before passing control to the PREDICTOR, and  $p_i^k$  denotes the i-th CONSTRUCTOR operation carried out at position k in the word string.

Subsequently, the *language model* probability assignment for the word at position k+1 in the input sentence is made using:

$$P_{\text{SLM}}(w_{k+1} | W_k) = \sum_{T_k \in S_k} P(w_{k+1} | W_k T_k) \cdot \rho(W_k T_k)$$

$$\rho(W_k T_k) = P(W_k T_k) / \sum_{T_k \in S_k} P(W_k T_k)$$

which ensures a proper probability normalization over strings  $W^*$  where  $S_k$  is the set of all parses built and retained by the model at the current stage k.

Neural networks are very suitable for modeling conditional discrete distribution with large vocabularies. These models work by first assigning a continuous feature vector with every token in the vocabulary, and then using a standard multi-layered neural net to get the conditional distribution at the output, given the input feature vectors. Training is achieved by searching for parameters  $\Theta$  of the neural network and the values of feature vectors that maximize the penalized log-likelihood of the training corpus:

$$L = \frac{1}{N} \sum_t \log p(y^t | x_1^t, x_2^t, \dots, x_{n-1}^t; \Theta) - R(\Theta)$$

where  $p(y^t | x_1^t, x_2^t, \dots, x_{n-1}^t; \Theta)$  is the probability of word  $y^t$  (network output at time t), N is the training data size and  $R(\Theta)$  is a regularization term, L-2 norm squared of the parameters in our case.

We have used a neural net to model the SCORER component of the SLM. By the SCORER we refer to the model  $P(w_{k+1} | W_k T_k)$ . The neural net SCORER's parameters can be obtained by training it on the events extracted from the gold standard (usually one best) parses obtained from an external source (humans or an automatic parser). However, there would be a mismatch during evaluation since the partial parses during that phase are not provided and have to be hypothesized by the SLM itself. We have called the SCORER trained in this manner the *mismatched* SCORER.

On the other hand, one can train the model on partial parses hypothesized by the baseline SLM, thus maximizing the proper log-likelihood function. We have called this procedure the *matched* training of the SCORER.

Experimental results have shown considerable improvement in both perplexity and WER when using a neural net based SLM, specially in the case of matched SCORER training. On the UPenn section of the WSJ corpus, perplexity reductions of 12% and 19% over the baseline SLM (with a perplexity of 132) have been observed when using the mismatched and matched neural net models respectively.

For the WER experiments, the neural net bases models were used to re-rank an N-best list output by a speech recognizer on the WSJ DARPA'93 HUB1 test set (with a 1-best WER of 13.7%). The mismatched and matched neural net models reduced the SLM baseline WER of 12.6% to 12.0% and 11.8% (for relative improvements of 4.8% and 6.3%) respectively.

In summary, neural network models showed to be capable of taking advantage of the richer probabilistic dependencies extracted through syntactic analysis. In our case the use of a neural net for the SCORER component of the Structured Language Model resulted in considerable improvements in both perplexity and Word Error Rate (WER) with the best results achieved when using a training procedure matched with the evaluation.

## **C. Signal, Speech, and Language Interpretation Lab, University of Washington**

Prosody can be thought of as the "punctuation" in spoken language, particularly the indicators of phrasing and emphasis in speech. Most theories of prosody have a symbolic (phonological) representation for these events, but a relatively flat structure. In English, for example, two levels of prosodic phrases are usually distinguished: intermediate (minor) and intonational (major) phrases [Beckman & Pierrehumbert, 1986]. While there is evidence that both phrase-level emphasis (or, prominence) of words and prosodic phrases (perceived groupings of words) provide information for syntactic disambiguation [Price et al., 1991], the most important of these cues seems to be the prosodic phrases or the boundary events marking them. While prior work has looked at the use of prosody in automatic parsing of isolated sentences, a key component of our work involved sentence detection as well, since our goal is to handle continuous conversational speech. Hence, the focus of our work has been on automatically recognizing sentence boundaries and sentence-internal prosodic phrase structure and investigating methods for integrating that structure in parsing.

To support these efforts, we also worked on analysis of acoustic cues to prosodic structure. The most important (and best understood) acoustic correlates of prosody are continuous-valued, including fundamental frequency (F0), energy, and duration cues. Particularly at phrase boundaries, cues include significant falls or rises in fundamental frequency accompanied by word-final duration lengthening, energy drops, and optionally a silent pause. In addition, however, there is evidence of spectral cues to prosodic events, so some of our work explored these cues, which also have implications for improving speech recognition.

Our approach to integrating prosody in parsing is to use symbolic boundary events that have categorical perceptual differences, including prosodic break labels that build on linguistic notions of intonational phrases and hesitation phenomena but also higher level structure. These events are predicted from a combination of the continuous acoustic features, rather than using the acoustic features directly, because the intermediate representation simplifies training with high-level (sparse) structures. Just as most automatic speech recognition (ASR) systems use a small inventory of phones as an intermediate level between words and acoustic feature to have robust word models (especially for unseen words), the small set of word boundary events are also useful as a mechanism for generalizing the large number of continuous-valued acoustic features to different parse structures. This approach is currently somewhat controversial because of the high cost of hand labeling, and to some extent because of its association with a particular linguistic theory. However, the specific subset of labels used in this work are relatively theory neutral and language independent, and a key contribution of this work is the use of weakly supervised learning to reduce the cost of prosodic labeling.

An alternative approach, as in [Noth et al, 2000], is to assign categorical "prosodic" labels defined in terms of syntactic structure and presence of a pause, without reference to human perception, and automatically learn the association of other prosodic cues with these labels. While this approach has been very successful in parsing speech from human-computer dialogs, we expect that it will be problematic for conversational speech because of the longer utterance and potential confusion between fluent and disfluent pauses.

### **C.1 Data, Annotation and Development**

The usefulness of speech databases for statistical analysis is substantially increased when the database is labeled for a range of linguistic phenomena, providing the opportunity to improve our understanding of the factors governing systematic suprasegmental and segmental variation in word forms. Prosodic labels and phonetic alignments are available for some read speech corpora, but only a few limited samples of spontaneous conversational speech have been prosodically labeled. To fill this gap, a subset of the

Switchboard corpus of spontaneous telephone-quality dialogs was labeled using a simplification of the ToBI system for transcribing pitch accents, boundary tones and prosodic constituent structure [Pitrelli et al., 1994]. The aim was to cover phenomena that would be most useful for automatic transcription and linguistic analysis of conversational speech. This effort was initiated under another research grant, but completed with partial support from this NSF ITR grant.

The corpus is based on about 4.7 hours of hand-labeled conversational speech (excluding silence and noise) from 63 conversations of the Switchboard corpus and 1 conversation from the CallHome corpus. The orthographic transcriptions are time-aligned to the waveform using segmentations available from Mississippi State University (<http://www.cavs.msstate.edu/hse/ies/projects/switchboard/index.html>) and a large vocabulary speech recognition system developed at the JHU Center for Language and Speech Processing for the 1997 Language Engineering Summer Workshop ([www.clsp.jhu.edu/ws97](http://www.clsp.jhu.edu/ws97)) [Byrne et al., 1997]. All conversations were analyzed using a high quality pitch tracker [Talkin, 1995] to obtain F0 contours, then post-processed to eliminate errors due to crosstalk. The prosody transcription system included: i) breaks (0-4), which indicate the depth of the boundary after each word; ii) phrase prominence (none, \*, \*?); and iii) tones, which indicate syllable prominence and tonal boundary markers. It also provides a way to deal with the disfluencies that are common in spontaneous conversational speech (a p diacritic associated with the prosodic break), and to indicate labeler uncertainty about a particular transcription. The annotation does not include accent tone type, primarily to reduce transcription costs and because we hypothesized that the phrase tones would be relevant to dialog act representations which may be relevant for question-answering. For further information on the corpus and an initial distributional analysis, see [Ostendorf et al. 2001].

The prosodically labeled subset of Switchboard overlaps with the subset of that corpus annotated with Treebank parses, but there is a mismatch in the orthographic transcriptions because the Treebank parses were based on an earlier version of transcripts and the prosodic annotation was based on the higher quality corrections done by Prof. Picone's group (ISIP) at Mississippi State. In addition, we made use of the DARPA EARS metadata annotations that overlapped with the Treebank parses, which were again based on the higher quality transcriptions. To be able to use all of these resources, we used an alignment of words provided by the ISIP team, and mapped the Treebank parse information to the more recent word transcriptions, which could then be aligned with the EARS metadata annotations. Differences in transcriptions were handled by: dropping the parse information for deletions, transferring it as is for word substitutions, and treating it as "missing" information for insertions in the corrected transcripts. While most of the differences between the Treebank and corrected word transcriptions involved simple substitutions (or deletions) that had little or no impact on the parse (e.g. "a" vs. "the"), there were some cases where the transfer introduced noise into the collection of parses. The most frequent such cases were in disfluent regions, where transcribers tend to have more difficulties, including missed word fragments or repetitions ("I I" vs. "I I I"). An additional difference between the Treebank parses and the EARS metadata annotations is the marking of sentence boundaries. Since speakers frequently begin sentences with conjunctions, the metadata conventions often split up constituents marked as compound sentences in Treebank. Because the metadata labelers listened to the speech and the Treebank labelers did not, we chose to use the metadata constituents, which in most cases involved simply dropping a top-level (S) node, but in some cases involved adding a top-level node called "SUGROUP".

## **C.2 Automatic Labeling of Prosodic Structure**

An important part of the effort was development of an automatic prosodic labeling system that would provide cues to improve parsing. In addition, the resulting system was inspected to analyze possible dependencies between prosodic and parse structures in conversational speech. In the experiments, we used decision tree classifiers with different combinations of acoustic, punctuation, parse, and disfluency cues. While more sophisticated techniques, such as HMMs and maximum entropy models, have been

used for related tasks of sentence boundary detection (see [Liu et al., 2005] for a brief survey), we chose decision trees because they are easy to inspect for learning about the prosody-syntax relationship and because this simplified the weakly supervised learning experiments, which were the focus of our efforts.

For the prosody/syntax analyses, we designed trees to predict prosodic labels from syntactic structure, as well as trees to predict prosodic structure from a combination of syntactic and acoustic cues. For purposes of providing information to a parser, we designed trees to predict prosodic constituents from acoustic cues and part-of-speech (POS) tags, but as an intermediate step in designing these trees we also used syntactic cues in designing trees as part of the weakly supervised training. More specifically, a small set of labeled data was used to train prosody models based on both text and acoustic cues, which were then used in combination to automatically label a large set of data that had not been hand-annotated with prosodic structure, and finally new (separate) acoustic-based prosody models were designed from this larger data set for use in parsing new data.

Experiments were conducted on the Switchboard corpus, using the prosodically annotated subset described above for initial training and evaluation (independent subsets for each). Then the full Switchboard training set was incorporated using various methods for weakly supervised learning, as described below. The prosodic constituent labels were merged into 3 classes: major intonational phrase boundary (4), hesitation boundary (1p, 2p), and all other fluent word boundaries. We grouped minor intonational phrase boundaries (3) with the default word boundary class, because preliminary experiments showed that they were almost never predicted by the decision trees (even with sampled training to account for the low frequency) and because they were most often confused with the default word class in 4-class prediction experiments. The simple 3-class system also has the advantage that it is relatively theory neutral and language independent in that essentially all languages have a notion of fluent and disfluent segmentation.

The acoustic cues included normalized F0, energy and duration cues based on those used in [Kim et al., 2004] and similar to those used in other metadata detection studies [Shriberg et al., 2000]. Text-based cues -- including punctuation, parse structure and disfluency markers -- were taken from the annotations available from the Linguistic Data Consortium, aligned to the word transcriptions as described above. The parse features included right-to-left and left-to-right relative depth, part-of-speech, label of the left and right syntactic constituents, and parenthetical markers. In addition, disfluency interruption points and flags for filled pauses and sentence-initial conjunctions were used as features. Punctuation as inserted by a human transcriber (including incomplete sentences) and estimated speaker turn boundaries (defined simply as a word boundary with a silence of length greater than 4s) were also used.

The baseline system trained on hand-labeled data has error rates of 9.6% when all available cues are used (both syntax and prosody) and 16.7% when just acoustic and POS cues are used. This can be compared to an error rate of 30% when the default class is assigned to all word boundaries. We considered three different weakly supervised training techniques for adding data without prosodic labels (but with hand-labeled syntactic structure) into the training set: EM, co-training, and self-training. The co-training algorithm used classifiers designed on either acoustic or syntactic cues, and it differed slightly from the standard method in that we used an information-theoretic distance on the tree posteriors to determine when to omit samples with conflicting classifier decisions. The self-training algorithm used bagging with uniform class sampling to deal with data skew [Liu et al., 2004]. In all cases, only 1-2 iterations were needed. Both the co-training and EM approaches gave improved performance over the baseline, with the

EM algorithm giving the best results of 14.2% error for the acoustic-only trees, which corresponds to a 15% reduction in error rate over supervised training. The self-training strategy actually hurt performance.<sup>1</sup>

From analysis of the resulting trees, we find that punctuation and repair points are correlated with prosodic breaks and hesitations, as expected. Silence duration is the most useful individual acoustic feature, but alone it is not a reliable predictor of prosodic phrases, since there are many prosodic phrases that do not occur at silences and because silences are frequently associated with hesitations. Aside from syntactic structure, the most important text features for predicting prosodic constituents are punctuation, disfluency edit point markers, filler words (sentence-initial coordinating conjunctions, discourse markers, filled pauses), and turn boundaries. Some important syntactic features include depth of subtrees on the left and right sides of the boundary, previous and next syntactic constituent tag, length of closing phrase, and part-of-speech tags. These features were relevant when associated with the target word boundary, but frequently also with the next or previous word boundary. Surprisingly, the label of the joining constituent is not useful. This analysis provided input into the parse reranking work described in the next section.

Due to the success of the weakly supervised training on prosodic phrase boundary detection, we have recently started investigating use of the same technique for training models of prosodic prominence. Initial results show only a 4% reduction in error rate for the system based on acoustic cues, from 22% to 21% error. Despite the high error rate, however, the automatically annotated prominence appears to be useful in topic identifications in preliminary experiments associated with a separate NSF project (IIS-0121396).

### **C.3 Prosody and Parsing**

Parsing speech can be useful for a number of tasks, including information extraction and question answering from audio transcripts. However, parsing conversational speech presents a different set of challenges than parsing text: sentence boundaries are not well-defined, punctuation is absent, and presence of disfluencies (edits and restarts) impact the structure of language. Most prior work on parsing conversational speech has focused on handling disfluencies [Hindle 1983; Mayfield 1995; Charniak & Johnson, 2001], but experiments relied on hand-marked sentence boundaries and made use of punctuation as in text-based parsers. While utterance-level segmentation may be reasonable to assume in current human-computer dialog systems, it is not realistically available in recognized conversational speech. Hence, our work looked at the problem of parsing text with disfluencies and without punctuation.

We have investigated three main issues in the use of prosody in parsing: the impact of automatic sentence segmentation, the usefulness of interruption points, and the usefulness of automatically detected sub-sentence prosodic constituent boundaries (described above). In all cases, we use a two-stage architecture where metadata (constituent boundaries) are first detected with a combination of prosodic and simple text features, and then these symbolic events (or their posterior probabilities) are used in parsing. Our approach focuses on categorical boundary events, which are predicted from a combination of acoustic features, rather than using the acoustic features directly. As argued earlier, the intermediate representation simplifies training with sparse structures. Key research issues include whether the metadata should be treated as "words" or as features on words, whether edits should be represented with an independent component, and how to represent uncertainty of the metadata classifiers. Our work has begun investigating all of these questions, but some remain unanswered and are being pursued in ongoing work.

---

<sup>1</sup> The results reported here are in some cases worse than those reported in an earlier progress report, because they are based on a larger data set. Due to a data processing bug, several files were omitted from earlier studies. In addition, because of the larger amount of data used and the richer feature sets, the trees are much larger than those described in prior reports.

The data used in this work is the Treebank portion of the Switchboard corpus of conversational telephone speech, which includes sentence-like unit boundaries (SUs) as well as the reparandum and interruption point of disfluencies. The data consists of 816 hand-transcribed conversation sides (566K words), of which we reserve 128 conversation sides (61K words) for evaluation testing according to the 1993 NIST evaluation choices. In all cases, training was based on hand-labeled SUs. Parses were evaluated using SU boundaries rather than the standard punctuation-based units that the Treebank is based on, so the gold standard parses and parse evaluation metric were modified to incorporate the SUs.

The most exhaustive series of experiments looked at the impact of automatic segmentation on parsing. For this particular effort, we chose to work with the complete word sequence, i.e. including all of the words within edit regions, to allow experimentation with multiple parsers. In initial work [Kahn et al., 2004], we used the structured language model (SLM) as a parser with a simple pause-based segmentation and automatically detected SUs (69% vs. 35% slot error rate, respectively), showing a significant improvement in parsing performance when using the automatic SUs. We then confirmed the findings with results on the same segmentation alternatives but using two state-of-the-art statistical parsers trained on conversational speech: the Bikel [Bikel, 2004]<sup>2</sup> and Charniak [Charniak & Johnson, 2001]<sup>3</sup> parsers. Figure 1 shows the relative performance of each parser for the two different segmentations, comparing to the oracle performance for each using hand-annotated SUs. In order to have a single number for performance, we use the F-measure calculated from bracket precision and recall. (Trends with separate precision and recall measures and with bracket crossing statistics follow the same patterns.) For all three parsers, there is a significant gain in performance due to using the explicit SU detection algorithm, with more than half of the performance loss associated with the pause-based segmenter recovered when moving to the more sophisticated SU detection system. As SU detection improves, we would expect further performance gains. Also important is the observation that the impact of increased SU error is not lessened by using a better parser: the best oracle results were achieved with the Bikel parser, but it was much more sensitive to SU errors than the SLM.

In trying to understand the role of IPs in parsing, we ran several experiments, including some with and without human-annotated punctuation. Without punctuation, there was a small increase in parsing performance of the SLM using IPs. When the SUs are automatically detected. We were not able to confirm these gains with other parsers; however, recent work in [Johnson, Charniak & Lease, 2004] shows a benefit to edit detection from using IPs which presumably would lead to improved parsing in their two-stage processing strategy [Johnson & Charniak, 2004]. Including punctuation and IPs in experiments with the SLM showed an

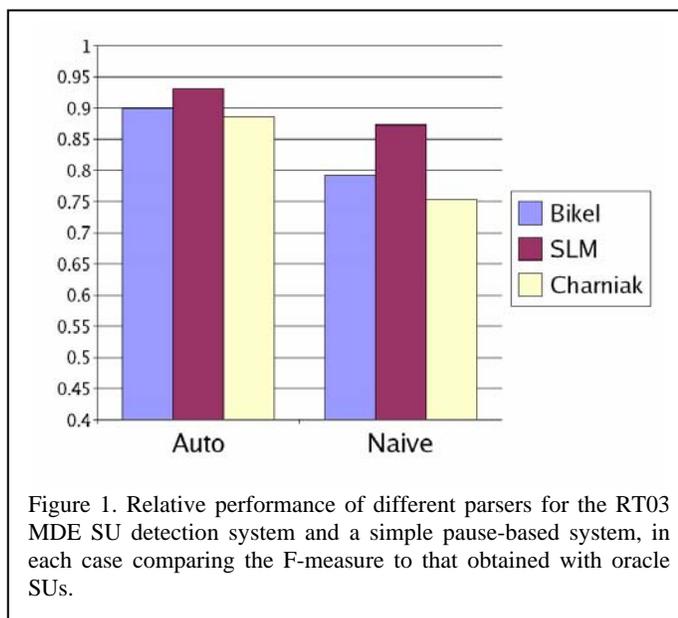


Figure 1. Relative performance of different parsers for the RT03 MDE SU detection system and a simple pause-based system, in each case comparing the F-measure to that obtained with oracle SUs.

<sup>2</sup> <http://www.cis.upenn.edu/~dbikel/download.html> (Version 0.9.9). For this work, we trained the Bikel parser on the Switchboard Treebank parses with the Collins settings.

<sup>3</sup> <ftp://ftp.cs.brown.edu/pub/nlparser/> (August 2004)

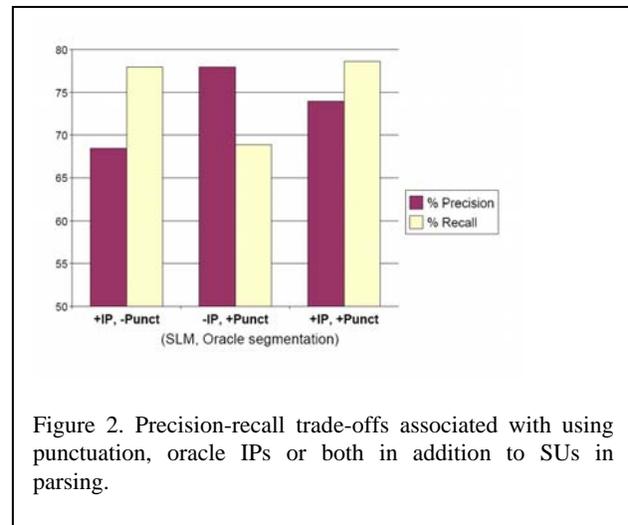
interesting trade-off of bracket precision and recall using the two different cues, as illustrated in Figure 2. We saw the improved precision associated with using both punctuation and IPs as possible evidence that sub-sentence prosodic constituents might be useful.

In all of the above work, metadata is incorporated as "word" tokens, similar to the standard mechanism for parsers to incorporate punctuation. For sentence segmentation with a reasonably reliable segmenter, this may make sense, but certainly for sub-sentence prosodic constituents there is the potential for the gains associated with adding prosody to be offset by a loss from the extra words blocking part of the history that might be used in a statistical model of word dependence. We conjecture that this may in part explain the negative results obtained in [Gregory et al., 2004], since our analyses of the prosody prediction trees provides some evidence that sub-sentence prosodic constituents may be useful in parsing. (The direct use of acoustic features may also be problematic.) In addition, the use of metadata events as "words" requires a hard decision in the first stage of detection, and many results in speech processing suggest that soft decisions (e.g. using class posteriors) are more effective.

To address these problems, we developed an extension to the SLM that uses prosodic constituents as hidden conditioning variables, similar to headword conditioning in the SLM. However, since our subsequent work obtained much better baseline performance with other parsers, we decided to explore a parse reranking framework [Johnson et al., ms. in prep.] as an alternative method for incorporating automatically detected prosodic constituents. The approach uses a maximum entropy reranking model and introduces new features based on counts of syntactic constituent types weighted by the posterior probability of different prosodic events. Experiments with this new approach are in progress, now under other funding, and we anticipate having results in early 2005. This series of experiments will also look at the question of whether a separate stage of edit detection benefits parsing compared to simply incorporating the edit structure in the parser with the same status as other constituents.

#### C.4 Prosody and Acoustic Modeling

Most research on the use of prosody in automatic speech processing has focused on F0, energy and duration correlates to prosodic structure. However, there is evidence from long standing acoustic, articulatory and perceptual studies of speech suggesting that there are spectral correlates as well. For that reason, we conducted an analysis of our prosodically labeled conversational speech data using acoustic parameters and clustering techniques that are standard in speech recognition. We found that prosodic factors are associated with acoustic differences that can be learned in standard speech recognition systems. Both prosodic phrase structure and phrasal prominence seem to provide distinguishing cues, with some phones being affected much more than others (as one would expect from the linguistics literature). We hypothesized that we would find that constituent onsets were important at all levels (syllable, word and prosodic phrase). Instead, we found that onset is more important for syllables, but constituent-final position is more important at higher levels. Prosodic prominence had a smaller affect than phrase structure in terms of increasing likelihood of the training data, but seemed to result in more separable models when it did play a role.



Finally, we found evidence that segmental cues can help distinguish fluent from disfluent phrase boundaries, in that segments associated with these categories are frequently placed in different clusters. These differences can be leveraged in a “multiple pronunciation” acoustic model to aid in detecting fluent vs. disfluent prosodic boundaries, though additional prosodic cues are necessary to separate these from unmarked word boundaries. A limitation of this work was that it was based on hand-labeled data, and therefore did not take advantage of the full training data set needed for designing a state-of-the-art recognition system. However, with our recent developments in prosodic annotation, we will be able to assess the usefulness on a much larger corpus in the future.

#### **D. Institute for Signal and Information Processing, Mississippi State University**

Hidden Markov models (HMMs) with Gaussian emission densities are the prominent modeling technique in speech recognition. HMMs suffer from an inability to learn discriminative information and are prone to overfitting and overparameterization. Recent work in machine learning has focused on models, such as the support vector machine (SVM), that automatically control generalization and parameterization as part of the overall optimization process. SVMs, however, require ad hoc (and unreliable) methods to couple it to probabilistic speech recognition systems. We have applied a probabilistic Bayesian learning machine termed the relevance vector machine (RVM) as the core statistical modeling unit in a speech recognizer. The RVM is shown to provide superior performance compared to HMMs and SVMs in terms of both accuracy and sparsity on a continuous alphanum digit task.

Computer speech recognition is typically posed as a statistical pattern recognition problem based on Bayesian methods in which the acoustic model and language model are treated as separate statistical models. The focus of our work has been the acoustic model, which maps sequences of feature vectors to probabilities that these vectors were produced by a given linguistic unit, such as phone. In most state-of-the-art recognition systems, a hidden Markov model (HMM) is used as the acoustic model. The popularity of the HMM representation is based on an HMM's ability to simultaneously model the temporal progression of speech (speech is usually seen as a “left-to-right” process) and the acoustic variability of the speech observations. The temporal variation is modeled via an underlying Markov process while the acoustic variability is modeled by an emission distribution at each state in the Markov chain.

The most commonly used emission distribution is the Gaussian mixture model (GMM). While combinations of HMMs and Gaussian mixture models have been extremely successful, there are two fundamental limitations of these approaches: (1) the parametric form of the underlying distribution is assumed to be Gaussian, (2) the maximum likelihood (ML) approach, which is typically used to estimate model parameters, does not explicitly improve the discriminative ability of the model. The ML approach maximizes the probability of the correct model while implicitly ignoring the probability of the incorrect model. Ideally, a training approach should force the model toward in-class training examples while simultaneously driving the model away from out-of-class training examples. Methods such as maximum mutual information and minimum classification error have been developed to incorporate discriminative training directly into the standard HMM/GMM framework. However, their success has been limited due primarily to their considerable computational costs.

The weaknesses of the HMM/GMM system have led researchers to explore other models, such as hybrid connectionist systems, which merge the power of artificial neural networks (ANNs) and HMMs. HMM/ANN hybrids have shown promise in terms of performance but have not yet found widespread use due to some serious problems. First, ANNs are prone to overfitting the training data if allowed to blindly converge. To avoid overfitting, a cross-validation set is often used to define a stopping point for training. This is wasteful of data and resources — a serious consideration in speech where the amount of labeled training data is limited. ANNs also typically converge much slower than HMMs. Most importantly, the

HMM/ANN hybrid systems have not shown the substantial improvements in recognition accuracy over HMM/GMM systems that would be necessary to force a paradigm shift.

The approaches developed in this work draw significantly from the HMM/ANN hybrid systems. However, we seek methods which are automatically immune to overfitting without the artificial imposition of a cross-validation set as well as methods which can automatically learn the appropriate model structure as part of the overall optimization process. One such model that has come to the forefront of pattern recognition research is the support vector machine (SVM). The support vector paradigm is based upon structural risk minimization (SRM) in which the learning process is posed as one of optimizing some risk function. The optimal learning machine is the one whose free parameters are set such that the risk is minimized.

Finding a minimum of the risk function is typically impossible due to the unknown distribution. Instead, it has been shown that a relationship exists between the actual risk, which is related to the empirical risk (i.e. the training set error which can be measured) and the Vapnik-Chervonenkis (VC) dimension. The VC dimension is a measure of the capacity of a learning machine to learn any training set and is typically closely related to the complexity of the learning machine's structure. With this result, we can guarantee both a small empirical risk (training error) and good generalization — an ideal situation for a learning machine.

In their most basic form SVMs use the SRM principle to impose an order on the optimization process by ranking candidate separating hyperplanes based on the margin they induce. For separable data, the optimal linear hyperplane is the one that maximizes the margin. The true power of the SVM, however, lies in how it deals with nonlinear class separating surfaces. Providing for a nonlinear decision region is accomplished using kernels. The optimization process yields a decision function where the sign of can be used to classify examples as either in-class or out-of-class. The decision function is formed from only those training vectors that lie on the margin or in overlap regions. In practice, the proportion of the training set that becomes support vectors is small, making the classifier sparse. Consequently, the training process, along with the training set, directly optimize the complexity of the learning machine. In contrast, ANN systems often make *a priori* assumptions about the form of the model.

SVMs have had great success on static classification tasks. However, it is only recently, that these techniques have been applied to continuous speech recognition. While the SVMs provide an excellent classification paradigm, they suffer from two serious drawbacks that hamper their effectiveness in speech recognition. First, while sparse, the size of the SVM models (number of non-zero weights) tends to scale linearly with the quantity of training data. For a large speaker independent corpus this effect is prohibitive. Second, the SVMs are binary classifiers. In speech recognition this is an important disadvantage since there is significant overlap in the feature space which can not be modeled by a yes/no decision boundary. Further, the combination of disparate knowledge sources (such as linguistic models, pronunciation models, acoustic models, etc.) requires a method for combining the scores produced by each model so that alternate hypotheses can be compared. Thus, we require a probabilistic classification which reflects the amount of uncertainty in our predictions.

We have investigated a Bayesian model termed the relevance vector machine (RVM) which is similar in form to the SVM but which addresses these two problems. The essence of an RVM is a fully probabilistic model with an automatic relevance determination prior over each model parameter. Thus, sparseness in the RVM model is explicitly sought in a probabilistic framework.

## D.1 Sparse Bayesian Methods

Supervised learning in speech recognition implemented via a maximum likelihood approach is the dominant approach for finding values of the parameters in our model that best match the training data. Our expectation in data modeling is that given sufficient training data, the model would generalize to unseen test sets. Two levels of inference must be implemented to accomplish this. First, assuming that a particular model is true, we seek to infer the values for the parameters of the model that best fit the data at hand. This is exactly the training process used in the ANN hybrids. Second, we must decide which model is most appropriate given the data at hand, i.e. model comparison.

A simple approach to model comparison might dictate that we simply choose the model that fits the data best — the maximum likelihood solution. However, a more complex model can always fit the data better. Jaynes describes an extreme interpretation of this problem where we would always choose the so-called *Sure Thing* hypothesis, under which exactly the training set and only the training set is possible. Though it is the maximum likelihood solution, the *Sure Thing* hypothesis is intuitively displeasing and is counter to our desire for a solution which generalizes. We avoid choosing the *Sure Thing* hypothesis by expressing an *a priori* preference for simpler solutions using the principle of Occam's Razor. MacKay and others have formalized this preference mechanism through the use of Bayesian methods. These provide a natural and quantitative embodiment of Occam's razor. The first level of inference requires that we find the best-fit parameters. The second level of inference requires the comparison of competing hypotheses. If we assume that the competing hypotheses are *a priori* equiprobable then the best hypothesis is chosen by evaluating the evidence. The evidence is computed by marginalization across the model parameters.

The evidence provides a natural trade-off between the best-fit likelihood and the Occam factor. This concept is closely related to other methods such as the Minimum Description Length and the Bayesian Information Criteria where the model is directly penalized by the number of parameters used. A similar idea was also incorporated into SVM models, which penalize the models with too large a capacity (VC dimension). However, while the SVM models are forced to estimate the penalty via cross-validation schemes, Bayesian techniques automatically determine and apply the penalty in a fully probabilistic framework.

Assuming we have no prior knowledge that would cause us to favor a particular prior, we can find the optimal value for by evaluating the evidence. If we did have prior knowledge, we would simply repeat the inference over using the prior. At some level of the inference, we will arrive at a point where our prior knowledge is too weak to apply and then we evaluate the evidence. This extreme application of the Bayesian prior as a control parameter is known as the method of *automatic relevance determination* (ARD). It is so named because the prior over the input unit weights in a neural network can 'shut-off' those input dimensions which are irrelevant to the problem at hand. ARD is at the heart of the relevance vector machines that will be described next.

## D.2 Relevance Vector Machines

An application of the evidence framework to kernel machines is the relevance vector machine (RVM). As with SVMs, RVMs use a weighted linear combination of basis functions. Due to the large number of parameters in this model — one per observation — we must guard against overfitting of the model to the training data. SVMs use a control parameter to implicitly balance the trade-off between training error and generalization. RVMs take a Bayesian approach and explicitly define an ARD prior distribution over the weights.

Each weight in the RVM model has an individual hyperparameter,  $\alpha_i$ , that is iteratively reestimated as part of the optimization process. As the hyperparameter grows larger, the prior on  $w_i$  becomes infinitely peaked around zero, forcing  $w_i$  to go to zero and, thus, contributing nothing to the summation. This process

automatically embodies the principle of Occam’s Razor because it explicitly seeks the simplest model that satisfies the data constraints. In practice, the majority of the weights are pruned, resulting in an exceedingly sparse model with generalization abilities on par with SVMs. To complete the Bayesian specification of the model, we have to specify a prior probability. In practice we use a non-informative (flat) prior to indicate a lack of preference.

The parameter estimation process is an iterative reduction process. That is, initially each vector of the system is allocated one parameter. As the procedure continues, vectors are pruned from the model when they are found to be irrelevant with respect to the remaining parameters. Integral to this iterative reestimation process is the computation of the inverse Hessian matrix. This operation requires the inversion of an  $M \times M$  Hessian matrix where  $M$  is initially set to the size of training set. For larger training sets (on the order of a few thousand), this computation is prohibitive both in time and memory.

### D.3 Experiments

RVMs have had significant success in several classification tasks. These tasks have, however, involved relatively small quantities of static data. Speech recognition, on the other hand, involves processing a very large amount of temporally evolving signals. In order to gain insight into the effectiveness of RVMs for speech recognition, we explored two tasks. We first experimented on the Deterding static vowel classification task which is a common benchmark used for new classifiers. Second, we applied the techniques described above to a complete small vocabulary recognition task. Comparison with SVM models are given below. For each task, the RVMs outperformed the SVM models both in terms of model sparsity and error rate.

In our first pilot experiment, we applied SVMs and RVMs to a publicly available vowel classification task, Deterding Vowels. This was a good data set to evaluate the efficacy of static classifiers on speech classification data since it has been used as a standard benchmark for several nonlinear classifiers for several years. In this evaluation, the speech data was collected at a 10 kHz sampling rate and low pass filtered at 4.7 kHz. The signal was then transformed to 10 log-area parameters, giving a 10 dimensional input space. A window duration of 50 msec was used for generating the features. The training set consisted of 528 frames from eight speakers and the test set consisted of 462 frames from a different set of seven speakers. The speech data consisted of 11 vowels uttered by each speaker in a h\*d context. Though it appears to be a simple task, the small training set and significant confusion in the vowel data make it a very challenging task.

Table 1 shows the results for a range of nonlinear classification schemes on the Deterding vowel data. From the table, the SVM and RVM are both superior to nearly all other techniques. The RVM achieves performance rivaling the best performance reported on this data (30% error rate) while exceeding the error performance of SVMs and the best neural network classifier. Importantly, the RVM classifiers achieve superior performance to the SVM classifiers while utilizing nearly an order of magnitude fewer parameters. While we do not expect the superior error performance to be typical (on pure classification tasks) we do expect the superior sparseness to be typical. This sparseness property is particularly important when attempting to build systems which are practical to train and test.

Approach	Error Rate	# Parameters
K-Nearest Neighbor	44%	
Gaussian Node Network	44%	
SVM: Polynomial Kernels	49%	
SVM: RBF Kernels	35%	83 SVs
Separable Mixture Models	30%	
RVM: RBF Kernels	30%	13 RVs

Table 1. Performance comparison of SVMs and RVMs to other nonlinear classifiers on static vowel classification data.

A hybrid recognition architecture was also developed that is a parallel of our SVM hybrid. Each phone-level classifier (either an SVM or RVM dichotomous classifier) is trained as a one-vs-all classifier. The classifiers are used to predict the probability of an acoustic segment. For the SVM hybrid, a sigmoid posterior fit is used to map the SVM distance to a probability. The RVM output is naturally probabilistic so no link function is needed.

The HMM system is used to generate alignments at the phone level. Each phone instance is treated as one segment. Since each segment could span a variable duration, we divide the segment into three regions in a set ratio and construct a composite vector from the mean vectors of the three regions. In our experiments empirical evidence showed that a 3-4-3 proportion generally gave optimal performance. The classifiers in our hybrid systems operate on composite vectors. For decoding, the segmentation information is obtained from a baseline HMM system—a cross-word triphone system with 8 Gaussian mixtures per state. Composite vectors are generated for each of the segments and posterior probabilities are hypothesized that are used to find the best word sequence using the Viterbi decoder. The HMM system also outputs a set of N-best hypotheses. The posterior probabilities for each hypothesis are determined and the most likely entry of the N-best list is produced.

The performance of RVMs on the static classification of vowel data gave us good reason to expect the performance on continuous speech would be appreciably better than that of the SVM system in terms of sparsity and on par with the SVM system in terms of accuracy. Our initial tests of this hypothesis have been on a telephone alphadigit task. Recent work on both alphabet and alphadigit systems has taken a focus on resolving the high rates of recognizer confusion for certain word sets. In particular, the E-set (B,C, D, E, G, P, T, V, Z, THREE) and A-set (A, J, K, H, EIGHT). The problems occur mainly because the acoustic differences between the letters of the sets are minimal. For instance, the letters B and D differ primarily in the first 10-20 ms during the consonant portion of the letter.

The OGI Alphadigit Corpus is a telephone database collected from approximately 3000 subjects. Each subject was a volunteer responding to a posting on the USEnet. The subjects were given a list of either 19 or 29 alphanumeric strings to speak. The strings in the lists were each six words long, and each list was “set up to balance phonetic context between all letter and digit pairs.” There were 1102 separate prompting strings which gave a balanced coverage of vocabulary and contexts. The training, cross-validation and test sets consisted of 51544, 13926 and 3329 utterances respectively, each balanced for gender. The data sets have been chosen to make them speaker independent.

The hybrid SVM and RVM systems have been benchmarked on the OGI alphadigit corpus with a vocabulary of 36 words. A total of 29 phone models, one classifier per model, were used to cover the pronunciations. Each classifier was trained using the segmental features derived from 39-dimensional frame-level feature vectors comprised of 12 cepstral coefficients, energy, delta and acceleration coefficients. The full training set has as many as 30k training examples per classifier. However, the training routines employed for the RVM models are unable to utilize such a large set as mentioned earlier. The training set was, thus, reduced to 10,000 training examples per classifier (5,000 in-class and 5,000 out-of class).

The test set was an open-loop speaker independent set with 3329 sentences. The composite vectors are also normalized to the range -1 to 1 to assist in convergence of the SVM classifiers. Both the SVM and RVM hybrid systems use identical RBF kernels with the width parameter set to 0.5. The trade-off parameter for the SVM system was set to 50. The sigmoid posterior estimate for the SVM was constructed using a held-out set of nearly 14000 utterances. The results of the RVM and SVM systems are shown in Table 2. The important columns to notice in terms of performance are the error rate, average number of parameters and testing time. In all three, the RVM system outperforms the SVM system. It achieves a slightly better error rate of 14.8% compared to 15.5%. This error rate is obtained in over an order of magnitude fewer parameters. This naturally translates to well over an order of magnitude better runtime performance. However, the RVM does require significantly longer to train. Fortunately, that added training time is done off-line.

#### D.4 Summary

This work is the first application of sparse Bayesian methods to continuous speech recognition. By using an automatic relevance determination mechanism, we are able to achieve state-of-the-art performance in extremely sparse models. Further, this is accomplished while maintaining a purely probabilistic framework. We also achieve performance better than the popular SVM kernel classifier while using an order of magnitude fewer parameters for both a static classification task and a continuous speech task. However, this runtime efficiency comes at a large up front cost during training. Thus, most of our work at this point is focused on more efficient training schemes so that we can move to larger vocabulary tasks. To this end, we have developed an iterative subset refinement approach which attempts to optimize the global criteria by locally optimizing the model on small subsets of the total training set. The subset models are incrementally used to generate a model of the full training set.

We are continuing our work on learning machines in speech recognition, and are now exploring new nonlinear statistical models under separate funding. This ITR project was our first opportunity to explore such risky and innovative methods.

Approach	Word Error Rate	Avg # Parameters	Training Time	Testing Time
SVM: RBF Kernels	15.5%	994	3 hours	1.5 hours
RVM: RBF Kernels	14.8%	72	5 days	5 minutes

Table 2. Performance comparison of SVMs and RVMs on Alphadigit recognition data. The RVMs yield a large reduction in the parameter count while attaining superior performance.

## E. References

- M. Beckman and J. Pierrehumbert, "Intonational structure in Japanese and English," *Phonology Yearbook*, 3, 255-309, 1986.
- W. Byrne et al., "Pronunciation Modelling Using a Hand-Labelled Corpus for Conversational Speech Recognition," in *Proc. ICASSP*, 1998.
- D. Bikel, *On the Parameter Space of Lexicalized Statistical Parsing Models*, Ph.D. thesis, University of Pennsylvania, 2004.
- E. Charniak and M. Johnson, "Edit detection and parsing for transcribed speech," in *Proc. NAACL* 2001.
- M. Gregory, M. Johnson, and E. Charniak, "Sentence-internal prosody does not help parsing the way punctuation does," in *Proc. HLT-NAACL*, 2004, pp. 81-88.
- D. Hindle, "Deterministic parsing of syntactic non-fluencies," in *Proc. ACL*, 1983, pp. 123-128.
- M. Johnson and E. Charniak, "A {TAG}-based noisy channel model of speech repairs," in *Proc. ACL*, 2004, pp. 33-39.
- M. Johnson, E. Charniak and M. Lease, "An improved model for recognizing disfluencies in conversational speech," in *Proc. NIST Rich Transcription Workshop*, 2004, to appear.
- J. Kahn, M. Ostendorf, and C. Chelba, "Parsing conversational speech using acoustic segmentation," in *Proc. HLT-NAACL*, comp. vol., 2004, pp. 125-128.
- J. Kim, S. Schwarm and M. Ostendorf, "Detecting structural metadata with decision trees and transformation-based learning," in *Proc. HLT-NAACL*, pp. 137-144, May 2004.
- Y. Liu, E. Shriberg, A. Stolcke and M. Harper, "Using machine learning to cope with imbalanced classes in natural speech: Evidence from sentence boundary and disfluency detection," in *Proc. ICSLP*, 2004.
- Y. Liu et al., "Structural metadata research in the EARS program," in *Proc. ICASSP*, to appear, 2005.
- L. Mayfield, M. Gavalda, Y. Seo, B. Suhm, W. Ward, and A. Waibel, "Parsing real input in {JANUS}: a concept-based approach," in *Proc. TMI 95*, 1995.
- E. Noth, A. Batliner, A. Kiebling, R. Kompe, and H. Niemann, "Verbmobil: The Use of Prosody in the Linguistic Components of a Speech Understanding System," *IEEE Trans. on Speech and Audio Processing*, 8(5):519-532, 2000.
- M. Ostendorf, "Linking Speech Recognition and Language Processing Through Prosody," *CC-AI*, vol. 15, no. 3, pp. 279-303, 1998.
- M. Ostendorf, I. Shafran, S. Shattuck-Hufnagel, B. Byrne and L. Carmichael, "A prosodically labeled database of spontaneous speech," *Proc. of the ISCA Workshop on Prosody in Speech Recognition and Understanding*, pp. 119-121, October 2001.

- J. Pitrelli, M. Beckman, and J. Hirschberg, "Evaluation of prosodic transcription labeling reliability in the ToBI framework," In Proc. of the International Conference on Spoken Language Processing, 1, 123-126, 1994.
- P. J. Price, M. Ostendorf, S. Shattuck-Hufnagel and C. Fong, "The Use of Prosody in Syntactic Disambiguation," Journal of the Acoustical Society of America, vol. 90, no. 6, December 1991, pp. 2956-2970.
- I. Shafran, M. Ostendorf, and R. Wright, "Prosody and phonetic variability: lessons learned from acoustic model clustering," Proc. of the ISCA Workshop on Prosody in Speech Recognition and Understanding, pp. 127-131, October 2001.
- E. Shriberg et al., "Prosody-based automatic segmentation of speech into sentences and topics," Speech Communication, 32(1-2), pp. 127-154, 2000.
- D. Talkin, "Pitch Tracking," in Speech Coding and Synthesis, ed. W. B. Kleijn and K. K. Paliwal, Elsevier Science B.V., 1995.
- F. Jelinek, *Statistical Methods for Speech Recognition*, MIT Press, Cambridge, Massachusetts, USA, 1997.
- F. Jelinek, "Up From Trigrams! The Struggle for Improved Language Models," *Proceedings of the 2<sup>nd</sup> European Conference on Speech Communication and Technology (Eurospeech)*, pp. 1037-1040, Genova, Italy, September 1991.
- L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257-285, February 1989.
- L. R. Bahl, P. F. Brown, P. V. de Souza, and R. L. Mercer, "Maximum mutual information estimation of hidden Markov model parameters for speech recognition," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, pp. 49-52, Tokyo, Japan, October 1986.
- P. Woodland and D. Povey, "Very Large Scale MMIE Training for Conversational Telephone Speech Recognition," *Proceedings of the 2000 Speech Transcription Workshop*, University of Maryland, MD, USA, May 2000.
- S. Renals, *Speech and Neural Network Dynamics*, Ph. D. dissertation, University of Edinburgh, UK, 1990.
- A. J. Robinson, *Dynamic Error Propagation Networks*, Ph.D. dissertation, Cambridge University, UK, February 1989.
- J. Tebelskis, *Speech Recognition using Neural Networks*, Ph. D. dissertation, Carnegie Mellon University, Pittsburg, USA, 1995.
- A. Ganapathiraju, *Support Vector Machines for Speech Recognition*, Ph.D. Dissertation, Mississippi State University, Mississippi State, Mississippi, USA, 2002.
- M.D. Richard and R.P. Lippmann, "Neural Network Classifiers Estimate Bayesian a Posteriori Probabilities", *Neural Computation*, vol. 3, no. 4, pp. 461-483, 1991.

- H.A. Bourlard and N. Morgan, *Connectionist Speech Recognition — A Hybrid Approach*, Kluwer Academic Publishers, Boston, USA, 1994.
- G.D. Cook and A.J. Robinson, "The 1997 ABBOT System for the Transcription of Broadcast News," *Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop*, Lansdowne, VA, USA, 1998.
- V.N. Vapnik, *Statistical Learning Theory*, John Wiley, New York, NY, USA, 1998.
- C. Cortes, V. Vapnik. Support Vector Networks, *Machine Learning*, vol. 20, pp. 293-297, 1995.
- C.J.C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, [http:// svm.research.bell-labs.com/SVMdoc.html](http://svm.research.bell-labs.com/SVMdoc.html), AT&T Bell Labs, November 1999.
- B. Schölkopf, C. Burges and A. Smola, *Advances in Kernel Methods: Support Vector Machines*, MIT Press, Cambridge, MA, USA, December 1998.
- O. L. Mangasarian, W. Nick Street and W. H. Wolberg: "Breast cancer diagnosis and prognosis via linear programming", *Operations Research*, 43(4), pp. 570-577, July-August 1995.
- Y. Le Cun, L.D. Jackel, L. Bottou, C. Cortes, J. Denker, H. Drucker, I. Guyon, U.A. Muller, E. Sackinger, P. Simard, and V. Vapnik, "Learning algorithms for classification: A comparison on handwritten digit recognition," in *Neural Networks: The Statistical Mechanics Perspective*, J.H. Kwon and S. Cho, eds., pp. 261--276, World Scientific, 1995.
- T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," *Technical Report 23, LS VIII, University of Dortmund*, Germany, 1997.
- P. Clarkson and P. Moreno, "On the Use of Support Vector Machines for Phonetic Classification," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, Phoenix, Arizona, USA, 1999.,
- A. Ganapathiraju, J. Hamaker and J. Picone, "Support Vector Machines for Speech Recognition," *Proceedings of the International Conf. on Spoken Language Processing*, Sydney, Australia, Nov. 1998.
- A. Ganapathiraju, J. Hamaker and J. Picone, "A Hybrid ASR System Using Support Vector Machines," submitted to the *International Conf. of Spoken Language Processing*, Beijing, China, October, 2000.
- A. Ganapathiraju, J. Hamaker and J. Picone, "Hybrid HMM/SVM Architectures for Speech Recognition," *Proceedings of the Department of Defense Hub 5 Workshop*, College Park, Maryland, USA, May 2000.
- M. Ostendorf, V. Digalakis, and O. Kimball, "From HMM's to Segment Models: A Unified View of Stochastic Modeling for Speech Recognition," *IEEE Transactions on Speech and Audio Processing*, vol. 4, no. 5, pp. 360-378, 1996.
- M. Ostendorf and S. Roukos, "A Stochastic Segment Model for Phoneme-Based Continuous Speech Recognition," *IEEE Trans. on Acoustics Speech and Signal Processing*, vol. 37, pp. 1857- 1867, 1989.
- M. Tipping, "Sparse Bayesian Learning and the Relevance Vector Machine," *Journal of Machine Learning*, vol. 1, pp. 211-244, June 2001.

M. Tipping, "The Relevance Vector Machine," in *Advances in Neural Information Processing Systems*, S. Solla, T. Leen and K.-R. Muller, eds., pp. 652-658, MIT Press, 2000.

D. J. C. MacKay, *Bayesian Methods for Adaptive Models*, Ph. D. thesis, California Institute of Technology, Pasadena, California, USA, 1991.

D. J. C. MacKay, "Probable networks and plausible predictions --- a review of practical Bayesian methods for supervised neural networks," *Network: Computation in Neural Systems*, 6, pp. 469-505, 1995.

E. T. Jaynes, "Bayesian Methods: General Background," *Maximum Entropy and Bayesian Methods in Applied Statistics*, J. H. Justice, ed., pp. 1-25, Cambridge University Press, Cambridge, UK, 1986.

S. F. Gull, "Bayesian Inductive Inference and Maximum Entropy," *Maximum Entropy and Bayesian Methods in Science and Engineering, vol. 1, Foundations*, G. J. Erickson and C. R. Smith, eds., Kluwer Publishing, 1988.

J. Rissanen, "Modeling by Shortest Data Description," *Automatica*, 14, pp. 465-471, 1978.

G. Schwartz, "Estimating the Dimension of a Model," *Annals of Statistics*, vol. 6, no. 2, pp. 461-464, 1978.

J. Hamaker, J. Picone, and A. Ganapathiraju, "A Sparse Modeling Approach to Speech Recognition Based on Relevance Vector Machines," *Proceedings of the International Conference of Spoken Language Processing*, vol. 2, pp. 1001-1004, Denver, Colorado, USA, September 2002.

E. Osuna, R. Freund, and F. Girosi, "Support Vector Machines: Training and Applications," *MIT AI Memo 1602*, March 1997.

A. Faul and M. Tipping, "Analysis of Sparse Bayesian Learning," *Proceedings of the Conference on Neural Information Processing Systems*, preprint, 2001.

J. Hamaker, *Sparse Bayesian Methods for Continuous Speech Recognition*, Ph.D. Dissertation, Mississippi State University, Mississippi State, Mississippi, USA, 2003.

M. E. Tipping and A. Faul, "Fast Marginal Likelihood Maximisation for Sparse Bayesian Models," submitted to *Artificial Intelligence and Statistics '03*, 2003.

D. Deterding, M. Niranjan and A. J. Robinson, "Vowel Recognition (Deterding data)," Available at <http://www.ics.uci.edu/pub/machine-learning-databases/undocumented/connectionist-bench/vowel>, 2000.

P. Loizou and A. Spanias, "High-Performance Alphabet Recognition," *IEEE Transactions on Speech and Audio Processing*, pp. 430-445, 1996.

R. Cole, "Alphadigit Corpus v1.0," <http://www.cse.ogi.edu/CSLU/corpora/alphadigit>, Center for Spoken Language Understanding, Oregon Graduate Institute, USA, 1998.

J. Hamaker and J. Picone, "Iterative Refinement of Relevance Vector Machines for Speech Recognition," submitted to *Eurospeech'03*, Geneva, Switzerland, September 2003.

**08/15/03 — 08/14/04: RESEARCH AND EDUCATIONAL ACTIVITIES**

The overall goal of this ITR was to create a strong synergy between speech recognition (ASR) and natural language processing (NLP). In line with the primary goal of the ITR program, this project created close collaborations between groups who did not previously work together. The PIs collaborated on a number of new initiatives as offshoots of this project, including applications in parsing, information retrieval, and homeland security. In this annual report, we briefly describe some of the significant findings of our research below.

## **F. Laboratory for Linguistic Information Processing, Brown University**

Learning general functional dependencies, i.e. functions between arbitrary input and output spaces, is one of the main goals in supervised machine learning. Recent progress has to a large extent focused on designing flexible and powerful input representations, for instance by using kernel-based methods such as Support Vector Machines. We have addressed the complementary issue of problems involving complex outputs such as multiple dependent output variables and structured output spaces. In the context of this project we have mainly dealt with the problem of label sequence learning, a class of problems where dependencies between labels take the form of nearest neighbor dependencies along a chain or sequence of labels. The latter is a natural generalization of categorization or multiclass-classification that has many applications in the context of natural language processing and information extraction. Special cases include part-of-speech tagging, named entity recognition, and speech-accent prediction. More specifically, we have developed and empirically investigated several extensions of state-of-the-art categorization algorithms such as AdaBoost, Support Vector Machines, and Gaussian Process classification. We have designed and implemented several scalable learning algorithms that combine standard optimization techniques employed in the context of the above mentioned methods with dynamic programming techniques that account for the nearest neighbor dependencies. Experimental evaluations on a wide variety of tasks have shown the competitiveness of these methods compared to existing techniques like Hidden Markov Models and Conditional Random Fields.

## **G. Center for Language and Speech Processing, Johns Hopkins University**

The Structured Language Model (SLM) aims at making a prediction of the next word in a given word string by making a syntactical analysis of the preceding words. However, it faces the data sparseness problem because of the large dimensionality and diversity of the information available in the syntactic parses. A neural network model is better suited to tackle the data sparseness problem and its use has been shown to give significant improvements in perplexity and word error rate over the baseline SLM (Emami et al, 2003).

In this work we have investigated a new method of training the neural net based SLM. Our model makes use of a neural network for that component of the SLM that is responsible for predicting the next word given the previous words and their partial syntactic structure. We have investigated both a mismatched and a matched training scenario. In matched training, the neural network is trained on partial parses similar to those that are likely to be encountered during evaluation. On the other hand in the mismatched scenario, faster training time is achieved but at the cost of mismatch between training and evaluation and hence, possible degradation in performance.

The Structured Language Model works by assigning a probability  $P(W,T)$  to every sentence  $W$  and every possible binary parse  $T$  of  $W$ . The joint probability  $P(W,T)$  of a word sequence  $W$  and a complete parse  $T$  is broken into:

$$P(W, T) = \prod_{k=1}^{n+1} P(W_k | W_{k-1} T_{k-1}) \cdot P(t_k | W_{k-1} T_{k-1}, W_k) \cdot \prod_{i=1}^{N_k} P(p_i^k | W_{k-1} T_{k-1}, w_k, t_k, p_1^k \cdots p_{i-1}^k)$$

where  $W_{k-1}T_{k-1}$  is the word-parse (k-1)-prefix,  $t_k$  is the tag assigned to  $w_k$  by the TAGGER,  $N_k - 1$  is the number of operations the CONSTRUCTOR executes at sentence position k before passing control to the PREDICTOR, and  $P_i^k$  denotes the i-th CONSTRUCTOR operation carried out at position k in the word string.

Subsequently, the *language model* probability assignment for the word at position k+1 in the input sentence is made using:

$$P_{\text{SLM}}(w_{k+1} | W_k) = \sum_{T_k \in S_k} P(w_{k+1} | W_k T_k) \cdot \rho(W_k T_k)$$

$$\rho(W_k T_k) = P(W_k T_k) / \sum_{T_k \in S_k} P(W_k T_k)$$

which ensures a proper probability normalization over strings  $W^*$  where  $S_k$  is the set of all parses built and retained by the model at the current stage k.

Neural networks are very suitable for modeling conditional discrete distribution with large vocabularies. These models work by first assigning a continuous feature vector with every token in the vocabulary, and then using a standard multi-layered neural net to get the conditional distribution at the output, given the input feature vectors. Training is achieved by searching for parameters  $\Theta$  of the neural network and the values of feature vectors that maximize the penalized log-likelihood of the training corpus:

$$L = \frac{1}{N} \sum_t \log p(y^t | x_1^t, x_2^t, \dots, x_{n-1}^t; \Theta) - R(\Theta)$$

where  $p(y^t | x_1^t, x_2^t, \dots, x_{n-1}^t; \Theta)$  is the probability of word  $y^t$  (network output at time t), N is the training data size and  $R(\Theta)$  is a regularization term, L-2 norm squared of the parameters in our case.

We have used a neural net to model the SCORER component of the SLM. By the SCORER we refer to the model  $P(w_{k+1} | W_k T_k)$ . The neural net SCORER's parameters can be obtained by training it on the events extracted from the gold standard (usually one best) parses obtained from an external source (humans or an automatic parser). However, there would be a mismatch during evaluation since the partial parses during that phase are not provided and have to be hypothesized by the SLM itself. We have called the SCORER trained in this manner the *mismatched* SCORER.

On the other hand, one can train the model on partial parses hypothesized by the baseline SLM, thus maximizing the proper log-likelihood function. We have called this procedure the *matched* training of the SCORER.

## **H. Signal, Speech, and Language Interpretation Lab, University of Washington**

Prosody can be thought of as the "punctuation" in spoken language, particularly the indicators of phrasing and emphasis in speech. Most theories of prosody have a symbolic (phonological) representation for these events, but a relatively flat structure. In English, for example, two levels of prosodic phrases are usually distinguished: intermediate (minor) and intonational (major) phrases [Beckman & Pierrehumbert, 1986]. While there is evidence that both phrase-level emphasis (or, prominence) of words and prosodic phrases (perceived groupings of words) provide information for syntactic disambiguation [Price et al., 1991], the most important of these cues seems to be the prosodic phrases or the boundary events marking them. While prior work has looked at the use of prosody in automatic parsing of isolated sentences, a key component of our work involved sentence detection as well, since our goal is to handle continuous conversational speech. Hence, the focus of our work has been on automatically recognizing sentence boundaries and sentence-internal prosodic phrase structure and investigating methods for integrating that structure in parsing.

To support these efforts, we also worked on analysis of acoustic cues to prosodic structure. The most important (and best understood) acoustic correlates of prosody are continuous-valued, including fundamental frequency (F0), energy, and duration cues. Particularly at phrase boundaries, cues include significant falls or rises in fundamental frequency accompanied by word-final duration lengthening, energy drops, and optionally a silent pause. In addition, however, there is evidence of spectral cues to prosodic events, so some of our work explored these cues, which also have implications for improving speech recognition.

Our approach to integrating prosody in parsing is to use symbolic boundary events that have categorical perceptual differences, including prosodic break labels that build on linguistic notions of intonational phrases and hesitation phenomena but also higher level structure. These events are predicted from a combination of the continuous acoustic features, rather than using the acoustic features directly, because the intermediate representation simplifies training with high-level (sparse) structures. Just as most automatic speech recognition (ASR) systems use a small inventory of phones as an intermediate level between words and acoustic feature to have robust word models (especially for unseen words), the small set of word boundary events are also useful as a mechanism for generalizing the large number of continuous-valued acoustic features to different parse structures. This approach is currently somewhat controversial because of the high cost of hand labeling, and to some extent because of its association with a particular linguistic theory. However, the specific subset of labels used in this work are relatively theory neutral and language independent, and a key contribution of this work is the use of weakly supervised learning to reduce the cost of prosodic labeling.

An alternative approach, as in [Noth et al, 2000], is to assign categorical "prosodic" labels defined in terms of syntactic structure and presence of a pause, without reference to human perception, and automatically learn the association of other prosodic cues with these labels. While this approach has been very successful in parsing speech from human-computer dialogs, we expect that it will be problematic for conversational speech because of the longer utterance and potential confusion between fluent and disfluent pauses.

### **H.1 Data, Annotation and Development**

The usefulness of speech databases for statistical analysis is substantially increased when the database is labeled for a range of linguistic phenomena, providing the opportunity to improve our understanding of the factors governing systematic suprasegmental and segmental variation in word forms. Prosodic labels and phonetic alignments are available for some read speech corpora, but only a few limited samples of spontaneous conversational speech have been prosodically labeled. To fill this gap, a subset of the

Switchboard corpus of spontaneous telephone-quality dialogs was labeled using a simplification of the ToBI system for transcribing pitch accents, boundary tones and prosodic constituent structure [Pitrelli et al., 1994]. The aim was to cover phenomena that would be most useful for automatic transcription and linguistic analysis of conversational speech. This effort was initiated under another research grant, but completed with partial support from this NSF ITR grant.

The corpus is based on about 4.7 hours of hand-labeled conversational speech (excluding silence and noise) from 63 conversations of the Switchboard corpus and 1 conversation from the CallHome corpus. The orthographic transcriptions are time-aligned to the waveform using segmentations available from Mississippi State University (<http://www.cavs.msstate.edu/hse/ies/projects/switchboard/index.html>) and a large vocabulary speech recognition system developed at the JHU Center for Language and Speech Processing for the 1997 Language Engineering Summer Workshop ([www.clsp.jhu.edu/ws97](http://www.clsp.jhu.edu/ws97)) [Byrne et al., 1997]. All conversations were analyzed using a high quality pitch tracker [Talkin, 1995] to obtain F0 contours, then post-processed to eliminate errors due to crosstalk. The prosody transcription system included: i) breaks (0-4), which indicate the depth of the boundary after each word; ii) phrase prominence (none, \*, \*?); and iii) tones, which indicate syllable prominence and tonal boundary markers. It also provides a way to deal with the disfluencies that are common in spontaneous conversational speech (a p diacritic associated with the prosodic break), and to indicate labeler uncertainty about a particular transcription. The annotation does not include accent tone type, primarily to reduce transcription costs and because we hypothesized that the phrase tones would be relevant to dialog act representations which may be relevant for question-answering. For further information on the corpus and an initial distributional analysis, see [Ostendorf et al. 2001].

The prosodically labeled subset of Switchboard overlaps with the subset of that corpus annotated with Treebank parses, but there is a mismatch in the orthographic transcriptions because the Treebank parses were based on an earlier version of transcripts and the prosodic annotation was based on the higher quality corrections done by Prof. Picone's group (ISIP) at Mississippi State. In addition, we made use of the DARPA EARS metadata annotations that overlapped with the Treebank parses, which were again based on the higher quality transcriptions. To be able to use all of these resources, we used an alignment of words provided by the ISIP team, and mapped the Treebank parse information to the more recent word transcriptions, which could then be aligned with the EARS metadata annotations. Differences in transcriptions were handled by: dropping the parse information for deletions, transferring it as is for word substitutions, and treating it as "missing" information for insertions in the corrected transcripts. While most of the differences between the Treebank and corrected word transcriptions involved simple substitutions (or deletions) that had little or no impact on the parse (e.g. "a" vs. "the"), there were some cases where the transfer introduced noise into the collection of parses. The most frequent such cases were in disfluent regions, where transcribers tend to have more difficulties, including missed word fragments or repetitions ("I I" vs. "I I I"). An additional difference between the Treebank parses and the EARS metadata annotations is the marking of sentence boundaries. Since speakers frequently begin sentences with conjunctions, the metadata conventions often split up constituents marked as compound sentences in Treebank. Because the metadata labelers listened to the speech and the Treebank labelers did not, we chose to use the metadata constituents, which in most cases involved simply dropping a top-level (S) node, but in some cases involved adding a top-level node called "SUGROUP".

## **H.2 Automatic Labeling of Prosodic Structure**

An important part of the effort was development of an automatic prosodic labeling system that would provide cues to improve parsing. In addition, the resulting system was inspected to analyze possible dependencies between prosodic and parse structures in conversational speech. In the experiments, we used decision tree classifiers with different combinations of acoustic, punctuation, parse, and disfluency cues. While more sophisticated techniques, such as HMMs and maximum entropy models, have been

used for related tasks of sentence boundary detection (see [Liu et al., 2005] for a brief survey), we chose decision trees because they are easy to inspect for learning about the prosody-syntax relationship and because this simplified the weakly supervised learning experiments, which were the focus of our efforts.

For the prosody/syntax analyses, we designed trees to predict prosodic labels from syntactic structure, as well as trees to predict prosodic structure from a combination of syntactic and acoustic cues. For purposes of providing information to a parser, we designed trees to predict prosodic constituents from acoustic cues and part-of-speech (POS) tags, but as an intermediate step in designing these trees we also used syntactic cues in designing trees as part of the weakly supervised training. More specifically, a small set of labeled data was used to train prosody models based on both text and acoustic cues, which were then used in combination to automatically label a large set of data that had not been hand-annotated with prosodic structure, and finally new (separate) acoustic-based prosody models were designed from this larger data set for use in parsing new data.

Experiments were conducted on the Switchboard corpus, using the prosodically annotated subset described above for initial training and evaluation (independent subsets for each). Then the full Switchboard training set was incorporated using various methods for weakly supervised learning, as described below. The prosodic constituent labels were merged into 3 classes: major intonational phrase boundary (4), hesitation boundary (1p, 2p), and all other fluent word boundaries. We grouped minor intonational phrase boundaries (3) with the default word boundary class, because preliminary experiments showed that they were almost never predicted by the decision trees (even with sampled training to account for the low frequency) and because they were most often confused with the default word class in 4-class prediction experiments. The simple 3-class system also has the advantage that it is relatively theory neutral and language independent in that essentially all languages have a notion of fluent and disfluent segmentation.

The acoustic cues included normalized F0, energy and duration cues based on those used in [Kim et al., 2004] and similar to those used in other metadata detection studies [Shriberg et al., 2000]. Text-based cues -- including punctuation, parse structure and disfluency markers -- were taken from the annotations available from the Linguistic Data Consortium, aligned to the word transcriptions as described above. The parse features included right-to-left and left-to-right relative depth, part-of-speech, label of the left and right syntactic constituents, and parenthetical markers. In addition, disfluency interruption points and flags for filled pauses and sentence-initial conjunctions were used as features. Punctuation as inserted by a human transcriber (including incomplete sentences) and estimated speaker turn boundaries (defined simply as a word boundary with a silence of length greater than 4s) were also used.

The baseline system trained on hand-labeled data has error rates of 9.6% when all available cues are used (both syntax and prosody) and 16.7% when just acoustic and POS cues are used. This can be compared to an error rate of 30% when the default class is assigned to all word boundaries. We considered three different weakly supervised training techniques for adding data without prosodic labels (but with hand-labeled syntactic structure) into the training set: EM, co-training, and self-training. The co-training algorithm used classifiers designed on either acoustic or syntactic cues, and it differed slightly from the standard method in that we used an information-theoretic distance on the tree posteriors to determine when to omit samples with conflicting classifier decisions. The self-training algorithm used bagging with uniform class sampling to deal with data skew [Liu et al., 2004]. In all cases, only 1-2 iterations were needed. Both the co-training and EM approaches gave improved performance over the baseline, with the

EM algorithm giving the best results of 14.2% error for the acoustic-only trees, which corresponds to a 15% reduction in error rate over supervised training. The self-training strategy actually hurt performance.<sup>4</sup>

From analysis of the resulting trees, we find that punctuation and repair points are correlated with prosodic breaks and hesitations, as expected. Silence duration is the most useful individual acoustic feature, but alone it is not a reliable predictor of prosodic phrases, since there are many prosodic phrases that do not occur at silences and because silences are frequently associated with hesitations. Aside from syntactic structure, the most important text features for predicting prosodic constituents are punctuation, disfluency edit point markers, filler words (sentence-initial coordinating conjunctions, discourse markers, filled pauses), and turn boundaries. Some important syntactic features include depth of subtrees on the left and right sides of the boundary, previous and next syntactic constituent tag, length of closing phrase, and part-of-speech tags. These features were relevant when associated with the target word boundary, but frequently also with the next or previous word boundary. Surprisingly, the label of the joining constituent is not useful. This analysis provided input into the parse reranking work described in the next section.

Due to the success of the weakly supervised training on prosodic phrase boundary detection, we have recently started investigating use of the same technique for training models of prosodic prominence. Initial results show only a 4% reduction in error rate for the system based on acoustic cues, from 22% to 21% error. Despite the high error rate, however, the automatically annotated prominence appears to be useful in topic identifications in preliminary experiments associated with a separate NSF project (IIS-0121396).

### **H.3 Prosody and Parsing**

Parsing speech can be useful for a number of tasks, including information extraction and question answering from audio transcripts. However, parsing conversational speech presents a different set of challenges than parsing text: sentence boundaries are not well-defined, punctuation is absent, and presence of disfluencies (edits and restarts) impact the structure of language. Most prior work on parsing conversational speech has focused on handling disfluencies [Hindle 1983; Mayfield 1995; Charniak & Johnson, 2001], but experiments relied on hand-marked sentence boundaries and made use of punctuation as in text-based parsers. While utterance-level segmentation may be reasonable to assume in current human-computer dialog systems, it is not realistically available in recognized conversational speech. Hence, our work looked at the problem of parsing text with disfluencies and without punctuation.

We have investigated three main issues in the use of prosody in parsing: the impact of automatic sentence segmentation, the usefulness of interruption points, and the usefulness of automatically detected sub-sentence prosodic constituent boundaries (described above). In all cases, we use a two-stage architecture where metadata (constituent boundaries) are first detected with a combination of prosodic and simple text features, and then these symbolic events (or their posterior probabilities) are used in parsing. Our approach focuses on categorical boundary events, which are predicted from a combination of acoustic features, rather than using the acoustic features directly. As argued earlier, the intermediate representation simplifies training with sparse structures. Key research issues include whether the metadata should be treated as "words" or as features on words, whether edits should be represented with an independent component, and how to represent uncertainty of the metadata classifiers. Our work has begun investigating all of these questions, but some remain unanswered and are being pursued in ongoing work.

---

<sup>4</sup> The results reported here are in some cases worse than those reported in an earlier progress report, because they are based on a larger data set. Due to a data processing bug, several files were omitted from earlier studies. In addition, because of the larger amount of data used and the richer feature sets, the trees are much larger than those described in prior reports.

The data used in this work is the Treebank portion of the Switchboard corpus of conversational telephone speech, which includes sentence-like unit boundaries (SUs) as well as the reparandum and interruption point of disfluencies. The data consists of 816 hand-transcribed conversation sides (566K words), of which we reserve 128 conversation sides (61K words) for evaluation testing according to the 1993 NIST evaluation choices. In all cases, training was based on hand-labeled SUs. Parses were evaluated using SU boundaries rather than the standard punctuation-based units that the Treebank is based on, so the gold standard parses and parse evaluation metric were modified to incorporate the SUs.

The most exhaustive series of experiments looked at the impact of automatic segmentation on parsing. For this particular effort, we chose to work with the complete word sequence, i.e. including all of the words within edit regions, to allow experimentation with multiple parsers. In initial work [Kahn et al., 2004], we used the structured language model (SLM) as a parser with a simple pause-based segmentation and automatically detected SUs (69% vs. 35% slot error rate, respectively), showing a significant improvement in parsing performance when using the automatic SUs. We then confirmed the findings with results on the same segmentation alternatives but using two state-of-the-art statistical parsers trained on conversational speech: the Bikel [Bikel, 2004]<sup>5</sup> and Charniak [Charniak & Johnson, 2001]<sup>6</sup> parsers. Figure 1 shows the relative performance of each parser for the two different segmentations, comparing to the oracle performance for each using hand-annotated SUs. In order to have a single number for performance, we use the F-measure calculated from bracket precision and recall. (Trends with separate precision and recall measures and with bracket crossing statistics follow the same patterns.) For all three parsers, there is a significant gain in performance due to using the explicit SU detection algorithm, with more than half of the performance loss associated with the pause-based segmenter recovered when moving to the more sophisticated SU detection system. As SU detection improves, we would expect further performance gains. Also important is the observation that the impact of increased SU error is not lessened by using a better parser: the best oracle results were achieved with the Bikel parser, but it was much more sensitive to SU errors than the SLM.

In trying to understand the role of IPs in parsing, we ran several experiments, including some with and without human-annotated punctuation. Without punctuation, there was a small increase in parsing performance of the SLM using IPs. When the SUs are automatically detected. We were not able to confirm these gains with other parsers; however, recent work in [Johnson, Charniak & Lease, 2004] shows a benefit to edit detection from using IPs which presumably would lead to improved parsing in their two-stage processing strategy [Johnson & Charniak, 2004]. Including punctuation and IPs in experiments with the SLM showed an

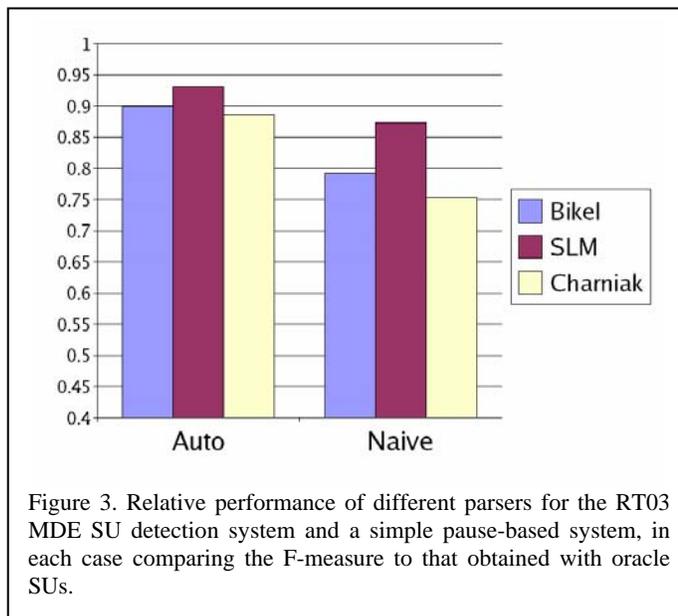


Figure 3. Relative performance of different parsers for the RT03 MDE SU detection system and a simple pause-based system, in each case comparing the F-measure to that obtained with oracle SUs.

<sup>5</sup> <http://www.cis.upenn.edu/~dbikel/download.html> (Version 0.9.9). For this work, we trained the Bikel parser on the Switchboard Treebank parses with the Collins settings.

<sup>6</sup> <ftp://ftp.cs.brown.edu/pub/nlparser/> (August 2004)

interesting trade-off of bracket precision and recall using the two different cues, as illustrated in Figure 2. We saw the improved precision associated with using both punctuation and IPs as possible evidence that sub-sentence prosodic constituents might be useful.

In all of the above work, metadata is incorporated as "word" tokens, similar to the standard mechanism for parsers to incorporate punctuation. For sentence segmentation with a reasonably reliable segmenter, this may make sense, but certainly for sub-sentence prosodic constituents there is the potential for the gains associated with adding prosody to be offset by a loss from the extra words blocking part of the history that might be used in a statistical model of word dependence. We conjecture that this may in part explain the negative results obtained in [Gregory et al., 2004], since our analyses of the prosody prediction trees provides some evidence that sub-sentence prosodic constituents may be useful in parsing. (The direct use of acoustic features may also be problematic.) In addition, the use of metadata events as "words" requires a hard decision in the first stage of detection, and many results in speech processing suggest that soft decisions (e.g. using class posteriors) are more effective.

To address these problems, we developed an extension to the SLM that uses prosodic constituents as hidden conditioning variables, similar to headword conditioning in the SLM. However, since our subsequent work obtained much better baseline performance with other parsers, we decided to explore a parse reranking framework [Johnson et al., ms. in prep.] as an alternative method for incorporating automatically detected prosodic constituents. The approach uses a maximum entropy reranking model and introduces new features based on counts of syntactic constituent types weighted by the posterior probability of different prosodic events. Experiments with this new approach are in progress, now under other funding, and we anticipate having results in early 2005. This series of experiments will also look at the question of whether a separate stage of edit detection benefits parsing compared to simply incorporating the edit structure in the parser with the same status as other constituents.

## I. Institute for Signal and Information Processing, Mississippi State University

Hidden Markov models (HMMs) with Gaussian emission densities are the prominent modeling technique in speech recognition. HMMs suffer from an inability to learn discriminative information and are prone to overfitting and overparameterization. Recent work in machine learning has focused on models, such as the support vector machine (SVM), that automatically control generalization and parameterization as part of the overall optimization process. SVMs, however, require ad hoc (and unreliable) methods to couple it to probabilistic speech recognition systems. We have applied a probabilistic Bayesian learning machine termed the relevance vector machine (RVM) as the core statistical modeling unit in a speech recognizer. The RVM is shown to provide superior performance compared to HMMs and SVMs in terms of both accuracy and sparsity on a continuous alphadigit task.

Computer speech recognition is typically posed as a statistical pattern recognition problem based on Bayesian methods in which the acoustic model and language model are treated as separate statistical models. The focus of our work has been the acoustic model, which maps sequences of features vectors to probabilities that these vectors were produced by a given linguistic unit, such as phone. In most state of

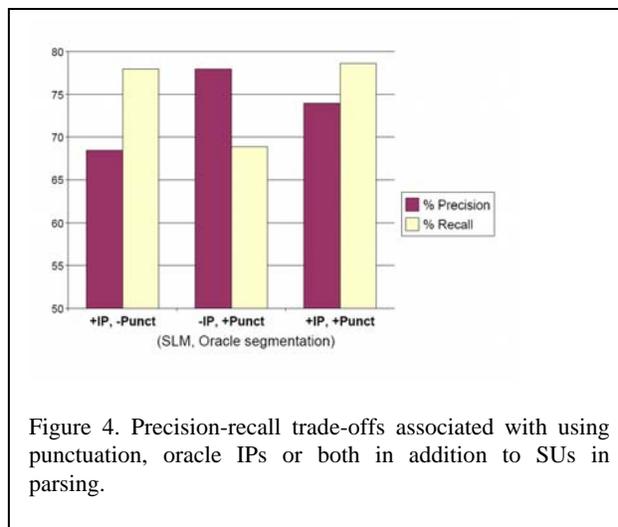


Figure 4. Precision-recall trade-offs associated with using punctuation, oracle IPs or both in addition to SUs in parsing.

the art recognition systems, a hidden Markov model (HMM) is used as the acoustic model. The popularity of the HMM representation is based on an HMMs ability to simultaneously model the temporal progression of speech (speech is usually seen as a “left-to-right” process) and the acoustic variability of the speech observations. The temporal variation is modeled via an underlying Markov process while the acoustic variability is modeled by an emission distribution at each state in the Markov chain.

The most commonly used emission distribution is the Gaussian mixture model (GMM). While combinations of HMMs and Gaussian mixture models have been extremely successful, there are two fundamental limitations of these approaches: (1) the parametric form of the underlying distribution is assumed to be Gaussian, (2) the maximum likelihood (ML) approach, which is typically used to estimate model parameters, does not explicitly improve the discriminative ability of the model. The ML approach maximizes the probability of the correct model while implicitly ignoring the probability of the incorrect model. Ideally, a training approach should force the model toward in-class training examples while simultaneously driving the model away from out-of-class training examples. Methods such as maximum mutual information and minimum classification error have been developed to incorporate discriminative training directly into the standard HMM/GMM framework. However, their success has been limited due primarily to their considerable computational costs.

The weaknesses of the HMM/GMM system have led researchers to explore other models, such as hybrid connectionist systems, which merge the power of artificial neural networks (ANNs) and HMMs. HMM/ANN hybrids have shown promise in terms of performance but have not yet found widespread use due to some serious problems. First, ANNs are prone to overfitting the training data if allowed to blindly converge. To avoid overfitting, a cross-validation set is often used to define a stopping point for training. This is wasteful of data and resources — a serious consideration in speech where the amount of labeled training data is limited. ANNs also typically converge much slower than HMMs. Most importantly, the HMM/ANN hybrid systems have not shown the substantial improvements in recognition accuracy over HMM/GMM systems that would be necessary to force a paradigm shift.

The approaches developed in this work draw significantly from the HMM/ANN hybrid systems. However, we seek methods which are automatically immune to overfitting without the artificial imposition of a cross-validation set as well as methods which can automatically learn the appropriate model structure as part of the overall optimization process. One such model that has come to the forefront of pattern recognition research is the support vector machine (SVM). The support vector paradigm is based upon structural risk minimization (SRM) in which the learning process is posed as one of optimizing some risk function. The optimal learning machine is the one whose free parameters are set such that the risk is minimized.

Finding a minimum of the risk function is typically impossible due to the unknown distribution. Instead, it has been shown that a relationship exists between the actual risk, which is related to the empirical risk (i.e. the training set error which can be measured) and the Vapnik-Chervonenkis (VC) dimension. The VC dimension is a measure of the capacity of a learning machine to learn any training set and is typically closely related to the complexity of the learning machine’s structure. With this result, we can guarantee both a small empirical risk (training error) and good generalization — an ideal situation for a learning machine.

In their most basic form SVMs use the SRM principle to impose an order on the optimization process by ranking candidate separating hyperplanes based on the margin they induce. For separable data, the optimal linear hyperplane is the one that maximizes the margin. The true power of the SVM, however, lies in how it deals with nonlinear class separating surfaces. Providing for a nonlinear decision region is accomplished using kernels. The optimization process yields a decision function where the sign of can be used to classify examples as either in-class or out-of-class. The decision function is formed from only

those training vectors that lie on the margin or in overlap regions. In practice, the proportion of the training set that becomes support vectors is small, making the classifier sparse. Consequently, the training process, along with the training set, directly optimize the complexity of the learning machine. In contrast, ANN systems often make *a priori* assumptions about the form of the model.

SVMs have had great success on static classification tasks. However, it is only recently, that these techniques have been applied to continuous speech recognition. While the SVMs provide an excellent classification paradigm, they suffer from two serious drawbacks that hamper their effectiveness in speech recognition. First, while sparse, the size of the SVM models (number of non-zero weights) tends to scale linearly with the quantity of training data. For a large speaker independent corpus this effect is prohibitive. Second, the SVMs are binary classifiers. In speech recognition this is an important disadvantage since there is significant overlap in the feature space which can not be modeled by a yes/no decision boundary. Further, the combination of disparate knowledge sources (such as linguistic models, pronunciation models, acoustic models, etc.) requires a method for combining the scores produced by each model so that alternate hypotheses can be compared. Thus, we require a probabilistic classification which reflects the amount of uncertainty in our predictions.

We have investigated a Bayesian model termed the relevance vector machine (RVM) which is similar in form to the SVM but which addresses these two problems. The essence of an RVM is a fully probabilistic model with an automatic relevance determination prior over each model parameter. Thus, sparseness in the RVM model is explicitly sought in a probabilistic framework.

## 1.1 Sparse Bayesian Methods

Supervised learning in speech recognition implemented via a maximum likelihood approach is the dominant approach for finding values of the parameters in our model that best match the training data. Our expectation in data modeling is that given sufficient training data, the model would generalize to unseen test sets. Two levels of inference must be implemented to accomplish this. First, assuming that a particular model is true, we seek to infer the values for the parameters of the model that best fit the data at hand. This is exactly the training process used in the ANN hybrids. Second, we must decide which model is most appropriate given the data at hand, i.e. model comparison.

A simple approach to model comparison might dictate that we simply choose the model that fits the data best — the maximum likelihood solution. However, a more complex model can always fit the data better. Jaynes describes an extreme interpretation of this problem where we would always choose the so-called *Sure Thing* hypothesis, under which exactly the training set and only the training set is possible. Though it is the maximum likelihood solution, the *Sure Thing* hypothesis is intuitively displeasing and is counter to our desire for a solution which generalizes. We avoid choosing the *Sure Thing* hypothesis by expressing an *a priori* preference for simpler solutions using the principle of Occam's Razor. MacKay and others have formalized this preference mechanism through the use of Bayesian methods. These provide a natural and quantitative embodiment of Occam's razor. The first level of inference requires that we find the best-fit parameters. The second level of inference requires the comparison of competing hypotheses. If we assume that the competing hypotheses are *a priori* equiprobable then the best hypothesis is chosen by evaluating the evidence. The evidence is computed by marginalization across the model parameters.

The evidence provides a natural trade-off between the best-fit likelihood and the Occam factor. This concept is closely related to other methods such as the Minimum Description Length and the Bayesian Information Criteria where the model is directly penalized by the number of parameters used. A similar idea was also incorporated into SVM models, which penalize the models with too large a capacity (VC dimension). However, while the SVM models are forced to estimate the penalty via cross-validation

schemes, Bayesian techniques automatically determine and apply the penalty in a fully probabilistic framework.

Assuming we have no prior knowledge that would cause us to favor a particular prior, we can find the optimal value for by evaluating the evidence. If we did have prior knowledge, we would simply repeat the inference over using the prior. At some level of the inference, we will arrive at a point where our prior knowledge is too weak to apply and then we evaluate the evidence. This extreme application of the Bayesian prior as a control parameter is known as the method of *automatic relevance determination* (ARD). It is so named because the prior over the input unit weights in a neural network can ‘shut-off’ those input dimensions which are irrelevant to the problem at hand. ARD is at the heart of the relevance vector machines that will be described next.

Approach	Error Rate	# Parameters
K-Nearest Neighbor	44%	
Gaussian Node Network	44%	
SVM: Polynomial Kernels	49%	
SVM: RBF Kernels	35%	83 SVs
Separable Mixture Models	30%	
RVM: RBF Kernels	30%	13 RVs

Table 3. Performance comparison of SVMs and RVMs to other nonlinear classifiers on static vowel classification data.

## I.2 Relevance Vector Machines

An application of the evidence framework to kernel machines is the relevance vector machine (RVM). As with SVMs, RVMs use a weighted linear combination of basis functions. Due to the large number of parameters in this mode  $l$ — one per observation — we must guard against overfitting of the model to the training data. SVMs use a control parameter to implicitly balance the trade-off between training error and generalization. RVMs take a Bayesian approach and explicitly define an ARD prior distribution over the weights.

Each weight in the RVM model has an individual hyperparameter,  $\lambda$ , that is iteratively reestimated as part of the optimization process. As the hyperparameter grows larger, the prior on  $w$  becomes infinitely peaked around zero, forcing  $w$  to go to zero and, thus, contributing nothing to the summation. This process automatically embodies the principle of Occam’s Razor because it explicitly seeks the simplest model that satisfies the data constraints. In practice, the majority of the weights are pruned, resulting in an exceedingly sparse model with generalization abilities on par with SVMs. To complete the Bayesian specification of the model, we have to specify a prior probability. In practice we use a non-informative (flat) prior to indicate a lack of preference.

The parameter estimation process is an iterative reduction process. That is, initially each vector of the system is allocated one parameter. As the procedure continues, vectors are pruned from the model when they are found to be irrelevant with respect to the remaining parameters. Integral to this iterative reestimation process is the computation of the inverse Hessian matrix. This operation requires the inversion of an  $M \times M$  Hessian matrix where  $M$  is initially set to the size of training set. For larger training sets (on the order of a few thousand), this computation is prohibitive both in time and memory.

## I.3 Experiments

RVMs have had significant success in several classification tasks. These tasks have, however, involved relatively small quantities of static data. Speech recognition, on the other hand, involves processing a very large amount of temporally evolving signals. In order to gain insight into the effectiveness of RVMs for speech recognition, we explored two tasks. We first experimented on the Deterding static vowel classification task which is a common benchmark used for new classifiers. Second, we applied the

techniques described above to a complete small vocabulary recognition task. Comparison with SVM models are given below. For each task, the RVMs outperformed the SVM models both in terms of model sparsity and error rate.

In our first pilot experiment, we applied SVMs and RVMs to a publicly available vowel classification task, Deterding Vowels. This was a good data set to evaluate the efficacy of static classifiers on speech classification data since it has been used as a standard benchmark for several nonlinear classifiers for several years. In this evaluation, the speech data was collected at a 10 kHz sampling rate and low pass filtered at 4.7 kHz. The signal was then transformed to 10 log-area parameters, giving a 10 dimensional input space. A window duration of 50 msec was used for generating the features. The training set consisted of 528 frames from eight speakers and the test set consisted of 462 frames from a different set of seven speakers. The speech data consisted of 11 vowels uttered by each speaker in a h\*d context. Though it appears to be a simple task, the small training set and significant confusion in the vowel data make it a very challenging task.

Table 1 shows the results for a range of nonlinear classification schemes on the Deterding vowel data. From the table, the SVM and RVM are both superior to nearly all other techniques. The RVM achieves performance rivaling the best performance reported on this data (30% error rate) while exceeding the error performance of SVMs and the best neural network classifier. Importantly, the RVM classifiers achieve superior performance to the SVM classifiers while utilizing nearly an order of magnitude fewer parameters. While we do not expect the superior error performance to be typical (on pure classification tasks) we do expect the superior sparseness to be typical. This sparseness property is particularly important when attempting to build systems which are practical to train and test.

## J. References

- M. Beckman and J. Pierrehumbert, "Intonational structure in Japanese and English," *Phonology Yearbook*, 3, 255-309, 1986.
- W. Byrne et al., "Pronunciation Modelling Using a Hand-Labelled Corpus for Conversational Speech Recognition," in *Proc. ICASSP*, 1998.
- D. Bikel, *On the Parameter Space of Lexicalized Statistical Parsing Models*, Ph.D. thesis, University of Pennsylvania, 2004.
- E. Charniak and M. Johnson, "Edit detection and parsing for transcribed speech," in *Proc. NAACL* 2001.
- M. Gregory, M. Johnson, and E. Charniak, "Sentence-internal prosody does not help parsing the way punctuation does," in *Proc. HLT-NAACL*, 2004, pp. 81-88.
- D. Hindle, "Deterministic parsing of syntactic non-fluencies," in *Proc. ACL*, 1983, pp. 123-128.
- M. Johnson and E. Charniak, "A {TAG}-based noisy channel model of speech repairs," in *Proc. ACL*, 2004, pp. 33-39.
- M. Johnson, E. Charniak and M. Lease, "An improved model for recognizing disfluencies in conversational speech," in *Proc. NIST Rich Transcription Workshop*, 2004, to appear.
- J. Kahn, M. Ostendorf, and C. Chelba, "Parsing conversational speech using acoustic segmentation," in *Proc. HLT-NAACL*, comp. vol., 2004, pp. 125-128.
- J. Kim, S. Schwarm and M. Ostendorf, "Detecting structural metadata with decision trees and transformation-based learning," in *Proc. HLT-NAACL*, pp. 137-144, May 2004.
- Y. Liu, E. Shriberg, A. Stolcke and M. Harper, "Using machine learning to cope with imbalanced classes in natural speech: Evidence from sentence boundary and disfluency detection," in *Proc. ICSLP*, 2004.
- Y. Liu et al., "Structural metadata research in the EARS program," in *Proc. ICASSP*, to appear, 2005.
- L. Mayfield, M. Gavalda, Y. Seo, B. Suhm, W. Ward, and A. Waibel, "Parsing real input in {JANUS}: a concept-based approach," in *Proc. TMI 95*, 1995.
- E. Noth, A. Batliner, A. Kiebling, R. Kompe, and H. Niemann, "Verbmobil: The Use of Prosody in the Linguistic Components of a Speech Understanding System," *IEEE Trans. on Speech and Audio Processing*, 8(5):519-532, 2000.
- M. Ostendorf, "Linking Speech Recognition and Language Processing Through Prosody," *CC-AI*, vol. 15, no. 3, pp. 279-303, 1998.
- M. Ostendorf, I. Shafran, S. Shattuck-Hufnagel, B. Byrne and L. Carmichael, "A prosodically labeled database of spontaneous speech," *Proc. of the ISCA Workshop on Prosody in Speech Recognition and Understanding*, pp. 119-121, October 2001.

- J. Pitrelli, M. Beckman, and J. Hirschberg, "Evaluation of prosodic transcription labeling reliability in the ToBI framework," In Proc. of the International Conference on Spoken Language Processing, 1, 123-126, 1994.
- P. J. Price, M. Ostendorf, S. Shattuck-Hufnagel and C. Fong, "The Use of Prosody in Syntactic Disambiguation," Journal of the Acoustical Society of America, vol. 90, no. 6, December 1991, pp. 2956-2970.
- I. Shafran, M. Ostendorf, and R. Wright, "Prosody and phonetic variability: lessons learned from acoustic model clustering," Proc. of the ISCA Workshop on Prosody in Speech Recognition and Understanding, pp. 127-131, October 2001.
- E. Shriberg et al., "Prosody-based automatic segmentation of speech into sentences and topics," Speech Communication, 32(1-2), pp. 127-154, 2000.
- D. Talkin, "Pitch Tracking," in Speech Coding and Synthesis, ed. W. B. Kleijn and K. K. Paliwal, Elsevier Science B.V., 1995.
- F. Jelinek, *Statistical Methods for Speech Recognition*, MIT Press, Cambridge, Massachusetts, USA, 1997.
- F. Jelinek, "Up From Trigrams! The Struggle for Improved Language Models," *Proceedings of the 2<sup>nd</sup> European Conference on Speech Communication and Technology (Eurospeech)*, pp. 1037-1040, Genova, Italy, September 1991.
- L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257-285, February 1989.
- L. R. Bahl, P. F. Brown, P. V. de Souza, and R. L. Mercer, "Maximum mutual information estimation of hidden Markov model parameters for speech recognition," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, pp. 49-52, Tokyo, Japan, October 1986.
- P. Woodland and D. Povey, "Very Large Scale MMIE Training for Conversational Telephone Speech Recognition," *Proceedings of the 2000 Speech Transcription Workshop*, University of Maryland, MD, USA, May 2000.
- S. Renals, *Speech and Neural Network Dynamics*, Ph. D. dissertation, University of Edinburgh, UK, 1990.
- A. J. Robinson, *Dynamic Error Propagation Networks*, Ph.D. dissertation, Cambridge University, UK, February 1989.
- J. Tebelskis, *Speech Recognition using Neural Networks*, Ph. D. dissertation, Carnegie Mellon University, Pittsburg, USA, 1995.
- A. Ganapathiraju, *Support Vector Machines for Speech Recognition*, Ph.D. Dissertation, Mississippi State University, Mississippi State, Mississippi, USA, 2002.
- M.D. Richard and R.P. Lippmann, "Neural Network Classifiers Estimate Bayesian a Posteriori Probabilities", *Neural Computation*, vol. 3, no. 4, pp. 461-483, 1991.

- H.A. Bourlard and N. Morgan, *Connectionist Speech Recognition — A Hybrid Approach*, Kluwer Academic Publishers, Boston, USA, 1994.
- G.D. Cook and A.J. Robinson, "The 1997 ABBOT System for the Transcription of Broadcast News," *Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop*, Lansdowne, VA, USA, 1998.
- V.N. Vapnik, *Statistical Learning Theory*, John Wiley, New York, NY, USA, 1998.
- C. Cortes, V. Vapnik. Support Vector Networks, *Machine Learning*, vol. 20, pp. 293-297, 1995.
- C.J.C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, [http:// svm.research.bell-labs.com/SVMdoc.html](http://svm.research.bell-labs.com/SVMdoc.html), AT&T Bell Labs, November 1999.
- B. Schölkopf, C. Burges and A. Smola, *Advances in Kernel Methods: Support Vector Machines*, MIT Press, Cambridge, MA, USA, December 1998.
- O. L. Mangasarian, W. Nick Street and W. H. Wolberg: "Breast cancer diagnosis and prognosis via linear programming", *Operations Research*, 43(4), pp. 570-577, July-August 1995.
- Y. Le Cun, L.D. Jackel, L. Bottou, C. Cortes, J. Denker, H. Drucker, I. Guyon, U.A. Muller, E. Sackinger, P. Simard, and V. Vapnik, "Learning algorithms for classification: A comparison on handwritten digit recognition," in *Neural Networks: The Statistical Mechanics Perspective*, J.H. Kwon and S. Cho, eds., pp. 261--276, World Scientific, 1995.
- T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," *Technical Report 23, LS VIII, University of Dortmund*, Germany, 1997.
- P. Clarkson and P. Moreno, "On the Use of Support Vector Machines for Phonetic Classification," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, Phoenix, Arizona, USA, 1999.,
- A. Ganapathiraju, J. Hamaker and J. Picone, "Support Vector Machines for Speech Recognition," *Proceedings of the International Conf. on Spoken Language Processing*, Sydney, Australia, Nov. 1998.
- A. Ganapathiraju, J. Hamaker and J. Picone, "A Hybrid ASR System Using Support Vector Machines," submitted to the *International Conf. of Spoken Language Processing*, Beijing, China, October, 2000.
- A. Ganapathiraju, J. Hamaker and J. Picone, "Hybrid HMM/SVM Architectures for Speech Recognition," *Proceedings of the Department of Defense Hub 5 Workshop*, College Park, Maryland, USA, May 2000.
- M. Ostendorf, V. Digalakis, and O. Kimball, "From HMM's to Segment Models: A Unified View of Stochastic Modeling for Speech Recognition," *IEEE Transactions on Speech and Audio Processing*, vol. 4, no. 5, pp. 360-378, 1996.
- M. Ostendorf and S. Roukos, "A Stochastic Segment Model for Phoneme-Based Continuous Speech Recognition," *IEEE Trans. on Acoustics Speech and Signal Processing*, vol. 37, pp. 1857- 1867, 1989.
- M. Tipping, "Sparse Bayesian Learning and the Relevance Vector Machine," *Journal of Machine Learning*, vol. 1, pp. 211-244, June 2001.

- M. Tipping, "The Relevance Vector Machine," in *Advances in Neural Information Processing Systems*, S. Solla, T. Leen and K.-R. Muller, eds., pp. 652-658, MIT Press, 2000.
- D. J. C. MacKay, *Bayesian Methods for Adaptive Models*, Ph. D. thesis, California Institute of Technology, Pasadena, California, USA, 1991.
- D. J. C. MacKay, "Probable networks and plausible predictions --- a review of practical Bayesian methods for supervised neural networks," *Network: Computation in Neural Systems*, 6, pp. 469-505, 1995.
- E. T. Jaynes, "Bayesian Methods: General Background," *Maximum Entropy and Bayesian Methods in Applied Statistics*, J. H. Justice, ed., pp. 1-25, Cambridge University Press, Cambridge, UK, 1986.
- S. F. Gull, "Bayesian Inductive Inference and Maximum Entropy," *Maximum Entropy and Bayesian Methods in Science and Engineering, vol. 1, Foundations*, G. J. Erickson and C. R. Smith, eds., Kluwer Publishing, 1988.
- J. Rissanen, "Modeling by Shortest Data Description," *Automatica*, 14, pp. 465-471, 1978.
- G. Schwartz, "Estimating the Dimension of a Model," *Annals of Statistics*, vol. 6, no. 2, pp. 461-464, 1978.
- J. Hamaker, J. Picone, and A. Ganapathiraju, "A Sparse Modeling Approach to Speech Recognition Based on Relevance Vector Machines," *Proceedings of the International Conference of Spoken Language Processing*, vol. 2, pp. 1001-1004, Denver, Colorado, USA, September 2002.
- E. Osuna, R. Freund, and F. Girosi, "Support Vector Machines: Training and Applications," *MIT AI Memo 1602*, March 1997.
- A. Faul and M. Tipping, "Analysis of Sparse Bayesian Learning," *Proceedings of the Conference on Neural Information Processing Systems*, preprint, 2001.
- J. Hamaker, *Sparse Bayesian Methods for Continuous Speech Recognition*, Ph.D. Dissertation, Mississippi State University, Mississippi State, Mississippi, USA, 2003.
- M. E. Tipping and A. Faul, "Fast Marginal Likelihood Maximisation for Sparse Bayesian Models," submitted to *Artificial Intelligence and Statistics '03*, 2003.
- D. Deterding, M. Niranjan and A. J. Robinson, "Vowel Recognition (Deterding data)," Available at <http://www.ics.uci.edu/pub/machine-learning-databases/undocumented/connectionist-bench/vowel>, 2000.
- P. Loizou and A. Spanias, "High-Performance Alphabet Recognition," *IEEE Transactions on Speech and Audio Processing*, pp. 430-445, 1996.
- R. Cole, "Alphadigit Corpus v1.0," <http://www.cse.ogi.edu/CSLU/corpora/alphadigit>, Center for Spoken Language Understanding, Oregon Graduate Institute, USA, 1998.
- J. Hamaker and J. Picone, "Iterative Refinement of Relevance Vector Machines for Speech Recognition," submitted to *Eurospeech'03*, Geneva, Switzerland, September 2003.