

# Derivative Kernels for Noise Robust ASR

A. Ragni and M. J. F. Gales

*Cambridge University Engineering Department  
Trumpington St., Cambridge CB2 1PZ, U.K.  
{ar527,mjfg}@eng.cam.ac.uk*

**Abstract**—Recently there has been interest in combining generative and discriminative classifiers. In these classifiers features for the discriminative models are derived from the generative kernels. One advantage of using generative kernels is that systematic approaches exist to introduce complex dependencies into the feature-space. Furthermore, as the features are based on generative models standard model-based compensation and adaptation techniques can be applied to make discriminative models robust to noise and speaker conditions. This paper extends previous work in this framework in several directions. First, it introduces derivative kernels based on context-dependent generative models. Second, it describes how derivative kernels can be incorporated in structured discriminative models. Third, it addresses the issues associated with large number of classes and parameters when context-dependent models and high-dimensional feature-spaces of derivative kernels are used. The approach is evaluated on two noise-corrupted tasks: small vocabulary AURORA 2 and medium-to-large vocabulary AURORA 4 task.

## I. INTRODUCTION

Most automatic speech recognition (ASR) systems use generative models, hidden Markov models (HMM), as the acoustic models. Likelihoods from these models are combined with the prior, the language model, using Bayes' rule to yield the sentence posterior. Although successful, it is widely known that the underlying models are not correct. This has lead to interest in discriminative classifiers which directly model sentence posteriors/decision boundaries given a set of features extracted from the observation sequence. There are several options how features can be extracted from observation sequences. This includes event detectors [1], generative kernels [2] and other parametric and non-parametric approaches [3]. Event detectors make use of multiple parallel feature streams which operate at different levels of granularity such as word, multi-phone and phone. This flexibility enables a wide range of short and long-spanning dependencies. However, the current applications of event detectors do not attempt to improve the underlying acoustic models, the recognition results from these models are used to derive indicator features. Additionally, the issues associated with adapting feature streams to noise and speaker conditions are not easy to handle.

Generative kernels derive features from generative models and have several advantages. First, the use of competing log-likelihoods, first and higher order derivatives of log-likelihood offers a systematic approach of adding new acoustic features. In contrast to log-likelihoods the derivatives do not inherit conditional independence assumptions from generative models and enable other short and long-spanning dependencies. Second, the generative kernels can be adapted to noise and speaker

conditions using model-based compensation and adaptation approaches [4]. Third, since generative kernels derive features from generative models the parameters of these models can be re-estimated to extract more discriminative features. Previous work with generative kernels has examined several feature configurations. The use of log-likelihood features extracted from whole-word and context-dependent HMMs was investigated in [5] and [6] respectively. However, the features in these approaches inherited the underlying HMM conditional independence assumptions. Derivative features have been examined in [4] and [7]. However, the generative models used in these approaches were whole-word models and small vocabulary recognition tasks were considered.

This paper extends the previous work with derivative features to handle medium/large vocabulary speech recognition tasks. This requires three fundamental issues to be addressed. The first issue is the large number of context-dependent discriminative classes. The approach based on phonetic decision tree clustering [6] to ensure that sufficient amount of training data exists for robust parameter estimation is adopted. With derivative features parameter tying introduces another issue. When more than one distinct generative model is used to extract features the discriminative parameters become sensitive to the order of components in these models. A simple approach is proposed where discriminative parameters associated with derivatives are tied within the states. The third issue is how to handle high-dimensional derivative features during training and decoding. In this paper general on-the-fly variants of training and decoding with generative kernels are proposed.

## II. COMBINED GENERATIVE AND DISCRIMINATIVE CLASSIFIERS

Generative models are well known for their ability to handle variable length sequences, adaptability to varying noise/speaker conditions, efficient learning and inference algorithms. For discriminative classifiers these issues are not easy to handle. This section provides details on a combined approach which offers the benefits of generative models with the additional power of discriminative classifiers.

Consider the framework shown in Figure 1 where the shaded part corresponds to generative models and the rest to discriminative classifiers. The generative part is a standard model-based HMM compensation/adaptation framework. Given noise and speaker-dependent observation sequence  $\mathbf{O}$  the parameters of the canonical HMMs are compensated to the target conditions using model-based techniques. The

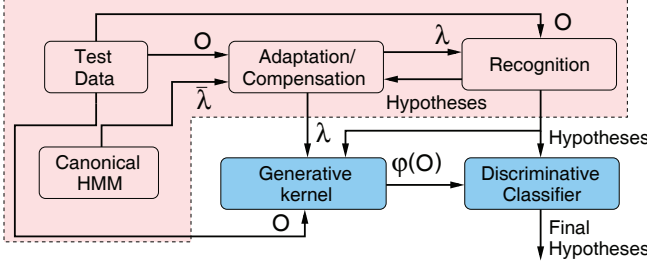


Fig. 1. Combined generative and discriminative framework

discriminative part makes use of these compensated HMMs and observation sequences to extract a set of features. These features handle the mapping from variable length sequences to a fixed dimension and incorporate a range of short and long-spanning dependencies. The advantage of this framework is that the features extracted from the compensated HMMs will be automatically adapted to target noise and speaker conditions. By compensating the features it is then possible to train noise and speaker-independent discriminative classifiers.

In this work vector Taylor series (VTS) model-based compensation is applied to map HMM parameters to target noise conditions. The first-order VTS scheme described in [8] is used. The mismatch function between the static part of clean  $\mathbf{x}_t^s$  and noise-corrupted  $\mathbf{o}_t^s$  observation is given by

$$\mathbf{o}_t^s = \mathbf{x}_t^s + \mathbf{h} + \mathbf{C} \log(1 + \exp(\mathbf{C}^{-1}(\mathbf{n}_t^s - \mathbf{x}_t^s - \mathbf{h}))) \quad (1)$$

where  $\mathbf{C}$  is a discrete cosine transformation matrix,  $\mathbf{1}$  is a unit vector,  $\mathbf{n}_t^s$  and  $\mathbf{h}$  are additive and convolutional noise vectors. Applying the first-order VTS expansion and taking expectation with respect to static parameters of component  $\theta^{jm}$  yields the following form of updated mean and covariance

$$\mu_{jm}^s = \mathbf{C} \log(\exp(\mathbf{C}^{-1}(\bar{\mu}_{jm}^s - \mu_h)) + \exp(\mathbf{C}^{-1}\mu_n^s)) \quad (2)$$

$$\Sigma_{jm}^s = \mathbf{J}_{jm} \bar{\Sigma}_{jm}^s \mathbf{J}_{jm}^T + (\mathbf{I} - \mathbf{J}_{jm}) \Sigma_n^s (\mathbf{I} - \mathbf{J}_{jm})^T \quad (3)$$

where  $\mu_h$  and  $\mu_n$ ,  $\Sigma_n$  are convolutional and additive noise parameters estimated using maximum likelihood (ML) training [9],  $\mathbf{I}$  is identity matrix,  $\mathbf{J}_{jm}$  is a component-specific Jacobian matrix computation of which is fully described in [8].

Several options exist to estimate the canonical HMM parameters. One approach is to train HMMs on clean data. Another approach is to adaptively train HMMs using ML [10], [11] or minimum phone error (MPE) [9], [12] training on multi-style data collected in various noise and speaker conditions. This allows more data to be used in estimating canonical model parameters.

### III. DERIVATIVE KERNELS

Generative kernels in Figure 1 extract features from generative models. Several types of generative kernels have been proposed in literature [13], [2], [7]. As this work primarily deals with the feature-spaces of those kernels no particular form of the kernel such as polynomial, Gaussian etc. is assumed. The simplest example of generative kernels are *log-likelihood kernels*. The base (b) feature vector in equation (4)

is a log-likelihood of generative model computed for class  $\omega_i$  from observation sequence  $\mathbf{O}$

$$\phi_b^0(\mathbf{O}|\omega_i) = [\log(p(\mathbf{O}|\omega_i))] \quad (4)$$

Another example is shown in equation (5) where, in addition to the correct class  $\omega_i$ , log-likelihoods of competing classes are also appended (a)

$$\phi_a^0(\mathbf{O}|\omega_i) = \begin{bmatrix} \log(p(\mathbf{O}|\omega_1)) \\ \log(p(\mathbf{O}|\omega_2)) \\ \vdots \\ \log(p(\mathbf{O}|\omega_K)) \end{bmatrix} \quad (5)$$

The appended feature vector in equation (5) may include log-likelihoods of all classes [5] or a subset of them [6]. The features derived from the log-likelihood kernels inherit conditional independence assumptions of the underlying generative models. In contrast to log-likelihood kernels features derived from *derivative kernels* have different conditional independence assumptions. Consider the  $\rho$ -th order base derivative kernel where the feature vector has the following form

$$\phi_b^\rho(\mathbf{O}|\omega_i) = \begin{bmatrix} \log(p(\mathbf{O}|\omega_i)) \\ \nabla_{\lambda} \log(p(\mathbf{O}|\omega_i)) \\ \vdots \\ \nabla_{\lambda}^\rho \log(p(\mathbf{O}|\omega_i)) \end{bmatrix} \quad (6)$$

In addition to correct class log-likelihood the feature vector in equation (6) incorporates derivatives up to the order  $\rho$  with respect to generative model parameters. Consider the first-order derivatives taken with respect to component  $\theta^{jm}$  output distribution parameters  $\lambda_{jm} = \{\mu_{jm}, \Sigma_{jm}\}$

$$\nabla_{\lambda_{jm}} \log(p(\mathbf{O}|\omega_i)) = \sum_{t=1}^T P(\theta_t^{jm}|\mathbf{O}) \nabla_{\lambda_{jm}} \log(p(\mathbf{o}_t|\theta_t^{jm})) \quad (7)$$

These derivatives are functions of component posterior probabilities,  $P(\theta_t^{jm}|\mathbf{O})$ , which depend on the whole observation sequence. This means that the use of derivatives introduces additional dependencies into the features. Higher-order derivatives offer more complex dependencies [7].

Since not all first and higher order derivatives are equally discriminative a subset of them are normally used. In [14] the derivatives with respect to the mean vectors (1m) were found to be the most discriminative first-order derivatives. The feature vector in this case has the following form

$$\phi_b^{1m}(\mathbf{O}|\omega_i) = \begin{bmatrix} \log(p(\mathbf{O}|\omega_i)) \\ \sum_{t=1}^T P(\theta_t^{1,1}|\mathbf{O}) \Sigma_{1,1}^{-1/2}(\mathbf{o}_t - \mu_{1,1}) \\ \vdots \\ \sum_{t=1}^T P(\theta_t^{1,M}|\mathbf{O}) \Sigma_{1,M}^{-1/2}(\mathbf{o}_t - \mu_{1,M}) \\ \sum_{t=1}^T P(\theta_t^{2,1}|\mathbf{O}) \Sigma_{2,1}^{-1/2}(\mathbf{o}_t - \mu_{2,1}) \\ \vdots \\ \sum_{t=1}^T P(\theta_t^{N,M}|\mathbf{O}) \Sigma_{N,M}^{-1/2}(\mathbf{o}_t - \mu_{N,M}) \end{bmatrix} \quad (8)$$

where  $N$  is the number of states and  $M$  is the number of components in every state. Note that consistently with other work in this area standard deviation rather than variance normalisation is performed [4].

In order to illustrate the advantages of using first and higher derivatives consider the following example [7]. A discrete

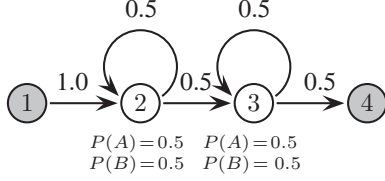


Fig. 2. Example discrete HMM topology, transition and output probabilities

HMM with the topology shown in Figure 2 is used to model two classes  $\omega_1$  and  $\omega_2$ . The data for the two classes are

$\omega_1$  : AAAA, BBBB  
 $\omega_2$  : AABB, BBAA

When ML training is used to estimate HMM parameters then the state transition and output probabilities shown in Figure 2 are obtained. Since all estimated distributions yield equal probabilities the HMM is not capable of distinguishing between the two classes. The situation is different with derivative kernels. Table III shows values of selected derivatives. When the first and second order derivatives are computed with

TABLE I  
FEATURE VECTOR VALUES FOR SECOND-ORDER GENERATIVE KERNEL

Feature	Class $\omega_1$		Class $\omega_2$	
	AAAA	BBBB	AABB	BBAA
$\nabla_{2A}$	0.50	-0.50	0.33	-0.33
$\nabla_{2A} \nabla_{2A}^T$	-3.83	0.17	-3.28	-0.61
$\nabla_{2A} \nabla_{3A}^T$	-0.17	-0.17	-0.06	-0.06

respect to output symbol A in state 2 (line 1 and 2) then all training examples may be correctly classified provided non-linear decision boundaries can be modelled. With the cross-state second order derivative  $\nabla_{2A} \nabla_{3A}^T$  (line 3) a linear decision boundary is sufficient. This second order derivative is capable of capturing whether label changes or not on transition from state 2 to state 3.

#### IV. CLASSIFICATION WITH DERIVATIVE KERNELS

##### A. Isolated word kernels and classification

The derivative kernels from Section III can be directly applied for isolated word classification tasks. One option to extend them to classify sequences rather than isolated words is to use acoustic code-breaking [15]. In this approach recognition of continuous speech is broken down into classification of a sequence of isolated speech segments. Given a word-level hypothesis with alignment an isolated discriminative classifier is sequentially applied to every segment. One classifier used for this task is the support vector machine (SVM). For these

SVMs the feature vector of the first-order derivative kernel has the following form [14]

$$\phi(\tilde{\mathbf{O}}, \omega_i, \omega_j) = \begin{bmatrix} \log(p(\tilde{\mathbf{O}}|\omega_i)) - \log(p(\tilde{\mathbf{O}}|\omega_j)) \\ \nabla_{\lambda} \log(p(\tilde{\mathbf{O}}|\omega_i)) \\ \nabla_{\lambda} \log(p(\tilde{\mathbf{O}}|\omega_j)) \end{bmatrix} \quad (9)$$

where  $\omega_i$  and  $\omega_j$  are two classes and  $\tilde{\mathbf{O}}$  is a segment of observation sequence. Note that the feature vector in equation (9) is a *joint feature vector* which incorporates features simultaneously from two classes. For multi-class classification with SVMs schemes such as majority voting [14], [4] and tree-based reductions [16] have been examined. However, with large number of words the number of binary SVMs required in these approaches becomes large. One option to address this issue is to use a multi-class SVM [5]. However, with high-dimensional derivative features and large number of classes the total dimensionality of the joint feature vector becomes huge. This makes constraint satisfaction of maximum margin training computationally infeasible. Therefore in [5] log-likelihood rather than derivative kernels were used.

The use of acoustic code-breaking approach is suboptimal in several ways. The first issues is that the discriminative model is defined on a word-level which is not useful for medium/large vocabulary tasks. The use of subword models is complicated as phone boundaries are hard to reliably estimate. Another issue with acoustic code-breaking is that isolated segments rather than continuous sequences are modelled.

##### B. Continuous speech kernels and classification

This paper extends derivative kernels to classify continuous speech by using structured discriminative models [7], [5], [6]. The model considered in this work has a log-linear form

$$P(\mathbf{W}|\mathbf{O}) = \frac{\exp(\alpha^T \phi(\mathbf{O}, \mathbf{W}, \theta))}{\sum_{\mathbf{W}'} \exp(\alpha^T \phi(\mathbf{O}, \mathbf{W}', \theta'))} \quad (10)$$

where  $\alpha$  are discriminative parameters and  $\theta$  is an alignment. An important decision to make is at which level  $\theta$  will segment the data as this defines the level of independence assumptions in the model. Segmenting data at the word level is not useful for medium/large vocabulary acoustic models. Similar to the work in [6] here the data is segmented at the phone level. Figure 3 illustrates the structure of the model in equation (10) by a lattice typically used in discriminative HMM training [17]. Here every word arc is segmented into a sequence of

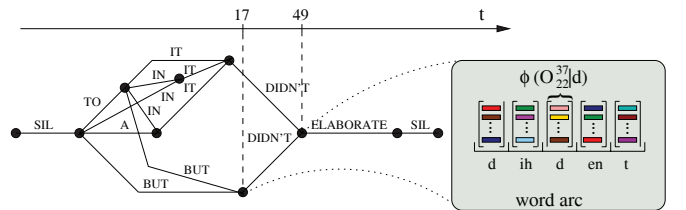


Fig. 3. Structure modelling approach in continuous discriminative models

phone arcs. This allows generative models attached to phone arcs to be directly used in derivative kernels. The features extracted are shown in Figure 3 as the column vectors. Note that for simplicity context-independent labels are shown.

Given observation sequence  $\mathbf{O}$  and hypothesised word sequence  $\mathbf{W}$  aligned by  $\theta$  the model in equation (10) assigns a *score* equal to the exponent of the dot-product below

$$\alpha^\top \phi(\mathbf{O}, \mathbf{W}, \theta) = \sum_{i=1}^{L_p} \alpha^\top \phi(\mathbf{O}_{t(w_i, \theta)}, w_i) + \sum_{j=1}^{L_w} \log(P(\mathbf{w}_j)) \quad (11)$$

The dot product in equation (11) is a summation of phone-level dot products and word-level language model probabilities. Features used at the phone-level are those extracted by derivative kernels from generative models

$$\phi(\mathbf{O}, w) = \begin{bmatrix} \delta(w, \omega_1) \phi_b^p(\mathbf{O}|\omega_1) \\ \vdots \\ \delta(w, \omega_{K_p}) \phi_b^p(\mathbf{O}|\omega_{K_p}) \end{bmatrix} \quad (12)$$

where  $w$  is one of  $K_p$  context-dependent classes. The vector in equation (12) is a high-dimensional *joint feature vector*, the use of delta functions ensures that only one  $\phi_b^p(\mathbf{O}|\omega_i)$  is active for each phone arc. The language model probabilities in equation (11) are obtained in this work from a  $n$ -gram model.

In this work the alignment  $\theta$  is produced by HMMs. However, this introduces suboptimality into the discriminative model since in training/decoding the most likely alignment

$$\hat{\theta}/\{\hat{\theta}, \hat{\mathbf{W}}\} = \arg \max_{\theta/\{\theta, \mathbf{W}\}} \{\alpha^\top \phi(\mathbf{O}, \mathbf{W}, \theta)\} \quad (13)$$

with respect to discriminative parameters  $\alpha$  theoretically must be used. The inference problem in equation (13) can be solved using the semi-Markov equivalent [18] of the Viterbi algorithm. In [19] the impact of using the most likely alignment was investigated with the appended log-likelihood features  $\phi_a^0$  on a digit string recognition task. Small improvements were observed over using alignments produced by HMMs. In this work the use of optimal alignment was not investigated and the alignment provided by HMMs was adopted.

## V. PARAMETER TYING AND ESTIMATION

### A. Parameter tying

When context-dependent generative models are used the number of possible classes becomes large. It is unlikely that the amount of training data available will be sufficient to robustly estimate parameters of all classes. The standard approach with generative models is to use *state-level* phonetic decision trees to cluster phonetically similar states together [20]. Since discriminative classes in this work are defined on a model rather than state level the trees created for generative models can not be re-used. Therefore, another set of *model-level* decision trees is created as described in [6].

When derivative features are used there is an additional issue to consider. The derivatives are computed with respect to the fixed order of components of generative model states. The clustering procedure applied in phonetic decision tree building, however, is insensitive to the order of components. This may

introduce accuracy issues when several generative models are tied without taking into account the order of components in these models. Consider an example on the left in equation (14) where for simplicity the class label is omitted and one-state HMM is assumed.

$$\begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_M \end{bmatrix}^\top \begin{bmatrix} \log(p(\mathbf{O})) \\ \nabla_{\lambda_1} \log(p(\mathbf{O})) \\ \vdots \\ \nabla_{\lambda_M} \log(p(\mathbf{O})) \end{bmatrix} \quad \begin{bmatrix} \alpha_0 \\ \alpha \\ \vdots \\ \alpha \end{bmatrix}^\top \begin{bmatrix} \log(p(\mathbf{O})) \\ \nabla_{\lambda_1} \log(p(\mathbf{O})) \\ \vdots \\ \nabla_{\lambda_M} \log(p(\mathbf{O})) \end{bmatrix} \quad (14)$$

When several generative models  $\lambda^{(1)}, \dots, \lambda^{(K)}$  are used to extract features then the estimate of the discriminative parameter  $\alpha_m$  will depend on the features extracted from the components  $\lambda_m^{(1)}, \dots, \lambda_m^{(K)}$ . It is then sufficient to have one outlier component to distort the expected range of features from the rest  $K - 1$  components. One option to overcome this is to ensure that a small number of generative models is used by any discriminative class. However, this can introduce robustness issues as fewer training examples will be available. The option considered in this work is to tie parameters associated with derivatives within states as shown on the right in equation (14). This is expected to help to deweight the contribution of outliers into the estimate of  $\alpha$  as more features will be used. With limited amount of training data this approach can also improve robustness by reducing the number of parameters by a factor of  $M$ .

### B. Parameter estimation

The standard criterion to train log-linear models is a conditional maximum likelihood (CML). For tasks such as ASR minimum Bayes' risk (MBR) training is a popular alternative approach. The objective function in MBR training is given by

$$\mathcal{F}_{\text{mbr}}^a(\alpha) = \sum_{r=1}^R \sum_{\mathbf{W}} P(\mathbf{W}|\mathbf{O}^{(r)}) \mathcal{A}(\mathbf{W}, \mathbf{W}_{\text{ref}}^{(r)}) \quad (15)$$

where the accuracy function  $\mathcal{A}(\mathbf{W}, \mathbf{W}_{\text{ref}})$  can be defined on a sentence, word, phone or frame level. In MBR training of log-linear models [6] gradient-based optimisation is performed

$$\nabla_{\alpha} \mathcal{F}_{\text{mbr}}^a(\alpha) = \sum_{r=1}^R \sum_{\mathbf{a} \in \mathbf{L}_{\text{den}}^{(r)}} \mathcal{C}(\mathbf{a}) P(\mathbf{a}|\mathbf{O}^{(r)}) \phi(\mathbf{O}_{t(\mathbf{a})}^{(r)}, w) \quad (16)$$

where  $\mathcal{C}(\mathbf{a})$  is phone arc  $\mathbf{a}$  contribution to the average accuracy,  $P(\mathbf{a}|\mathbf{O})$  is arc posterior probability and  $\phi(\mathbf{O}, w)$  is given by equation (12).

Storing high-dimensional features attached to every phone arc as in Figure 3 is impractical for medium/large vocabulary tasks. Therefore, in this paper *on-the-fly training* and *decoding* is performed. In contrast to the standard MBR training in the on-the-fly variant every lattice is passed through twice. The first pass extracts derivative features and accumulates dot products with discriminative parameters on every phone arc. These phone-level dot products multiplied by acoustic deweighting constant  $\gamma$  are then combined with language model probabilities in a lattice-based forward-backward algorithm [17]



to yield arc posterior probabilities and contributions. In the second pass the gradient in equation (16) is accumulated. Although every derivative is computed twice there is no need to keep features attached to phone arcs. The derivative feature vectors can be computed on-the-fly and freed once arc  $a$  is finished. For the on-the-fly decoding only one pass over lattices is required. This is then followed by the best path search as in the standard lattice rescoring case.

In this work regularised training is performed where the final objective function to maximise has the following form

$$\mathcal{F}(\alpha) = \mathcal{F}_{\text{mbr}}^a(\alpha) - \frac{1}{2}(\alpha - \alpha_0)^T \Sigma_\alpha^{-1}(\alpha - \alpha_0) \quad (17)$$

The second term in equation (17) originates from a Gaussian prior. The mean of the prior has the form

$$\alpha_0^{(\omega)} = [1 \quad 0 \quad \dots \quad 0]^T \quad (18)$$

which in equation (10) would yield the generative model performance. The weight matrix  $\Sigma_\alpha$  in this work has a diagonal form where a separate  $\sigma_1$  and  $\sigma_d$  weights are used for parameters associated with log-likelihood and derivatives.

## VI. RESULTS

This section describes experiments with derivative kernels in AURORA 2 and AURORA 4 task. For all systems the structured discriminative models were initialised with the sparse parameter vector in equation (18) to yield generative model performance on the first iteration. The acoustic deweighting constant  $\gamma$  was set to  $\frac{1}{48}$  in AURORA 2 and  $\frac{1}{16}$  in AURORA 4 task respectively. Similarly to other work in this area RProp optimisation was performed [1]. To prevent over-training a subset of test data was used to stop training, Set A for AURORA 2 and Set B for AURORA 4.

### A. AURORA 2

AURORA 2 is a connected digit string recognition task. The number of classes is 11 plus the `sil` and `sp` model. Whole-word HMMs with 16 states and 3 components/mixture trained using ML on the clean data were used as the generative models. The setup used follows the one described in [4]. The structured discriminative model was based on derivative features  $\phi_b^{1m}$  in equation (8), no language model was used. As a contrast a one-dimensional feature  $\phi_b^0$  in equation (4) was also used. The number of discriminative parameters was 21,554 and 13. The multi-style data was used for training.

The word error-rate (WER) averaged over 0-20 dB test data of the VTS-compensated HMMs (VTS) and discriminative classifiers is shown in Table II. The first block quotes

TABLE II  
AURORA2 RECOGNITION RESULTS BASED ON CLEAN-TRAINED HMMs

System	Test set			Avg
	A	B	C	
VTS	9.8	9.1	9.5	<b>9.5</b>
SVM	7.5	7.4	8.1	<b>7.6</b>
$\phi_b^0$	8.1	7.4	8.2	<b>7.8</b>
$\phi_b^{1m}$	7.0	6.6	7.6	<b>7.0</b>

the acoustic code-breaking results with binary SVMs [4] as described in Section IV-A. As can be seen from Table II the use of isolated discriminative classifier with derivative kernels yields large gains over the VTS. The second block in Table II shows the performance of structured discriminative models. The use of one-dimensional features  $\phi_b^0$  gives results comparable to the performance of SVMs but with significantly fewer parameters. The derivative features  $\phi_b^{1m}$  improve the result of  $\phi_b^0$  by 10% relatively, however, the number of additional parameters is approximately half of those in the HMMs. Comparing the performance of the isolated SVMs and continuous derivative kernels it can be seen that modelling whole sentences rather than isolated segments gives consistent gains. The same was observed with the appended log-likelihood kernels in [5].

### B. AURORA 4

AURORA 4 is a noise-corrupted medium/large vocabulary task based on the Wall Street Journal (WSJ) data. Two configurations of canonical HMMs were considered. The first repeats the previous setup where the HMMs were trained from clean data (SI-84 WSJ0 part, ~14 hours). In the second more advanced VTS-adaptive training (VAT) was used to obtain the canonical HMM [9], [11]. For both setups the HMMs were state-clustered triphones (~3140 states) with ~16 components/mixture. Model-based noise compensation was done in two cycles where multiple (4) iterations of VTS compensation were performed for the training and test data, the supervision hypothesis was updated after each cycle. The discriminative model was based on  $\phi_b^0$  and  $\phi_b^{1m}$  features. The parameters of discriminative classes were tied to yield 47 and 4020 classes. Evaluation was performed using the standard 5000-word WSJ0 bigram model on four noise-corrupted test sets based on NIST Nov'92 WSJ0 test set.<sup>1</sup> The language model scale was set to 16.

The first configuration investigated the usefulness of derivative kernels based on context-dependent HMMs with different number of discriminative classes. Table III shows AURORA 4 recognition results. The first block gives baseline performance

TABLE III  
AURORA 4 RECOGNITION RESULTS BASED ON CLEAN-TRAINED HMMs

Classes	System	State tied $\alpha$	Test set				Avg
			A	B	C	D	
	VTS		7.1	15.3	12.1	23.1	<b>17.9</b>
47	$\phi_b^0$	yes	7.6	14.6	11.8	22.2	<b>17.2</b>
	$\phi_b^{1m}$	yes	7.5	14.1	11.3	21.6	<b>16.6</b>
		no	7.4	14.3	11.7	21.9	<b>16.9</b>
4020	$\phi_b^0$	yes	6.6	14.2	10.7	21.8	<b>16.7</b>
	$\phi_b^{1m}$	yes	6.8	13.7	10.6	21.3	<b>16.2</b>
		no	6.7	13.5	10.2	21.1	<b>16.0</b>

of the VTS-compensated HMMs. The second block gives results for 47-class discriminative model based on  $\phi_b^0$  and

<sup>1</sup>Test set A is clean, set B has 6 types of noise added, set C has the channel distortion introduced and set D has both the additive noise and the channel distortion. Average SNR in noise-corrupted data is 10 dB.

$\phi_b^{1m}$  features. The first line in the second block shows that the use of one-dimensional  $\phi_b^0$  features gives gains over VTS, the number of additional parameters is only 47. The next two lines show that when  $\phi_b^{1m}$  features are used the arbitrary ordering of components in HMMs has a clear impact on discriminative model performance. The third block in Table III shows results for 4020-class discriminative model in the three cases described above. As in the case of 47 classes the use of  $\phi_b^0$  features yields gains over the VTS. The addition of derivatives further improves the results. As discussed in Section V with large number of classes the impact of arbitrary component ordering is expected to be small. The results in Table III confirm this by showing that tying parameters instead has lead to a small drop in classification accuracy. However, since the number of parameters in the tied case is less by a factor of 16 the within-state tying is useful for making compact discriminative models based on derivative kernels. Comparing the second and third block consistent gains can be observed from using more discriminative classes.

The second configuration used a VTS adaptively trained HMM system (VAT). Note in this configuration both the generative and discriminative models were trained on multi-style data. The following table shows the performance of baseline VAT, MPE-trained VAT (MPE-VAT) from [21] and 4020-class discriminative models based on  $\phi_b^0$  and  $\phi_b^{1m}$  features. Comparing the VAT in Table IV (line 1) and the

TABLE IV  
AURORA 4 RECOGNITION RESULTS BASED ON VAT HMMs AND  
COMPARISON TO MPE-VAT HMMs

System	Test set				Avg
	A	B	C	D	
VAT	8.6	13.8	12.0	20.1	<b>16.0</b>
MPE-VAT	7.2	12.8	11.5	19.7	<b>15.3</b>
VAT+ $\phi_b^0$	7.7	13.1	11.0	19.5	<b>15.3</b>
VAT+ $\phi_b^{1m}$	7.4	12.6	10.7	19.0	<b>14.8</b>

VTS in Table III (line 1) gain around 2% absolute can be observed from the adaptive training of generative parameters. The VAT+ $\phi_b^0$  discriminative model gives gains over the VAT. The use of derivative features again improves the performance further. Comparing Tables III and IV shows that the use of derivative features in 4020-class VTS+ $\phi_b^{1m}$  model achieves similar performance to more advanced VAT system. Further, by comparing lines 2 and 3 in Table IV the VAT+ $\phi_b^0$  model gives the same level of performance as the MPE-VAT though the number of discriminatively trained parameters in the former is just 4020. Finally, comparing the performance of VAT+ $\phi_b^{1m}$  with the MPE-VAT system shows that derivative features on average yield 0.5% absolute improvement, the largest gain comes from the most noisy conditions.

## VII. CONCLUSION

This paper has described a structured discriminative model based on derivative kernels which is suitable for noise-robust medium/large vocabulary speech recognition. Here the generative models are adapted to noise and speaker conditions using

model-based techniques. The adapted models are then used to extract features from observation sequences. Previous work on small vocabulary tasks with whole word/phone models has been extended to allow context-dependent models to be used. At the phone level, in addition to log-likelihood, the first-order derivatives with respect to HMM mean vectors are used as the features. Parameter tying and estimation with large number of discriminative classes and high-dimensional features is described. The performance of structured discriminative model was evaluated on two noise-corrupted tasks: AURORA 2 and AURORA 4. Consistent gains have been observed over VTS-compensated clean-trained ML, VTS adaptively trained ML and MPE HMM systems.

## ACKNOWLEDGMENT

Anton Ragni is jointly funded by EPSRC, HTK and Toshiba Research Europe Limited.

## REFERENCES

- [1] G. Zweig and P. Nguyen, "A segmental CRF approach to large vocabulary continuous speech recognition," in *Proc. ASRU*, 2009.
- [2] N. Smith, "Using Augmented Statistical Models and Score Spaces for Classification," Ph.D. dissertation, Cambridge University, 2003.
- [3] G. Zweig *et al.*, "Speech recognition with segmental conditional random fields: a summary of the JHU CLSP 2010 summer workshop," in *Proc. ICASSP*, 2011.
- [4] M. Gales and F. Flego, "Discriminative classifiers with adaptive kernels for noise robust speech recognition," *Computer Speech and Language*, vol. 24, pp. 648–662, 2010.
- [5] S.-X. Zhang, A. Ragni, and M. Gales, "Structured log-linear models for noise robust speech recognition," *IEEE Sig. Proc. Lett.*, 2010.
- [6] A. Ragni and M. Gales, "Structured discriminative models for noise robust speech recognition," in *Proc. ICASSP*, 2011.
- [7] M. Layton, "Augmented statistical models for classifying sequence data," Ph.D. dissertation, Cambridge University, 2006.
- [8] A. Acero, L. Deng, T. Kristjansson, and J. Zhang, "HMM Adaptation using Vector Taylor Series for Noisy Speech Recognition," in *Proc. ICSLP*, Beijing, China, 2000.
- [9] H. Liao and M. Gales, "Joint uncertainty decoding for robust large vocabulary speech recognition," Cambridge Uni., Tech. Rep. CUED/F-INFENG/TR552, November 2006.
- [10] T. Anastasakos, J. McDonough, R. Schwartz, and J. Makhoul, "A compact model for speaker-adaptive training," in *Proc. ICSLP*, 1996.
- [11] O. Kalinli, M. Seltzer, and A. Acero, "Noise adaptive training using a vector Taylor series approach for noise robust automatic speech recognition," in *Proc. ICASSP*, 2009.
- [12] F. Flego and M. Gales, "Discriminative adaptive training with VTS and JUD," in *Proc. ASRU*, 2009.
- [13] T. Jaakkola and D. Haussler, "Exploiting generative models in discriminative classifiers," in *Adv. in Neural Inf. Proc. Systems 11*, 1999.
- [14] N. Smith and M. Gales, "Speech recognition using SVMs," in *Adv. in Neural Inf. Proc. Systems*, 2001.
- [15] V. Venkataramani, S. Chakrabartty, and W. Byrne, "Support vector machines for segmental minimum Bayes risk decoding of continuous speech," in *Proc. ASRU*, 2003.
- [16] M. Gales, A. Ragni, H. AlDamarki, and C. Gautier, "Support vector machines for noise robust ASR," in *Proc. ASRU*, 2009.
- [17] D. Povey, "Discriminative training for large vocabulary speech recognition," Ph.D. dissertation, Cambridge University, 2004.
- [18] S. Sarawagi and W. Cohen, "Semi-markov conditional random fields for information extraction," in *Proc. NIPS*, 2005.
- [19] S.-X. Zhang and M. Gales, "Structured support vector machines for noise robust continuous speech recognition," in *Proc. Interspeech*, 2011.
- [20] S. Young, J. Odell, and P. Woodland, "Tree-based state tying for high accuracy acoustic modelling," in *Proc. ARPA Workshop HLT*, 1994.
- [21] F. Flego and M. Gales, "Factor analysis based VTS and JUD noise estimation and compensation," Cambridge University, Tech. Rep. CUED/F-INFENG/TR653, 2011.