



**HAL**  
open science

## Improving phone duration modelling using support vector regression fusion

Alexandros Lazaridis, Iosif Mporas, Todor Ganchev, George Kokkinakis,  
Nikos Fakotakis

► **To cite this version:**

Alexandros Lazaridis, Iosif Mporas, Todor Ganchev, George Kokkinakis, Nikos Fakotakis. Improving phone duration modelling using support vector regression fusion. *Speech Communication*, 2010, 53 (1), pp.85. 10.1016/j.specom.2010.07.005 . hal-00699049

**HAL Id: hal-00699049**

**<https://hal.science/hal-00699049>**

Submitted on 19 May 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Accepted Manuscript

Improving phone duration modelling using support vector regression fusion

Alexandros Lazaridis, Iosif Mporas, Todor Ganchev, George Kokkinakis, Nikos Fakotakis

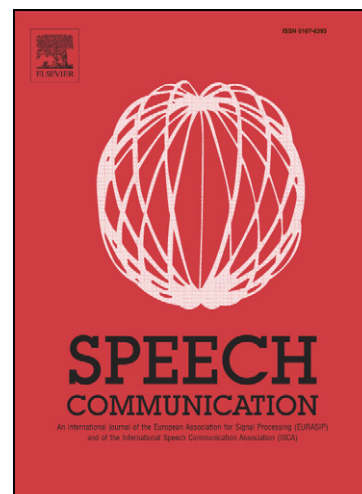
PII: S0167-6393(10)00132-9  
DOI: [10.1016/j.specom.2010.07.005](https://doi.org/10.1016/j.specom.2010.07.005)  
Reference: SPECOM 1912

To appear in: *Speech Communication*

Received Date: 30 October 2009  
Revised Date: 7 July 2010  
Accepted Date: 21 July 2010

Please cite this article as: Lazaridis, A., Mporas, I., Ganchev, T., Kokkinakis, G., Fakotakis, N., Improving phone duration modelling using support vector regression fusion, *Speech Communication* (2010), doi: [10.1016/j.specom.2010.07.005](https://doi.org/10.1016/j.specom.2010.07.005)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



# IMPROVING PHONE DURATION MODELLING USING SUPPORT VECTOR REGRESSION FUSION

*Alexandros Lazaridis, Iosif Mporas, Todor Ganchev, George Kokkinakis  
and Nikos Fakotakis*

Artificial Intelligence Group, Wire Communications Laboratory,  
Department of Electrical and Computer Engineering,  
University of Patras, 26500 Rion-Patras, Greece  
Tel. +30 2610 996496, Fax. +30 2610 997336  
alaza@upatras.gr

## ABSTRACT

In the present work, we propose a scheme for the fusion of different phone duration models, operating in parallel. Specifically, the predictions from a group of dissimilar and independent to each other individual duration models are fed to a machine learning algorithm, which reconciles and fuses the outputs of the individual models, yielding more precise phone duration predictions. The performance of the individual duration models and of the proposed fusion scheme is evaluated on the American-English KED TIMIT and on the Greek WCL-1 databases. On both databases, the SVR-based individual model demonstrates the lowest error rate. When compared to the second-best individual algorithm, a relative reduction of the mean absolute error (MAE) and the root mean square error (RMSE) by 5.5% and 3.7% on KED TIMIT, and 6.8% and 3.7% on WCL-1 is achieved. At the fusion stage, we evaluate the performance of twelve fusion techniques. The proposed fusion scheme, when implemented with SVR-based fusion, contributes to the improvement of the phone duration prediction accuracy over the one of the best individual model, by 1.9% and 2.0% in terms of relative reduction of the MAE and RMSE on KED TIMIT, and by 2.6% and 1.8% on the WCL-1 database.

***Index Terms***— Duration modelling, parallel fusion scheme, phone duration prediction, support vector regression, text-to-speech synthesis

## 1. INTRODUCTION

In Text-to-Speech synthesis (TTS) there are two major issues concerning the quality of the synthetic speech, namely the intelligibility and the naturalness (Dutoit, 1997; Klatt, 1987). The former refers to

the capability of a synthesized word or phrase to be comprehended by the average listener. The latter represents how close to the human natural speech, the synthetic speech is perceived. One of the most important factors for achieving intelligibility and naturalness in synthetic speech is the accurate modelling of prosody.

Prosody can be regarded as the implicit channel of information in the speech signal that conveys linguistic, paralinguistic and extralinguistic information related to communicative functions. Such functions are the linguistic functions of prominence (stress and accent), the phrasing, the discourse segmentation, the information about expression of emphasis, attitude, assumptions, the emotional state of the speaker, the information about the identify of the speaker (particular with respect to habitual factors). These functions provide to the listener clues supporting the recovery of the verbal message (Clark and Yallop, 1995; Laver, 1980; Laver, 1994). The accurate modelling and control of prosody in a text-to-speech system leads to synthetic speech of higher quality.

Prosody is shaped by the relative level of the fundamental frequency, the intensity and last but not least by the duration of the pronounced phones (Dutoit, 1997; Furui, 2000). The duration of the phones controls the rhythm and the tempo of speech (Yamagishi et al., 2008) and the flattening of the prosody in a speech waveform would result in a monotonous, neutral, toneless and without rhythm synthetic speech, sounding unnatural, unpleasant to the listener or sometimes even scarcely intelligible (Chen et al., 2003). Thus, the accurate modelling of phones' duration is essential in speech processing.

Several areas of speech technology, among which TTS, automatic speech recognition (ASR) and speaker recognition benefit from duration modelling. In TTS, the correct segmental duration contributes to the naturalness of synthetic speech (Chen et al., 1998; Klatt, 1976). In hidden Markov model (HMM)-based ASR, state duration models improve the speech recognition performance (Bourlard et al., 1996; Jennequin and Gauvain, 2007; Levinson, 1986; Mitchell et al., 1995; Pols et al., 1996). Finally, significant improvement of the performance in the speaker recognition task was achieved by Ferrer et al. (2003), when duration-based speech parameters were used for the characterization of the speaker's voice.

Various approaches for segment duration modelling and many factors influencing the segmental duration have been studied in the literature (Bellegarda et al., 2001; Crystal and House, 1988; Edwards and Beckman, 1988; Riley, 1992; Shih and Ao, 1997; van Santen, 1994). The features related to these factors can be extracted from several levels of linguistic information, such as the phonetic, the

morphological and the syntactic level. With respect to the way duration models are built, the duration prediction approaches can be divided in two major categories: the rule-based (Klatt, 1976) and the data-driven methods (Campbell, 1992; Chen et al., 1998; Lazaridis et al., 2007; Monkowski et al., 1995; Rao and Yegnanarayana, 2005; Riley, 1992; Takeda et al., 1989; van Santen, 1992).

The rule-based methods use manually produced rules, extracted from experimental studies on large sets of utterances, or based on previous knowledge. The extraction of these rules requires labour of expert phoneticians. In the most prominent attempt in the rule-based duration modelling category, proposed by Klatt (1976), rules which were derived by analyzing a phonetically balanced set of sentences, were used in order to predict segmental duration. These rules were based on linguistic information such as positional and prosodic factors. Initially a set of intrinsic (starting) values was assigned on each phone which was modified each time according to the extracted rules. Models of this type and similar to this were developed in many languages such as French (Bartkova and Sorin, 1987), Swedish (Carlson and Granstrom, 1986), German (Kohler, 1988) and Greek (Epitropakis et al., 1993; Yiourgalis and Kokkinakis, 1996), as well as in several dialects such as American English (Allen et al., 1987; Olive and Liberman, 1985) and Brazilian Portuguese (Simoes, 1990). The main disadvantage of the rule-based approaches is the difficulty to represent and tune manually all the linguistic factors, such as the phonetic, the morphological and the syntactic ones, which influence the segmental duration in speech. As a result, it is very difficult to collect all the appropriate (or even enough) rules without long-term devotion to this task (Klatt, 1987). Consequently the rule-based duration models are restricted to controlled experiments, where only a limited number of contextual factors are involved in order to be able to deduce the interaction among these factors and extract the corresponding rules (Rao and Yegnanarayana, 2007).

Data-driven methods for the task of phone duration modelling were developed after the construction of large databases (Kominek and Black, 2003). Data-driven approaches overcame the problem of the extraction of manual rules by employing either statistical methods or artificial neural network (ANN) based techniques which automatically produce phonetic rules and construct duration models from large speech corpora. Their main advantage is that this process is automated and thus significantly reduces the efforts that have to be spent by phoneticians.

Several machine learning methods have been used in the phone duration modelling task. The linear regression (LR) (Takeda et al., 1989) models are based on the assumption that among the features

which affect the segmental duration there is linear independency. These models achieve reliable predictions even with small amount of training data but do not model the dependency among the features. On the other hand, decision tree models (Monkowski et al., 1995) and in particular classification and regression tree (CART) models (Riley, 1992), which are based on binary splitting of the feature space, can represent the dependencies among the features but cannot insert constraints of linear independency for reliable predictions (Iwahashi and Sagisaka, 2000). Another technique which has been used on the phone duration modelling task is the sums-of-products (SOP), where the segment duration prediction is based on a sum of factors and their product terms that affect the duration (van Santen, (1992, 1994)). The advantage of these models is that they can be trained with a small amount of data. Bayesian networks models have also been introduced on the phone duration prediction task. These models incorporate a straightforward representation of the problem domain information and despite their time consuming training phase, they can make accurate predictions even when unknown values come across in some features (Goubanova and King, 2008; Goubanova and Taylor, 2000). Furthermore, instance-based algorithms (Lazaridis et al. 2007) have been used in phone duration modelling. In instance-based approaches the training data are stored and a distance function is employed during the prediction phase in order to determine which member of the training set is closer to the test instance and predict the phone duration. In a recent study (Yamagishi et al., 2008), the gradient tree boosting (GTB) (Friedman, 2001; Friedman, 2002) approach was proposed for the phone duration modelling task as an alternative to the conventional approach using regression trees. The GTB algorithm is a meta-algorithm which is based on the construction of multiple regression trees and consequently taking advantage of them.

On the task of syllable duration modelling various neural networks have been used, including feedforward neural networks (Campbell, 1992; Rao and Yegnanarayana, 2007) and recurrent neural networks (RNN) (Chen et al., 1998). Furthermore, in the case of syllable duration prediction the SVM regression model has been used in order to perform the function estimation from the training instances using non-linear mapping of the data onto a high-dimensional feature space (Rao and Yegnanarayana, 2005). Iwahashi and Sagisaka (2000) proposed a scheme for statistical modelling of prosody control in speech synthesis. It is based on a combination of regression trees and linear regression models. It offers a mechanism for evading the disadvantages inherent to one algorithm by benefiting from the

advantages provided by another algorithm. This can be explained by the observation that different algorithms perform better in different conditions.

As a result, the task of phone duration modelling based on the data-driven approaches gives the ability to overcome the time consuming labour of the manual extraction of the rules which are needed in the rule-based approaches. However, as shown by van Santen and Olive (1990), these methods are not always satisfactory for the task of phone duration prediction.

All previous studies on phone and syllable duration modelling are restricted to the use of a single linear or non-linear regression algorithm. The only exception to this trend is the work of Iwahashi and Sagisaka (2000), where a hierarchical structure for syllable duration prediction using the outputs of a phone duration model was used. However, this structure is restricted to the post-processing of a single duration prediction model, and no extension to a parallel regression fusion of the duration predictions of multiple models has been studied.

In the present work, aiming at improving the accuracy of the prediction of the segmental durations (here phone durations), we propose a fusion scheme based on the use of multiple dissimilar phone duration predictors which operate on a common input, and whose predictions are combined using a regression fusion method. The proposed scheme is based on the observation that predictors implemented with different machine learning algorithms perform differently in dissimilar conditions. Hence, we suppose that an appropriate combination of their outputs could result in a new set of more precise phone duration predictions. Thus, an appropriate fusion scheme that can learn how to combine the outputs of a number of individual predictors in a beneficial manner, will contribute to the reduction of the overall prediction error, when compared to the error of each individual predictor.

Based on this assumption, we investigate various implementations of the proposed fusion scheme and study its accuracy for duration prediction on different levels of granularity: vowels/consonants, phonetic category and individual phones. In this connection, initially, we investigate the performance of eight linear and non-linear regression algorithms, five of them already examined in previous studies (Iwahashi and Sagisaka, 2000; Lee and Oh, 1999; Riley, 1992; Takeda et al., 1989; Yamagishi et al., 2008) as individual predictors. These are based on linear regression and decision trees – model trees, regression trees and pruning decision trees. Furthermore, another two of them– the meta-learning algorithms, additive regression and bagging, using REPTrees as base classifier –are modifications of algorithms that were already studied in the phone duration prediction task (Yamagishi et al., 2008), and

finally, the support vector regression (SVR) algorithm, which to our best knowledge has not yet been employed on the phone duration prediction task. Next, the durations predicted by the individual duration models are fed as inputs to a machine learning algorithm referred to as fusion model, which uses these predictions and produces the final phone duration prediction. For the purpose of fusion, we evaluate twelve different (linear and non-linear) regression fusion techniques, which are the linear regression, decision trees, support vector regression, neural networks, meta-learning and lazy-learning algorithms, and finally average linear combination and best-case fusion.

The present study was inspired by the work of Kominek and Black (2004), where a family of acoustic models, providing multiple estimates for each boundary point, was used for segmenting a speech database, creating synthetic speech of higher quality using a corpus-based unit selection TTS system. This approach was found more robust than a single estimate, since by taking consensus values large labelling errors are less prevalent in the synthesis catalogue, which improves the resulting synthetic speech. To the extent of our knowledge, a parallel regression fusion of individual models has not yet been studied on the phone duration prediction or on the syllable duration prediction tasks. Furthermore, although SVR models have been used for syllable duration prediction (Rao and Yegnanarayana, 2005), to this end, they have not been employed on the phone duration prediction task.

The remainder of this article is organized as follows. In Section 2 we outline the proposed fusion scheme. In Section 3 we briefly outline the individual phone duration modelling algorithms, the algorithms used in the fusion scheme, the speech databases and the experimental setup used in the evaluation. The experimental results are presented and discussed in Section 4 and finally this work is concluded in Section 5.

## 2. FUSION SCHEME FOR DURATION MODELLING AND PREDICTION

Phone duration modelling, which mainly relies on regression algorithms, suffers from specific types of errors. The most commonly occurring type of error is the bias (systematic) error (Freedman et al., 2007). This error is a constant shift of the predicted phone durations from the real ones and can be estimated as the difference between the real and predicted mean durations. Other prediction errors that may occur in the phone duration modelling task are small miss-predictions and gross errors (outliers) (Freedman et al., 2007). Small miss-predictions in phone duration, i.e. less than 20 milliseconds, apart from the cases of short-duration phones such as schwas and flaps, do not significantly affect the quality



of the synthetic speech signal. However, larger than 20 ms errors have been reported to degrade the quality of synthetic speech (Wang et al., 2004). Here, we assume that an appropriate combination of the predictions of a number of dissimilar phone duration prediction models will improve the overall phone duration prediction accuracy. This is because different phone duration models will err in a dissimilar manner, and the fusion of their outputs, through a machine learning technique, would be able to learn and compensate some of these errors. Especially, we suppose that such a fusion scheme, apart from improving the overall accuracy of duration prediction, will be able to reduce the amount of gross errors.

In Fig. 1, we present the block diagram of the proposed fusion scheme, which relies on the combination of predictions that are produced by multiple dissimilar phone duration models, which operate on a common input. As the figure presents, the predictions of the individual models are introduced into the fusion stage, where a machine learning algorithm uses them for obtaining more precise phone duration predictions. The training and the operational phases of the proposed fusion scheme is discussed in the following subsections.

---

**Figure 1**

---

### 2.1 Training of the fusion scheme

The training of the proposed fusion scheme is an off-line two-step procedure, which relies on two non-overlapping datasets: the training and the development data. The training process can be summarized concisely as follows. During the first training step, the individual phone duration models are created using the training dataset. Subsequently, at the second step, these models are employed to process the development dataset. The outcome of this processing is a set of predictions, which together with the ground truth labels (manually annotated tags) serve as input for training the adjustable parameters of the fusion algorithm. This procedure can be formalized as follows:

Let us define a set of  $N$  individual phone duration prediction models,  $DM_n$ , with  $1 \leq n \leq N$ . The input feature vector,  $X_j^p$ , for the  $j$ th instance ( $1 \leq j \leq J$ ) of the phone  $p$ , which is used for training the  $N$  individual phone duration models,  $DM_n$ , is defined as:

$$X_j^p = [\theta_1, \theta_2, \dots, \theta_m, \dots, \theta_M]^T, \quad j=1,2,\dots,J, \quad (1)$$

where  $X \in \mathbb{C}^M$  ( $\mathbb{C}$  is the feature space) and  $\theta_m$  is the  $m$ th feature ( $1 \leq m \leq M$ ) of the feature vector  $X_j^p$ .

Once the individual duration models are trained, they are fed with the development dataset. The outcome of its processing is the set of phone duration predictions,  $y_j^{p,n}$ , of the  $n$ th duration model for the  $j$ th instance of the phone  $p$ , to be predicted:

$$y_j^{p,n} = f_{DM_n}^p(X_j^p), \quad j=1,2,\dots,J, \quad (2)$$

where  $y_j^{p,n} \in \mathbb{R}^N$ . The vector,  $Y_j^p$ , formed by appending the individual phone duration predictions,  $y_j^{p,n}$ , is

$$Y_j^p = \{y_j^{p,n}\}^T, \quad j=1,2,\dots,J, \quad (3)$$

where  $1 \leq n \leq N$ , for the  $j$ th instance of the phone  $p$ , together with the ground truth labels, are used in the training of the fusion algorithm. Once the fusion stage is trained, the proposed composite phone duration modelling scheme, shown in Fig. 1, is ready for operation.

## 2.2 Operation of the fusion scheme

In the operational mode, the input vector,  $X_j^p$ , for the  $j$ th instance of the phone  $p$  of the test dataset, appears as input to the  $N$  individual phone duration prediction models,  $DM_n$ , with  $1 \leq n \leq N$  (refer to Fig.1). Their outputs,  $y_j^{p,n}$ , as computed in eq. 2, form the vector of predictions  $Y_j^p$ , which serves as input for the fusion stage. At the fusion stage the vector  $Y_j^p$  is processed by the fusion algorithm, which computes the final phone duration prediction for the  $j$ th instance as:

$$O_j^p = g^p(Y_j^p), \quad j=1,2,\dots,J, \quad (4)$$

with  $O_j^p \in \mathbb{R}$ .

The fusion of multiple different predictions is expected to contribute to the reduction of the types of errors described above (first paragraph in section 2), and consequently to contribute to the decrease of the overall error rate. This expectation is based on the observation that different predictors, which rely on different machine learning algorithms, err in a dissimilar manner. Employing an appropriate fusion scheme, which is capable to learn the proper mapping between a set of noisy predictions and the true phone duration values, could turn out beneficial in terms of improved accuracy.

### 3. EXPERIMENTAL SETUP

To investigate the practical usefulness of the proposed approach, we trained several individual phone duration models, and then employed them in the fusion scheme described in Section 2. The various individual phone duration models and the fusion algorithms involved in the phone duration prediction fusion scheme, as well as the speech databases used in the experiments and the experimental protocol that was followed, are described in the following subsections.

#### 3.1 Individual phone duration models

In the present work we consider eight different machine learning algorithms for phone duration modelling, the outputs of which are then fed to the fusion model. These algorithms are well known and have successfully been used over the years, in different modelling tasks. One exception is the support vector regression (SVR) based modelling, which to this end has not been employed on the phone duration modelling task. In brief, the eight individual phone duration modelling algorithms that we consider here are:

- (i) the linear regression (LR) (Witten and Frank, 1999) using Akaike's Information Criterion (AIC) (Akaike, 1974) in backward stepwise selection (BSS) (Kohavi and John, 1997) procedure eliminating unnecessary variables of the training data,
- (ii) the m5p model tree, using a linear regression function on each leaf, and the m5pR regression tree, using a constant value on each leaf node instead (Quinlan, 1992; Wang and Witten, 1997).
- (iii) two additive regression algorithms (Friedman, 2002) and two bagging algorithms (Breiman, 1996) were used, by using two different regression trees (m5pR and REPTrees) (Kaariainen and Malinen, 2004; Quinlan, 1992; Wang and Witten, 1997) as base classifiers in each case. The latter four algorithms are meta-learning algorithms (Vilalta and Drissi, 2002) using regression trees as base classifiers.

During the training process, the additive regression algorithm builds a regression tree in each iteration, using the residuals of the previous tree as training data. The regression trees are combined together creating the final prediction function. In these two cases of additive regression meta-classification, the shrinkage parameter,  $\nu$ , indicating the learning rate, was set equal to 0.5 and the number of the regression trees,  $nt\_num$ , was set equal to ten. These values

were selected after grid-search experiments ( $v=\{0.1, 0.3, 0.5, 0.7, 0.9\}$ ,  $rt-num=\{5, 10, 15, 20\}$ ) on a randomly selected subset of the training data, of size approximately equal to 20% of the size of the training set.

In the bagging algorithm, the dataset was split in multiple subsets using a regression tree for each of them. The final prediction value is the average of the values predicted from each regression tree. In a similar manner, the number of the regression trees ( $rt-num$ ) was set equal to ten after a number of grid-search experiments ( $rt-num=\{5, 10, 15, 20\}$ ) on the randomly selected subset of the training data mentioned above.

- (iv) Finally, the support vector regression (SVR) model (Platt, 1999), which employs the sequential minimal optimization (SMO) algorithm for training a support vector classifier (Smola and Scholkopf, 1998), was used. Many kernel functions have been used in SVR such as the polynomial, the radial basis function (RBF) and the Gaussian functions (Scholkopf and Smola, 2002). In our experiments the RBF kernel was used as mapping function. The  $\epsilon$  and  $C$  parameters, where  $\epsilon \geq 0$  is the maximum deviation allowed during training and  $C > 0$  is the penalty parameter for exceeding the allowed deviation, were set equal to  $10^{-3}$  and  $10^{-1}$  respectively. This was done after a grid search ( $\epsilon=\{10^{-1}, 10^{-2}, \dots, 10^{-5}\}$ ,  $C=\{0.05, 0.1, 0.3, 0.5, 0.7, 1.0, 10, 100\}$ ) on the randomly selected subset of the training data mentioned above.

Our motivation to select these algorithms was based on previous research (Iwahashi and Sagisaka, 2000; Lee and Oh, 1999; Riley, 1992; Takeda et al., 1989; Yamagishi et al., 2008), where these algorithms were reported successful on the segmental duration modelling task. Along with the phone duration prediction task, many of these algorithms have also been used in syllable duration prediction task, supporting different languages and databases (cf. supra Section 1).

### 3.2 Fusion algorithms

The outputs of the eight individual duration models, outlined in Section 3.1, serve as the input for the fusion stage, which combines and disambiguates their predictions (refer to Fig. 1). The fusion stage can be implemented through different linear or nonlinear machine learning techniques. In this work in order to select the most advantageous fusion method, we evaluate twelve different algorithms for numerical prediction. These include the eight algorithms outlined in Section 3.1, as well as (i) the radial basis function neural network (RBFNN) with Gaussian kernel (Park and Sandberg, 1993), (ii) the instance-

based algorithm (IBK) (Aha and Kibler, 1991), which is a  $k$ -nearest neighbours classifier, (iii) the average linear combination, where the mean value of all individual prediction estimations of the first stage duration models is regarded as the final phone duration prediction value, and (iv) the best-case selection method, where for each instance the best prediction is selected among the predictions of all first-stage individual phone duration models. The selection of the best duration model per instance can be performed for different categories of clusters, such as voiced/unvoiced, vowels/consonants, phonetic category, individual phones, diphones and triphones. The last two methods, namely the average linear combination and the best case selection are well-known fusion schemes, and therefore in the present study they serve as intuitive reference points. In all the implementations listed above, the fusion method takes as input the predictions of the individual phone duration models,  $y_j^{p,n}$ , obtained at the first stage and combines them for obtaining the final prediction  $O_j^p(X_j^p)$ .

The functioning of the first eight machine learning techniques was already outlined in Section 3.1, therefore here we focus on the other four techniques. In the case of RBFNN, the  $k$ -means algorithm is used as a first step in the training process, for the estimation of the centres of the radial basis units in the hidden layer on the network. The outputs of the hidden layers are combined with linear regression. The number of clusters ( $num-cl$ ) for the  $k$ -means to generate and the minimum standard deviation ( $cl-std$ ) for the clusters were set equal to 135 and  $10^{-2}$  respectively. These parameters were determined after a grid search ( $num-cl=\{5,10, \dots, 200\}$ ,  $cl-std=\{0.001, 0.01, 0.1, 0.5\}$ ) on a randomly selected subset of the training set, consisting of approximately 20% of the set.

In the case of IBK a linear nearest neighbours search algorithm was used, employing the Euclidean distance as a distance function. Leave-one-out cross-validation was used to select the best value for  $k$ , under the restriction,  $k \leq 35$ , i.e. an upper limit of 35 nearest neighbours. The predictions from the  $k$  nearest neighbours were weighted according to the inverse distance.

In the best-case selection fusion scheme, we relied on the root mean square error (RMSE) of each phone prediction algorithm over the development data as the criterion for the selection of the best model for each case (Chen et al., 1998; Goubanova and King, 2008; Yamagishi et al., 2008). Specifically, the duration model prediction with the lowest RMSE for the cluster (vowels/consonants, phonetic category and individual phones) of each instance was selected.

### 3.3 Databases and feature set

In the evaluation experiments we used two databases: the American-English speech database CSTR US KED TIMIT (CSTR, 2001) and the Modern Greek speech prosodic database, WCL-1 (Zervas et al., 2008). KED TIMIT consists of 453 phonetically balanced sentences (3400 words approximately) uttered by a Native American male speaker. The WCL-1 prosodic database consists of 5500 words distributed in 500 paragraphs, each one of which may be a single word, a short sentence, a long sentence, or a sequence of sentences uttered by a female professional radio actress. The final corpus includes 390 declarative sentences, 44 exclamation sentences, 36 decision questions and 24 “wh” questions.

For the experiments on the KED TIMIT database, we adopted the phone set provided with the database (CSTR, 2001) which consists of 44 phones. For the experiments using the WCL-1 database we adopted the phone set provided with the database (Zervas et al., 2008) consisting of 34 phones. In all experiments, the manually labelled phone durations were used as the ground truth (reference) durations. In this work, a number of features, which have been reported successful in the literature (Crystal and House, 1988; Campbell, 1992; Klatt, 1987; Goubanova and King, 2008; Riley, 1992; van Santen, 1994), are considered for the task of phone duration modelling. From each utterance we computed 33 features, and for some of them we also used their temporal neighbours, defined on the level of the respective feature, i.e. phone-level, syllable-level, word-level. These features are summarized in the following:

- (i) eight phonetic features: the phone class (consonants/non-consonants), the phone types (vowels, diphthongs, schwa, consonants), the vowel height (high, mid or low), the vowel frontness (front, central or back), the lip rounding (rounded/unrounded), the manner of production (plosive, fricative, affricate, approximant, lateral, nasal), the place of articulation (labial, labio-dental, dental, alveolar, palatal, velar, glottal), the consonant voicing. Along with the aforementioned features, which concern each current instance, the two previous and the two next instances (temporal context information) were also used,
- (ii) three segment-level features: the phone name with the temporal context information of the neighbouring instances (previous, next), the position of the phone in the syllable and the onset-coda type (if the specific phone is before or after the vowel in the syllable).

- (iii) thirteen syllable-level features: the position type of the syllable in the word (single, initial, middle or final) with the temporal context information of the neighbouring instances (previous, next), the number of all the syllables in the word, the number of the accented syllables and the number of the stressed syllables since the last and to the next phrase break (i.e. the break index tier of ToBI with values, 0, 1, 2, 3, 4,) (cf. below, item (v)), syllable's onset-coda size (the number of phones before and after the vowel of the syllable) with the temporal context information of previous and next instances, the onset-coda type (if the consonant before and after the vowel in the syllable is voiced or unvoiced) with the temporal context information of previous and next instances, the position of the syllable in the word and the onset-coda consonant type (the manner of production of the consonant before and after the vowel in the syllable).
- (iv) two word-level features: the part-of-speech (noun, verb, adjective, etc) and the number of syllables of the word.
- (v) one phrase-level feature: the syllable break (i.e. the phrase break after the syllable) with the temporal context information of the neighbouring (two previous, two next) instances. The syllable break feature is implemented based on the break index tier of ToBI (0, 1, 2, 3, 4). The break index specifies an inventory of numbers expressing the strength of a prosodic juncture. The prosodic association of words is shown using the break index tier, by labelling the end of each word for the subjective strength of its association with the next word on a scale from 0 (strongest perceived conjoining) to 4 (most disjoint), defined as follows (Beckman and Ayers, 1994; Silverman et al., 1992; Huang et al., 2001):
- a. 0 for cases of clear phonetic marks of clitic groups,
  - b. 1 most phrase-medial word boundaries,
  - c. 2 a strong disjuncture marked by a pause or virtual pause, but with no tonal marks; i.e. a well-formed tune continues across the juncture or a disjuncture that is weaker than expected at what is tonally a clear intermediate or full intonation phrase boundary,
  - d. 3 intermediate intonation phrase boundary; i.e. marked by a single phrase tone affecting the region from the last pitch accent to the boundary,
  - e. 4 full intonation phrase boundary; i.e. marked by a final boundary tone after the last phrase tone.

(vi) six accentual features: the ToBI accents and boundary tones with the temporal context information of the neighbouring (previous, next) instances, the last-next accent (the number of the syllables since the last and to the next accented syllable) and we also included the stressed-unstressed syllable feature (if the syllable is stressed or not) and the accented-unaccented syllable feature (if the syllable is accented or not) with the temporal context information of the neighbouring (two previous, two next) instances.

The overall size of the feature vector, which was used for the individual phone duration models, including the aforementioned features and their temporal context information as reported above (one or two previous and next instances on the level of the respective feature, phone-level, syllable-level, word-level) is 93.

In all experiments we followed an experimental protocol based on 10-fold cross-validation. Specifically, in each fold the training data were split in two portions, the training dataset and the development dataset. The former, amounting to approximately 60% of the full dataset, was used for the training of the individual phone duration predictors, and the latter, amounting to approximately 30% of the full dataset, for the training of the fusion algorithm. Furthermore, the test dataset, amounting to approximately 10% of the full dataset, was used for evaluating the performance of the eight individual duration prediction algorithms, as well as the performance of the fusion scheme.

### 3.4 Performance metrics

The experimental results were evaluated using the two most commonly used figures of merit, namely the mean absolute error (MAE) and the root mean squared error (RMSE), between the predicted duration and the actual (reference) duration of each phone (Chen et al., 1998; Goubanova and King, 2008; van Santen, 1992; Yamagishi et al., 2008). Due to the squaring of values in the RMSE, large errors (outliers) are weighted heavily, which makes this figure of metric more sensitive to outliers than the MAE (Witten and Frank, 1999). This sensitivity of the RMSE makes it a more illustrative measurement concerning the outliers, e.g. the gross errors, in comparison to the MAE.

## 4. EXPERIMENTAL RESULTS

In the present work, we consider clustering of the instances on the basis of (i) vowels/consonants categorization, (ii) phonetic categories and (iii) individual phones. This offers different degree of detail and allows us to gain insights about the advantages and disadvantages of each algorithm. The same



clustering of the instances is used in the best-case selection fusion model. In this fusion model, as mentioned in Section 3.2, the criterion for the selection of the best model for each case is the RMSE of each phone duration prediction algorithm over the development data, as in Chen et al. (1998), Goubanova and King (2008) and Yamagishi et al. (2008). Specifically, the phone duration model prediction with the lowest RMSE per cluster (vowels/consonants, phonetic categories, individual phone) of each instance is selected.

#### *4.1 Duration prediction with individual phone duration models*

As a first step we examined the performance of the eight individual algorithms on both databases using the entire feature set described in Section 3. The RMSE, the MAE and the standard deviation of the absolute error (STD of AE) for all individual algorithms specified in Section 3.1 are shown in Table 1, where Table 1 (a) presents the results obtained on KED TIMIT and Table 1 (b) the ones on the WCL-1 database. The results of the best performing model, among the eight individual prediction models, are in bold. As can be seen, on both databases, the proposed support vector regression (SMOreg) model, implemented with the SMO regression algorithm, outperforms all the other models. Specifically, on the KED TIMIT database the SMOreg model outperformed the second-best model, i.e. the meta-classifier additive regression using m5pR model, by approximately 5.5% and 3.7% in terms of MAE and RMSE respectively. On the WCL-1 database the SMOreg model outperformed the second-best model, i.e. the Linear Regression model, by approximately 6.8% and 3.7% in terms of MAE and RMSE respectively. This advantage of the SMOreg models, on both databases, is owed to the advantage of SVMs to cope better with high-dimensional feature spaces (Vapnik, 1995; Vapnik, 1998), when compared to the other classification and regression algorithms.

##### Table 1 (a)

##### Table 1 (b)

In Table 2 we present the performance per phonetic category as well as for the vowel/consonant categorization of the eight individual phone duration models, implemented by different algorithms, on the KED TIMIT (Table 2 (a)) and the WCL-1 (Table 2 (b)) databases. As can be seen, the SMOreg

model demonstrates the lowest RMSE on both databases, in all cases except for the Affricates on KED TIMIT, where the lowest RMSE is observed for the REPTrees and the SMOreg achieves the second-best performance.

**Table 2 (a)**

**Table 2 (b)**

In Table 3 the phone duration prediction results obtained on the level of individual phones are presented. Specifically, Table 3 (a) shows the RMSE for the 44 phone set of the KED TIMIT database and Table 3 (b) for the 34 phone set of the WCL-1 database. The results for the best performing algorithm are in bold. As shown in the tables, despite the fact that the SMOreg model demonstrates the highest overall performance on both databases (refer to Table 1), in one phonetic category (Affricates in Table 2 (a)) and in some particular phones, such as *ch*, *ay*, etc (Table 3), other models offer a higher phone duration prediction accuracy. For instance, on the KED TIMIT database, the highest accuracy for the phone *ch* is observed for the Linear Regression model, while for the phone *ay* the highest accuracy is for the m5p model (refer to Table 3 (a)). These specific results, and other similar cases shown in the Table 3, are in support of our observation that different algorithms perform better in different phonetic categories and phones. This indicates that an appropriate fusion of the outputs of the individual phone duration prediction models could be beneficial for reducing the overall error rate. Experimental results for various implementations of the fusion stage are presented in Section 4.2.

**Table 3 (a)**

**Table 3 (b)**

*4.2 Duration prediction with the proposed fusion scheme*

In the following, we report the evaluation results for the twelve fusion algorithms outlined in Section 3.2. In Table 4, we present the results obtained in the evaluation of the *average* linear combination and the *best-case selection* techniques for different cases: overall performance, vowel/consonant categories,

phonetic categories and individual phones. The results for the best individual phone duration prediction model (SMOreg) are duplicated from Table 1 for the purpose of direct comparison. We consider the *average* and the *best-case selection* techniques as intuitive reference points against which the performance of the other ten fusion algorithms, evaluated here, is compared. As can be seen in Table 4 (a), for KED TIMIT, and in Table 4 (b), for the WCL-1 database, in the cases of vowel/consonant categories and phonetic categories the best-case selection algorithm did not offer advantage over the best individual model, SMOreg, since the SMOreg outperforms the other individual predictors in all categories (Table 2). As discussed above the only exception is Affricates on KED TIMIT, where the additive regression with REPTrees algorithm is the best, but is not sufficient for significant advantage of the fusion scheme.

Concerning the clustering according to individual phones, the best-case selection method slightly outperformed the best individual duration model on both databases. In detail, the best-case selection outperformed the SMOreg model by approximately 0.2% and 0.9% in terms of MAE and RMSE on the KED TIMIT database, and by approximately 0.6% and 0.5% on the WCL-1 database. As can be seen in Table 3, this is owed to the fact that, for both databases, in approximately 35-40% of the phones the best performing algorithm is not the SMOreg. Consequently, the best-case selection fusion scheme outperforms the best individual phone duration prediction model (SMOreg).

Table 4 (a)

Table 4 (b)

In Table 5, we present results for the remaining ten fusion algorithms (refer to Section 3.2): the Linear Regression, the m5p model tree, the m5pR regression tree, the additive regression algorithms based on m5pR and REPTrees, the bagging algorithms based on m5pR and REPTree, the instance based learning (IBK), the support vector regression (SVR) which implements the sequential minimal optimization (SMO) algorithm, and the radial basis function neural network (RBFNN). The best fusion result is shown in bold. For the reason of comparison, in Table 5 we duplicate the results for the best individual phone duration model, SMOreg. As can be seen from the results on the KED TIMIT (Table 5 (a)) and the WCL-1 (Table 5 (b)) databases, the SMOreg fusion model outperformed all the other

fusion models that were evaluated here. It is also noteworthy to mention that only the SMOREg fusion model outperformed the best individual duration prediction model. Specifically, the SMOREg fusion model outperformed the individual SMOREg predictor by approximately 1.9% and 2.0% in terms of MAE and RMSE on KED TIMIT, and by approximately 2.6% and 1.8% on the WCL-1 database, respectively. Furthermore, we should point out that the SMOREg fusion model apart from reducing the overall error also reduced the outliers. Specifically, in comparison to the best individual predictor, i.e. the SMOREg model, the SMOREg fusion model reduced the standard deviation of the absolute error (STD of AE), by approximately 2.1% on KED TIMIT and by approximately 1.2% on the WCL-1 database, respectively.

**Table 5 (a)**

**Table 5 (b)**

Finally, in order to investigate the statistical significance of the difference between the results for the best individual phone duration model (SMOREg) and the results for the best fusion scheme (fusion with SMOREg algorithm) the Wilcoxon test (Wilcoxon, 1945) was carried out. The Wilcoxon test showed that on both databases, the difference between the results for the best individual model and these for the fusion scheme is statistically significant. Specifically, for a significance level of 0.05 the Wilcoxon test estimated a  $p$ -value equal to  $5.77e^{-09}$  and  $3.5e^{-11}$  on KED TIMIT and WCL-1 databases, respectively. Consequently, the fusion scheme contributes to the improvement of the accuracy of phone duration prediction, in comparison to best predictor among all evaluated individual phone duration prediction models.

## 5. SUMMARY AND CONCLUSIONS

In this work we studied the accuracy of various machine learning algorithms on the task of phone duration modelling. The experimental results showed that on this task, Support Vector Machines (SVM), as a regression model, outperforms various other machine learning techniques. Specifically, in terms of relative decrease of the mean absolute error and root mean square error, the SMO regression

model outperformed the second-best model by approximately 5.5% and 3.7% on KED TIMIT, and by approximately 6.8% and 3.7% on the WCL-1 database, respectively.

Furthermore, the proposed fusion scheme, which combines predictions from multiple individual phone duration models, operating on a common input, takes advantage of the observation that different prediction algorithms perform better in different situations. The experimental validation demonstrated that the fusion scheme improves the accuracy of phone duration prediction. The SVM-based fusion algorithm was found to outperform all other fusion techniques. Specifically, the fusion scheme based on the SVM regression algorithm outperformed the best individual predictor (SVM regression) by approximately 1.9% and 2.0% in terms of relative reduction of the mean absolute error and root mean square error respectively, on the KED TIMIT database, and by 2.6% and 1.8% on the WCL-1 database, respectively.

## 6. ACKNOWLEDGEMENTS

The authors are thankful to the anonymous reviewers for their valuable comments and corrections on an earlier version of our manuscript, which contributed to the significant improvement of the quality of this article.

## 7. REFERENCES

- Aha, D., Kibler, D., 1991. Instance-based learning algorithms. *Machine Learning*. 6, 37-66.
- Allen, J., Hunnicutt, S., Klatt, D.H., 1987. *From text to speech: The MITalk system*. Cambridge University Press, Cambridge.
- Akaike, H., 1974. A new look at the statistical model identification. *IEEE Trans. on Automatic Control*. 19(6), 716-723.
- Bartkova, K., Sorin, C., 1987. A model of segmental duration for speech synthesis in French. *Speech Communication*. 6, 245-260.
- Beckman, M.E., Ayers, G.M., 1994. *Guidelines for ToBI labelling*. The Ohio State University, [http://www.ling.ohio-state.edu/research/phonetics/E\\_ToBI](http://www.ling.ohio-state.edu/research/phonetics/E_ToBI)
- Bellegarda, J.R., Silverman, K.E.A., Lenzo, K., Anderson, V., 2001. Statistical prosody modelling: from corpus design to parameter estimation. *IEEE Trans. on Speech and Audio Processing*. 9(1), 52-66.
- Breiman, L., 1996. Bagging predictors. *Machine Learning*. 24(2), 123-140.

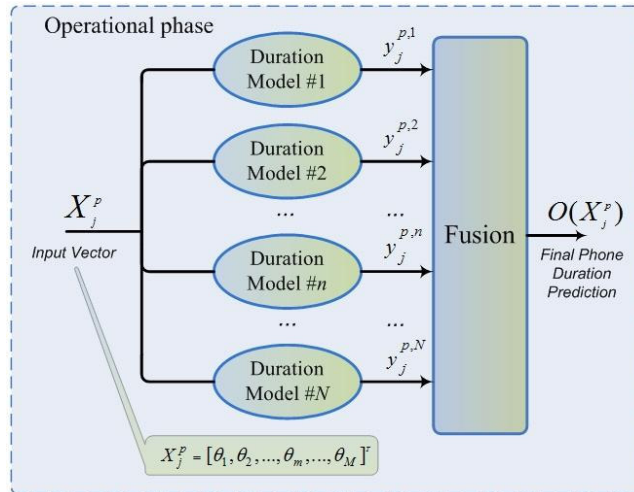
- Bourlard, H., Hermansky, H., Morgan, N., 1996. Towards increasing speech recognition error rates. *Speech Communication*. 18(3), 205-231.
- Campbell, W.N., 1992. Syllable based segment duration, in: Bailly, G., Benoit, C., Sawallis, T.R. (Eds.), *Talking Machines: Theories, Models and Designs*, Elsevier, Amsterdam, pp. 211-224.
- Carlson, R., Granstrom, B., 1986. A search for durational rules in real speech database. *Phonetica*. 43, 140-154.
- Chen, S.H., Hwang, S.H., Wang, Y.R., 1998. An RNN-based prosodic information synthesizer for Mandarin text-to-speech. *IEEE Trans. on Speech and Audio Processing*. 6(3), 226-239.
- Chen, S.H., Lai, W.H., Wang, Y.R., 2003. A new duration modeling approach for Mandarin speech. *IEEE Trans. on Speech and Audio Processing*. 11(4), 308-320.
- Clark, J., Yallop, C., 1995. *An introduction to phonetics and phonology*, second ed. Blackwell, Oxford.
- Crystal, T.H., House, A.S., 1988. Segmental durations in connected-speech signals: Current results. *Journal of the Acoustical Society of America*. 83(4), 1553-1573.
- CSTR, 2001. CSTR US KED TIMIT. University of Edinburgh, [http://www.festvox.org/dbs/dbs\\_kdt.html](http://www.festvox.org/dbs/dbs_kdt.html).
- Dutoit, T., 1997. *An Introduction to Text-To-Speech Synthesis*, Kluwer Academic Publishers, Dordrecht.
- Edwards, J., Beckman, M.E., 1988. Articulatory timing and the prosodic interpretation of syllable duration. *Phonetica*. 45, 156-174.
- Epitropakis, G., Tambakas, D., Fakotakis, N., Kokkinakis, G., 1993. Duration modelling for the Greek language. In *Proc. of EUROSPEECH-1993*, Berlin, Germany, pp. 1995-1998.
- Ferrer, L., Bratt, H., Gadde, V.R.R., Kajarekar, S.S., Shriberg, E., Sonmez, K., Stolcke, A., Venkataraman, A., 2003. Modeling duration patterns for speaker recognition. In *Proc. of EUROSPEECH-2003*, Geneva, Switzerland, pp. 2017-2020.
- Freedman, D., Pisani, R., Purves, R., 2007. *Statistics*, fourth ed. W.W. Norton & Company, New York London.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann. Statist.* 29(5), 1189-1232.
- Friedman, J.H., 2002. Stochastic gradient boosting. *Comput. Statist. Data Anal.* 38(4), 367-378.
- Furui, S., 2000. *Digital Speech Processing, Synthesis, and Recognition*, second ed. Marcel Dekker, New York.
- Goubanova, O., King, S., 2008. Bayesian networks for phone duration prediction. *Speech Communication*. 50(4), 301-311.
- Goubanova, O., Taylor, P., 2000. Using Bayesian belief networks for model duration in text-to-speech systems. In *Proc. of ICSLP-2000*, Beijing, China, pp. 427-430.
- Huang, X., Acero, A., Hon, H.W., 2001. *Spoken Language Processing: a guide to theory, algorithm, and system development*. Prentice Hall, Redmond.

- Iwahashi, N., Sagisaka, Y., 2000. Statistical modeling of speech segment duration by constrained tree regression. *IEICE Trans. Inform. Systems*. E83-D(7), 1550-1559.
- Jennequin, N., Gauvain, J.L., 2007. Modeling duration via lattice rescoring. In *Proc. of ICASSP-2007*, Honolulu, Hawaii, pp. 641-644.
- Kaariainen, M., Malinen, T., 2004. Selective Rademacher penalization and reduced error pruning of decision trees. *Journal of Machine Learning Research*. 5, 1107-1126.
- Klatt, D.H., 1976. Linguistic uses of segmental duration in English: Acoustic and perceptual evidence. *Journal of the Acoustic Society of America*. 59, 1209-1221.
- Klatt, D.H., 1987. Review of text-to-speech conversion for English. *Journal of the Acoustical Society of America*. 82(3), 737-793.
- Kohavi, R., John, G.H., 1997. Wrappers for feature subset selection. *Artificial Intelligence*. 97(1-2), 273-324.
- Kohler, K.J., 1988. Zeistrukturierung in der Sprachsynthese. *ITG-Tagung Digitalc Sprachverarbeitung*. 6, 165-170.
- Kominek, J., Black, A.W., 2003. CMU ARCTIC databases for speech synthesis. CMU-LTI-03-177, Language Technologies Institute, School of Computer Science, Carnegie Mellon University.
- Kominek, J., Black, A.W., 2004. A family-of-models approach to HMM-based segmentation for unit selection speech synthesis. In *Proc. of INTERSPEECH-2004*, Jeju Island, Korea, pp. 1385-1388.
- Laver, J., 1980. *The phonetic description of voice quality*. Cambridge University Press, Cambridge.
- Laver, J., 1994. *Principles of Phonetics*. Cambridge University Press, Cambridge.
- Lazaridis, A., Zervas, P., Kokkinakis, G., 2007. Segmental duration modeling for Greek speech synthesis. In *Proc. of IEEE ICTAI-2007*, Patras, Greece, pp. 518-521.
- Lee, S., Oh, Y., 1999. CART-based modelling of Korean segmental duration. In *Proc. of Oriental COCOSDA-1999*, Taipei, Taiwan, pp. 109-112.
- Levinson, S., 1986. Continuously variable duration hidden Markov models for speech analysis. In *Proc. of ICASSP-1986*, Tokyo, Japan, pp.1241-1244.
- Mitchell, C., Harper, M., Jamieson, L., Helzermam, R., 1995. A parallel implementation of a hidden Markov model with duration modeling for speech recognition. *Digital Signal Processing*. 5, 43-57.
- Monkowski, M.D., Picheny, M.A., Rao, P.S., 1995. Context dependent phonetic duration models for decoding conversational speech. In *Proc. of ICASSP-1995*, Detroit, Michigan, USA, pp. 528-531.
- Olive, J.P., Liberman, M.Y., 1985. Text to speech - an overview. *Journal of the Acoustical Society of America*. 78 (Suppl. 1):S6.
- Park, J., Sandberg, I.W., 1993. Approximation and radial-basis-function networks. *Neural Computation*. 5(2), 305-316.

- Platt, J., 1999. Fast training of support vector machines using sequential minimal optimization, in: Scholkopf, B., Burges, C., Smola, A. (Eds.), *Advances in kernel methods: Support vector learning*, MIT Press, Cambridge, pp. 185-208.
- Pols, L.C.W., Wang, X., ten Bosch, L.F.M., 1996. Modelling of phone duration (using the TIMIT database) and its potential benefit for ASR. *Speech Communication*. 19, 161-176.
- Quinlan, R.J., 1992. Learning with continuous classes. In *Proc. of 5th Australian Joint Conference on Artificial Intelligence*, Singapore, pp. 343-348.
- Rao, K.S., Yegnanarayana, B., 2005. Modeling syllable duration in Indian languages using support vector machines. In *Proc. of ICISIP-2005, India*, pp. 258-263.
- Rao, K.S., Yegnanarayana, B., 2007. Modeling durations of syllables using neural networks. *Computer Speech & Language*. 21(2), 282-295.
- Riley, M., 1992. Tree-based modelling for speech synthesis, in: Bailly, G., Benoit, C., Sawallis, T.R. (Eds.), *Talking Machines: Theories, Models and Designs*, Elsevier, Amsterdam, pp. 265-273.
- Scholkopf, B., Smola, A.J., 2002. *Learning with Kernels*. MIT Press, Cambridge.
- Shih, C., Ao, B., 1997. Duration study for the Bell Laboratories Mandarin text-to-speech system, in: van Santen, J., Sproat, R., Olive, J., Hirschberg, J. (Eds.), *Progress in Speech Synthesis*, Springer, New York, pp. 383-399.
- Silverman, K., Beckman, M., Pitrelli, J., Ostendorf, M., Wightman, C., Price, P., Pierrehumbert, J., Hirschberg, J., 1992. ToBI: A standard for labeling English Prosody. In *Proc. of ICSLP-1992, Banff, Alberta, Canada*, pp. 867-870.
- Simoes, A.R.M., 1990. Predicting sound segment duration in connected speech: An acoustical study of Brazilian Portuguese. In *Proc. of Workshop on Speech Synthesis, AuTrans, France*, pp.173-176.
- Smola, A.J., Scholkopf, B., 1998. A tutorial on support vector regression. Royal Holloway College, London, U.K., NeuroCOLT Tech. Rep. TR 1998-030.
- Takeda, K., Sagisaka, Y., Kuwabara, H., 1989. On sentence-level factors governing segmental duration in Japanese. *Journal of Acoustic Society of America*. 86(6), 2081-2087.
- van Santen, J.P.H., 1992. Contextual effects on vowel durations. *Speech Communication*. 11(6), 513-546.
- van Santen, J.P.H., 1994. Assignment of segmental duration in text-to-speech synthesis. *Computer Speech & Language*. 8(2), 95-128.
- van Santen, J., Olive, J., 1990. The analysis of contextual effects on segmental duration. *Computer, Speech & Language*. 4(4), 359-390.
- Vapnik, V., 1995. *The Nature of Statistical Learning Theory*. Springer, New York.
- Vapnik, V., 1998. *Statistical Learning Theory*. Wiley, New York.



- Vilalta, R., Drissi, Y., 2002. A perspective view and survey of meta-learning. *Artificial Intelligence Review*. 18(2), 77-95.
- Wang, Y., Witten, I.H., 1997. Inducing model trees for continuous classes. In *Poster papers of the 9th European Conference on Machine Learning*, Prague, Czech Republic, pp. 128-137.
- Wang, L., Zhao, Y., Chu, M., Zhou, J., Cao, Z., 2004. Refining segmental boundaries for TTS database using fine contextual-dependent boundary models. In *Proc. of ICASSP-2004*, Montreal, Quebec, Canada, pp. 641-644.
- Wilcoxon, F., 1945. Individual comparisons by ranking methods. *Biometrics*. 1, 80-83.
- Witten, H.I., Frank, E., 1999. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufman Publishing, San Francisco.
- Yamagishi, J., Kawai, H., Kobayashi, T., 2008. Phone duration modeling using gradient tree boosting. *Speech Communication*. 50(5), 405-415.
- Yiourgalis, N., Kokkinakis, G., 1996. A TtS system for the Greek language based on concatenation of formant coded segments. *Speech Communication*. 19(1), 21-39.
- Zervas, P., Fakotakis, N., Kokkinakis, G., 2008. Development and evaluation of a prosodic database for Greek speech synthesis and research. *Journal of Quantitative Linguistics*. 15(2), 154-184.



**Fig. 1.** Block diagram of the proposed fusion scheme, which exploits multiple dissimilar phone duration predictors, operating on a common input.

**Table 1.** Mean Absolute Error (MAE), standard deviation of absolute error (STD of AE) and Root Mean Square Error (RMSE) (in milliseconds) for the eight individual phone duration prediction algorithms on: (a) the KED TIMIT database, and (b) the WCL-1 database.

(a) results on the KED TIMIT database

<i>Individual models (KED TIMIT database)</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
SMOreg	<b>14.95</b>	<b>14.11</b>	<b>20.56</b>
Add. Reg. m5pR (Yamagishi et al., 2008)	15.82	14.34	21.35
Add. Reg. REPTrees	16.29	15.06	22.19
Bagging m5pR (Lee and Oh, 1999)	16.51	14.76	22.14
m5p (Iwahashi and Sagisaka, 2000)	16.62	14.77	22.23
Bagging REPTrees	16.69	15.89	23.04
m5pR (Riley, 1992)	16.93	15.16	22.72
Linear Regression (Takeda et al., 1989)	17.15	15.16	22.89

(b) results on the WCL-1 database

<i>Individual models (WCL-1 database)</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
SMOreg	<b>16.78</b>	<b>18.81</b>	<b>25.21</b>
Linear Regression (Takeda et al., 1989)	18.00	19.02	26.19
Add. Reg. REPTrees	18.08	19.97	26.94
Add. Reg. m5pR (Yamagishi et al., 2008)	18.13	19.16	26.38
Bagging m5pR (Lee and Oh, 1999)	18.14	19.63	26.72
m5p (Iwahashi and Sagisaka, 2000)	18.31	20.08	27.17
Bagging REPTrees	18.93	20.32	27.77
m5pR (Riley, 1992)	19.07	20.10	27.71

**Table 2.** Root Mean Square Error (in milliseconds) per phonetic category for the eight individual phone duration prediction algorithms on: (a) the KED TIMIT database, and (b) the WCL-1 database.

(a) results on the KED TIMIT database

<i>KED TIMIT database</i>	<i>LR</i>	<i>m5p</i>	<i>m5pR</i>	<i>Additive Regression</i>		<i>Bagging</i>		<i>SMOreg</i>
				<i>m5pR</i>	<i>REPTrees</i>	<i>m5pR</i>	<i>REPTrees</i>	
Vowel	24.56	24.18	25.46	23.67	24.87	24.78	26.34	<b>22.72</b>
Consonant	21.72	20.86	20.74	19.69	20.24	20.24	20.60	<b>19.02</b>
<i>Phonetic category</i>	<i>LR</i>	<i>m5p</i>	<i>m5pR</i>	<i>Additive Regression</i>		<i>Bagging</i>		<i>SMOreg</i>
				<i>m5pR</i>	<i>REPTrees</i>	<i>m5pR</i>	<i>REPTrees</i>	
Vowel	24.56	24.18	25.46	23.67	24.87	24.78	26.34	<b>22.72</b>
Affricate	22.44	24.41	23.48	22.86	<b>21.72</b>	22.96	23.34	21.88
Approximant	22.23	22.44	23.09	21.77	22.56	22.59	24.07	<b>20.42</b>
Fricative	22.51	21.67	21.10	20.19	20.69	20.63	20.96	<b>19.63</b>
Lateral	21.16	20.98	21.18	20.29	21.16	20.52	21.89	<b>19.77</b>
Nasal	18.59	17.88	17.80	17.11	16.94	17.28	17.57	<b>16.53</b>
Plosive	23.39	22.07	21.62	20.26	20.97	21.04	20.80	<b>19.61</b>

(b) results on the WCL-1 database

<i>WCL-1 database</i>	<i>LR</i>	<i>m5p</i>	<i>m5pR</i>	<i>Additive Regression</i>		<i>Bagging</i>		<i>SMOreg</i>
				<i>m5pR</i>	<i>REPTrees</i>	<i>m5pR</i>	<i>REPTrees</i>	
Vowel	24.22	24.68	26.04	24.51	25.18	24.91	26.62	<b>23.12</b>
Consonant	27.86	29.25	29.13	27.97	28.44	28.27	28.77	<b>26.57</b>
<i>Phonetic category</i>	<i>LR</i>	<i>m5p</i>	<i>m5pR</i>	<i>Additive Regression</i>		<i>Bagging</i>		<i>SMOreg</i>
				<i>m5pR</i>	<i>REPTrees</i>	<i>m5pR</i>	<i>REPTrees</i>	
Vowel	24.22	24.68	26.04	24.51	25.18	24.91	26.62	<b>23.12</b>
Affricate	24.72	27.61	24.62	22.22	20.74	25.83	23.57	<b>20.73</b>
Fricative	25.67	27.04	26.93	25.79	26.23	25.57	26.45	<b>23.95</b>
Liquid	19.46	19.19	19.55	18.84	17.83	18.47	18.02	<b>16.38</b>
Nasal	22.44	22.94	23.11	22.27	22.15	22.18	22.27	<b>20.62</b>
Plosive	34.22	36.33	36.03	34.62	35.64	35.48	36.14	<b>33.69</b>

**Table 3 (a).** Root Mean Square Error (in milliseconds) per phone for the eight individual phone duration prediction algorithms on the KED TIMIT database

<i>KED TIMIT database</i>	<i>IPA Symbols</i>	<i>LR</i>	<i>m5p</i>	<i>m5pR</i>	<i>Additive Regression</i>		<i>Bagging</i>		<i>SMOreg</i>
					<i>m5pR</i>	<i>REPTrees</i>	<i>m5pR</i>	<i>REPTrees</i>	
<i>aa</i>	ɒ	27.81	25.57	28.01	<b>24.57</b>	27.27	26.71	29.22	25.64
<i>ae</i>	æ	31.40	30.97	31.67	29.75	31.19	31.64	33.13	<b>29.23</b>
<i>ah</i>	ʌ	19.67	22.34	22.27	20.31	21.50	20.88	22.27	<b>19.38</b>
<i>ao</i>	ɔ:	32.79	<b>29.21</b>	32.95	30.54	32.11	32.66	33.29	29.65
<i>aw</i>	aʊ	33.46	<b>32.89</b>	37.35	34.55	38.49	37.25	40.02	33.07
<i>ax</i>	ə	16.10	15.66	16.06	15.16	15.56	15.54	15.93	<b>14.80</b>
<i>ay</i>	aɪ	37.12	<b>32.78</b>	38.37	34.43	34.04	36.64	37.23	34.51
<i>b</i>	b	23.89	22.36	23.42	22.24	23.03	22.52	<b>21.19</b>	21.33
<i>ch</i>	tʃ	<b>19.69</b>	23.34	21.17	20.48	20.43	19.89	22.36	20.57
<i>d</i>	d	20.77	19.36	19.66	19.12	20.05	19.32	20.54	<b>18.26</b>
<i>dh</i>	ð	17.56	16.03	15.72	15.19	<b>14.57</b>	15.16	15.30	15.14
<i>dx</i>	d/t	11.08	10.38	11.30	9.99	<b>8.78</b>	9.54	8.86	9.63
<i>eh</i>	ɛ	20.94	20.39	22.50	21.41	21.44	21.32	22.61	<b>19.05</b>
<i>el</i>	əl	21.39	27.24	21.52	20.79	<b>18.98</b>	19.97	21.05	22.32
<i>em</i>	əm	13.61	15.31	10.51	10.44	<b>10.13</b>	10.28	13.99	11.58
<i>en</i>	ən	22.26	24.60	25.01	23.18	<b>20.67</b>	22.44	21.80	21.01
<i>er</i>	ɜr	28.41	29.28	28.73	27.09	27.77	28.15	29.87	<b>25.29</b>
<i>ey</i>	eɪ	27.76	26.99	29.43	28.12	29.90	28.72	31.36	<b>26.73</b>
<i>f</i>	f	22.84	23.90	21.52	20.09	21.05	21.08	22.43	<b>18.91</b>
<i>g</i>	g	18.23	17.14	18.73	17.04	17.65	17.88	17.62	<b>16.22</b>
<i>hh</i>	h	19.13	18.79	18.82	18.52	18.73	18.28	18.73	<b>17.54</b>
<i>ih</i>	ɪ	19.38	19.76	20.16	19.09	19.81	19.82	20.86	<b>17.53</b>
<i>iy</i>	i:	23.04	23.06	23.87	22.05	24.93	23.27	25.39	<b>20.99</b>
<i>jh</i>	dʒ	24.36	25.22	25.14	24.56	<b>22.68</b>	25.08	24.07	22.85
<i>k</i>	k	22.18	21.82	20.62	18.65	18.63	19.94	18.93	<b>17.64</b>
<i>l</i>	l	21.13	20.18	21.14	20.24	21.39	20.58	21.98	<b>19.47</b>
<i>m</i>	m	16.07	15.32	16.20	15.45	16.19	15.81	17.04	<b>14.38</b>
<i>n</i>	n	18.69	17.65	17.29	16.70	<b>16.18</b>	16.80	16.32	16.19
<i>ng</i>	ŋ	22.38	20.86	20.13	<b>19.90</b>	20.88	20.61	22.41	20.91
<i>ow</i>	Oʊ	28.12	28.98	28.93	27.20	28.85	27.73	30.68	<b>25.54</b>
<i>oy</i>	ɔɪ	<b>25.45</b>	30.16	34.58	28.81	30.61	33.13	34.72	31.19
<i>p</i>	p	25.06	24.90	22.50	21.05	21.32	21.94	21.25	<b>20.45</b>
<i>r</i>	r/ɹ	19.20	18.84	20.18	19.28	20.11	19.92	21.18	<b>18.25</b>
<i>s</i>	s	26.37	24.47	24.31	23.46	24.45	24.36	24.54	<b>23.21</b>
<i>sh</i>	ʃ	19.71	21.72	19.28	18.30	20.53	18.49	20.27	<b>16.41</b>
<i>t</i>	t	28.18	25.60	25.06	23.64	25.14	24.72	24.93	<b>23.37</b>
<i>th</i>	θ	24.09	26.39	29.14	25.58	<b>21.31</b>	25.21	22.59	22.05
<i>uh</i>	ʊ	20.64	20.61	23.10	20.45	25.35	22.68	26.16	<b>19.88</b>
<i>uw</i>	u:	27.65	27.73	29.05	28.00	30.35	29.40	33.64	<b>24.97</b>
<i>v</i>	v	17.26	17.31	16.72	16.93	17.15	16.66	17.34	<b>16.26</b>
<i>w</i>	w	20.28	20.09	22.35	19.81	20.93	20.89	22.59	<b>19.12</b>
<i>y</i>	j	18.36	19.08	18.85	18.80	19.42	19.22	20.56	<b>16.34</b>
<i>z</i>	z	22.38	20.42	19.94	19.07	19.37	19.24	19.10	<b>18.99</b>
<i>zh</i>	ʒ	25.60	28.40	25.25	<b>22.62</b>	26.38	23.95	27.28	24.66

**Table 3 (b).** Root Mean Square Error (in milliseconds) per phone for the eight individual phone duration prediction algorithms on the WCL-1 database

WCL-1 database	IPA Symbols	LR	m5p	m5pR	Additive Regression		Bagging		SMOreg
					m5pR	REPTrees	m5pR	REPTrees	
<i>a</i>	a	24.25	25.85	25.76	24.07	24.57	24.83	26.07	<b>22.71</b>
<i>b</i>	b	21.05	24.66	24.41	21.84	22.05	22.33	22.53	<b>20.20</b>
<i>c</i>	tʃ	24.16	28.40	25.43	22.48	<b>20.29</b>	26.62	23.62	20.85
<i>D</i>	ð	<b>22.64</b>	23.08	25.08	24.24	24.39	24.00	26.25	<b>22.64</b>
<i>d</i>	d	<b>19.33</b>	20.40	23.54	21.01	21.44	21.39	24.61	20.10
<i>e</i>	e	25.05	25.11	26.69	25.62	26.79	25.71	26.48	<b>24.05</b>
<i>f</i>	f	30.13	34.11	30.41	29.95	33.13	<b>29.50</b>	31.94	29.56
<i>G</i>	ʒ	30.72	37.75	37.14	31.68	31.56	33.82	33.99	<b>29.89</b>
<i>g</i>	g	34.43	38.79	34.94	35.14	40.05	34.30	37.80	<b>33.85</b>
<i>h</i>	ŋ	24.73	25.91	26.39	24.69	24.87	23.78	25.88	<b>23.50</b>
<i>i</i>	i	24.17	24.30	25.54	24.27	24.68	24.68	26.98	<b>23.09</b>
<i>j</i>	dʒ	25.65	26.18	23.13	21.75	<b>20.52</b>	24.39	23.48	21.50
<i>K</i>	c	45.75	44.50	45.47	43.94	46.86	45.73	45.82	<b>43.28</b>
<i>k</i>	k	<b>42.27</b>	46.15	44.65	43.31	44.58	43.90	47.61	43.61
<i>ks</i>	ks	<b>22.50</b>	24.00	42.80	39.97	26.32	42.34	27.10	23.16
<i>L</i>	ʌ	<b>24.98</b>	32.86	32.93	29.20	28.64	29.98	29.43	26.64
<i>l</i>	l	19.95	19.34	20.63	19.85	20.17	19.65	20.96	<b>18.31</b>
<i>m</i>	m	22.90	22.82	23.74	22.56	23.46	22.67	23.96	<b>22.27</b>
<i>N</i>	ɲ	26.99	33.30	36.39	33.22	24.13	34.11	24.37	<b>21.26</b>
<i>n</i>	n	21.76	22.14	21.58	21.21	21.07	20.96	20.88	<b>19.38</b>
<i>o</i>	o	23.70	23.81	25.99	24.18	25.08	24.36	25.85	<b>22.72</b>
<i>p</i>	p	29.65	32.71	31.57	28.65	30.51	29.80	30.04	<b>28.24</b>
<i>Q</i>	θ	<b>23.08</b>	25.82	25.22	23.49	24.99	23.82	26.83	23.85
<i>r</i>	r	18.64	17.53	17.30	17.06	14.94	16.41	14.53	<b>13.85</b>
<i>s</i>	s	26.93	27.65	27.28	26.10	24.75	25.47	25.11	<b>23.49</b>
<i>t</i>	t	34.70	36.84	34.98	34.31	36.09	34.96	36.04	<b>34.07</b>
<i>u</i>	u	23.24	<b>22.63</b>	27.51	24.78	25.11	25.37	29.45	23.12
<i>v</i>	v	<b>23.87</b>	24.11	26.09	25.80	34.70	25.56	27.08	24.86
<i>w</i>	ps	<b>20.83</b>	25.71	40.93	42.47	25.92	42.98	29.62	23.66
<i>X</i>	ç	22.75	24.38	26.33	24.45	23.44	25.03	25.95	<b>21.44</b>
<i>x</i>	x	<b>20.33</b>	24.82	26.45	23.35	21.87	24.98	25.06	21.58
<i>Y</i>	ʝ	<b>26.68</b>	28.56	29.65	28.08	27.41	28.37	28.39	26.82
<i>y</i>	ɣ	20.77	20.35	22.98	21.03	21.35	21.38	21.64	<b>19.68</b>
<i>z</i>	z	23.05	22.38	23.31	22.64	23.13	22.98	24.68	<b>21.58</b>

A

**Table 4.** Mean Absolute Error (MAE), standard deviation of absolute error (STD of AE) and Root Mean Square Error (RMSE) (in milliseconds) for the fusion scheme, implemented with the average linear combination and the best-case selection fusion algorithms on: (a) the KED TIMIT database, and (b) the WCL-1 database.

(a) results on the KED TIMIT database

<i>Fusion algorithms on the KED TIMIT database</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
Overall ( <i>average</i> linear combination)	15.32	14.01	20.76
Vowel/consonant (best-case selection)	14.95	14.11	20.56
Phonetic category (best-case selection)	14.94	14.11	20.54
Phone (best-case selection)	14.92	13.87	20.37
No fusion – best individual model, SMOreg	14.95	14.11	20.56

(b) results on the WCL-1 database

<i>Fusion algorithms on the WCL-1 database</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
Overall ( <i>average</i> linear combination)	16.91	18.72	25.29
Vowel/consonant (best-case selection)	16.78	18.81	25.21
Phonetic category (best-case selection)	16.78	18.81	25.21
Phone (best-case selection)	16.68	18.65	25.08
No fusion – best individual model, SMOreg	16.78	18.81	25.21

**Table 5.** Mean Absolute Error (MAE), standard deviation of absolute error (STD of AE) and Root Mean Square Error (RMSE) (in milliseconds) for the various fusion techniques on: (a) the KED TIMIT database, and (b) the WCL-1 database.

(a) results on the KED TIMIT database

<i>KED TIMIT database</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
SMOreg	<b>14.66</b>	<b>13.82</b>	<b>20.14</b>
IBK	15.19	14.69	21.02
Linear Regression	15.49	14.45	21.18
RBFNN	15.53	14.49	21.24
m5p	15.56	14.60	21.34
Add. Regr. m5pR	15.72	14.94	21.69
Add. Regr. REPTrees	15.79	14.94	21.74
Bagging m5pR	15.81	15.09	21.86
Bagging REPTrees	15.88	15.15	21.95
m5pR	15.97	15.28	22.10
No fusion – best individual model, SMOreg	14.95	14.11	20.56

(b) results on the WCL-1 database

<i>WCL-1 database</i>	<i>MAE (ms)</i>	<i>STD of AE (ms)</i>	<i>RMSE (ms)</i>
SMOreg	<b>16.35</b>	<b>18.59</b>	<b>24.76</b>
IBK	16.98	18.85	25.47
RBFNN	17.34	19.51	26.10
Add. Regr. m5pR	17.69	19.84	26.58
Bagging m5pR	17.72	19.84	26.60
m5p	17.84	20.51	27.18
m5pR	17.91	20.00	26.85
Bagging REPTrees	17.99	20.45	27.23
Add. Regr. REPTrees	18.00	20.56	27.32
Linear Regression	18.32	20.19	27.26
No fusion – best individual model, SMOreg	16.78	18.81	25.21