# Efficient computation of the frame-based extended union model and its application in speech recognition against partial temporal corruptions

Arthur Chan [1], Manhung Siu [*]

*Department of EEE, The Hong Kong University of Science and Technology, Clearwater Bay, Hong Kong*

## Abstract

The extended union model (EUM) was recently proposed and shown to be effective in handling short time temporal corruption. Because of the computational complexity, the EUM probability can only be computed over groups of consecutive observations (called segments) and recognition can only be performed under $N$-best re-scoring paradigm. In this paper, we introduce a hidden variable called "pattern of corruption" and re-formulate the extended union model as marginalizing over possible patterns of corruption with likelihood computed via the missing feature theory. We then introduce a recursive relationship between the EUM probabilities of two successive observation sequences that can greatly simplify the EUM probability computation. This makes it possible to compute the EUM probability over a long sequence. Using this recursive relationship, the EUM probability over frames, called the "frame-based EUM" can easily be computed. To simplify the EUM-based recognition, we propose an approximated, dynamic programming-based EUM recognition algorithm, called the Frame-based EUM Viterbi algorithm (FEVA), that performs recognition directly instead of via $N$-best re-scoring. Experimental results on digit recognition under added impulsive noises show that both the frame-base EUM and the FEVA outperform the segment-based EUM.
© 2004 Elsevier Ltd. All rights reserved.

---

[*] Corresponding author. Tel.: +852 2358 8027; fax: +852 2358 1485.
  *E-mail addresses:* archan@cs.cmu.edu (A. Chan), msiu@ieee.org (M. Siu).
[1] Currently with CMU.

## 1. Introduction

As speech recognition matures and is being fielded in real world environments, robustness against unknown noisy environments becomes vital. Some real-life noises, such as door slams or telephone rings, are short and impulsive in nature. Their frequencies of occurrence and durations are difficult to predict.

Because these corrupted observations do not match well with the clean speech models, they would bias the observation likelihood and cause recognition errors. The problem of handling corrupted observations is recently addressed by missing feature theory (Lippmann and Carlson, 1997), which proposed to omit the corrupted portion of the observation sequence from the likelihood computation. Other authors suggested that the corrupted components can be estimated from the clean components either in the spectral domain (Cooke et al., 1997) or even in the general spectral–temporal domain (Raj et al., 1998). It is well known that prior knowledge such as location of corruptions and the number of corruptions can greatly improve the robustness of the algorithms (de Veth et al., 1999).

The probabilistic union model, first introduced by Ming and Smith (1999), is a general approach to reduce the model's sensitivity to unknown corruptions by using the probabilistic union operation. It was found to be very useful in dealing with partial corruptions. It was adapted to handle corruptions in sub-band-level (Ming et al., 2002) or word-level (Sicilia-Garcia et al., 2002). In Ming and Smith (2001), it was generalized into the extended union model (EUM) to make it applicable to robust speech recognition against *partial temporal corruptions*, or corruptions which affect only a subset of the observations This class of oise includes door-slam, keyboard click sounds and packet loss in digital transmission network.

The EUM is a probabilistic union of all possible observation subsets. For a $T$-frame observation sequence in which $M$ observations are corrupted by short time noise at unknown positions, there are $C_M^T$ possible subsets of $T–M$ observations. One of these subsets contains no corrupted observation. Ideally, if the identities of the corrupted observations are known and the joint observation probability of the observation sequence is computed by excluding the corrupted observation as suggested in Lippmann and Carlson (1997), then, this joint observation probability would be the same as the probability of the clean subset. In reality, the identities of the corrupted observations are often not known. The EUM probability of the $T$-frame observation sequence is defined as the probability of the union of all the $C_M^T$ possible subsets and can be computed as a sum over the probabilities of each subset. By assuming that the probability of the corrupted observations would be small compared to the clean observations, this sum would be dominated by the probability of the clean subset making the EUM probability a close approximation of ideal probability.

Using the EUM probability as a recognition criterion would require solving two computationally intensive problems. First, one needs to compute the EUM probability of a given state sequence. As will be discussed in Section 2, because the EUM is combinatorial in nature, direct computation of the EUM probability of a long observation sequence (large $T$) is infeasible computationally. Second, in speech recognition, one needs to evaluate the EUM probabilities over all possible state sequences.

Two approximations were proposed in Ming and Smith (2001) to handle the computational issues. First, the observation sequence can be partitioned into a sequence of fixed-length segments [2]. Since the number of segments would be substantially smaller than the number of frames, direct computation of the EUM probability would be feasible. Second, recognition can be performed under the *N*-best re-scoring paradigm (Schwartz and Chow, 1990) to reduce the number of state sequences to be evaluated. By using these approximations, recognition using the EUM probability was reported to perform better than using the joint observation likelihood when the test utterances were corrupted by short-time impulsive noises. The segment length holds the balance between performance and computation. Using a long segment reduces computation. But, it could either ignore short corruptions or mark a segment corrupted even though only few observations (frames) in the segment are corrupted. Using short segments would increase modeling resolution but is computationally expensive. The choice of the segment length have to be optimized for different noise conditions empirically (Section 4.2; Ming and Smith, 2001) which imposes some constraints on the general applicability of the approach. Long segments are typically used to reduce computation. To alleviate the side effect of long segments, (Ming, 2001) suggested shifting the segment boundaries (the first segment would not always begin at the first observation) to minimize the number of segments affected by the corruptions.

Single-observation segments in the EUM (or the frame-based EUM) would be ideal if one can solve the computational problem. In this paper, we propose an efficient algorithm to compute the EUM probability that makes the frame-based EUM feasible. This algorithm centers on a recursive relationship between the EUM probability of a $t + 1$-frame sequence and the EUM probability of the first $t$-frame sub-sequence. By means of direct computation, computation complexity in obtaining the EUM probability would be exponential on the sequence length $T$. By using this recursion, the complexity is reduced to a linear function of $T$.

To take full advantage of the EUM framework, it is desirable to use the EUM probability directly in Viterbi recognition instead of under the *N*-best re-scoring paradigm. In this paper, we derive another recursive relationship between the EUM probabilities across different state sequences and devise an approximate dynamic-programming-based algorithm called the frame-based extended union model Viterbi algorithm (FEVA), which is an approximation to the full EUM search. The FEVA provides an alternative solution to the two-stage recognition scheme described in Ming and Smith (2001).

This paper is organized as follows. In Section 2, we review the EUM formulation and express it as a marginalization over different patterns of corruption. We also describe how the EUM can be used in speech recognition. In Section 3, we describe the first recursive relation mentioned above and discuss how the computation of the EUM probability can be simplified. The FEVA will then be described in Section 4. In Section 5, we evaluate the proposed algorithms under partial temporal corruptions. We summarize and conclude the paper in Section 6.

---

[2] Here, the term "segment" is not a representation of an acoustic unit as in segmental modeling. It is simply a sequence of consecutive frames. The usage of term follows (Ming and Smith, 2001).

## 2. The extended union model for speech recognition

Denote a sequence of $T$ observations as $O_1^T = (o_1, \ldots, o_T)$. The conditional likelihood of this sequence for a given state sequence $Q_1^T$ can be expressed as

$$p(O_1^T \mid Q_1^T) = p(o_1, \ldots, o_T \mid Q_1^T)$$
$$= \prod_{t=1}^{T} p(o_t \mid q_t).$$

Suppose out of the $T$ frames, $M$ are corrupted at unknown locations. Denote $l_i = \{t_{i,1}, \ldots, t_{i,M}\}$ as a possible set of $M$ observation time indexes such that $o_{t_{i,k}}$ is corrupted. For example, to mark $o_2$, $o_4$, $o_7$ as corrupted, $l = \{2,4,7\}$. By rules in combinatorics, there are $c_M^T$ possible $l_i$, i.e. $1 \leqslant i \leqslant c_M^T$. We call the $l_i$'s the **patterns of corruption**. If the pattern of corruption is known, then according to the missing feature theory (Lippmann and Carlson, 1997), the corrupted observations should be ignored in the calculation of the joint likelihood. That is,

$$p(O_1^T \mid Q_1^T, l_i) \doteq \prod_{t=1, t \notin l_i}^{T} p(o_t \mid q_t).$$

This joint probability of only the known clean frames can be used as the optimization function in recognition.

In reality, $l_i$ is not known. Instead, *all possible* patterns of corruption can be considered by using their **union**. Denote the union of the observations across all possible corruption patterns as $O_{\cup(T)}^{\circ M}$ [3] in which the $M$ denotes the number of corrupted frames and $T$ denotes the total number of frames in the sequence

$$O_{\cup(T)}^{\circ M} = O_1^T \wedge (l_1 \cup \cdots \cup l_{c_M^T})$$
$$= (O_1^T \wedge l_1) \cup \cdots \cup (O_1^T \wedge l_{c_M^T}),$$

in which $\wedge$ denotes the "and" operator (or joint) and $\cup$ denotes the "or" or union operator. Because $O_1^T \wedge l_1$ represents the observations in which a particular subset is corrupted, $O_{\cup(T)}^{\circ M}$ is the union of all possible subsets. This is called an extended union model (EUM) of order $M$, which describes the union event of all possible patterns with $M$ corruptions.

Because the patterns of corruption are mutually exclusive, the $p(O_{\cup(T)}^{\circ M})$ can be expressed as a sum of the probability of the observations and the individual patterns. That is,

---

[3] Strictly speaking, because the EUM probability is calculated from frame 1 to frame $t$, a more vigorous notation would be $(O_1^t)_{\cup(t)}^{\circ M}$. Since we typically assumes that the beginning frame is 1 and that the ending index is always the same as the total number of frames considered, we adopt a simpler notation $O_{\cup(t)}^{\circ M}$ in this paper.

$$
\begin{aligned}
p(O^{\circ M}_{\cup(T)} \mid Q_1^T) &= p(O_1^T \wedge (l_1 \cup \ldots l_i \ldots \cup l_{c_M^T}) \mid Q_1^T) \\
&= \sum_{i=1}^{c_M^T} p(O_1^T, l_i \mid Q_1^T) \\
&= \sum_{i=1}^{c_M^T} p(O_1^T \mid l_i, Q_1^T) p(l_i \mid Q_1^T) \\
&= \sum_{i=1}^{c_M^T} p(O_1^T \mid l_i, Q_1^T) p(l_i) \\
&= \frac{1}{c_M^T} \sum_{i=1}^{c_M^T} p(O_1^T \mid l_i, Q_1^T).
\end{aligned}
\tag{1}
$$

The above formulation assumes that the pattern of corruption is not dependent on the state sequence which is equivalent to assuming that the noise is independent of the speech signal. Furthermore, with no prior knowledge of the noise, the prior probability of the pattern of corruption, $p(l_i)$, is assumed to be uniform.

Why would the EUM probability be robust to corruptions? One can view this from two perspectives. First, corrupted observations do not match the model well and would result in lower likelihood than the clean observations. This implies that $p(O^{\circ M}_{\cup(T)} \mid Q_1^T)$ would be dominated by the term $p(O_1^T \mid l_i, Q_1^T)$ which contains well matched clean observations. Then, the effect of the corruptions is reduced in the decision process. The idea is closely related to robust statistics: if we regard the process of likelihood computation as gathering statistics from individual frames, Eq. (1) can be interpreted as jackknife estimates (Duda et al., 2001) (similar to cross validation or leavel-one-out) of likelihood that leaves $M$ observations out.

The second perspective is suggested in Ming and Smith (2001). They observed that the likelihood function, which involved an "and" operation on the observations which is good for determining the correctness of all the frames against a model but is sensitive to corruption of individual frames. A probabilistic general union model, on the other hand, involves an "or" operation on observations that would average out the observation probabilities of corrupted frames and clean frames. The EUM works as a trade-off between the above two cases. As $M$ becomes smaller, the EUM approaches the joint likelihood function. As $M$ becomes larger, the EUM approaches a general union model. By properly adjusting the value of $M$, the EUM should be able to cope with any pattern of corruption.

As we can see in Eq. (1), the computation of $p(O^{\circ M}_{\cup(T)} \mid Q_1^T)$ involves $c_M^T$ terms. Table 1 shows the EUM of order $0, 1, 2, 3$ for a four-frame sequence. (We ignore the state information in the table and only show the observation information there.) Because of the number of terms depends on $c_M^T$, we can see that direct implementation can be computationally intractable for large $T$ ($T > 100$) which is the typical length of a speech utterance. In Section 3, we describe how this can be solved more efficiently.

In a HMM-based speech recognition system, instead of considering only one state sequence, we need to find the most likely state sequence, $\hat{Q}_1^T$, given the observations $O_1^T$, i.e.,

Table 1
Example of the EUM probability of different orders in the case for 4 frames, order $0, 1, 2, 3$ represent $0, 1, 2, 3$ frames are missing

| $M$ | $p(O^{\circ M}_{\cup(T)})$ | Number of terms |
|---|---|---|
| 0 | $p(o_1 o_2 o_3 o_4)$ | $c^{T=4}_{M=0} = 1$ |
| 1 | $p(o_1 o_2 o_3 \cup o_1 o_2 o_4 \cup o_1 o_3 o_4 \cup o_2 o_3 o_4)$ | $c^{T=4}_{M=1} = 4$ |
| 2 | $p(o_1 o_2 \cup o_1 o_3 \cup o_1 o_4 \cup o_2 o_3 \cup o_2 o_4 \cup o_3 o_4)$ | $c^{T=4}_{M=2} = 6$ |
| 3 | $p(o_1 \cup o_2 \cup o_3 \cup o_4)$ | $c^{T=4}_{M=3} = 4$ |

Note that the state information was ignored here.

$$\hat{Q}^T_1 = \arg \max_{Q^T_1} p(Q^T_1 \mid O^T_1). \tag{2}$$

By using Bayes rule and using the fact that $p(O^T_1)$ is independent of the state sequence, Eq. (2) can be re-written as maximizing the joint observation and state likelihood, i.e.,

$$\begin{aligned}
\hat{Q}^T_1 &= \arg \max_{Q^T_1} p(O^T_1, Q^T_1) \\
&= \arg \max_{Q^T_1} p(O^T_1 \mid Q^T_1) p(Q^T_1).
\end{aligned} \tag{3}$$

Under the EUM framework, the observation likelihood is replaced by the EUM probability of order $M$ that is robust to $M$ corruptions, i.e.,

$$\hat{Q}^T_1 = \arg \max_{Q^T_1} p(O^{\circ M}_{\cup(T)} \mid Q^T_1) p(Q^T_1). \tag{4}$$

Eq. (4) poses another computation difficulty: The need to search for the state sequence with highest EUM probability. While the Viterbi algorithm can efficiently find the best state sequence to solve Eq. (3), it is not clear how to search for the best state sequence that optimized the EUM probability in Eq. (4).

   To reduce the number of state sequences to be searched, (Ming and Smith, 2001) suggested applying the $N$-best re-scoring paradigm. This simplifies the problem of evaluating the EUM probabilities of all state sequences into evaluating $N$ hypotheses only at the expense of a more complex decoding process and a dependence on the quality of the $N$-best list. Furthermore, to reduce $T$ in the EUM, the observations sequence was partitioned into fixed-length segments of length $D$, denoted as $(z_1, \ldots, z_{\lceil T/D \rceil})$. Because the number of segments, $\lceil T/D \rceil$ is significantly smaller than $T$, the best state sequence $\hat{Q}^T_1$, which optimized the segment-based extended union probability with $M$ segments being corrupted can be found. The segment based EUM recognition equation is given by

$$\hat{Q}^T_1 = \arg \max_{Q^T_1 \in (N\text{-best})} p(Z^{\circ M}_{\cup(\lceil T/D \rceil)}, Q^T_1). \tag{5}$$

Partitioning the observations to segments may have some drawbacks. If part of the segment was corrupted, the whole segment could either be regarded as corrupted or completely clean. This could result in either excessive data loss or reduced robustness. Another drawback is that setting

the segment length can be difficult. To reduce the side effect of long segments, Ming (2001) suggested the use of a "variable" segmentation that shifts the beginning of the segments while keeping the segment length fixed. By selecting the segment origin that minimizing the number of segments affected by corruption, this reduces the chance that multiple long segments to be corrupted and reduces the order of the EUM.

In general, the EUM over individual frames is desirable but is hindered by its computation complexity. We will discuss how this problem can be solved in the next section.

## 3. Fast computation of the EUM likelihood

As we have mentioned in the previous section, computing the EUM probability of a given state sequence using Eq. (1) directly is computationally expensive for any moderate $T$. However, many of the computations are repeated. Similar to the dynamic programming approach applied in Viterbi and forward–backward algorithm, a recursive relationship can be established between the EUM probability of $t$ observations and $t - 1$ observations. By creating such a recursion, the computation complexity can be reduced from exponential to linear order.

Using the definition of the EUM in Section 2 Eq. (1), the EUM probabilities of a $t$-frame sequence with $m$ corruptions given a state sequence $Q_1^t$ can be expressed as

$$p(O_{\cup(t)}^{\circ m} \mid Q_1^t) = \frac{1}{c_m^t} \sum_{i=1}^{c_m^t} p(O_1^t \mid l_i, Q_1^t), \tag{6}$$

where $l_i$ is the $i$-th pattern of corruption. Within the $c_m^t$ patterns of corruption, some patterns mark observation $o_t$ (the last observation) as corrupted and some do not. Thus, we can divide the $c_m^t$ patterns into two subsets, $L_1$ and $L_2$. The first set, $L_1$, includes only those patterns with $o_t$ marked as corrupted. The second set, $L_2$, is the complement of $L_1$.

Because the sets $L_1$ and $L_2$ are mutually exclusive, Eq. (6) can be rewritten as

$$p(O_{\cup(t)}^{\circ m} \mid Q_1^t) = \frac{1}{c_m^t} \left( \sum_{l_i \in L_1} p(O_1^t \mid l_i, Q_1^t) + \sum_{l_i \in L_2} p(O_1^t \mid l_i, Q_1^t) \right). \tag{7}$$

For $L_1$, since $o_t$ is marked as corrupted (and ignored in likelihood computation), there are only $m - 1$ corruptions out of the previous $t - 1$ observation. Thus,

$$\sum_{l_i \in L_1} p(O_1^t \mid l_i, Q_1^t) = \sum_{i=1}^{c_{m-1}^{t-1}} p(O_1^{t-1} \mid l_i, Q_1^t)$$
$$= c_{m-1}^{t-1} p(O_{\cup(t-1)}^{\circ m-1} \mid Q_1^{t-1}). \tag{8}$$

Notice that the conditioning on $Q_1^t$ is changed to $Q_1^{t-1}$ because the knowledge of $q_t$ is not needed when $o_t$ is marked as corrupted and being ignored in likelihood computation.

The corruption patterns in $L_2$ do not mark $o_t$ as corrupted. This means that there are $m$ corruptions out of the first $t - 1$ observations. Also, observation $o_t$ is included in the likelihood computation. Thus,

$$\sum_{l_i \in L_2} p(O_1^t \mid l_i, Q_1^t) = \sum_{i=1}^{c_m^{t-1}} p(o_t \mid l_i, Q_1^t) p(O_1^{t-1} \mid l_i, Q_1^t) \tag{9}$$

$$= p(o_t \mid Q_1^t) \sum_{i=1}^{c_m^{t-1}} p(O_1^{t-1} \mid l_i, Q_1^t) \tag{10}$$

$$= p(o_t \mid q_t) \sum_{i=1}^{c_m^{t-1}} p(O_1^{t-1} \mid l_i, Q_1^{t-1}) \tag{11}$$

$$= c_m^{t-1} p(o_t \mid q_t) p(O_{\cup(t-1)}^{\circ m} \mid Q_1^{t-1}). \tag{12}$$

Because of the conditional independence assumption, the likelihood of the current observation $o_t$ can be decomposed as expressed in Eq. (9). By definition, all the corruption patterns in $L_2$ do not include $t$, implying that the likelihood of $o_t$ is always included in the sum and is independent of the corruption patterns. Thus, Eq. (9) can be rewritten into Eq. (10). By using the conditional independence assumption again, $p(o_t \mid Q_1^t)$ is simplified into $p(o_t \mid q_t)$ and $p(O_1^{t-1} \mid l_i, Q_1^t)$ into $p(O_1^{t-1} \mid l_i, Q_1^{t-1})$ in Eq. (11). Finally, using Eq. (6), the sum in Eq. (11) is converted into the EUM probability at time $t - 1$ in Eq. (12). Combining Eqs. (7)–(12) and using the fact that

$$\frac{c_{m-1}^{t-1}}{c_m^t} = \frac{m}{t}$$

and

$$\frac{c_m^{t-1}}{c_m^t} = \frac{t-m}{t},$$

we have the following recursion:

$$
\begin{aligned}
p(O_{\cup(t)}^{\circ m} \mid Q_1^t) &= \left(\frac{c_{m-1}^{t-1}}{c_m^t}\right) p(O_{\cup(t-1)}^{\circ m-1} \mid Q_1^{t-1}) + \left(\frac{c_m^{t-1}}{c_m^t}\right) p(O_{\cup(t-1)}^{\circ m} \mid Q_1^{t-1}) p(o_t \mid q_t) \\
&= \left(\frac{m}{t}\right) p(O_{\cup(t-1)}^{\circ m-1} \mid Q_1^{t-1}) + \left(\frac{t-m}{t}\right) p(O_{\cup(t-1)}^{\circ m} \mid Q_1^{t-1}) p(o_t \mid q_t).
\end{aligned}
\tag{13}
$$

We call Eq. (13) the EUM formula. It has two recursive variables, $m$ and $t$. To compute $p(O_{\cup(t)}^{\circ m} \mid Q_1^t)$, one would need $p(O_{\cup(t-1)}^{\circ m-1} \mid Q_1^{t-1})$ at time $t - 1$ which in turn depends on $p(O_{\cup(t-2)}^{\circ m-1} \mid Q_1^{t-2})$.

In addition to a recursive formulation on the condition observation likelihood $p(O_{\cup(t)}^{\circ m} \mid Q_1^t)$, it is sometimes useful to compute the unconditional observational likelihood $p(O_{\cup(t)}^{\circ m})$. Such a recursive relationship can easily be derived using Eq. (13) and the fact that

$$p(O_{\cup(t)}^{\circ m}, Q_1^t) = p(O_{\cup(t)}^{\circ m} \mid Q_1^t) p(Q_1^t).$$

A summary of the recursion equations on both conditional and unconditional observation likelihood is given in Fig. 1.

In general, using the EUM formula for computing $p(O_{\cup(T)}^{\circ M} \mid Q_1^T)$, we need to compute $p(O_{\cup(t)}^{\circ m} \mid Q_1^t)$ for time $t \leqslant T$ and all order $m \leqslant M$. Each update of $p(O_{\cup(t)}^{\circ m} \mid Q_1^t)$ requires three mul-

Denote the number of states as $N$ and

$$\kappa_t(m) = p(\mathring{O}^m_{\cup(t)}|Q^t_1),$$

$$\hat{\alpha}_t(j,m) = p(\mathring{O}^m_{\cup(t)}, Q^{t-1}_1, q_t = j).$$

**Initialization:**

$$\kappa_1(0) \quad = b_{q_1}(o_1)$$

$$\kappa_1(1) \quad = 1$$

$$\hat{\alpha}_1(j,0) \; = b_j(o_1) \qquad\qquad\qquad\qquad\qquad\qquad 0 \le j \le N$$

$$\hat{\alpha}_1(j,1) \; = 1$$

$$\hat{\alpha}_t(j,m) = \kappa_t(m) = 0 \qquad\qquad\qquad\qquad 0 \le j \le N, 1 \le t \le T, m > \min(M,t)$$

**Recursion:** $(1 \le j \le N)$

$$\kappa_t(m) \quad = (\tfrac{m}{t})\kappa_{t-1}(m-1) + (1 - (\tfrac{m}{t}))\kappa_{t-1}(m)b_{q_t}(o_t) \qquad 0 \le m \le \min(M,t), 1 < t \le T$$

$$\hat{\alpha}_t(j,m) = \sum_i a_{i,j} \left[ (\tfrac{m}{t})\hat{\alpha}_t(i,m-1) + (1 - (\tfrac{m}{t}))\hat{\alpha}_t(i,m)b_j(o_t) \right] \; 0 \le m \le \min(M,t), 1 < t \le T$$

**Termination:**

$$p(\mathring{O}^M_{\cup(T)}) = \sum_j \hat{\alpha}_T(j,M)$$

Fig. 1. Summary of recursive equations for EUM probability computation.

tiplications and one addition. For a given state sequence, there is only one state to evaluate at each time $t$. For each state, the number of EUM scores to compute is $M$. Therefore, the computation, is of O(MT). Consider direct computation using Eq. (1) instead of using the EUM formula. The number of possible unions is $C^M_T$ and each has to be evaluated separately resulting in O($C^M_T$).

As we can see in Fig. 1, to compute of the unconditional EUM probability at time $t$ for state $j$ for EUM order $m$, a sum of $N$ scores would be needed. Therefore, the computation in this case for all state and all time would be O($MTN^2$) which is $M$ times more in computation as in the computing the forward probability in HMM. Direct computation of this EUM probability would be too expensive.

The same principle as in Eq. (13) can be applied to the segment-based EUM as described in Ming and Smith (2001) and can potentially be useful for other EUM applications as well.

## 4. The frame-based EUM viterbi algorithm

The proposed recursion in Section 3 simplifies the computation of the EUM probability of a given state sequence and thus removes the constraint on segment length. This makes the frame-based EUM possible. However, in order to use it for speech recognition, $N$-best re-scoring paradigm that reduces the number of state sequences is still needed. While the $N$-best rescoring paradigm has been used extensively for testing out new models, its usefulness is dependent on the quality of the $N$-best. If a long $N$-best list is needed, it can be expensive to rescore. Insights from other works such as (Siu and Chan, 2002) on robust speech recognition suggested that it is possible to devise a dynamic algorithm which directly optimizes the EUM probability.

For an HMM-based recognition system, efficient solution for Eq. (2) can be obtained using the Viterbi algorithm by computing the partial path score $\delta_t(j)$, which is the best likelihood of $O_1^t$ and state sequence up to time $t - 1$ such that $q_t = j$. Mathematically,

$$\delta_t(j) = \max_{Q_1^{t-1}} p(O_1^t, Q_1^{t-1}, q_t = j). \tag{14}$$

A recursive relation on $\delta_t(j)$ in the Viterbi algorithm is given by

$$\delta_t(j) = \max_i \left[\delta_{t-1}(i)a_{ij}\right] b_j(o_t), \quad 1 \leqslant t \leq T, \tag{15}$$

where $b_j(o_t)$ is the observation probability of $o_t$ at state $j$.

In Eq. (13), we have a recursive relationship across time for a given state sequence. In order to make the EUM-based recognition possible, a recursion across different state sequences similar to Eq. (15) would be needed.

### 4.1. Computing the EUM probability

We define the partial path score $\delta_t(j,m)$ as the EUM likelihood at time $t$ for state $j$ with the best state sequence up to time $t - 1$ for an order $m$ EUM. Mathematically, $\delta_t(j,m)$ can be expressed as

$$\delta_t(j,m) = \max_{Q_1^{t-1}} p(O_{\cup(t)}^{\circ m}, Q_1^{t-1}, q_t = j). \tag{16}$$

If we denote $v_m(t) = m/t$, then, an approximated algorithm is derived as follows:

$$\delta_t(j,m) = \max_{Q_1^{t-1}} p(O_{\cup(t)}^{\circ m}, Q_1^{t-1}, q_t = j) = \max_i \left[\max_{Q_1^{t-2}} \left(p(O_{\cup(t)}^{\circ m}, Q_1^{t-2}, q_{t-1} = i, q_t = j)\right)\right]$$

$$= \max_i \left[\max_{Q_1^{t-2}} \left(p(O_{\cup(t)}^{\circ m} \mid Q_1^{t-2}, q_{t-1} = i, q_t = j) p(Q_1^{t-2}, q_{t-1} = i, q_t = j)\right)\right]$$

$$= \max_i a_{ij} \left[\max_{Q_1^{t-2}} \left(p(Q_1^{t-2}, q_{t-1} = i) p(O_{\cup(t)}^{\circ m} \mid Q_1^{t-2}, q_{t-1} = i, q_t = j)\right)\right], \tag{17}$$

$$\delta_t(j,m) = \max_i a_{ij}\left[\max_{Q_1^{t-2}}\left(p(Q_1^{t-2}, q_{t-1}=i)\left\{v_m(t)p(O_{\cup(t-1)}^{\circ\, m-1} \mid Q_1^{t-2}, q_{t-1}=i)\right.\right.\right.$$
$$\left.\left.\left.+(1-v_m(t))p(O_{\cup(t-1)}^{\circ\, m} \mid Q_1^{t-2}, q_{t-1}=i)b_j(o_t)\right\}\right)\right], \tag{18}$$

$$\delta_t(j,m) = \max_i a_{ij}\left[\max_{Q_1^{t-2}}\left(v_m(t)p(O_{\cup(t-1)}^{\circ\, m-1}, Q_1^{t-2}, q_{t-1}=i)\right.\right.$$
$$\left.\left.+(1-v_m(t))p(O_{\cup(t-1)}^{\circ\, m}, Q_1^{t-2}, q_{t-1}=i)b_j(o_t)\right)\right], \tag{19}$$

$$\delta_t(j,m) \approx \max_i a_{ij}\left[v_m(t)\max_{Q_1^{t-2}}p(O_{\cup(t-1)}^{\circ\, m-1}, Q_1^{t-2}, q_{t-1}=i)\right.$$
$$\left.+(1-v_m(t))\max_{Q_1^{t-2}}p(O_{\cup(t-1)}^{\circ\, m}, Q_1^{t-2}, q_{t-1}=i)b_j(o_t)\right], \tag{20}$$

$$\delta_t(j,m) = \max_i a_{ij}\left[\left(v_m(t)\delta_{t-1}(i,m-1) + (1-v_m(t))\delta_{t-1}(i,m)b_j(o_t)\right)\right]. \tag{21}$$

By means of Bayes rule and the Markov assumption, we can rewrite Eq. (16) into Eq. (17) such that $\delta_t(j,m)$ depends on quantities at time $t-1$. Eq. (18) is obtained by applying the recursion from Eq. (13) such that the EUM probability at time $t$ can be expressed as a sum of the EUM probabilities at time $t-1$. By applying the Bayes rule Eq. (18) can be rewritten into Eq. (19) which is a maximization over a sum of two terms at time $t-1$. Notice that $\delta_t(j,m)$ is now expressed as the score of a **single** state sequence $Q_1^{t-2}$ that maximizes the sum of the two terms. However, each term by itself is not necessarily the best score at time $t-1$ at state $i$. This implies that to compute $\delta_t(j,m)$ exactly, one has to store the partial path scores for all state sequence $Q_1^{t-2}$ rendering the recursion useless. The maximum of the sum, however, can be approximated by the sum of the maximum in Eq. (20). This in effect is an optimistic estimation of the EUM path score. Then, a recursion on $\delta_t(j,m)$ can be obtained in Eq. (21).

There are two recursion variables $m$ and $i$ in Eq. (21). $\delta_t(j,m)$ depends on the $\delta$'s at $t-1$ from different states $i$ of orders $m$ and $m-1$. One can view the above recursion as expanding the path metric $\delta_t(j)$ into a $M+1$ dimensional vector. Hence, the recursion can be simply implemented by expanding the search space $M+1$ times and applying Eq. (21). We called this algorithm the frame-base EUM Viterbi algorithm (FEVA).

## 4.2. Back-tracking

In standard Viterbi algorithm, determination of the best state sequence requires only state information in the previous time $t-1$. In the FEVA, back-tracking requires both the best previous state and the knowledge of the number of previous order $M_{t-1}$ because of the two recursion variables in Eq. (21). Define $\psi_t(j,m)$ as the best previous state for an order $m$ EUM at time $t$ and state $j$.

$$\psi_t(j,m) = \arg\max_i a_{ij}[(1 - v_m(t))\delta_{t-1}(i,m)b_j(o_t) + v_m(t)\delta_{t-1}(i,m-1)], \tag{22}$$

where $v_m(t) = m/t$. However, because of the approximation in Eqs. (20) and Eq. (21), determination of the exact EUM order information is very difficult. This is because the $m$th order partial score at time $t$ would depends on the sum of the $m$th and $m - 1$th order EUM scores at time $t - 1$ that can have different state histories. This situation is analogous to the difference between the forward–backward algorithm and the Viterbi algorithm. To search for an approximate best EUM state sequence, we propose the following approximation:

$$M_t(m) = \begin{cases} m & \text{if } (1 - v_m(t))\delta_{t-1}(\hat{i},m)b_j(o_t) > v_m(t)\delta_{t-1}(\hat{i},m-1), \\ m-1 & \text{otherwise,} \end{cases} \tag{23}$$

where $\hat{i} = \psi_t(j,m)$. This approximation would mean that the score computed is not that of the recognized path. This type of approach, however is also used in other recognition system. For example, (Lamere et al., 2003) used the forward probability as the partial path metric while the most likely state was used in back-tracking.

### 4.3. Implementation of the frame-based EUM viterbi algorithm

Similar to the $N$-best recognition algorithm proposed in Schwartz and Chow (1990), one can directly implement Eqs. (21)–(23) by augmenting the state to store $M + 1$ scores at each time instance.

To illustrate the implementation of FEVA, Fig. 2 shows a simple left-to-right HMM with two states, state $a$ and state $b$ and Fig. 3 shows the computation of the path score and the augmented state space for FEVA for $M = 2$. Each box represents an augmented state that stores the partial path score $\delta_t(j,m)$, $0 \leqslant m \leqslant M$. The original Viterbi path which is the same path as the EUM order 0 is stored as $m = 0$. During the update (from time $t$ to $t + 1$), the scores from all possible incoming arcs are calculated using Eq. (21). For example, if we ignore the transition probabilities, for state $b$ at time $t + 1$ at $m = 1$, the score is computed as

$$\delta_{t+1}(b,1) = \max_{j=a,b} \left( \frac{1}{t}\delta_t(j,0) + \frac{t-1}{t}\delta_t(j,1)p(o_t \mid q_j) \right).$$
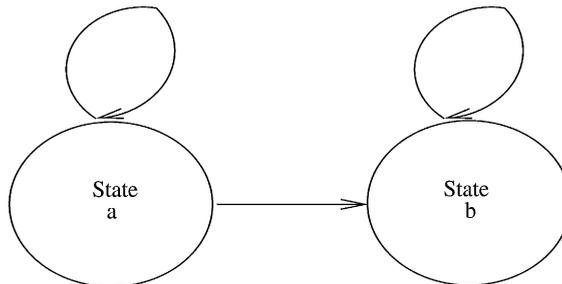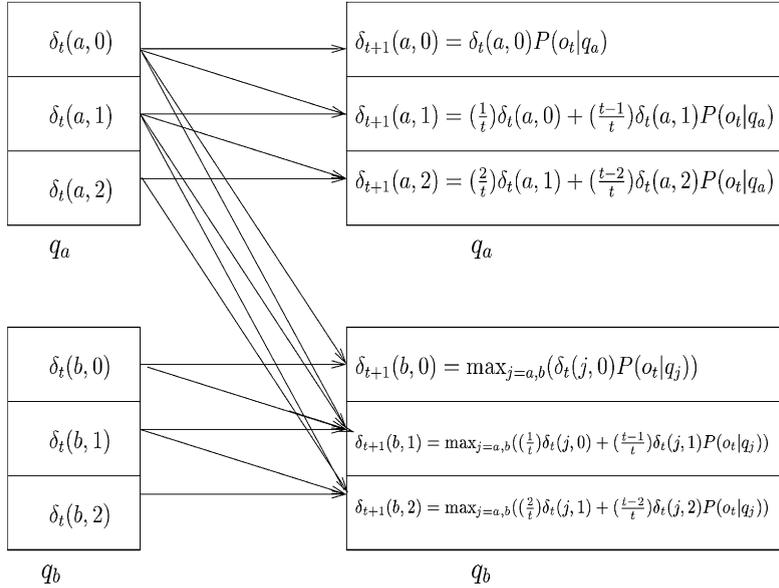


Fig. 2. Topology of a simple HMM.

Fig. 3. FEVA path score update for the topology shown in Fig. 2 with EUM order $M = 2$. In here, $\delta_t(j, m)$ is the EUM score at time $t$ for state $j$ and order $m$.

It is possible that Eq. (21) requires computation of union probabilities with either $m < 0$ or $t < 0$ or $m > t$. In all these cases, $p(O^{\circ m}_{\cup(t)})$ is set to 0.

Backtracking pointers can be used to record the state with maximum partial scores for each state as described in Eq. (23). Because the dynamic range between $v_m(t)$ and the EUM probability of the previous frame can be large, a simplified implementation would ignored the contribution of $v_m(t)$.

The computation and memory size of this implementation would be $M$ times larger than a standard Viterbi algorithm. With the use of caching on the observation likelihood, the computation could be reduced to a closer level to the Viterbi algorithm.

Both the $N$-best re-scoring and the FEVA are approximations to the "exact" search. In $N$-best, the search space is reduced first by using a weaker model. This can be particularly problematic under noisy environment because the impact of the noise can significantly degrade the quality of the $N$-best. In FEVA, while the likelihood is approximated, the more robust search is applied to the whole space considering all possible state sequences. For the $N$-best re-scoring to work well, the quality of the $N$-best is one key which in turn depends on the length of the $N$-best which has implication on computation complexity.

## 5. Experimental results

In this section, we describe the experimental results of both the frame-based EUM and the FEVA. Similar to (Ming and Smith, 2001), experiments were conducted to evaluate recognition performance under partial temporal corruptions (PTC). A description of the corpus and

experimental setup is given in Section 5.1. Using *N*-best re-scoring paradigm, we compare recognition performance of the segment-based EUM, the frame-based EUM and the FEVA under wide ranges of noisy conditions in Section 5.2.

### 5.1. Corpus and conditions

Before we describe the experimental conditions in detail, three terms, "frame", "segment" and "block" are used in this section that can be confusing. "Frame" denotes an acoustic feature vector (MFCC). "Segment" denotes a sequence of consecutive frames under the EUM and is used in the context of the segment-based EUM. Finally, "block" denotes a short interval of noise that is added to speech to simulate the noisy environment. By using the above definitions, we hope to avoid ambiguity.

The experiments were carried out using the TI-digit (Leonard, 1984) on the task of connected digit recognition. 8668 training utterances and 8668 testing utterances were used. Each utterance is a digit string with a maximum length of 7 digits and up to 3 s long. Speech observations (MFCC frames) were generated at a rate of 100 frames per second using a 25 ms window. Thirty-nine dimensional front-end features were used including 12 MFCC, log frame energy, and their first and second-order derivatives.

Each digit was represented by an 18-state left-to-right HMM. Silence was represented as a 5-state left-to-right HMM. Short pause was modeled using a 3-state *t*-model as described in Young et al. (1999). We followed the training recipe described in Hirsch and Pearce (2000) to estimate the parameters of the HMMs and obtained baseline word accuracy of 98.83% in connected digit recognition. [4]

We tested the algorithm in two conditions: clean and noise-corrupted. The original testing utterances in the TIDIGITS corpus were used in the clean condition. For the noise-corrupted cases, partially temporal corruptions (PTC) discussed in Ming and Smith (2001) were added. In our experiments, noise blocks were extracted from four telephone ring tones, one whistling noise and two message sound clips in ICQ. These noise blocks were randomly added to the testing utterances. To test the effectiveness of the proposed algorithms, we varied the following attributes of the noise corruptions.

- *Signal-to-noise ratio* (*S*) is computed as the ratio between the power of the speech signals and the noise added. In our experiments, the signal-to-noise ratio was computed based the following equation:

$$\text{SNR} = 10 \log_{10}\left(\frac{\text{power of the speech signal}}{\text{power of the added noise}}\right), \tag{24}$$

  with silences and short-pauses excluded in the calculation of the power of the speech signal.
- *Rate of corruption* (*C*) is the probability that a block of speech would be corrupted by additive noise.

---

[4] The results was un-tuned, by tuning different penalties, higher accuracy above 99.05% can be reached.

- *Block length* (*L*) is the duration of the block measured in milli-seconds (ms). In our experiments, block lengths are always multiples of 25 ms.

In our experiments, each speech utterance was partitioned into non-overlapping blocks. For each block, a Bernoulli random variable with $p = C$ was generated to decide whether noise should be added. This simulates the general scenario in which short-time noise occurs in a time-varying and random manner. Because the noise is added in time domain, the use of overlapping window in cepstral generation and derivative computation results in corruption of multiple frames.

## 5.2. Digit recognition under partial temporal corruptions

We first tested the frame-based EUM in the *N*-best re-scoring framework. Fifty *N*-best hypotheses were generated using HTK (Young et al., 1999). The segment-based EUM and the frame-based EUM re-scored these 50-best utterances separately.

Because the length of each utterance varies in connected digit recognition, we followed the convention of (Ming and Smith, 2001) to use the notation of relative order, which is the ratio of the order of the EUM against the length of an utterance. For example, a relative order of 0.2 would mean that 20% of the frames (or segment) in an utterance is assumed to be corrupted. Similar to (Ming and Smith, 2001), we chose a relative order of 0.1% or 10% for all the experiments.

The frame-based and segmented-based EUMs were first compared in two conditions, clean and noisy with PTC added at $S = -10$ dB, $C = 20\%$, $L = 25$ ms. The results are summarized in Table 2. For the segment-based experiments, a segment length of 6 frames were used. In the frame-based experiments, the segment length was effectively 1 frame. The frame-based EUM was possible computationally because of the efficient recursive formula described in Section 3. Under the noisy environments, both the segment-based and the frame-based EUM outperformed the Viterbi algorithm. Furthermore, the frame-based approach outperformed the segment-based approach in noisy condition by 12% showing that a finer resolution of segment is beneficial for the EUM. In clean condition in which no corruption was present, both approaches hurt performance with similar amounts of degradation. This is consistent with the reported results in Ming and Smith (2001).

Significance tests were performed on results on Table 2 using the matched pair test (Gillick and Cox, 1989). For the clean results, the differences between the segment-based, frame-based EUM and FEVA are not significant at a *p* value of 5%. Under the noisy conditions, the improvement of the frame-based EUM over the segment-based is significant (using a *p* value of 5% as above) and so is the improvement of FEVA over frame-based EUM.

Table 2
Comparison of different algorithms, frame-based EUM, segment-based EUM, FEVA and baseline Viterbi algorithm, in clean and noisy conditions

| Conditions | BL | Segment | Frame | FEVA |
|---|---|---|---|---|
| Clean | 98.83 | 98.50 | 98.58 | 98.44 |
| −10 dB, $C = 20\%$, $L = 25$ ms | 84.61 | 88.49 | 89.90 | 90.70 |

In terms of computation, the time (elapsed) used in rescoring the segment-based and fame-based EUM were similar at around 16 h. These included the initial HMM recognition, $N$-best generation, state-level alignment and the actual EUM re-scoring. The time (elapsed) used for FEVA is approximately 12 h which was slightly faster than re-scoring. It should be noted that for both schemes, computation was not optimized. For example, no pruning was applied in either search and the code was not optimized for speed. Furthermore, the computation time would be sensitive to the experimental setup. For example, the re-scoring time would depend on the number of $N$-best generated while FEVA computation would depend on the value of $M$.

In the second set of experiments, the frame-based and segmented-based EUMs were tested under different partial temporal corruption conditions. Columns 3 and 4 in Tables 3–5 show the comparison under various SNRs, rates of corruption and lengths of corruption.

In Table 3, the SNR was varied between −10 and 10 dB. Lower SNR caused more degradation but also increased the usefulness of the EUM. For the segment-based EUM, it was better than the baseline Viterbi only when the SNR fell to −10 dB. However, the frame-based EUM was useful even at 0 dB. In Table 4, the rate of corruption was varied between 20% and 40% at SNR = −10 dB and the corruption block length set at $L = 25$ ms. As the rate of corruption increased, the ASR performance degraded but increased the effectiveness of the frame-based EUM. This effect is similar to what was observed with changing SNR. Furthermore, consistent with results in Table 3, the frame-based EUM is always better than the segment-based EUM across the different rates of corruption.

The results for varying the length of corruption from 25 to 200 ms are shown in Table 5 in which we kept $C = 20\%$ and SNR = −10 dB. Because the expected number of observations to be corrupted remains constant, shorter length of corruption means more bursts of corruptions distributed across the whole sentence. This caused more damage to the utterances because more phonetic units would be affected. Since we are measuring word error rates, the more word corrupted increases the word error rate. On the other hand, a long segment of corruption would mean that

Table 3

Performances of different algorithms in PTC with varying SNR and $C = 20\%$, $L = 25$ ms at $M/N = 0.1$

| SNR | Viterbi | Seg-based EUM | Frame-based EUM |
|---|---|---|---|
| ∞ | 98.83 | 98.50 | 98.58 |
| 10 | 98.83 | 98.47 | 98.53 |
| 0 | 97.67 | 97.47 | 97.77 |
| −10 | 84.61 | 88.49 | 89.90 |

Table 4

Performances of different algorithms in PTC with varying rates of corruption and SNR = −10 dB, $L = 25$ ms with $M/N = 0.1$

| Rate of corruption (%) | Viterbi | Seg-based EUM | Frame-based EUM |
|---|---|---|---|
| 20 | 84.61 | 88.49 | 89.90 |
| 30 | 68.92 | 74.28 | 76.82 |
| 40 | 59.06 | 64.17 | 68.44 |

Table 5
Performances of different algorithms with varying corruption durations and SNR = −10 dB, $C$ = 20% with $M/N$ = 0.1

| Length of corruption (ms) | Viterbi | Seg-based EUM | Frame-based EUM |
|---|---|---|---|
| 25 | 84.61 | 88.49 | 89.90 |
| 50 | 87.86 | 90.38 | 91.31 |
| 75 | 90.82 | 92.88 | 93.72 |
| 100 | 93.16 | 94.92 | 95.55 |
| 125 | 95.31 | 96.55 | 96.69 |
| 150 | 95.48 | 96.65 | 96.86 |
| 175 | 95.35 | 96.47 | 96.59 |
| 200 | 95.22 | 96.17 | 96.38 |

only very few regions of corruption occurred. While one or even two phonetic units may be completely corrupted, other units would not be affect, thus reducing the damage (in terms of word error) from the impulse noise.

It is interesting to note that in the last experiment, as the assumption of long segment length became more accurate when the segment-length was longer than 100 ms, the performances of the frame-based and the segment-based EUMs became more similar. However, the frame-based EUM was still slightly better the segment-based EUM in these cases. One reason may be that while the added noise block is long and consistent with the segment assumption, it may fall between two segments thus, corrupting parts of two consecutive segments.

Comparing the EUM and the Viterbi algorithm, results shown in Tables 3–5 suggest that the EUM model is more powerful when the condition is worse (low SNR or high rate of corruption). Comparing the frame-base EUM and the segment-based EUM with a 6-frame segment, the frame-based EUM is consistently better than the segment-based EUM because of its finer resolution.

Our next set of experiments focused on the evaluation of the FEVA. Table 2 also shows the performance of using the FEVA under clean to −10 dB noisy conditions. Compared to the frame-based EUM with $N$-best re-scoring, the FEVA gave comparable performance with better performance at −10 dB and slightly worse performance under clean environment. However, the FEVA was better than the segment-based EUM. Results from Tables 3 and 4 suggest that the most interesting case for using the EUM with acceptable recognition performance is at SNR = −10 dB and rate $C$ = 20%. Similar to the experiments reported in Table 5, we varied the noise block length at this setting. To make comparison easier, we tabulate the average performance across different block lengths in Table 6. The results show that the performance of the FEVA was comparable but no better than $N$-best re-scoring on the frame-based EUM probably because of the approximations taken. However, the FEVA outperformed the segment-based

Table 6
Average performance of different algorithms in PTC across different block length of 25–200 ms at SNR = −10 dB with $M/N$ = 0.1

| Viterbi | Seg-based EUM | Frame-based EUM | FEVA |
|---|---|---|---|
| 92.23 | 94.06 | 94.63 | 94.70 |

EUM and allowed the use of EUM probability directly in recognition without the need for an extra $N$-best generation pass. Furthermore, its performance is not dependent on the quality of $N$-best which can be badly degraded under difficult noisy environments.

## 6. Conclusions and future work

Union model is one of the emerging acoustic modeling techniques that move beyond traditional HMM and has been shown to be effective in handling corrupted observations. The application of union model to speech recognition is hindered by two computational problems, namely, the computation of the EUM probability for a given state sequence and the search for the optimal EUM probability state sequence.

This paper analyses the properties of the EUM probability and addresses these two computation problems by proposing two recursive algorithms that avoid repeated computations.

The first algorithm is based on the EUM formula which exploits the recursive relationship between the EUM probabilities of a length $t$ sequence and a length $t + 1$ sequence. The proposed algorithm reduces the computation complexity of the EUM probability from an exponential order on sequence length to a linear order, making it possible to compute the EUM probabilities over long sequences. Because of this, it is no longer needed to artificially divide a speech utterance into segments. Instead, the EUM probability over individual speech frames, which we call the frame-based EUM, is feasible. The frame-based EUM not only removes the need to select the segment length which can be difficult and dependent on the noise environments, it also outperforms the segment-based EUM experimentally over different impulsive noise conditions. It is also important to note that even at long durations of impulsive noise close to the segment-length, the frame-based EUM is consistently better than the segment-based EUM.

To solve the EUM recognition problem, this paper describes an approximate Viterbi-like algorithm, called the frame-based EUM Viterbi algorithm (FEVA). The FEVA can efficiently perform recognition using the EUM probability as the criterion. Experimentally, the FEVA performed as well as the $N$-best re-scoring without the need for $N$-best re-scoring which can add complexity to the recognition paradigm. Results from $N$-best re-scoring could also be sensitive to the quality of the $N$-best generation. The derivation of the FEVA is also valuable for further development of the EUM.

The results we presented in the experiment section assumed a fixed relative order. It would be more desirable to have ways to select the order of the EUM. Work in Jancovic and Ming (2001) made used of duration model to re-score results with different orders. In Siu and Chan (2002), the problem was treated as a model-selection problem. Similar approaches may be applicable to the frame-based EUM and would be one focus of our future study.

## Acknowledgement

# References

Cooke, M., Morris, A., Green, P., 1997. Missing data techniques for robust speech recognition. In: Proc. of the IEEE Inter. Conf. on Acoust., Speech and Signal Proc., pp. 863–866.

de Veth, J., de Wet, F., Cranen, B., Boves, L., 1999. Missing feature theory in asr: make sure you miss the right type of features. In: Proc. of the Workshop on Robust Methods in Adverse Conditions. Tampere, Finland, pp. 231–234.

Duda, R.O., Hart, P.E., Stork, D.G., 2001. Pattern Classification, second ed. Wiley, New York.

Gillick, L., Cox, S., 1989. Some statistical issues in the comparison of speech recognition algorithms. In: Proc. of the IEEE Inter. Conf. on Acoust., Speech and Signal Proc., pp. 532–535.

Hirsch, H., Pearce, D., 2000. The aurora experimental framework for the performance evaluation of speech recognitions system under noisy conditions. In: International Workshop on Automatic Speech Recognition, ASR2000. Paris, France, pp. 181–188.

Jancovic, P., Ming, J., 2001. A probabilistic union model with automoatic order selection for noisy speech recognition. Journal of Acoustic Society of America 110, 1641–1648.

Lamere, P., Kwok, P., Walker, W., Gouvilla, E., Singh, R., Raj, B., Woolf, P., 2003. Design of the cmu sphinx-4 decoder. In: Proc. of the European Conference on Speech Comm. and Tech..

Leonard, R.G., 1984. A database for speaker-independent digit recognitionProc. of the IEEE Inter. Conf. on Acoust., Speech and Signal Proc., vol. 9, pp. 328–331.

Lippmann, R., Carlson, B., 1997. Robust speech recognition with time-varying filtering interruptions and noise. In: IEEE Automatic Speech Recognition and Understanding (ASRU) Workshop, pp. 37–40.

Ming, J., 2001. An improved union model for continous speech recognition with partial duration corruption. In: IEEE Automatic Speech Recognition and Understanding (ASRU) Workshop, pp. 25–28.

Ming, J., Jancovic, P., Smith, F., 2002. Robust speech recognition using probabilistic union models. IEEE Transactions on Speech and Audio Processing 10, 403–414.

Ming, J., Smith, F.J., 1999. Union: A new approach for combining sub-band observations for noisy speech recognition. In: International Workshop on Robust Methods for Speech Recognition in Adverse Conditions, pp. 175–178.

Ming, J., Smith, F.J., 2001. Union: a model for partial temporal corruption of speech. Computer Speech and Language 15, 217–231.

Raj, B., Singh, R., Stern, R.M., 1998. Inference of missing spectrographic for robust speech recognition. In: Proc. of the Inter. Conf. on Spoken Language Processing, pp. 1491–1494.

Schwartz, R., Chow, Y.L., 1990. The *N*-best algorithm: an efficient and exact procedure for finding the n most likely sentence hypotheses. In: Proc. of the IEEE Inter. Conf. on Acoust., Speech and Signal Proc. Albuquerque, pp. 81–84.

Sicilia-Garcia, E., Ming, J., Smith, F.J., 2002. Individual word language models and the frequency approach. In: Proc. of the Inter. Conf. on Spoken Language Processing, pp. 897–900.

Siu, M., Chan, A., 2002. Robust speech recognition against short-time noise. In: Proc. of the Inter. Conf. on Spoken Language Processing, vol. 2, pp. 1049–1052.

Young, S.J., Odell, J.J., Olasen, D., Woodland, P.C., 1999. The HTK Book (for HTK Version 3.0).