

Confidence Measures for Large Vocabulary Continuous Speech Recognition

Present by Tzan-Hwei Chen

Reference (1/3)

- [1] H. Jiang , “Confidence Measures for Speech Recognition : A Survey”, Speech communication 2005 .
- [2] F. Wessel , R. Schluter, K. Macherey, and H. Ney, “Confidence Measures for Large Vocabulary Continuous Speech Recognition”, IEEE SAP 2001.
- [3] R. Sarikaya, Y. Gao, M. Picheny and H. Erdogan, “Semantic Confidence Measurement for Spoken Dialog Systems.”, IEEE SAP 2005.
- [4] M. Afify, F. Liu, H. Jiang and O. Siohan, “A New Verification-based Fast-Match for Large Vocabulary Continuous Speech Recognition” IEEE SAP 2005

Reference (2/3)

- [5] F. Wessel , R. Schluter and H. Ney, “Using Posterior Word Probabilities for Improved Speech Recognition”, ICASSP 2000.
- [6] F. Wessel , R. Schluter and H. Ney, “Explicit Word Error Minimization Using Word Hypothesis Posterior Probabilities”, ICASSP 2001.
- [7] A. Kobayashi, K. Onoe, S. Sato and T. Imai, “Word Error Minimization Using an Integrate Confidence Measure”, INTERSPEECH 2005.
- [8] F. K. Soong and W. K. Lo, “Generalized Posterior Probability for Minimum Error Verification of Recognized Sentences”, ICASSP 2005

Reference (3/3)

- [9]J. Razik, O. Mella, D. Fohr, J.-P Haton, “Local Word Confidence Measure Using Word Graph and N-Best List.”
- [10]T, Fabian, R. Lieb, G. Ruske, M. Thomaе, “Impact of Word Graph Density on the Quality of Posterior Probability Based Confidence Measures.

Introduction (1/3)

- It is extremely important to be able to make an appropriate and reliable judgement based on the error-prone ASR result.
- researchers have proposed to compute a score (preferably 0~1), called confidence measure (CM) to indicate reliability of any recognition decision made by ASR system.
- Under a different name, namely ***utterance verification***, Rose et al.(1995) first formally cast the CM problem in speech recognition as a statistical hypothesis testing problem

Introduction (2/3)

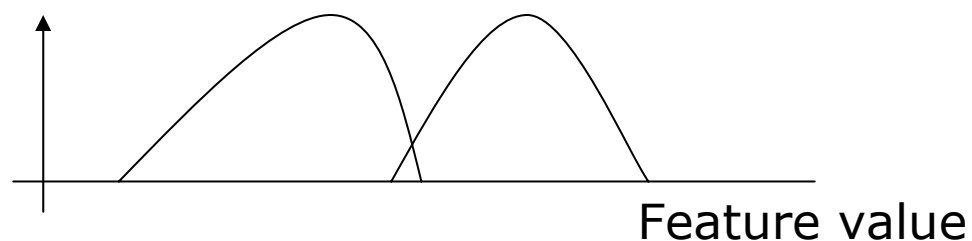
- First of all, we can backtrack some early research on CM to rejection in word-spotting systems.
- Other early CM-related works lie in automatic detection of new words in LVCSR.
- From the past few years, the CM is applied to more and more research areas, ex :
 - To improve speech recognition
 - The algorithm about look-head in LVCSR
 - Guide the system to perform unsupervised learning
 - ...

Introduction (3/3)

- all methods proposed for computing CMs can be roughly classified into three major categories :
 - Predictor features.
 - Posterior probability
 - Utterance verification (UV)

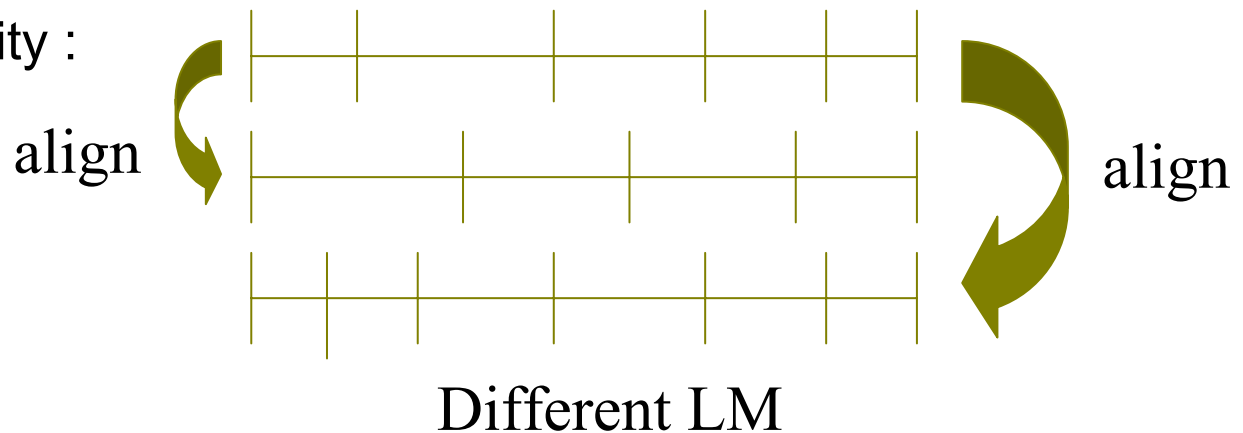
How to compute CMs (1/16)

- Predictor features :
 - Collected during decoding procedure and may include acoustic as well as language information
 - Any feature can be called a predictor if its p.d.f of correctly recognized words is clear distinct form that of misrecognized words



How to compute CMs (2/16)

- Predictor features (cont) : Some common predictor features
 - Pure normalized likelihood score related : acoustic score per frame.
 - N-best related : count in the N-best list, N-best homogeneity score, ...
 - Acoustic stability :



How to compute CMs (3/16)

- Predictor features (cont) : Some common predictor features

- Hypothesis density :

$$D(t') = |\{a' : [w_{a'}, s_{a'}, e_{a'}] \in WG \wedge s_{a'} \leq t' \leq e_{a'}\}|$$

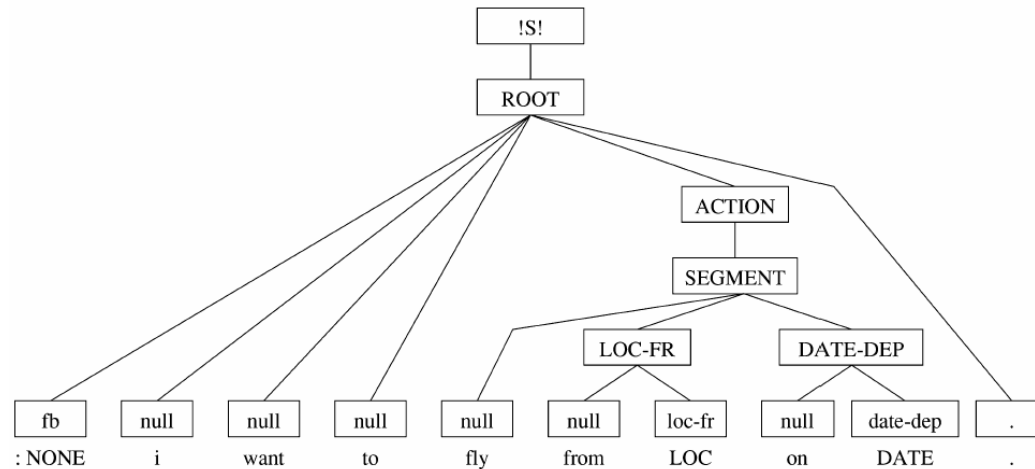
$$HD(a : [w_a; s_a, e_a]) = \frac{1}{e_a - s_a + 1} \sum_{t=s_a}^{e_a} D(t)$$

- Posterior probability
- Log-likelihood-ratio
- Duration related : phone duration, word duration
- LM score

How to compute CMs (4/16)

- Predictor features (cont) : Some common predictor features

– Parsing related :



Parse Tree Probabilities

	[!S!	[ROOT	:NONE_fb	i_null	want_null	to_null
Node	1	0.999516	1	0.984263	0.988049	0.994018
Extension	1	0.999522	0.984119	0.992231	0.995098	0.995098
	[ACTION	[SEGMENT	fly_flights	[LOC-FR	from_null	LOC_loc-fr]
Node	0.997064	1	0.990373	1	0.702338	0.991334
Extension	0.995984	0.995123	1	0.998472	0.999981	0.994569
	LOC-FR]	[DATE-DEP	on_null	DATE_date-dep	DATE-DEP]	SEGMENT]
Node	1	1	0.544135	0.998505	1	1
Extension	0.998472	0.990492	0.974586	0.999671	0.990492	0.995123
	ACTION]	...	ROOT]	!S!]		
Node	0.997064	1	0.999516	1		
Extension	0.995984	0.999519	0.999522	1		

How to compute CMs (5/16)

- Predictor features (cont) : some people attempt to combine several different features for a better performance
 - Line discriminant function
 - Generalized linear model
 - Neural networks
 - Decision tree
 - Support vector machine
 - Boosting
 - Naïve bayes classifier

How to compute CMs (6/16)

- Posterior probability :
 - The conventional ASR system

$$\begin{aligned}\hat{W} &= \arg \max_{W \in \Sigma} P(W | X) \\ &= \arg \max_{W \in \Sigma} \frac{P(X | W)P(W)}{P(X)} \text{ ignore} \\ &= \arg \max_{W \in \Sigma} P(X | W)P(W)\end{aligned}$$

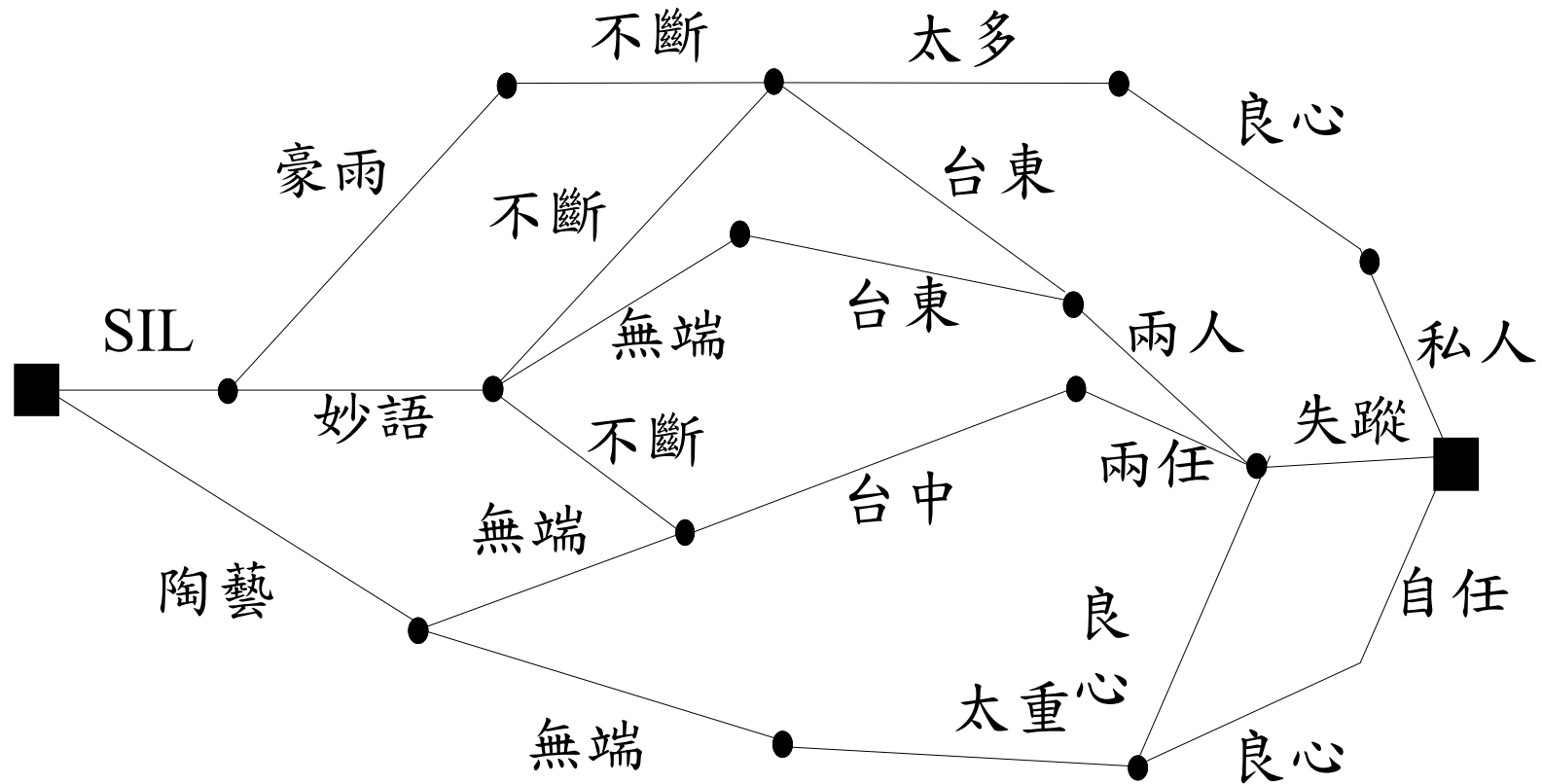
- In theory, we should compute $P(X)$ as follow :

$$P(X) = \sum_H P(X, H)$$

Impossible to estimate in a precise manner

How to compute CMs (7/16)

- Posterior probability (cont) :
 - Word graph is a compact and fairly accurate representation of all alternative competing hypotheses



How to compute CMs (8/16)

- Posterior probability (cont) :
 - The any arc in word graph can be computed as :

$$p([w, s, t] | x_1^T) = \frac{1}{p(x_1^T)} \sum_W \sum_{W'} p(x_1^{s-1} | W) p(x_s^t | w) \cdot p(x_{t+1}^T | W) \cdot p(W_w W') \quad (1)$$

- some issue are addressed and the word posterior probability is generalized
 - Reduced search space
 - Relaxed time registration
 - Optimal acoustic and language model weights

How to compute CMs (9/16)

- Posterior probability (cont) :
 - some approaches

$$p([w, s, e] | X) = \sum_{\substack{M, [w, s, t]_1^M \\ \exists n, 1 \leq n \leq M \\ w_n = w \\ (s_n, t_n) \cap (s, t) \neq \emptyset}} \frac{\prod_{m=1}^M p^\alpha(x_{s_m}^{e_m} | w_m) p^\beta(w_m | w_1^M)}{p(X)} \quad (2)$$

$$C_{\text{sec}}([w, s, e]) = \sum_{\substack{[w]_{s'}^{e'} \\ (s, e) \cap (s', e') \neq \emptyset}} p([w, s', e'] | X) \quad (3)$$

$$C_{\text{med}}([w, s, e]) = \sum_{\substack{[w]_{s'}^{e'} \\ s' \leq [s + \frac{e-s}{2}] \leq e'}} p([w, s', e'] | X) \quad (4)$$

$$C_{\text{max}}([w, s, e]) = \max_{t \in \{s \dots e\}} \sum_{\substack{[w]_{s'}^{e'} \\ s' \leq t \leq e'}} p([w, s', e'] | X) \quad (5)$$

How to compute CMs (10/16)

- Posterior probability (cont) :
 - The drawbacks of above methods – needed the additional pass.
 - So in [9], they proposed the “local word confidence measure”

$$C([w, s, e]) = \frac{\max(p(x_s^e | w))^\alpha p(w)^\beta}{\sum_{[w', s', e'] \in E} \max(p(x_{s'}^{e'} | w'))^\alpha p(w')^\beta} \quad (6)$$

$$C([w, s, e]) = \frac{\max(p(x_s^e | w))^\alpha \sum_{w'} p(w | w_h)^\beta}{\sum_{[w', s', e'] \in E} \max(p(x_{s'}^{e'} | w'))^\alpha \sum_{w'_h} p(w' | w'_h)^\beta} \quad (7)$$

$$C([w, s, w]) = \frac{\max(p(x_s^e | w))^\alpha \sum_{w_h} \sum_{w_f} \{p(w | w_h) p(w_f | w)\}^\beta}{\sum_{[w', s', e'] \in E} \max(p(x_{s'}^{e'} | w'))^\alpha \sum_{w'_h} \sum_{w'_f} \{p(w' | w'_h) p(w'_f | w')\}^\beta} \quad (8)$$

How to compute CMs (11/16)

- Posterior probability (cont) :
 - Impact of Word Graph Density on the Quality of Posterior Probability

Verbmobil		NaDia	
WGD	CER [%]	WGD	CER [%]
17.7	19.7	30.9	10.9
206.5	16.7	383.5	10.2
736	16	1170.4	10.3

Baseline

27.3

15.4

$$CER = \frac{\text{the number of classification errors}}{\text{the total number of recognized words}}$$

$$baseline = \frac{\text{sub. + ins.}}{\text{the total number of recognized words}}$$

How to compute CMs (12/16)

- Utterance verification

- The CM problem is formulated as a statistical hypothesis testing problem.

- Under the framework of UV, we propose two complementary hypotheses

H_0 (Null Hypothesis): X is correctly recognized and truly comes from model λ_w

H_1 (Alternative Hypothesis): X is wrongly recognized and is NOT from model λ_w

- Then we test H_0 against H_1

$$\text{likelihood ratio testing(LRT)} = \frac{P(X | H_0)}{P(X | H_1)} \underset{H_1}{\overset{H_0}{>}} \tau \quad (9)$$

How to compute CMs (13/16)

- Utterance verification (cont)
 - As pointed out by Lee(2001), the above LRT score can be transformed to a CM based on a monotonic 1-1 mapping function.
 - The major difficulty with LRT is how to model the alternative hypothesis.
 - In practice, the same HMM structure is adopted to model the alternative hypothesis.
 - A discriminative training procedure plays a crucial role in improving modeling performance.

How to compute CMs (14/16)

- Utterance verification (cont) – MCE training

frame based distance:

$$r_{ij}(y_t) = \log(a_{ij}^c b_j^c(y_t)) - \log(a_{ij}^a b_j^a(y_t)) \quad (10)$$

segment based distance is obtained by averaging the frame based distances as

$$R_u(Y^u) = \frac{1}{t_{f_u} - t_{i_u} + 1} \sum_{t=t_{i_u}}^{t=t_{f_u}} r_{q_{t-1}q_t}(y_t) \quad (11)$$

where the indicator function $\delta(u)$ is defined as

$$\delta(u) = \begin{cases} -1, & u \in \text{correct} \\ 1, & u \in \text{imposter} \end{cases}$$

A gradient update is performed on the expected cost $E\{F_u(Y^u, \Lambda^u)\}$

$$\Lambda_{n+1}^u = \Lambda_n^u - \varepsilon \nabla E\{F_u(Y^u, \Lambda^u)\} \quad (12)$$

How to compute CMs (15/16)

- Incorporation of high-level information for CM
 - LSA :

$$\begin{array}{c} s_1 s_2 \dots s_n \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{array} \begin{array}{c} \mathbf{A} \end{array} = \begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_m \end{array} \begin{array}{c} \mathbf{U} \end{array} \begin{array}{c} \mathbf{\Sigma} \end{array} \begin{array}{c} s_1 s_2 \dots s_n \\ \mathbf{V}^T \end{array}$$

- The key property of LSA is that words whose vectors are “close” correspond to semantically similar words.
- These similarities can be used to provide an estimate of the likelihood of the words co-occurring within the same utterance.

How to compute CMs (16/16)

- Incorporation of high-level information for CM (cont)
 - Inter-word mutual information :

Assume $N(x, y)$ be the co - occurrence times of word x and word y
in all training documents, the joint probability $P(x, y)$ is :

$$P(x, y) = \frac{N(x, y)}{\sum_{x, y} N(x, y)}$$

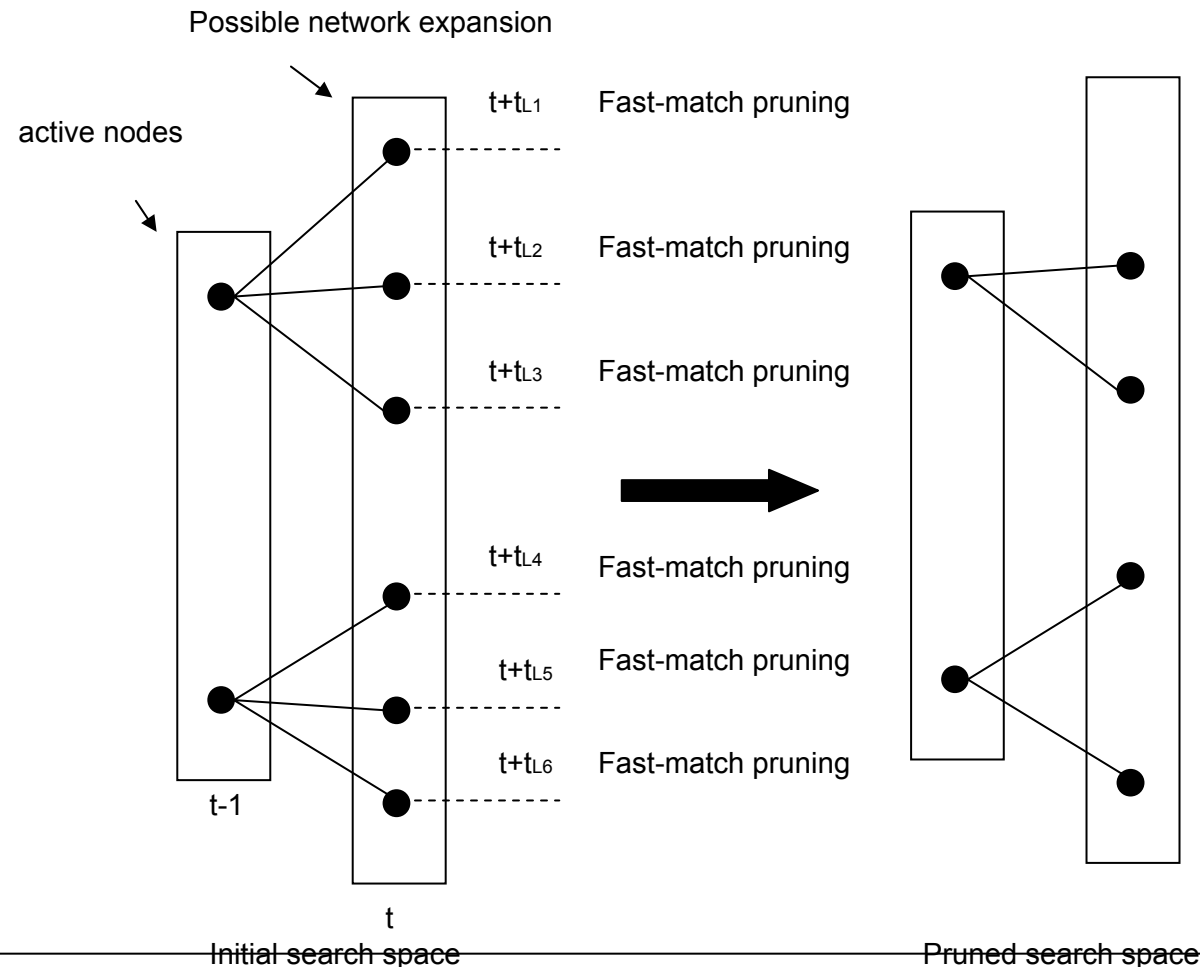
Mutual information between any two words x and y can be calculated as follows

$$MI = \log\left(\frac{p(x, y)}{p(x)p(y)}\right)$$

CM of each recognized word is calculated as the average mutual information of this word
with all other recognized words

Some recent applications (1/9)

- Verification-Based Fast-Match
 - Introduction



Some recent applications (2/9)

- Verification-Based Fast-Match (cont)
 - For a fast-match the null and alternative hypothesis testing for phoneme α can be written as :

H_0 : α starts at time t

H_1 : α does not start at time t

$$LRT = \frac{P(X | H_0)}{P(X | H_1)} \underset{H_1}{\overset{H_0}{>}} \tau$$

$$P(X | H_0) \equiv P(x_t^{t+d_\alpha} | \alpha)$$

$$P(X | H_1) \equiv P(x_t^{t+d_\alpha} | \bar{\alpha})$$

Some recent applications (3/9)

- Word error minimization :
 - Statistical decision theory aims at minimizing the expected of making error

$$w_1^{N*} = \arg \max_{w_1^N} P(w_1^N | x_1^T) \quad (13)$$

- To assume the boundary time of word sequence is given [5]:

$$\begin{aligned} p(w_1^N | x_1^T) &= p([w, s, t]_1^N | x_1^T) \\ &= \prod_{n=1}^N p([w_n, s_n, t_n] | [w, s, t]_1^{n-1}, x_1^T) \\ &= \prod_{n=1}^N p([w_n, s_n, t_n] | x_1^T) \end{aligned} \quad (14)$$

Some recent applications (4/9)

- Word error minimization (cont) :
 - Minimizing the expected SER does not necessarily minimizing the expected WER
 - In [7]

$$w_1^{N*} = \arg \min_{w_1^N} \text{WER}(w_1^N | x_1^T) \quad (15)$$

$$\text{WER}(w_1^N | x_1^T) = 1.0 - \frac{1}{N} \sum_i \{P(w_i = \textit{correct}) \times P(w_i | x_1^T)\} \quad (16)$$

Some recent applications (5/9)

- Word error minimization (cont) :

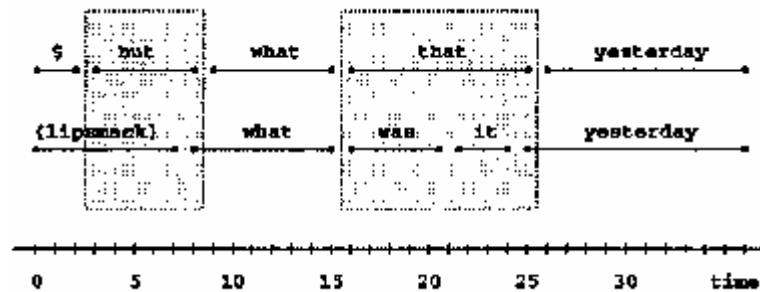
- In this very general framework

$$w_1^{N*} = \arg \min_{w_1^N} \left\{ \sum_{v_1^M} C(w_1^N, v_1^M) \cdot p(v_1^M | x_1^T) \right\} \quad (17)$$

- The easiest way to overcome this mismatch it to use the same cost function – Levensthein distance
- In (Stolcke et. al 1997), the pairwise alignment is restricted to N-best list.
- Let us assume that sub. were the one type of error.
 - A dynamic programming alignment would thus not be necessary.

Some recent applications (6/9)

- Word error minimization (cont) :
 - With these considerations and with the fact that a word graph contains the start and end time :

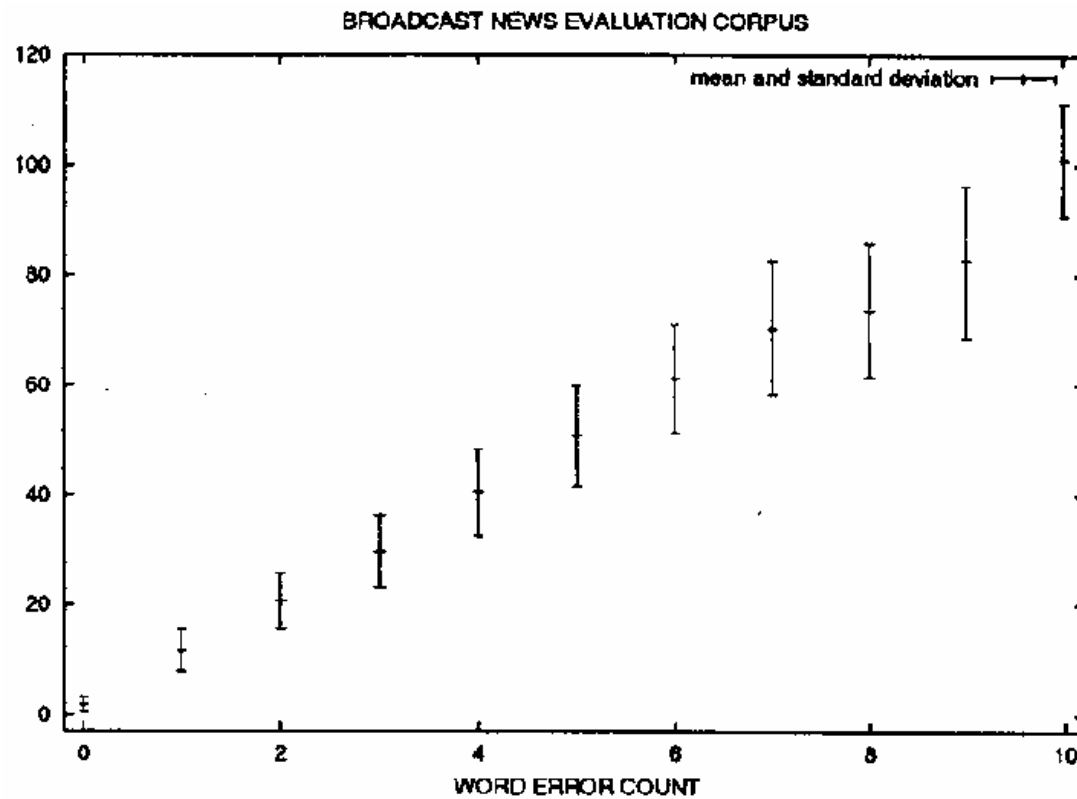


- We can now define a new cost function using the notion of time frame errors

$$C(w_1^N, v_1^M) = \sum_{n=1}^N \frac{\sum_{t=s_n}^{t=e_n} 1 - \delta(w_n, v_t)}{1 + \alpha(e_n - s_n - 1)} \quad (18)$$

Some recent applications (7/9)

- Word error minimization (cont) :
 - Correlation analysis



Some recent applications (8/9)

- Word error minimization (cont) :

Time Frame Error decoding [6]

$$\begin{aligned}
 \{[w, s, e]_1^N\}_{opt} &= \arg \min_{[w, s, e]_1^N} \left\{ \sum_{[v, s', e']_1^M} \ell([w, s, e]_1^N, [v, s', e']_1^M) P([v, s', e']_1^M | x_1^T) \right\} \\
 &= \arg \min_{[w, s, e]_1^N} \left\{ \sum_{[v, s', e']_1^M} \sum_{n=1}^N \frac{\sum_{t=s_n}^{t=e_n} 1 - \delta(w_n, v_t)}{1 + \alpha(e_n - s_n - 1)} P([v, s', e']_1^M | x_1^T) \right\} \\
 &= \arg \min_{[w, s, e]_1^N} \left\{ \sum_{n=1}^N \frac{\sum_{t=s_n}^{t=e_n} \sum_{[v, s', e']_1^M} 1 \times P([v, s', e']_1^M | x_1^T) - \sum_{t=s_n}^{t=e_n} \sum_{[v, s', e']_1^M} \delta(w_n, v_t) P([v, s', e']_1^M | x_1^T)}{1 + \alpha(e_n - s_n - 1)} \right\} \\
 &= \arg \min_{[w, s, e]_1^N} \left\{ \sum_{n=1}^N \frac{\sum_{t=s_n}^{t=e_n} \left[1 - \sum_{[v, s', e']_1^M} \delta(w_n, v_t) P([v, s', e']_1^M | x_1^T) \right]}{1 + \alpha(e_n - s_n - 1)} \right\}
 \end{aligned}$$

Some recent applications (9/9)

- Word error minimization (cont) :

$$\begin{aligned} & \sum_{[v,s',e']_1^M} \delta(w_n, v_t) P([v, s', e']_1^M | x_1^T) \\ &= \sum_{[v,s',e']_1^M} \sum_{v_m: s'_m \leq t \leq e'_m} \delta(w_n, v_t) P([v, s', e']_1^M | x_1^T) \\ &= \sum_{[v,s',e'], v: s' \leq t \leq e'} \delta(w_n, v) P([v, s', e'] | x_1^T) \\ &= p(w_n | t, x_1^T) \end{aligned}$$

$$\frac{\sum_{t=s_n}^{t=e_n} [1 - p(w_n | t, x_1^T)]}{1 + \alpha(e_n - s_n - 1)}$$

Can be interpreted as the normalized Probability of a word being incorrect.

Summary

- Almost all CMs in acoustic level fundamentally rely almost entirely on a single information source.
- We believe it is critical to improve performance of CMs by taking this segmentation issue into account.