

# Adding Word Duration Information to Bigram Language Models

George Doddington

Yufeng Wu, Aravind Ganapathiraju, Joseph Picone

National Institute of Standards and Technology  
 doddington@nist.gov  
 http://www.nist.gov

Institute for Signal and Information Processing  
 Mississippi State University  
 {wu, ganapath, picone}@isip.msstate.edu  
 http://www.isip.msstate.edu



### Motivation

reference: found out that that was n't  
 baseline: and all that was an  
 duration model found out that was an

- suprasegmental information plays an important role in human speech
- usually model suprasegmental information jointly with segmental measures

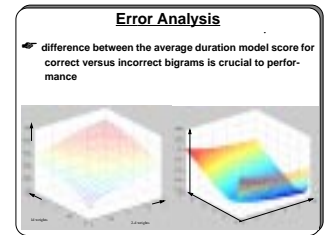
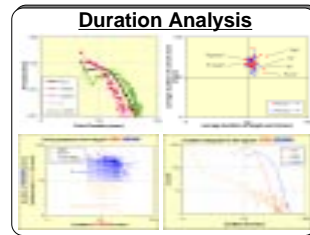
### Bigram Duration Model

- Modeling duration within the context of the bigram language model

$$Pr(w_{i-1}^{\tau_{i-1}} w_i^{\tau_i}) = Pr(w_{i-1}^{\tau_{i-1}} w_i^{\tau_i}) Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}})$$

$$= \left[ \frac{Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})}{Pr(w_{i-1}^{\tau_{i-1}})} \right] \left[ \frac{Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}})}{Pr(w_i^{\tau_i})} \right]$$

$$Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}}) = \left[ \frac{Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})}{Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})} \right] \left[ \frac{Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})}{Pr(w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})} \right]$$

$$Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}}) = \left[ \frac{Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}} | w_{i-2}^{\tau_{i-2}})}{Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}})} \right] \left[ \frac{Pr(w_i^{\tau_i} | w_{i-1}^{\tau_{i-1}})}{Pr(w_i^{\tau_i})} \right]$$


### Suprasegmental Information

- each feature represented as a single scalar attribute of each word
- incorporating word duration into a bigram language model:

$$Pr(F_{i-1}^{\tau_{i-1}}) = \frac{Pr(w_{i-1}^{\tau_{i-1}} | F_{i-1}^{\tau_{i-1}})}{Pr(w_{i-1}^{\tau_{i-1}})} Pr(w_{i-1}^{\tau_{i-1}})$$

where,  $\tau$  is duration,  $w$  is word identity and  $F$  is the new feature vector

### Back-off Weighting

- sparsity of training data for many duration bigrams
- combine bigram-specific models with word-specific and word-independent models in a back-off framework

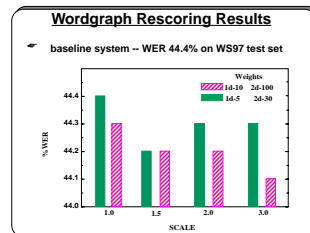
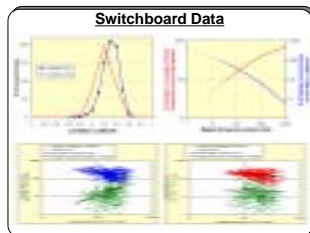
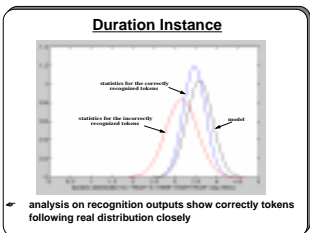
$$Pr_{smoothed}(w_{i-1}^{\tau_{i-1}} w_i^{\tau_i}) = \frac{\Omega_{bigram} Pr(w_{i-1}^{\tau_{i-1}} w_i^{\tau_i}) + \Omega_{word} Pr(w_{i-1}^{\tau_{i-1}}) Pr(w_i^{\tau_i}) + \Omega_{global} Pr(w_i^{\tau_i})}{\Omega_{bigram} + \Omega_{word} + \Omega_{global}}$$

### N-best Rescoring Results

- baseline - 32.4% WER on 637 SWB utterances
- rescoring the 100-best hypotheses
- N-best error rate of 21.2%

	[ weight 1d, weight 2d]			
scale	0.1 0.1	0.1 0.5	0.5 0.1	
0.01	32.5	32.4	32.3	
0.05	32.4	32.3	32.2	
0.1	32.3	32.3	32.2	

### Summary



### Future Work