

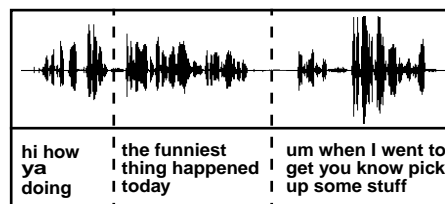
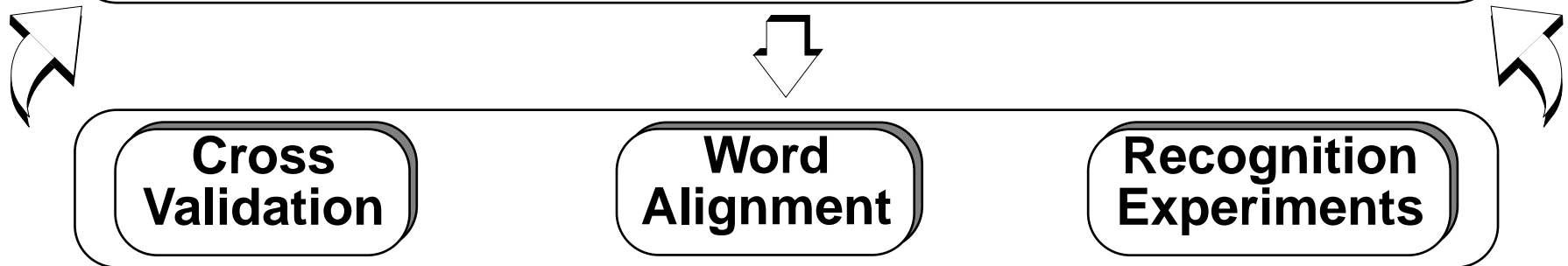
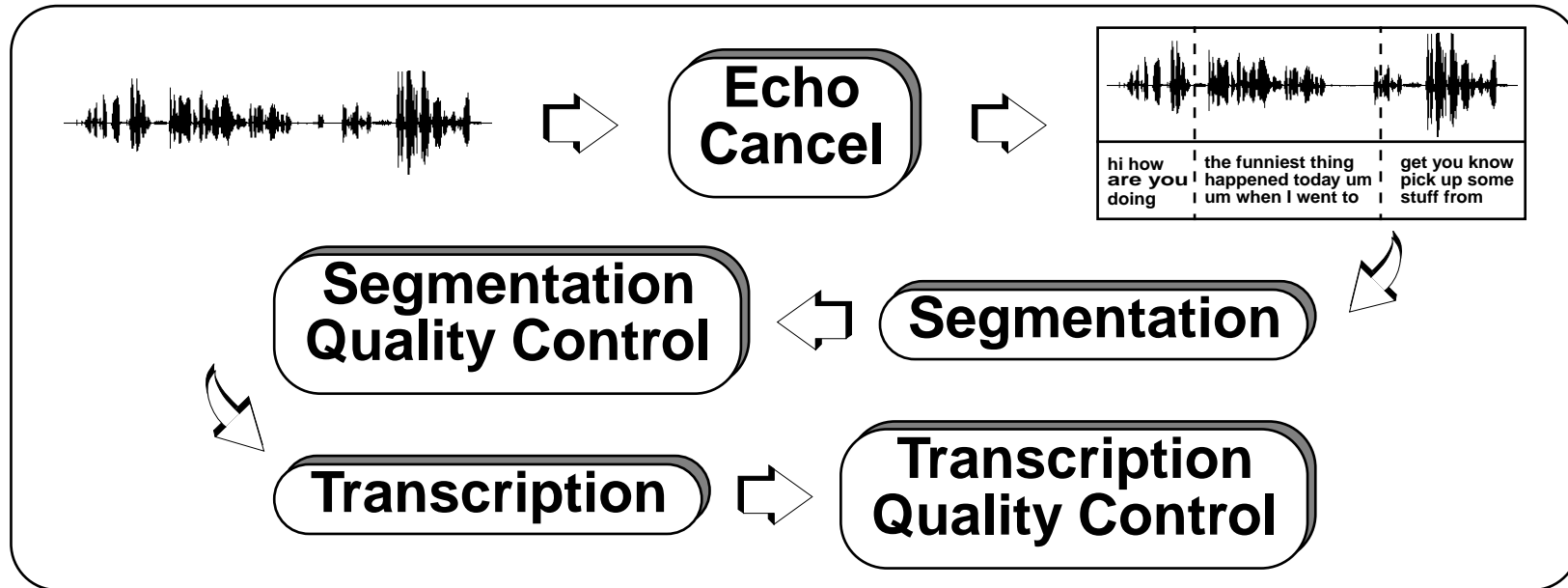
# Introduction to SWITCHBOARD

- ➔ **Challenging and popular LVCSR benchmark**
- ➔ **Spontaneous telephone conversations**
- ➔ **240 hours, 2430 conversations, 3+ million words, 500+ speakers (male and female)**
- ➔ **Low bandwidth, channel noise, echo**
- ➔ **Speaking rates, dialects, coarticulation, speaking styles, accents, dysfluencies**
- ➔ **Poor quality acoustic models, large mismatch**

# Motivation

- ➔ **Reduce acoustic model mismatch**
- ➔ **Segmentation and transcription must capture both acoustic and linguistic properties**
- ➔ **Automatic (energy-based) segmentation — unnatural breakpoints**
- ➔ **Linguistic structure-based segmentation — corrupted acoustic context**
- ➔ **Dysfluencies make transcription difficult (Current LDC transcription WER ~ 8%)**

# Approach



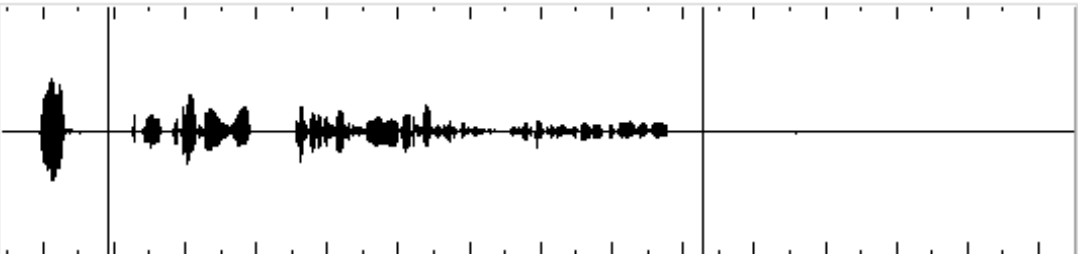
# Guidelines

- ➡ **Segment boundaries with at least 1 sec of silence between speech**
- ➡ **Segment along phrase / sentence / train-of-thought boundaries**
- ➡ **Merge utterances split at counterintuitive points (e.g. middle of sentence)**
- ➡ **Limit maximum utterance duration to 15 sec**
- ➡ **Fix transcriptions taking into account dysfluencies and capitalization issues**

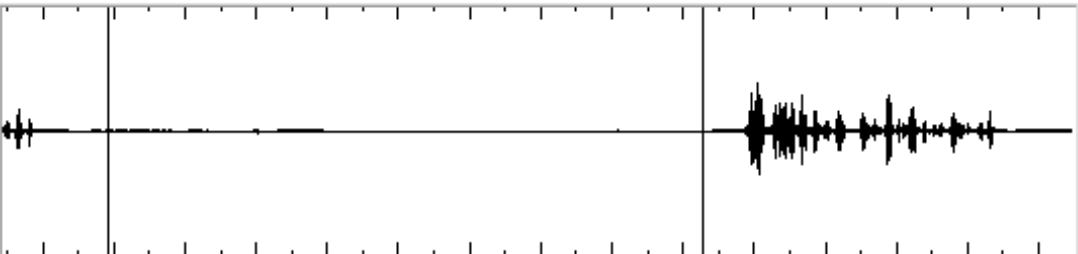
# Segmentation Tool

The Segmenter

W: [419.371125 , 434.524125] B: [420.908250 , 429.266125 :: 8.357875] T: [426.947625]



A-0069 A-0070 A-0071  
B-0090 B-0091 B-0092 B-0093




ID: sw2121A-ms98-a-0069 Endpoints: [419,371125 , 420,908250]  
yeah

ID: sw2121A-ms98-a-0070 Endpoints: [420,908250 , 429,266125]  
[noise] yeah we could also uh push for legislation for uh rapid transit systems uh  
this country seems to be a little behind on that

ID: sw2121A-ms98-a-0071 Endpoints: [429,266125 , 434,524125]  
[silence]

The Segmenter



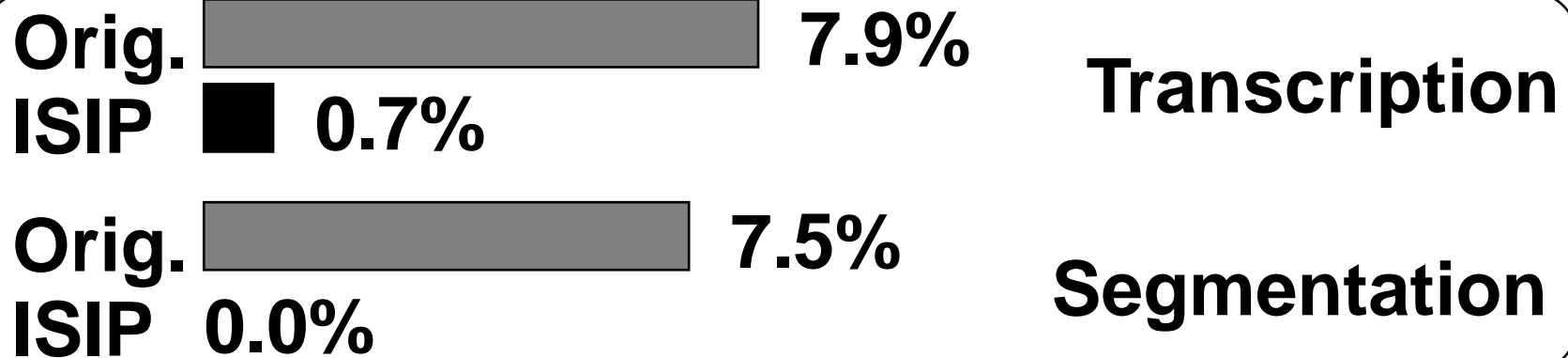
Load Config Log  
Save Help Quit  
Channel 0 [A] 1 [B]  
Lock Segments Transcriptions  
Volume  
Amplitude Gain  
Quick Play -2 -1 0 +1 +2  
Delete sw2121A-ms98-a-0067  
Insert sw2121A-ms98-a-0068  
Clear sw2121A-ms98-a-0069  
Set sw2121A-ms98-a-0070  
Play sw2121A-ms98-a-0071  
Merge sw2121A-ms98-a-0072  
Split sw2121A-ms98-a-0073  
Word sw2121A-ms98-a-0074  
Lexicon sw2121A-ms98-a-0075  
Verify sw2121A-ms98-a-0076  
sw2121A-ms98-a-0077  
sw2121A-ms98-a-0078  
sw2121A-ms98-a-0079  
sw2121A-ms98-a-0080  
sw2121A-ms98-a-0081  
sw2121A-ms98-a-0082  
sw2121A-ms98-a-0083

# Issues and Concerns

- ➡ **Large number of dysfluencies (pauses, laughter, partially pronounced words etc.)**
- ➡ **Affirmative statements (yes/no) and pause fillers (um/hmm) cover ~ 30% of utterances**
- ➡ **Marking boundaries near noise or echo**
- ➡ **Consistency in capitalization (“I” vs “i”) and handling proper nouns**
- ➡ **Marking asides, background noise / music and background speech**

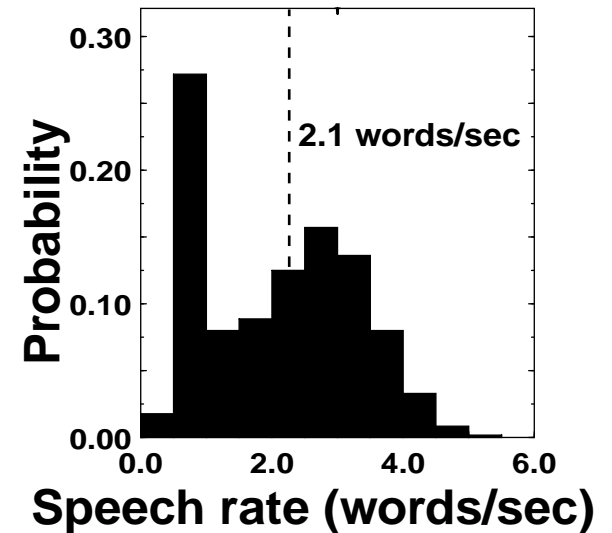
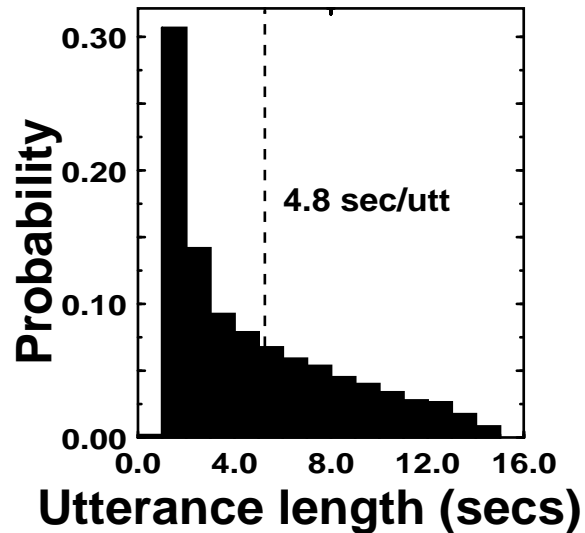
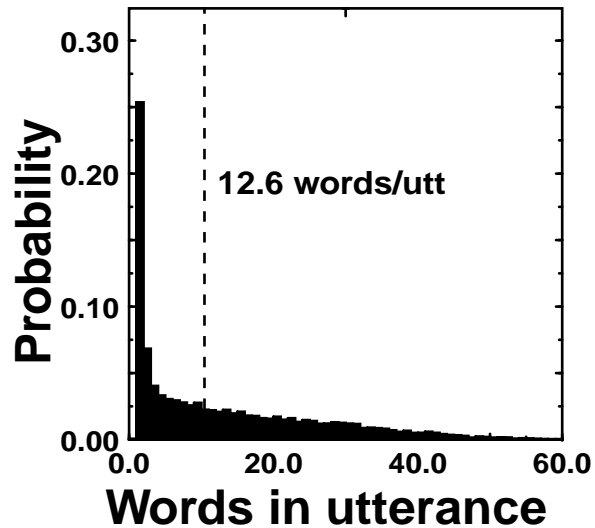
# Cross-Validation

- ➡ All validators segment / transcribe the same conversation
- ➡ Adjudicated reference transcription
- ➡ Word alignment review will further reduce error rate



**Cross-Validation Word Error Rates**

# The New SWITCHBOARD



- ➡ **Segmentation and transcription rate 20xRT**
- ➡ **Monosyllabic words constitute 53% of data on WS'97 subset (down from 67%)**
- ➡ **Lexicon updates — partial words, laughter words, alternate pronunciations**



# Effect on Recognition

- ➡ **Adapt existing acoustic models to resegmented speech data**
- ➡ **20 hours training data (27500 utterances) including silence**
- ➡ **Word-internal triphone system to bootstrap seed models (HTK)**
- ➡ **4 passes of re-estimation**
- ➡ **Lattice rescoring on WS'97 dev test set**

# Results

<b>Error Rate</b>	<b>ISIP</b>	<b>WS'97</b>
<b>Total WER</b>	<b>47.9%</b>	<b>49.8%</b>
<b>Correct words</b>	<b>55.8%</b>	<b>53.1%</b>
<b>Substitutions</b>	<b>31.6%</b>	<b>32.2%</b>
<b>Deletions</b>	<b>12.6%</b>	<b>14.8%</b>
<b>Insertions</b>	<b>3.7%</b>	<b>2.9%</b>

- ➡ **63% of total errors on monosyllabic words (down from 71%)**
- ➡ **Reduction in substitution and deletion errors**

# Analysis

- ➡ **1.9% absolute improvement in WER**
- ➡ **Monosyllabic words are the principal factor in error analysis**
- ➡ **Performance improvement attributed to better modeling of monosyllabic words**
- ➡ **Acoustically “complete” transcriptions (no partial words at utterance boundaries) help in improved acoustic modeling**
- ➡ **Longer utterance transcriptions facilitate LM application**

# Conclusions

- ➡ **Uniformity and accuracy are critical for the quality of training — segmentation and transcriptions**
- ➡ **Segmentation at natural boundaries allows better acoustic modeling**
- ➡ **Dysfluencies pose significant challenges to accurate transcription**
- ➡ **Acoustic models trained on corrected SWB data will result in major improvements in WER (e.g. 2% absolute improvement from adapting models)**