# Speech interfaces:
# A survey and some current projects

Dan Ellis & Nelson Morgan
International Computer Science Institute
Berkeley CA
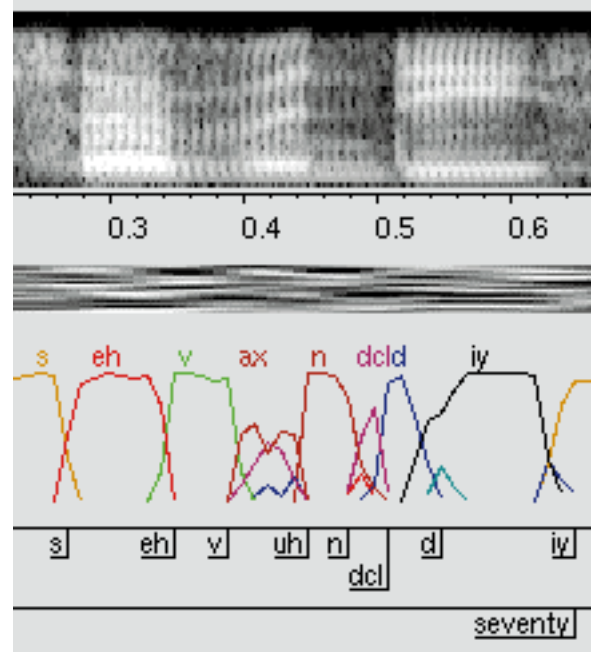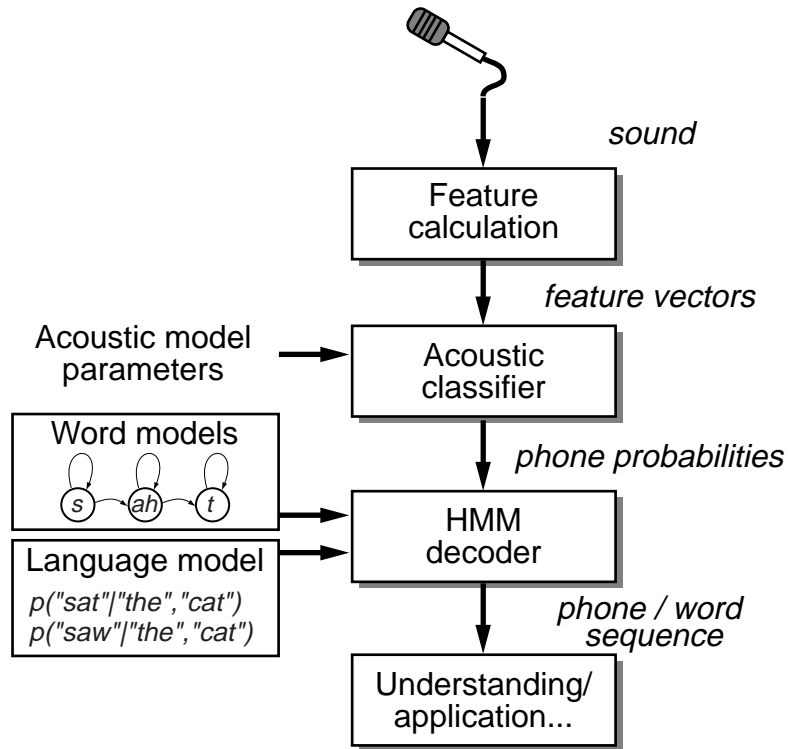{dpwe,morgan}@icsi.berkeley.edu

## Outline

**1** **Speech recognition: the state of the art**

**2** **Current projects at ICSI**

**3** **Conclusions**

**1**      # Speech recognition



- **Elements of a recognizer:**
  - feature design
  - acoustic modeling
  - pronunciation/language modeling    } data!

# How good is speech recognition?

- **Standard measure is word error rate (WER):**
  - dictation (close-mic): 2-5%
  - broadcast news: ~15%
  - telephone conversations: ~30%

F0: THE VERY EARLY RETURNS OF THE NICARAGUAN PRESIDENTIAL ELECTION SEEMED TO FADE BEFORE THE LOCAL MAYOR ON A LOT OF LAW

F4: AT THIS STAGE OF THE ACCOUNTING FOR SEVENTY SCOTCH ONE LEADER DANIEL ORTEGA IS IN SECOND PLACE THERE WERE TWENTY THREE PRESIDENTIAL CANDIDATES OF THE ELECTION

- **What are the problems?**
  - acoustic variability (noise, channel)
  - speech variability (accent, manner)
  - exploiting linguistic constraints
  - speech understanding...

# Frontiers of speech recognition

- **Acoustic modeling**
  - beyond head-mounted mics
  - background noise (mobile phones)
  - speech in mixtures (broadcast)
  - $\rightarrow$  robust feature design, better statistical models

- **Speaking styles**
  - coarticulation
  - pronunciation variability
  - speaking styles
  - $\rightarrow$  lump into acoustic model, more training data, better pron. models, context-dep. models

- **Linguistic constraints**
  - 'inferred' words
  - ambiguity
  - $\rightarrow$  higher-order n-grams (more training data), tree grammars

# Applications of speech recognition

- **Command & control**
  - more or less constrained

- **Dictation**
  - large vocabulary
  - known, co-operative user

- **Voice response systems**
  - dialog & speech understanding
  - robustness!
  - human factors (timing, barge-in etc.)

- **Information extraction & retrieval**
  - multimedia archive retrieval
  - live 'listener'

# Outline

**1** **Speech recognition: the state of the art**

**2** **Current projects at ICSI**

- Recognizer confidence measures
- Combining information sources
- The meeting recorder
- Audio content-based retrieval

**3** **Conclusions**

# Recognizer confidence measures

(Warner Warren, Andy Hatch, Eric Fosler + SRI)

- **Knowing which words are wrong can help**
  - hard to tell because recognition only *just* works
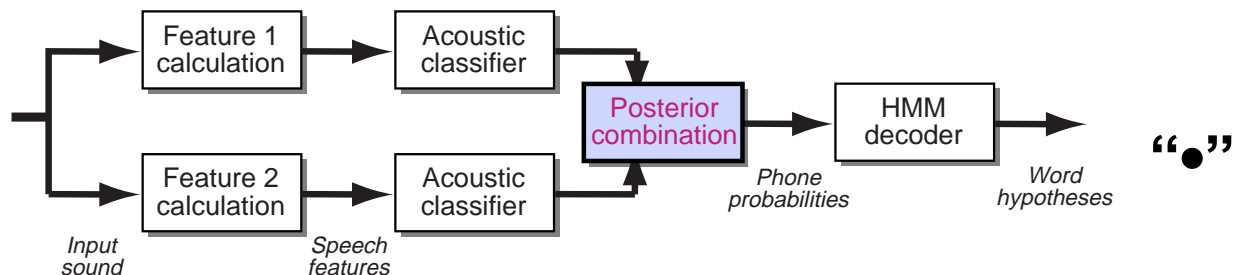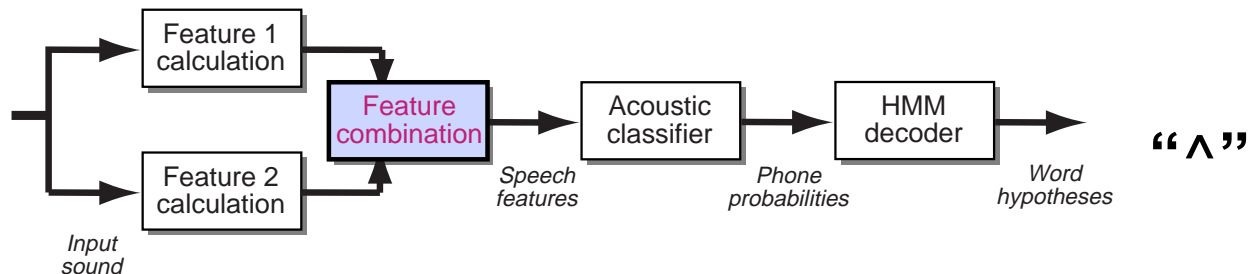
- **Average per-phone entropy + re-estimation:**



DET plot for word-level confidence estimation (AURORA)

- **Use for combining recognizer outputs**

# Combination schemes
## (Mike Shire, Barry Chen + Michael Jordan)



- **How best to combine different feature streams?**

| Features | Avg. WER |
|---|---|
| plp | 8.9% |
| msg | 9.5% |
| plp ^ msg | 8.1% |
| plp • msg | 7.1% |

# Tandem acoustic modeling

(with Hermansky et al., OGI)

- **ICSI pioneered 'hybrid-connectionist' ASR; Can it be combined with conventional models?**
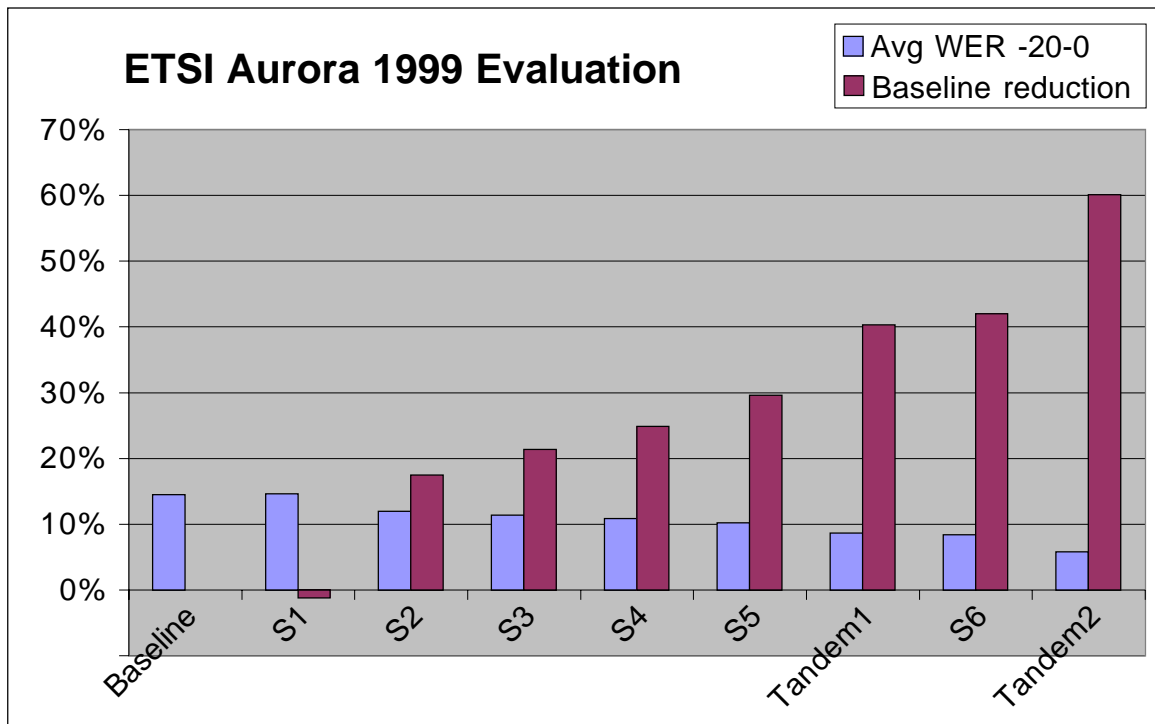


- **Result: better performance than either alone!**
  - neural net & Gaussian mixture models extract different information from training data

| System-features | Avg. WER |
|---|---|
| HTK-mfcc | 13.7% |
| Hybrid-mfcc | 9.3% |
| Tandem-mfcc | 7.4% |
| Tandem-plp+msg | 6.4% |

# Aurora "Distributed SR" evaluation

- **Organized by ETSI
  (European Telecoms. Standards Institute)**

**ETSI Aurora 1999 Evaluation**

Legend:
- ☐ Avg WER -20-0
- ☐ Baseline reduction

Chart (y-axis 0% to 70%, categories: Baseline, S1, S2, S3, S4, S5, Tandem1, S6, Tandem2)

- Tandem systems from OGI-ICSI-Qualcomm

# The meeting recorder project

(Adam Janin, Eric Fosler + UW, SRI, UPM, James Landay)

- **Idea: PDA records meetings to replace / enhance note-taking**

- **First task: Collect a training corpus**



- **Related to DARPA Communicator, SmartKom**

# Meeting recorder: Research areas

- **Audio recognition**
  - recognition from noisy microphones
  - speaker identification & tracking
  - nonspeech events

- **Indexing application**
  - understanding the structure of meetings
  - information retrieval
  - user interface

# Audio content-based retrieval
### (with Sheffield, Cambridge, BBC, Avideh Zakhor)

- **Idea: speech recognition output as indexes for broadcast news**
  - useful even with 15-30% WER

# Audio-video organization & retrieval

- **Proposed project:**



- **Synergy between audio & video features**

- **Query by terms or by examples**

- **Recovering temporal structuret**

# **Conclusions**

**3**

- **Speech recognition is now practical**
  - .. but still plenty of problems

- **Ongoing research in speech recognition**
  - recognition in demanding conditions
  - understanding / discourse a big issue

- **Multimodal information retrieval**
  - forgiving & fertile research area