# The Growing Impact of Speech Technology on Society
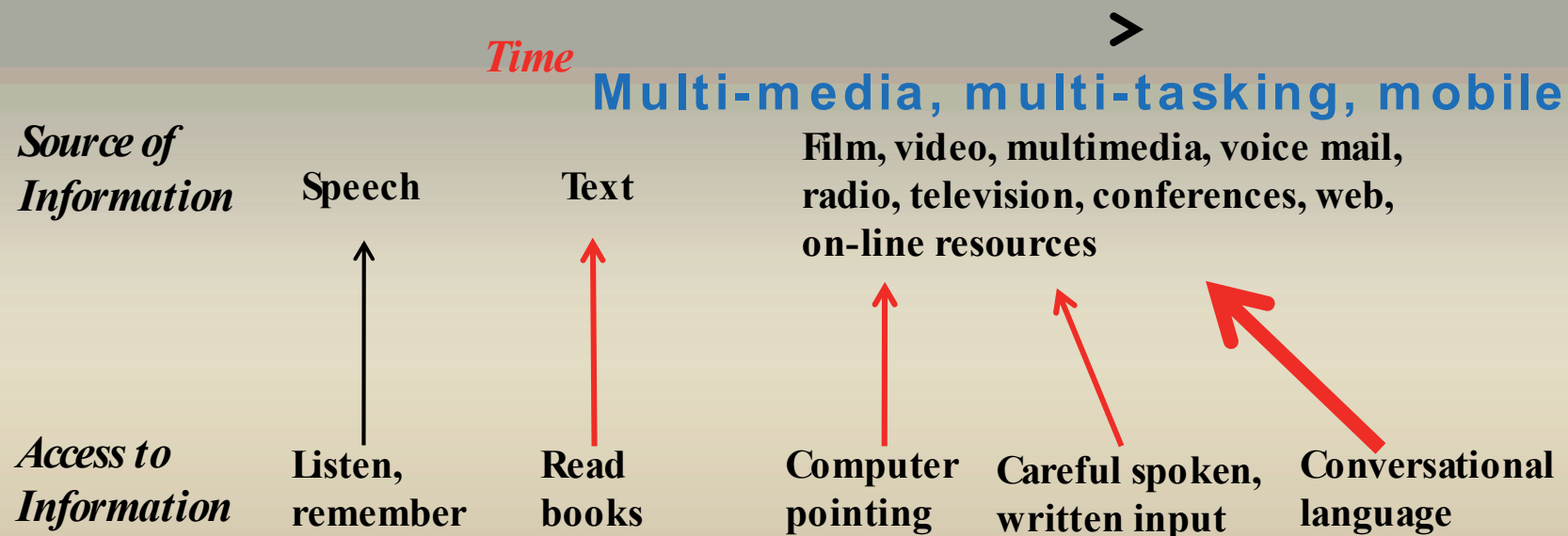
Patti Price, PPRICE Speech and Language Technology

**Intro: a few million years in about a minute**

Short review of speech recognition technology

Speech recognition performance

Social impact 1 (effect of society on speech)

Social impact 2 (people vs. technology)

Progress and challenges

# Speech in the Information Age

- **Speech & text were revolutionary because of information access**
- **New media and connectivity yield information overload**
- **Can speech technology help?**

>

*Time*

**Multi-media, multi-tasking, mobile**

*Source of Information*    Speech    Text    Film, video, multimedia, voice mail, radio, television, conferences, web, on-line resources

*Access to Information*    Listen, remember    Read books    Computer pointing    Careful spoken, written input    Conversational language

**Speech is social in ways our technology is not. Can it become a complementary partner in with humans?**
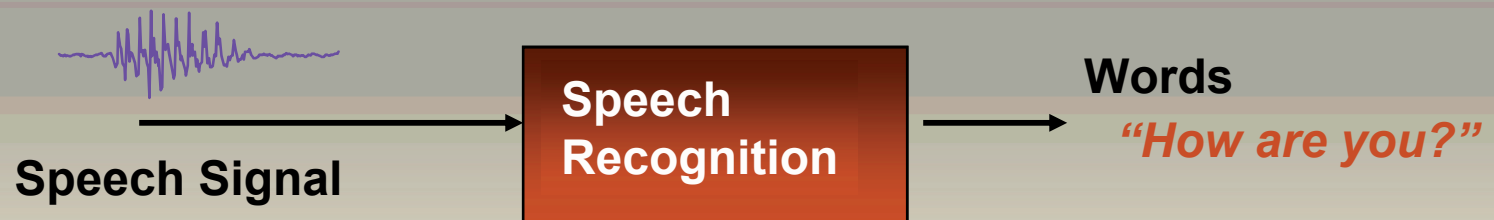
# The Growing Impact of
# Speech Technology on Society

Patti Price, PPRICE Speech and Language Technology

- Introduction: a few million years in about a minute

  **Short review of speech recognition technology**

  Speech recognition performance

  Social impact 1 (effect of society on speech)

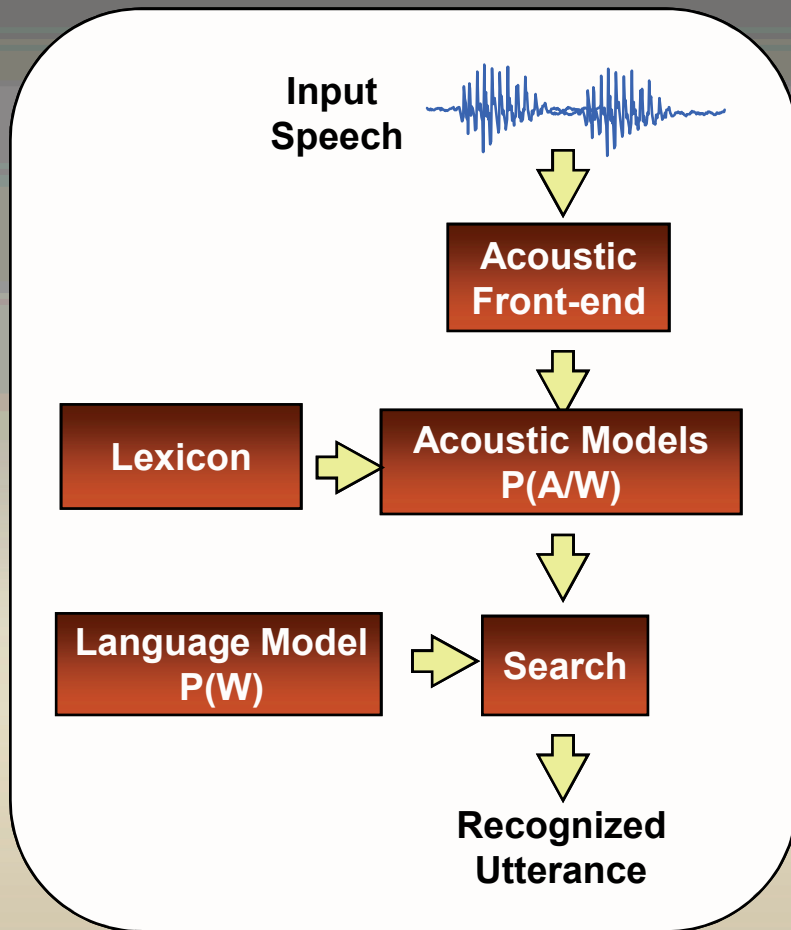  Social impact 2 (people vs. technology)

  Progress and challenges

# What is Speech Recognition?

**Goal:** **Automatically extract the string of words spoken from the speech signal**

**Speech Signal** → **Speech Recognition** → **Words** *"How are you?"*

Speech recognition does NOT determine

Who is talker (speaker recognition)

Speech output (speech synthesis or speech generation)

What the words mean (next two talks will address that)

# Recognition Architectures



Input Speech → Acoustic Front-end → Acoustic Models P(A/W) ← Lexicon; Acoustic Models → Search ← Language Model P(W); Search → Recognized Utterance

- The signal is converted to a sequence of feature vectors based on spectral and temporal measurements.

- Acoustic models represent sub-word units, such as phonemes, as a finite-state machine in which states model spectral structure and transitions model temporal structure.

- The language model predicts the next set of words, and controls which models are hypothesized.

- Search is crucial to the system, since many combinations of words must be investigated to find the most probable word sequence.

Probabilistic modeling requires training data, and match between test and training.

# The Growing Impact of Speech Technology on Society

Patti Price, PPRICE Speech and Language Technology

- Intro: a few million years in about a minute
- Short review of speech recognition technology

**Speech recognition performance**

Social impact 1 (effect of society on speech)

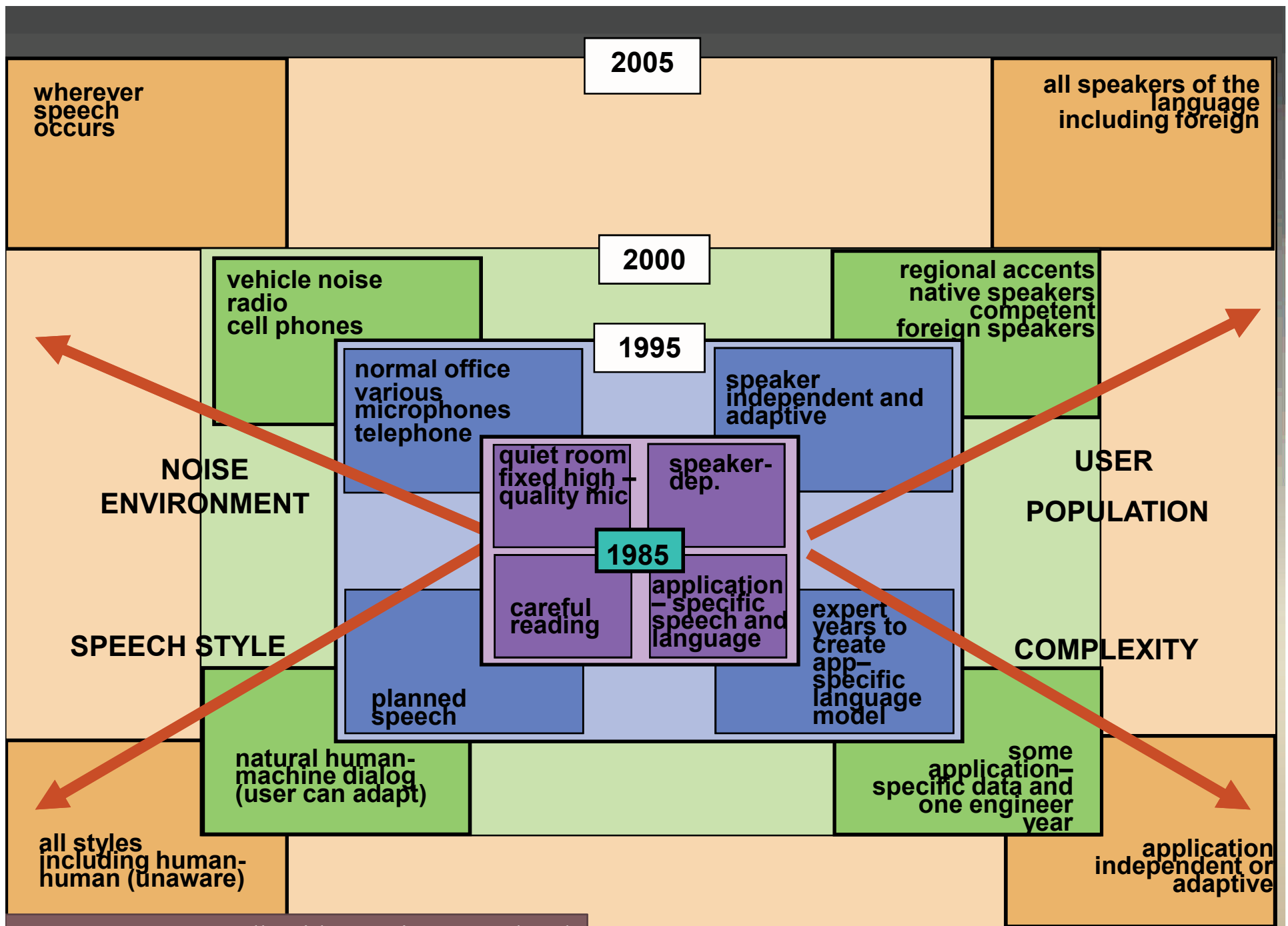Social impact 2 (people vs. technology)

Progress and challenges

# State of the Art

**It's easy to get 99% accuracy…**
**What that means depends on many factors…**

- **Common evaluations important**
- **Tasks become more challenging**
- **Word Error Rate (WER) < 10% is 'acceptable'**
- **Performance in field ~2x to 4x worse**

What was training set?

What was test set?

Were training and test independent?

Have other systems used same benchmark?

How large was the vocabulary and the sample size?

What speakers?

What style speech?

What kind of noise?

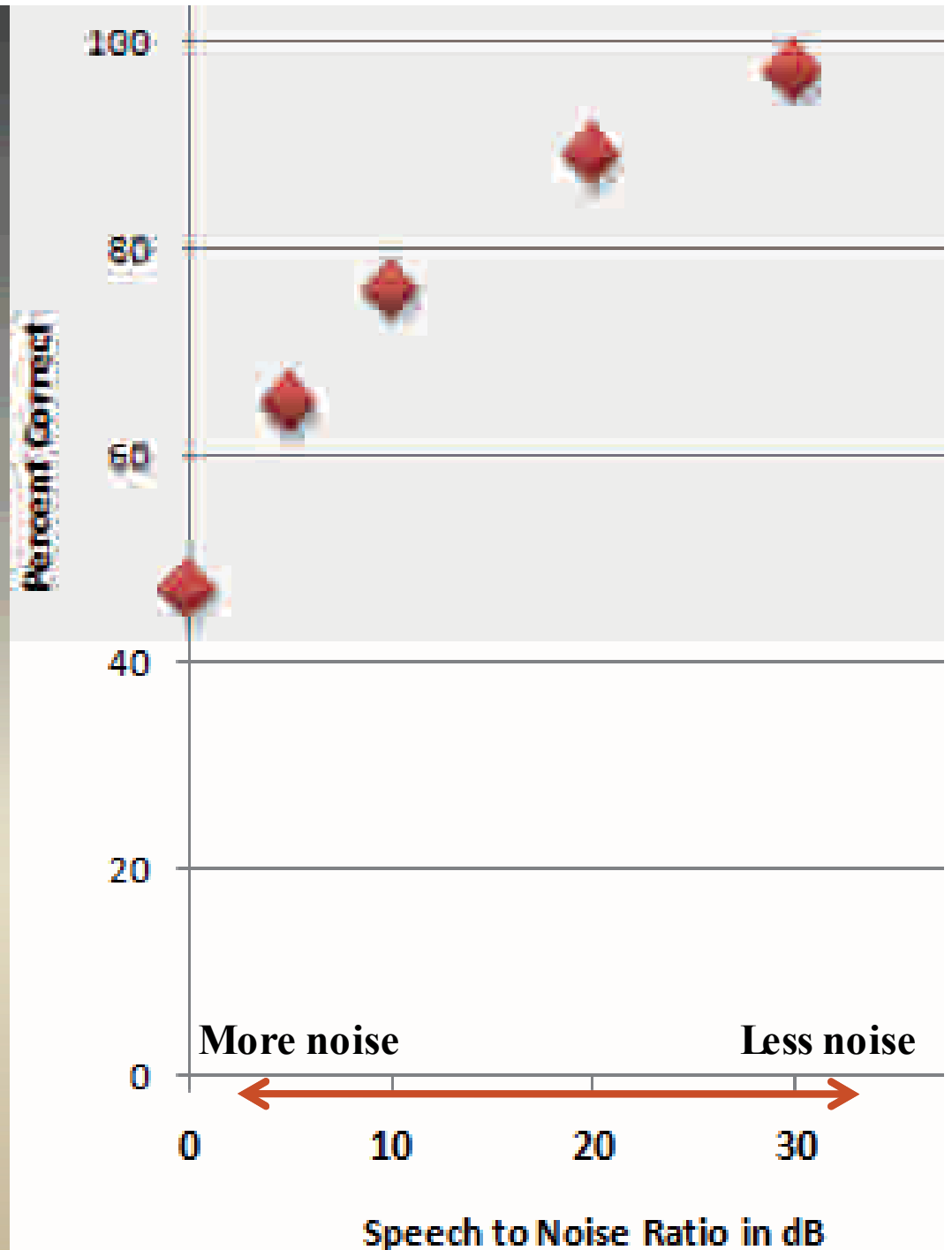From 2000 AAAS talk with Joe Picone, updated

# The Growing Impact of Speech Technology on Society

Patti Price, PPRICE Speech and Language Technology

- Intro: a few million years in about a minute
- Short review of speech recognition technology
- Speech recognition performance

**Social impact 1 (effect of society on speech)**

Social impact 2 (people vs. technology)

Progress and challenges

# Noise

We talk wherever we are, but noise degrades speech recognition, **especially speech-like noise**

Recognition of vowel-consonant-vowel consonants

Additive speech-shaped noise

Other noise also degrades speech recognition (speech, telephone channel, etc.)

The world is getting noisier

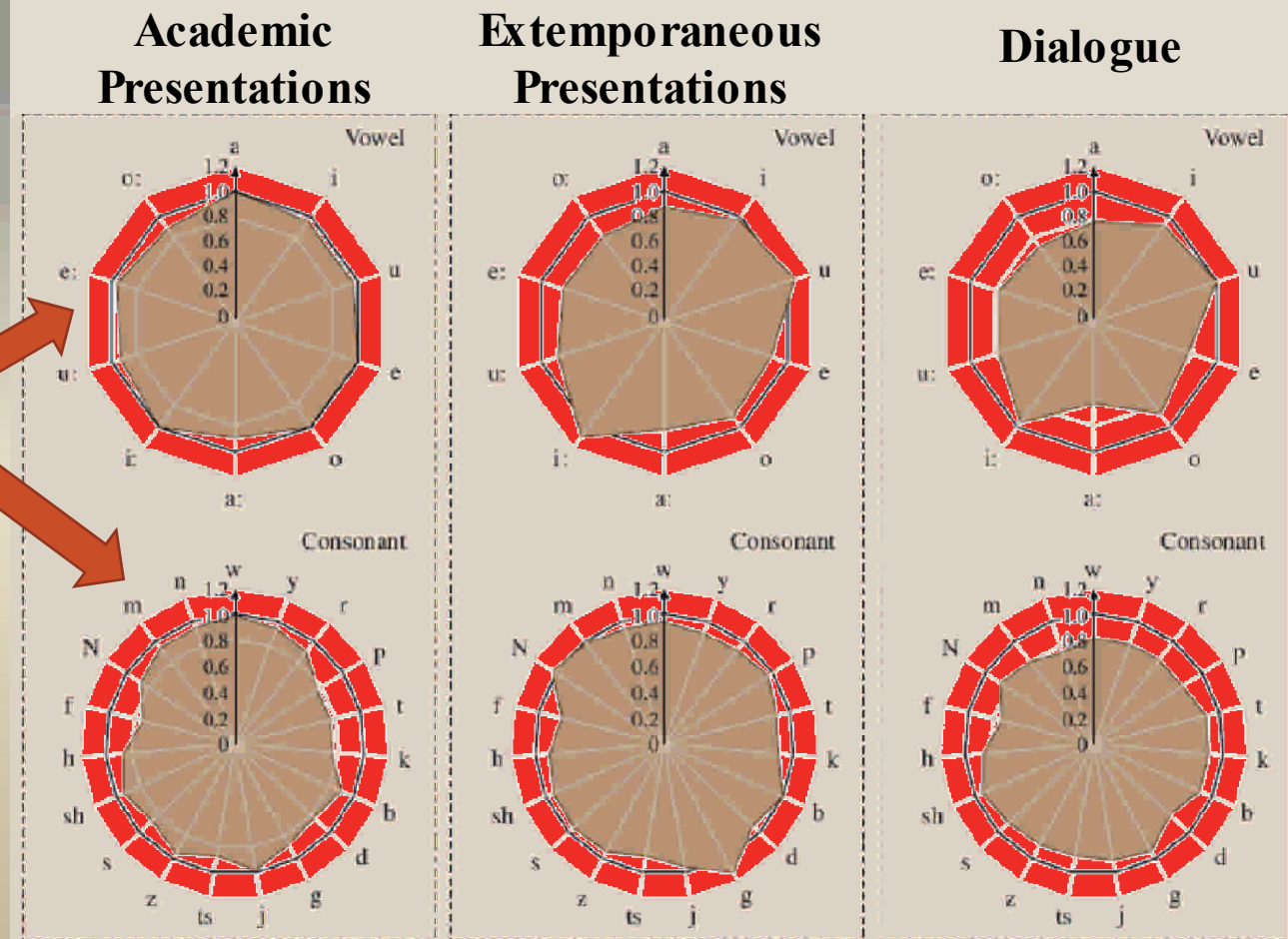Percent Correct

100
80
60
40
20
0

**More noise**       **Less noise**

0        10        20        30

Speech to Noise Ratio in dB

# Speech Style

**We don't always speak carefully…**

Orange outlines are read speech versions of originals

10 Japanese speakers

Beige inside portion is relative reduction for

vowels

and consonants

**Accuracy shrinks as reduction increases**



Academic Presentations

Extemporaneous Presentations

Dialogue

# Style Effects

**Non-speech**

**Filled pause**

**Reduction**

**Repetition**

*<laughter> uh I dunno I think in in in* in
today's increasingly global world it plays
a key role in *pinging breeple together*
bringing people together early on *in* in
their lives before they *sort of* embark
upon their *ca-* careers…

**Phrase
Correction**
(not in original)

**Repetition**

**False start**

**Hedge**

In today's increasingly global world it plays a key
role in bringing people together early on in their
lives before they embark upon their careers.
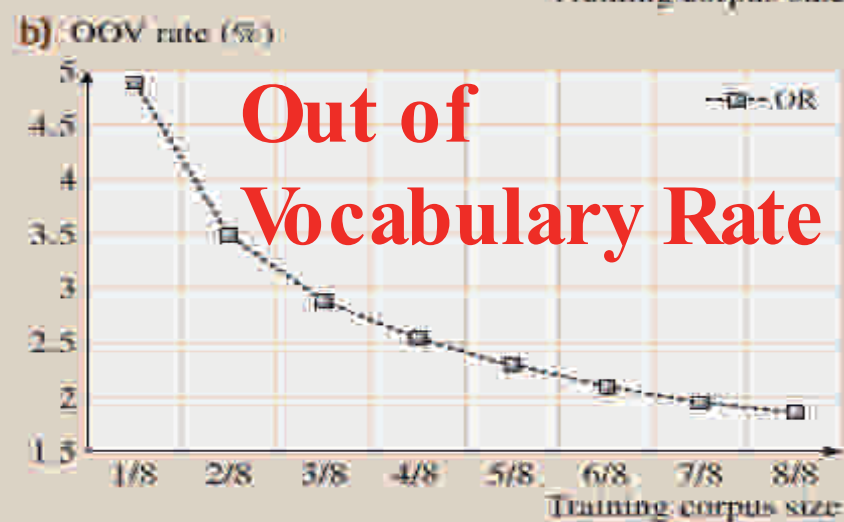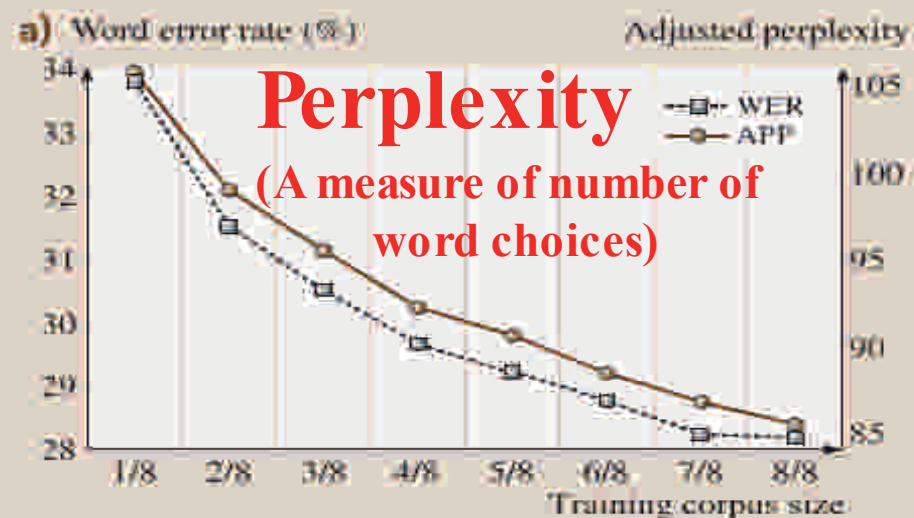
**Overlapping speech not illustrated**

**When disfluencies are
removed, spontaneous
speech had same
recognition error rates
as read speech.**
Butzberger et al. 1992

# Some Complexity Factors

**Perplexity**

**(A measure of number of word choices)**

**Out of Vocabulary Rate**

**Error rates rise with**
- a) **Less constraint on word choices**
- b) **More unknown words**

We wanted the knowledge navigator

We need to do better

Instead… we got
gethuman.com
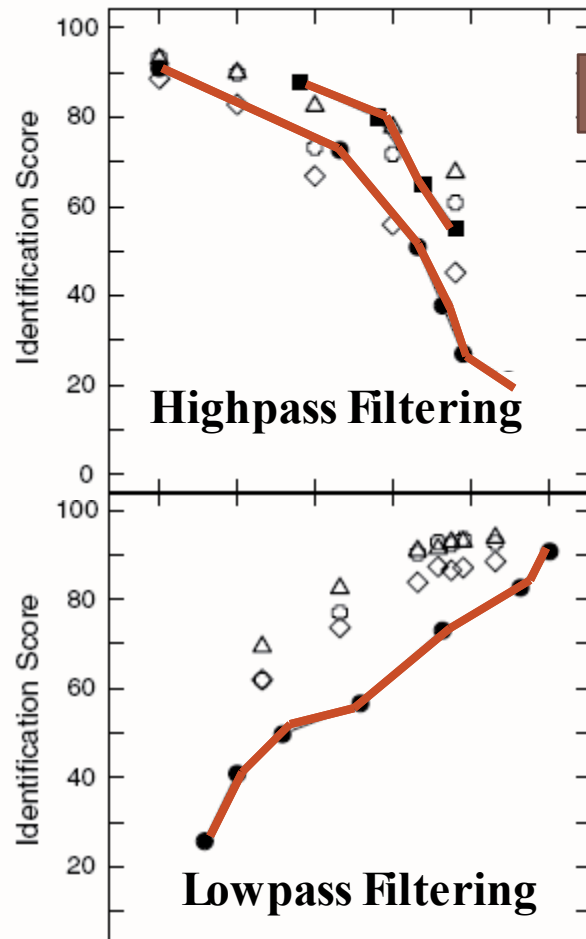And … Julie

From Apple 1987 visionary video

From Saturday Night Live, April 2006

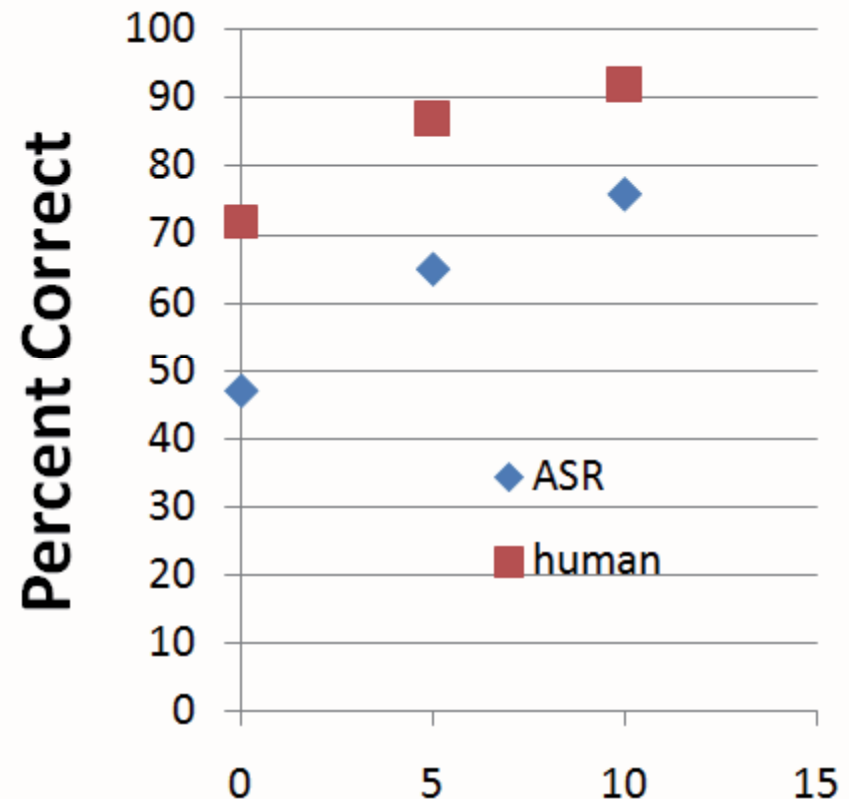# The Growing Impact of
# Speech Technology on Society

Patti Price, PPRICE Speech and Language Technology

- Intro: a few million years in about a minute
- Short review of speech recognition technology
- Speech recognition performance
- Social impact 1 (effect of society on speech)

**Social impact 2 (people vs. technology)**

Progress and challenges

# Human Recognition vs. ASR
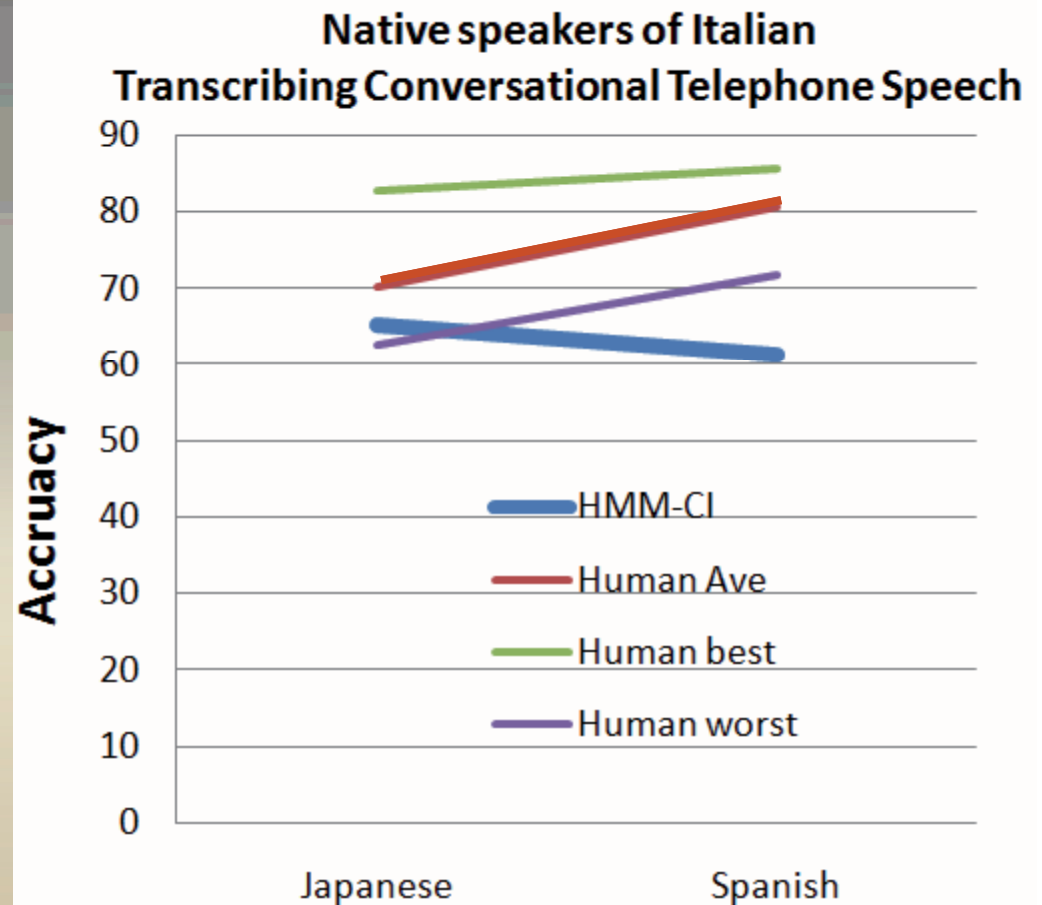
**Task: Label consonants in CVC and CV syllables**



**Highpass Filtering**

**Lowpass Filtering**

ASR
human

Percent Correct

**In most noise situations, humans are better than speech recognition but:**
- **About the same with high pass filtering**
- **ASR seems better with low pass filtering (data not same)**

# Human Recognition vs. ASR

Remove 'language model' but still use natural speech

Phonetic inventories are similar for Italian, Japanese and Spanish

Simple ASR is about the same as the worst of the 15 Italian transcribers

(Spanish and Italian are close in phonotactics)

**Native speakers of Italian**
**Transcribing Conversational Telephone Speech**

Accuracy (y-axis): 0, 10, 20, 30, 40, 50, 60, 70, 80, 90

Legend:
— HMM-CI
— Human Ave
— Human best
— Human worst

x-axis: Japanese, Spanish

Data from Shen et al., Interspeech 2008

# ASR is sometimes better than people…

Large vocabulary tasks where people may not know the vocabulary

   (e.g., thousands of names of companies in stock trading)

Small vocabulary tasks where memory plays a role

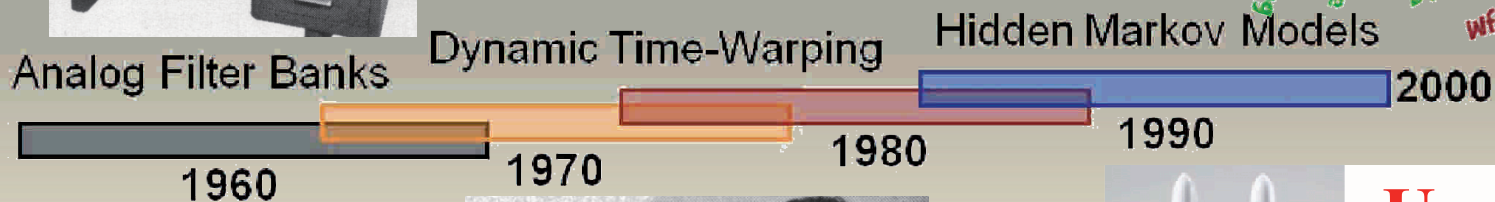   (e.g., transcribing sequence of 12 digits tracking numbers)

Artifacts of poorly designed experiments

   (e.g., testing on training data, correlational data that helps…)

**But generally people are more robust, flexible, adaptable… to situations that are normal human variability**

# The Growing Impact of Speech Technology on Society

Patti Price, PPRICE Speech and Language Technology

- Intro: a few million years in about a minute
- Short review of speech recognition technology
- Speech recognition performance
- Social impact 1 (effect of society on speech)
- Social impact 2 (people vs. technology)
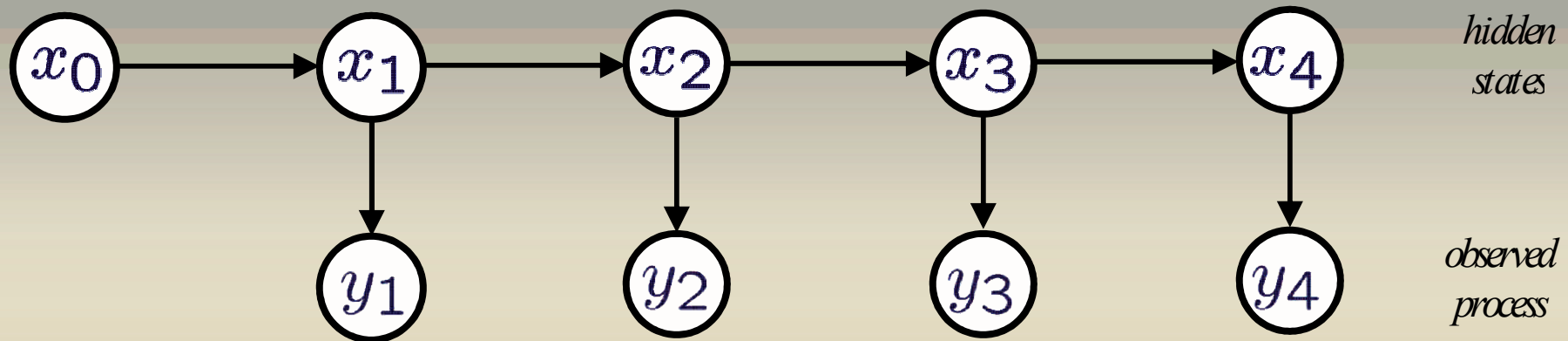- **Progress and challenges**

# Progress



Analog Filter Banks

Dynamic Time-Warping

Hidden Markov Models

1960     1970     1980     1990     2000

**Unsupervised training**
**New languages**
**User Interface**
**Distillation, NL**

# The Challenge of Hidden Markov Models

- Few realistic time series directly satisfy the assumptions of Markov processes:

*"Conditioned on the present, the past & future are independent"*
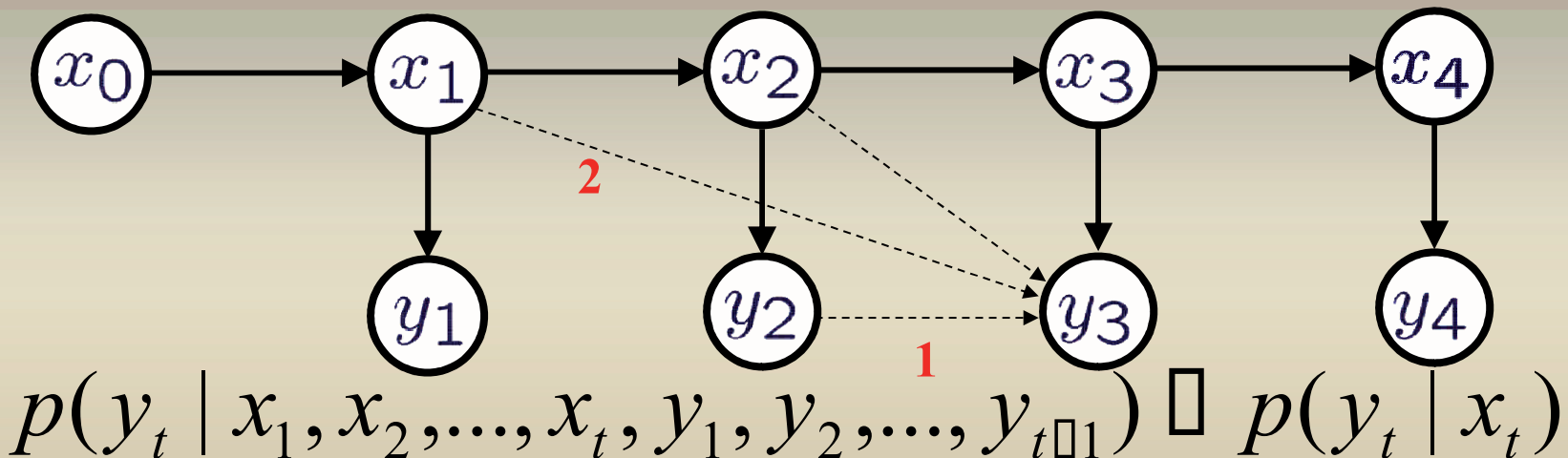


hidden states

observed process

# Conditional Independence Assumptions

1. Observations are conditionally independent given the HMM state

   We know this is false. It is a poor model.

2. Observations are conditionally independent of past states given the current state

   This is also false. We make up for it with trigrams, etc.



$$p(y_t \mid x_1, x_2, ..., x_t, y_1, y_2, ..., y_{t-1}) \qquad p(y_t \mid x_t)$$

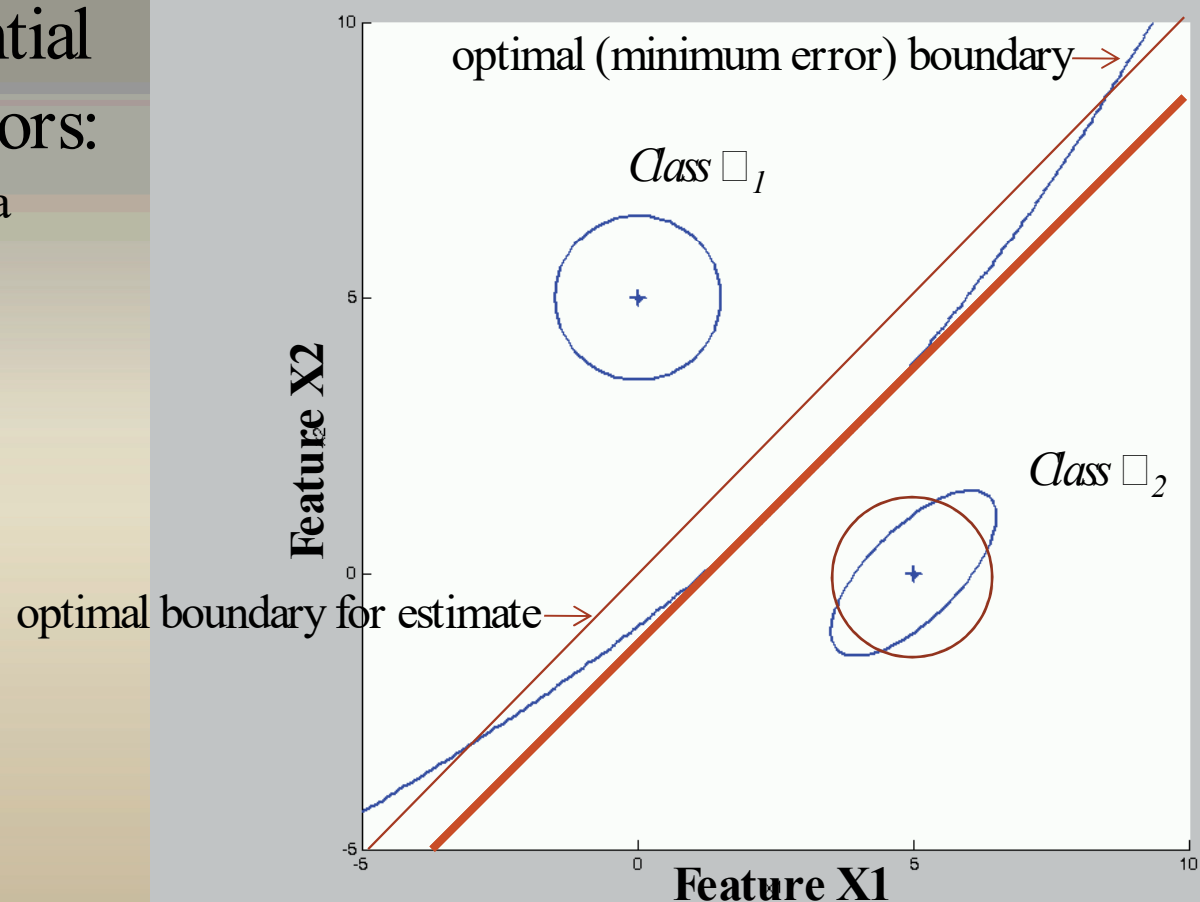**But that's just adjusting the boundary and not fixing the model...**

# Challenge: A Simple 2-Class Problem

Simple two-class, two-feature problem, with priors

$$p(\omega_1) = 1/3 \quad p(\omega_2) = 2/3$$

Using an exponential weight for the priors:

$$[p(\omega_1)]^a \quad [p(\omega_2)]^a$$



optimal (minimum error) boundary→

Class 1

Class 2

optimal boundary for estimate→

Feature X2

Feature X1

# Trends

Can our technology to help with information overload?

**150 years after Darwin's famous British A. A. S. debate**

• How do we get ourselves out of a niche as we learn more?

• Do we wait for a mutation? Wait for the old models to die?

• It's tough to start completely from scratch, but:

     • We can borrow useful mutations from others...

     • We can partner with our technology

---

**Limitations to speech technology arise from the evolution of speech as a social construct**

• Constrained by evolutionary history, production, perception, cognition
• Balancing needs of both speaker and hearer
• At the least, an existence proof, at best, a model we can improve on
• Speech technology lacks social skills; what do we do?

# 2009 Informal Survey

Sondra Ahlen

Fil Aleva

Francoise Beaufays

Joe Campbell

Rolf Carlson

Gerard Chollet

Mike Cohen

**Martin Cooke**

Deborah Dahl

**Vas Digalakis**

Farzad Ehsani

**Sadaoki Furui**

Juan Gilbert

John Makhoul

Bill Meisel

**Roger K. Moore**

Ariane Nabeth

Joe Picone

Alex Rudnicky

Paul Sawyer

Malcolm Slaney

Michel Stella

Gary Strong

Orith Toledo
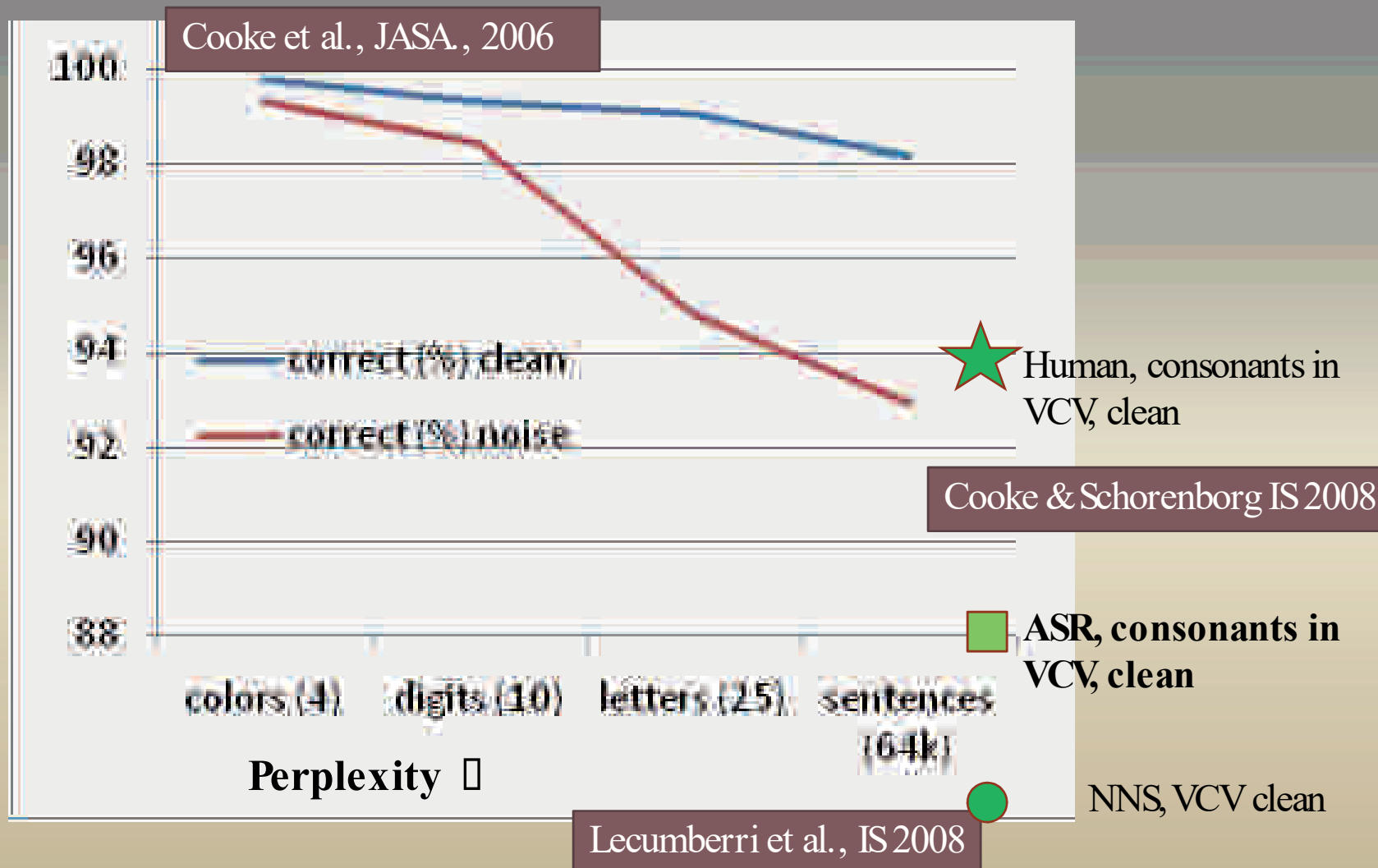
Carl Turner

Fuliang Weng

Steve Young

THANKS!!!

# Predictions Survey:

1997, 2003, 2009 speech conference attendees

Participants suggest year (or "never") when each statement might become true

Example: 1998, Kurzweil "By 2009 most routine business transactions take place between a human and a virtual personality (including an animated visual presence that looks like a human face)."

**In fact, in each sample the future gets farther away!**

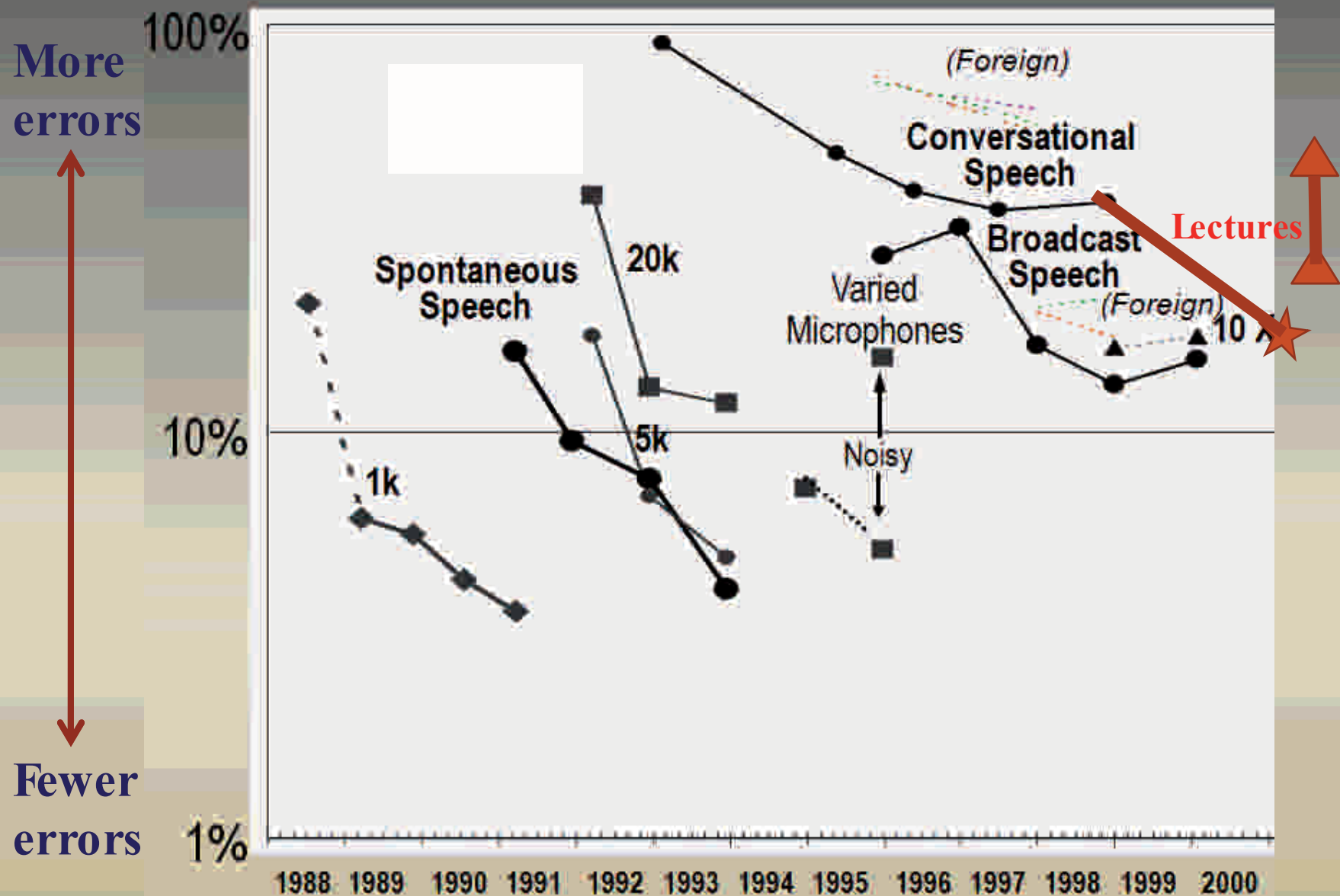Except for 'year when no speech research needed' (which is always "never")

# Human vs. ASR Recognition in Noise, Clean
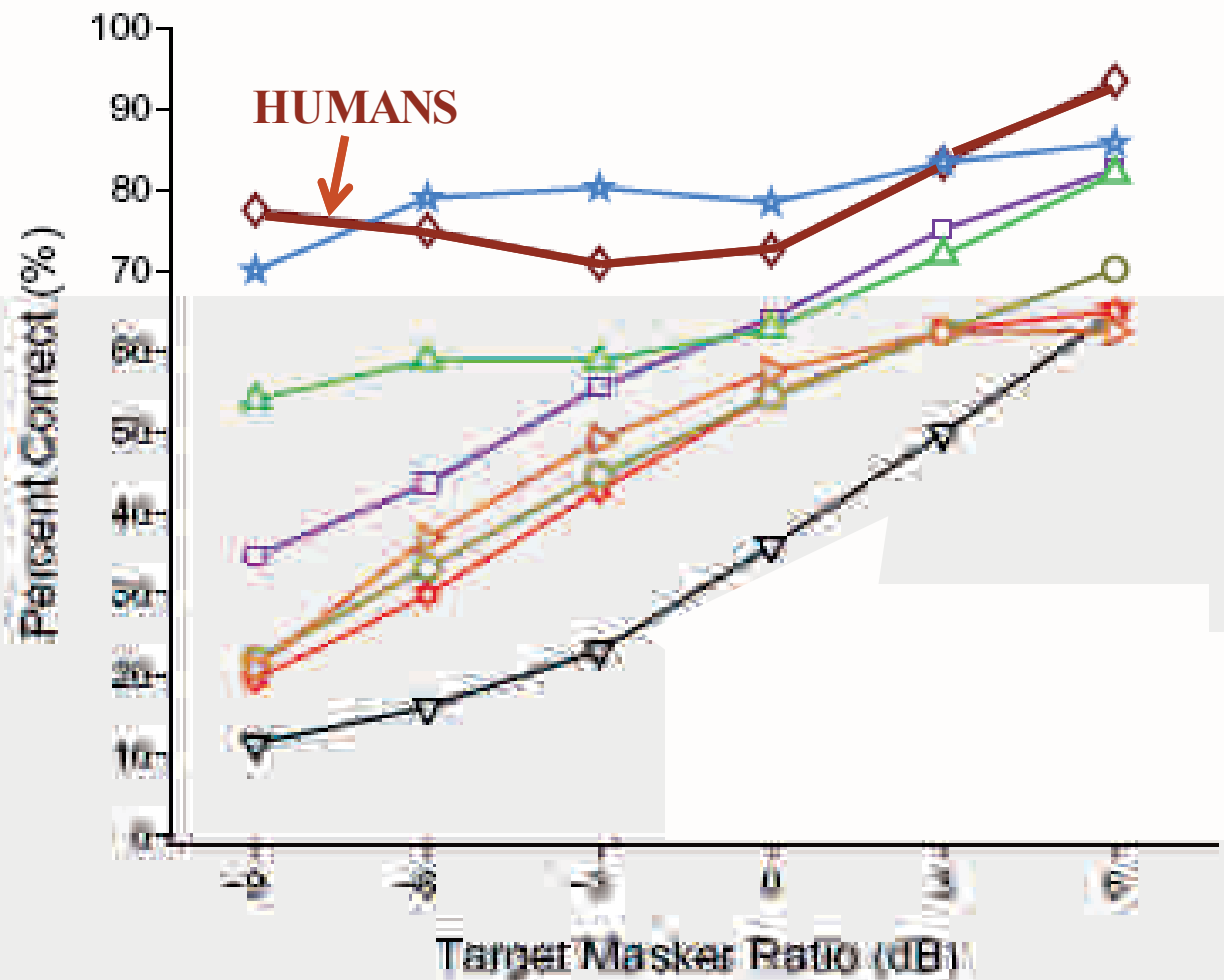
Noise is speech shaped noise at 6, 4 and 2 dB SNR



Cooke et al., JASA., 2006

correct (%) clean

correct (%) noise

colors (4)   digits (10)   letters (25)   sentences (64k)

**Perplexity**

Human, consonants in VCV, clean

Cooke & Schorenborg IS 2008

**ASR, consonants in VCV, clean**
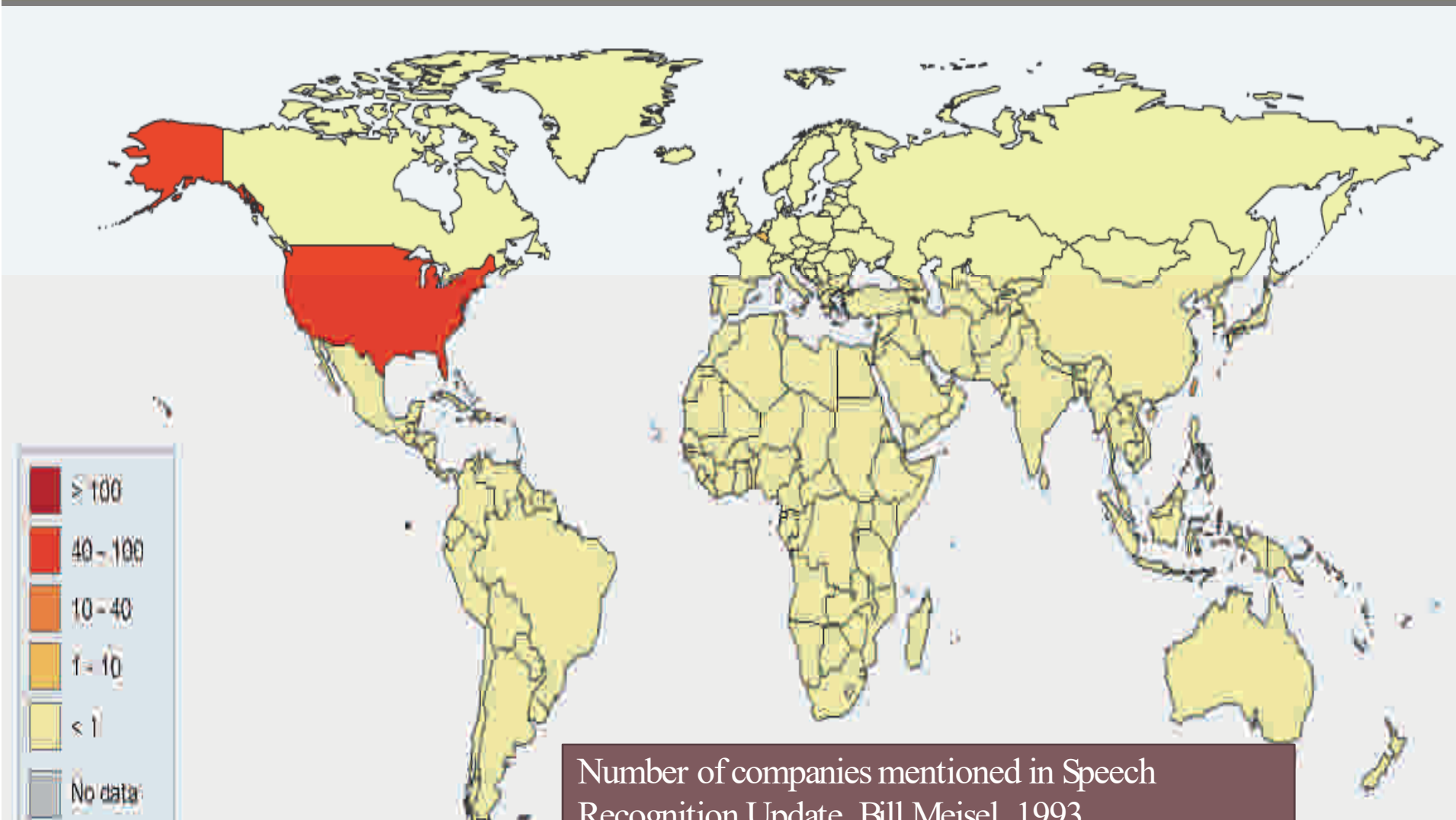
NNS, VCV clean

Lecumberri et al., IS 2008

# Error Rates

# Human Recognition vs. ASR

Task: Identify key words in sentences masked by similar sentences

Humans do fairly well, even when masker is louder target

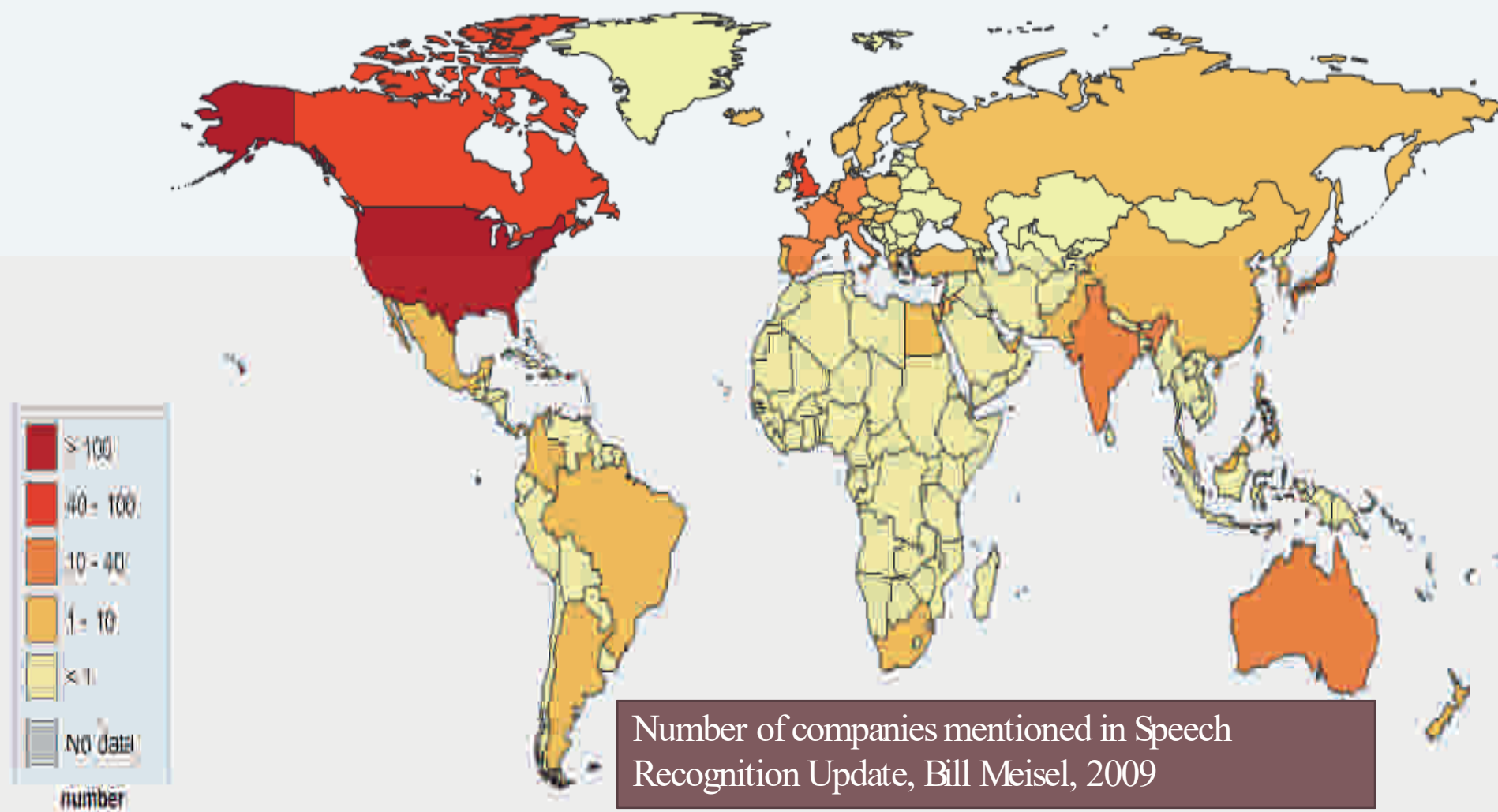Some systems took advantage of the fact that the target was always at a fixed level…

Cooke, Hershey and Rennie, Sp. Comm. 2010

HUMANS
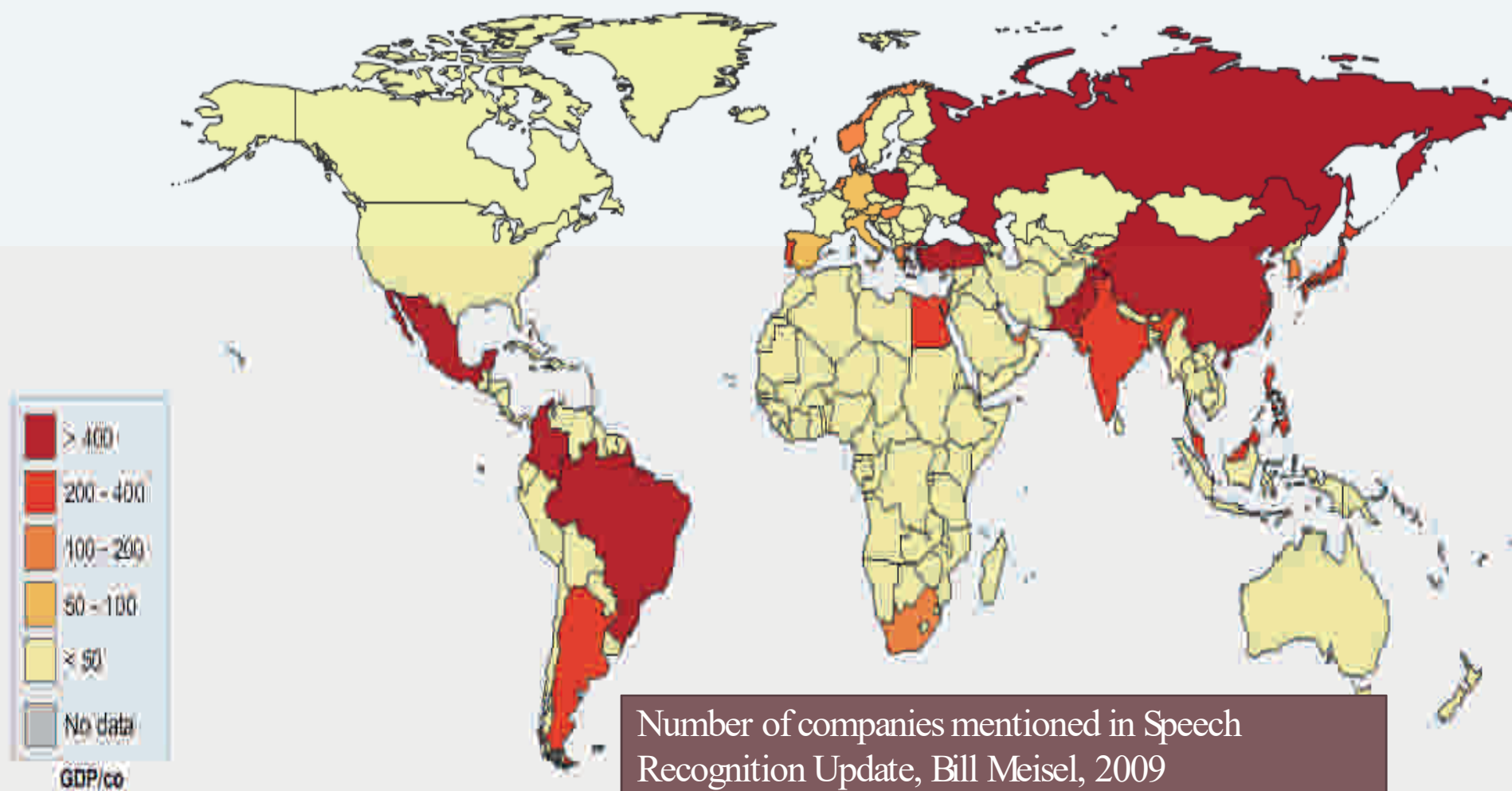
Percent Correct (%)

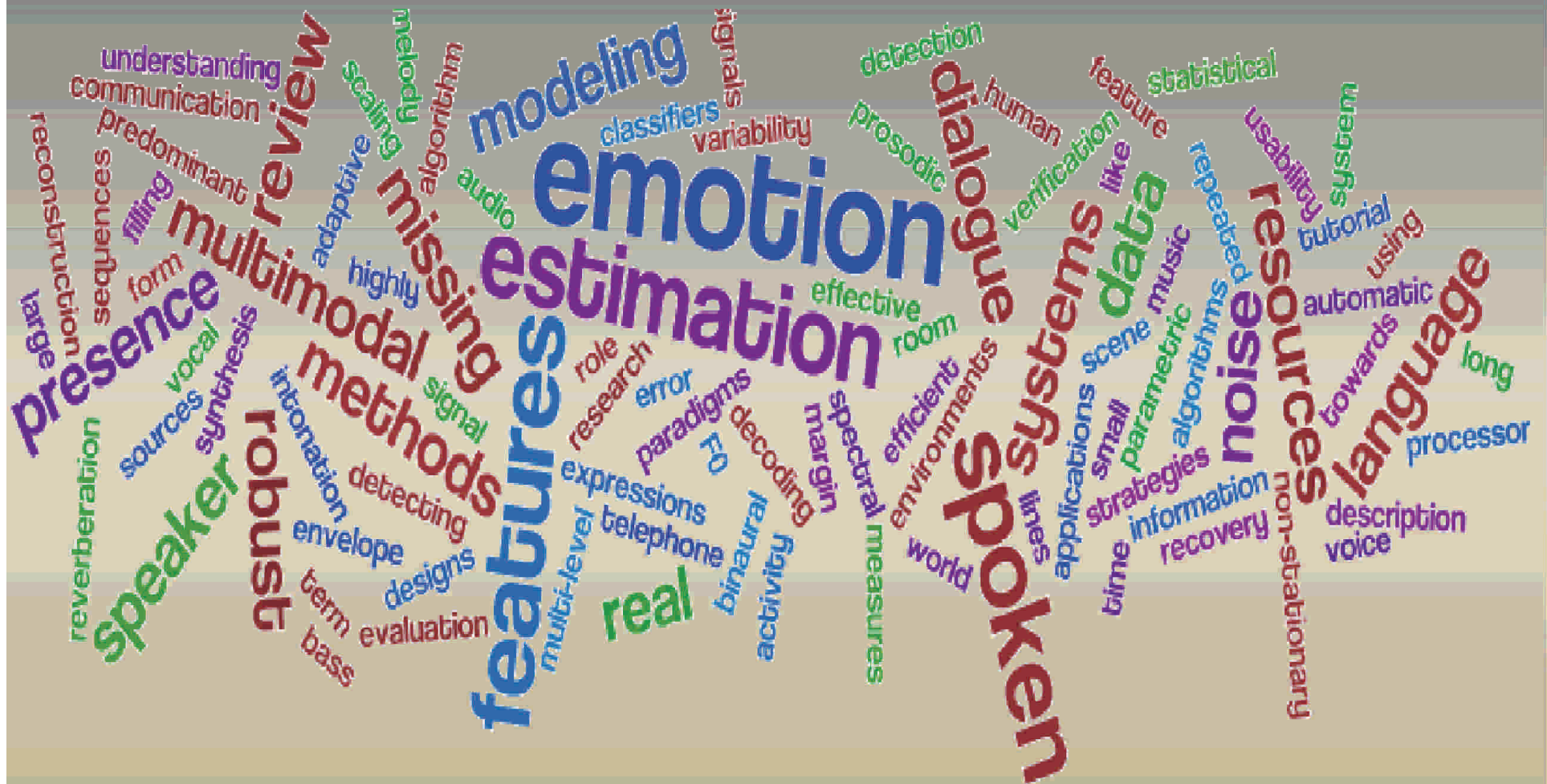100
90
80
70
60
50
40
30
20
10
0

Target Masker Ratio (dB)

# Companies 1993



Number of companies mentioned in Speech
Recognition Update, Bill Meisel, 1993

Companies 2009

Number of companies mentioned in Speech Recognition Update, Bill Meisel, 2009

# Companies 2009, GDP/ Company



Number of companies mentioned in Speech Recognition Update, Bill Meisel, 2009

# Challenges

# Most Cited, Most Downloaded

# Predictions