# Three Speaker Recognition Tasks
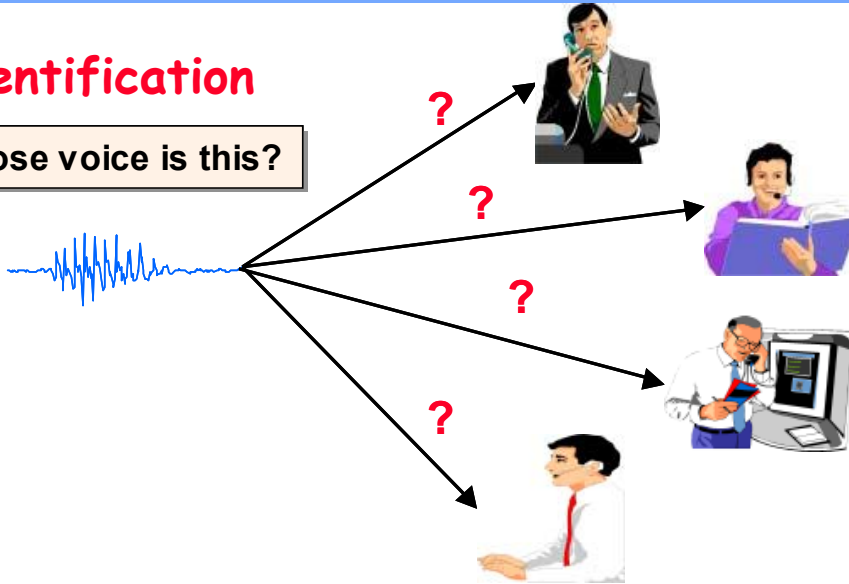
**Identification**

Whose voice is this?
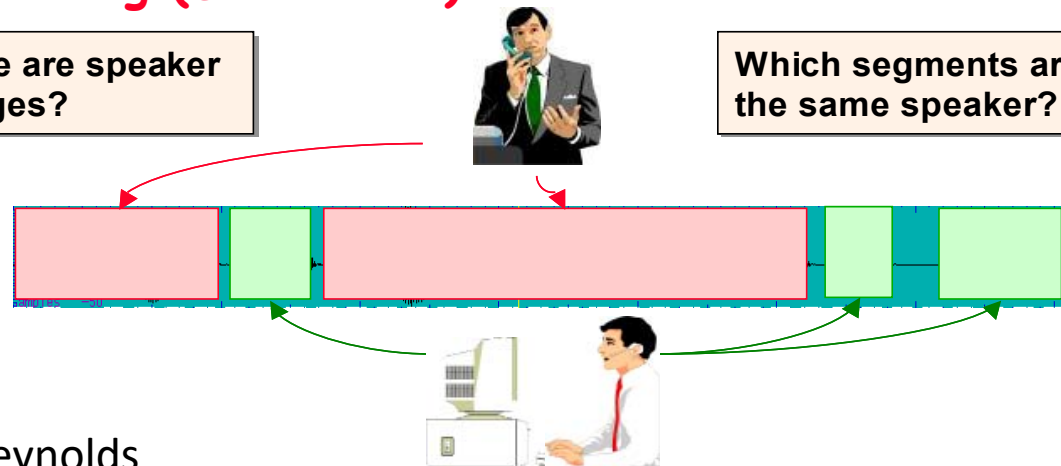
? ? ? ?

**Verification/Authentication/ Detection**

Is this Bob's voice?

?

**Segmentation and Clustering (Diarization)**

Where are speaker changes?

Which segments are from the same speaker?
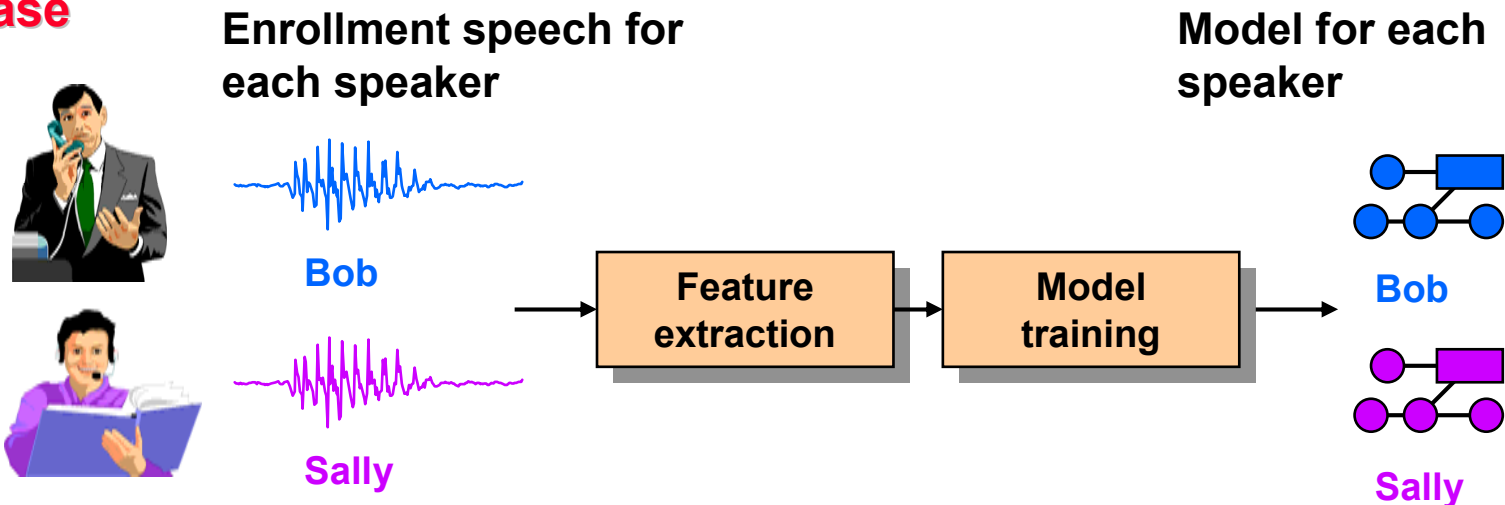
slide from Douglas Reynolds
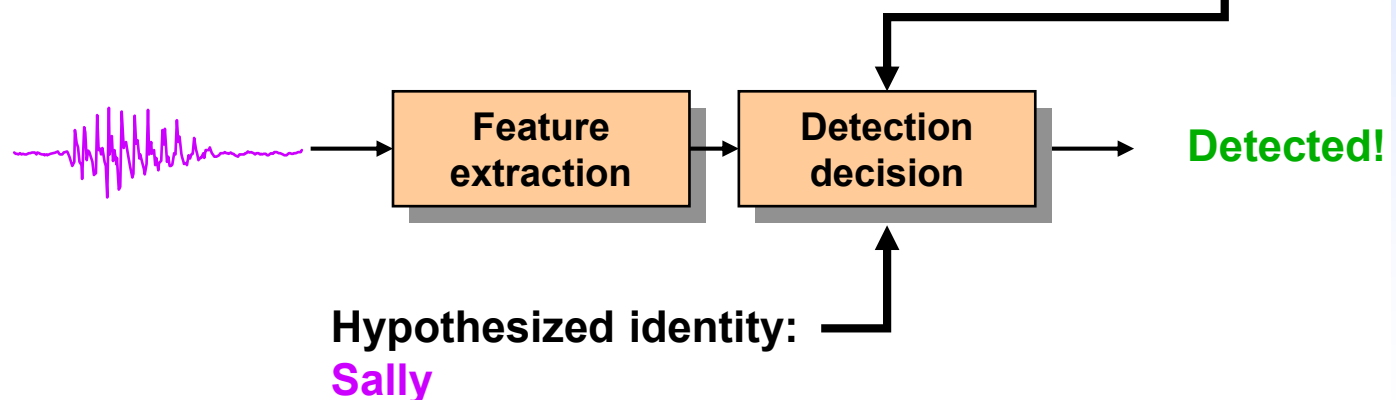
# Two kinds of speaker verification

- Text-dependent
  - Users have to say something specific
  - easier for system
- Text-independent
  - Users can say whatever they want
  - more flexible but harder

# Two phases to speaker detection



slide from Douglas Reynolds

# Detection: Likelihood Ratio

- Two-class hypothesis test:

  H0: X is **not** from the hypothesized speaker

  H1: X is from the hypothesized speaker

- **Choose the most likely hypothesis**

$$\Pr(H1 \mid X) \underset{<}{\overset{>}{\phantom{=}}} \Pr(H0 \mid X)$$

$$\frac{p(X \mid H1)\Pr(H1)}{p(X)} \underset{<}{\overset{>}{\phantom{=}}} \frac{p(X \mid H0)\Pr(H0)}{p(X)}$$

$$\frac{p(X \mid H1)}{p(X \mid H0)} \underset{<}{\overset{>}{\phantom{=}}} \frac{\Pr(H0)}{\Pr(H1)}$$
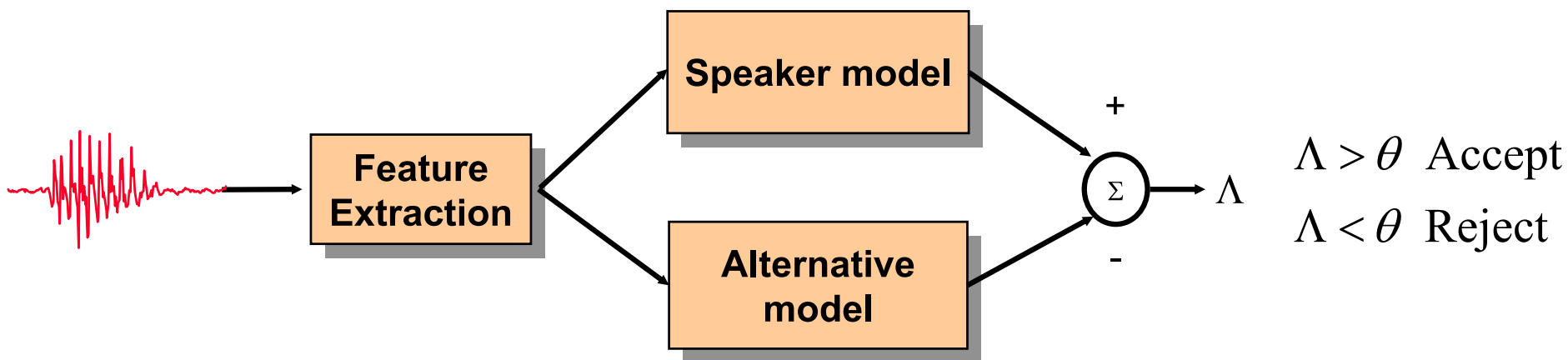
- **Likelihood ratio test:**

$$LR = \frac{p(X \mid H1)}{p(X \mid H0)} \qquad \begin{array}{l} LR > \theta \;\; \text{Accept } H1 \\ LR < \theta \;\; \text{Accept } H0 \end{array}$$

# Speaker ID
# Log-Likelihood Ratio Score

$$\text{LLR} = \Lambda = \log p(X|H1) - \log p(X|H0)$$



- Need *two* models
  - Hypothesized speaker model for H1
  - Alternative (background) model for H0
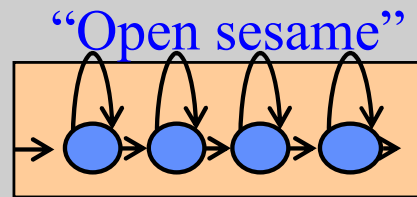
# How do we get H1?

- Pool speech from several speakers and train a single model:
  - a universal background model (UBM)
  - can train one UBM and use as H1 for all speakers
  - Should be trained using speech representing the expected impostor speech
  - Same type speech as speaker enrollment (modality, language, channel)
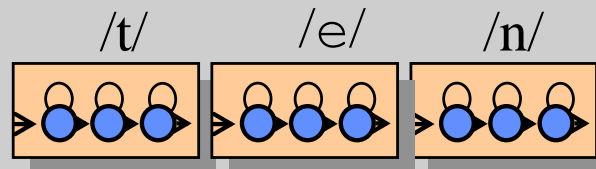
# How to compute P(H|X)?

- Gaussian Mixture Models (GMM)
  - The traditional best model for text-independent speaker recognition
- Support Vector Machines (SVM)
  - More recent use of discriminative model

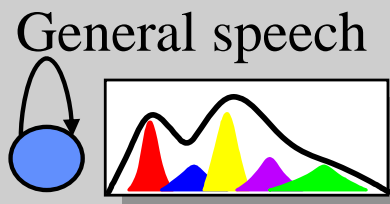# Form of GMM/HMM depends on application



**Fixed Phrase** ➡ Word/phrase models

"Open sesame"

**Prompted phrases/passwords** ➡ Phoneme models
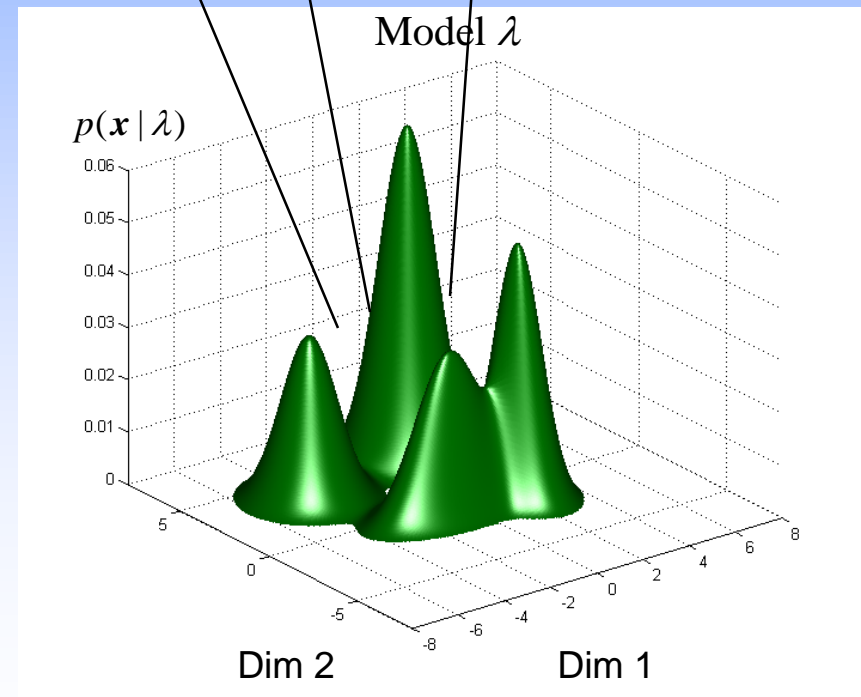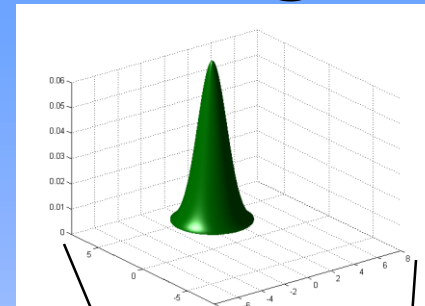
/t/   /e/   /n/

**Text-independent** ➡ single state HMM (GMM)

General speech
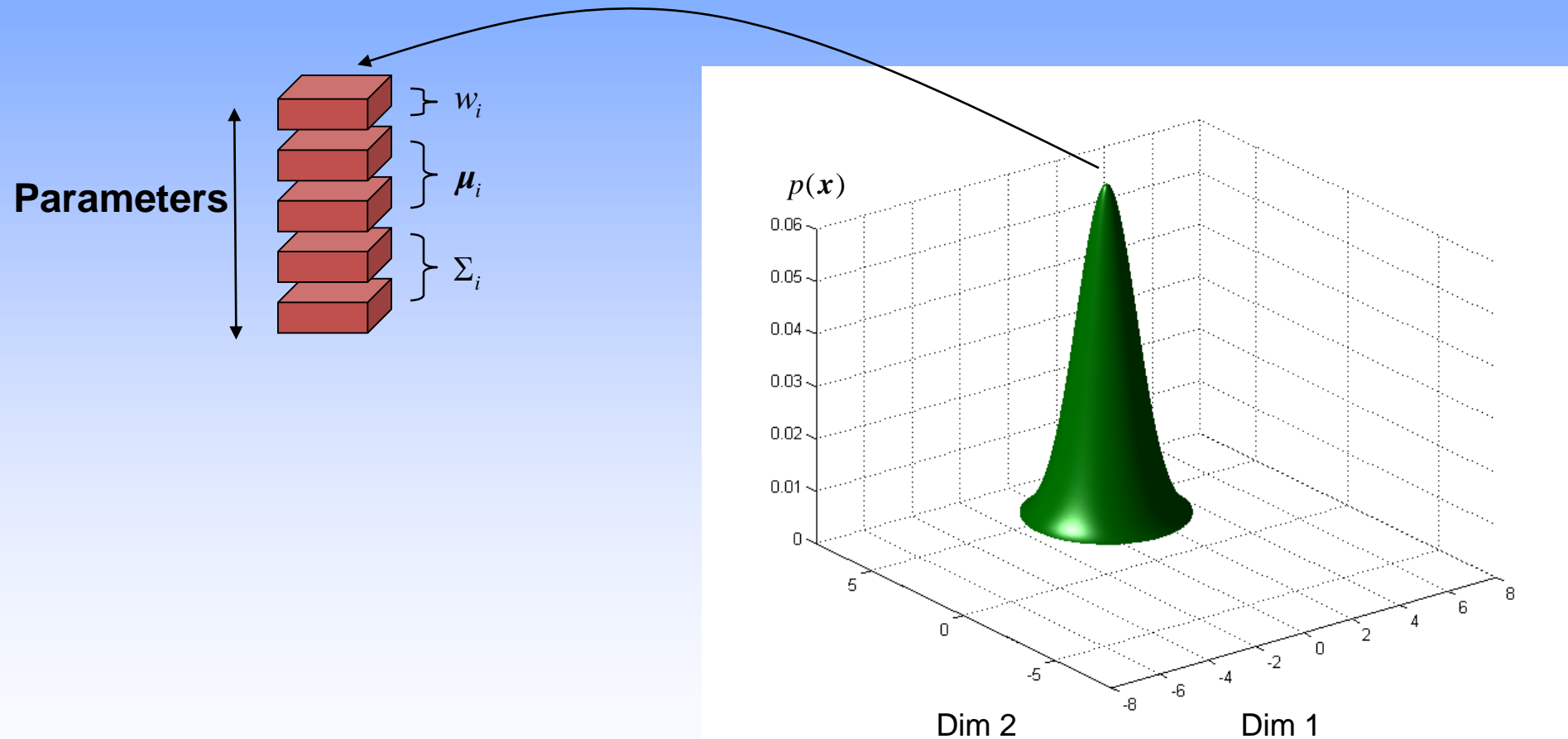
# GMMs for speaker recognition

- A Gaussian mixture model (GMM) represents features as the weighted sum of multiple Gaussian distributions

- Each Gaussian state i has a
  - Mean $\boldsymbol{\mu}_i$
  - Covariance $\Sigma_i$
  - Weight $w_i$

Model $\lambda$

$p(\boldsymbol{x}\,|\,\lambda)$

Dim 2     Dim 1

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# Recognition Systems
# Gaussian Mixture Models

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# Recognition Systems
# Gaussian Mixture Models



**Parameters**

**Model Components**

$p(\boldsymbol{x})$

Dim 2

Dim 1

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# GMM training

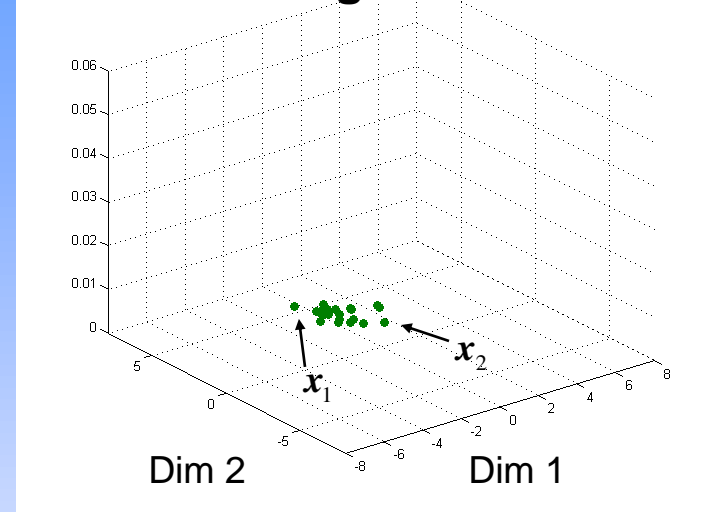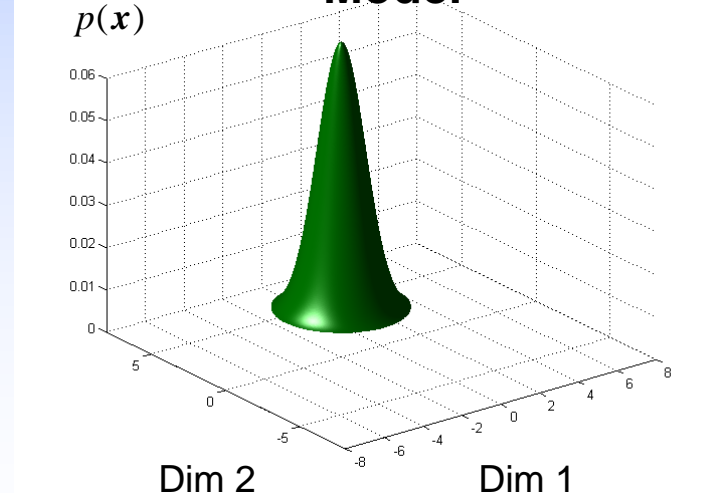- During training, the system learns about the data it uses to make decisions

  – A set of features are collected from a speaker (or language or dialect)
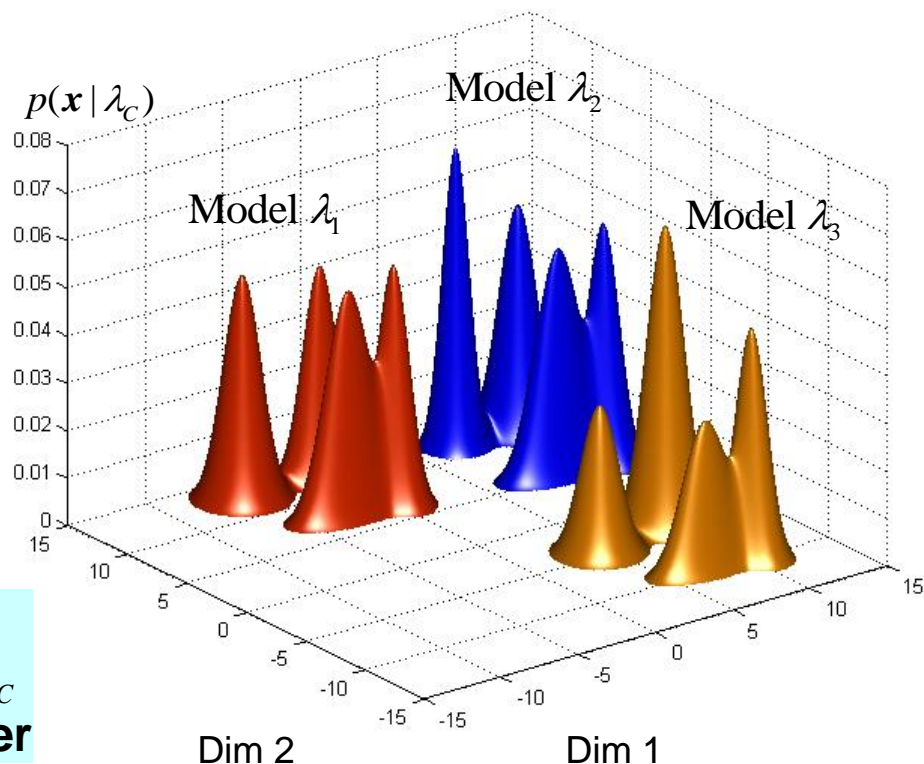


**Training Features**

Dim 2    Dim 1



**Model**

$p(\boldsymbol{x})$

Dim 2    Dim 1

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# Recognition Systems for Language, Dialect, Speaker ID

**Languages, Dialects, or Speakers**

**Parameters**

**Model Components**

$p(\boldsymbol{x}\,|\,\lambda_C)$

Model $\lambda_1$

Model $\lambda_2$

Model $\lambda_3$

Dim 2

Dim 1

**In LID, DID, and SID, we train a set of *target models* $\lambda_C$ for each dialect, language, or speaker**

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# Recognition Systems
# Universal Background Model

**Parameters**

**Model Components**

$p(\boldsymbol{x} \mid \lambda_{\bar{c}})$

Model $\lambda_{\bar{c}}$

Dim 2

Dim 1

**We also train a *universal background model* $\lambda_{\bar{c}}$ representing all speech**

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner
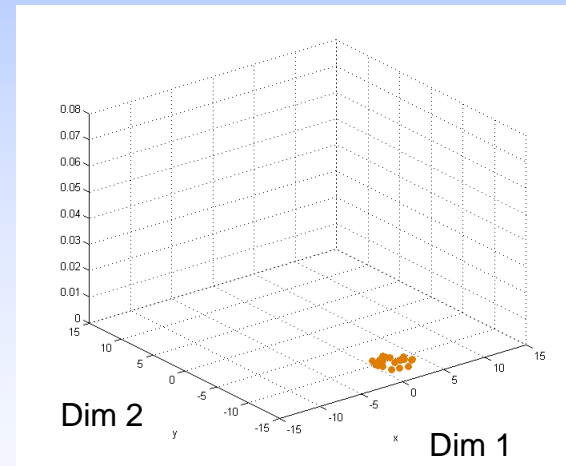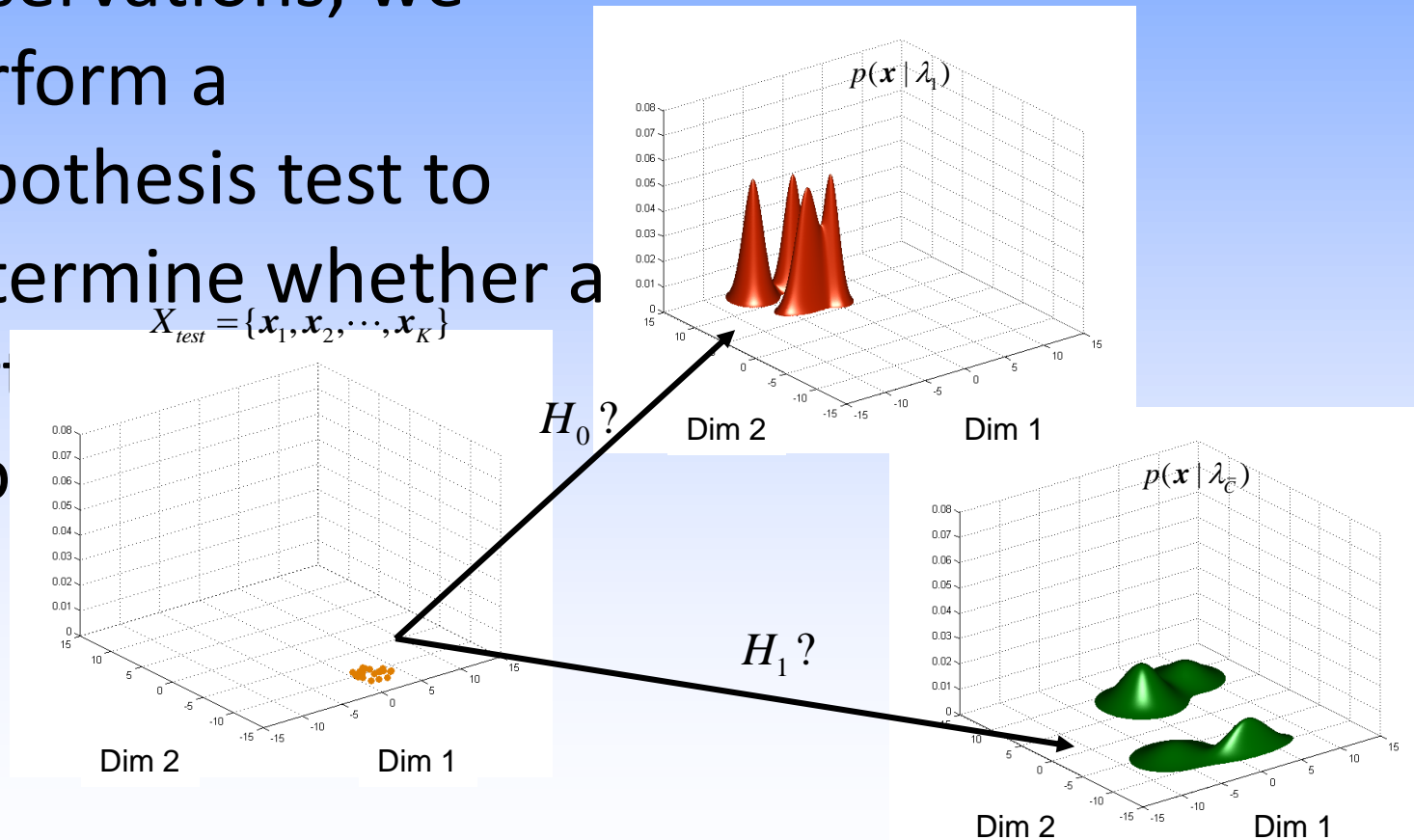
# Recognition Systems
# Hypothesis Test

- Given a set of test observations, we perform a hypothesis test to determine whether a certain class produced it

$H_0:$ $\quad X_{test}$ is from the hypothesized class

$H_1:$ $\quad X_{test}$ is not from the hypothesized class

$$X_{test} = \{ \boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_K \}$$



Dim 2

Dim 1

Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# Recognition Systems Hypothesis Test

$H_0:$   $X_{test}$ is from the hypothesized class
$H_1:$   $X_{test}$ is not from the hypothesized class

- Given a set of test observations, we perform a hypothesis test to determine whether a cert... pro...



$X_{test} = \{x_1, x_2, \cdots, x_K\}$

$p(x \mid \lambda_1)$

$H_0 ?$

Dim 2        Dim 1

$H_1 ?$

$p(x \mid \lambda_{\bar{C}})$

Dim 2        Dim 1

Dim 2        Dim 1

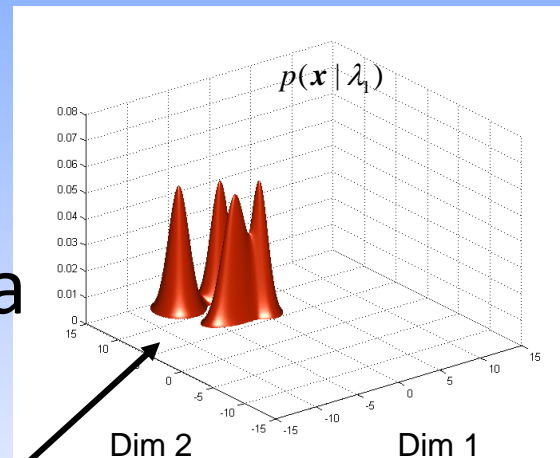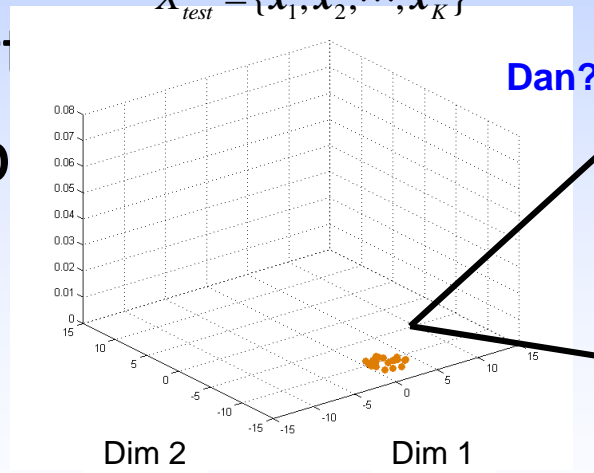Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner
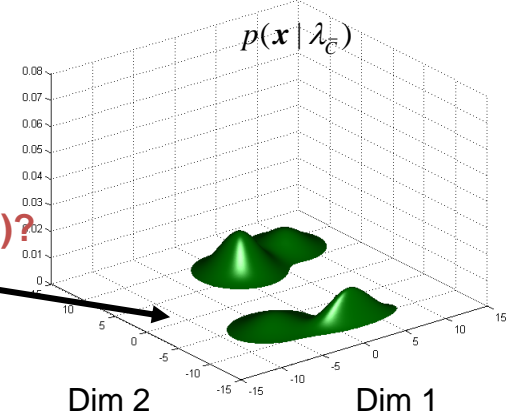
# Recognition Systems
# Hypothesis Test

- Given a set of test observations, we perform a hypothesis test to determine whether a cert~~~ pro~~~

$$X_{test} = \{x_1, x_2, \cdots, x_K\}$$

$p(x \mid \lambda_1)$

**Dan?**

Dim 2          Dim 1

$p(x \mid \lambda_{\bar{C}})$

**UBM (not Dan)?**

Dim 2          Dim 1

Dim 2          Dim 1

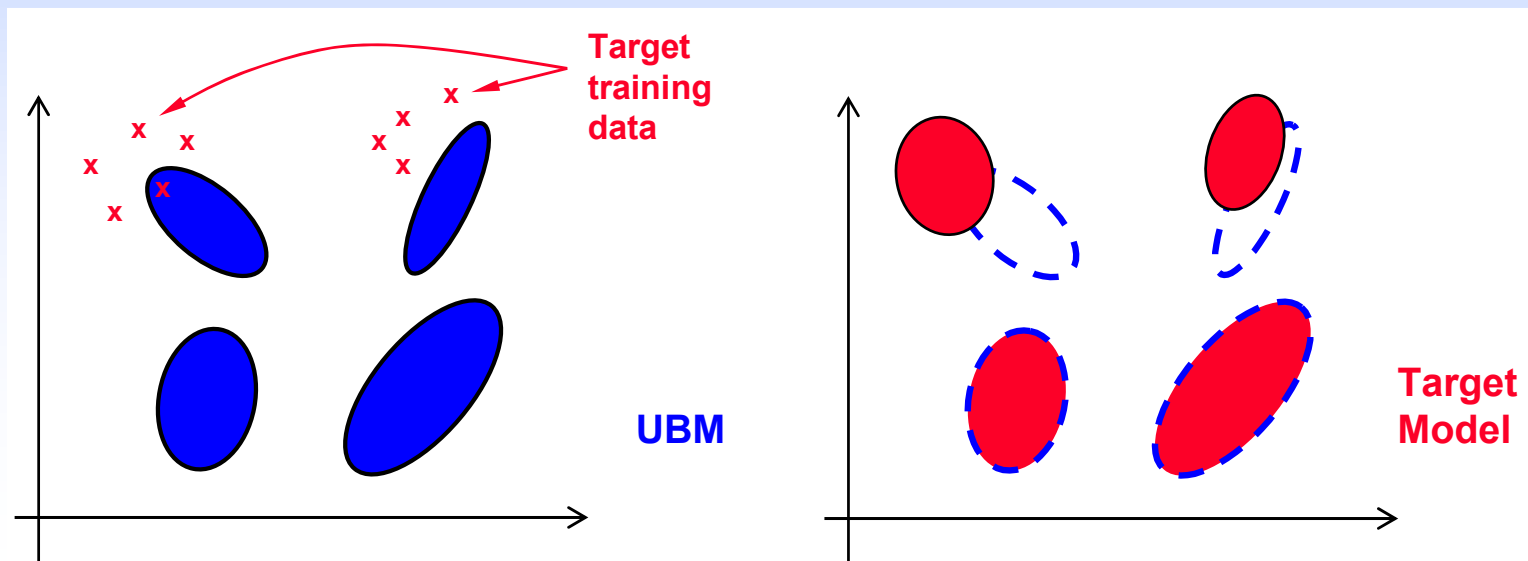Nicolas Malyska, Sanjeev Mohindra, Karen Lauro, Douglas Reynolds, and Jeremy Kepner

# More details on GMMs

- Instead of training speaker model on only speaker data

- Adapt the UBM to that speaker
  - takes advantage of all the data
  - MAP adaptation: new mean of each Gaussian is a weighted mix of the UBM and the speaker
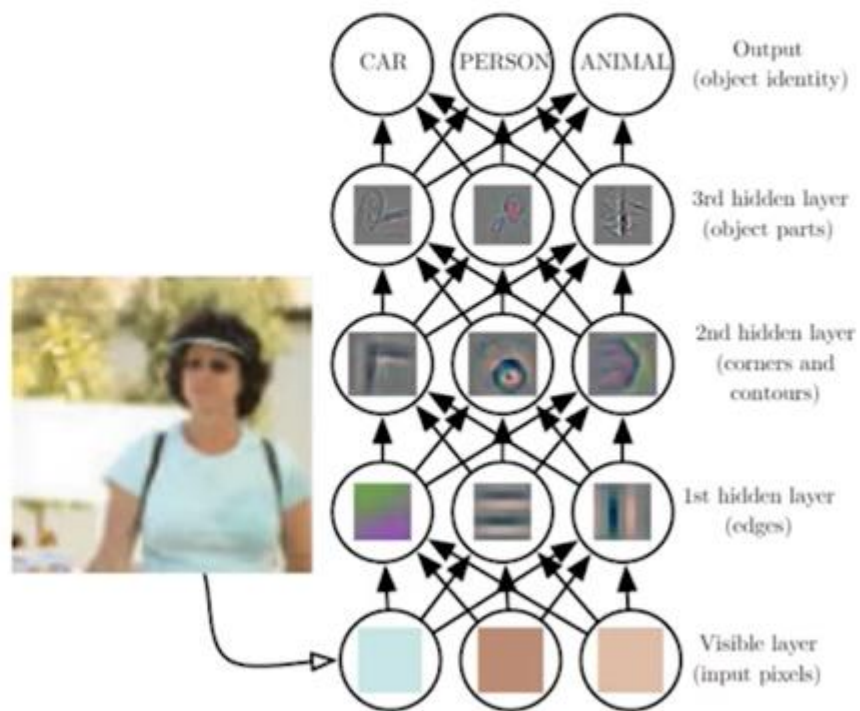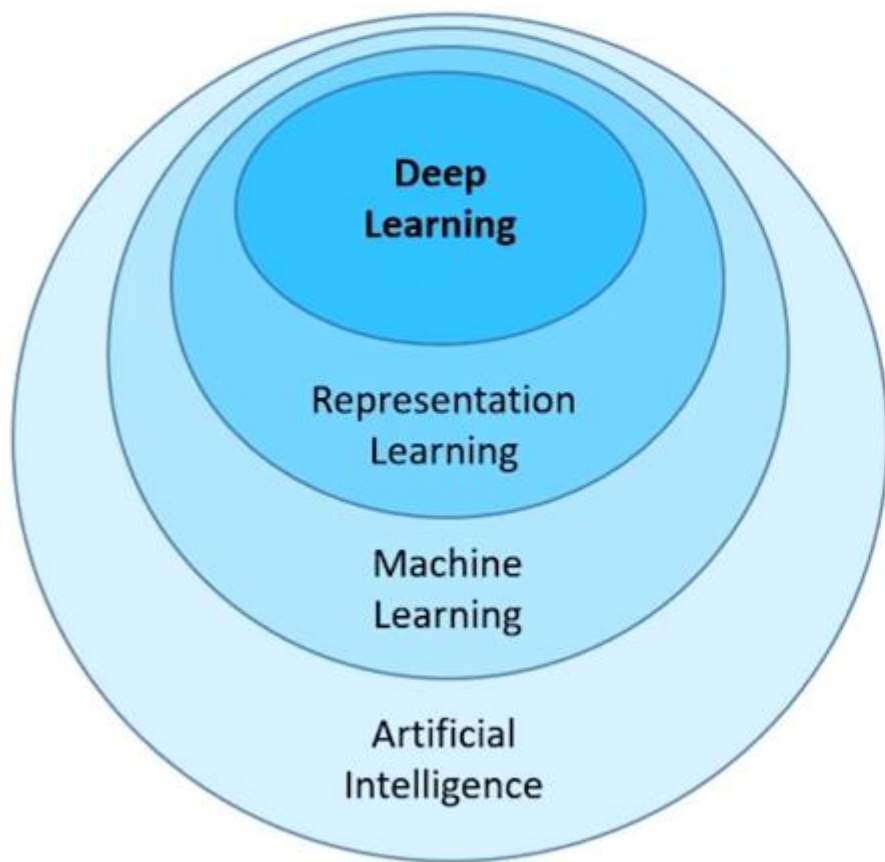
data:

**Target training data**

**UBM**

**Target Model**

# Gaussian mixture models

- Features are normal MFCC

  − can use more dimensions (20 + deltas)

- UBM background model: 512–2048 mixtures

- Speaker's GMM: 64–256 mixtures

- Often combined with other classifiers in mixture-of-experts

# Deep Learning is **Representation Learning**
## (aka Feature Learning)

# MIT Deep Learning Basics: Introduction and Overview

https://www.youtube.com/watch?v=O5xeyoRL95U