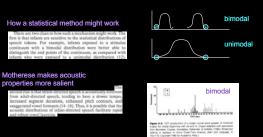**Slide 1**

# Psych 229:
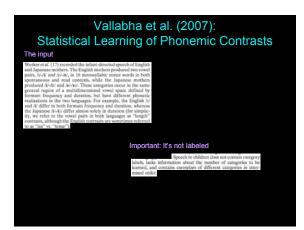# Language Acquisition

Lecture 11
Speech Perception

---

**Slide 2**

## Vallabha et al. (2007):
## Statistical Learning of Phonemic Contrasts

Testbed: Category emergence for English & Japanese vowel contrasts

Trajectory: 6 month olds have language-specific vowel distinctions

How a statistical method might work

bimodal

unimodal

There are two clues to how such a mechanism might work. The first is that infants are sensitive to the statistical distributions of speech tokens. For example, infants exposed to a stimulus continuum with a bimodal distribution were better able to distinguish the end points of the continuum, as compared with infants who were exposed to a unimodal distribution (12).

Motherese makes acoustic properties more salient

The second clue is that infant-directed speech is acoustically different from adult-directed speech, tending to have a slower tempo, increased segment durations, enhanced pitch contours, and exaggerated vowel formants (14–16). Thus, it is possible that the acoustic distributions of infant-directed speech facilitate rapid and robust vowel learning.

bimodal

---

**Slide 3**

## Vallabha et al. (2007):
## Statistical Learning of Phonemic Contrasts

The input

Werker et al. (17) recorded the infant-directed speech of English and Japanese mothers. The English mothers produced two vowel pairs, /i/-/ɪ/ and /ɛ/-/e/, in 16 monosyllabic nonce words in both spontaneous and read contexts, while the Japanese mothers produced /i/-/ii/ and /e/-/ee/. These categories occur in the same general region of a multidimensional vowel space defined by formant frequency and duration, but have different phonetic realizations in the two languages. For example, the English /i/ and /ɪ/ differ in both formant frequency and duration, whereas the Japanese /i/-/ii/ differ almost solely in duration (for simplicity, we refer to the vowel pairs in both languages as "length" contrasts, although the English contrasts are sometimes referred to as "lax" vs. "tense").

Important: It's not labeled

Speech to children does not contain category labels, lacks information about the number of categories to be learned, and contains exemplars of different categories in inter-mixed order.

---

**Slide 4**

## Vallabha et al. (2007):
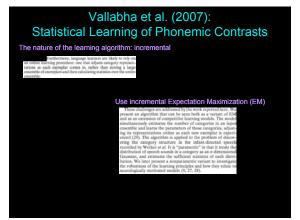## Statistical Learning of Phonemic Contrasts

A quick look at formants (F1, F2)

F1: depends on whether the sound is more open or closed. (Varies along y axis.) F1 increases as the vowel becomes more open and decreases as vowel closes.

F2: depends on whether the sound is made in the front or the back of the vocal cavity. (Varies along X axis). F2 increases the more forward the sound is.

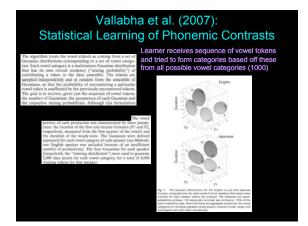Idea: As long as speakers use the same values for these formants, they will produce the same vowel.

---

**Slide 5**

## Vallabha et al. (2007):
## Statistical Learning of Phonemic Contrasts

The nature of the learning algorithm: incremental

Furthermore, language learners are likely to rely on an online learning procedure: one that adjusts category representations as each exemplar comes in, rather than storing a large ensemble of exemplars and then calculating statistics over the entire ensemble.

Use incremental Expectation Maximization (EM)

These challenges are addressed by the work reported here. We present an algorithm that can be seen both as a variant of EM and as an extension of competitive learning models. The model simultaneously estimates the number of categories in an input ensemble and learns the parameters of those categories, adjusting its representations online as each new exemplar is experienced (24). The algorithm is applied to the problem of discovering the category structure in the infant-directed speech recorded by Werker et al. It is "parametric" in that it treats the distribution of speech sounds in a category as an n-dimensional Gaussian, and estimates the sufficient statistics of each distribution. We later present a nonparametric variant to investigate the robustness of the learning principles and how they relate to neurologically motivated models (9, 27, 28).

---

**Slide 6**

## Vallabha et al. (2007):
## Statistical Learning of Phonemic Contrasts

A brief look at Expectation Maximization

Used for finding the maximum likelihood estimates of parameters in probabilistic models
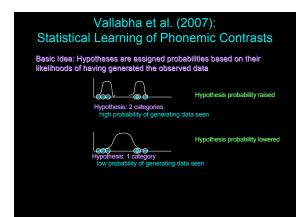
There are unknown (latent) variables in the model.

Algorithm alternates between doing an expectation step, which computes the expectation of the likelihood by using the latent variables, and a maximization step which computes the maximum likelihood estimates using the expected likelihood found in the previous step. It can then go back to the expectation step, using the results of the maximization step.
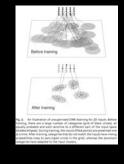
Learner receives sequence of vowel tokens and tried to form categories based off these from all possible vowel categories (1000)

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

For each token
- "responsibility" of each potential category is calculated
- more responsible categories get larger updates to their means & covariances
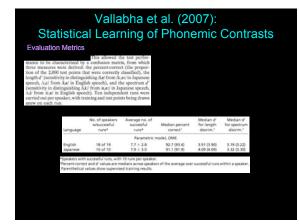- mixing probability (measure of success) of most "responsible" category [estimated] is updated a small amount

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Basic Idea: Hypotheses are assigned probabilities based on their likelihoods of having generated the observed data

Hypothesis probability raised

Hypothesis: 2 categories
high probability of generating data seen

Hypothesis probability lowered

Hypothesis: 1 category
low probability of generating data seen

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Before training

After training

- 50,000 data points to train on
- 2,000 data points tested on

Measure of Success

Each test point was classified with the category that had the greatest likelihood for that point. The run was considered "successful" if 95% of the test points were classified into four categories. For evaluation purposes, the categories were also assigned labels (e.g., the category to which most of the /i/ tokens were classified was labeled /i/).

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Evaluation Metrics

| Language | No. of speakers w/successful runs* | Average no. of successful runs* | Median percent correct† | Median d' for length discrim.‡ | Median d' for spectrum discrim.‡ |
|---|---|---|---|---|---|
| | | Parametric model, OME | | | |
| English | 18 of 19 | 7.7 ± 2.8 | 92.7 (93.4) | 3.91 (3.90) | 3.19 (3.22) |
| Japanese | 10 of 10 | 7.9 ± 3.0 | 91.1 (91.9) | 4.09 (4.09) | 3.32 (3.30) |

*Speakers with successful runs, with 10 runs per speaker.
†Percent-correct and d' values are medians across speakers of the average over successful runs within a speaker.
Parenthetical values show supervised training results.

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

What about inter-speaker variation within the same language?
Does that affect the categorization ability?

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Assumption of the OME mode: categories have Gaussian distribution

A model that doesn't do this: TOME



Hebbian learning: neural network, "cells that fire together wire together" - building associations

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

TOME process

TOME results



---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Success?



Discussion: Do we want perfect performance in these models, or do we want flawed performance since infants must go through stages of learning?

Relation to vowel category acquisition

A note on the implementational level

---

## Vallabha et al. (2007): Statistical Learning of Phonemic Contrasts

Now, back to speech acquisition - domain-specific vs. domain-general?

Gaussian distribution assumption = domain-general bias?

How important is biological plausibility in the learning algorithm?