

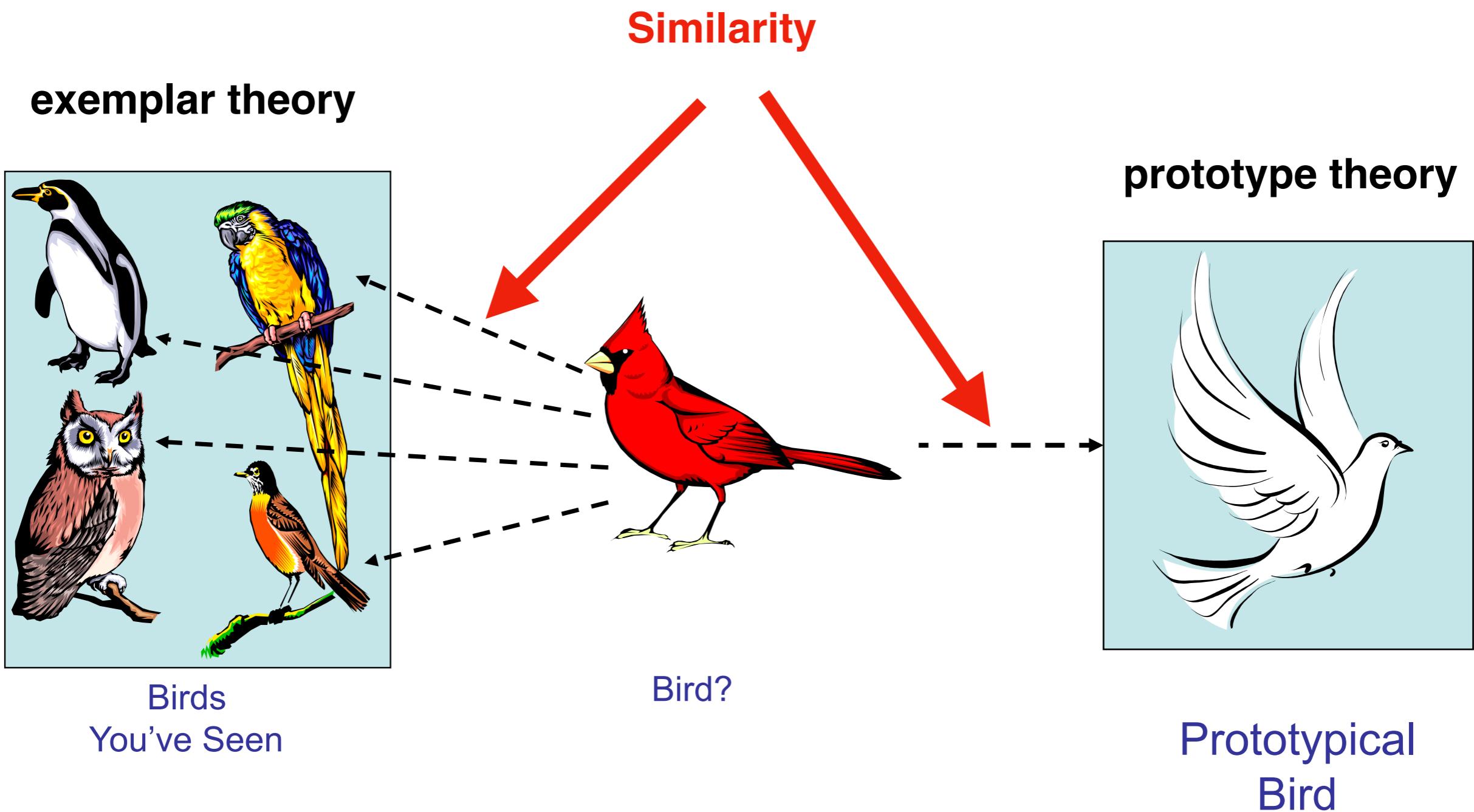
Categories and Concepts - Spring 2019

How categories influence perception

Brenden Lake

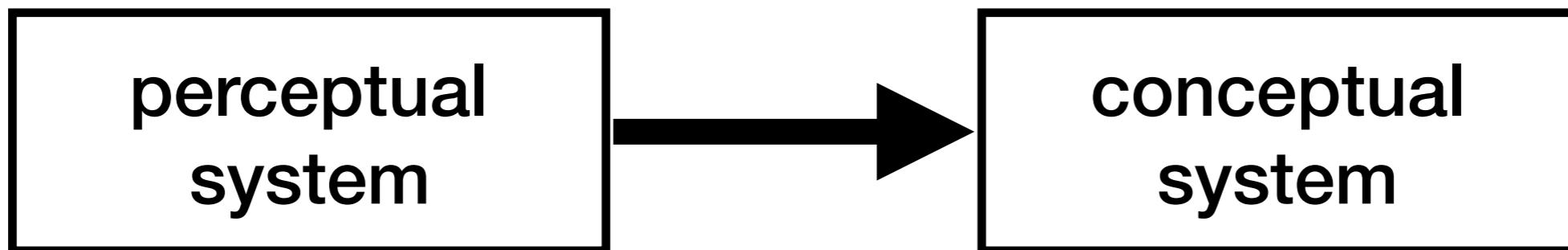
PSYCH-GA 2207

Traditional view: similarity → categorization

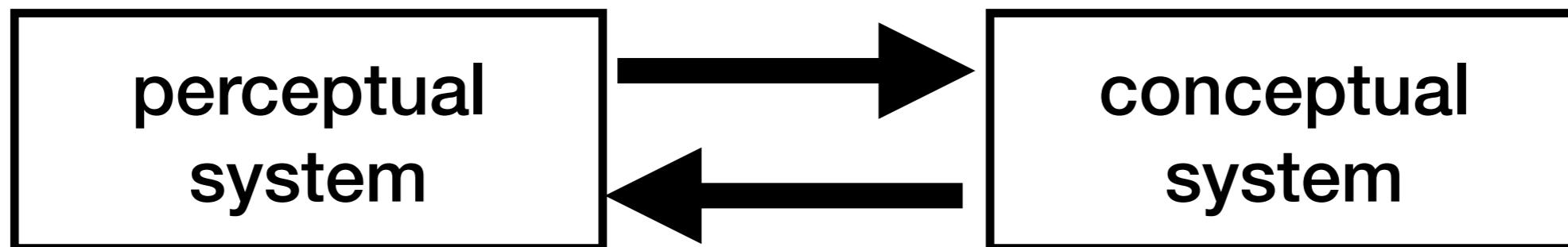


But does categorization also determine similarity?

Traditional view

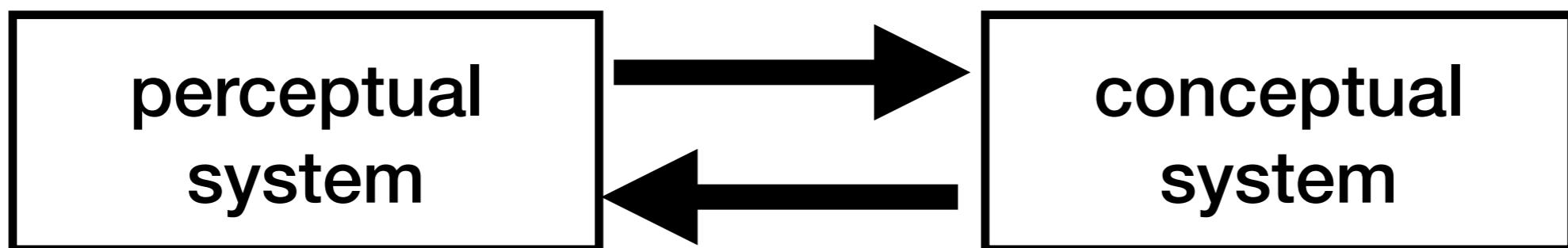


Interaction view



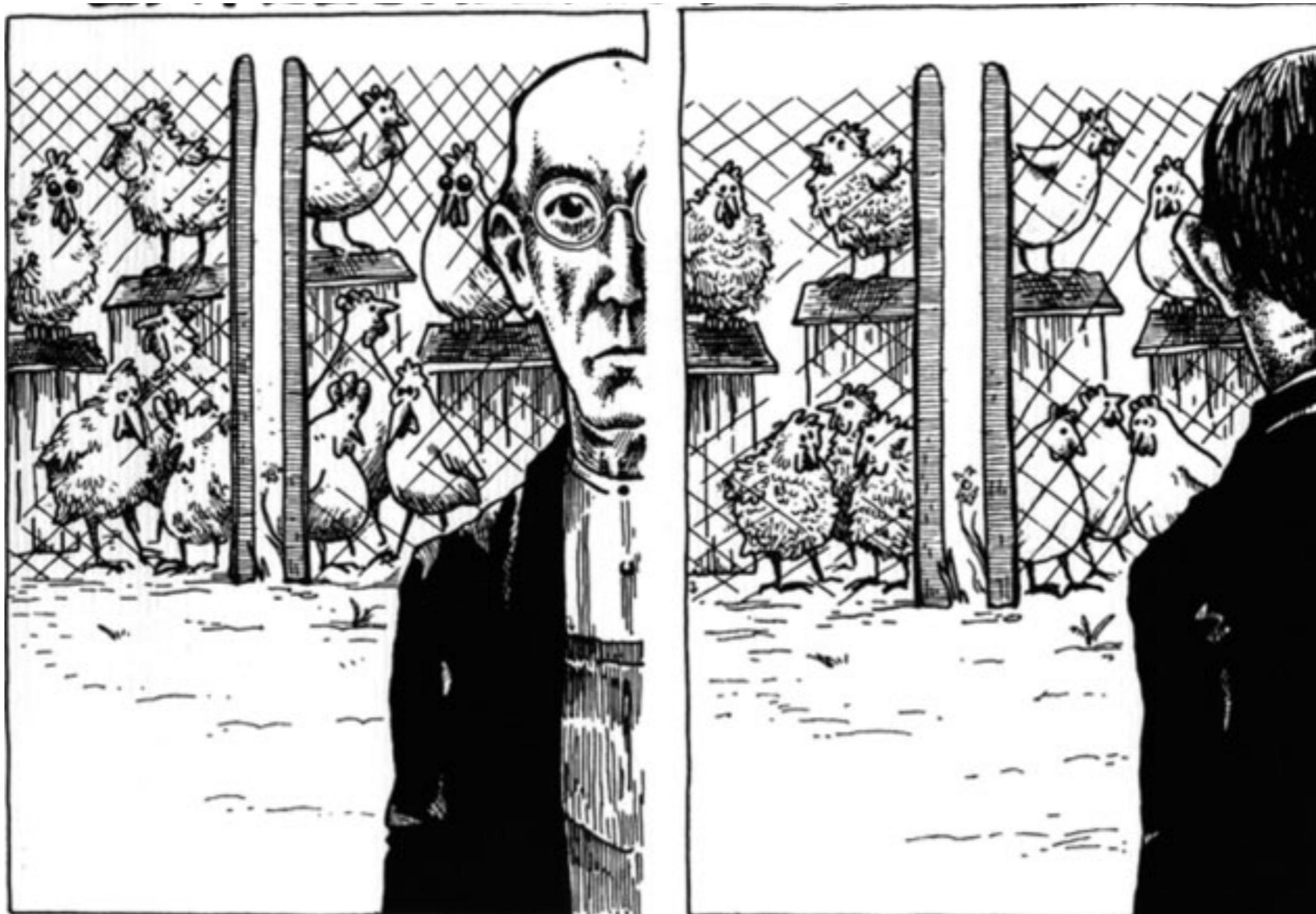
Categorical perception

- We tend to perceive our world in terms of the categories that we formed
- Our perceptions are warped such that differences between objects that belong in different categories are accentuated, and differences between objects that fall into the same categories are deemphasized





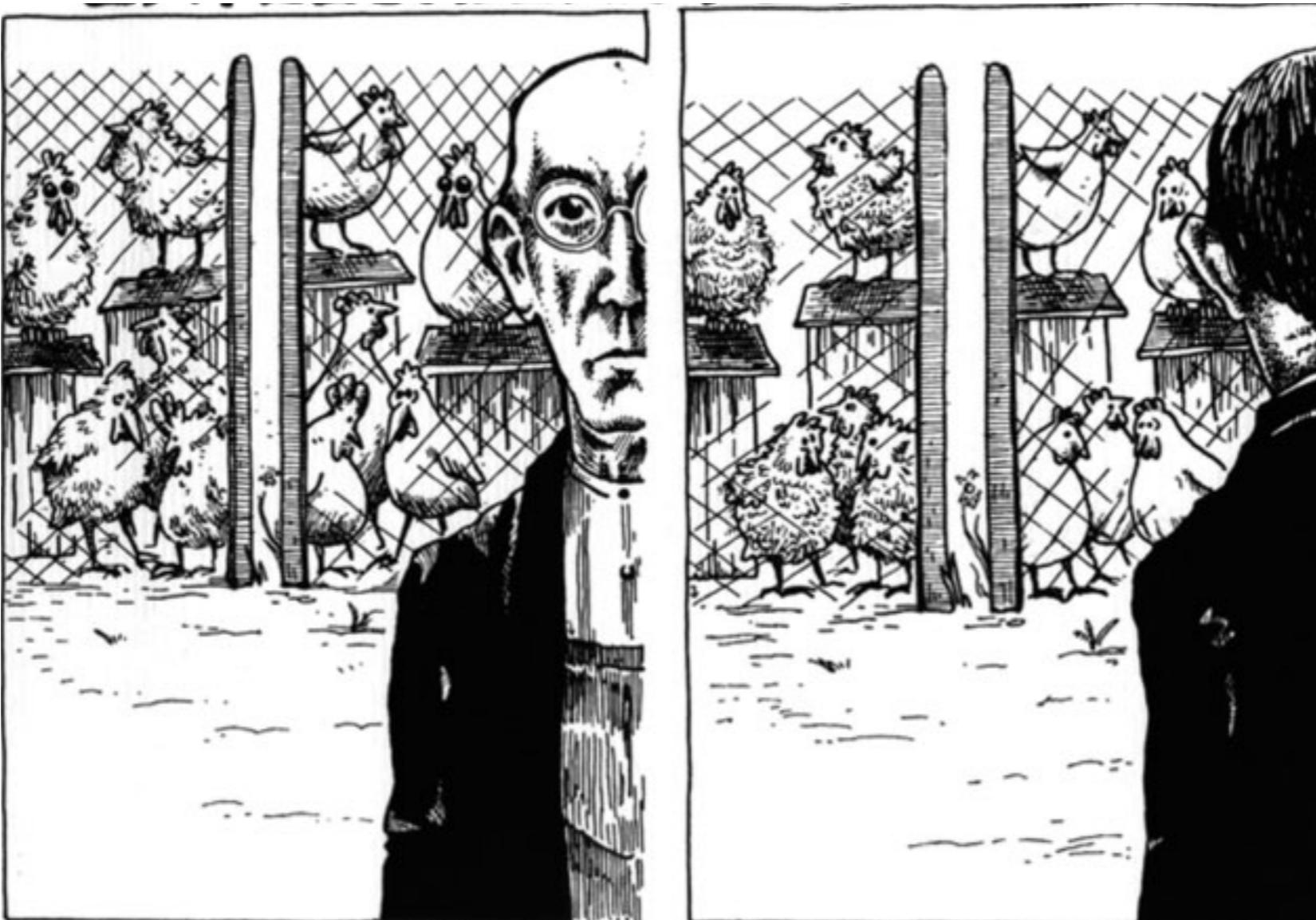
Categorical perception



From Goldstone and Hendrickson (2009)

acquired distinctiveness: differences between objects in different categories are emphasized

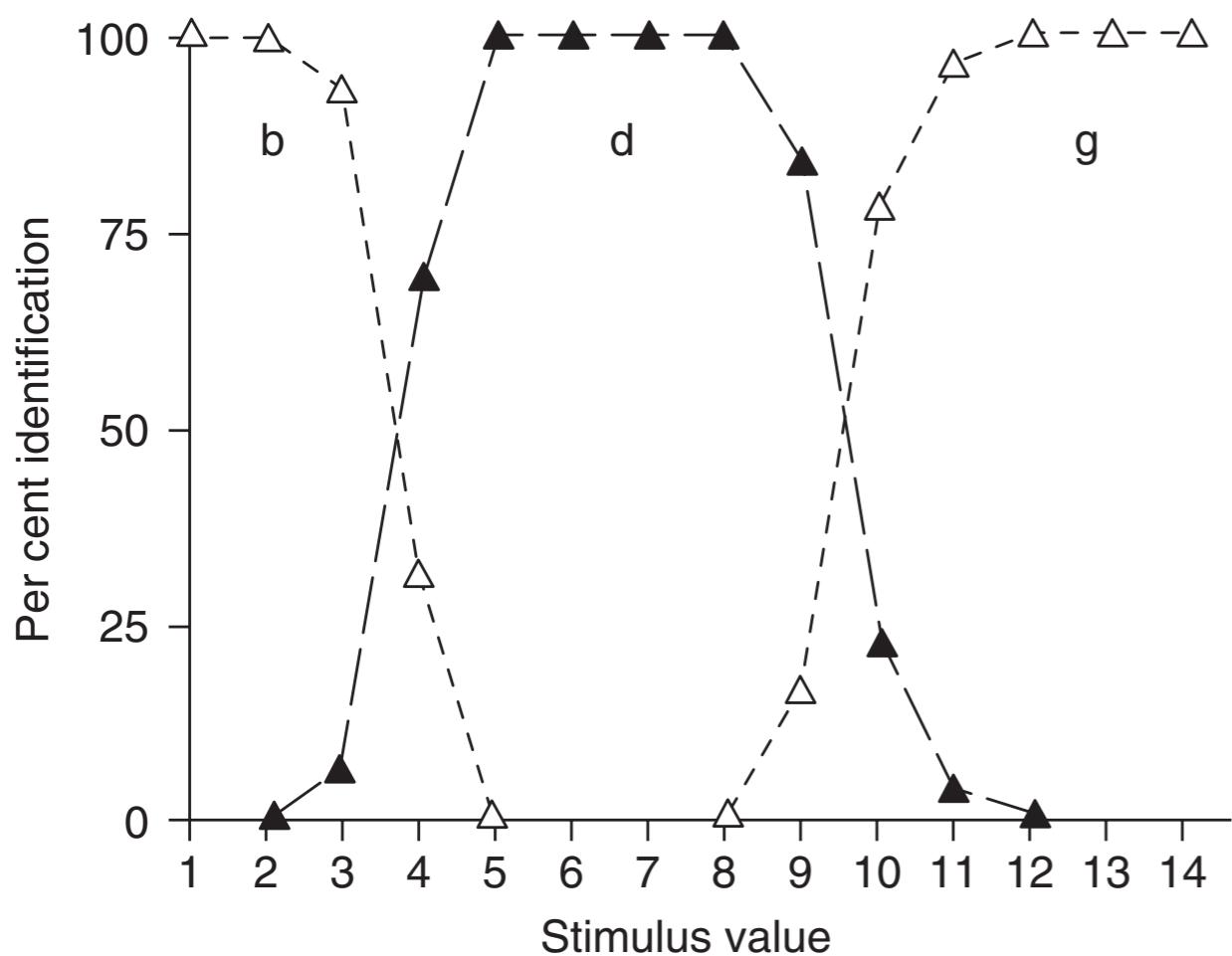
acquired similarity: differences between objects in the same categories are deemphasized



Categorical perception (CP) in speech

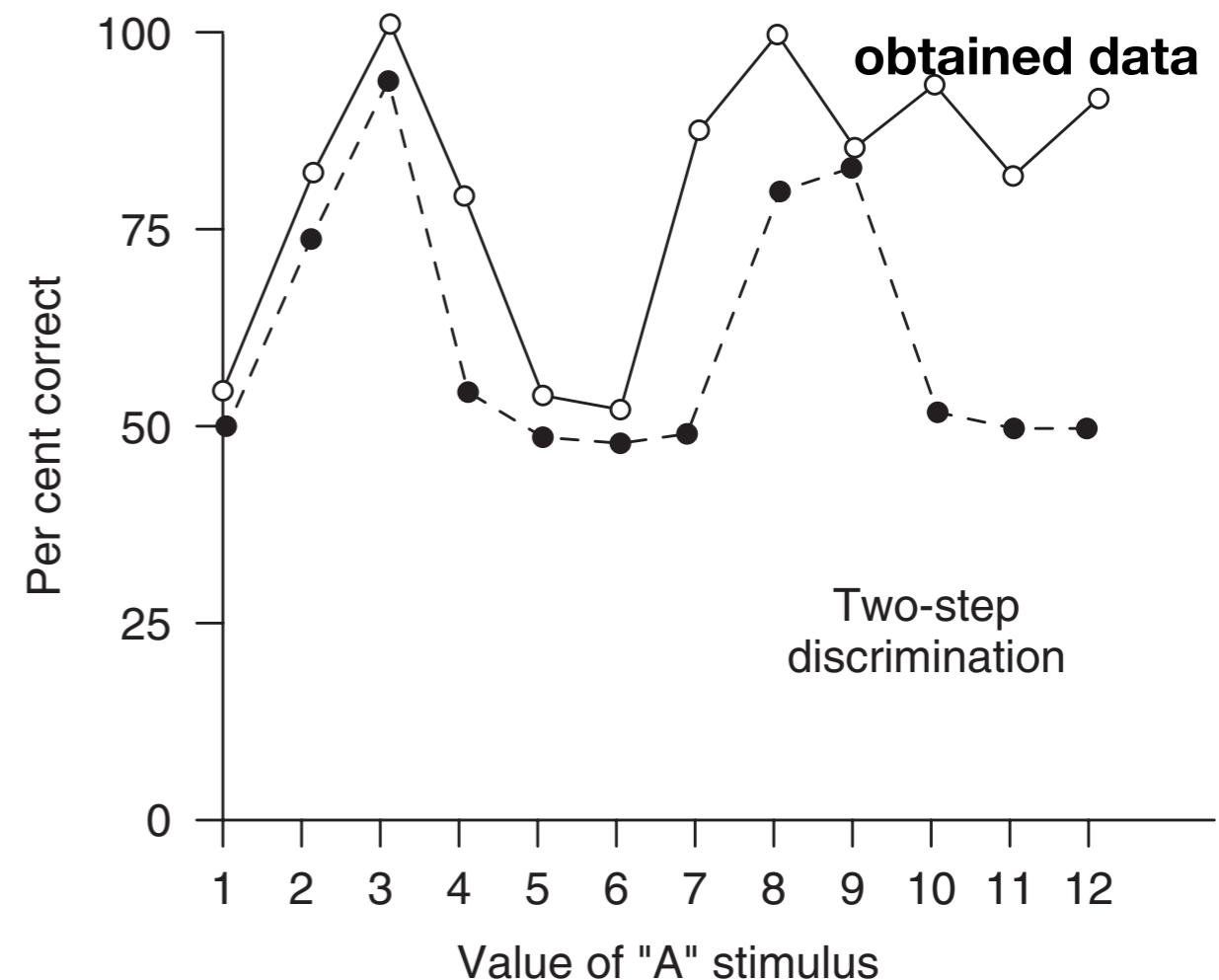
Categorization task

("ba" vs. "da" vs. "ga")



Discrimination task

(ABX; which is X identical to, A or B?)



From Liberman, A. M., Harris, K. S., Hoffman, H. S., & Griffith, B. C. (1957). The discrimination of speech sounds within and across phoneme boundaries. *Journal of experimental psychology*, 54(5), 358.

Linguistic Micro-Lectures

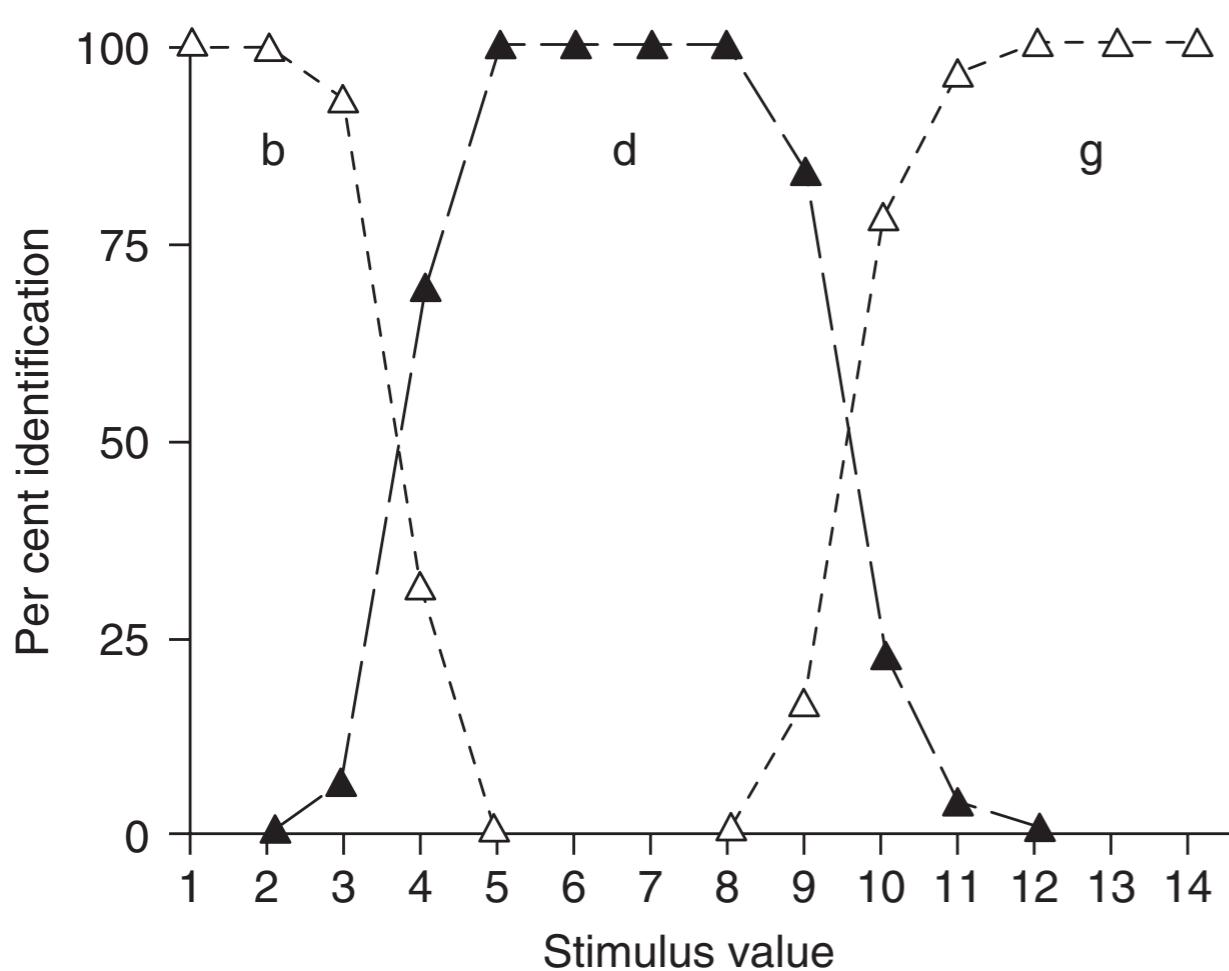


by Prof. Dr. Jürgen Handke
Marburg University, Germany

CP, but discrimination is not solely based on categorization

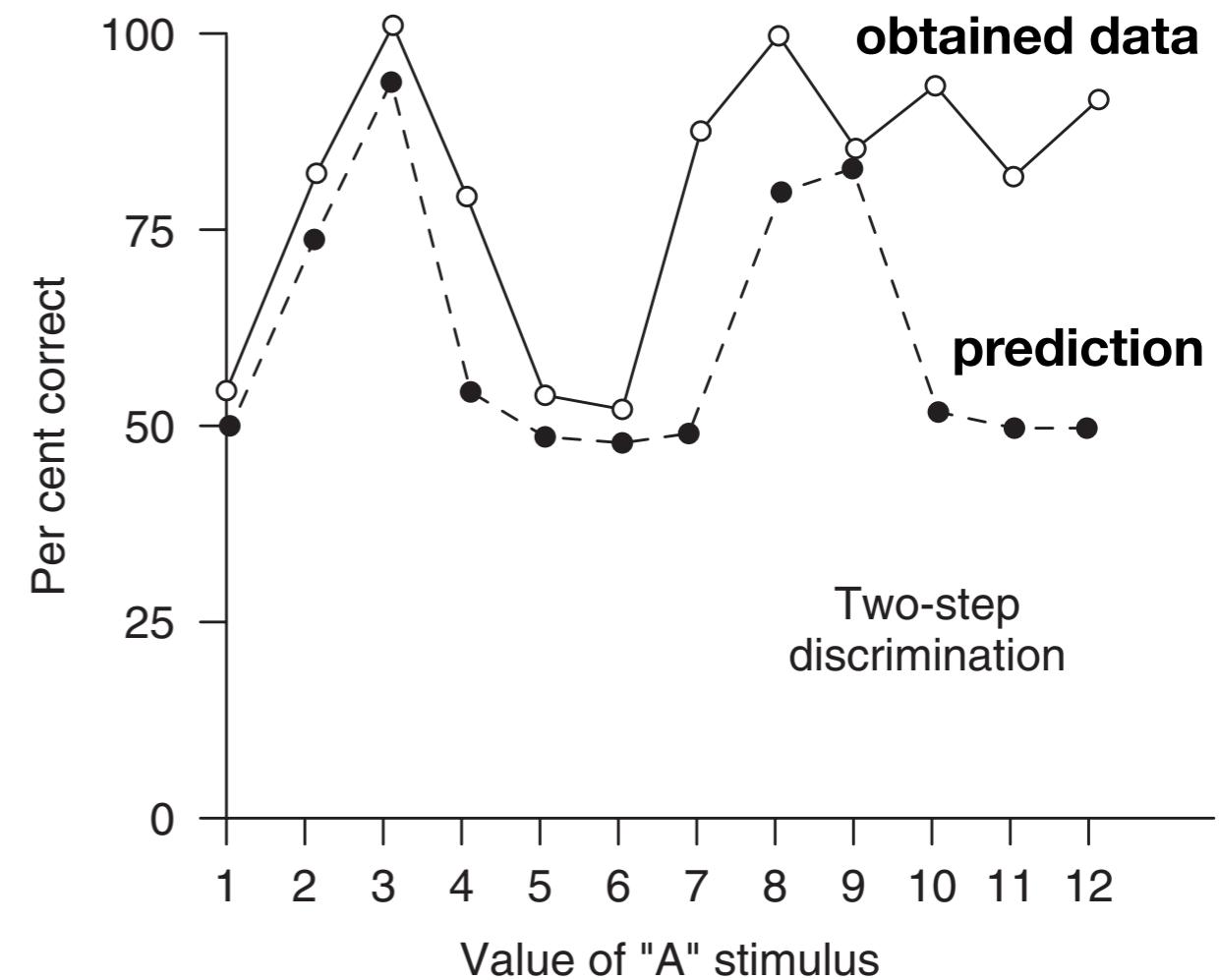
Categorization task

("ba" vs. "da" vs. "ga")



Discrimination task

(ABX; which is X identical to, A or B?)

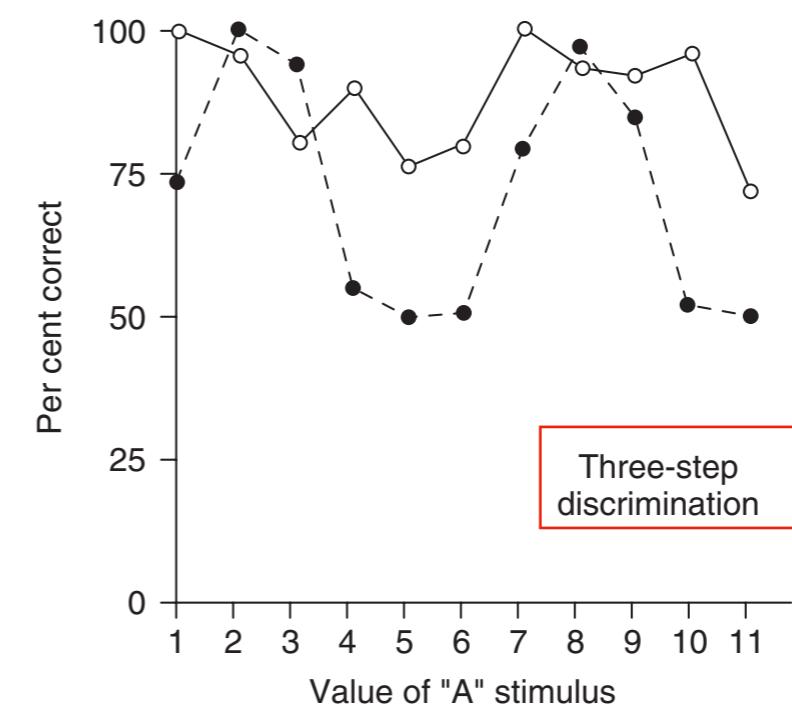
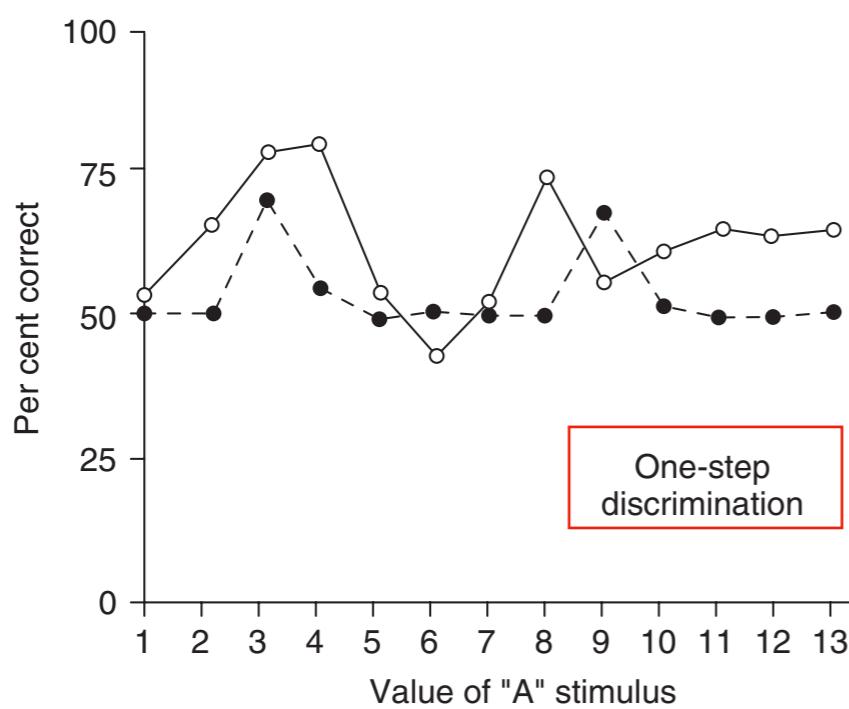
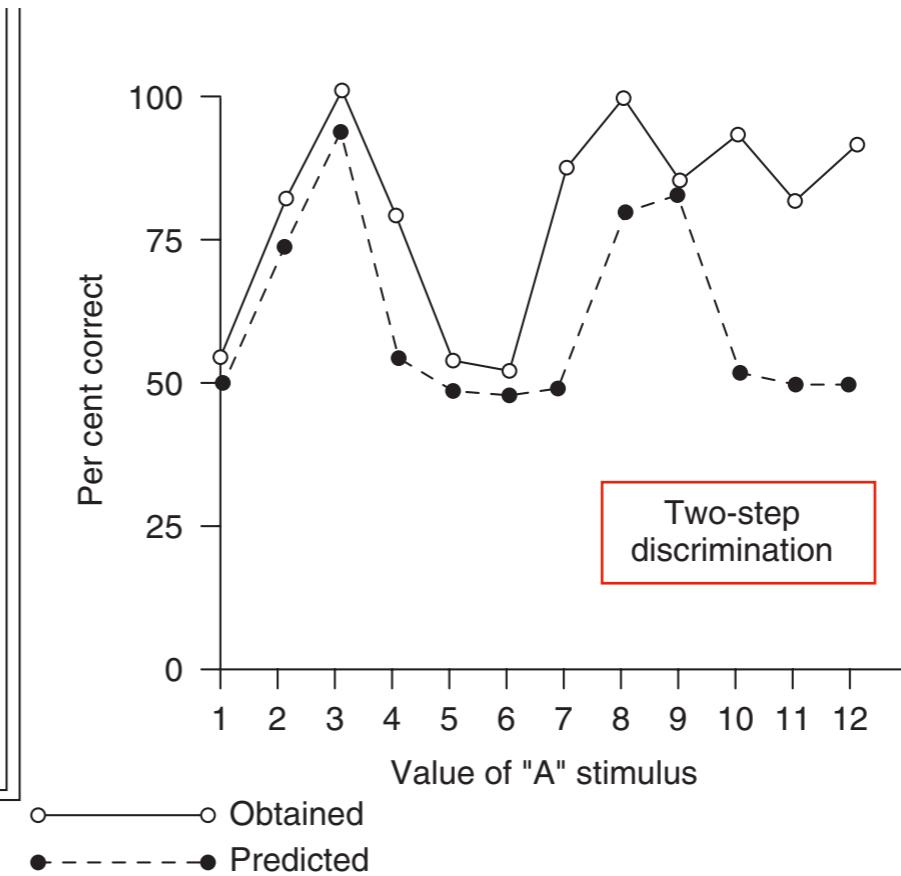
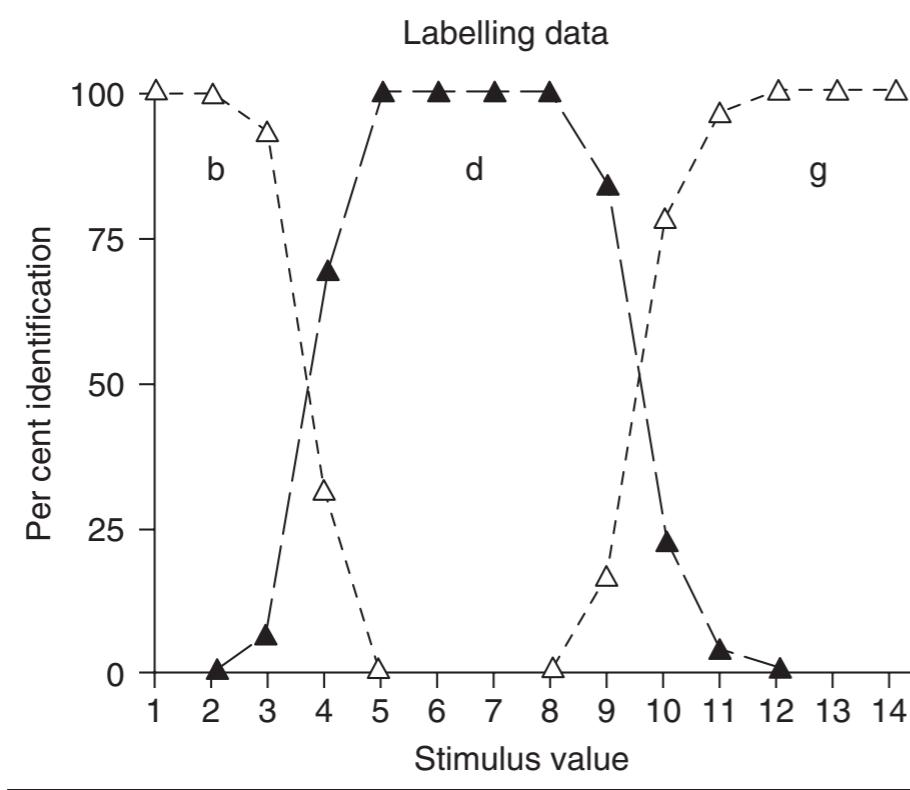


Strong CP model is wrong: discrimination cannot be predicted solely from a categorization task

$$P(\text{discrim } S_1 \text{ vs. } S_2) = \frac{1 + (P(A|S_1) - P(A|S_2))}{2}$$

- bounded between 0.5 and 1
- stimuli S_1 vs S_2
- categories A vs. B

CP, but discrimination is not solely based on categorization



Strong CP prediction:

$$P(\text{discrim } S_1 \text{ vs. } S_2) = \frac{1 + (P(A|S_1) - P(A|S_2))}{2}$$

Russian blues reveal effects of language on color discrimination

Jonathan Winawer^{*†‡}, Nathan Witthoft^{*‡}, Michael C. Frank^{*}, Lisa Wu[§], Alex R. Wade[¶], and Lera Boroditsky[‡]

^{*}Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139-4307; [§]Department of Neurology, David Geffen School of Medicine, University of California, Los Angeles, CA 90095-1769; [¶]Brain Imaging Center, Smith-Kettlewell Eye Research Institute, San Francisco, CA 94115; and [‡]Department of Psychology, Stanford University, Stanford, CA 94305

Communicated by Gordon H. Bower, Stanford University, Stanford, CA, March 7, 2007 (received for review September 22, 2006)

English and Russian color terms divide the color spectrum differently. Unlike English, Russian makes an obligatory distinction between lighter blues (“goluboy”) and darker blues (“siniy”). We investigated whether this linguistic difference leads to differences in color discrimination. We tested English and Russian speakers in a speeded color discrimination task using blue stimuli that spanned the siniy/goluboy border. We found that Russian speakers were faster to discriminate two colors when they fell into different linguistic categories in Russian (one siniy and the other goluboy) than when they were from the same linguistic category (both siniy or both goluboy). Moreover, this category advantage was eliminated by a verbal, but not a spatial, dual task. These effects were stronger for difficult discriminations (i.e., when the colors were perceptually close) than for easy discriminations (i.e., when the colors were further apart). English speakers tested on the identical stimuli did not show a category advantage in any of the conditions. These results demonstrate that (i) categories in language affect performance on simple perceptual color tasks and (ii) the effect of language is online (and can be disrupted by verbal interference).

categorization | cross-linguistic | Whorf

Different languages divide color space differently. For example, the English term “blue” can be used to describe all of the colors in Fig. 1. Unlike English, Russian makes an obligatory distinction between lighter blues (“goluboy”) and darker blues (“siniy”). Like other basic color words, “siniy” and “goluboy” tend to be learned early by Russian children (1) and share many of the usage and behavioral properties of other basic color words (2). There is no single generic word for “blue” in Russian that can be used to describe all of the colors in Fig. 1 (nor to adequately translate the title of this work from English to Russian). Does this difference between languages lead to differences in how people discriminate colors?

The question of cross-linguistic differences in color perception has a long and venerable history (e.g., refs. 3–14) and has been a cornerstone issue in the debate on whether and how much language shapes thinking (15). Previous studies have found cross-linguistic differences in subjective color similarity judgments and color confusability in memory (4, 5, 10, 12, 16). For

Most of the experiments have tested banal “weak” versions of the Whorfian hypothesis, namely that words can have some effect on memory or categorization. . . . In a typical experiment, subjects have to commit paint chips to memory and are tested with a multiple-choice procedure. In some of these studies, the subjects show slightly better memory for colors that have readily available names in their language. . . . All [this] shows is that subjects remembered the chips in two forms, a non-verbal visual image and a verbal label, presumably because two types of memory, each one fallible, are better than one. In another type of experiment subjects have to say which two of three color chips go together; they often put the ones together that have the same name in their language. Again, no surprise. I can imagine the subjects thinking to themselves, “Now how on earth does this guy expect me to pick two chips to put together? He didn’t give me any hints, and they’re all pretty similar. Well, I’d probably call these two ‘green’ and that one ‘blue,’ and that seems as good a reason to put them together as any.”

Because previous cross-linguistic comparisons have relied on memory procedures or subjective judgments, the question of whether language affects objective color discrimination performance has remained. Studies testing only color memory leave open the possibility that, when subjects make perceptual discriminations among stimuli that can all be viewed at the same time, language may have no influence. In studies measuring subjective similarity, it is possible that any language-congruent bias results from a conscious, strategic decision on the part of the subject (19). Thus, such methods leave open the question of whether subjects’ normal ability to discriminate colors in an objective procedure is altered by language.

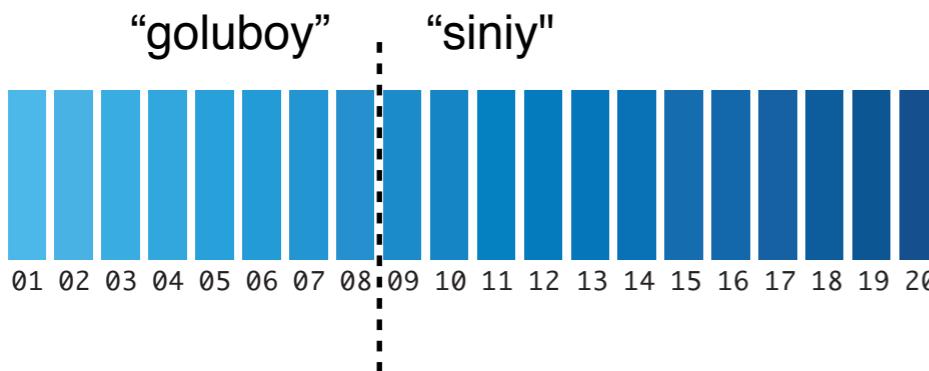
Here we measure color discrimination performance in two language groups in a simple, objective, perceptual task. Subjects were simultaneously shown three color squares arranged in a triad (see Fig. 1) and were asked to say which of the bottom two color squares was perceptually identical to the square on top.

This design combined the advantages of previous tasks in a way that allowed us to test for the effects of language on color

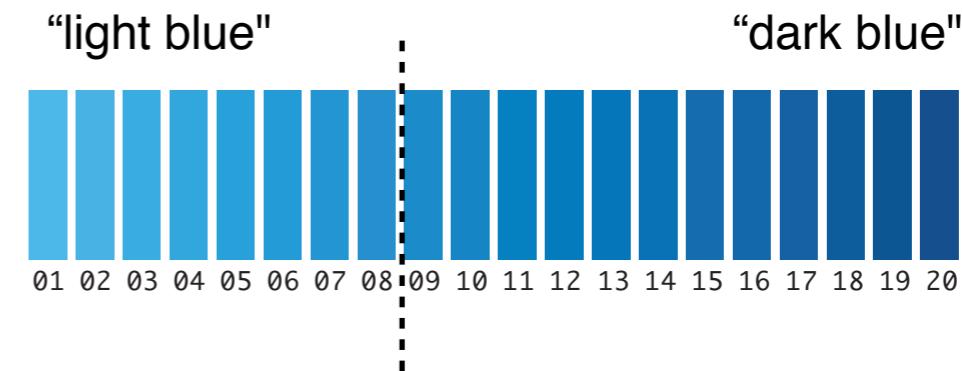


Effects of language on color discrimination

Russian (color word distinction)



English (NO color word distinction)



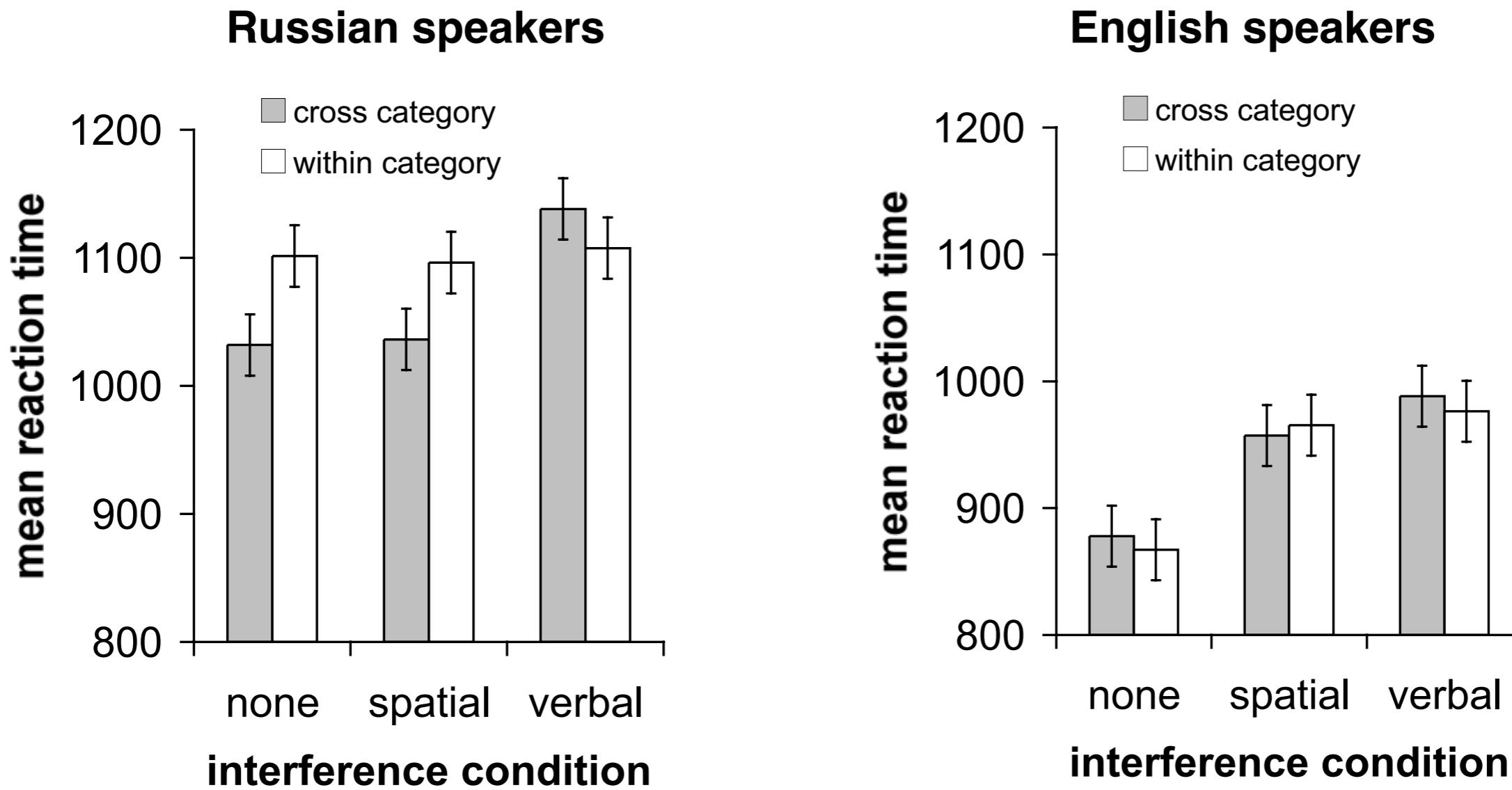
Discrimination task



which is the same as the above square?



Results: Effects of language on color discrimination



- Russian speakers show categorical perception, but the effect is eliminated with verbal (but not spatial) interference
- English speakers did not show categorical perception

Influences of Categorization on Perceptual Discrimination

Robert Goldstone



Robert Goldstone
Indiana University

Four experiments investigated the influence of categorization training on perceptual discrimination. Ss were trained according to 1 of 4 different categorization regimes. Subsequent to category learning, Ss performed a Same-Different judgment task. Ss' sensitivities (d' 's) for discriminating between items that varied on category-(ir)relevant dimensions were measured. Evidence for acquired distinctiveness (increased perceptual sensitivity for items that are categorized differently) was obtained. One case of acquired equivalence (decreased perceptual sensitivity for items that are categorized together) was found for separable, but not integral, dimensions. Acquired equivalence within a categorization-relevant dimension was never found for either integral or separable dimensions. The relevance of the results for theories of perceptual learning, dimensional attention, categorical perception, and categorization are discussed.

Psychologists have long been intrigued by the possibility that the concepts that people learn influence their perceptual abilities. It may be that the way people organize their world into categories alters the actual appearance of their world. The purpose of the present research is to investigate influences of concept learning on perception.

The notion that experience and expectations can influence perception can be traced back to the "New Look" movement of the 1940s and 50s (J. A. Bruner & Postman, 1949). Evidence suggests that experts perceive structures in X rays (Norman, Brooks, Coblenz, & Babcock, 1992), beers (Peron & Allen, 1988), and infant chickens (Biederman & Shiffrar, 1987) that are missed by novices. As the experts in these fields learn to distinguish among the concepts in their domain (types of fractures, brands of beer, or genders of chickens), they seem to acquire new ways of perceptually structuring the objects to be categorized.

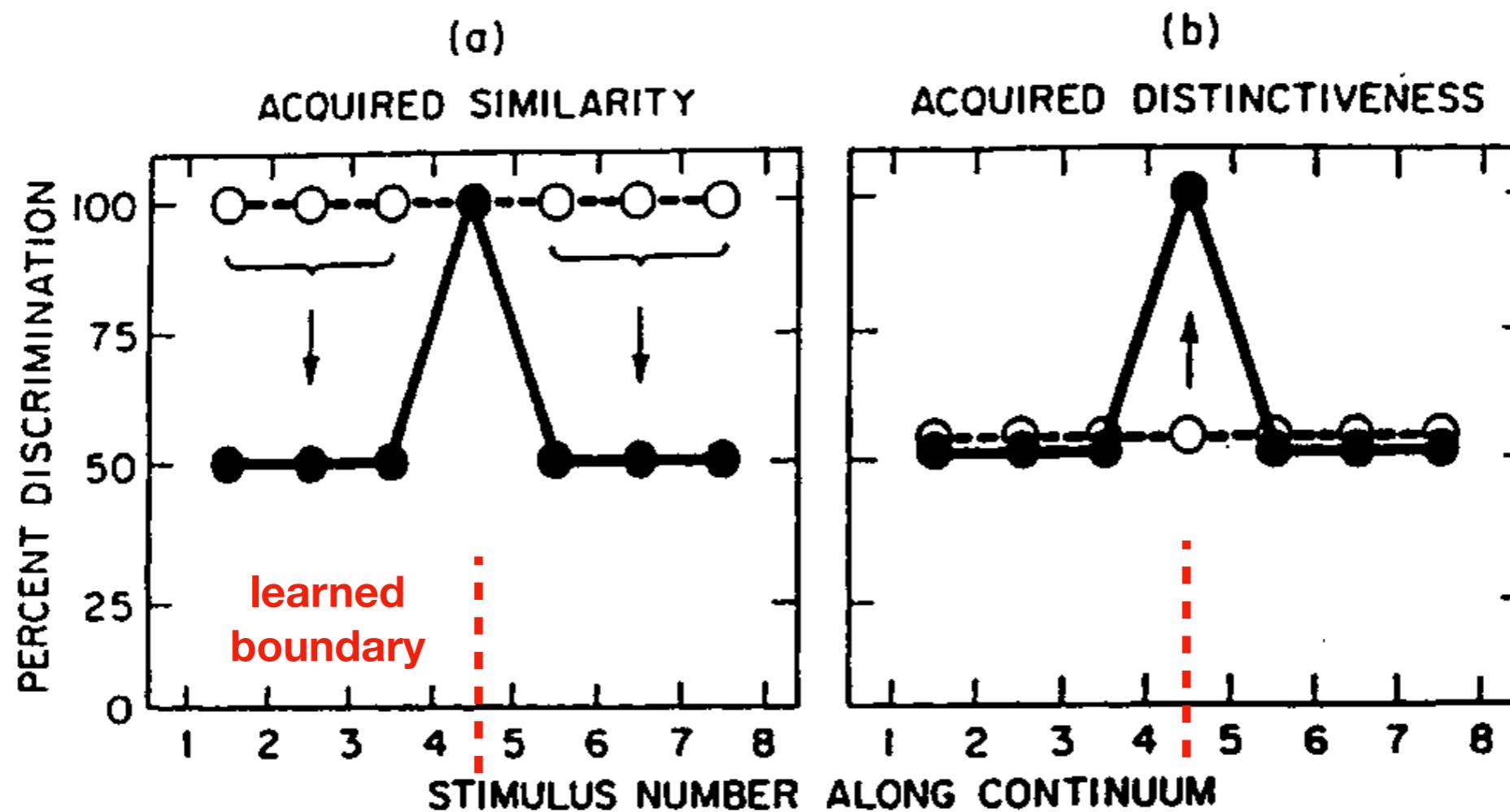
This suggestion—that categorization causes changes to perceptual abilities—is not implicated in most traditional accounts of concept learning. In J. S. Bruner, Goodnow, and Austin's (1956) classic studies of concept learning, subjects saw flash cards with shapes and were required to learn rules

learning has come a long way since J. S. Bruner et al.'s study (Estes, 1986; Kruschke, 1992; Medin & Schaffer, 1978; Nosofsky, 1986; Reed, 1972), vestiges of this earlier work are apparent in current research. Specifically, many researchers have investigated concept learning using stimuli that have clear-cut dimensions with clearly different values on these dimensions. Although such stimuli are mandatory in many cases for experimental control and precision, they do not require subjects to perceptually learn new dimensions or finer discriminations. In the present described concept learning tasks, subjects had to make fine discriminations along dimensions or isolate dimensions that normally are fused together. In both cases, the perceptual abilities required for the categorization task are not at a ceiling level before categorization training begins; consequently, experience with categorization may drive perceptual learning.

Evidence for an Influence of Learning on Perception
Perceptual Learning

acquired similarity: differences between objects in the same categories are deemphasized

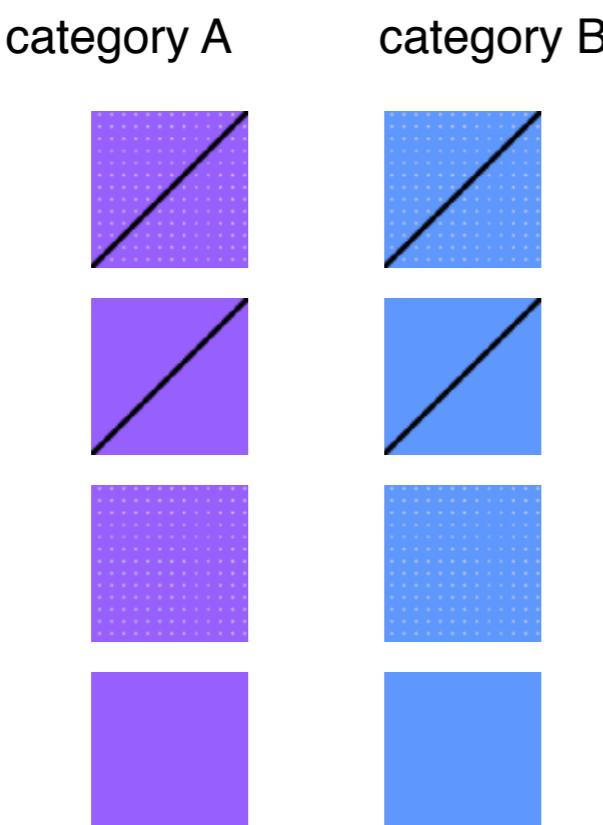
acquired distinctiveness: differences between objects in different categories are emphasized



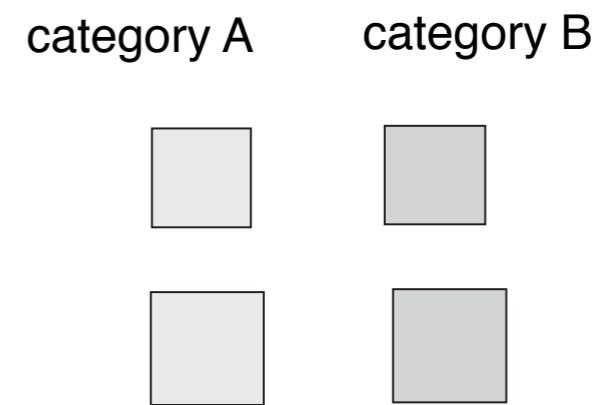
Category learning with hard perceptual discrimination

- Most category learning experiments have clear-cut stimuli with clear values on the dimensions.
- In present study, participants need to make fine discriminations between categories, showing how experience with categorization may drive perceptual learning

**Most experiments have easy perceptual discriminations
(e.g., Shepard, Honvland, & Jenkins)**

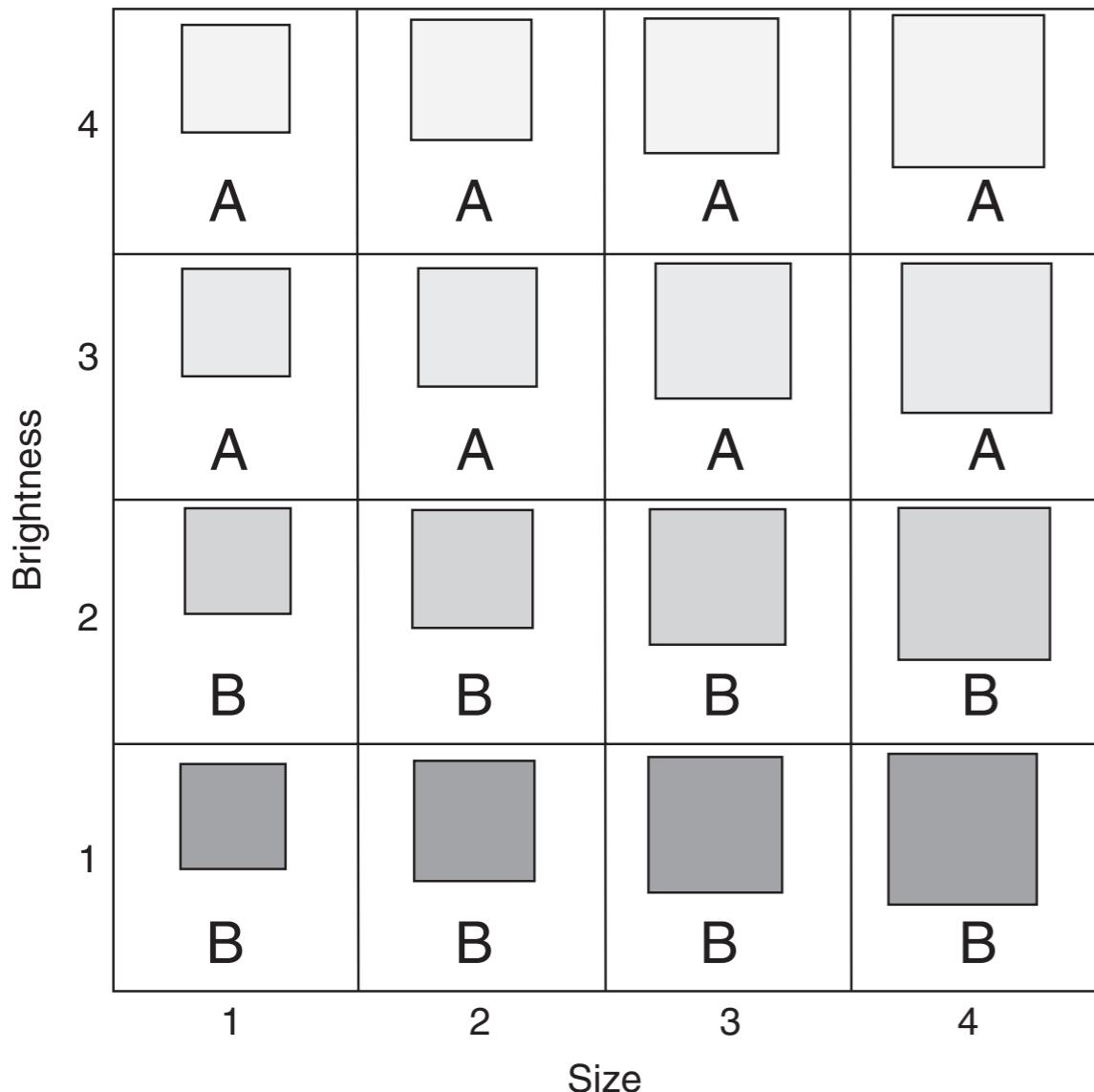


**Harder discrimination
(Goldstone)**

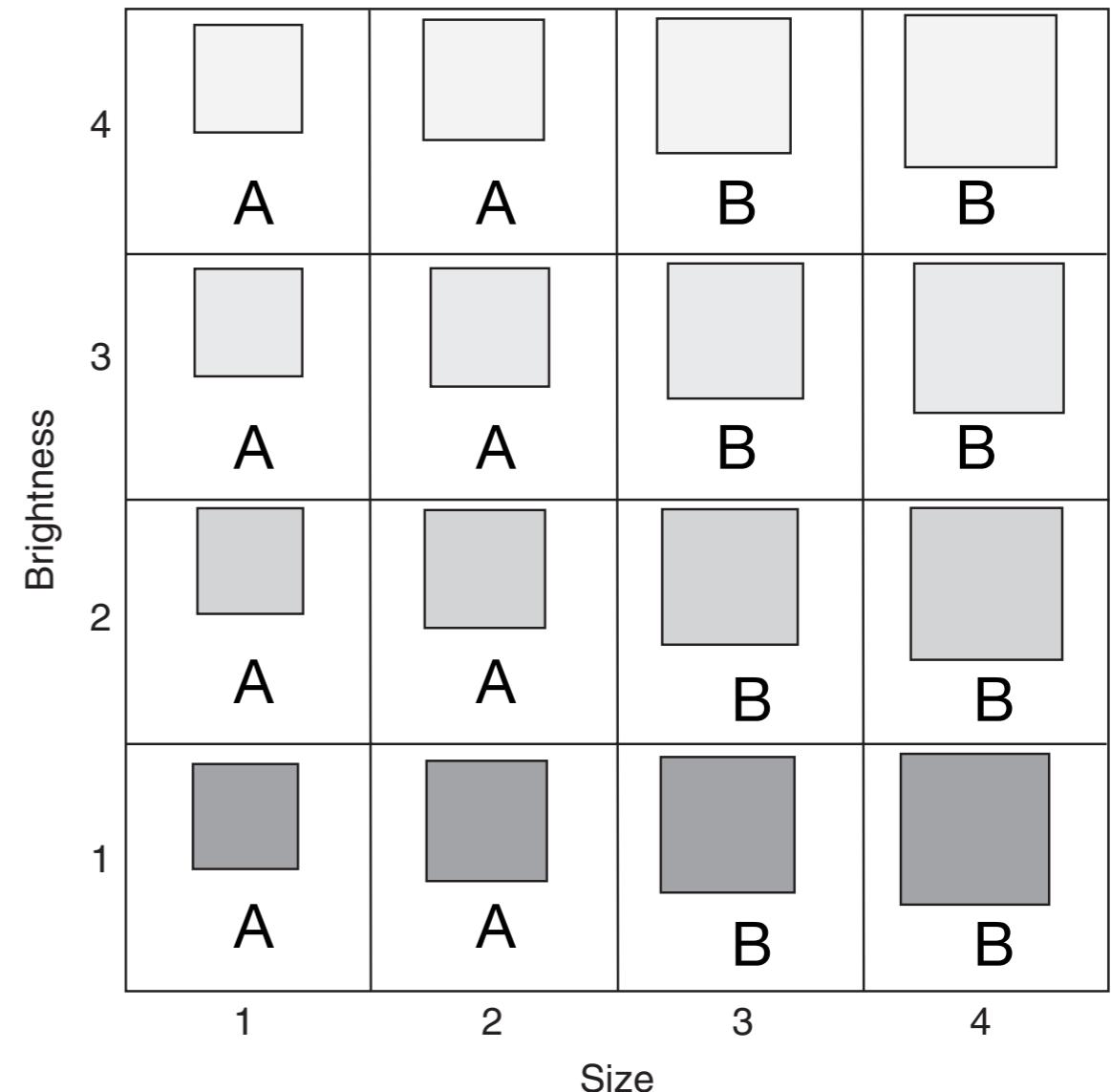


Goldstone Ex 2 : four conditions

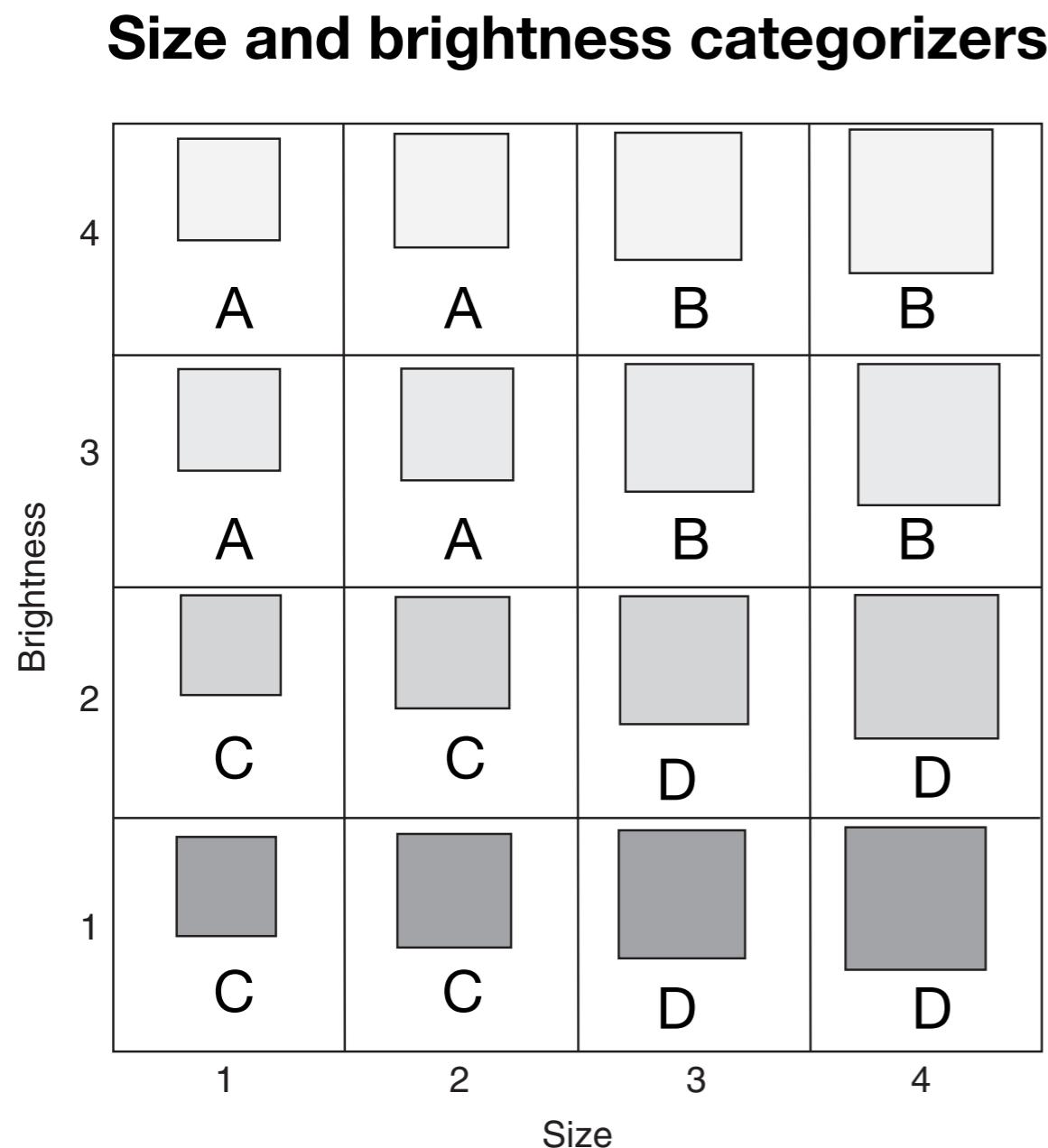
Brightness categorizers



Size categorizers



Goldstone Ex 2 : four conditions



**Control
(no categorization)**

Goldstone Ex 2 : method

Example categorization trial

- Phase 1: category learning (~60 min)
 - * 20 training runs over the 16 stimuli
- Phase 2: discrimination (~40 min)
 - * 576 trials
 - * judged “same” vs. “different”
 - * sensitivity measured as $d' = Z(\text{hit rate}) - Z(\text{false alarm rate})$, where Z is normal CDF

Category A or B?



Example discrimination trial (are stimuli physically identical?)

Stimulus 1
(1000 ms)



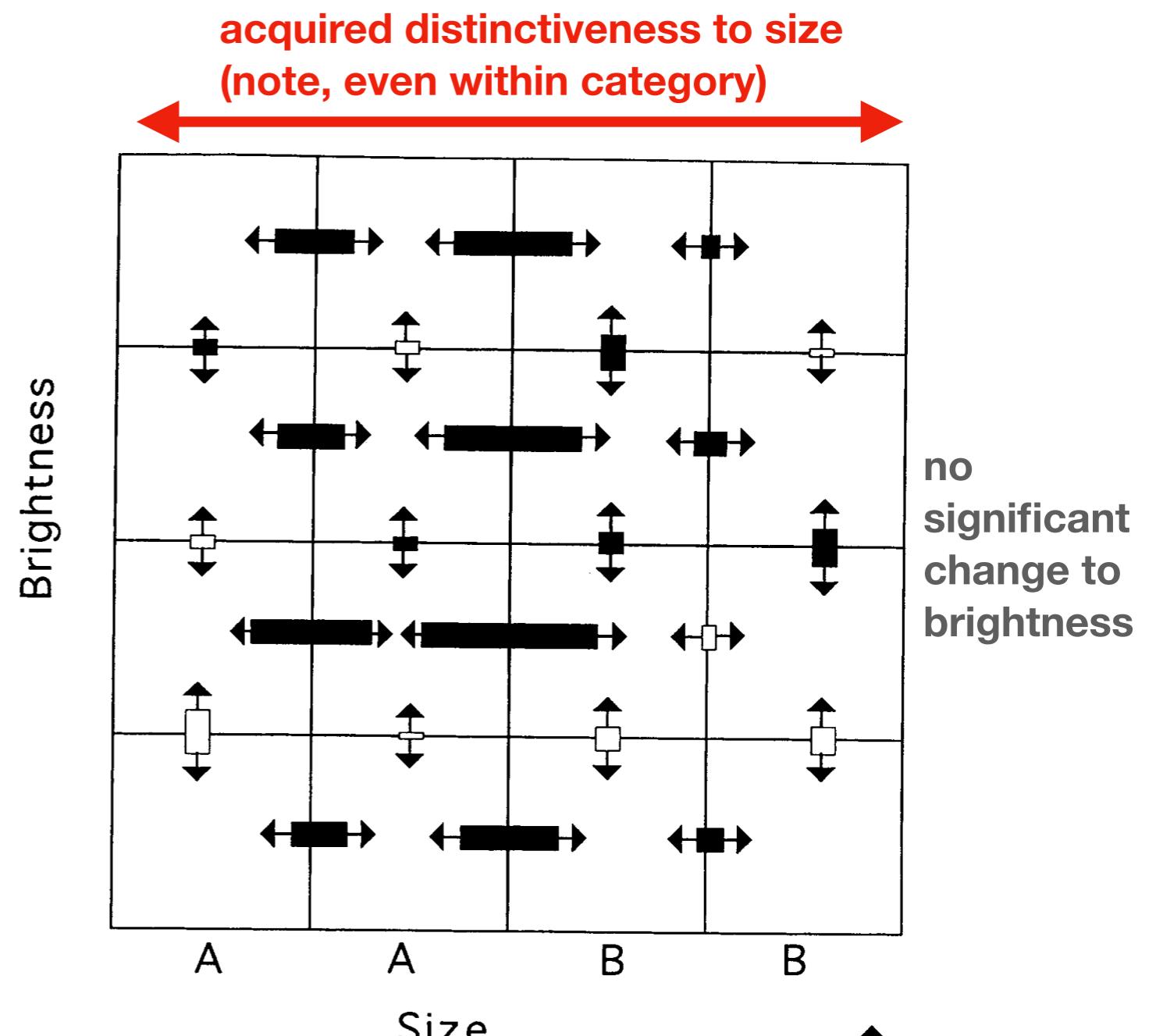
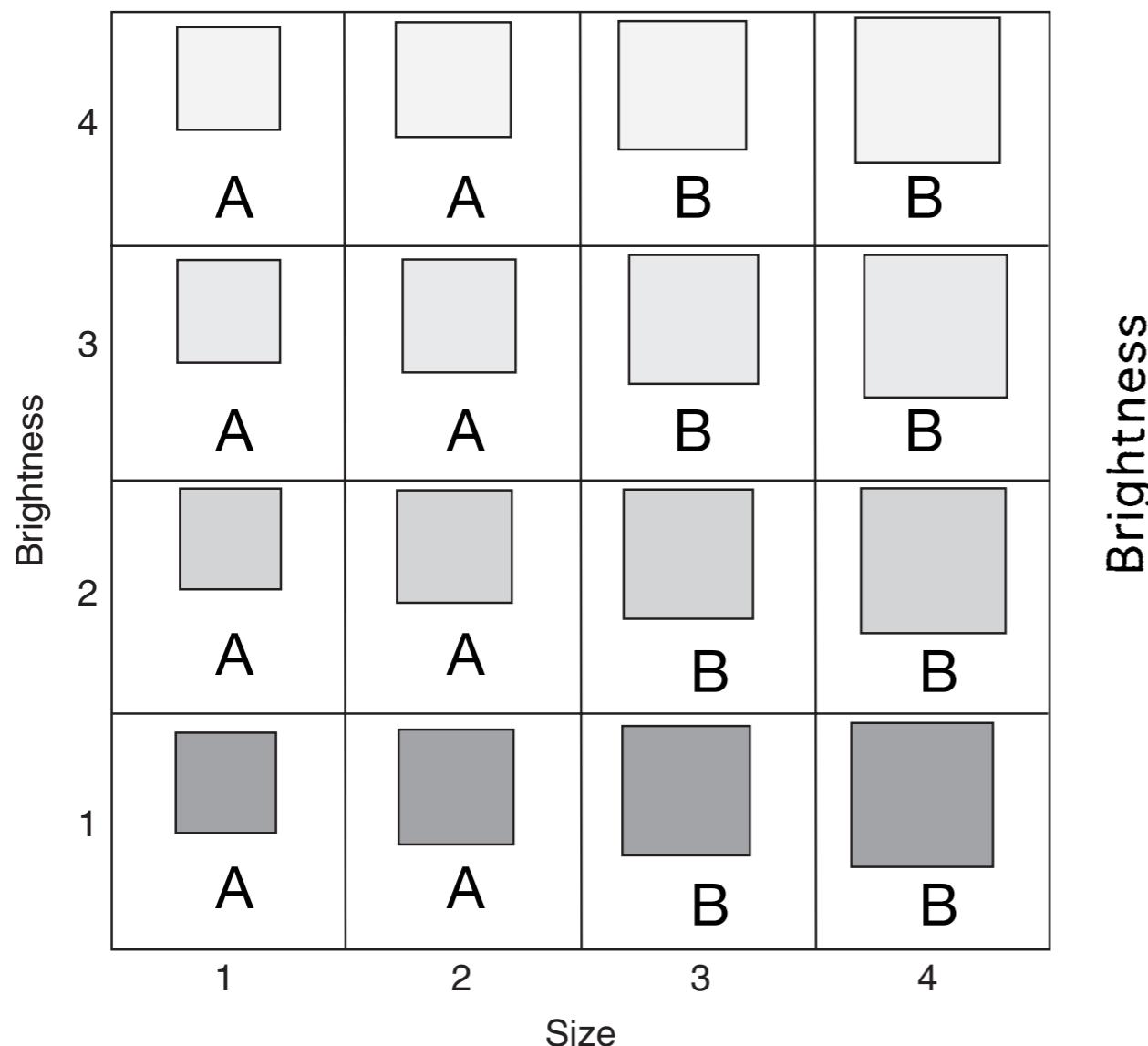
Stimulus 2
(1000 ms)



Response
“same” vs.
“different”

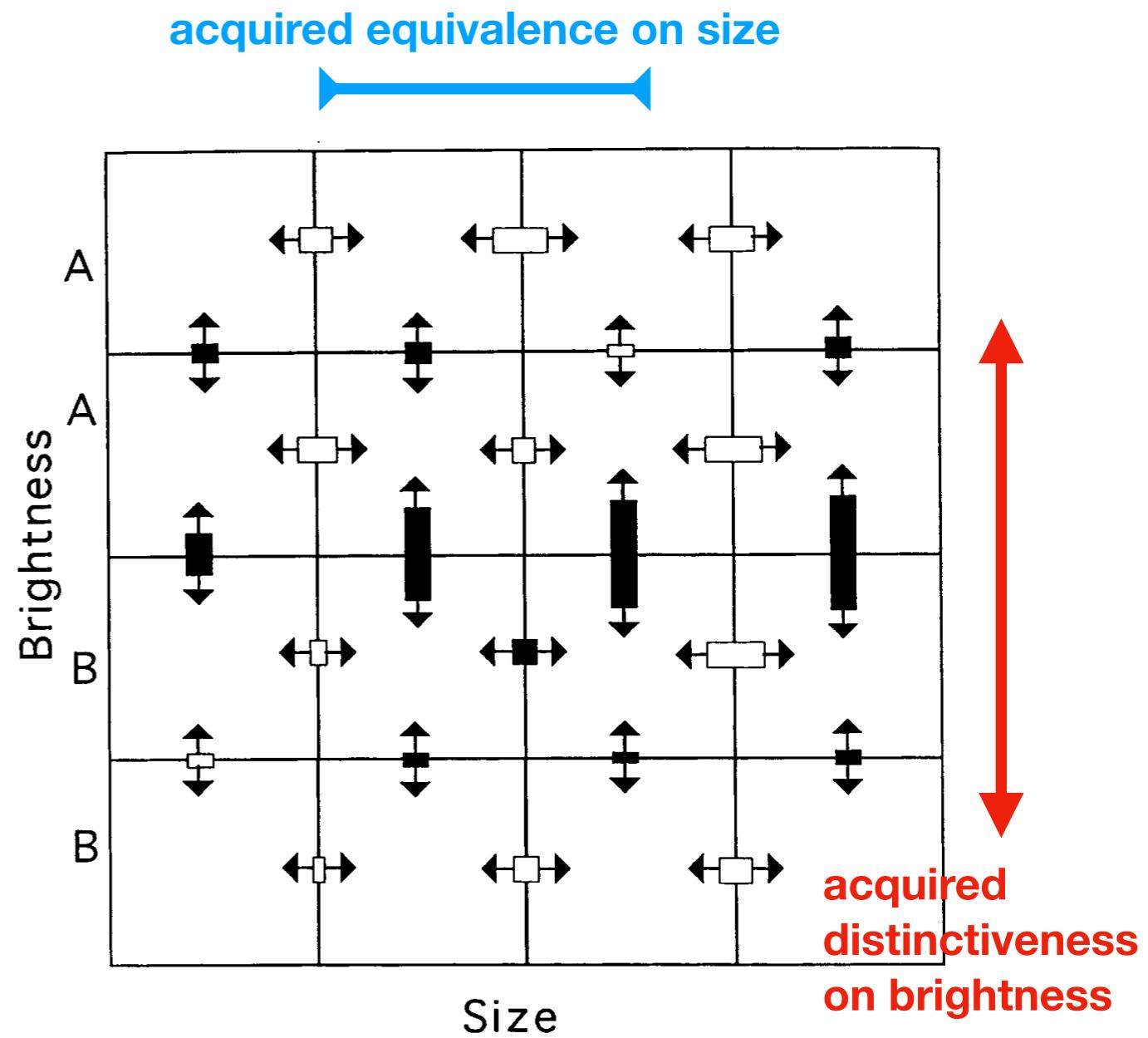
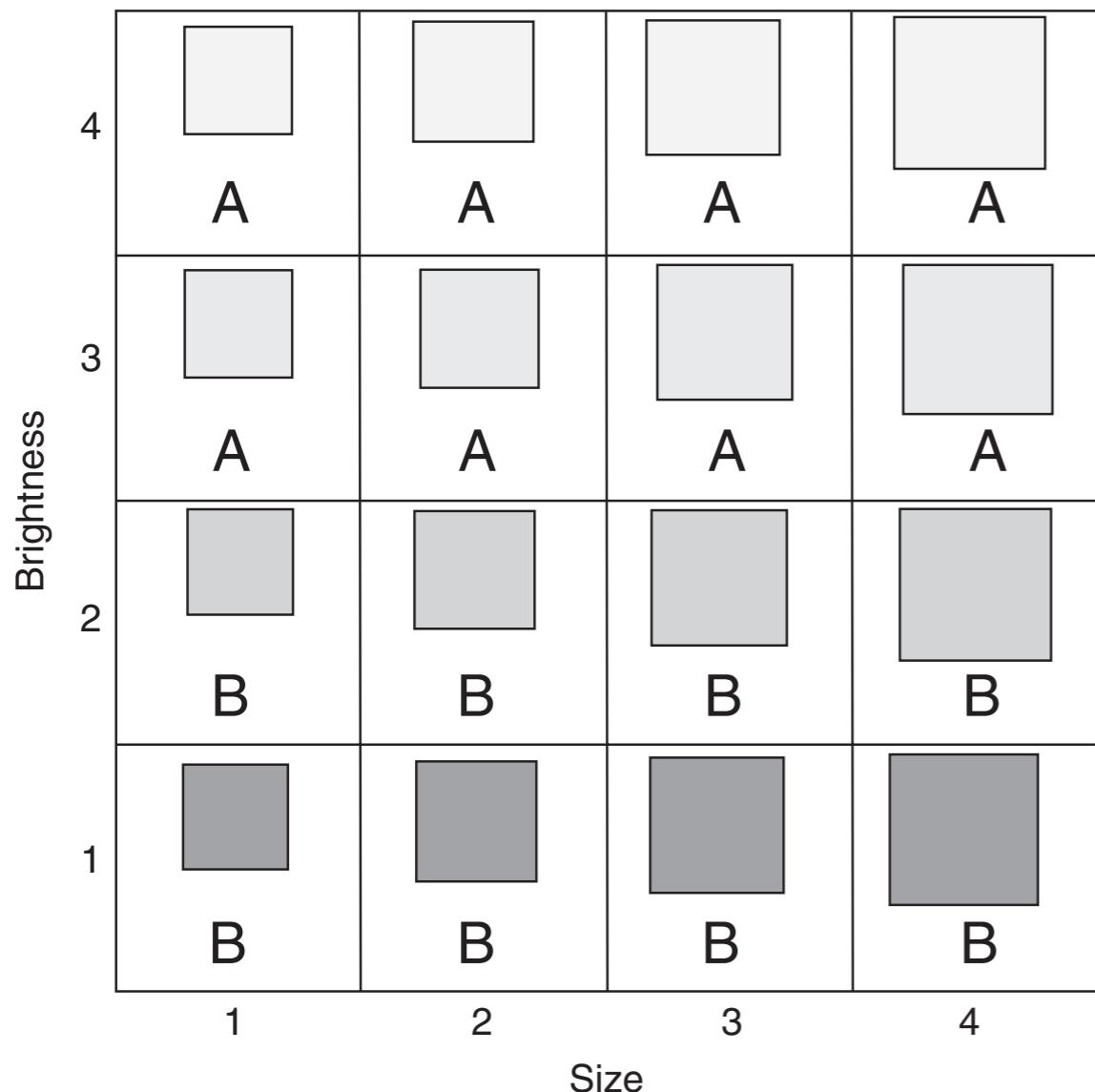
Goldstone Ex 2 : Results

Size categorizers

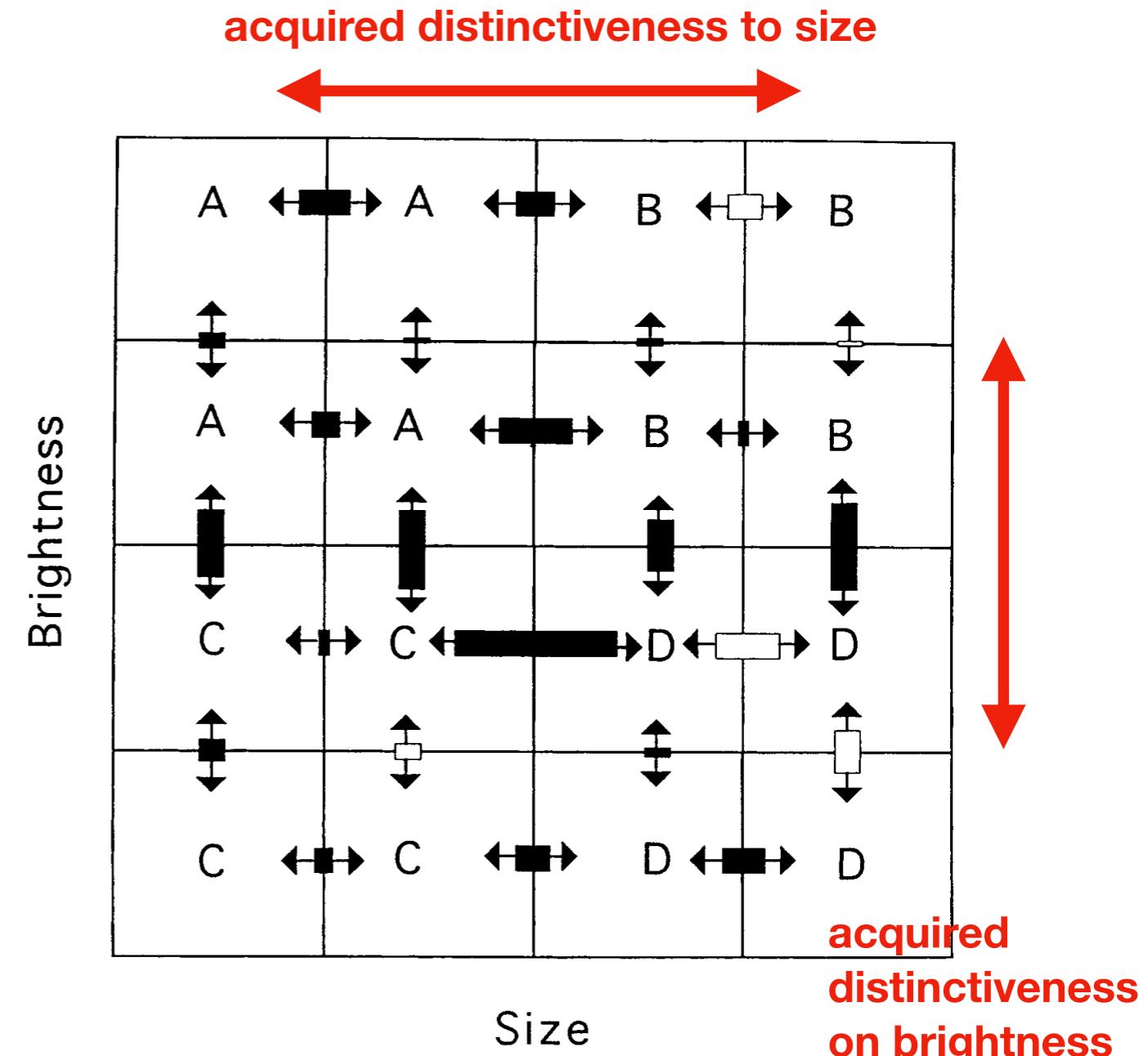
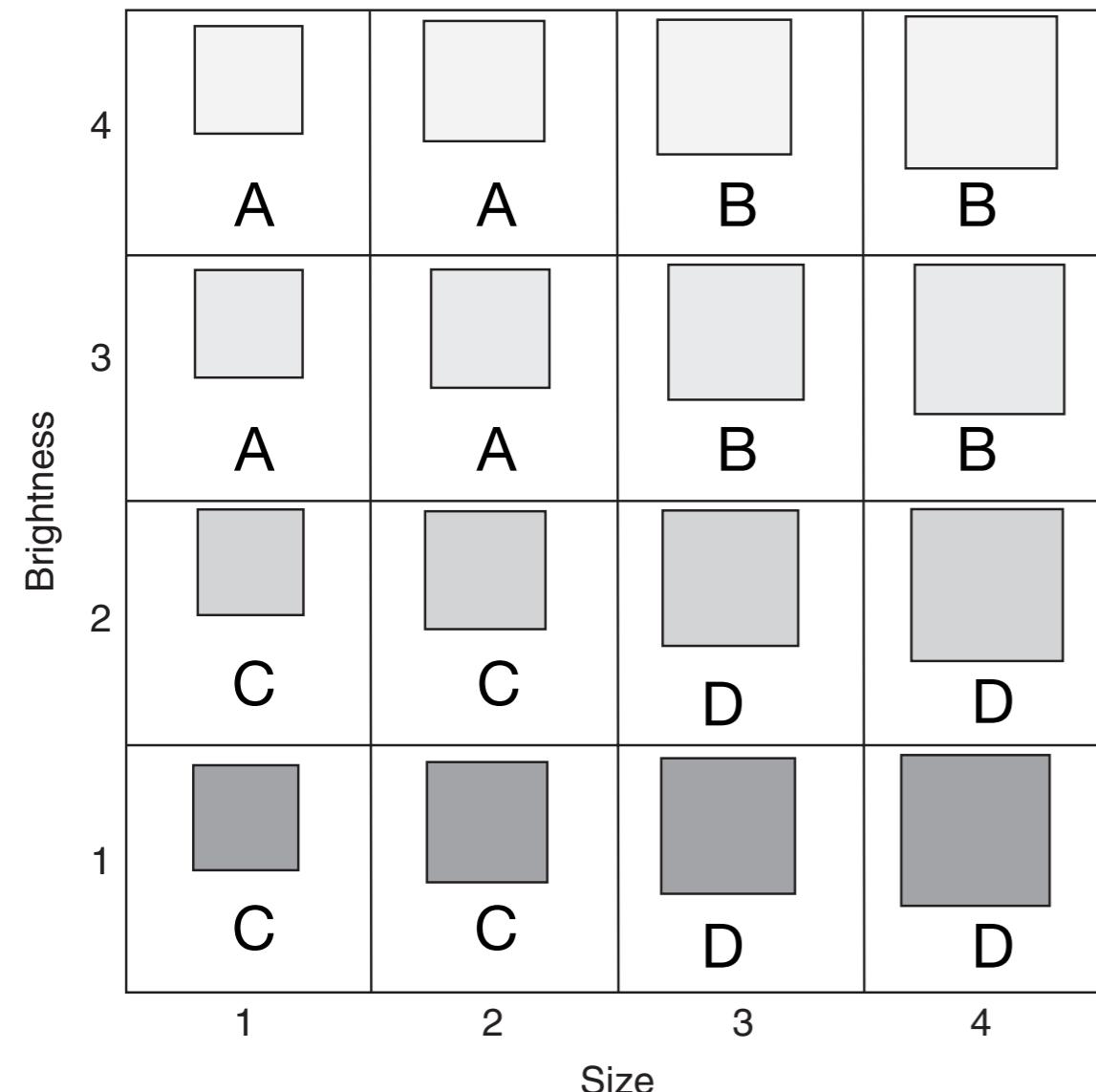


gain in sensitivity
over control
condition

Brightness categorizers



Size and brightness categorizers

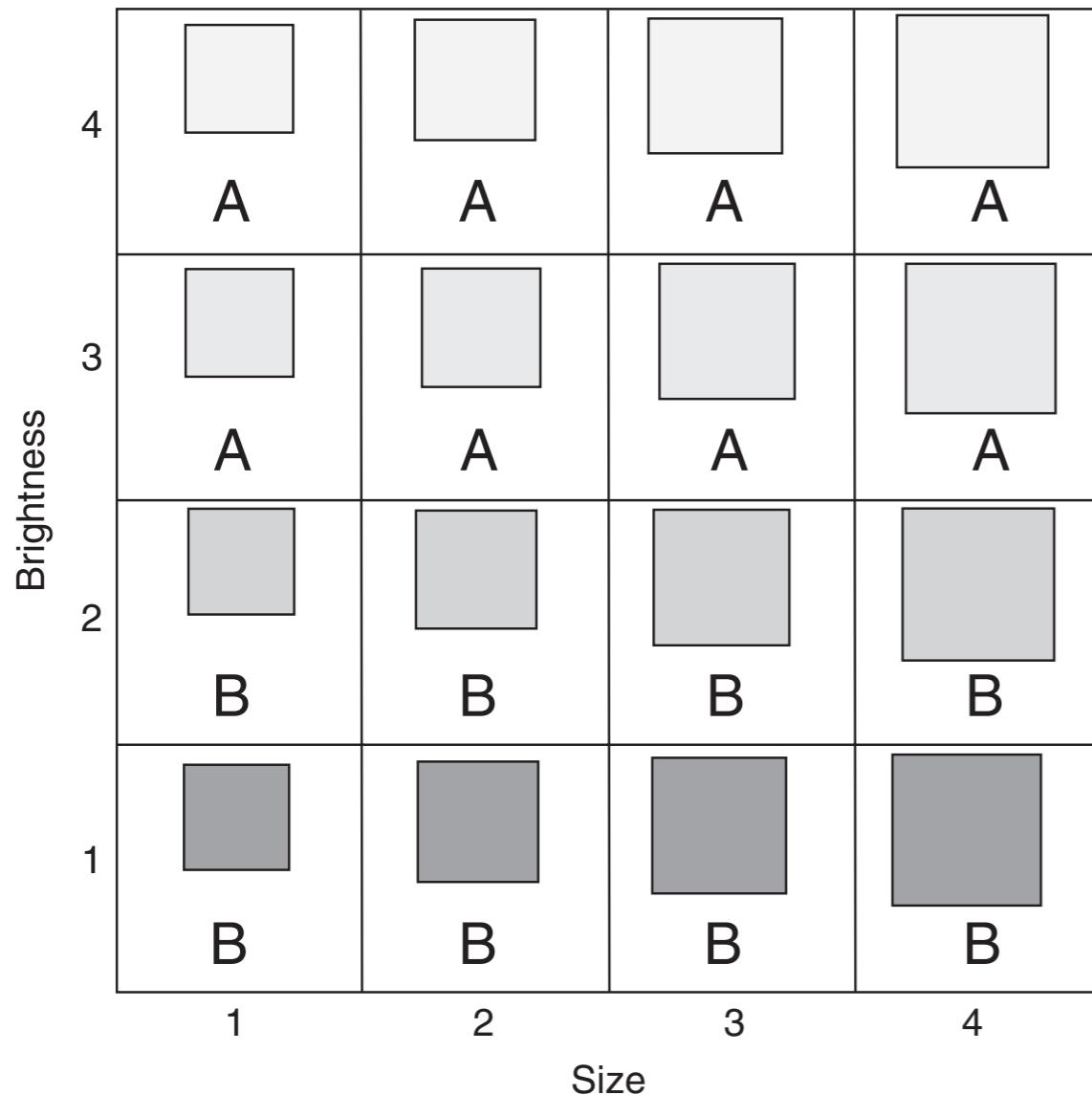


Acquired distinctiveness is attenuated compared to the other conditions, suggesting competition for attention

Separable vs. integral dimensions

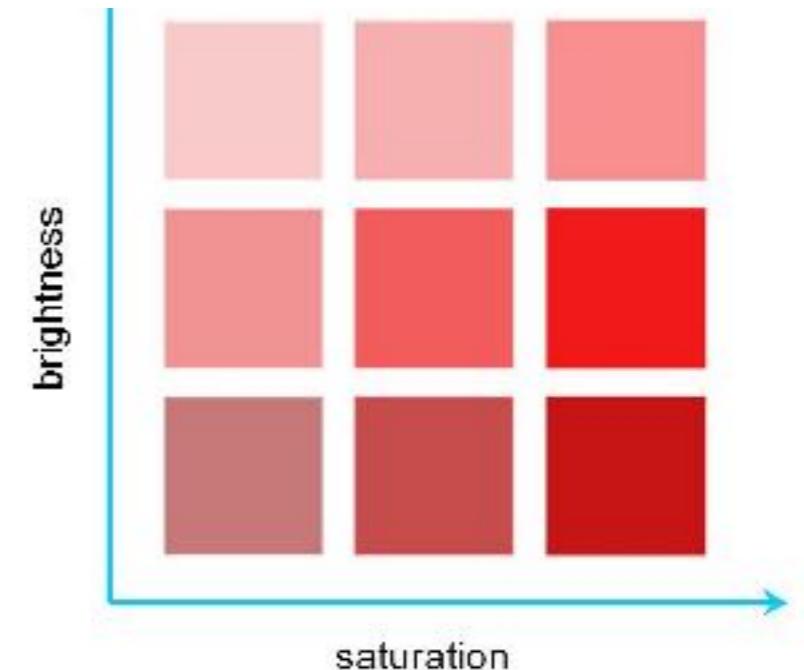
separable: dimensions are independent

e.g., brightness vs. size

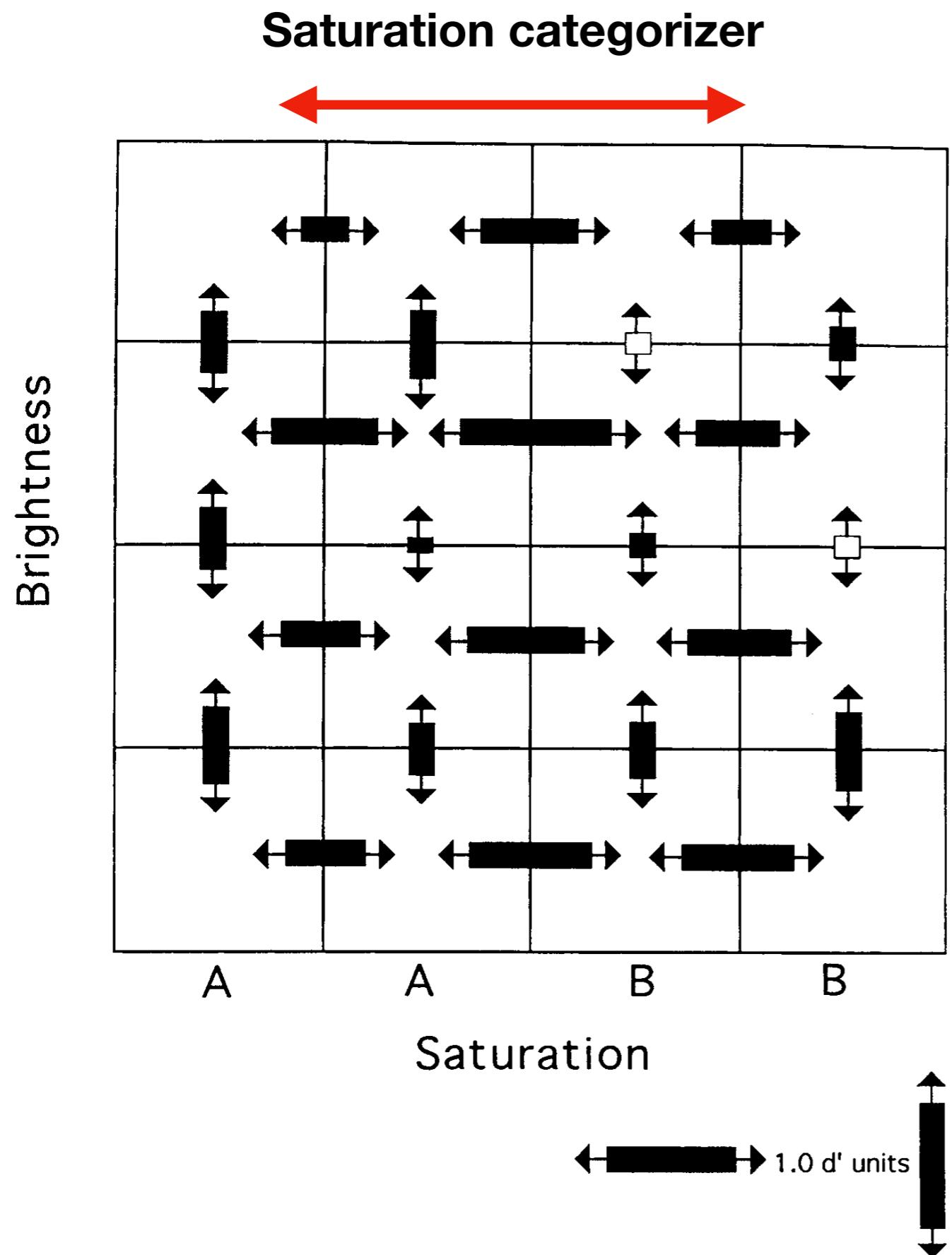


integral: dimensions interact

e.g., brightness vs. saturation



Goldstone Ex 4 : Results with integral dimensions



learning to categorize based on saturation causes acquired distinctive in brightness, and vice versa

Review: How does ALCOVE learn?

Learning is incremental fitting of the attention weights and association weights

Response rule

$$P(y \in A) = \frac{e^{\phi \mathbf{sim}(y, A)}}{e^{\phi \mathbf{sim}(y, A)} + e^{\phi \mathbf{sim}(y, B)}}$$

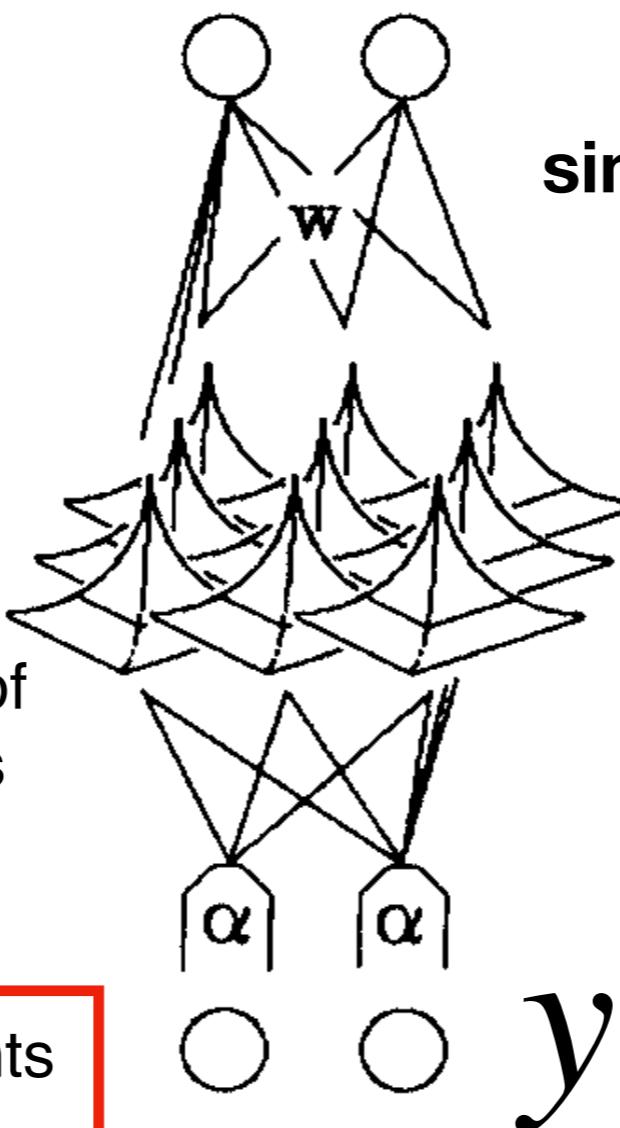
category A category B

$$\mathbf{sim}(y, x) = e^{\sum_{D_i} \alpha_i |x_i - y_i|}$$

x

denotes one of
the exemplars

α_i attention weights



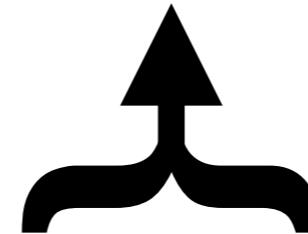
$$\mathbf{sim}(y, C) = \sum_{x \in X} w_{cx} \mathbf{sim}(y, x)$$

w_{cx} association weights
between exemplar x and
category c

y current stimulus

Review: Network before training

response $P(y \in A)$ 0.5



$$\mathbf{sim}(y, C) = \sum_{x \in X} w_{cx} \mathbf{sim}(y, x)$$

association weights w_{Ax}

0.0	0.0	0.0	0.0
1.0	0.11	0.11	0.01

$$\mathbf{sim}(y, x) = e^{\sum_{D_i} \alpha_i |x_i - y_i|}$$

w_{Ax}

$\mathbf{sim}(y, x)$

0.0	0.0	0.0	0.0
0.11	0.01	0.01	0.0

Exemplars \mathcal{X}

attention weights

α_{color}

0.33

$\alpha_{texture}$

0.33

α_{slash}

0.33

current stimulus

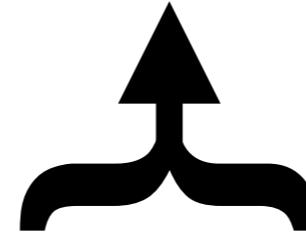
y



Review: Network after training

response $P(y \in A)$ 0.99

Response is now nearly perfect



$$\text{sim}(y, C) = \sum_{x \in X} w_{cx} \text{sim}(y, x)$$

association weights w_{Ax} 1.15 1.14 1.14 1.14

$\text{sim}(y, x)$	1.0	1.0	1.0	1.0

$$\text{sim}(y, x) = e^{\sum_{D_i} \alpha_i |x_i - y_i|}$$

w_{Ax} -1.14 -1.14 -1.14 -1.14

$\text{sim}(y, x)$	0.03	0.03	0.03	0.03

}

Exemplars \mathcal{X}

attention weights

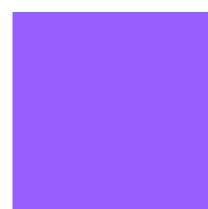
α_{color}
0.544

$\alpha_{texture}$
0.0

α_{slash}
0.0

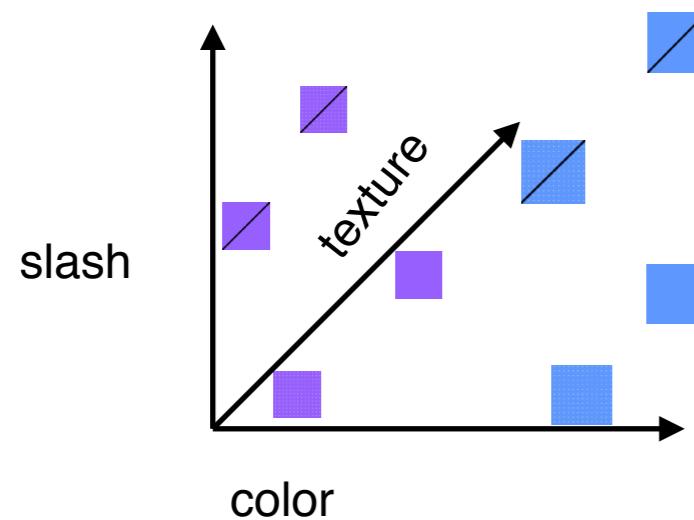
Attention only looks at color

current stimulus y

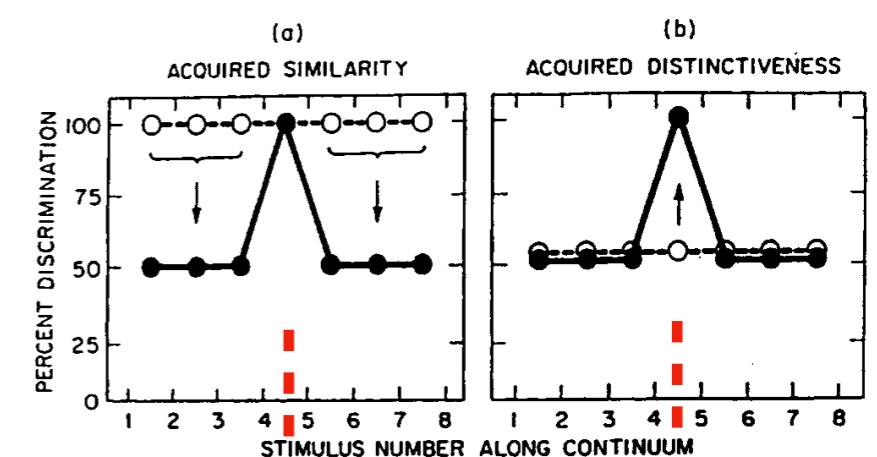


ALCOVE can explain acquired distinctiveness and acquired similarity, as applied across entire dimensions

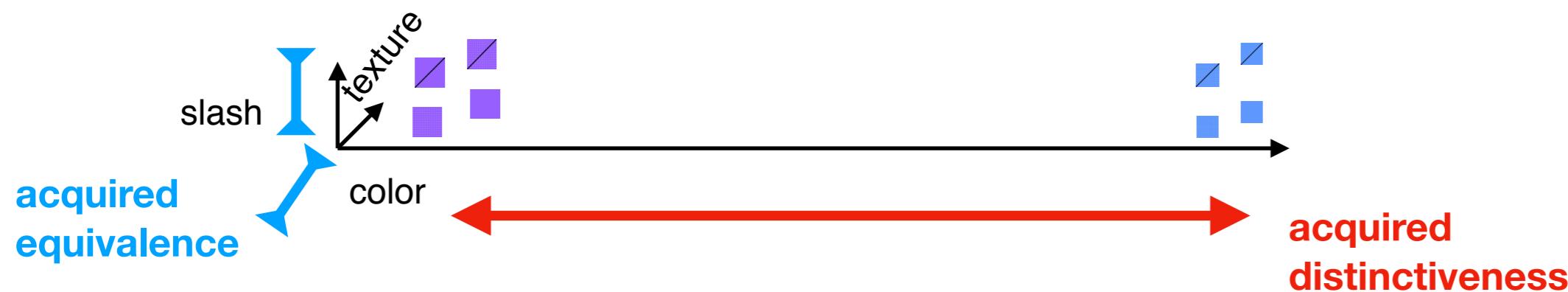
ALCOVE's similarity before training



But, ALCOVE's attention doesn't do this, but could if combined with category response



ALCOVE's similarity after training



$$P(\text{discrim } S_1 \text{ vs. } S_2) = \frac{1 + (P(A|S_1) - P(A|S_2))}{2}$$

Conclusion from Goldstone's results

- Goldstone mainly found **acquired distinctiveness**, such that differences between objects in different categories are emphasized
- ALCOVE does not necessarily predict discrimination judgements (it does categorization, not discrimination), but it is consistent with perceptual changes, especially when these effects operate across an *entire dimension*
- Different changes to perception depending on whether dimensions are integral or separable
- But this differs by domain: with speech and more complex stimuli, category learning can lead to both acquired distinctiveness and acquired similarity, and even acquired similarity within a class

Cross-Language Speech Perception: Evidence for Perceptual Reorganization During the First Year of Life*

JANET F. WERKER AND RICHARD C. TEES
University of British Columbia

Previous work in which we compared English infants, English adults, and Hindi adults on their ability to discriminate two pairs of Hindi (non-English) speech contrasts has indicated that infants discriminate speech sounds according to phonetic category without prior specific language experience (Werker, Gilbert, Humphrey, & Tees, 1981), whereas adults and children as young as age 4 (Werker & Tees, *in press*), may lose this ability as a function of age and/or linguistic experience. The present work was designed to (a) determine the generalizability of such a decline by comparing adult English, adult Salish, and English infant subjects on their perception of a new non-English (Salish) speech contrast, and (b) delineate the time course of the developmental decline in this ability. The results of these experiments replicate our original findings by showing that infants can discriminate nonnative speech contrasts without relevant experience, and that there is a decline in this ability during ontogeny. Furthermore, data from both cross-sectional and longitudinal studies show that this decline occurs within the first year of life, and that it is a function of specific language experience.

infants speech perception cross-language decline

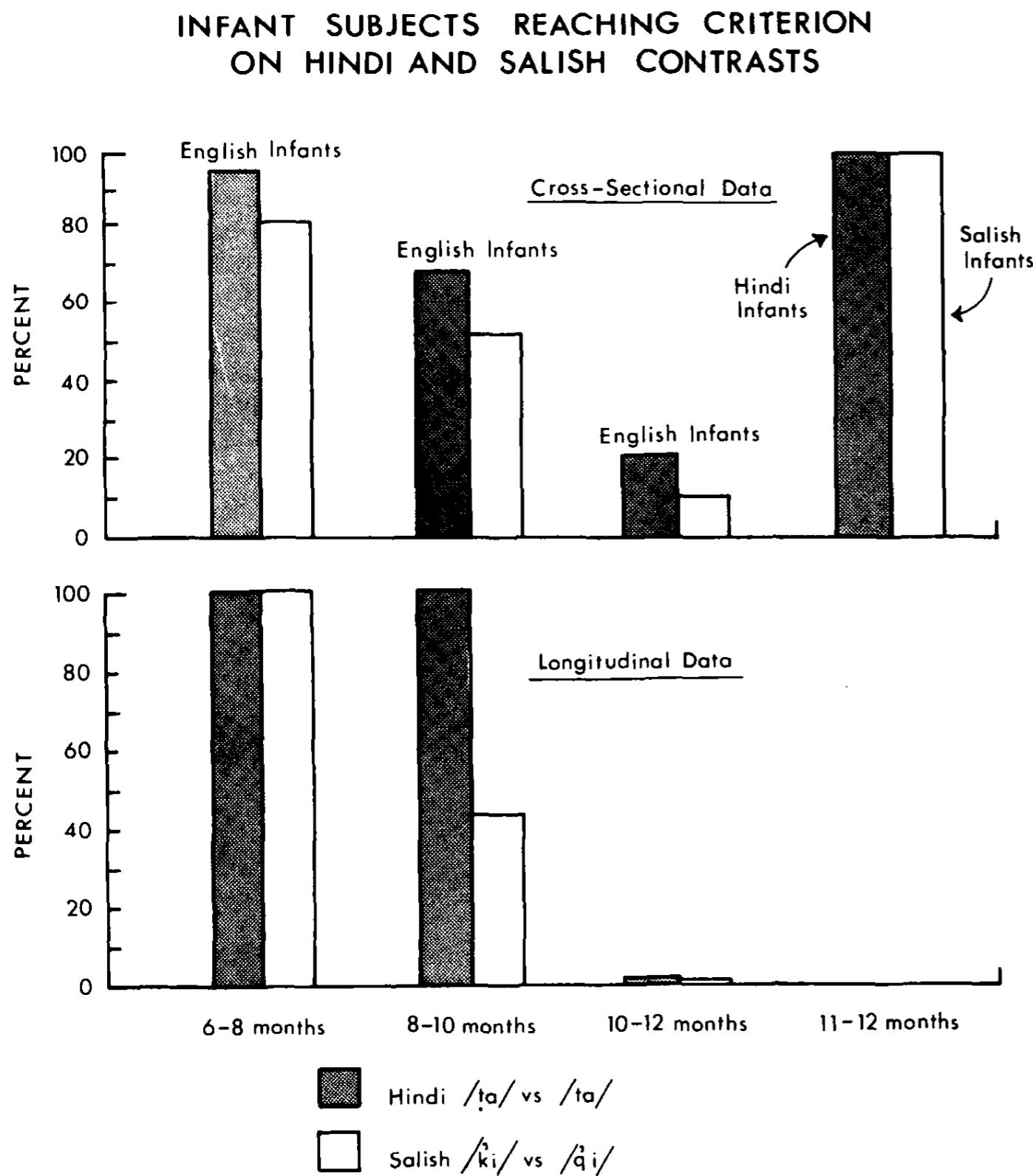


Janet Werker
University of British Columbia

While a large (but finite) number of sound segments occur in the languages of the world, only a subset is used phonemically (to differentiate meaning) in any particular language. Several researchers have predicted that human infants are born with the ability to discriminate the universal set of phonetic contrasts regardless of language experience, and that this ability declines as a function of specific linguistic experience (Eimas, 1978; Morse, 1978; Werker et al., 1981). Alternatively, it has been proposed that experience listening to a language may be necessary to facilitate the perception of the phonetic distinctions used in that language (Eilers, Gavin, & Wilson, 1979). Most relevant data support the first of these predictions, suggesting that rather than having to learn to differentiate phonetic features, young infants seem to respond to speech sounds according to the categories that could serve as the basis for adult phonemic

* This work was jointly supported by grants to Richard C. Tees from the Social Sciences and Humanities Research Council (410-81-0796), the National Research Council (PA0179) of Canada, and the National Institute of Mental Health (1R03NH35829), and by NICHD Grant HD12420 to Haskins Laboratories. We thank the infants and mothers who made this study possible. We also thank Kathy Searcy, Sue Tees, and Carole Bawden for their assistance. Special thanks to Al Liberman for making us welcome at Haskins Laboratories. Requests for reprints should be sent to Janet F. Werker, Department of Psychology, Dalhousie University, Halifax, Nova Scotia, B3H 4J1, or to Richard C. Tees, Department of Psychology, University of British Columbia, Vancouver, BC, V6T 1Y7, Canada.

Young infants can discriminate universal set of phonetic contrasts (6 mo), but they lose non-native contrasts between 8-12 mo



- Unlike in Goldstone's experiment, speech sound learning is characterized by declining discrimination in first year of life
- Procedure: Infants habituated to specific sound category, and rewarded for head turn when sound category changes
- Criterion based on succeeding on 8/10 change trials, and discriminating /ba/ vs /da/ before and after non-native contrast test
- Exp 2 (cross-sectional) with infants of different ages, and Exp 3 (longitudinal) with same infants, with similar results
- Infants lose non-native contrasts between 8-12 mo (acquired similarity)

Categorization Creates Functional Features

Philippe G. Schyns
University of Glasgow

Luc Rodet
University of Grenoble

Many theories of object recognition and categorization claim that complex objects are represented in terms of characteristic features. The origin of these features has been neglected in theories of object categorization. Do they form a fixed and independent set that exists before experience with objects, or are features progressively extracted and developed as an organism categorizes its world? This article maintains that features can be learned flexibly as a consequence of categorizing and representing objects. All 3 experiments reported in this article used categories of unfamiliar, computer-synthesized 2-dimensional objects ("Martian cells"). The results showed that varying the order of category learning induced the creation of different features that changed the perceptual appearance and the featural representation of identical category exemplars. Network simulations supported a flexible rather than a fixed-feature interpretation of the data.

Many theories of object recognition and categorization assume that objects are represented in memory as groups of components. To classify an object, one must first identify its components and then compare them to memory representations. For example, when a person sees a cup, he or she might first identify a container or a handle before categorizing the object properly. Of course, not all components of an object are necessary for its categorization, but many of them are probably identified during the recognition process. Componential accounts that embody this general approach include theories of object recognition and categorization (see, among others, Biederman, 1987; Marr & Nishihara, 1978; Rosch & Mervis, 1975; E. E. Smith & Medin, 1981; Treisman & Gelade, 1980).

Although most object categorization theories are componential, they have paid less attention to the origin of the components themselves. Do these components, or features, form a fixed and independent vocabulary that exists prior to the experience with objects, or are features progressively learned and developed as an organism categorizes and represents its world? Most current models of category learning leave aside the issue of feature learning and feature development. Their feature set is fixed and unaffected by the classification and learning processes.

Classification and learning processes, however, operate on a stable featural analysis (a perceptual organization) of

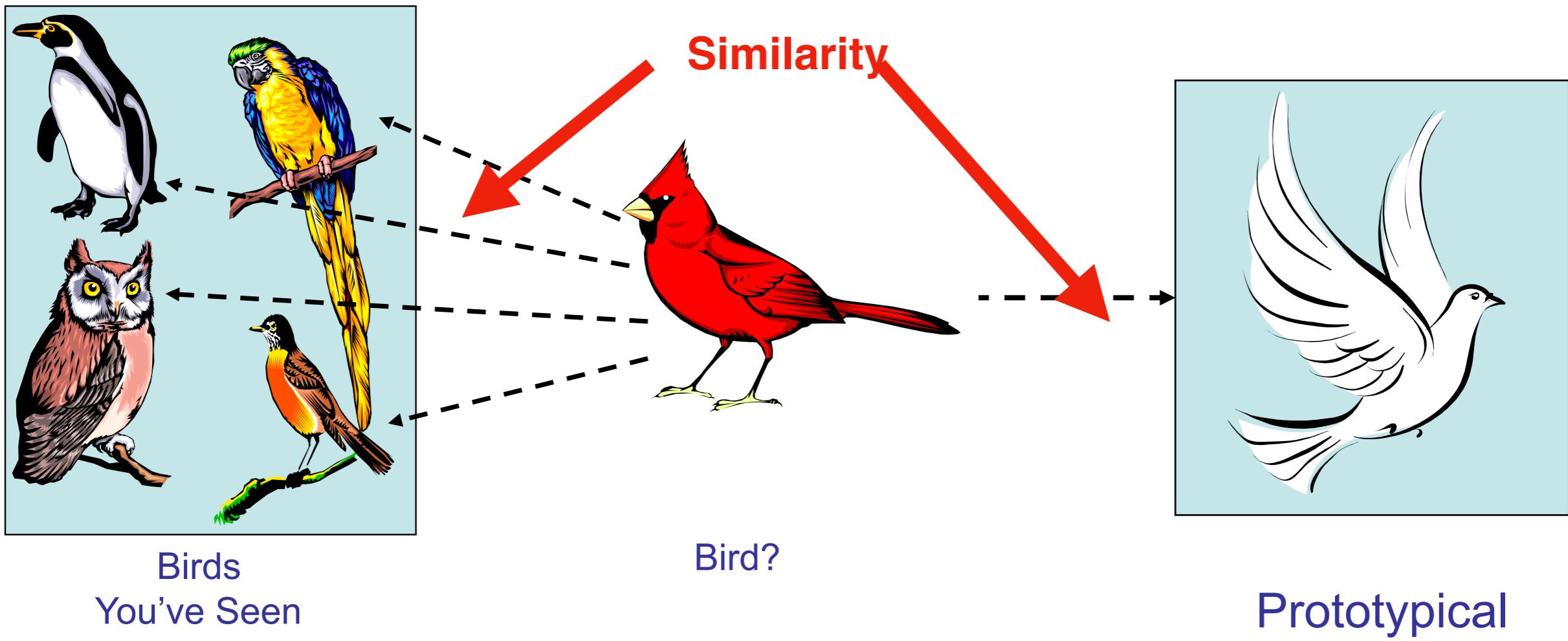
the ever-changing retinal input. Even though our sophisticated visual apparatus probably comes equipped with a priori ways of analyzing and organizing retinal images, there are occasions when a relevant perceptual analysis is not readily available. For example, complete novices reading chest X-rays (e.g., Christensen et al., 1981), sexing chickens (Biederman & Shiffrar, 1987), and categorizing dermatoses (Norman, Brooks, Coblenz, & Babcock, 1992) have little understanding of the relevant dimensional structure of these categories. Even when told what the signs of different diagnosis are, novices are not always able to see the features experts use to organize the input. If one takes a developmental perspective, it seems clear that infants and young children are not always able to analyze objects by using all the stimulus dimensions that are used by adults (C. Smith, Carey, & Wiser, 1985; L. B. Smith & Kemler, 1978; Ward, 1983).

Thus, there is suggestive evidence that features are flexible—that is, they adjust to the perceptual experience and the categorization history of the individual. Flexible features open the possibility that the same input is differently perceived and analyzed before being categorized. Hence, a complete theory of categorization and conceptual development should not only explain the ways in which object features are combined to form concepts, it should also explain the origin and the development of the features participating in the analysis of the input. The studies presented in this article investigate further the claim that a significant part of learning a category involves learning the



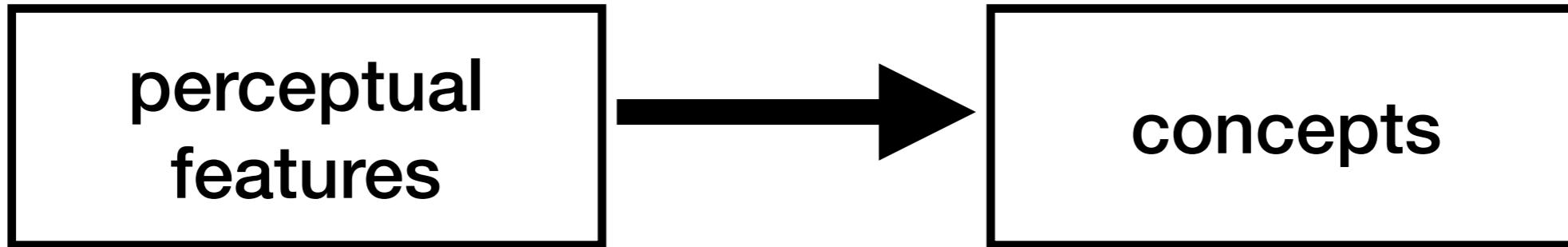
Philippe Schyns
University of Glasgow

Review: What counts as a feature?



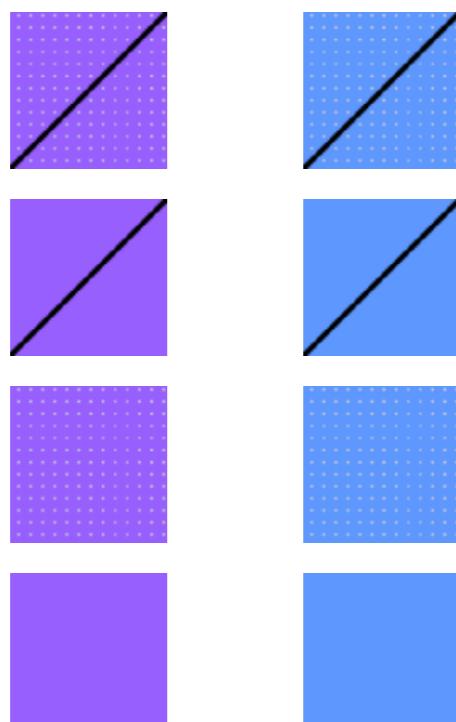
- What counts as a feature? (Murphy & Medin, 1985)
 - To change the importance of age, we could include features for "around 10 years ago," "around 100 years ago," "1000 years ago", etc.
 - To change importance of size, we could include "smaller than the earth," "smaller than a country", "smaller than a city," etc.
- It is difficult to establish the "respects for similarity" (Medin, Goldstone, & Gentner, 1993, *Psych Rev*)

Traditional “fixed feature” account of category learning

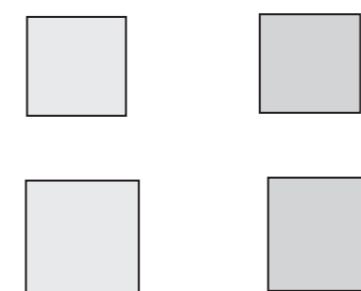


- Most category learning experiments have clear-cut stimuli with clear dimensions of variation
- Most models don't address where the features come from — instead they are provided. This can be characterized as a “fixed feature view”

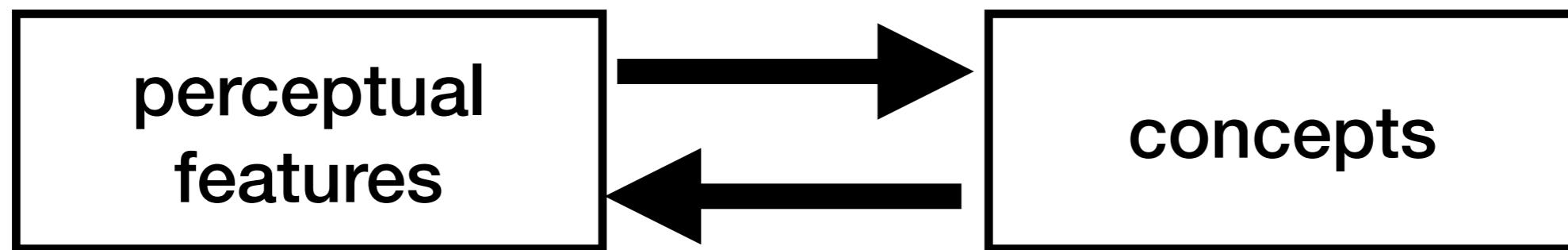
category A category B



category A category B



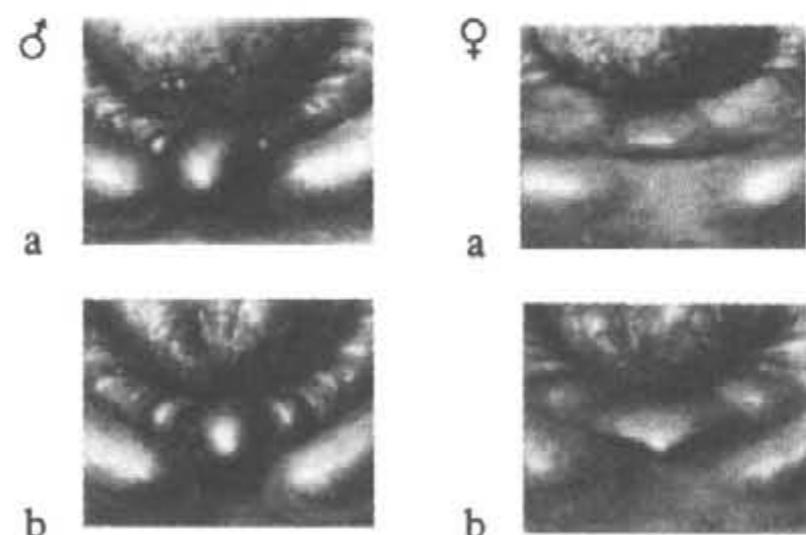
Can categorization create new features?



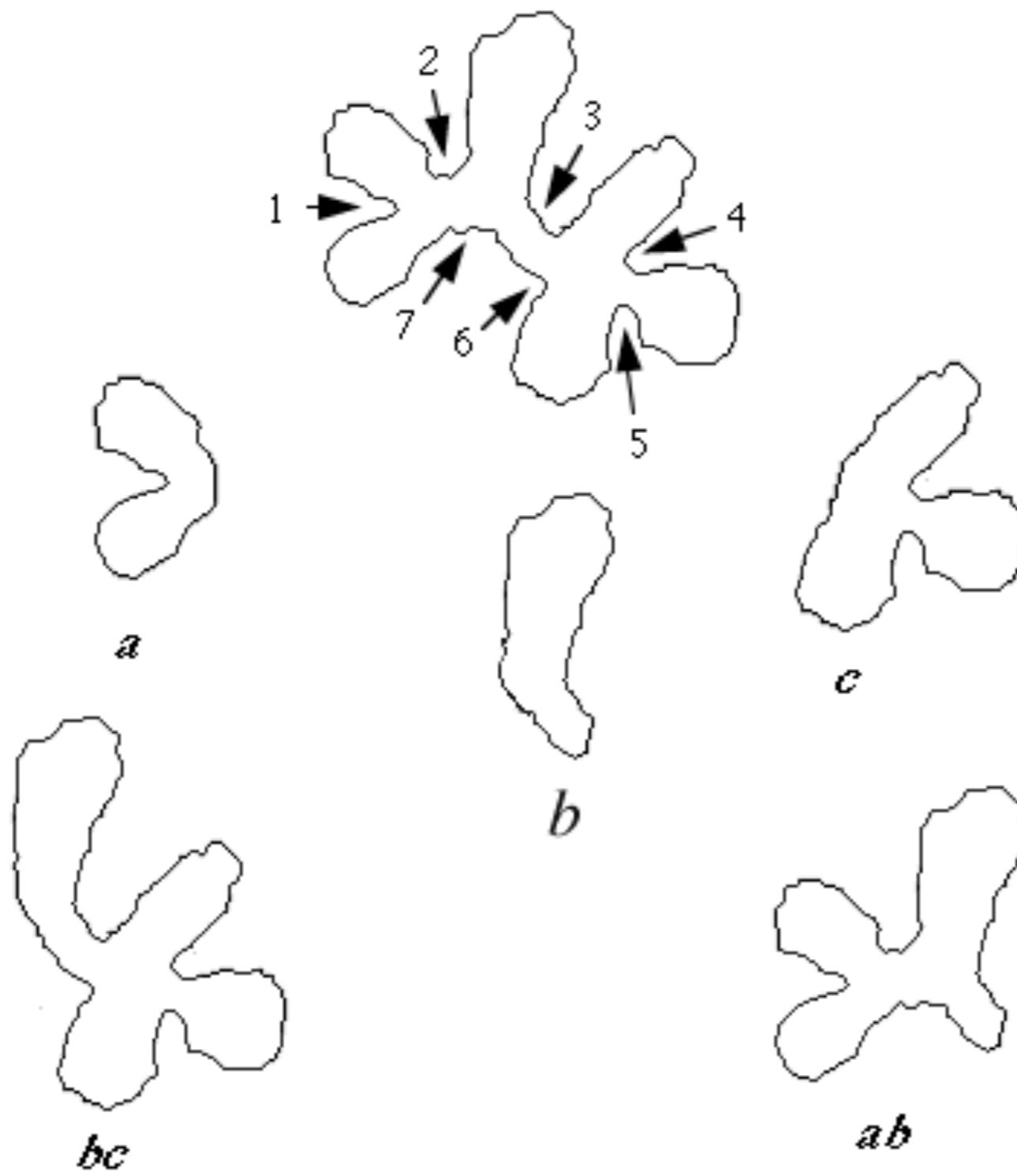
- Alternative is a “flexible feature view” : if a fragment of a stimulus categories objects, that fragment is instanced as a unit in the representational code
- Suggestive evidence reading chest X-Rays and sexing chickens seems to require experience and specialized features



*Figure 2. An accepted grasp for chick sexing. (Modified from “Chick Sexing” by J. H. Lunn, 1948, *American Scientist*, 36, pp. 280–287. Copyright 1948 by the American Scientist. Photograph by the University of Minnesota Photographic Laboratory. Adapted by permission.)*



Schyns and Rodet : Preliminary experiment



Schyns and Rodet : Preliminary experiment

study phase with no explicit task besides study stimuli

Group A-BC

shown 5 examples



A

shown 5 examples



BC



Group AB-C

shown 5 examples



AB

shown 5 examples



C

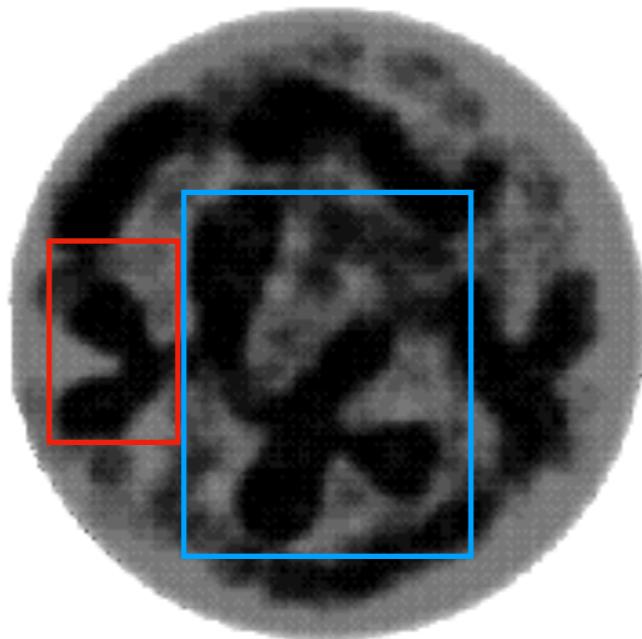


Schyns and Rodet : Pilot Results

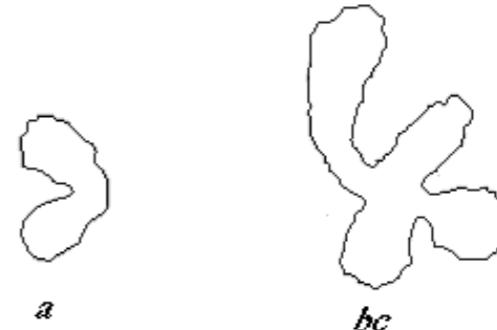
How do participants parse 5 novel test cells?

“draw outlines around the parts they saw during preexposure”

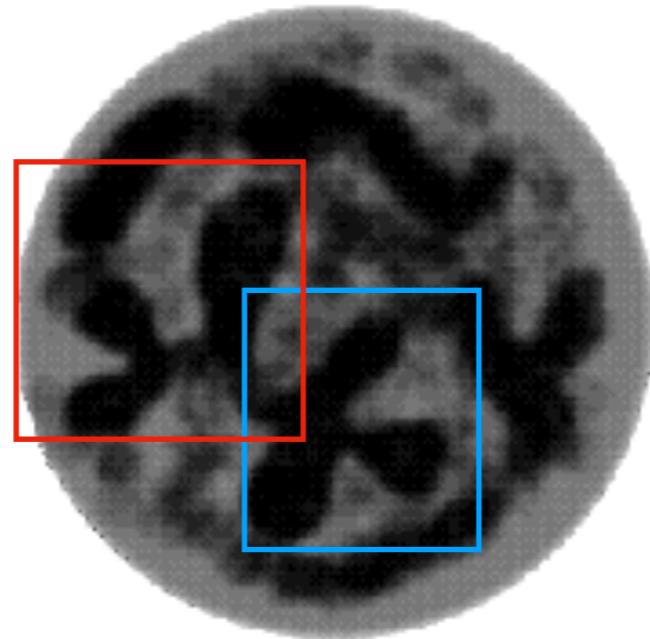
Group A-BC



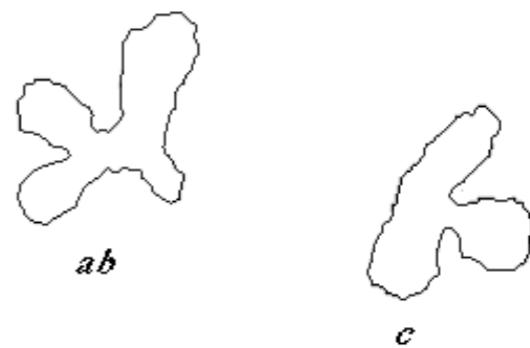
83% consistent with this parse



Group AB-C



100% consistent with this parse

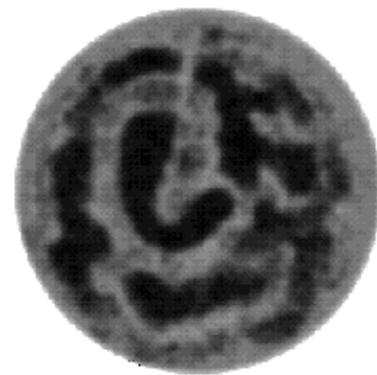


Schyns and Rodet : Experiment 2

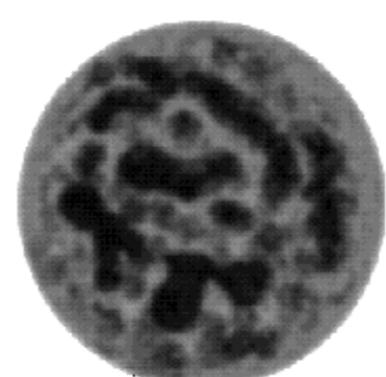
Explicit categorization task

Group X → Y → XY

Category 1
(10 examples)



Category 2
(10 examples)



Category 3
(10 examples)



**Hypothetical
feature library**



library is [X, Y]

Group XY → X → Y

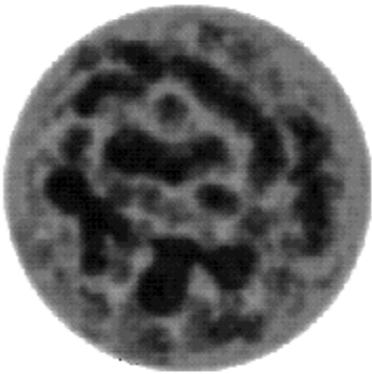
Category 1
(10 examples)



Category 2
(10 examples)



Category 3
(10 examples)



XY X Y



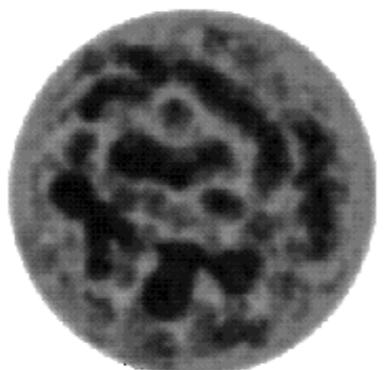
library is [XY, X, Y]

Schyns and Rodet : Experiment 2

How will they classify a novel stimulus X-Y?

Group X → Y → XY

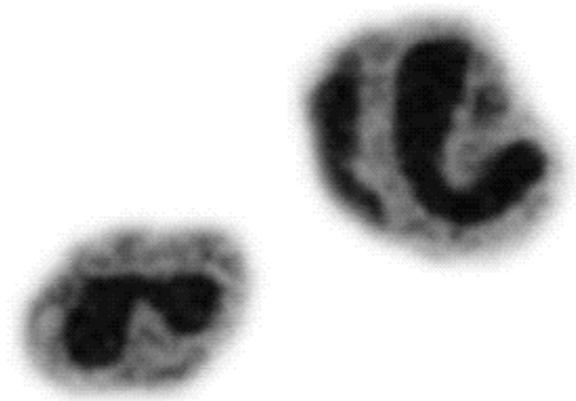
Category 1 Category 2 Category 3



“category representation”



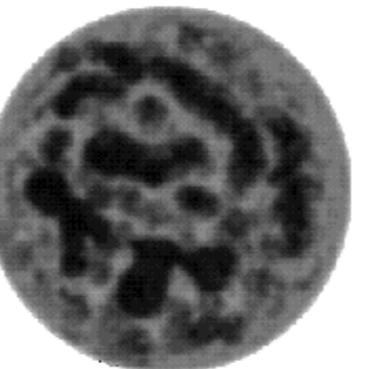
Stimulus X-Y



Hypothesis:
labeled “Category 3”

Group XY → X → Y

Category 1 Category 2 Category 3



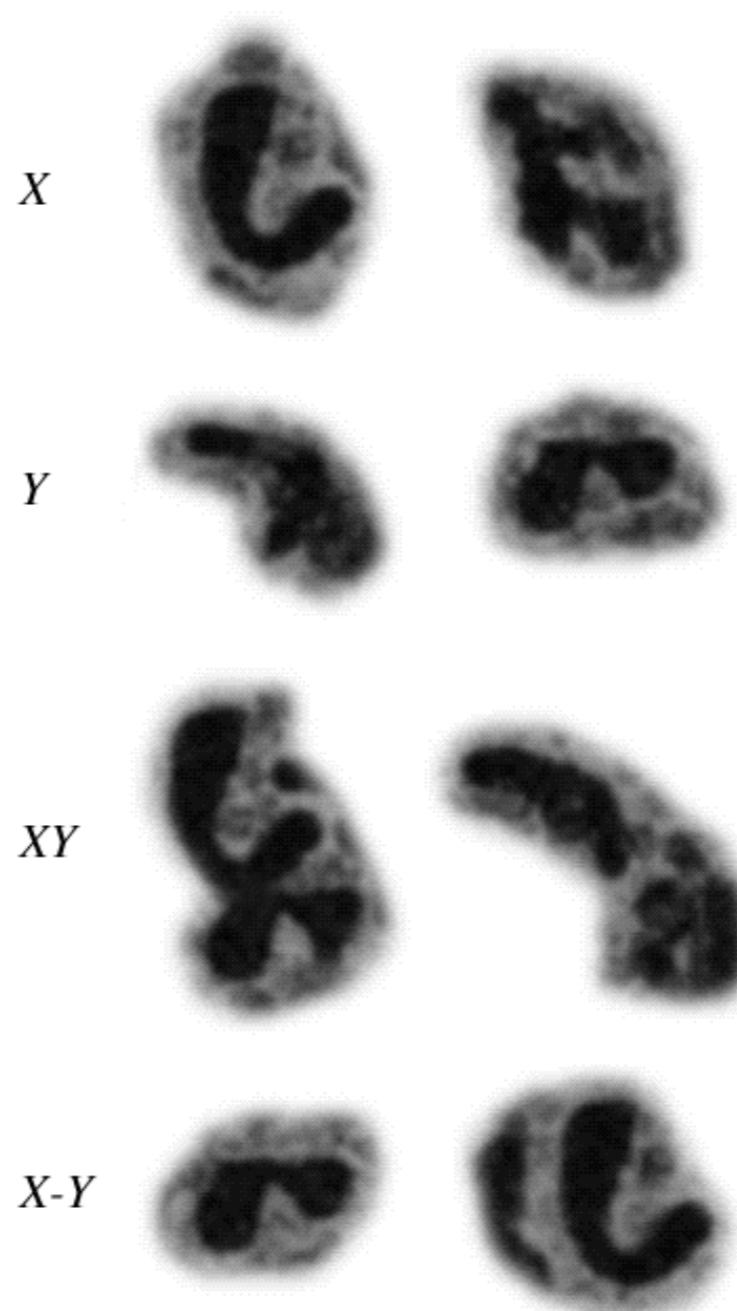
“category representation”



Hypothesis:
labeled “Category 1”
or “Category 2”

Schyns and Rodet : Experiment 2

Note stimuli were presented in two views, “as if two glimpses in a microscope” (does this matter??)

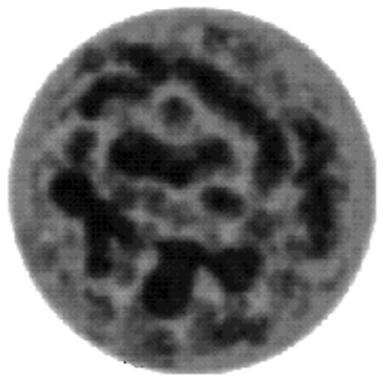


Schyns and Rodet : Experiment 2 Results

How will they classify a novel stimulus X-Y?

Group X → Y → XY

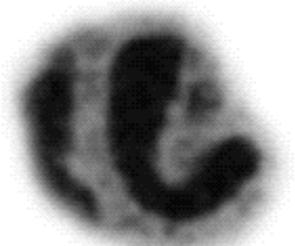
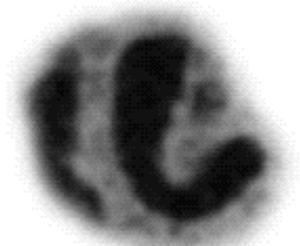
Category 1 Category 2 Category 3



“category representation”



Stimulus X-Y



**labeled “Category 3”
88% of time**

Group XY → X → Y

Category 1 Category 2 Category 3



“category representation”



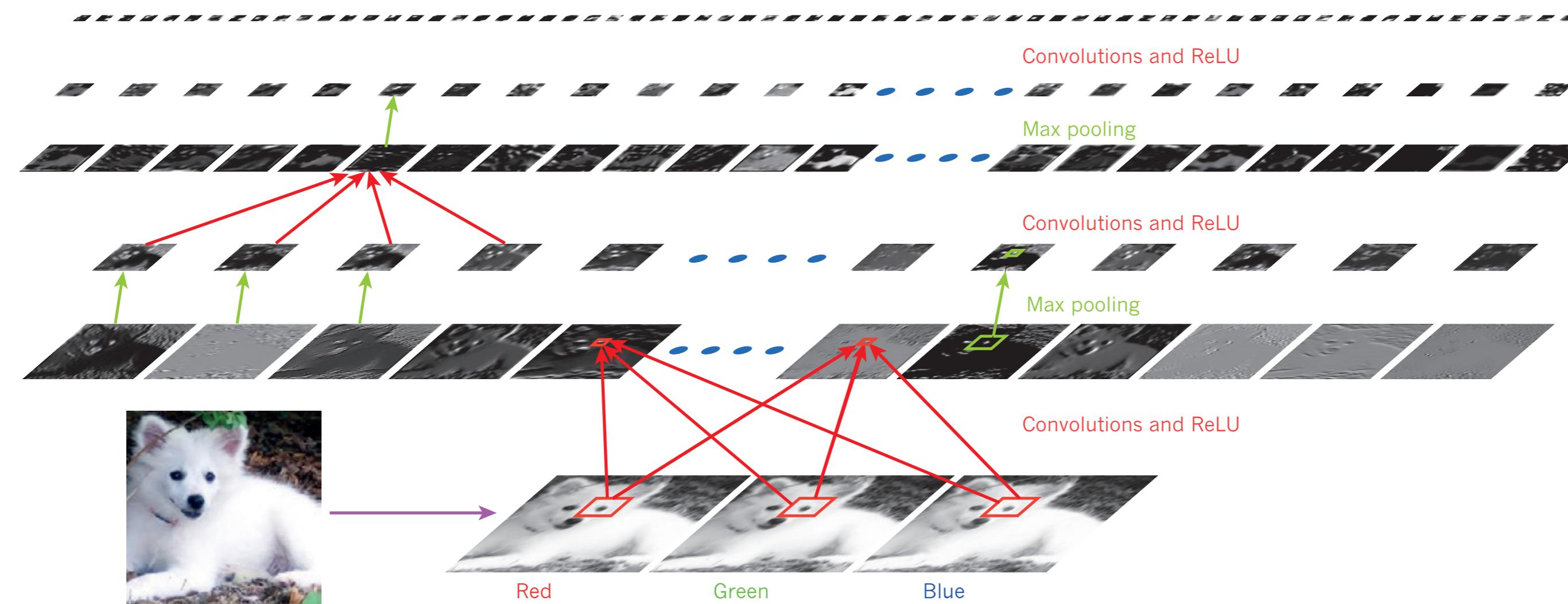
**labeled either
“Category 1”
or “Category 2”
79% of time**

Schyns and Rodet : Conclusions

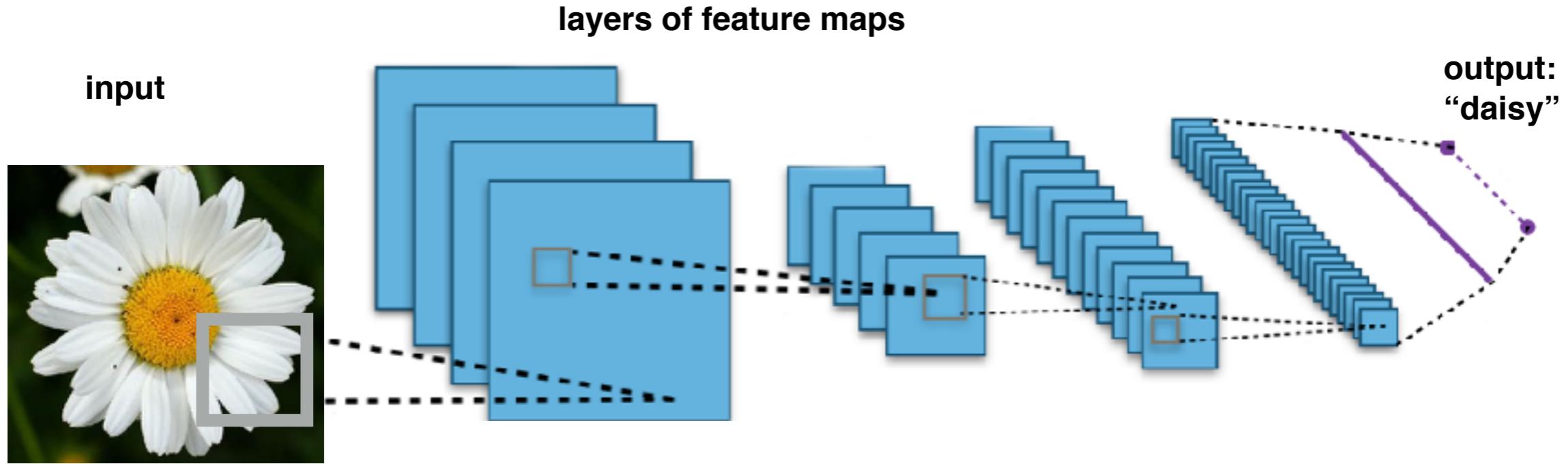
- features are flexible and develop with categorization experience, to influence perception of subsequent categorization examples
- Unlike theory-based concepts and the knowledge view, which also propose flexible features (or feature weighting), they emphasize the perceptual changes that accompany feature creation
- Beyond ALCOVE and almost every other model, more than “feature reweighting” will be required to understand how categorization influences perception

Review: Deep convolutional neural network

Samoyed (16); Papillon (5.7); Pomeranian (2.7); Arctic fox (1.0); Eskimo dog (0.6); white wolf (0.4); Siberian husky (0.4)



Review: Deep convolutional neural network



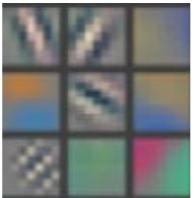
- These models learn from the raw input stimuli
- **Critically, these models learn their features.** They do not assume a fixed feature decomposition (although they assume features are translation invariant, and other inductive biases)
- They discover **functionally-relevant features** given the task at hand

Discovered functional features

Training data (ImageNet)

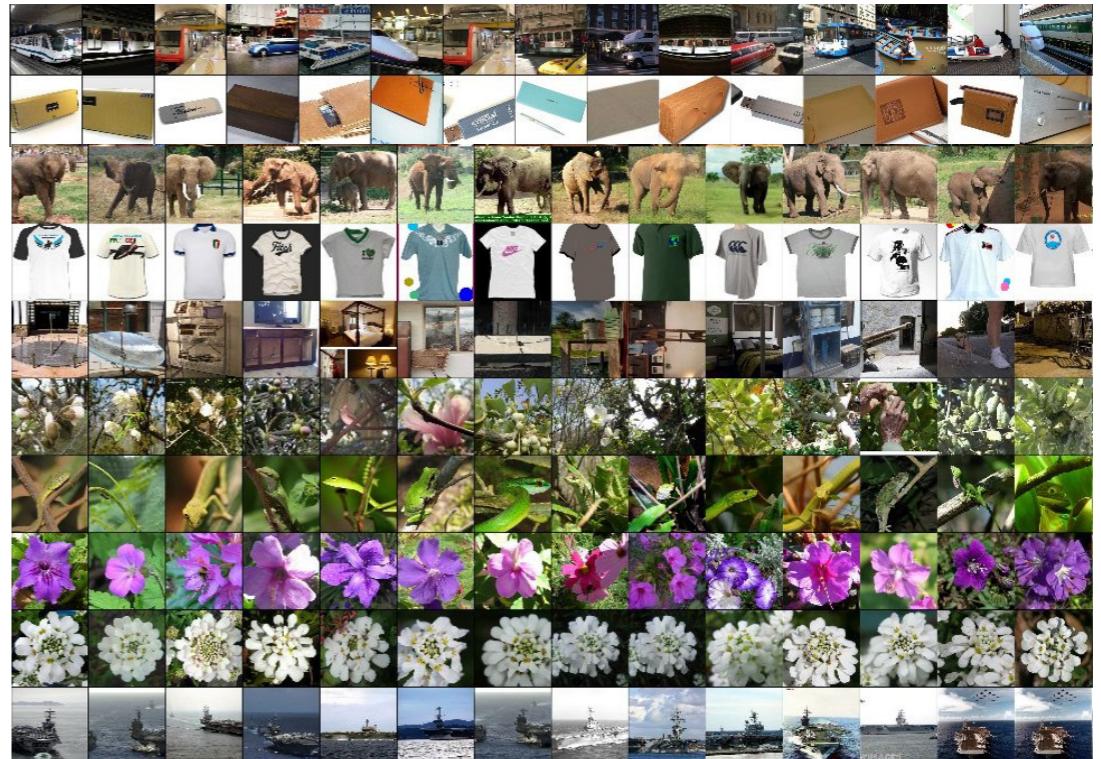
- Usually trained on ImageNet
- 1.2 million images with labels
- 1000 categories

Raw filters



Layer 1

Image patches that maximally activate each filter



(Zeiler & Fergus, 2014)

Discovered functional features

16 different filters

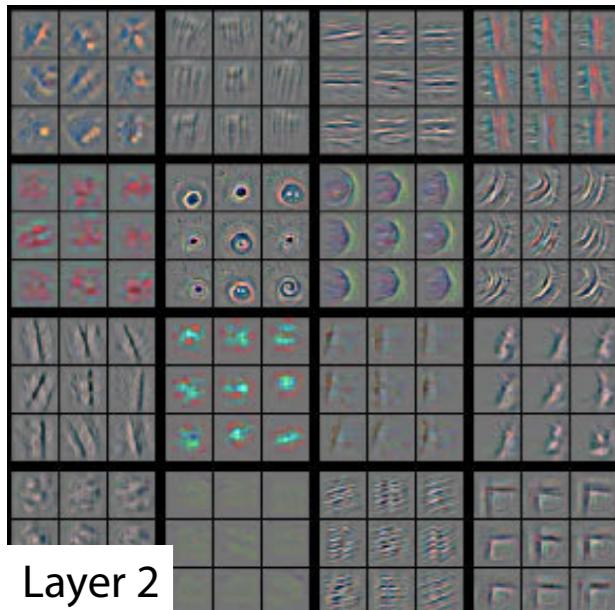
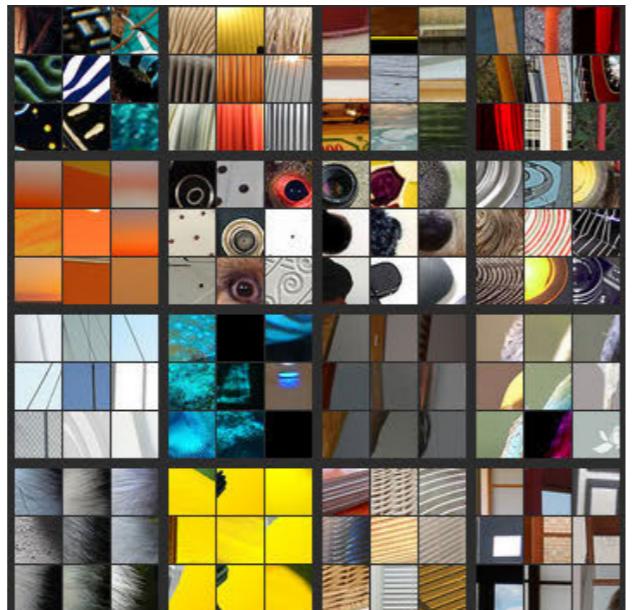


Image patches that
maximally activate
each filter



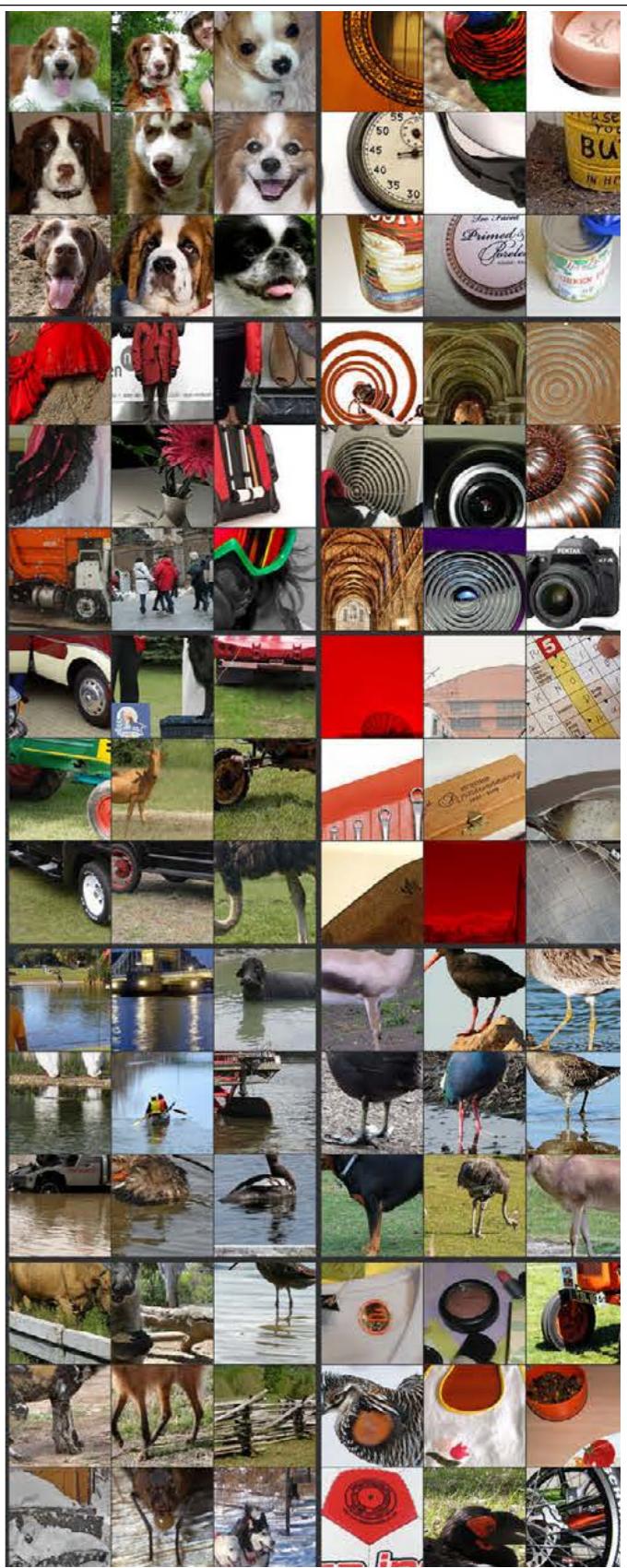
Layer 2

10 different filters



Layer 4

Image patches that
maximally activate
each filter

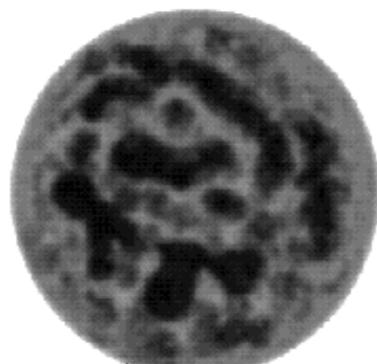


Potential explantation through modern neural network models?

How will they classify a novel stimulus X-Y?

Group X → Y → XY

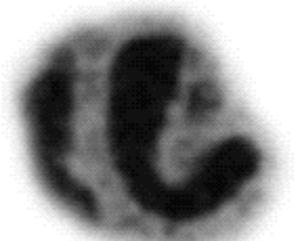
Category 1 Category 2 Category 3



“category representation”



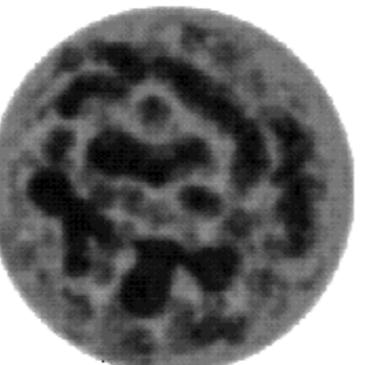
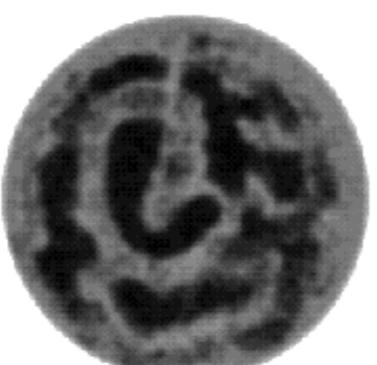
Stimulus X-Y



**labeled “Category 3”
88% of time**

Group XY → X → Y

Category 1 Category 2 Category 3



“category representation”



**labeled either
“Category 1”
or “Category 2”
79% of time**