# Part III

# Learning structured representations
# Hierarchical Bayesian models

Universal Grammar

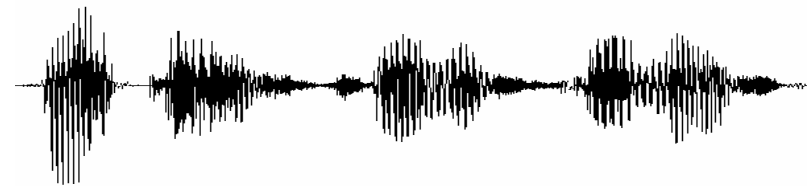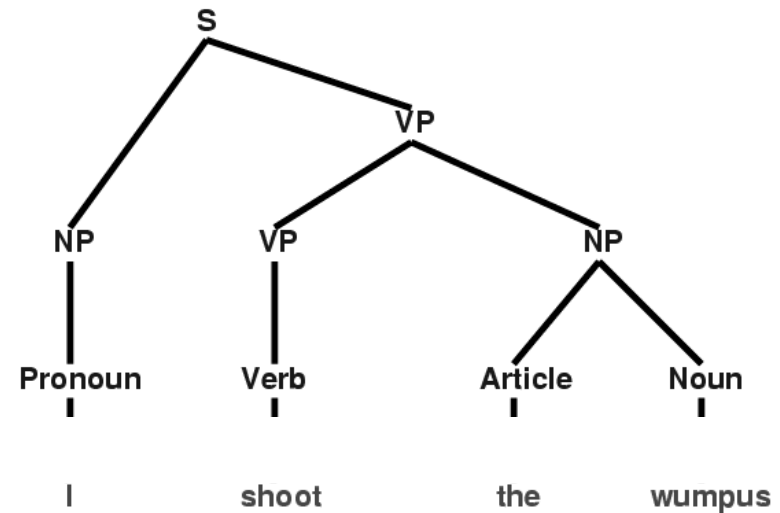Hierarchical phrase structure grammars (e.g., CFG, HPSG, TAG)

Grammar

$S \rightarrow NP\ VP$

$NP \rightarrow Det\ [Adj]\ Noun\ [RelClause]$

$RelClause \rightarrow [Rel]\ NP\ V$

$VP \rightarrow VP\ NP$

$VP \rightarrow Verb$

Phrase structure

Utterance

Speech signal

# Outline

- Learning structured representations
  - grammars
  - logical theories

- Learning at multiple levels of abstraction

# A historical divide

**Structured Representations**

**Innate knowledge**

vs

**Unstructured Representations**

**Learning**

(Chomsky,
 Pinker,
 Keil, ...)

(McClelland,
 Rumelhart, ...)

**Structured
Representations**

Chomsky
Keil

**Structure
Learning**
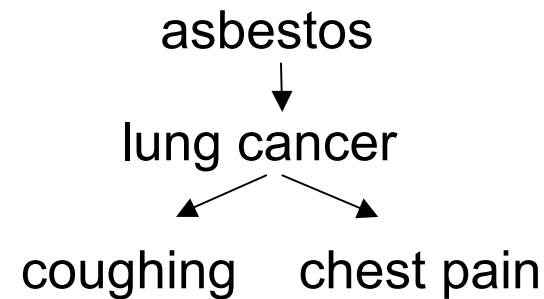
**Innate
Knowledge**

**Learning**

McClelland,
Rumelhart

**Unstructured
Representations**

# Representations

**Causal networks**

asbestos

$\downarrow$

lung cancer

coughing    chest pain

**Grammars**

$S \rightarrow NP\,VP$

$NP \rightarrow Det\,[\,Adj\,]\,Noun\,[\,RelClause\,]$

$RelClause \rightarrow [\,Rel\,]\,NP\,V$

$VP \rightarrow VP\,NP$

$VP \rightarrow Verb$

**Logical theories**

$\forall x\,y\;\mathsf{Sibling}(x,y) \leftarrow \mathsf{Sibling}(y,x)$

$\forall x\,y\;\mathsf{Ancestor}(x,y) \leftarrow \mathsf{Parent}(x,y)$

# Representations

Phonological rules

$$\begin{bmatrix} +\text{syllabic} \\ -\text{consonantal} \end{bmatrix} \rightarrow \begin{bmatrix} +\text{back} \end{bmatrix} / \begin{bmatrix} +\text{back} \\ +\text{syllabic} \\ -\text{consonantal} \end{bmatrix} \begin{bmatrix} +\text{consonantal} \end{bmatrix}^* \underline{\quad}$$

Semantic networks

# How to learn a R

- Search for R that maximizes

$$P(R|\mathbf{Data}) \propto P(\mathbf{Data}|R)P(R)$$

- Prerequisites
  - Put a prior over a hypothesis space of Rs.
  - Decide how observable data are generated from an underlying R.

# How to learn a R ~~a R~~ anything

- Search for R that maximizes

$$P(R|\textbf{Data}) \propto P(\textbf{Data}|R)P(R)$$

- Prerequisites
  - Put a prior over a hypothesis space of Rs.
  - Decide how observable data are generated from an underlying R.
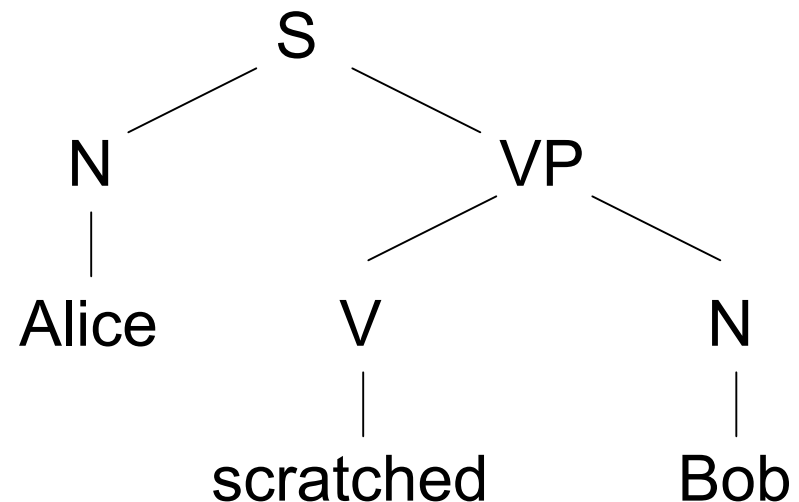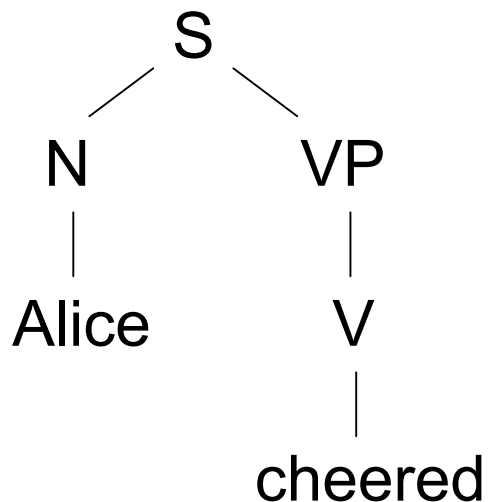
# Context free grammar

S → N VP    VP → V    N → "Alice"    V → "scratched"

VP → V N    N → "Bob"    V → "cheered"

```
        S                          S
       / \                        / \
      N   VP                     N   VP
      |   |                      |   / \
    Alice V                    Alice V   N
          |                          |   |
       cheered                   scratched Bob
```

# Probabilistic context free grammar

**1.0**
S → N VP

**0.6**
VP → V

**0.4**
VP → V N

**0.5**
N → "Alice"

**0.5**
N → "Bob"

**0.5**
V → "scratched"

**0.5**
V → "cheered"



S **1.0**
N **0.5**   VP **0.6**
Alice
V
cheered

probability = 1.0 * 0.5 * 0.6
= 0.3



S **1.0**
N **0.5**   VP **0.4**
Alice
V **0.5**   N **0.5**
scratched   Bob

probability = 1.0*0.5*0.4*0.5*0.5
= 0.05

# The learning problem

Grammar G:

$$S \xrightarrow{1.0} N\ VP \qquad VP \xrightarrow{0.6} V \qquad N \xrightarrow{0.5} \text{"Alice"} \qquad V \xrightarrow{0.5} \text{"scratched"}$$

$$VP \xrightarrow{0.4} V\ N \qquad N \xrightarrow{0.5} \text{"Bob"} \qquad V \xrightarrow{0.5} \text{"cheered"}$$

---

Data D:

Alice scratched.              Alice cheered.
Bob scratched.                Bob cheered.
Alice scratched Alice.        Alice cheered Alice.
Alice scratched Bob.          Alice cheered Bob.
Bob scratched Alice.          Bob cheered Alice.
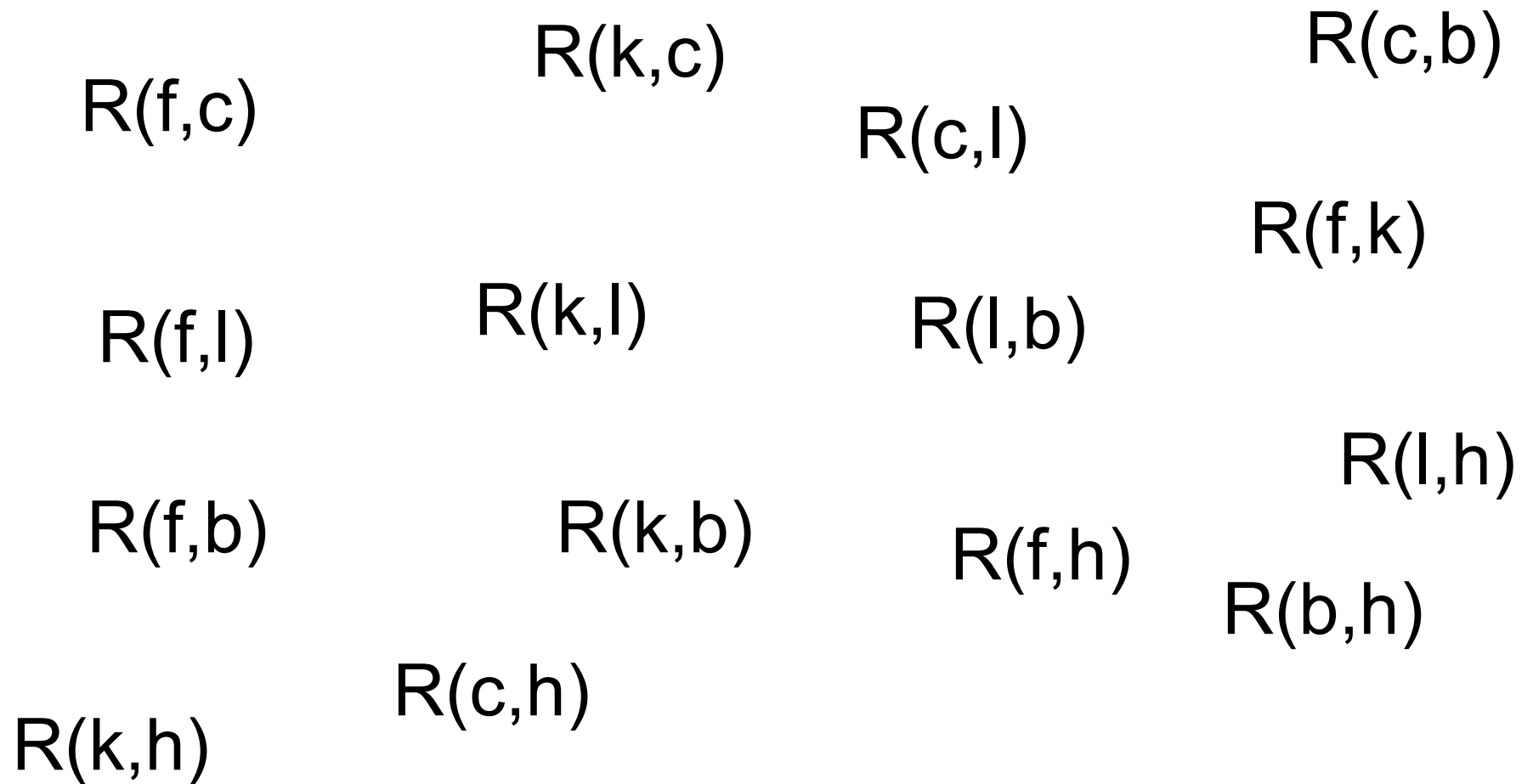Bob scratched Bob.            Bob cheered Bob.

# Grammar learning

- Search for G that maximizes

$$P(G|\text{Data}) \propto P(\text{Data}|G)P(G)$$

- Prior: $P(G) \propto 2^{-\text{length}(G)}$

- Likelihood: $P(\text{Data}|G)$
  – assume that sentences in the data are independently generated from the grammar.

(Horning 1969; Stolcke 1994)

# Experiment

```
S  --> NP VP
NP --> Det N
VP --> Vt NP
   --> Vc PP
   --> Vi
PP --> P NP
Det --> a
    --> the
Vt --> touches
   --> covers
Vc --> is
Vi --> rolls
   --> bounces
N  --> circle
   --> square
   --> triangle
P  --> above
   --> below
```

- Data: 100 sentences

```
the circle covers a square
a square is above the triangle
a circle bounces
```

⋮

(Stolcke, 1994)

## Generating grammar:

```
S  --> NP VP
NP --> Det N
VP --> Vt NP
   --> Vc PP
   --> Vi
PP --> P NP
Det --> a
    --> the
Vt --> touches
   --> covers
Vc --> is
Vi --> rolls
   --> bounces
N  --> circle
   --> square
   --> triangle
P  --> above
   --> below
```

## Model solution:

```
S  --> NP VP
NP --> Det N
VP --> VI
   --> X NP
X  --> VT
   --> VC P
Det --> a
    --> the
Vt --> touches
   --> covers
Vc --> is
Vi --> rolls
   --> bounces
N  --> circle
   --> square
   --> triangle
P  --> above
   --> below
```

# Predicate logic

- A compositional language

$$\forall x\, y\, \text{Sibling}(x, y) \leftarrow \text{Sibling}(y, x)$$

For all x and y, if y is the sibling of x then x is the sibling of y

$$\forall x\, y\, z\, \text{Ancestor}(x, z) \leftarrow \text{Ancestor}(x, y) \wedge \text{Ancestor}(y, z)$$

For all x, y and z, if x is the ancestor of y and y is the ancestor of z, then x is the ancestor of z.

# Learning a kinship theory

Theory T:

$$\forall x\, y\; \text{Sibling}(x, y) \leftarrow \text{Sibling}(y, x)$$

$$\forall x\, y\, z\; \text{Ancestor}(x, z) \leftarrow \text{Ancestor}(x, y) \wedge \text{Ancestor}(y, z)$$

$$\forall x\, y\; \text{Ancestor}(x, y) \leftarrow \text{Parent}(x, y)$$

$$\forall x\, y\, z\; \text{Uncle}(x, z) \leftarrow \text{Brother}(x, y) \wedge \text{Parent}(y, z)$$

---

Data D:

Sibling(victoria, arthur),   Sibling(arthur,victoria),

Ancestor(chris,victoria),   Ancestor(chris,colin),

Parent(chris,victoria),     Parent(victoria,colin),

Uncle(arthur,colin),        Brother(arthur,victoria)    …

(Hinton, Quinlan, …)

# Learning logical theories

- Search for T that maximizes

$$P(T|\textbf{Data}) \propto P(\textbf{Data}|T)P(T)$$

- Prior: $P(T) \propto 2^{-\text{length}(T)}$

- Likelihood: $P(\textbf{Data}|T)$
  - assume that the data include all facts that are true according to T

(Conklin and Witten; Kemp et al 08; Katz et al 08)

# Theory-learning in the lab

R(f,c)

R(k,c)

R(c,b)

R(c,l)

R(f,k)

R(f,l)

R(k,l)

R(l,b)

R(l,h)

R(f,b)

R(k,b)

R(f,h)

R(b,h)

R(c,h)

R(k,h)

(cf Krueger 1979)

# Theory-learning in the lab

Transitive:    R(f,k). R(k,c). R(c,l). R(l,b). R(b,h).

R(X,Z) ← R(X,Y), R(Y,Z).

---

| f,k | f,c | f,l | f,b | f,h |
|-----|-----|-----|-----|-----|
|     | k,c | k,l | k,b | k,h |
|     |     | c,l | c,b | c,h |
|     |     |     | l,b | l,h |
|     |     |     |     | b,h |

Learning time

Complexity

trans.

trans.

Theory length

trans.

(Kemp et al 08)

# Conclusion: Part 1

- Bayesian models can combine structured representations with statistical inference.

# Outline

- Learning structured representations
  - grammars
  - logical theories

- Learning at multiple levels of abstraction

# Vision



(Han and Zhu, 2006)

# Motor Control



symbolic representation
of tasks e.g. goal

mid-level representation
e.g. sequences of elements

low level dynamics
e.g. elements of
movements

(Wolpert et al., 2003)

# Causal learning

chemicals
↓
diseases
↓
symptoms

Schema

↓

Causal
models

asbestos
↓
lung cancer
↙    ↘
coughing   chest pain

mercury
↓
minamata disease
↓
muscle wasting

Contingency
Data

Patient 1: asbestos exposure, coughing, chest pain

Patient 2: mercury exposure, muscle wasting

(Kelley; Cheng; Waldmann)

Universal Grammar

Hierarchical phrase structure grammars (e.g., CFG, HPSG, TAG)

↓ P(grammar | UG)

Grammar

$S \rightarrow NP\ VP$

$NP \rightarrow Det\ [Adj]\ Noun\ [RelClause]$

$RelClause \rightarrow [Rel]\ NP\ V$

$VP \rightarrow VP\ NP$

$VP \rightarrow Verb$

P(phrase structure | grammar)

Phrase structure

```
               S
              / \
             /   VP
            /   / \
          NP  VP   NP
          |   |   / \
      Pronoun Verb Article Noun
          |   |    |     |
          I  shoot the  wumpus
```

P(utterance | phrase structure)

Utterance

↓ P(speech | utterance)

Speech signal

# Hierarchical Bayesian model

Universal Grammar

$\downarrow$ P(G|U)

Grammar

$\downarrow$ P(s|G)

Phrase structure

$\downarrow$ P(u|s)

Utterance

$U$

$\downarrow$

$G$

$s_1 \quad s_2 \quad s_3 \quad s_4 \quad s_5 \quad s_6$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

$u_1 \quad u_2 \quad u_3 \quad u_4 \quad u_5 \quad u_6$

A hierarchical Bayesian model specifies a joint distribution over all variables in the hierarchy:

$$P(\{u_i\}, \{s_i\}, G \mid U)$$

$$= P(\{u_i\} \mid \{s_i\}) \, P(\{s_i\} \mid G) \, P(G|U)$$

# Top-down inferences

Universal Grammar $\rightarrow$ U

$\downarrow$

Grammar $\rightarrow$ G

$\downarrow$

Phrase structure $\rightarrow$ $s_1$ $s_2$ $s_3$ $s_4$ $s_5$ $s_6$

$\downarrow$

Utterance $\rightarrow$ $u_1$ $u_2$ $u_3$ $u_4$ $u_5$ $u_6$

Infer $\{s_i\}$ given $\{u_i\}$, G:

$$P(\{s_i\} \mid \{u_i\}, G) \propto P(\{u_i\} \mid \{s_i\}) P(\{s_i\} \mid G)$$

# Bottom-up inferences

Universal Grammar        U

Grammar        G

Phrase structure     $s_1$    $s_2$    $s_3$    $s_4$    $s_5$    $s_6$

Utterance       $u_1$    $u_2$    $u_3$    $u_4$    $u_5$    $u_6$

Infer G given $\{s_i\}$ and U:

$$P(G| \{s_i\}, U) \; \alpha \; P(\{s_i\} | G) \, P(G|U)$$

# Simultaneous learning at multiple levels

Universal Grammar             U

Grammar             G

Phrase structure     $s_1$   $s_2$   $s_3$   $s_4$   $s_5$   $s_6$

Utterance         $u_1$   $u_2$   $u_3$   $u_4$   $u_5$   $u_6$

Infer G and $\{s_i\}$ given $\{u_i\}$ and U:

$$P(G, \{s_i\} \mid \{u_i\}, U) \; \alpha \; P(\{u_i\} \mid \{s_i\}) P(\{s_i\} \mid G) P(G \mid U)$$

# Word learning

Words in general

Individual words

Data

Whole-object bias
Shape bias

car          monkey          duck          gavagai

# A hierarchical Bayesian model



- Qualitative physical knowledge (symmetry) can influence estimates of continuous parameters ($F_H$, $F_T$).

- Explains why 10 flips of 200 coins are better than 2000 flips of a single coin: more informative about $F_H$, $F_T$ .

# Word Learning

"This is a dax."            "Show me the dax."

- 24 month olds show a shape bias
- 20 month olds do not

(Landau, Smith & Gleitman)

# Is the shape bias learned?

- Smith et al (2002) trained 17-month-olds on labels for 4 artificial categories:

"wib"

"lug"

- After 8 weeks of training 19-month-olds show the shape bias:

"zup"

"div"

"This is a dax."

"Show me the dax."

# Learning about feature variability



(cf. Goodman)

# Learning about feature variability



(cf. Goodman)

# A hierarchical model

Meta-constraints

Bags in general

Bag proportions

Data

M

Color varies across bags
but not much within bags

mostly
red

mostly
yellow

mostly
brown

mostly
green

…

mostly
blue?

…

# A hierarchical Bayesian model

Meta-constraints

M

Within-bag variability

Bags in general

$\alpha$ = 0.1

$\beta$ = [0.4, 0.4, 0.2]

Bag proportions

[1,0,0]   [0,1,0]   [1,0,0]   [0,1,0]   ...   [.1,.1,.8]

Data

[6,0,0]   [0,6,0]   [6,0,0]   [0,6,0]   ...   [0,0,1]

# A hierarchical Bayesian model

Meta-constraints      M     Within-bag
variability

$\alpha = 5$

Bags in general

$\beta = [0.4, 0.4, 0.2]$

Bag proportions   [.5,.5,0]   [.5,.5,0]   [.5,.5,0]   [.5,.5,0]   ⋯   [.4,.4,.2]

Data   [3,3,0]   [3,3,0]   [3,3,0]   [3,3,0]   ⋯   [0,0,1]

# Shape of the Beta prior

# A hierarchical Bayesian model

$$
\begin{array}{ll}
\alpha & \sim \text{Exponential}(\lambda) \\
\beta & \sim \text{Dirichlet}(\mathbf{1}) \\
\theta^i & \sim \text{Dirichlet}(\alpha\beta) \\
y^i & \sim \text{Multinomial}(\theta^i)
\end{array}
$$

Meta-constraints

M

Bags in general

$\alpha, \beta$

Bag proportions

$\theta^1 \qquad \theta^2 \qquad \theta^3 \qquad \theta^4 \qquad \cdots \qquad \theta^{\mathbf{n}}$

Data

$y^1 \qquad y^2 \qquad y^3 \qquad y^4 \qquad \qquad y^{\mathbf{n}}$

$$
p(\{y^i\}, \{\theta^i\}, \alpha, \beta | \lambda)
$$

# A hierarchical Bayesian model

$$
\begin{aligned}
\alpha &\sim \text{Exponential}(\lambda) \\
\boldsymbol{\beta} &\sim \text{Dirichlet}(\mathbf{1}) \\
\boldsymbol{\theta}^i &\sim \text{Dirichlet}(\alpha\boldsymbol{\beta}) \\
y^i &\sim \text{Multinomial}(\boldsymbol{\theta}^i)
\end{aligned}
$$

Meta-constraints

M

Bags in general

$\alpha, \boldsymbol{\beta}$

Bag proportions

$\theta^1 \quad \theta^2 \quad \theta^3 \quad \theta^4 \quad \cdots \quad \theta^n$

Data

$y^1 \quad y^2 \quad y^3 \quad y^4 \quad \cdots \quad y^n$

$$
p(\{\theta^i\}, \alpha, \boldsymbol{\beta} \mid \{y^i\}, \lambda)
$$

# Learning about feature variability



Meta-constraints

M

Categories in general

$\alpha, \beta$

Individual categories

$\theta^1$ $\theta^2$ $\theta^3$ $\theta^4$ $\theta^5$

Data

"wib"  "lug"  "zup"  "div"



| Category | 1 1 2 2 3 3 4 4 |
|----------|------------------|
| Shape    | 1 1 2 2 3 3 4 4 |
| Texture  | 1 2 3 4 5 6 7 8 |
| Color    | 1 2 3 4 5 6 7 8 |
| Size     | 1 2 1 2 1 2 1 2 |

"wib"   "lug"   "zup"   "div"

| Category | 1 1 2 2 3 3 4 4 | | 5 | ? | ? | ? |
|----------|-----------------|---|---|---|---|---|
| Shape    | 1 1 2 2 3 3 4 4 | | 5 | 5 | 6 | 6 |
| Texture  | 1 2 3 4 5 6 7 8 | | 9 | 10 | 9 | 10 |
| Color    | 1 2 3 4 5 6 7 8 | | 9 | 10 | 10 | 9 |
| Size     | 1 2 1 2 1 2 1 2 | | 1 | 1 | 1 | 1 |

"dax"

# Model predictions

# Where do priors come from?

Meta-constraints

M

Categories in general

$\alpha, \beta$

Individual categories

$\theta^1 \qquad \theta^2 \qquad \theta^3 \qquad \theta^4 \qquad \theta^5$

Data

# Knowledge representation



Mendeleev's Periodic Table of 1869[1]

# Children discover structural form

- Children may discover that
  - Social networks are often organized into cliques
  - The months form a cycle
  - "Heavier than" is transitive
  - Category labels can be organized into hierarchies

# A hierarchical Bayesian model

Meta-constraints

↓

Form

↓

Structure

↓

Data

M

↓

Tree

↓

mouse
squirrel
chimp
gorilla

↓

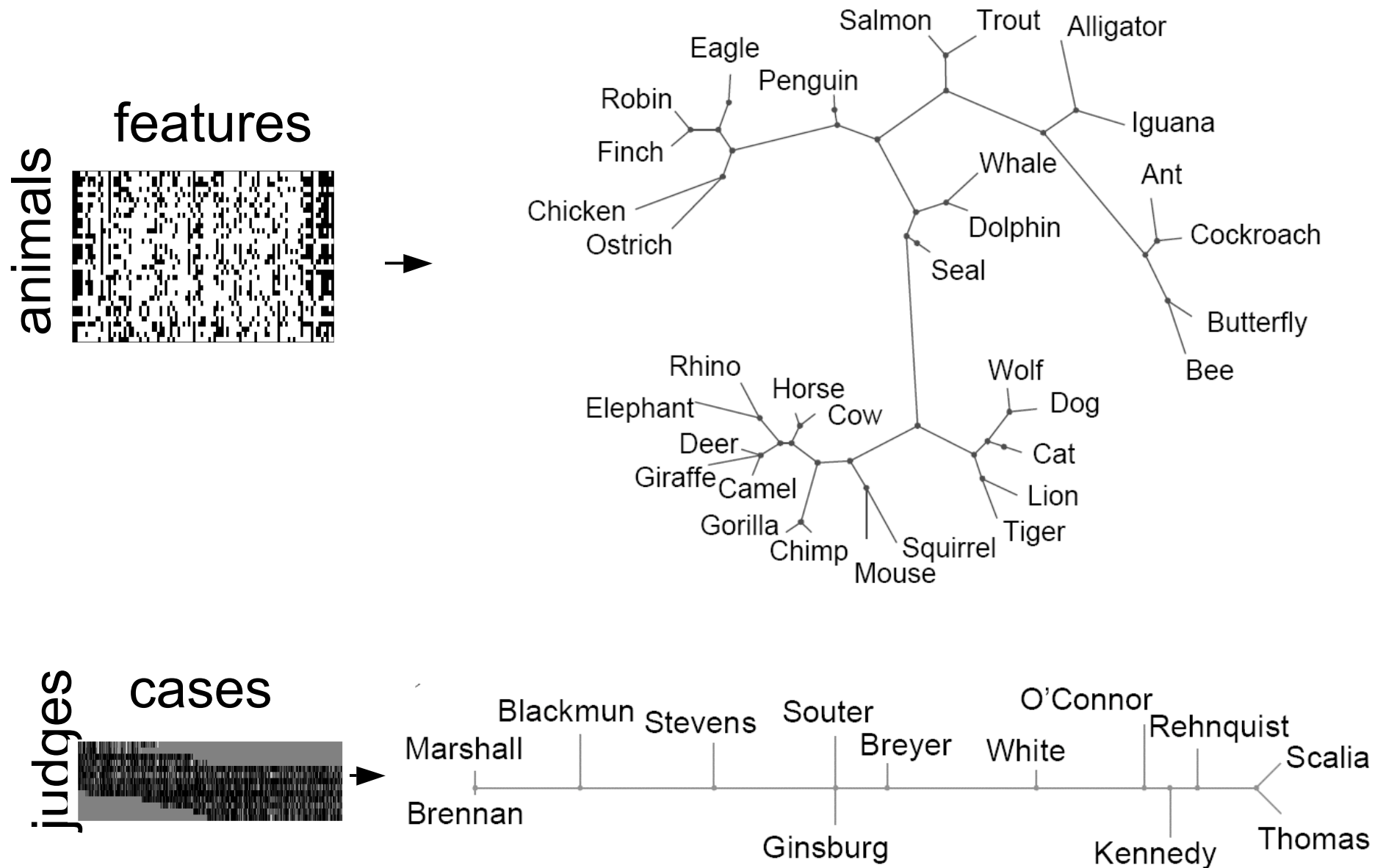|  | whiskers | hands | tail |
|---|---|---|---|
| mouse | ● | ○ | ● |
| squirrel | ● | ○ | ● |
| chimp | ○ | ● | ○ |
| gorilla | ○ | ● | ○ |

# A hierarchical Bayesian model

Meta-constraints          M

F: form

Tree

S: structure



D: data

| | whiskers | hands | tail |
|---|---|---|---|
| mouse | ● | ○ | ● |
| squirrel | ● | ○ | ● |
| chimp | ○ | ● | ○ |
| gorilla | ○ | ● | ○ |

$$P(S, F | D, n) \propto P(D|S)P(S|F,n)P(F)$$

# Structural forms



Partition     Order     Chain     Ring

Hierarchy     Tree     Grid     Cylinder

# P(S|F,n): Generating structures



- Each structure is weighted by the number of nodes it contains:

$$P(S|F) \propto \begin{cases} 0 & \text{if S inconsistent with F} \\ \theta(1-\theta)^{|S|} & \text{otherwise} \end{cases}$$

where $|S|$ is the number of nodes in $S$

# P(S|F, n): Generating structures from forms

- Simpler forms are preferred

# A hierarchical Bayesian model

Meta-constraints | M

F: form

Tree

S: structure

D: data

whiskers | hands | tail

mouse  ● ○ ●
squirrel  ● ○ ●
chimp  ○ ● ○
gorilla  ○ ● ○

$$P(S, F | D, n) \propto P(D|S)P(S|F, n)P(F)$$

# p(D|S): Generating feature data

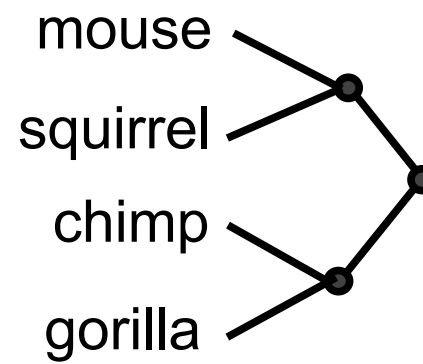- Intuition: features should be smooth over graph S

Relatively smooth          Not smooth

# p(D|S): Generating feature data



Let $f_i$ be the feature value at node $i$

$$p(f) \propto \exp\left(-\frac{1}{4}\sum_{i,j}\frac{(f_i - f_j)^2}{d_{ij}}\right)$$

(Zhu, Lafferty & Ghahramani)

# A hierarchical Bayesian model

Meta-constraints          M

F: form          Tree

S: structure

mouse
squirrel
chimp
gorilla

D: data

|  | whiskers | hands | tail |
|---|---|---|---|
| mouse | ● | ○ | ● |
| squirrel | ● | ○ | ● |
| chimp | ○ | ● | ○ |
| gorilla | ○ | ● | ○ |

$$P(S, F | D, n) \propto P(D|S)P(S|F,n)P(F)$$

# Feature data: results

# Developmental shifts



5 features

20 features

110 features

# Similarity data: results

# Relational data

Meta-constraints | M
Form | Cliques
Structure
Data

# Relational data: results

## Primates

"x dominates y"

## Bush cabinet

"x tells y"

## Prisoners

"x is friends with y"

# Why structural form matters

- Structural forms support predictions about new or sparsely-observed entities.

# Experiment: Form discovery

Cliques   (n = 8/12)

Chain   (n = 7/12)

# Universal Structure grammar     U

## Form



## Structure



mouse
squirrel
chimp
gorilla

## Data

|  | whiskers | hands | tail |
|---|---|---|---|
| mouse | ● | ○ | ● |
| squirrel | ● | ○ | ● |
| chimp | ○ | ● | ○ |
| gorilla | ○ | ● | ○ |

# A hypothesis space of forms

# Conclusions: Part 2

- Hierarchical Bayesian models provide a unified framework which helps to explain:

  – How abstract knowledge is acquired

  – How abstract knowledge is used for induction

# Outline

- Learning structured representations
  - grammars
  - logical theories

- Learning at multiple levels of abstraction

# Handbook of Mathematical Psychology, 1963