

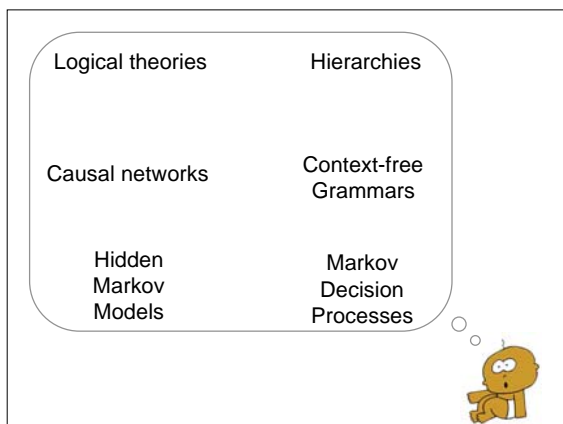
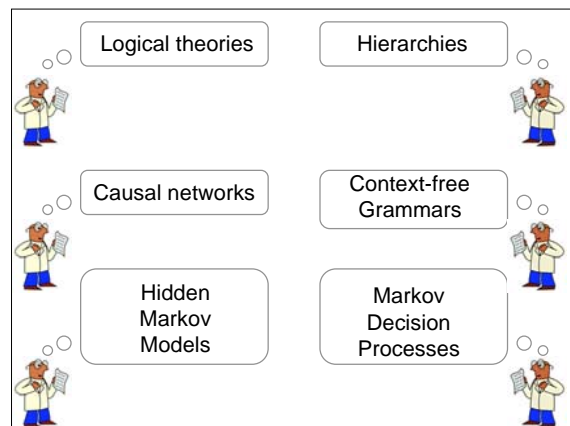
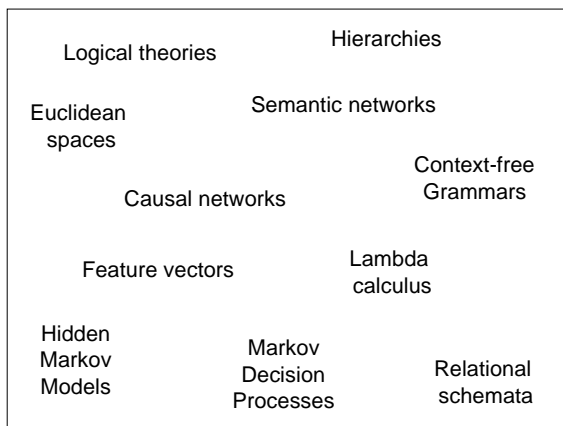
The development of structured representations

Charles Kemp

July 25, IPAM Summer School

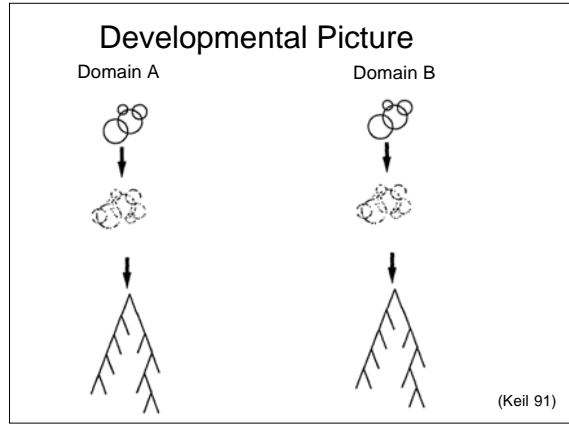
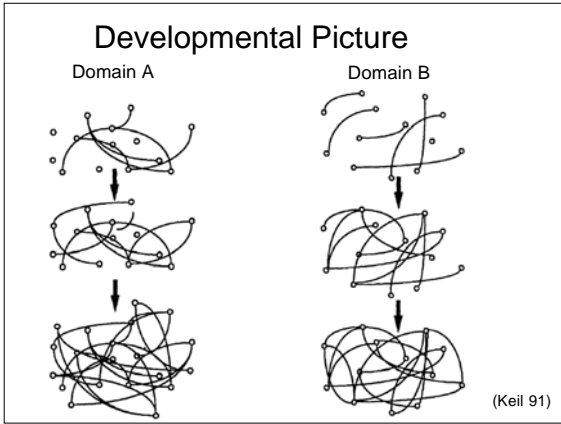
Acknowledgements

Josh Tenenbaum
Pooja Jotwani

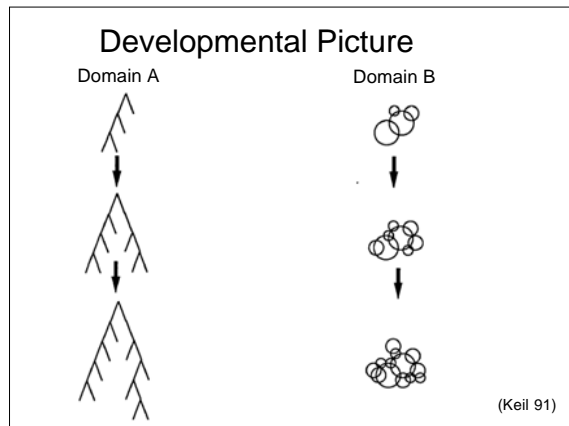


Hypotheses about representations

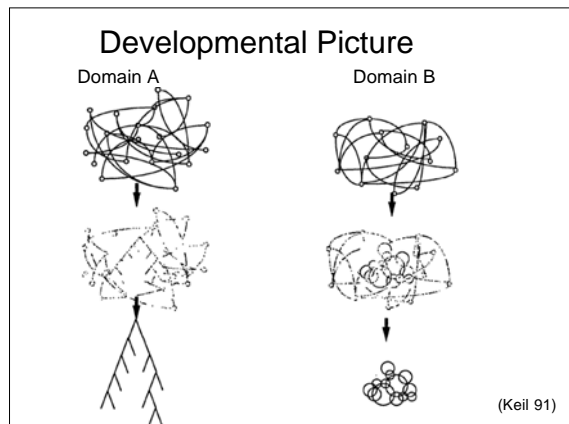
1. There is one kind of representation that will handle every domain.



- ### Hypotheses about representations
1. There is one kind of representation that will handle every domain.
 2. Children begin with innate, domain-specific representational constraints.



- ### Hypotheses about representations
1. There is one kind of representation that will handle every domain.
 2. Children begin with innate, domain-specific representational constraints.
 3. Children discover which kind of representation is best for each domain.



Scientists discover structural form

Mendeleev's Periodic Table of 1869¹

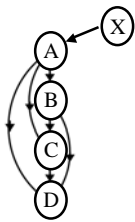
			Li	Be	B	C	N	O	F	Ne										
			Na	Mg	Al	Si	P	S	Cl	Ar										
			K	Ca	Sc	Ti	V	Cr	Mn	Fe	Cobalt	Nickel	Cu	Zn	Ga	Ge	As	Se	Br	Krypton
			Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Au	Hg	Indium	Sn	Antimony	Tellurium	Iodine	Xenon
			Cs	Ba	La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Hf	Ta	Tl	Pb	Bismuth	Po	Ast	Radon
			Fr	Ra	Ac	Th	Pa	Uranium					Ra	Actinium						

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

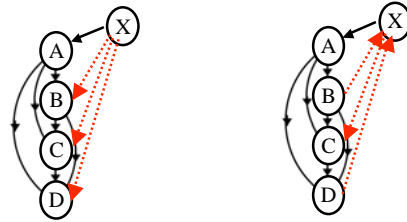
Why form discovery matters

- Structural forms provide inductive constraints

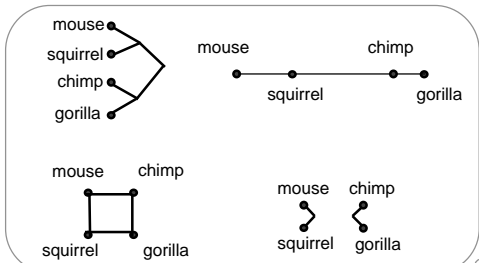


Why form discovery matters

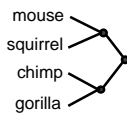
- Structural forms provide inductive constraints



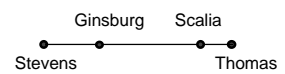
This talk: graph structures



BIOLOGY



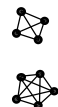
POLITICS



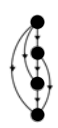
COLOR



FRIENDSHIP



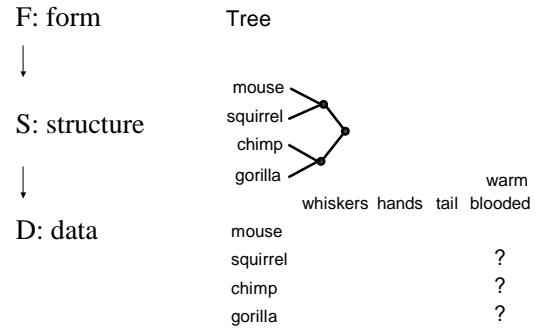
DOMINANCE



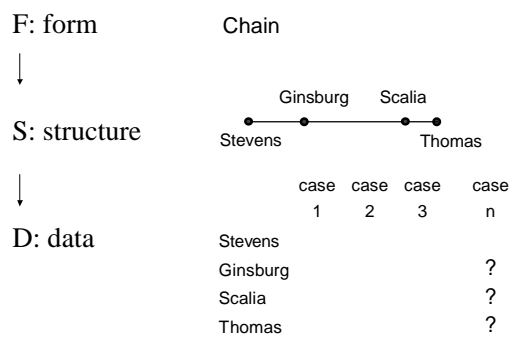
Outline

- Discovery of structural form
 - Feature data
 - Similarity
 - Relational data
- Form discovery in the lab

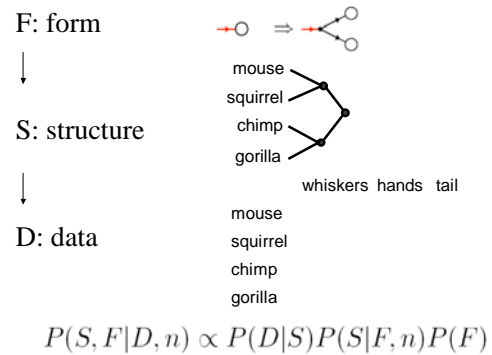
A hierarchical Bayesian framework



A hierarchical Bayesian framework



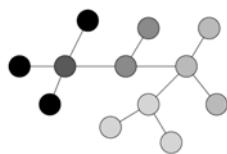
A hierarchical Bayesian framework



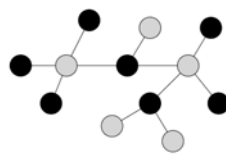
$p(D|S)$: Generating feature data

- Intuition: features should be smooth over graph S

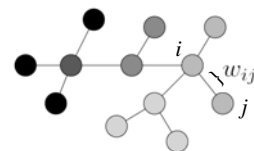
Relatively smooth



Not smooth



$p(f|W, \sigma)$: Generating a single feature



Let f_i be the feature value at node i

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

$$= \exp\left(-\frac{1}{2} f^T \Sigma^{-1} f\right)$$

(Zhu, Lafferty, Ghahramani 03)

$p(D|S,W,\sigma)$: Generating feature data

- The log likelihood for the feature model is

$$\log(p(D|W,\sigma)) = -\frac{mn}{2} \log(2\pi) - \frac{m}{2} \log |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1}DD^T)$$

where

- D is a matrix of objects (n) by features (m)
- W is a weighted graph
- σ specifies a prior on the variance of each feature value

$p(D|S)$: Generating feature data

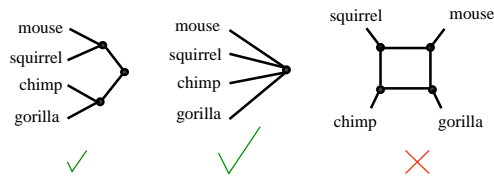
- Generating a weighted graph

$$w_{ij}|S, \beta \sim \text{Exponential}(\beta) \text{ if } s_{ij} = 1$$

- We integrate out W and σ using the Laplace approximation

$$p(D|S, \beta) = \int p(D|W, \sigma) p(W|S, \beta) p(\sigma) dW d\sigma$$

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



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

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

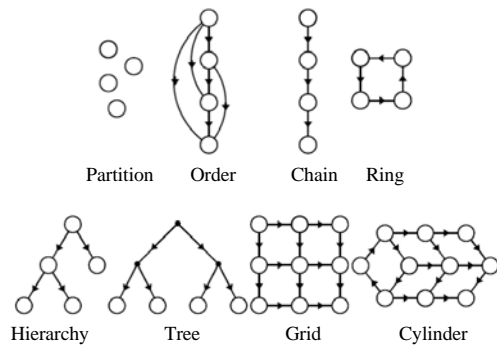
where $|S|$ is the number of nodes in S

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

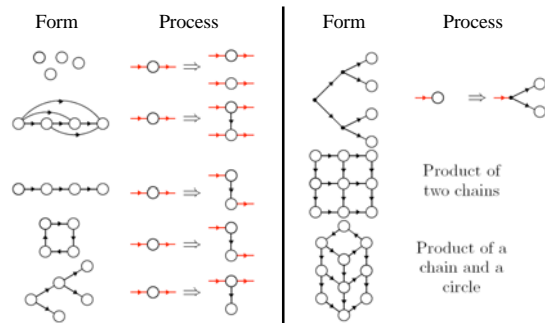
- Simpler forms generate fewer structures, and are therefore preferred by the prior

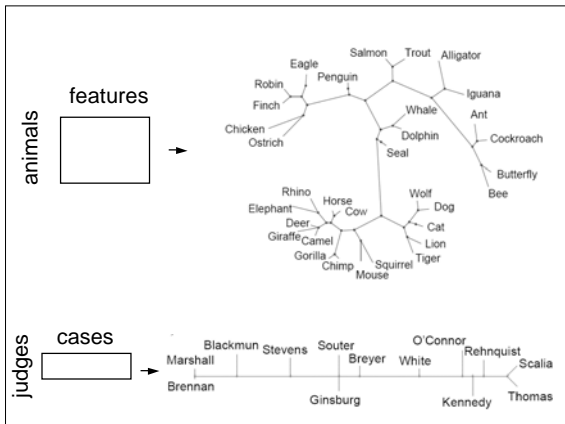
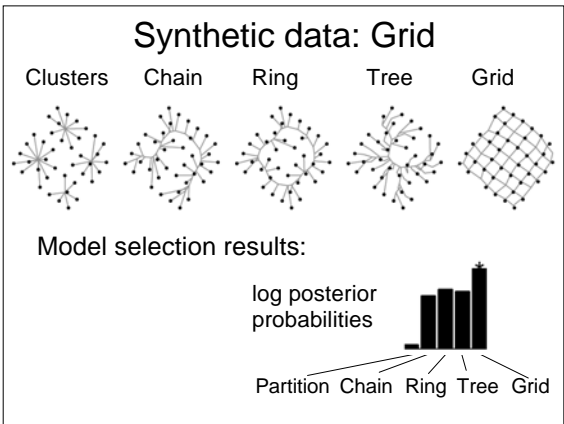
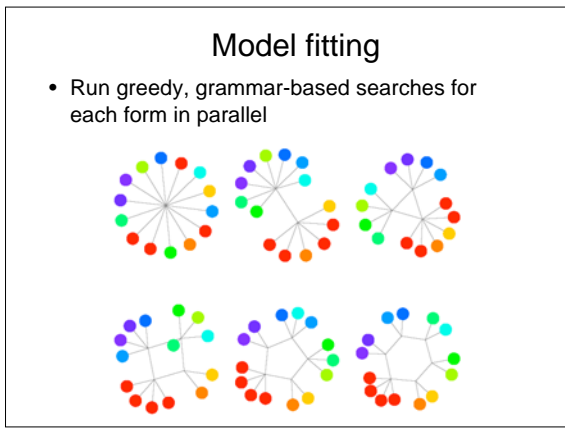
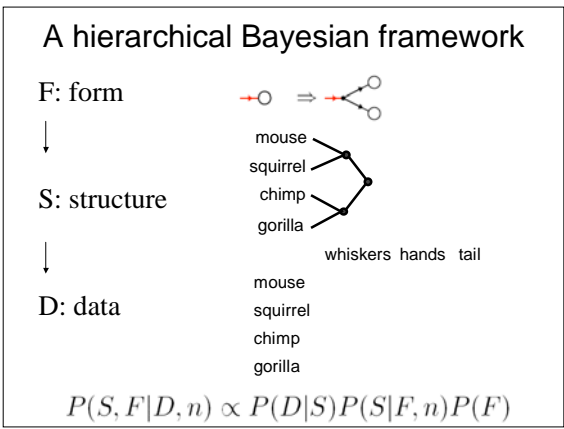
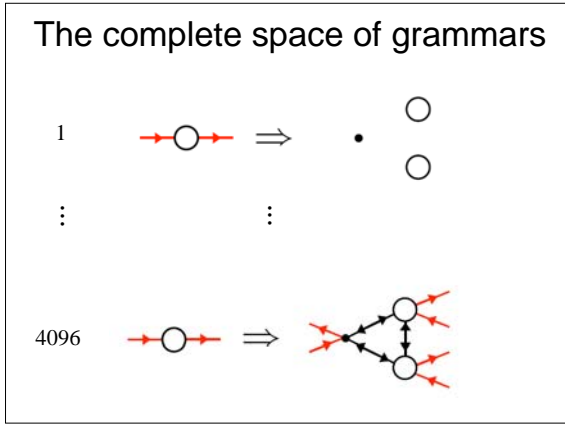
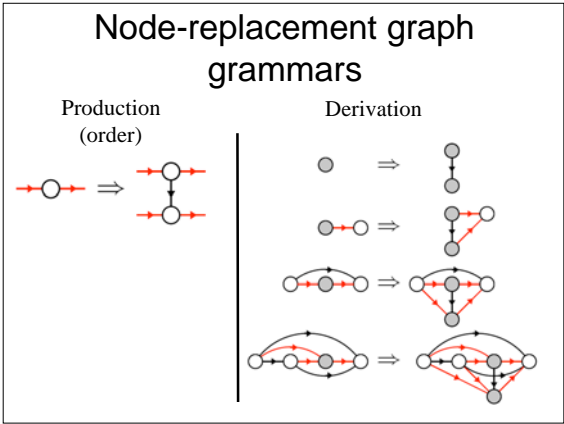


$P(F)$: Structural forms



Generating Graphs





Similarity data

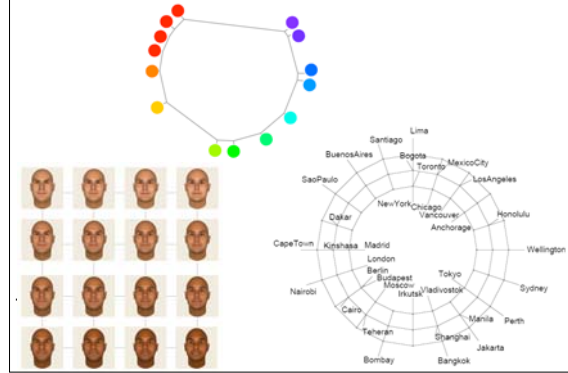
- The log likelihood for the feature model is

$$\log(p(D|W, \sigma)) = -\frac{mn}{2} \log(2\pi) - \frac{m}{2} \log |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} DD^T)$$

where D is a matrix of objects (n) by features (m)

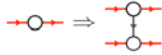
- The kernel trick: replacing DD^T with a similarity matrix lets us learn structural forms from similarity data

Similarity data: results

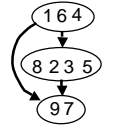


Relational Data

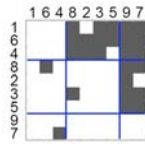
F: form



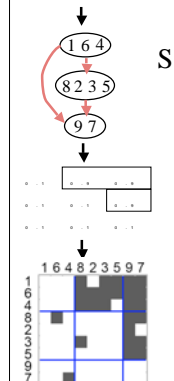
S: structure



D: data



Order F



S

$$\theta_{ab} | S \sim \begin{cases} \text{Beta}(\alpha_0, \beta_0), & \text{if } S_{ab} = 0 \\ \text{Beta}(\alpha_1, \beta_1), & \text{if } S_{ab} = 1 \end{cases}$$

$$D_{ij} | S, \theta \sim \text{Bernoulli}(\theta_{s_i s_j})$$

$P(D|S)$: Relational data

$$P(D|S) = \sum_{(\alpha_0, \beta_0, \alpha_1, \beta_1)} P(D|S, \alpha_0, \beta_0, \alpha_1, \beta_1) P(\alpha_0, \beta_0, \alpha_1, \beta_1)$$

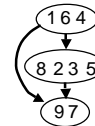
- The hyperparameters are drawn from a 4D grid where
 - $\alpha_0 + \beta_0$ and $\alpha_1 + \beta_1$ belong to $\{\frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2, 4, 8, 16, 32\}$
 - $\frac{\beta_0}{\alpha_0 + \beta_0}$ and $\frac{\beta_1}{\alpha_1 + \beta_1}$ belong to $\{0.05, 0.15, \dots, 0.95\}$
 - $\frac{\beta_0}{\alpha_0 + \beta_0} \leq \frac{\beta_1}{\alpha_1 + \beta_1}$

A hierarchical Bayesian framework

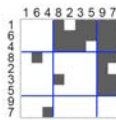
F: form



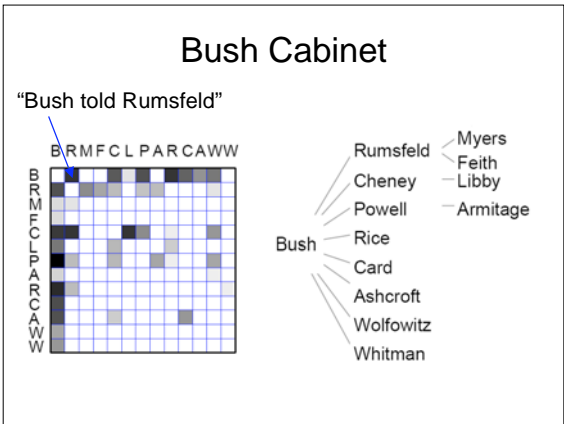
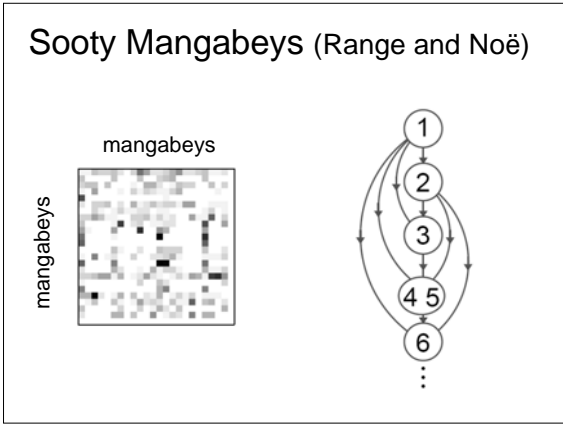
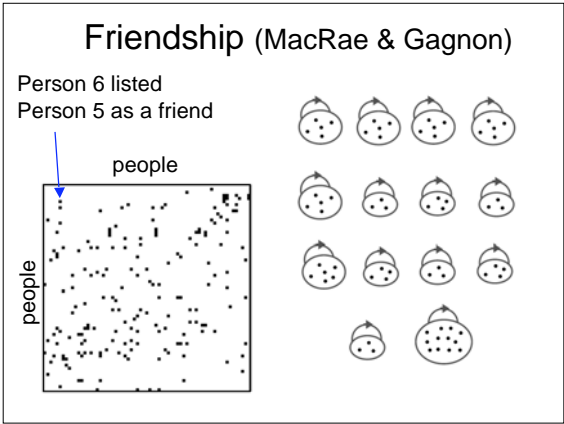
S: structure



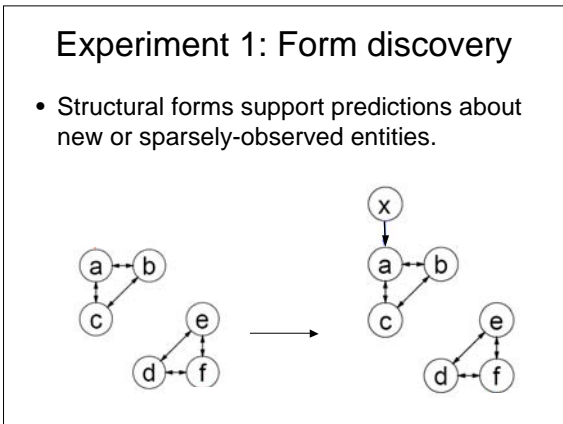
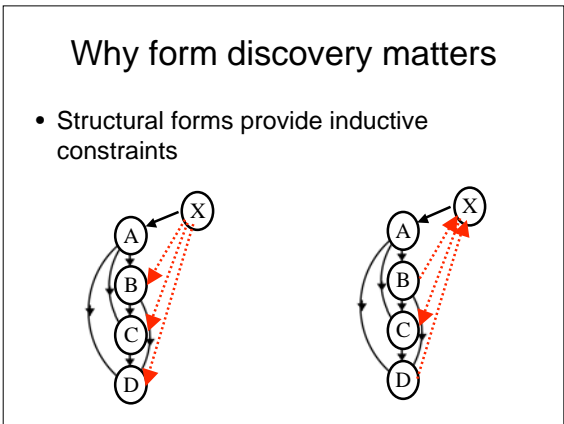
D: data

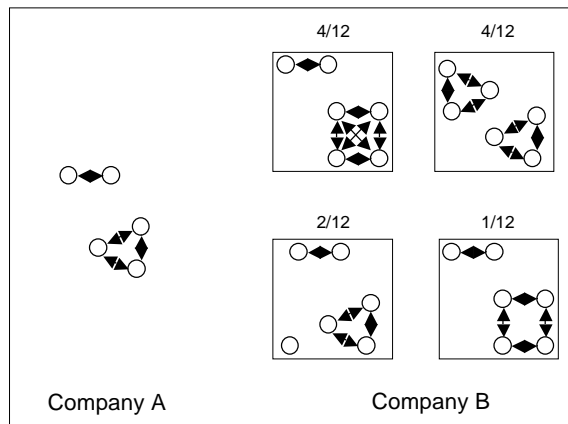
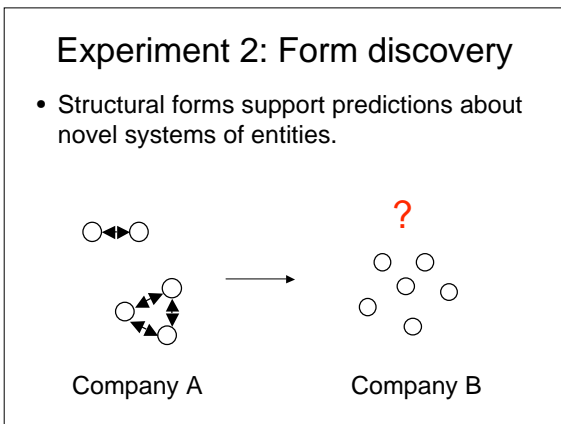
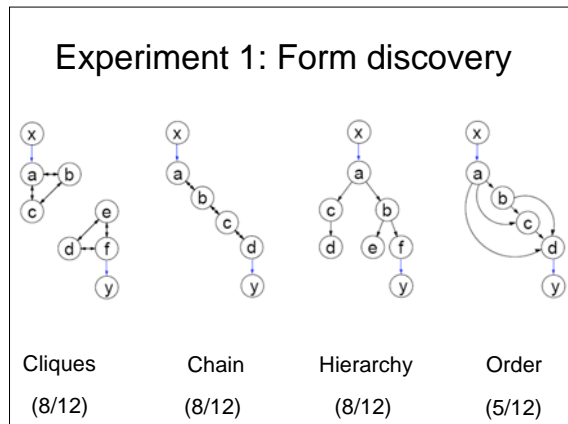
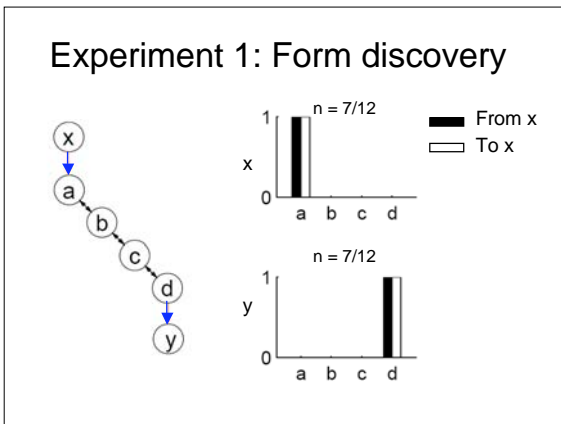
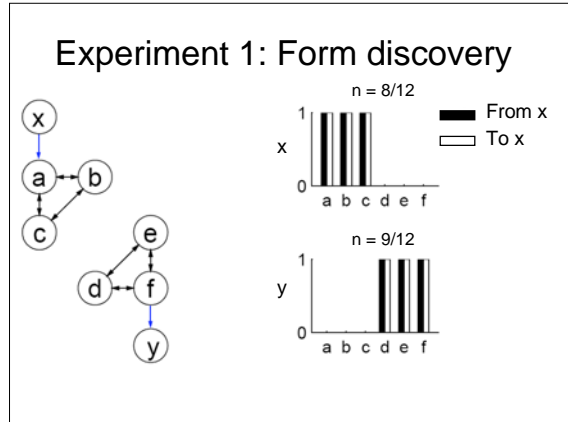
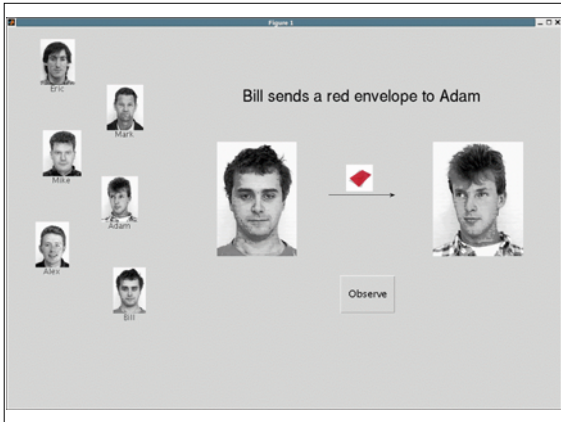


$$P(S, F|D, n) \propto P(D|S) P(S|F, n) P(F)$$

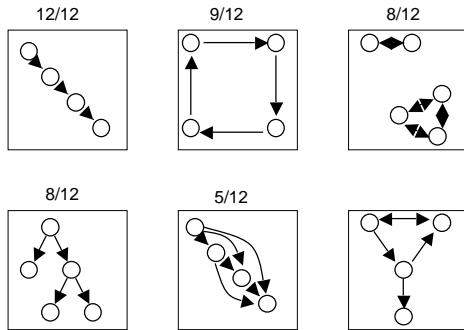


- ### Outline
- Discovery of structural form
 - Feature data
 - Similarity
 - Relational data
 - Form discovery in the lab

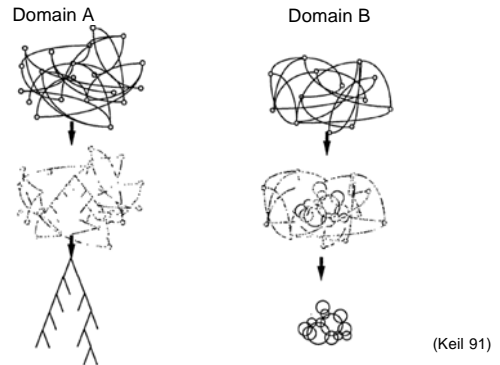




Experiment 2: Form discovery



Developmental Picture



Developmental predictions



Issues and questions

- How can we work with richer collections of structure grammars?
- Where do these structure grammars come from?
- What about representations other than graphs?

Conclusions

- Hierarchical models can help to explain:
 - how people acquire mental representations
 - How people learn what *kind* of representation is best for a domain

