

# TREE-STRUCTURED STICK BREAKING FOR HIERARCHICAL DATA

Ryan Prescott Adams  
University of Toronto

Zoubin Ghahramani  
University of Cambridge

Michael I. Jordan  
UC Berkeley

<http://www.cs.toronto.edu/~rpa>



CIFAR

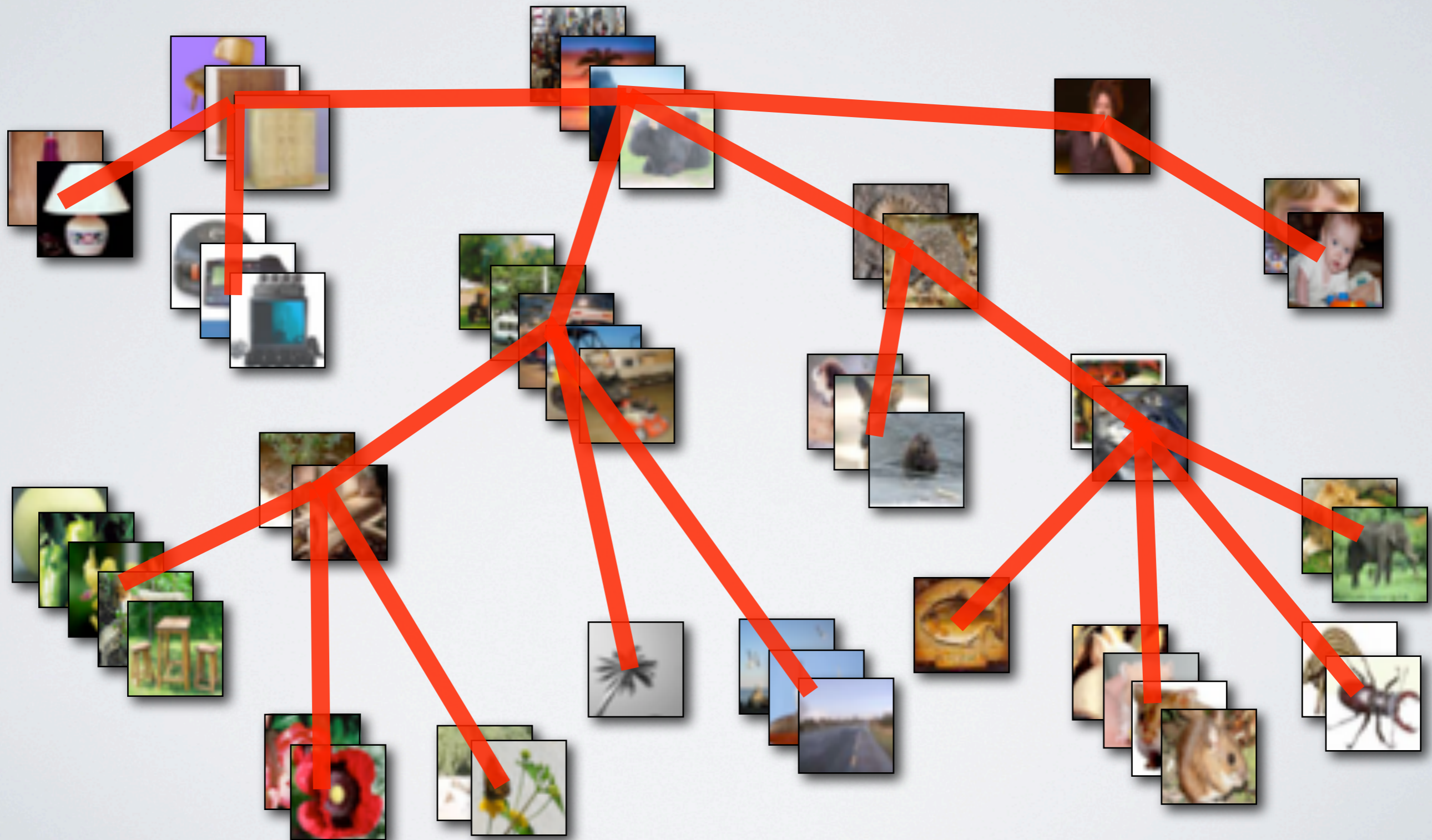
# THE BIG PICTURE

Discover a latent tree structure over data.



# THE BIG PICTURE

Discover a latent tree structure over data.



# HIERARCHICAL CLUSTERING

The main idea of our approach:

**Construct a mixture model in which the components have a tree-structured topology.**

- Bayesian nonparametric approach
- Unbounded width and depth.
- Data live at internal nodes of the tree
- Infinitely exchangeable urn scheme
- MCMC inference
- Applied to image and text

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$

0

1

Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$

0

1

Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set



# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set

# DIRICHLET PROCESS MIXTURES

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Standard stick breaking: natural numbers as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$

0

1

Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set



# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

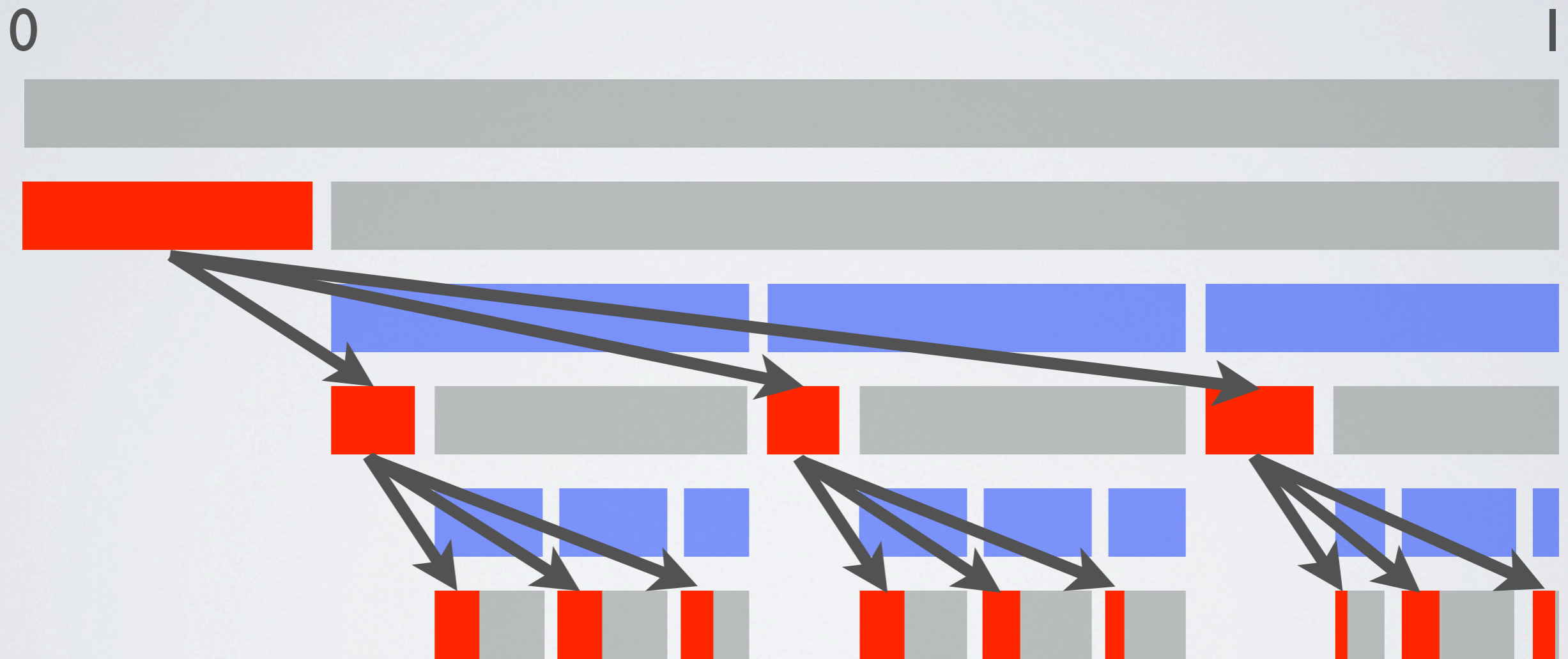
$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



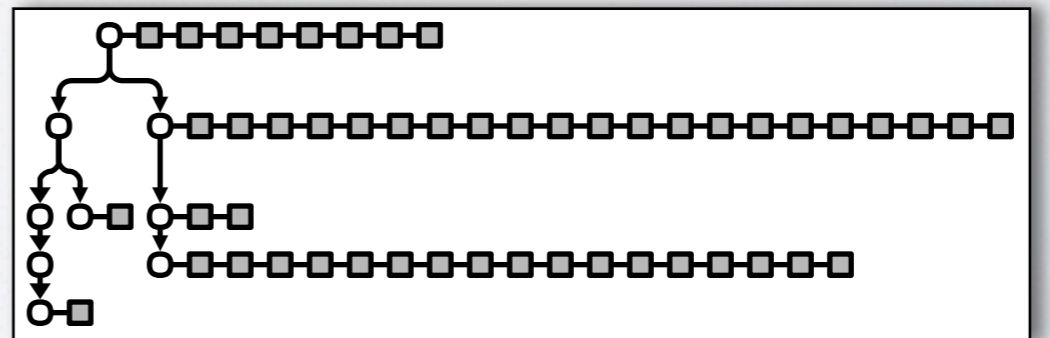
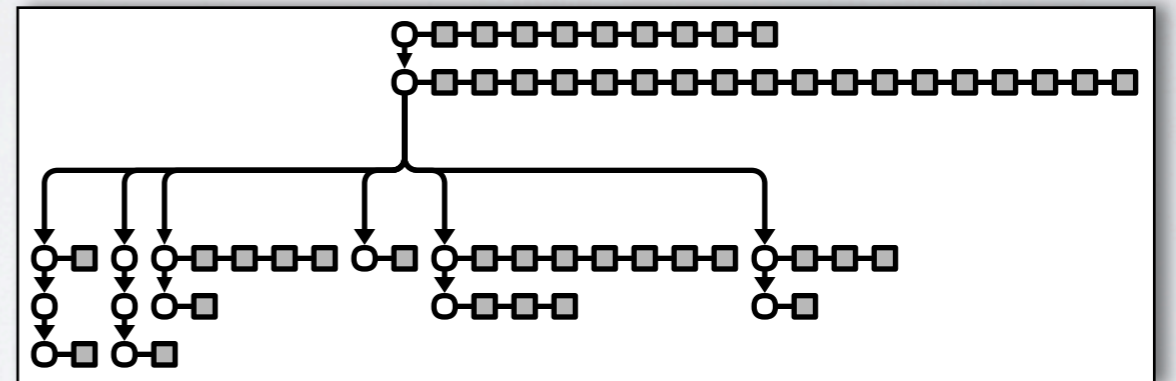
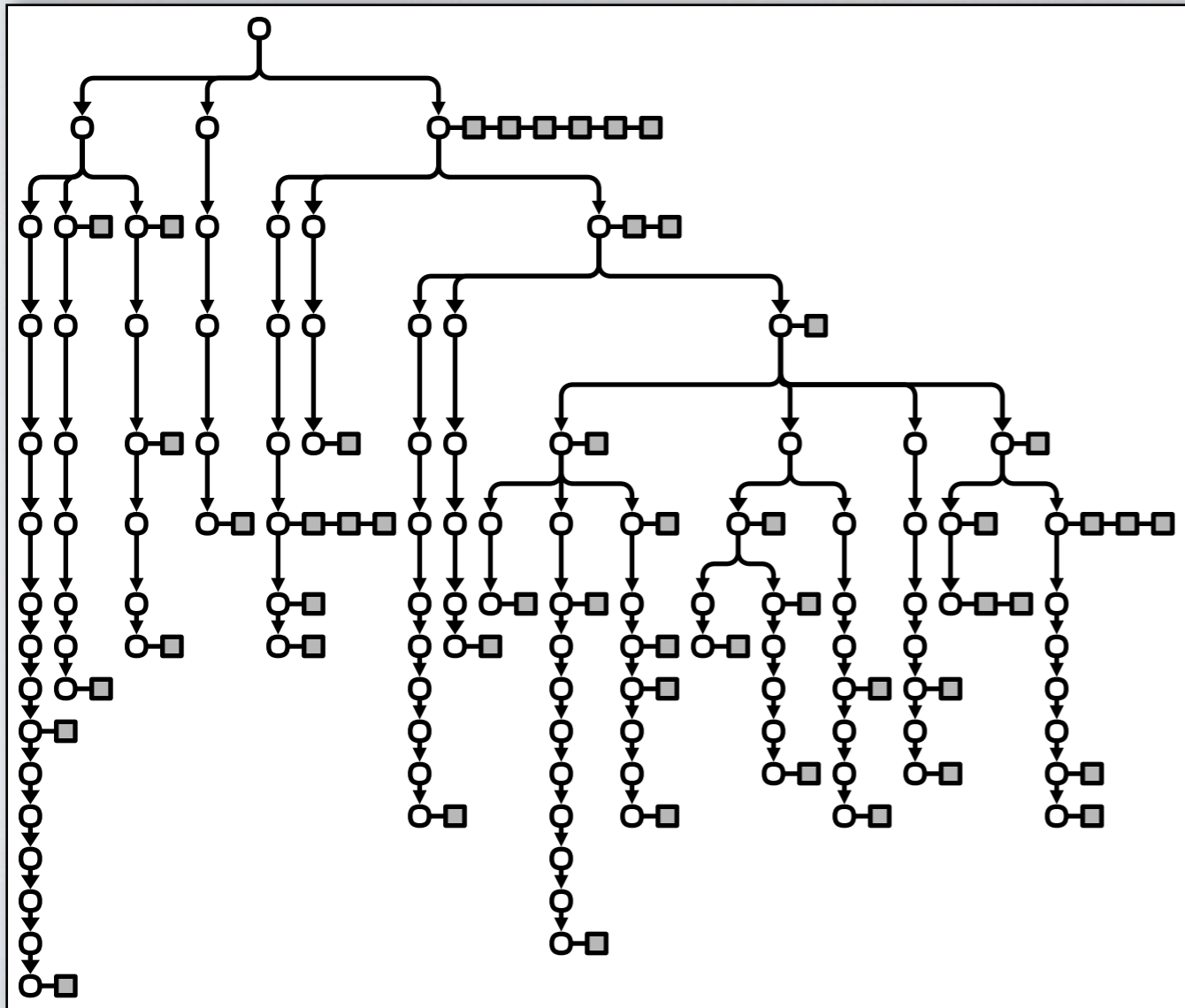
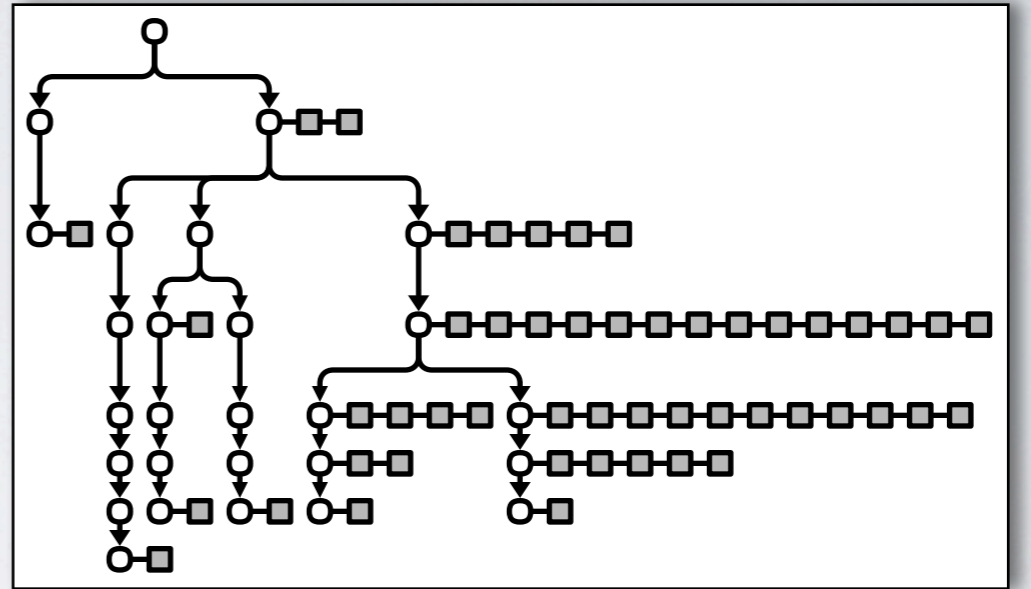
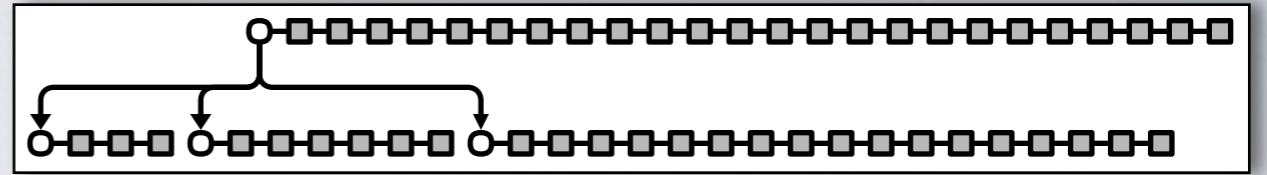
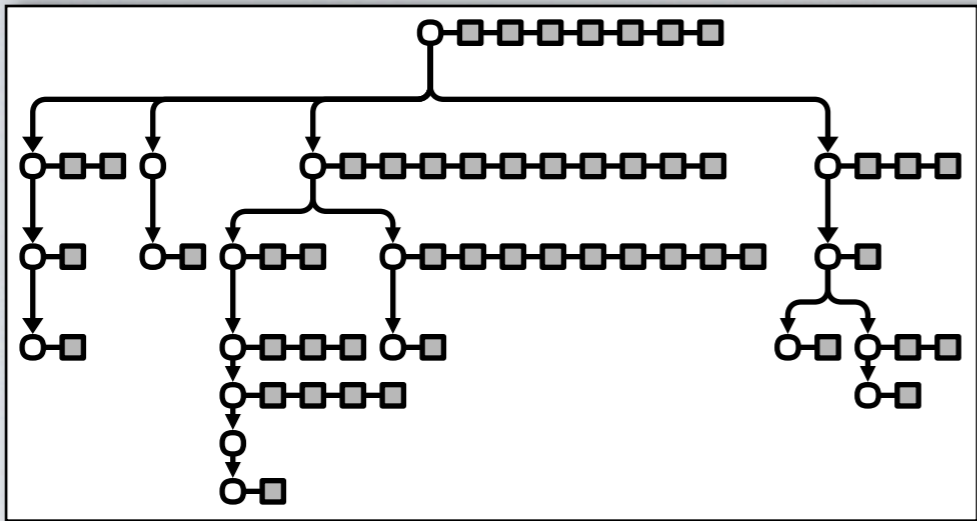
Tree-structured stick breaking: finite-length strings as index set

# TREE STRUCTURED PARTITIONS

$$p(x | \{\pi_k, \theta_k\}_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k f_X(x | \theta_k)$$



Tree-structured stick breaking: finite-length strings as index set

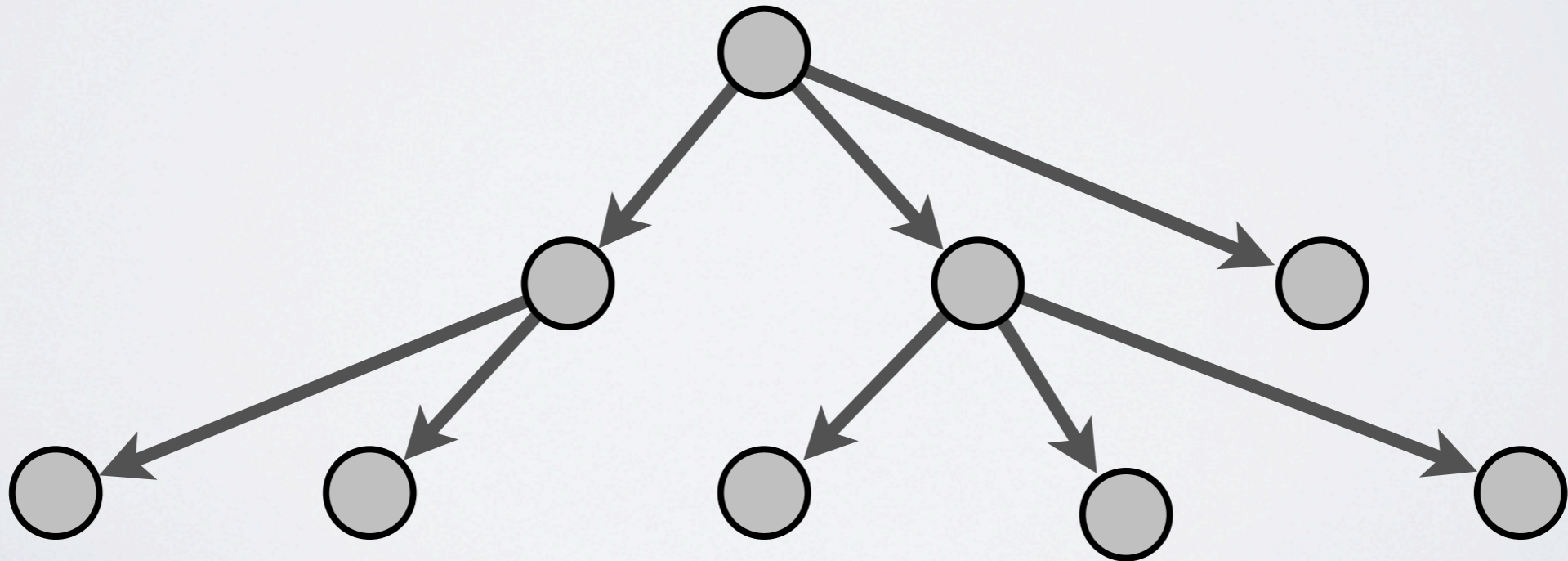


# PARAMETER DEPENDENCE

A fancy topology over the partitions is not useful unless the parameters are linked.

Tree topology = directed graphical model

Parents then provide the prior for the children.



# INFERENCE VIA MCMC

Simulate from the posterior over structure and parameters

Unknown quantities:

1. Assignments of data to nodes
2. Stick lengths
3. Node parameters
4. Stick-breaking hyperparameters

All moves are relatively standard, except #1.

To choose a “best” tree, we keep the assignments which maximized the complete data log likelihood.

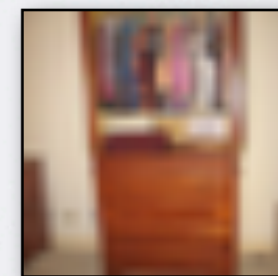
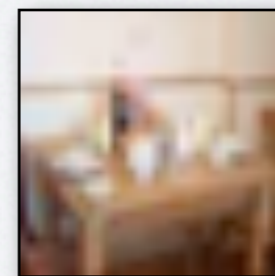
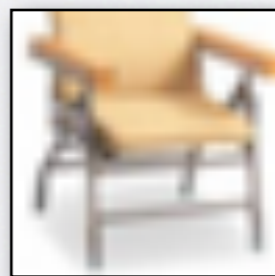
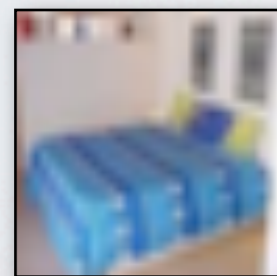
# CIFAR-100 IMAGE DATA

- 50,000 32x32 color images of 100 classes<sup>1</sup>
- Subset of the “80 Million Tiny Images” data set of Torralba, Fergus and Freeman<sup>2,3</sup>
- Sorted and labeled by Krizhevsky, Nair and Hinton
- The 100 classes aggregate into 20 “super classes”

*reptiles: crocodile, dinosaur, lizard, snake, turtle*



*furniture: bed, chair, couch, table, wardrobe*



1. <http://www.cs.toronto.edu/~kriz/cifar.html>

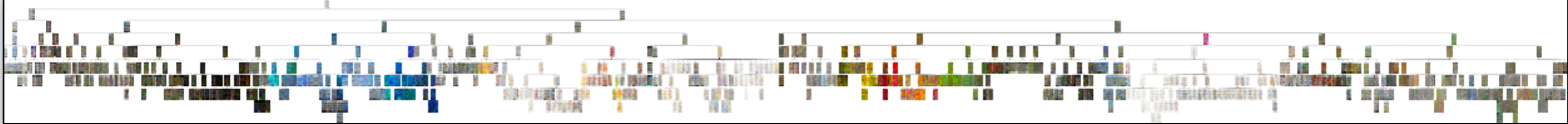
2. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2008

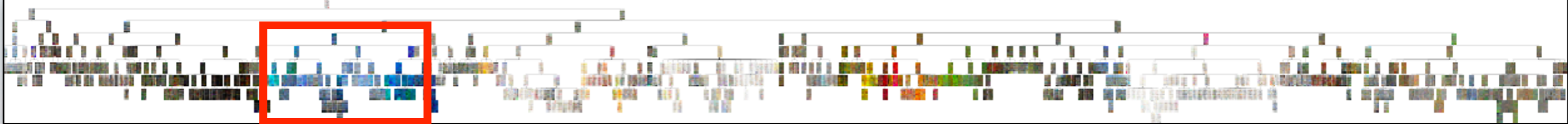
3. <http://groups.csail.mit.edu/vision/TinyImages/>

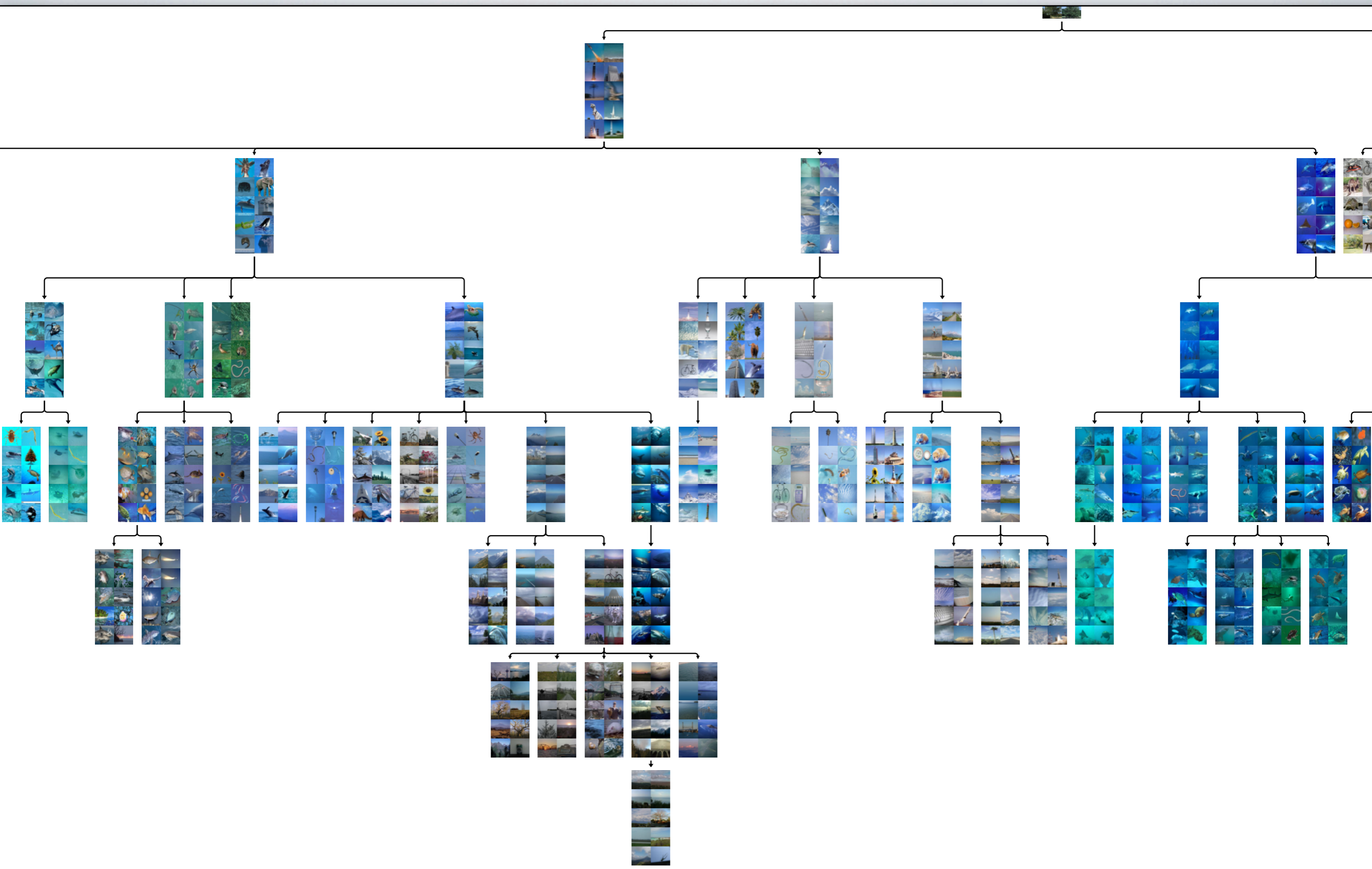


# CLUSTERING IMAGES

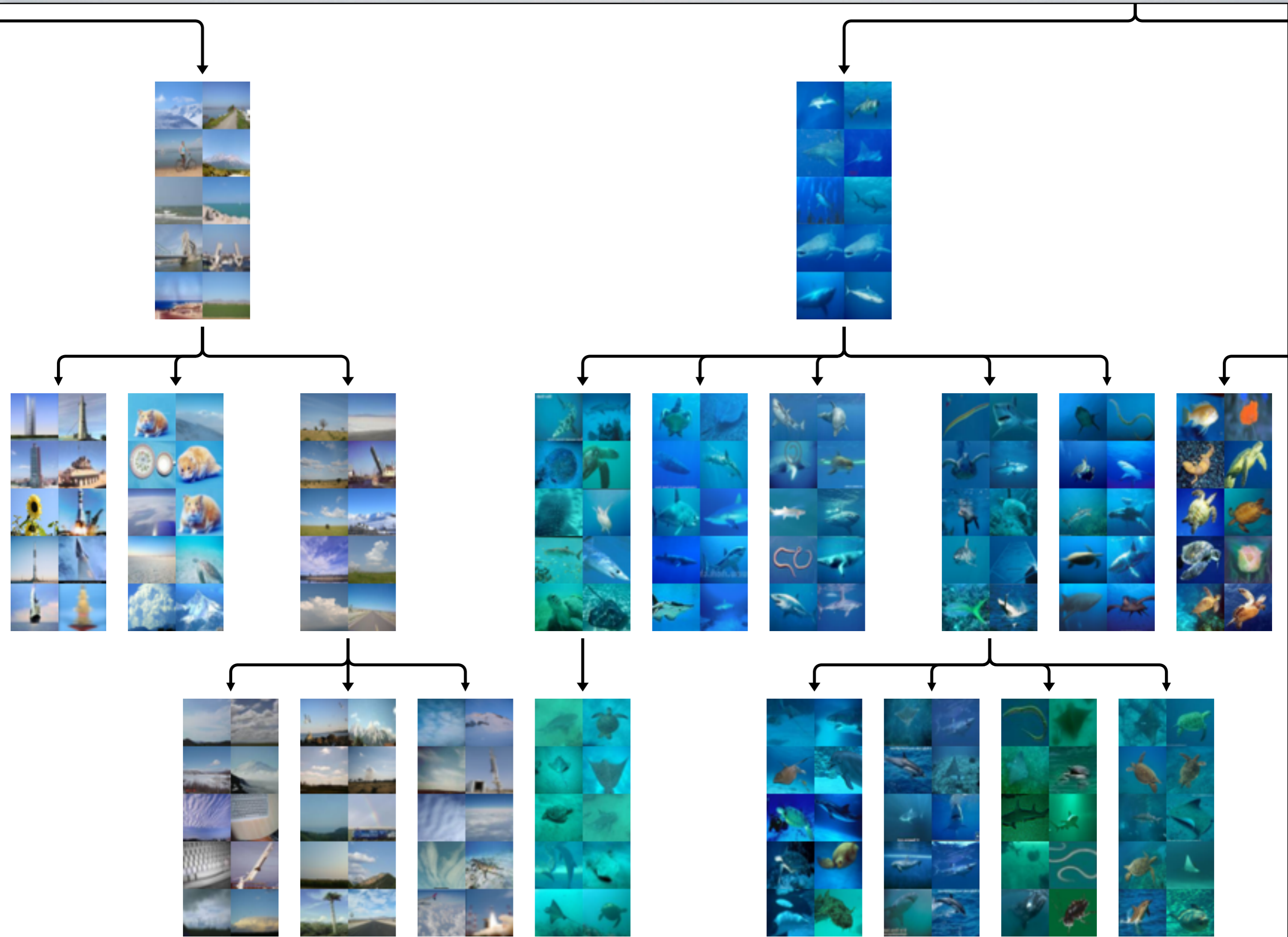
- Letting a mixture model see high-dimensional real-valued data (e.g., 1024 RGB pixels) can be a disaster.
  - Used 256-dimensional binary feature vectors.
  - Features were from a deep autoencoder trained for separate work on image retrieval by Alex Krizhevsky
  - Alex's codes are great!
  - Check out his tech report!
- 
- Each node owns a product of Bernoulli distributions, parameterized by logistic-transformed real values.
  - A child's parameters have a Gaussian prior with a mean determined by the parent's parameters.
  - The prior variances are shared and also inferred.

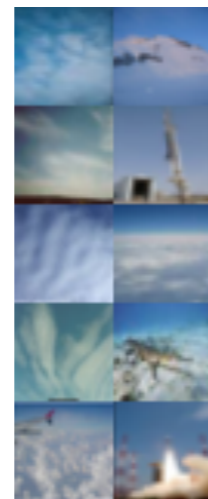
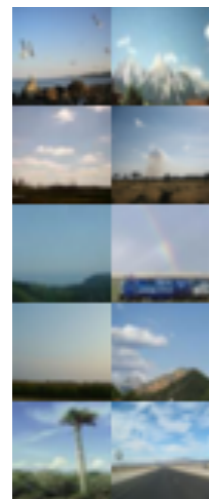
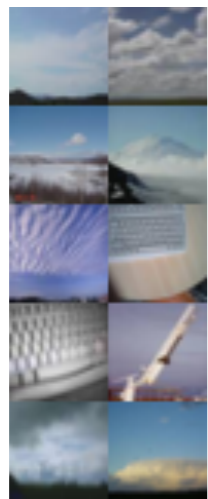
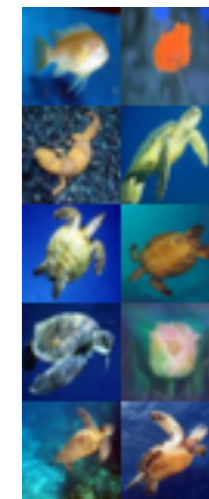
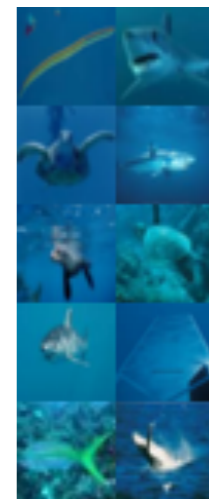
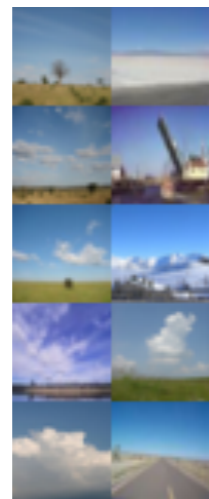
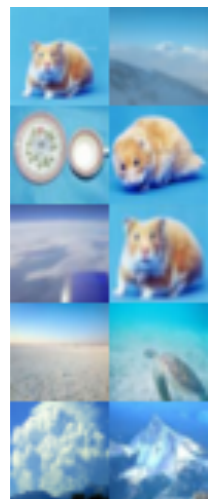
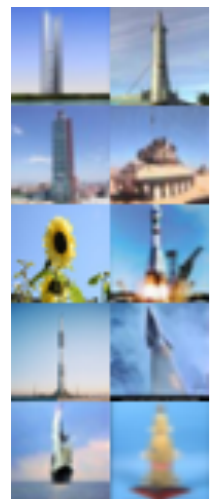
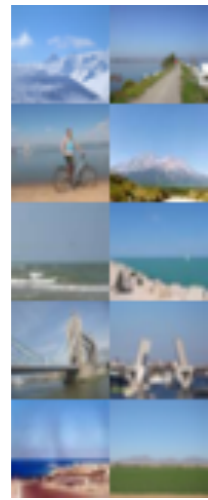


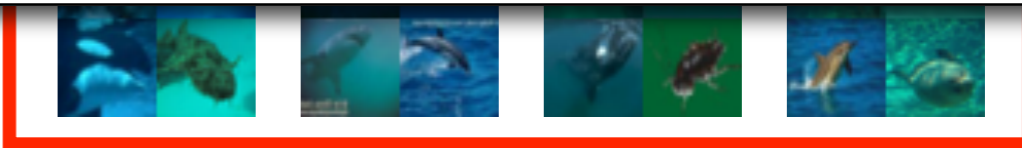
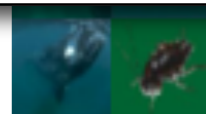
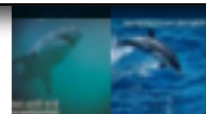
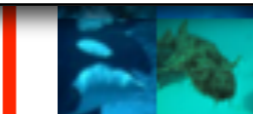
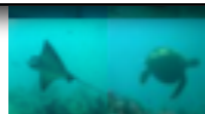
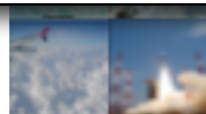
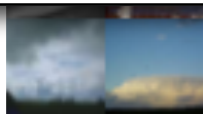
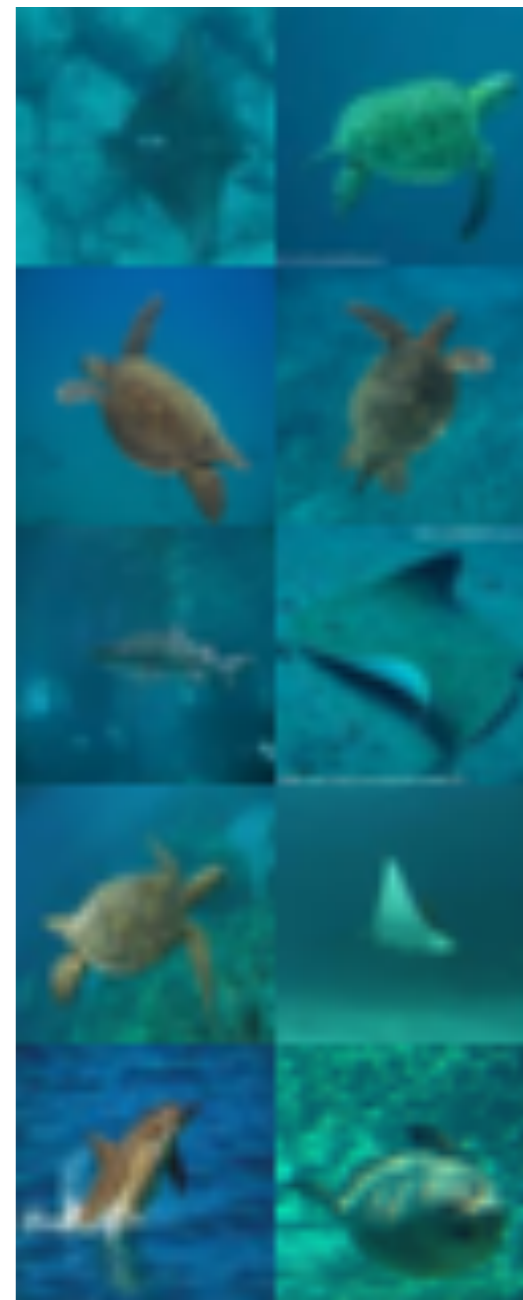




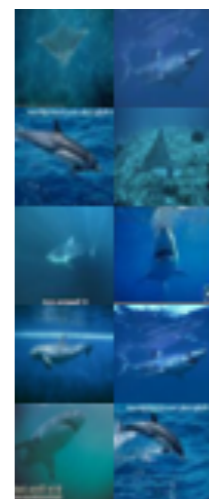
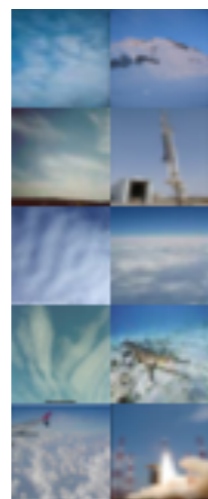
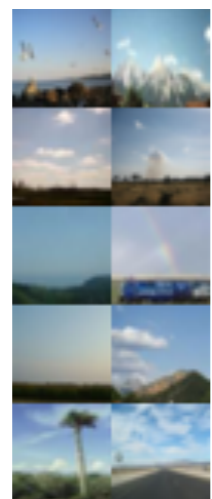
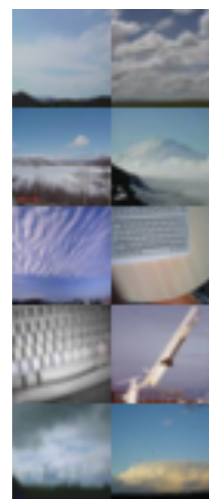
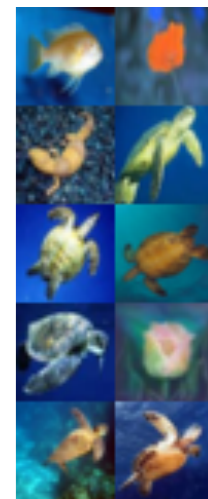
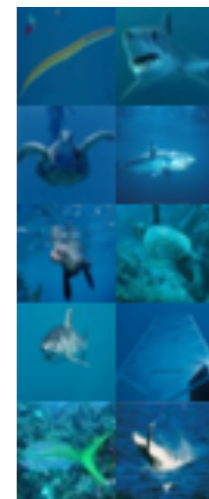
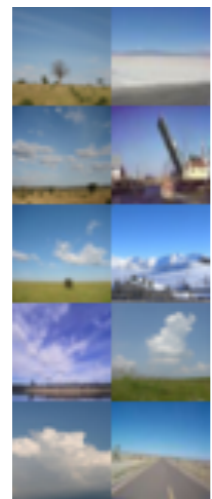
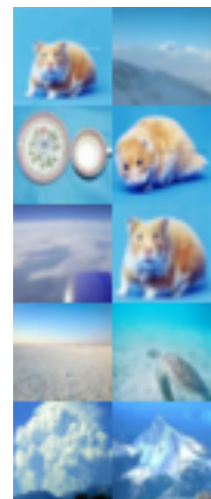
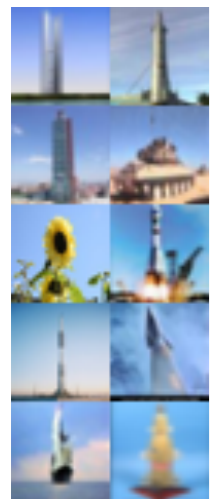
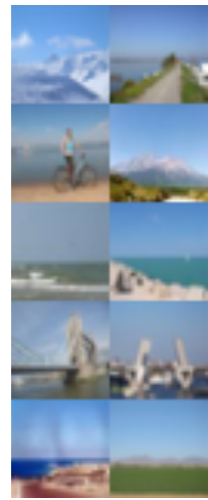




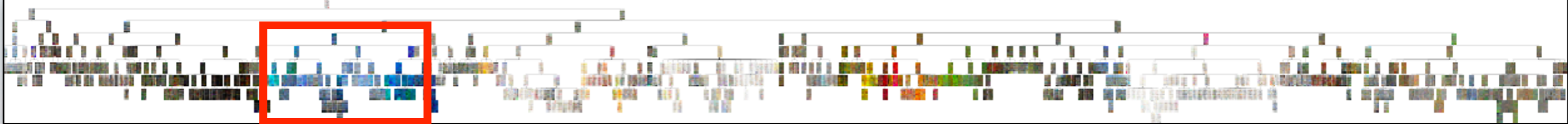


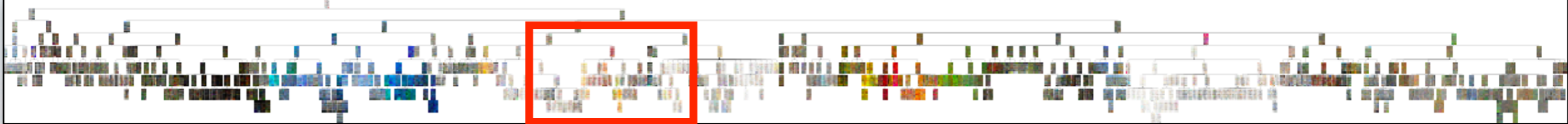


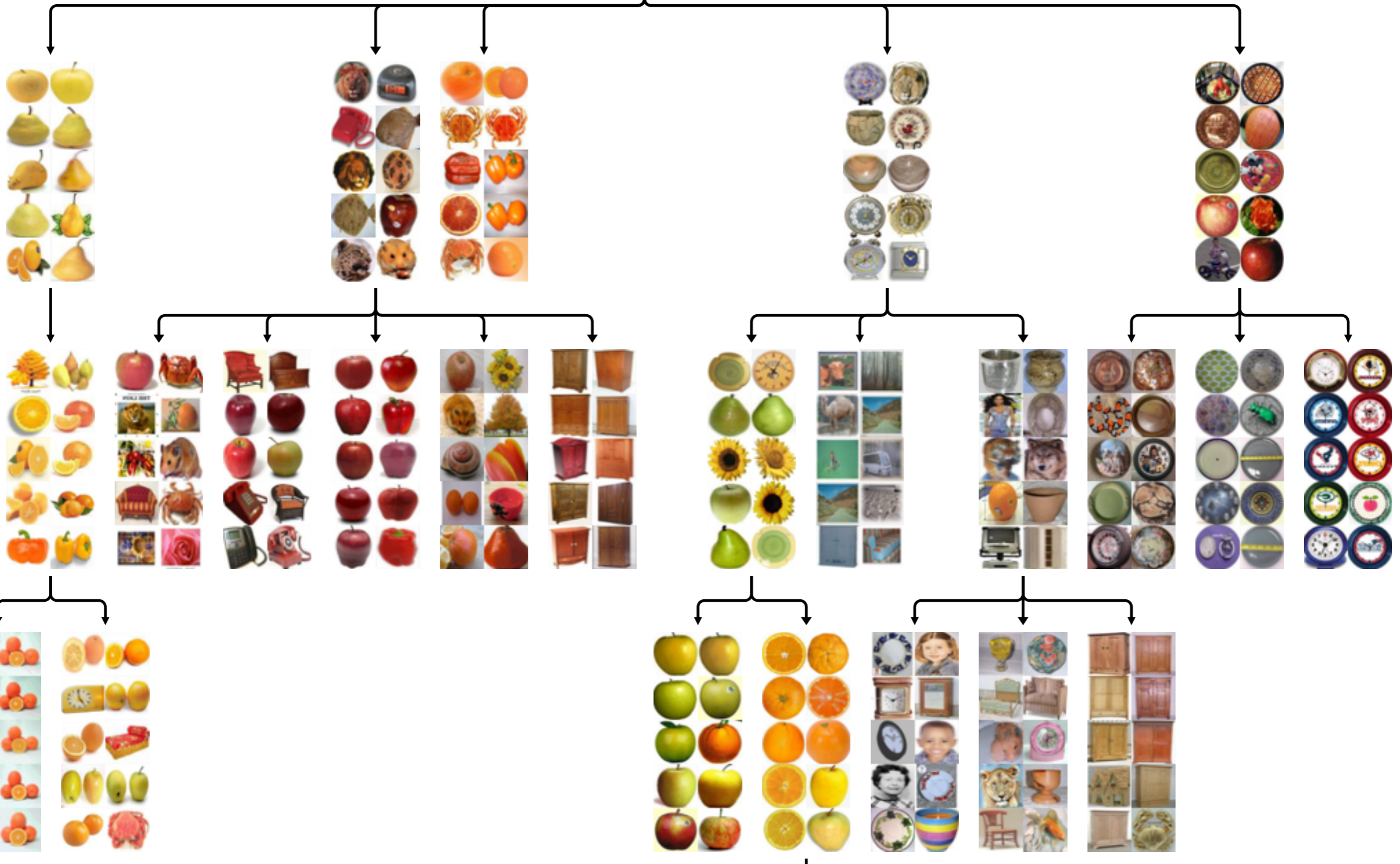


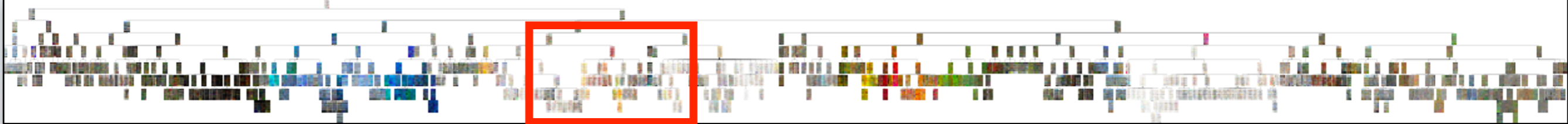




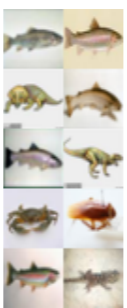
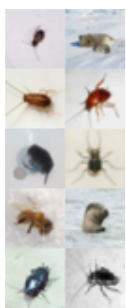
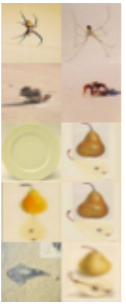
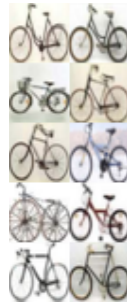
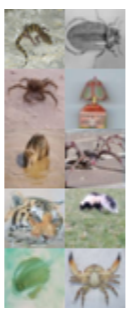
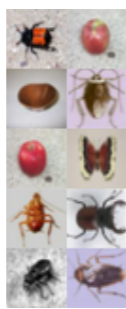
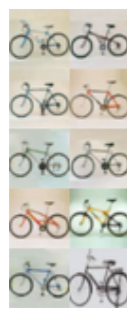
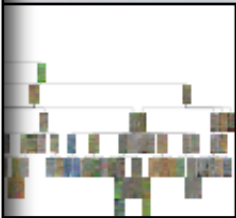
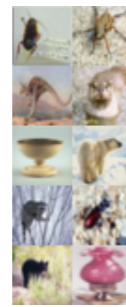
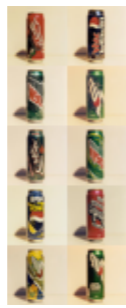




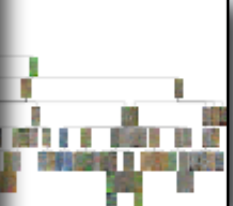
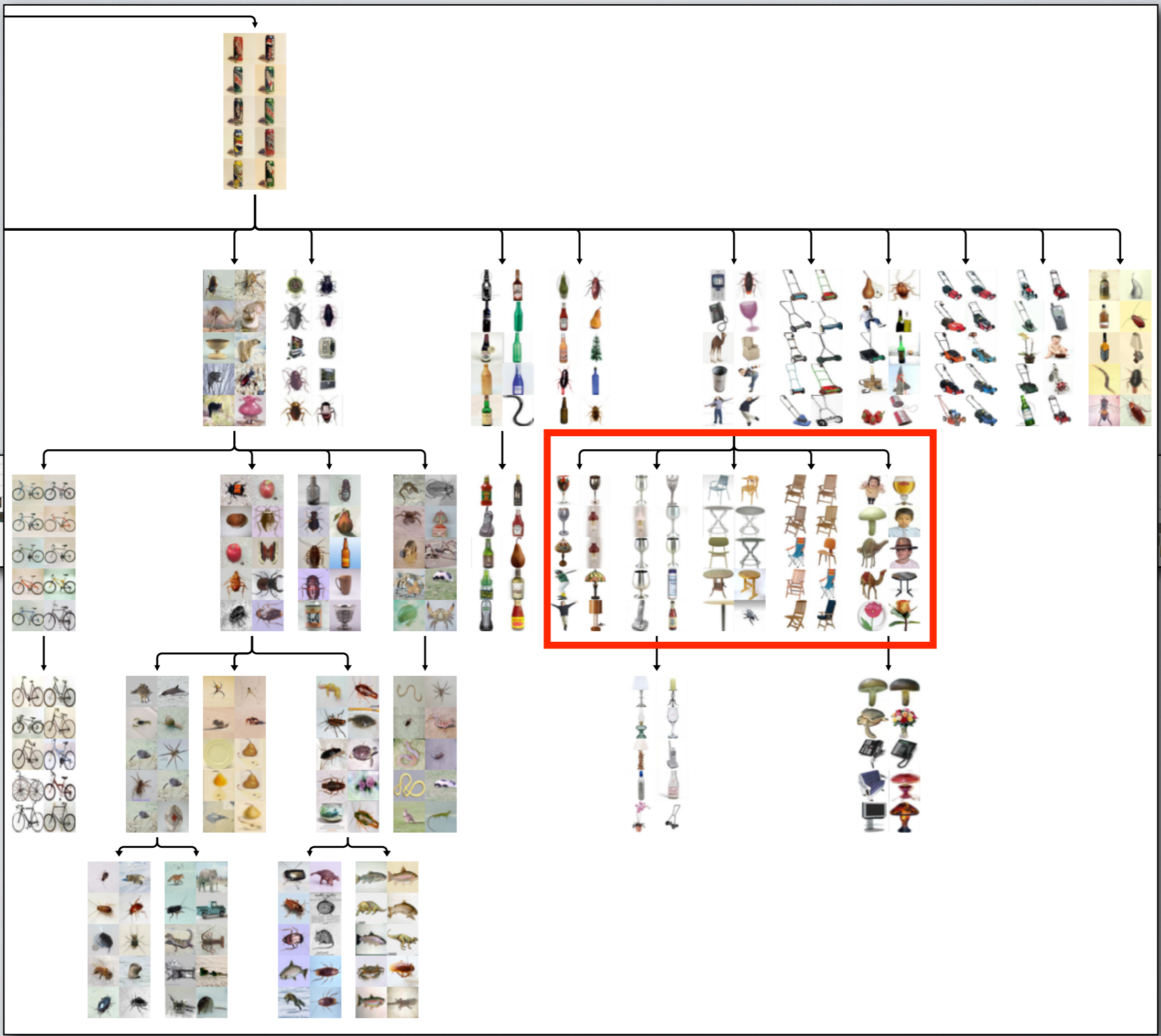


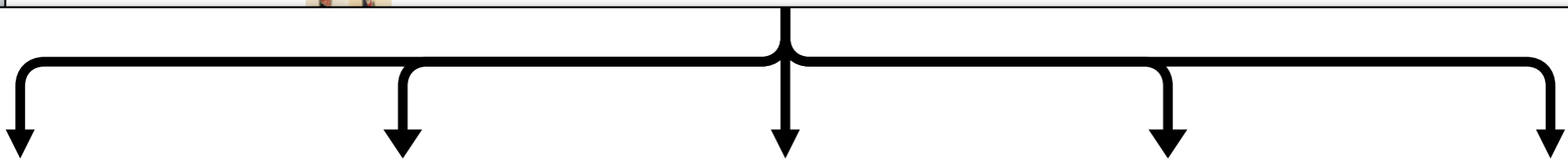


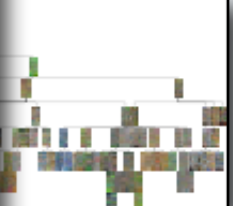
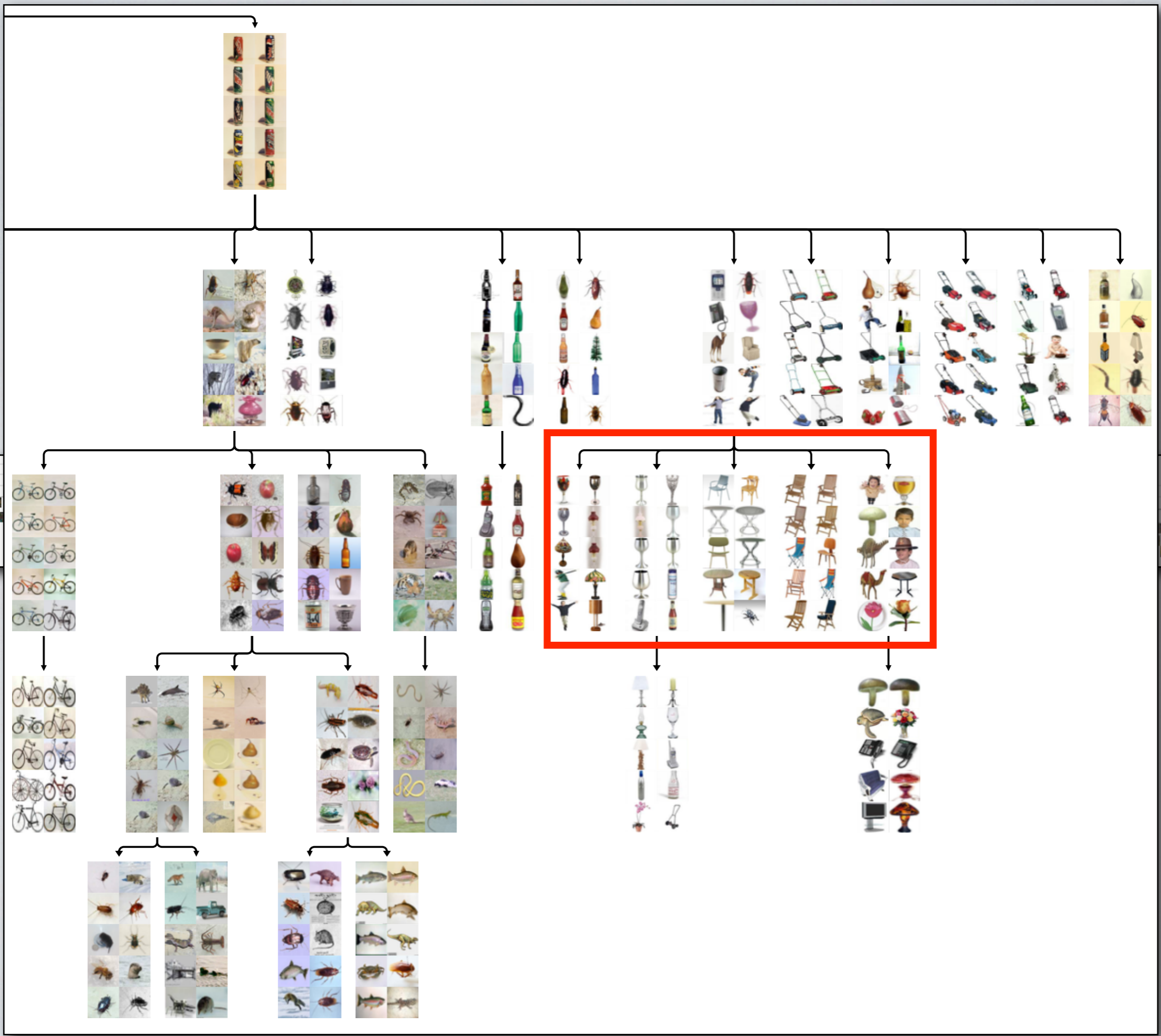






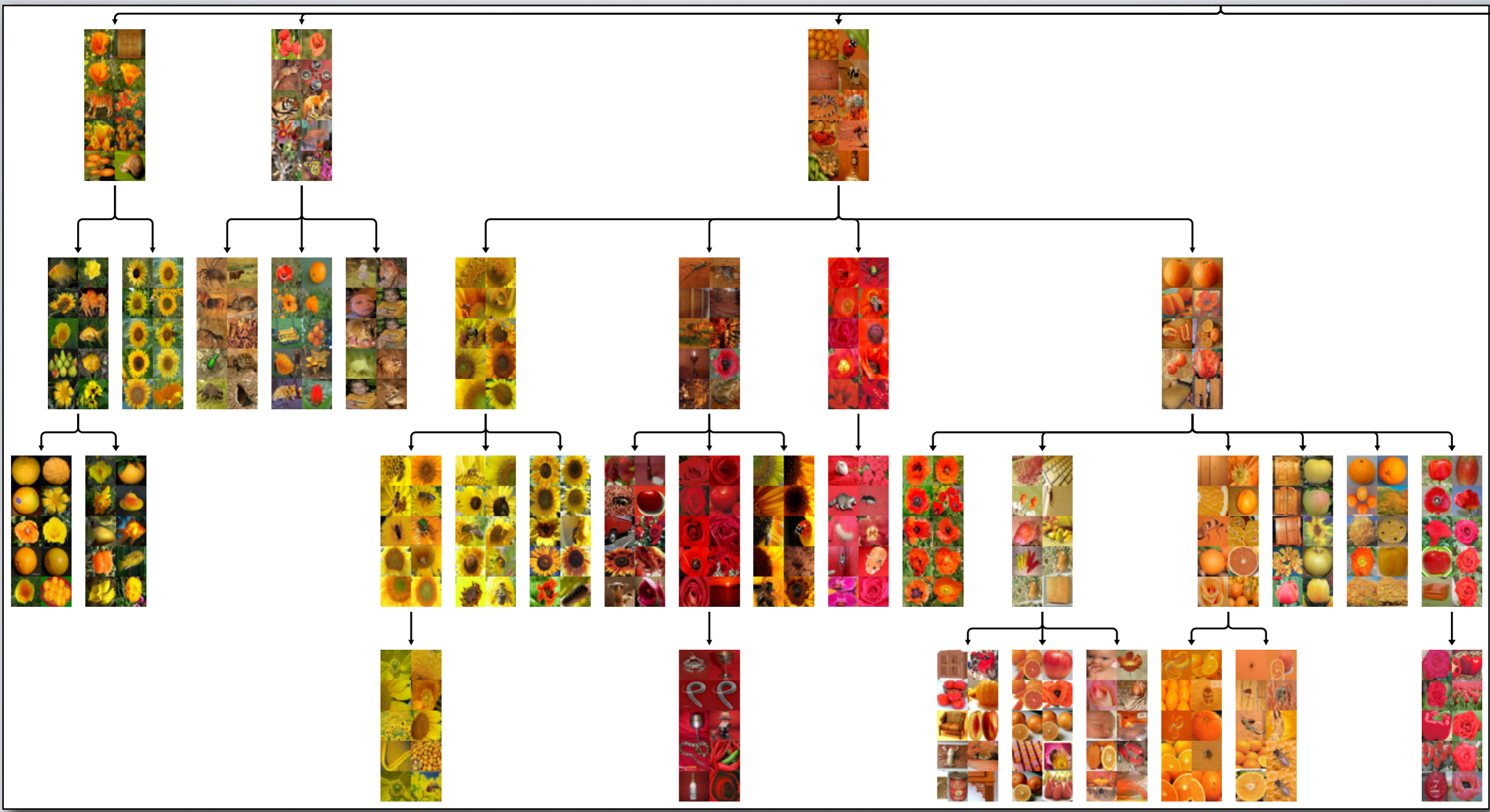


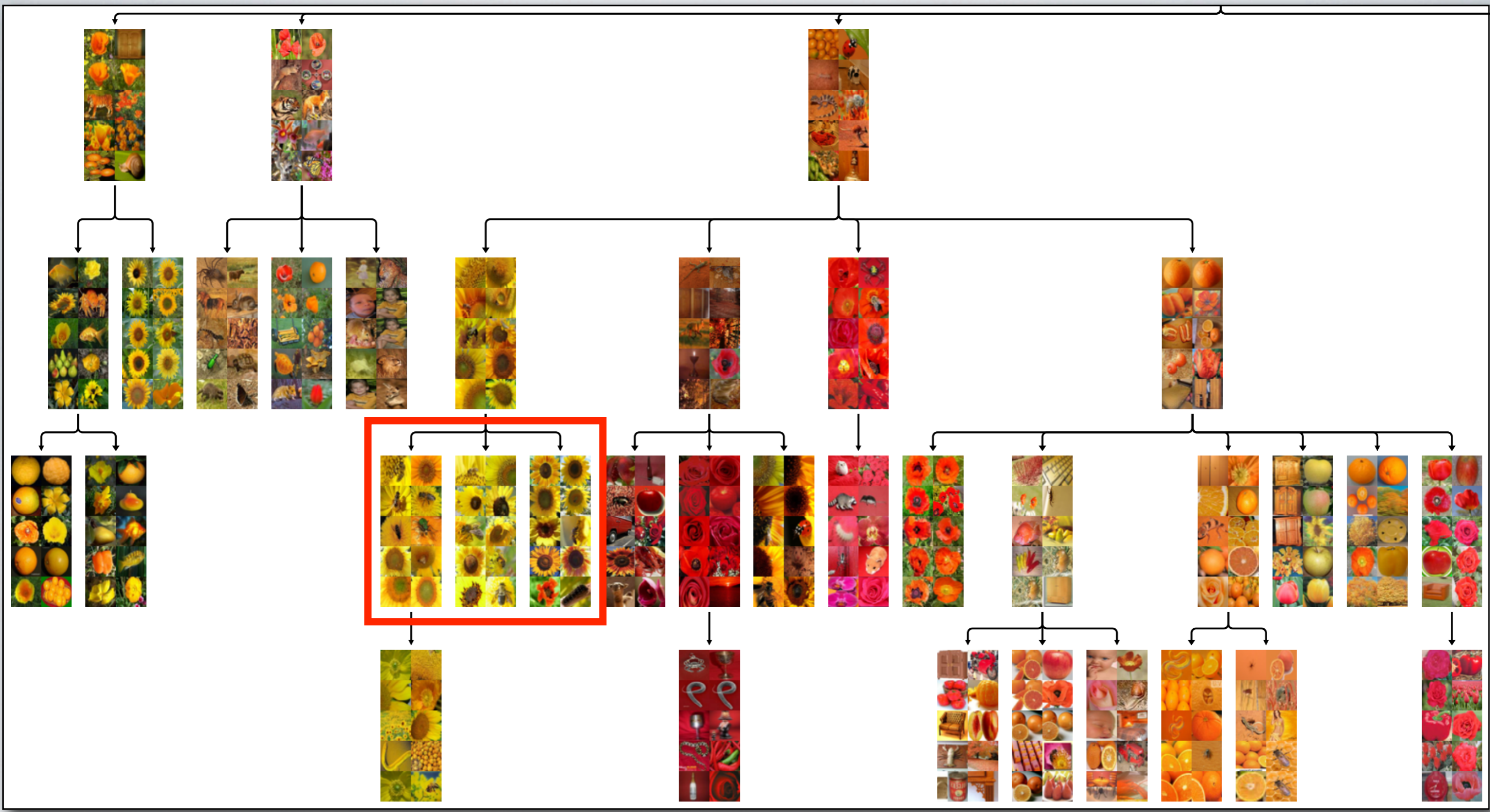


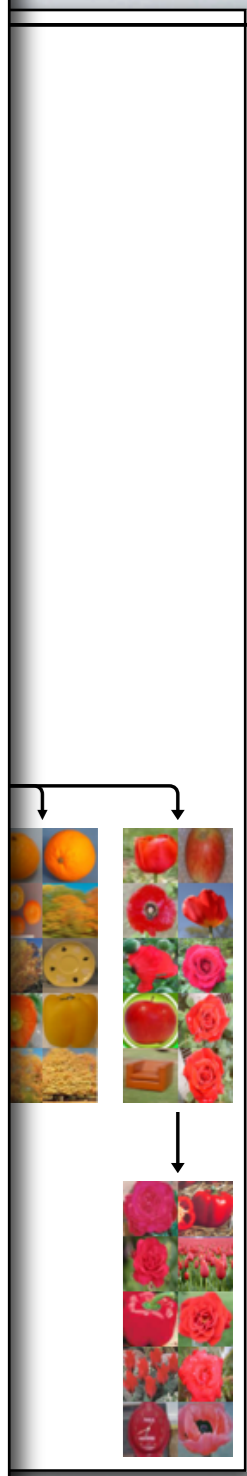
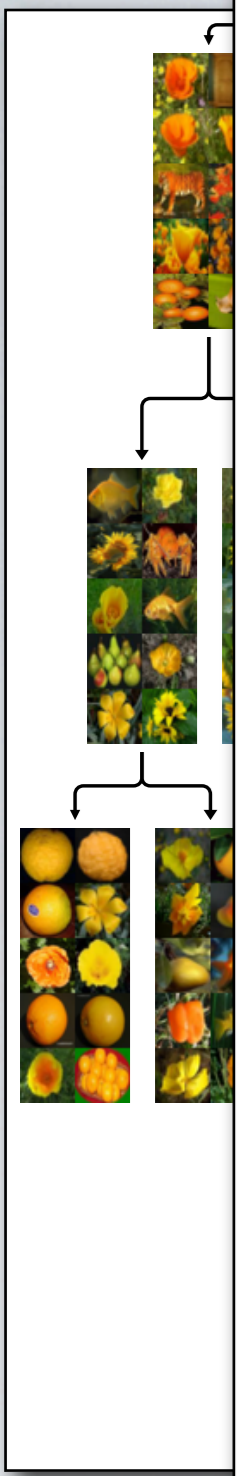
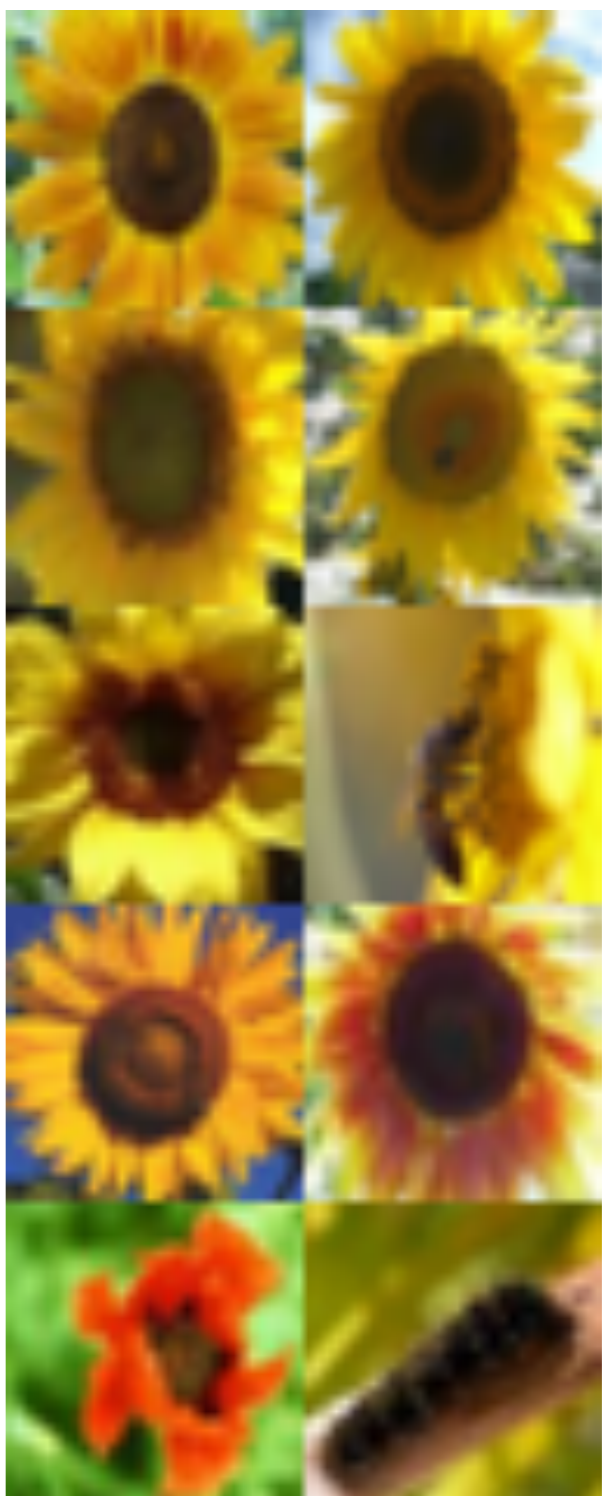
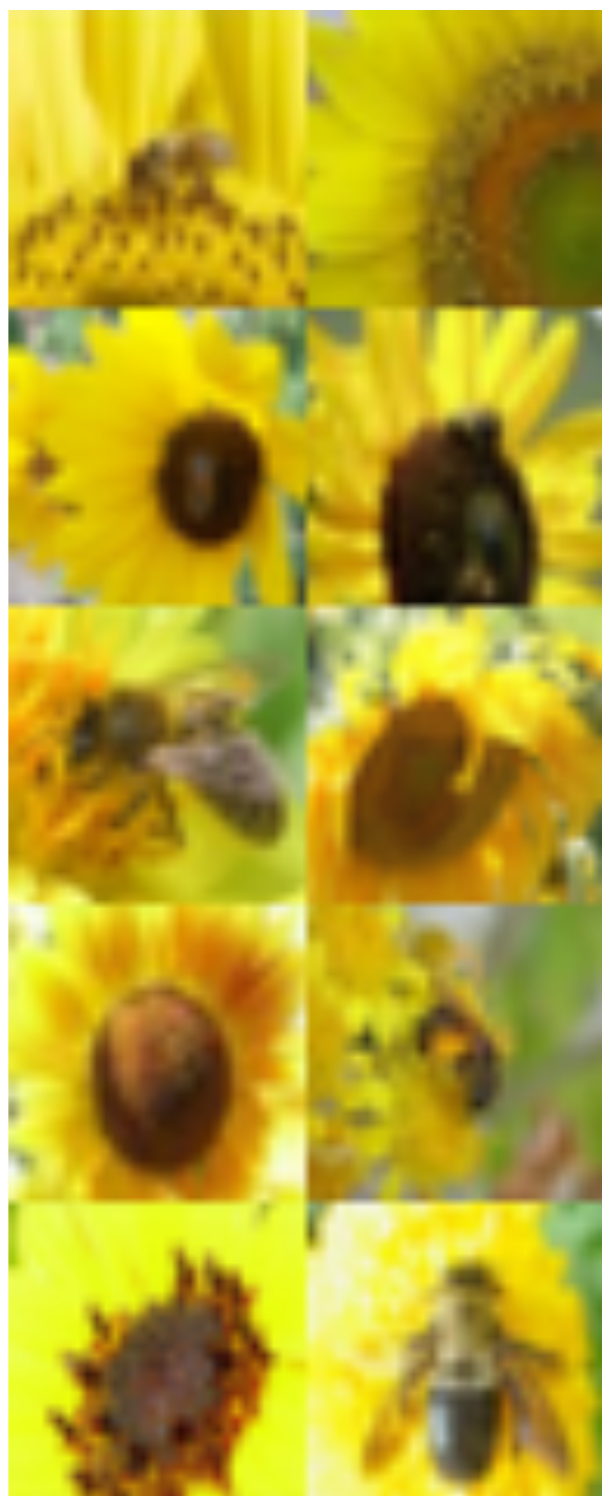
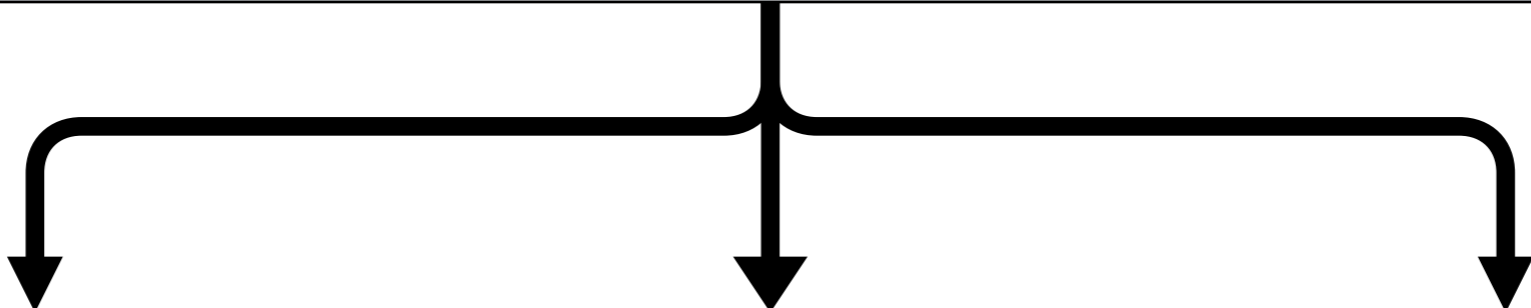




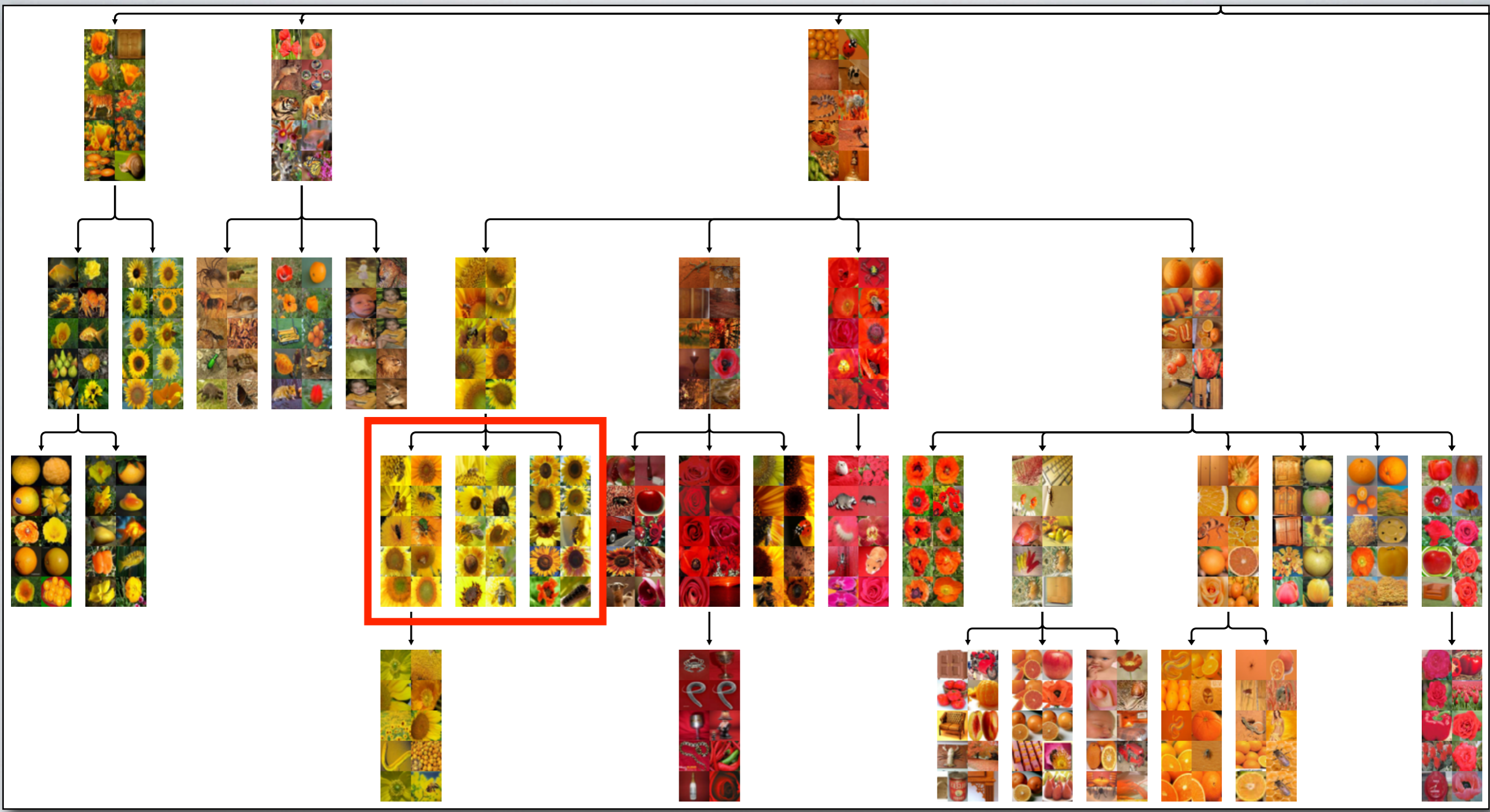








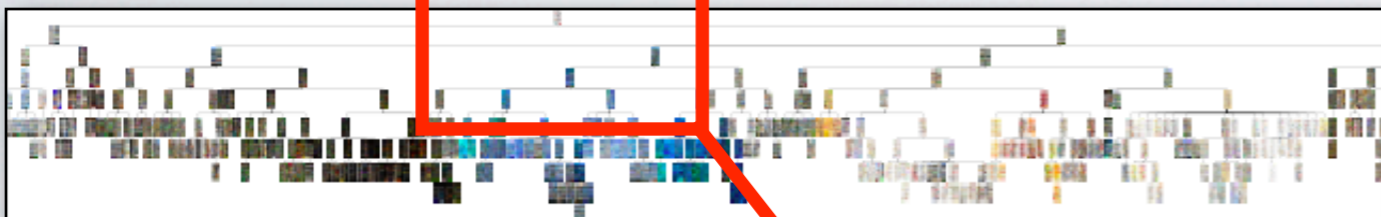




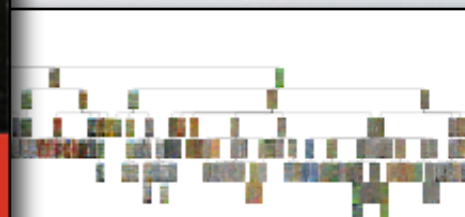


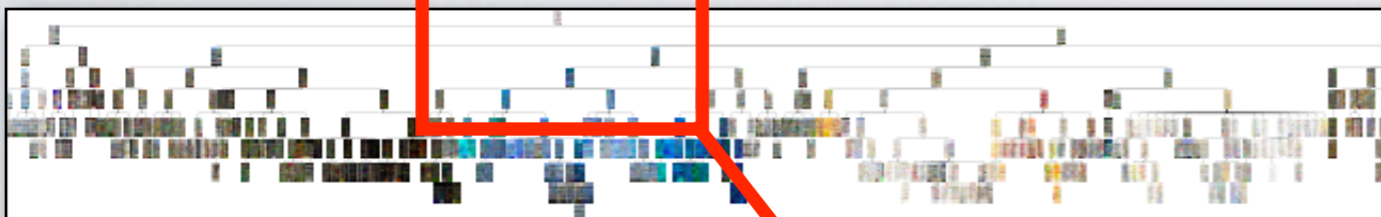
?



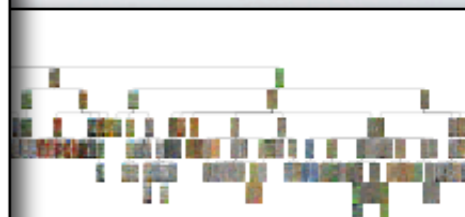


?





?



# DOCUMENT MODELING

- Hierarchical clustering with LDA-style topic model.
- A “topic” is a distribution over words.
- A document has a distribution over topics.
- Words are exchangeable in documents.
- Documents at a node share a single topic distribution.

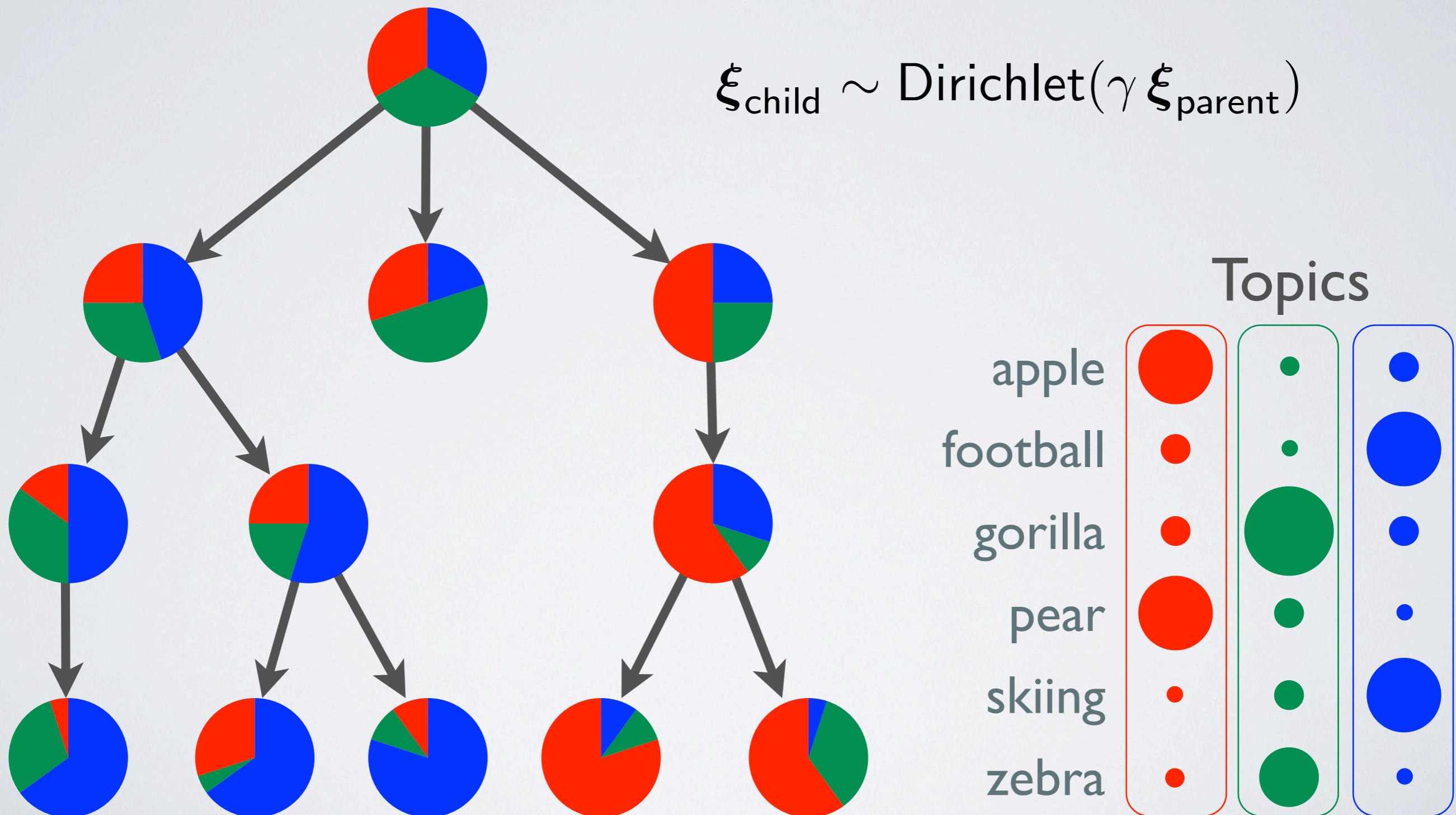
Difference with nested Chinese restaurant process<sup>1</sup>:  
nCRP has a *tree over topics*, and a document is an infinitely-long path down the tree.

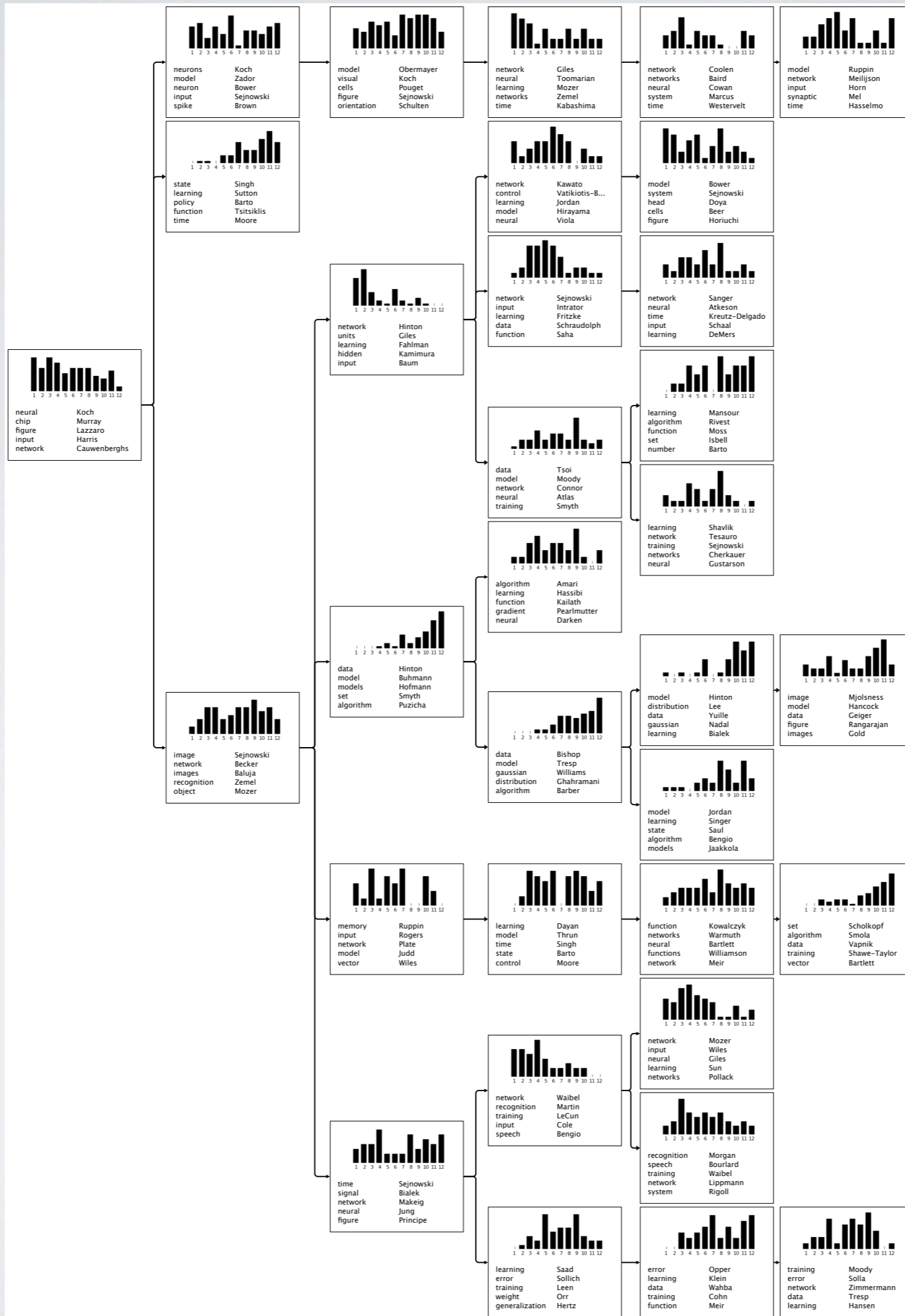
We applied the model and inference to NIPS 1-12.

1. Blei, Griffiths and Jordan. Journal of the ACM, 2010.

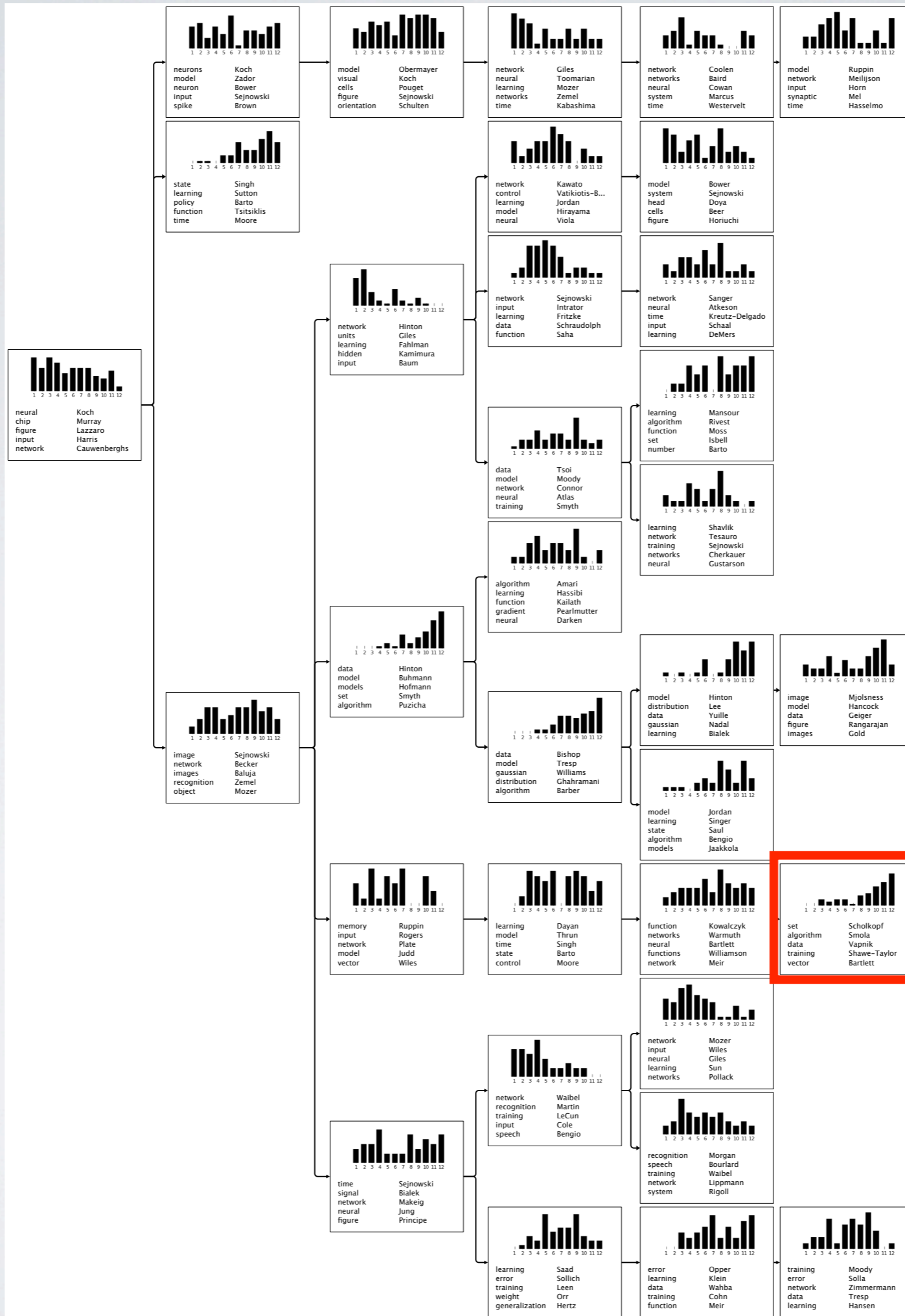
# DOCUMENT MODELING

Each node has a distribution over topics  $\xi$

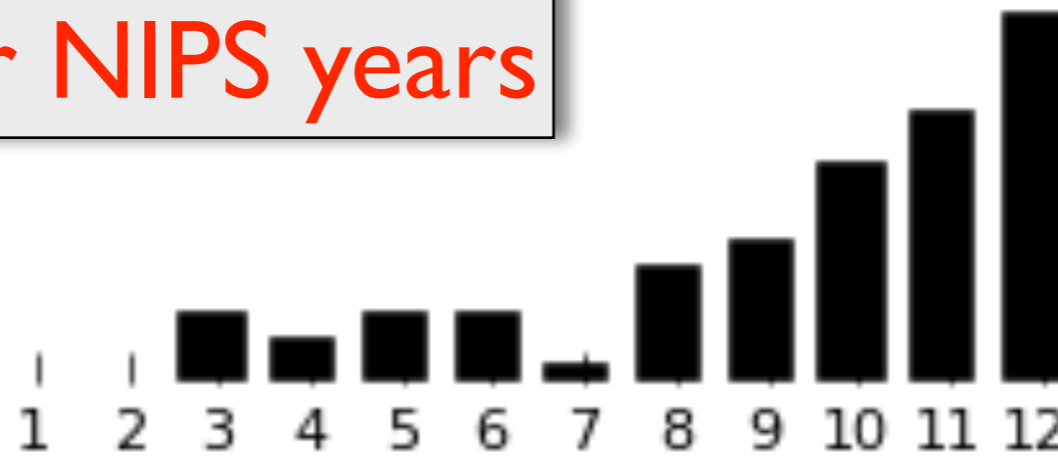








# Histogram over NIPS years

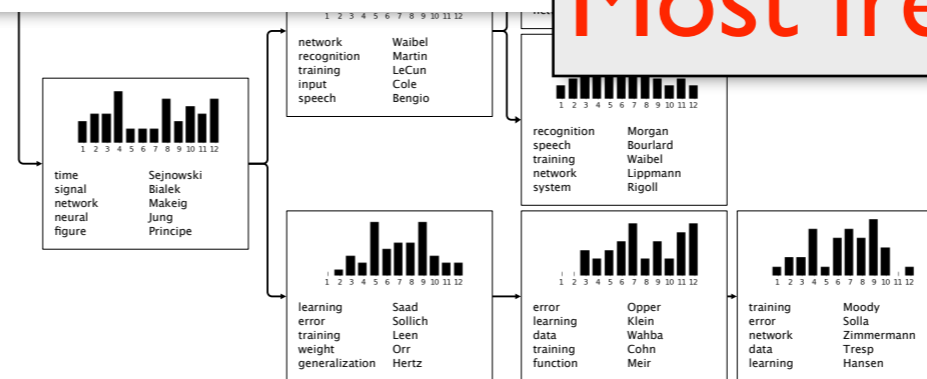


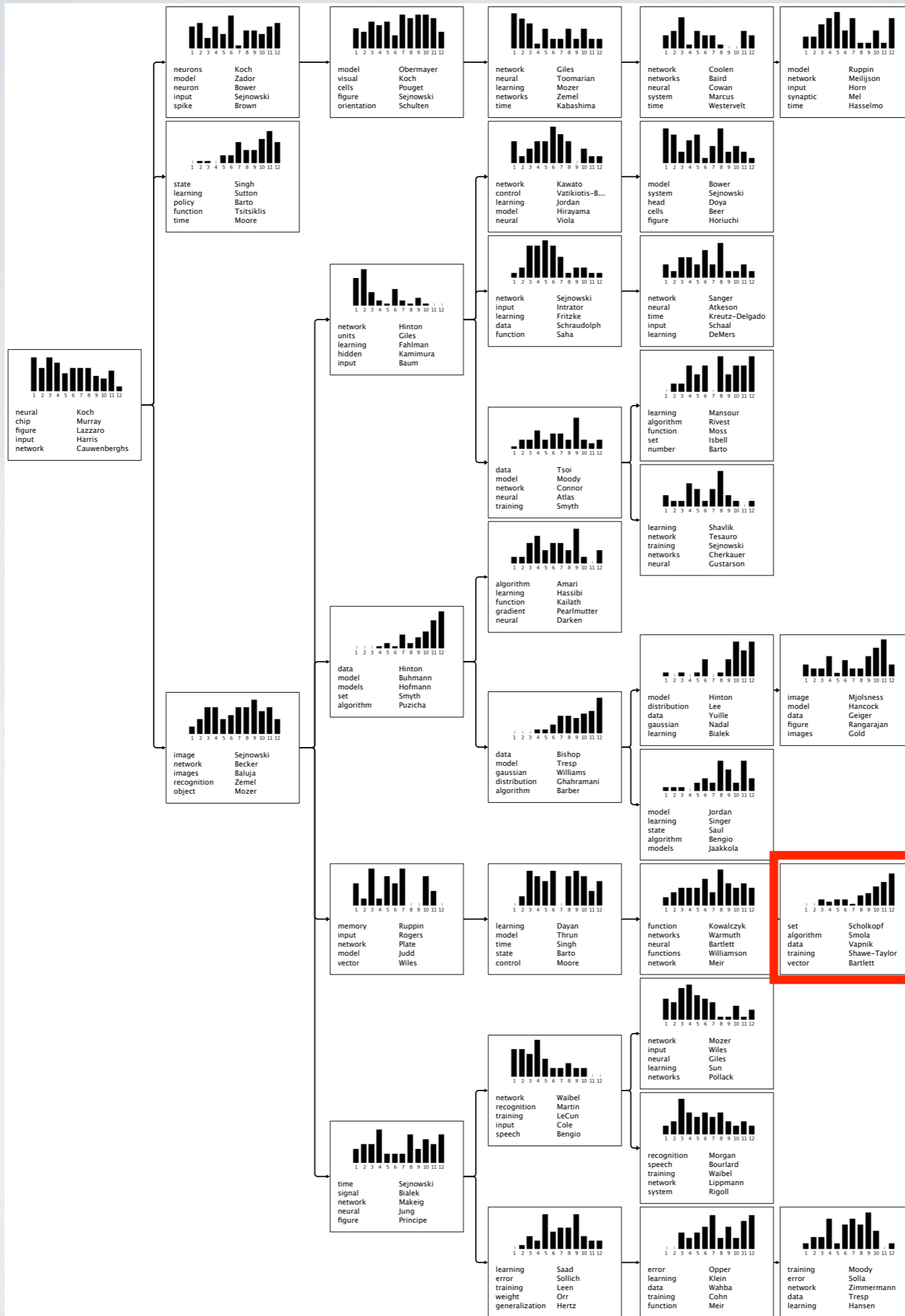
set  
algorithm  
data  
training  
vector

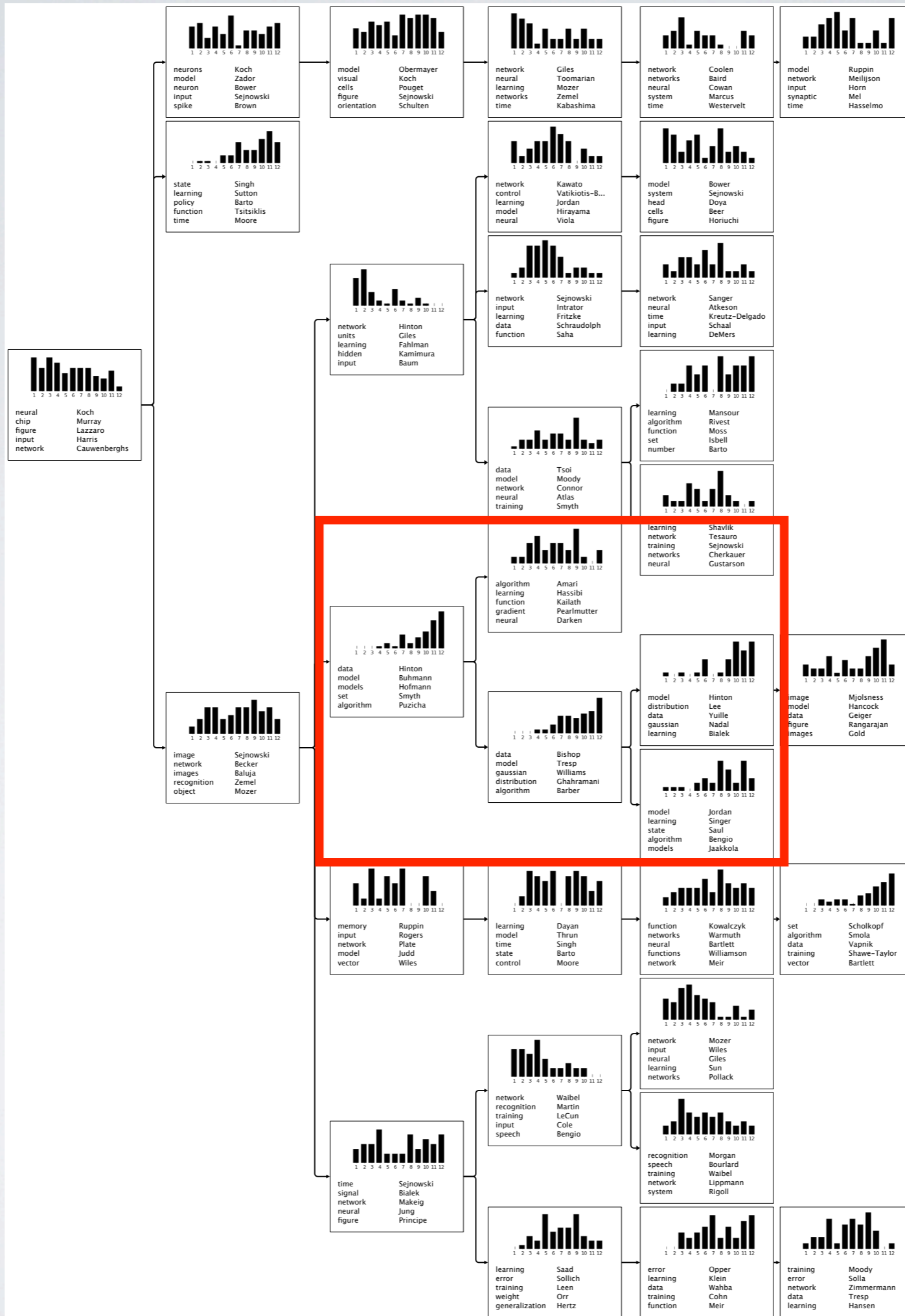
Scholkopf  
Smola  
Vapnik  
Shawe-Taylor  
Bartlett

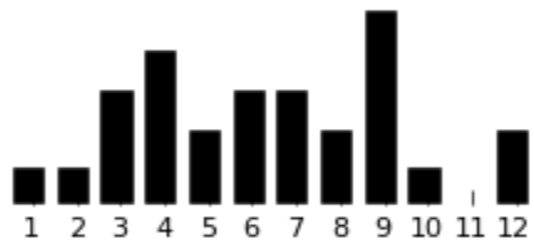
# Most likely words

# Most frequent authors



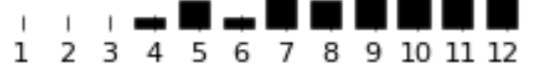







algorithm	Amari
learning	Hassibi
function	Kailath
gradient	Pearlmutter
neural	Darken

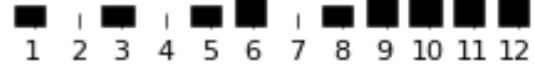
learning	Shavlik
network	Tesauro
training	Sejnowski
networks	Cherkauer
neural	Gustarson




data	Hinton
model	Buhmann
models	Hofmann
set	Smyth
algorithm	Puzicha



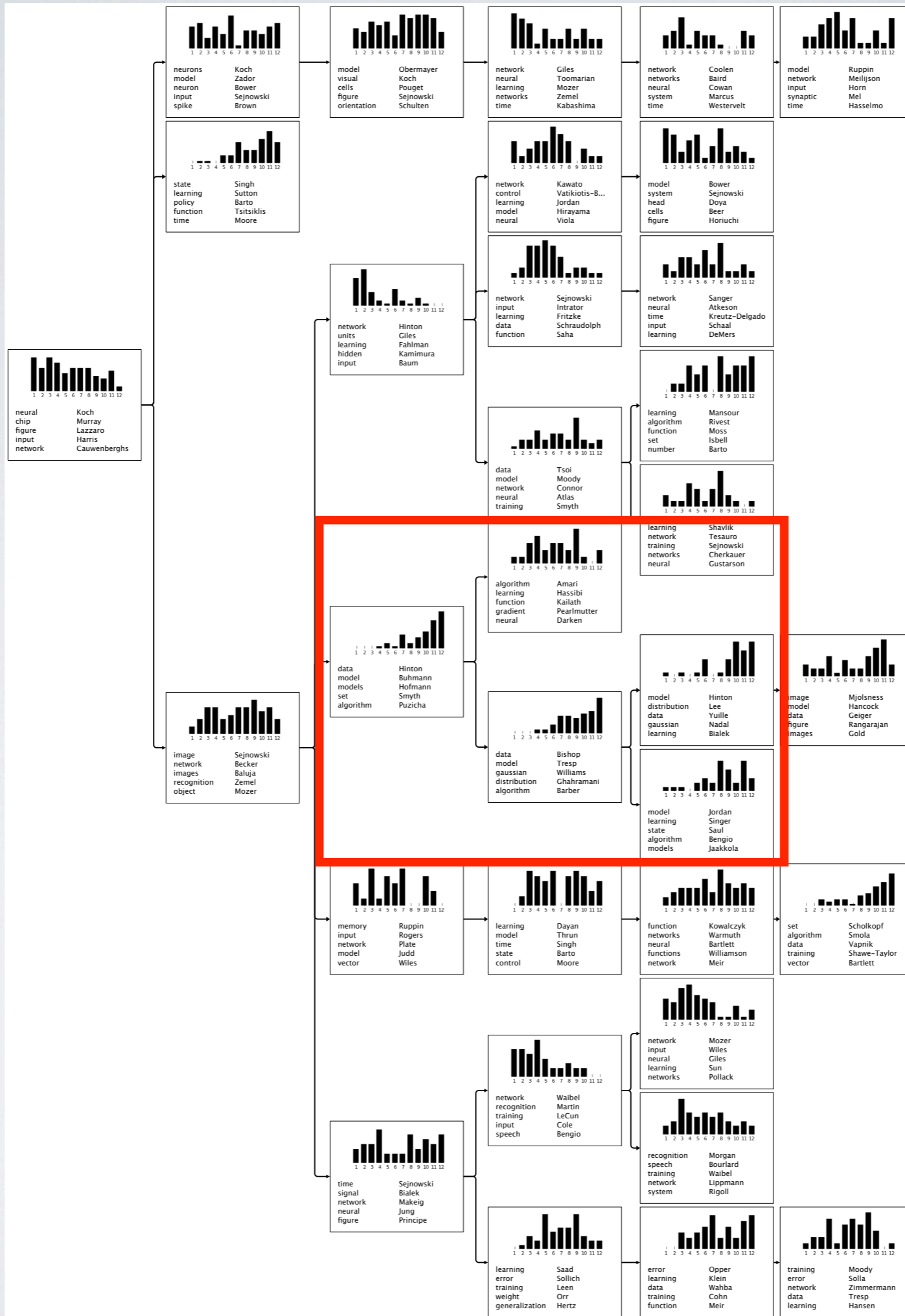
data	Bishop
model	Tresp
gaussian	Williams
distribution	Ghahramani
algorithm	Barber

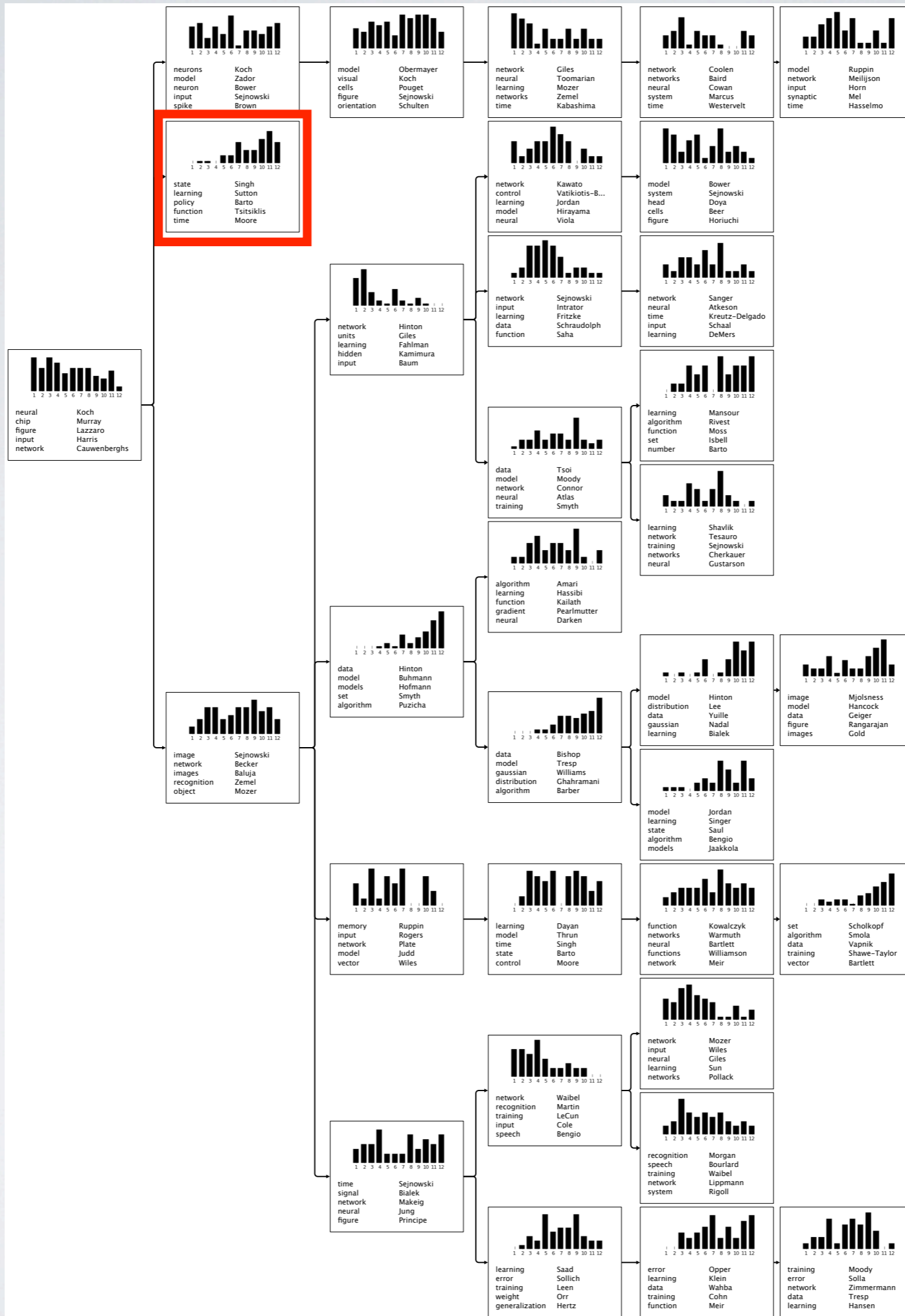


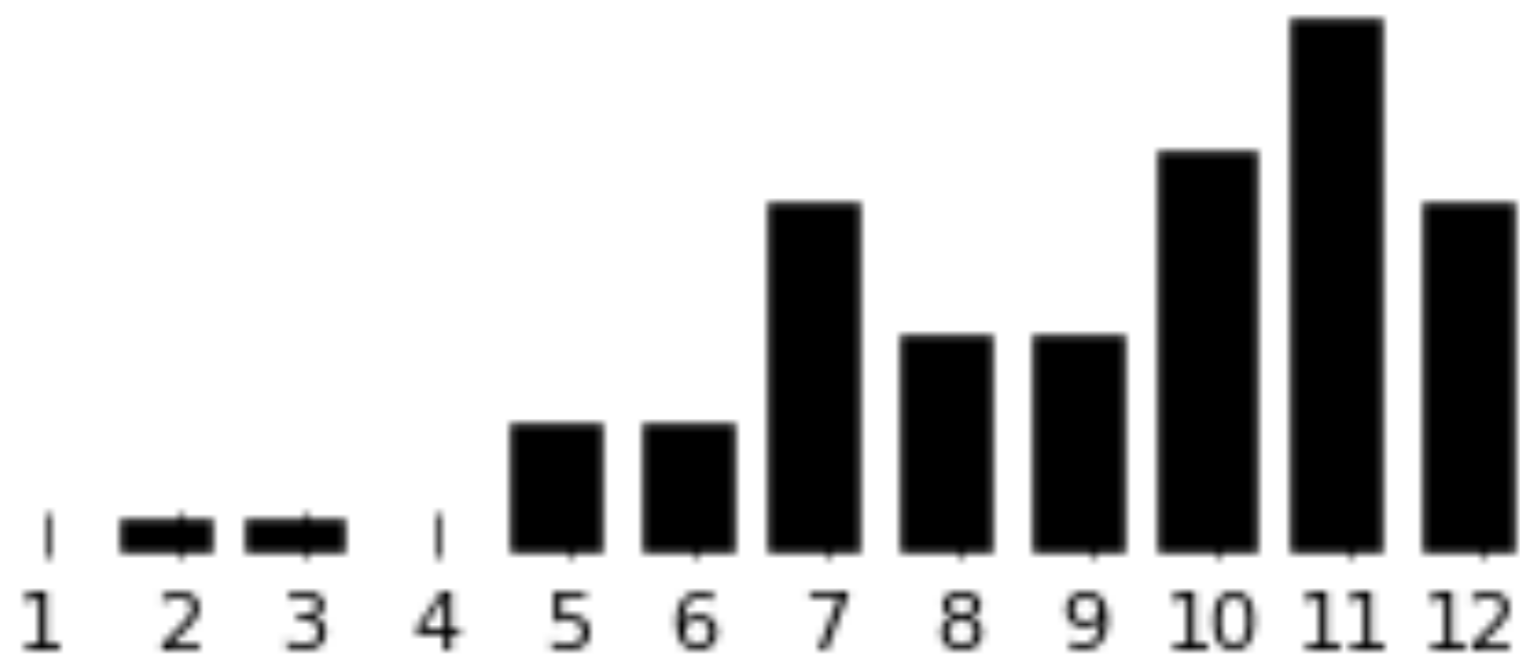
model	Hinton
distribution	Lee
data	Yuille
gaussian	Nadal
learning	Bialek



model	Jordan
learning	Singer
state	Saul
algorithm	Bengio
models	Jaakkola







state  
learning  
policy  
function  
time

Singh  
Sutton  
Barto  
Tsitsiklis  
Moore



# SUMMARY

- A generative prior on tree-structured measures:
  - Unbounded width and depth
  - Data live at internal nodes at finite depth
  - Data are infinitely exchangeable
  - Corresponds to a Blackwell-MacQueen urn model
- We perform inference via MCMC.
- Generates an interesting visualization of images.
- Compares well with LDA in modeling NIPS text.
- Thanks to: Alex Krizhevsky, Kurt Miller, Iain Murray, Yee Whye Teh, Hanna Wallach, Sinead Williamson