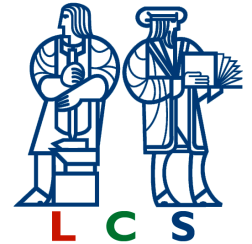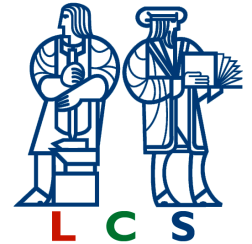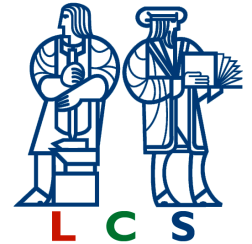# Project Overview

- ## Two themes:
  - Build models of how people learn and structure semantic concepts, and test those models with behavioral experiments.
  - Develop new machine learning algorithms that come closer to learning and structuring semantic concepts as people do.

- ## Two specific questions:
  - How do we combine unsupervised and supervised approaches to concept learning, using unsupervised learning to build better hypothesis spaces for supervised learning?
  - What is the large-scale structure of human semantic concepts in natural language, and how can this structure be learned from large text or hypertext corpora?

# Progress Through December 2002

- Combining unsupervised & supervised concept learning.
  - Task: Learn from very few examples + prior domain knowledge.
  - Approach: Bayesian inference over disjunctions of clusters in tree-structured taxonomy, with Occam's razor prior.
  - Results: Matches a wide range of human inductive inferences, better than best existing models.

- Learning semantics from large text corpora.
  - Task: Bottom-up extraction of semantically meaningful topics.
  - Approach: Embed Latent Dirichlet Allocation (LDA) model - capturing topic-specific semantic structure - in Hidden Markov Model - capturing topic-independent word order regularities.
  - Results: Discovers a factorization of word senses into syntactic and semantic structures, and returns cleaner semantics.

# Research Plan for the Next Six Months

- Combining unsupervised & supervised concept learning.
  - Applications to taxonomic domains: Amazon, WWW.
  - Sample task: Given the links on a web page as a random sample of all relevant pages, rank all other sites in open WWW directory by their estimated probability of relevance.
  - Challenge: Develop efficient, principled approximations to full Bayesian sum over all hypotheses, by exploiting tree structure.

- Learning semantics from large text corpora.
  - Applications to single-sentence information retrieval and learning perceptual grounding of words in images.
  - Sample tasks: Question answering, catalog browsing.
  - Challenges: Develop efficient approximate inference algorithms for LDA / HMM model, and richer models for representing images and overlapping semantic fields (shape, color, texture).