# Inferring Networks from Distances: the "Landscape" of Glycosidase Protein Structures

Sandhya Prabhakaran, Volker Roth
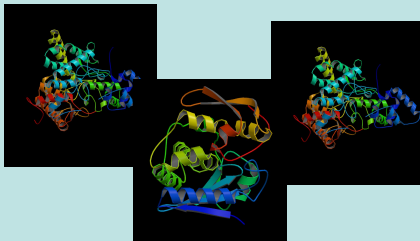
Department of Mathematics and Computer Science,
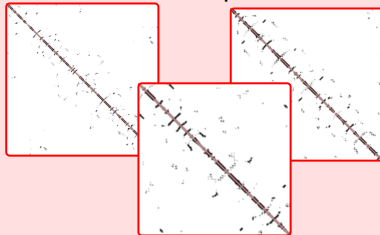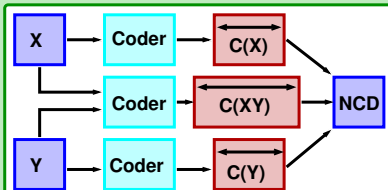University of Basel, Switzerland,

U N I
B A S E L

Nov 21, 2011

# Graphical Abstract

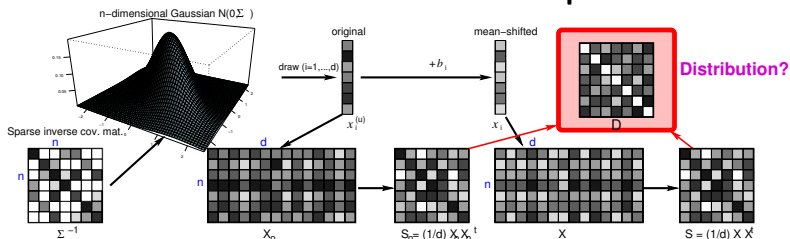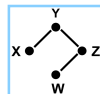# Bayesian Inference

- Idea: construct **probabilistic generative model**, **condition on observed data** and **infer distribution of parameters**.



- **We observe pairwise distances** $D$.
- $D \sim$ **singular Wishart**, parametrized by inverse covariance $\Psi$.
- **Link to Gaussian graphical models:**

  $\Psi_{ij} = 0 \rightsquigarrow (i,j)$ **conditionally independent**
  $\rightsquigarrow$ **no edge**.

- Infer $\Psi$ by **Markov Chain Monte Carlo sampling**.

# Compression Distance between Protein Structures

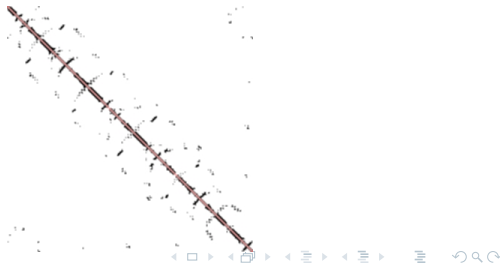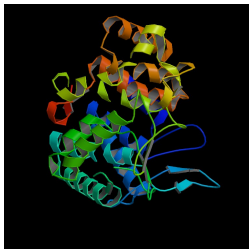- **Normalized Information Distance**: length of the shortest program that computes $x|y$ and $y|x$.

  $NID(x, y) \propto \max\{K(x|y), K(y|x)\} = K(xy) - \min\{K(x), K(y)\}$

- Approximation: **Normalized Compression Distance:**

  $$NCD(x, y) \propto C(xy) - \min\{C(x), C(y)\}.$$

- **In our application:** $x, y$ are vectorized **contact maps** for **Glycosidase enzymes** in **Escherichia coli**.