

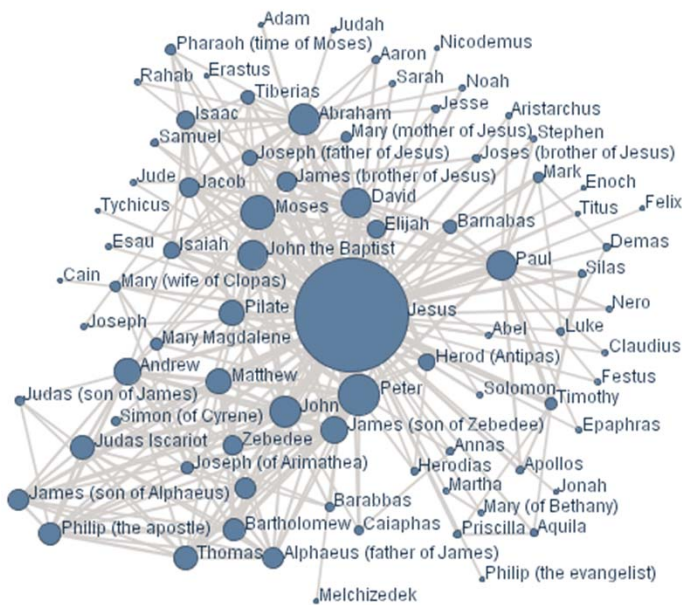
CSC 4020 Fundamentals of Machine Learning: Introduction to Probabilistic Graphical Models

Baoyuan Wu

April 12

What Are Graphical Models?

Graph



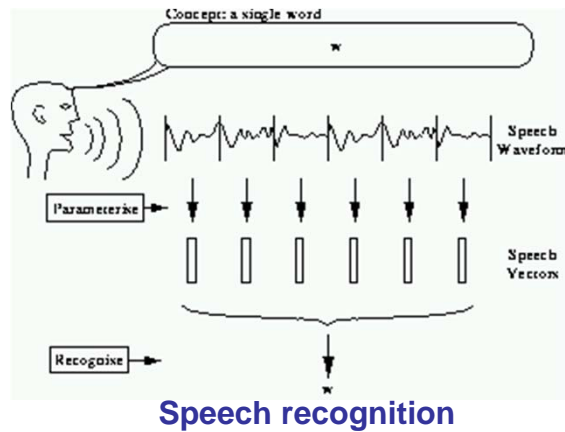
Model

\mathcal{M}

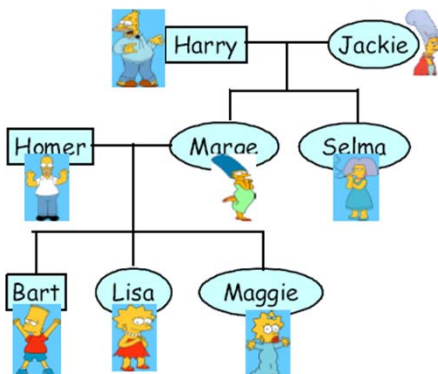
Data

$$\mathcal{D} \equiv \{X_1^{(i)}, X_2^{(i)}, \dots, X_m^{(i)}\}_{i=1}^N$$

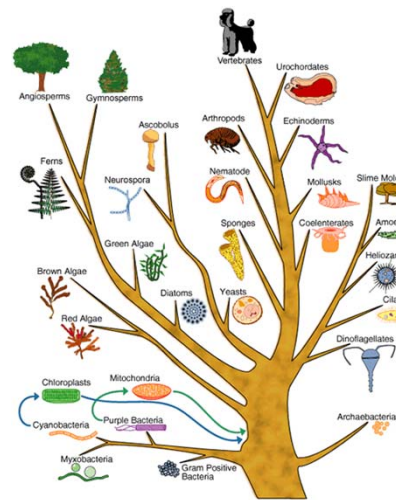
Reasoning under uncertainty!



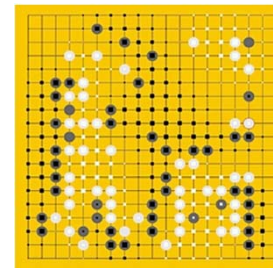
Computer vision



Pedigree



Evolution



Games



Robotic control



Planning

The Fundamental Questions

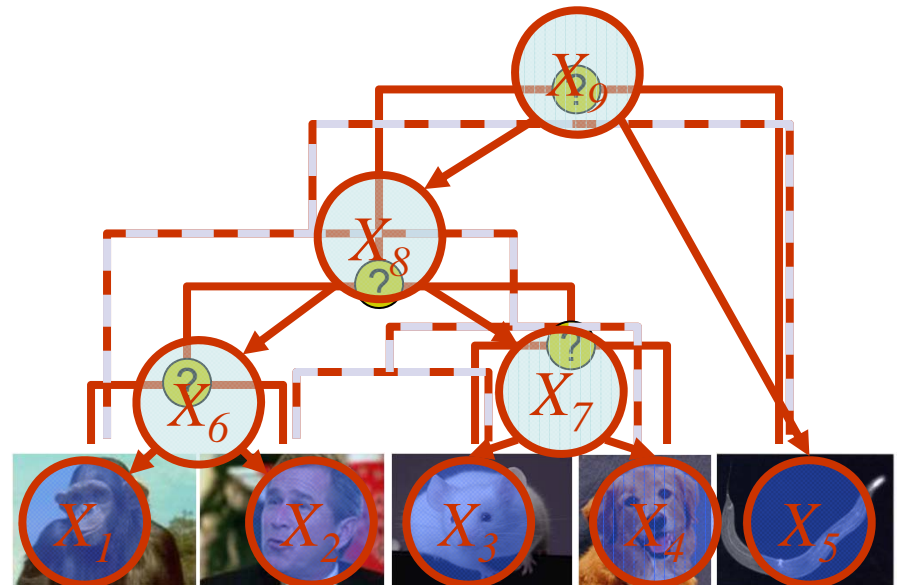
- Representation
 - How to capture/model uncertainties in possible worlds?
 - How to encode our domain knowledge/assumptions/constraints?

- Inference
 - How do I answers questions/queries according to my model and/or based given data?

e.g.: $P(X_i | \mathcal{D})$

- Learning
 - What model is "right" for my data?

e.g.: $\mathcal{M} = \arg \max_{\mathcal{M} \in \mathcal{M}} F(\mathcal{D}; \mathcal{M})$

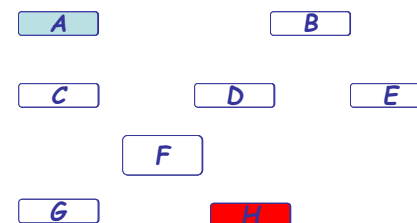


Recap of Basic Prob. Concepts

- Representation: what is the joint probability dist. on multiple variables?

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

- How many state configurations in total? --- 2^8
- Are they all needed to be represented?
- Do we get any scientific/medical insight?**

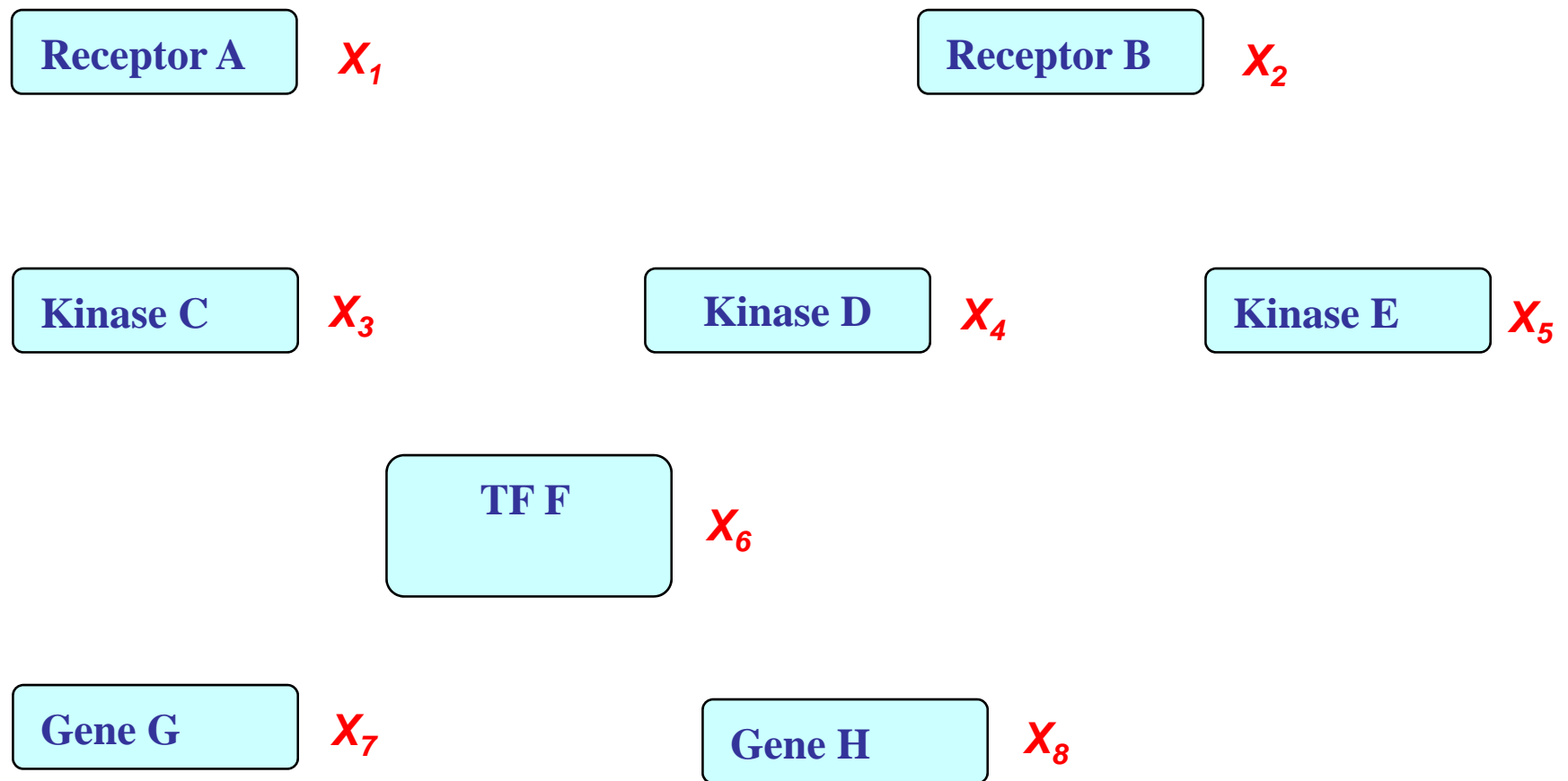


- Learning: where do we get all this probabilities?
 - Maximal-likelihood estimation? but how many data do we need?
 - Are there other est. principles?
 - Where do we put domain knowledge in terms of plausible relationships between variables, and plausible values of the probabilities?
- Inference: If not all variables are observable, how to compute the conditional distribution of latent variables given evidence?
 - Computing $p(H|A)$ would require summing over all 2^6 configurations of the unobserved variables

What is a Graphical Model?

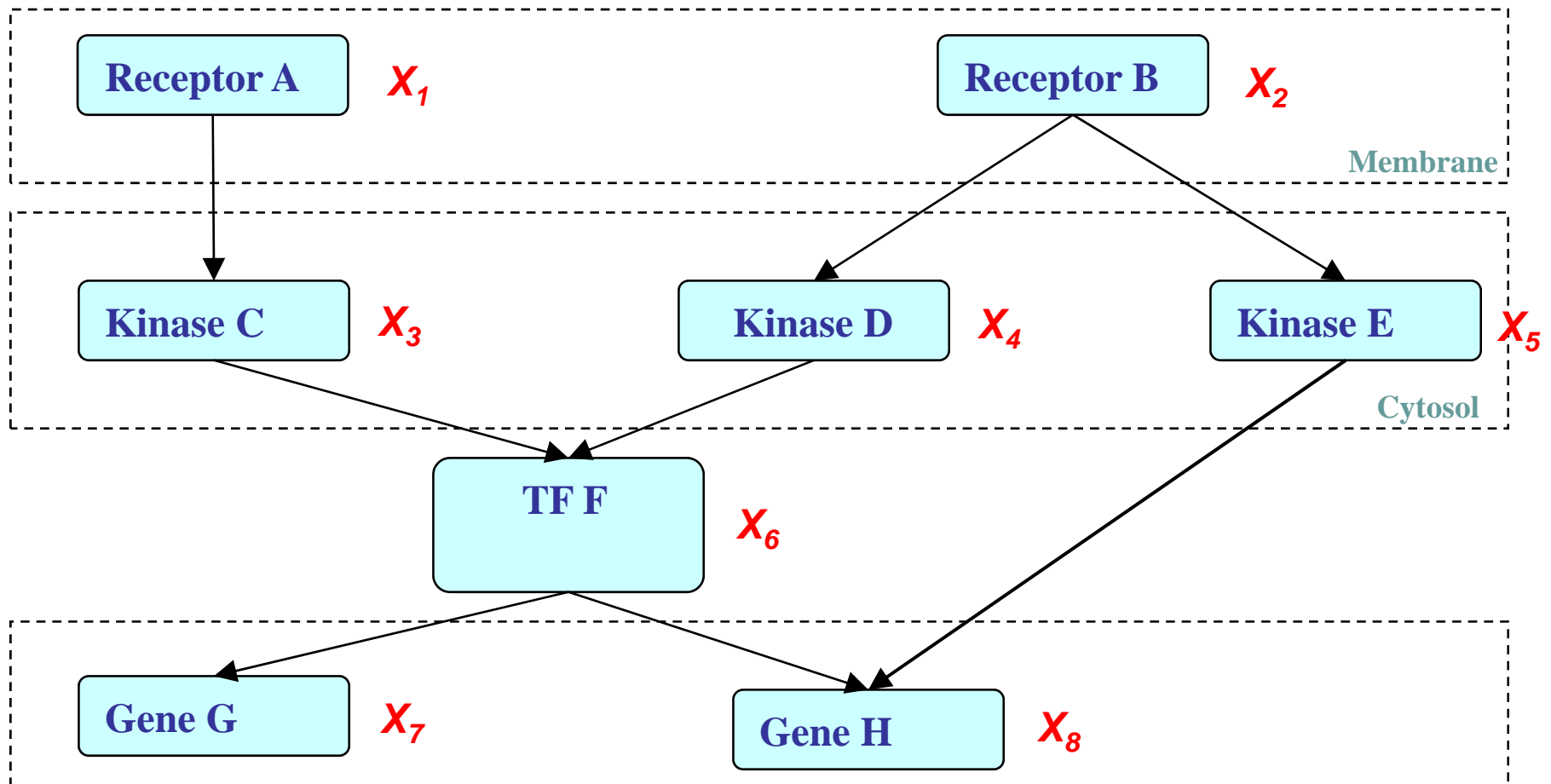
--- Multivariate Distribution in High-D Space

- A possible world for cellular signal transduction:



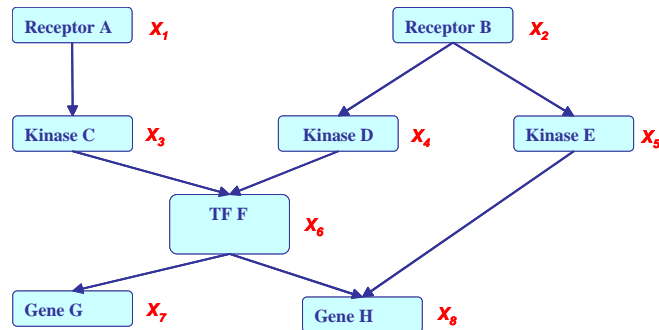
GM: Structure Simplifies Representation

- Dependencies among variables



Probabilistic Graphical Models

- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,

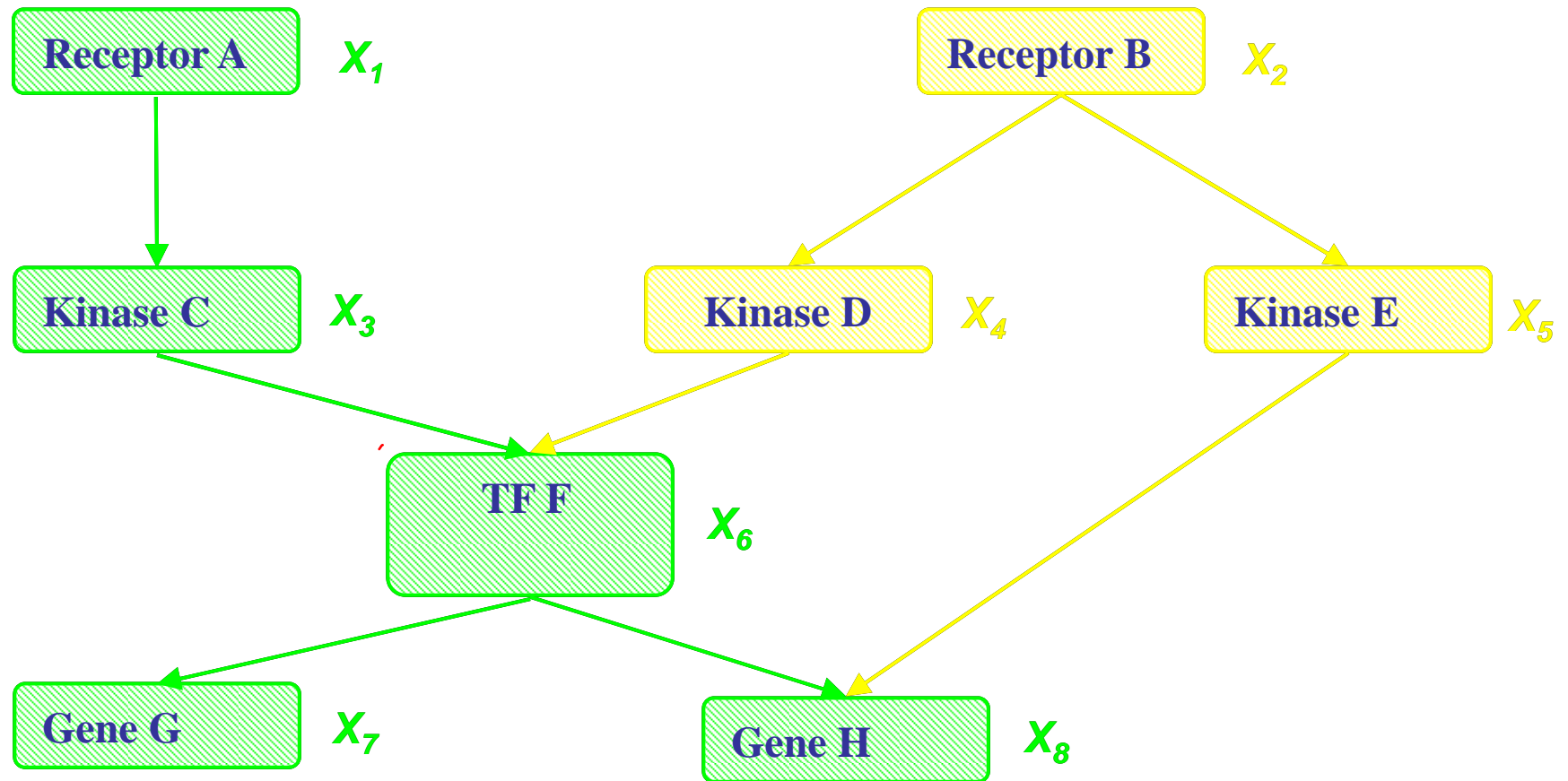


$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ &= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2) \\ &\quad P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6) \end{aligned}$$

Stay tune for what are these independencies!

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
1+1+2+2+2+4+2+4=18, a 16-fold reduction from 2^8 in representation cost !

GM: Data Integration

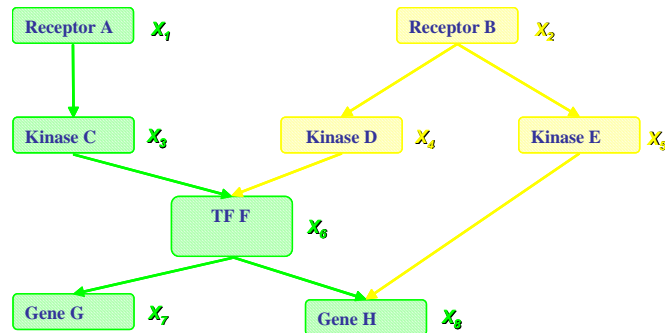


More Data Integration

- Text + Image + Network → Holistic Social Media
- Genome + Proteome + Transcriptome + Phenome + ... →
PanOmic Biology

Probabilistic Graphical Models

- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= P(X_2) P(X_4/X_2) P(X_5/X_2) P(X_1) P(X_3/X_1) \\
 &\quad P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)
 \end{aligned}$$

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
 $2+2+4+4+4+8+4+8=36$, an 8-fold reduction from 2^8 in representation cost !
 - Modular combination of heterogeneous parts – data fusion

Rational Statistical Inference

The Bayes Theorem:

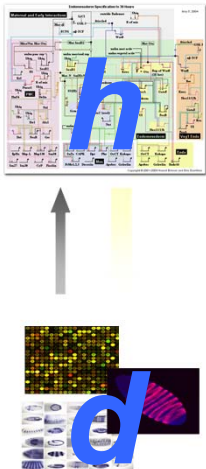
Posterior probability

Likelihood

Prior probability

$$p(h | d) = \frac{p(d | h) p(h)}{\sum_{h' \in H} p(d | h') p(h')}$$

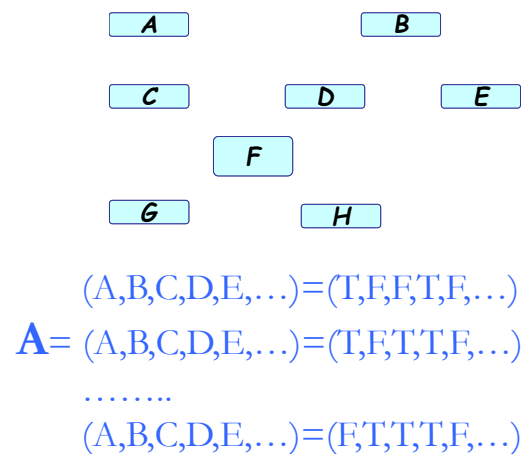
Sum over space of hypotheses



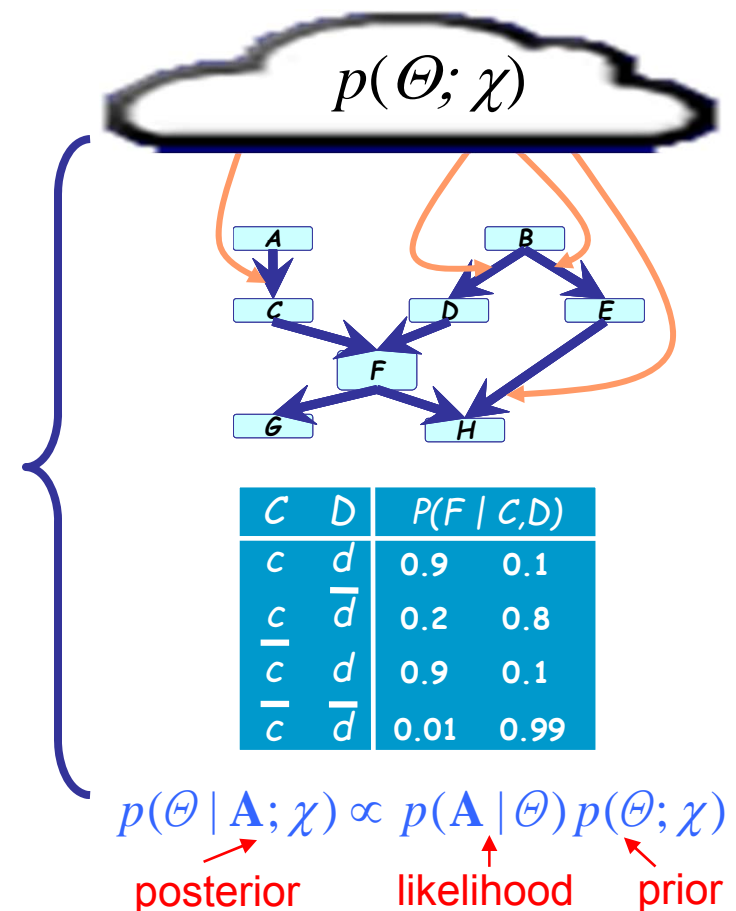
- This allows us to capture uncertainty about the model in a principled way
- But how can we specify and represent a complicated model?
 - Typically the number of genes need to be modeled are in the order of thousands!

GM: MLE and Bayesian Learning

- Probabilistic statements of Θ is conditioned on the values of the observed variables \mathbf{A}_{obs} and prior $p(\chi)$

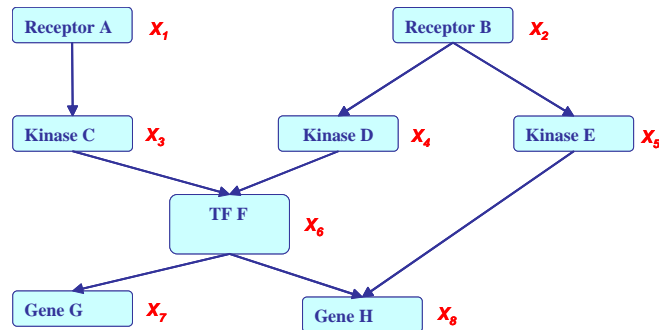


$$\Theta_{\text{Bayes}} = \int \Theta p(\Theta | \mathbf{A}, \chi) d\Theta$$



Probabilistic Graphical Models

- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2) \\
 &\quad P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)
 \end{aligned}$$

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
 $2+2+4+4+4+8+4+8=36$, an 8-fold reduction from 2^8 in representation cost !
 - Modular combination of heterogeneous parts – data fusion
 - Bayesian Philosophy
 - Knowledge meets data



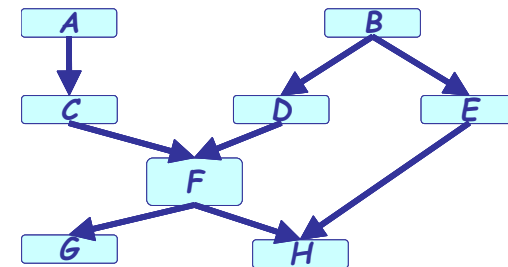
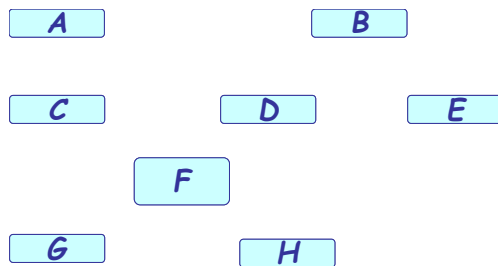
So What is a Graphical Model?

In a nutshell:

GM = Multivariate Statistics + Structure

What is a Graphical Model?

- The informal blurb:
 - It is a smart way to write/specify/compose/design exponentially-large probability distributions without paying an exponential cost, and at the same time endow the distributions with *structured semantics*



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$P(X_{1:8}) = P(X_1)P(X_2)P(X_3 | X_1X_2)P(X_4 | X_2)P(X_5 | X_2) \\ P(X_6 | X_3, X_4)P(X_7 | X_6)P(X_8 | X_5, X_6)$$

- A more formal description:
 - It refers to a family of distributions on a set of random variables that are compatible with all the probabilistic independence propositions encoded by a graph that connects these variables

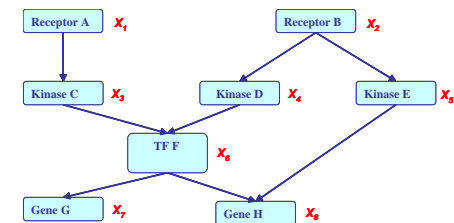
Two types of GMs

- Directed edges give causality relationships (**Bayesian Network** or **Directed Graphical Model**):

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2)$$

$$P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)$$

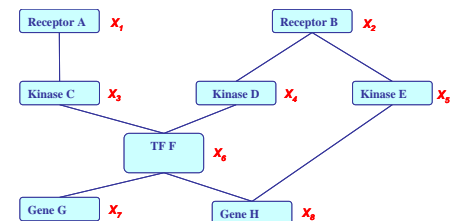


- Undirected edges simply give correlations between variables (**Markov Random Field** or **Undirected Graphical model**):

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= \frac{1}{Z} \exp\{E(X_1) + E(X_2) + E(X_3, X_1) + E(X_4, X_2) + E(X_5, X_2)$$

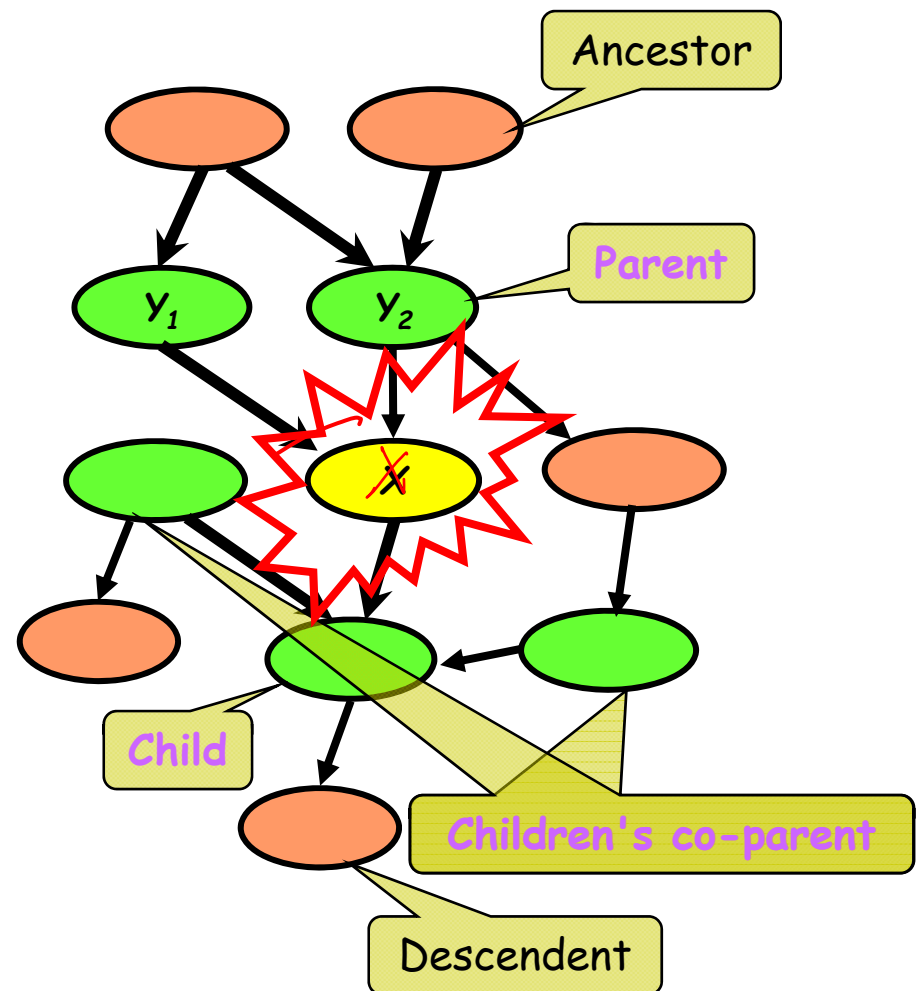
$$+ E(X_6, X_3, X_4) + E(X_7, X_6) + E(X_8, X_5, X_6)\}$$



Bayesian Networks

Structure: *DAG*

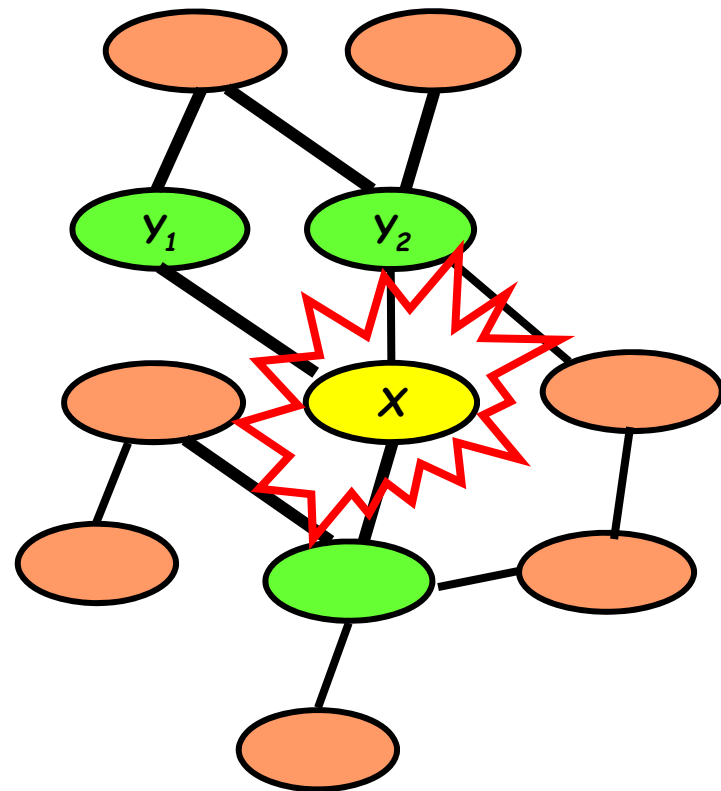
- Meaning: a node is **conditionally independent** of every other node in the network outside its **Markov blanket**
- Local conditional distributions (**CPD**) and the **DAG** completely determine the **joint** dist.
- Give **causality** relationships, and facilitate a **generative** process



Markov Random Fields

Structure: *undirected graph*

- Meaning: a node is **conditionally independent** of every other node in the network given its **Directed neighbors**
- Local contingency functions (**potentials**) and the **cliques** in the graph completely determine the **joint** dist.
- Give **correlations** between variables, but no explicit way to generate samples



Towards structural specification of probability distribution

- Separation properties in the graph imply independence properties about the associated variables
- For the graph to be useful, any conditional independence properties we can derive from the graph should hold for the probability distribution that the graph represents

- **The Equivalence Theorem**

For a graph G ,

Let \mathcal{D}_1 denote the family of all distributions that satisfy $I(G)$,

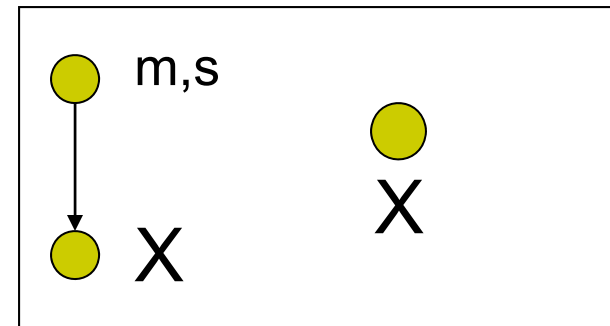
Let \mathcal{D}_2 denote the family of all distributions that factor according to G ,

Then $\mathcal{D}_1 \equiv \mathcal{D}_2$.

GMs are your old friends

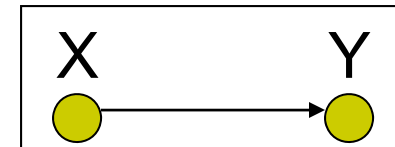
Density estimation

Parametric and nonparametric methods



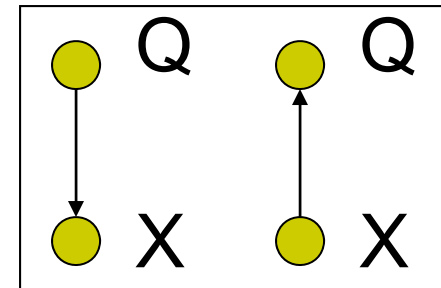
Regression

Linear, conditional mixture, nonparametric



Classification

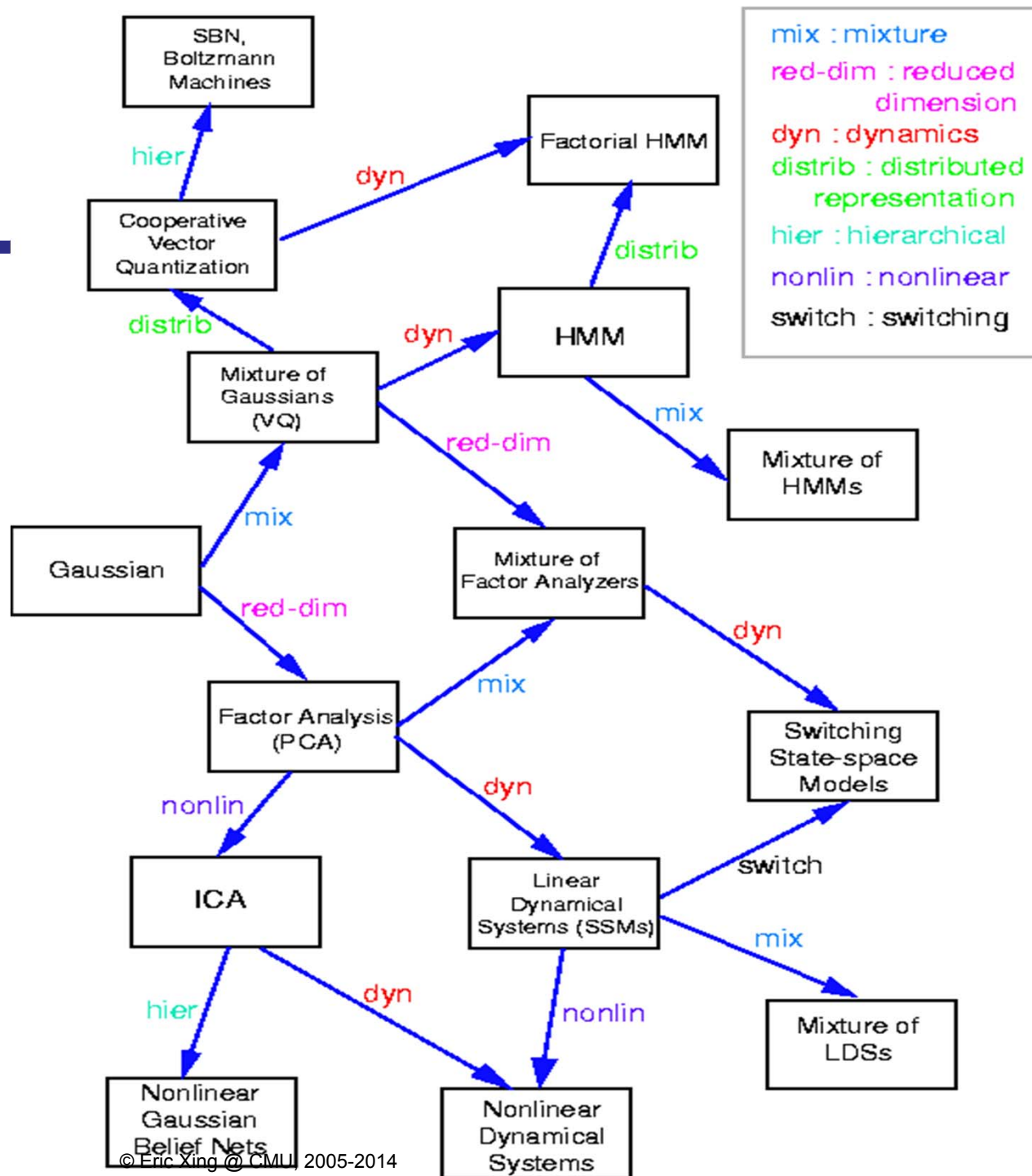
Generative and discriminative approach



Clustering

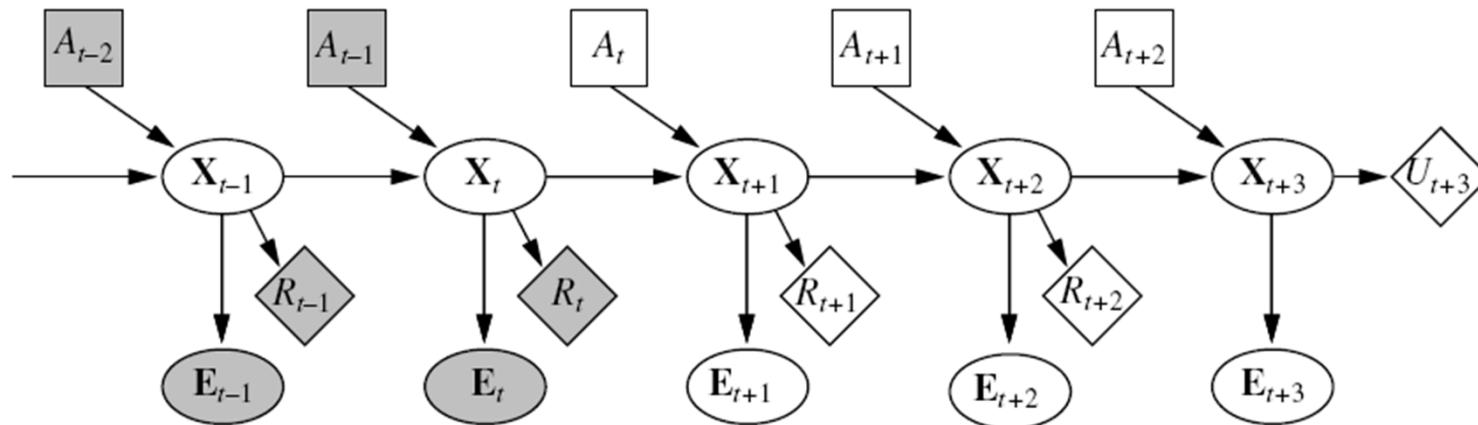
An (incomplete) genealogy of graphical models

(Picture by Zoubin
Ghahramani and
Sam Roweis)

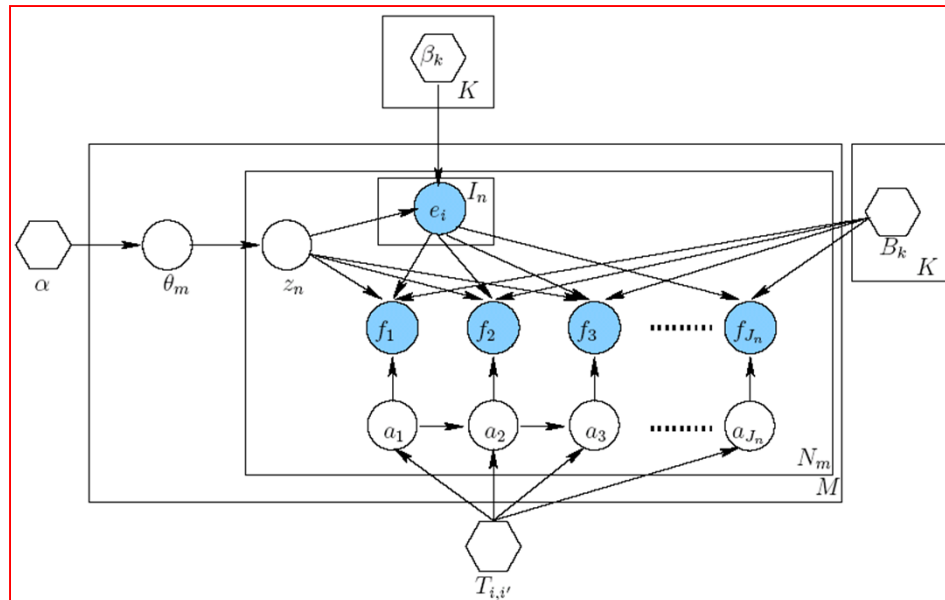
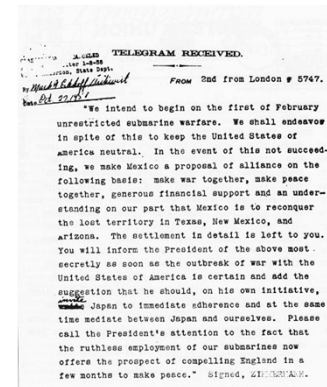
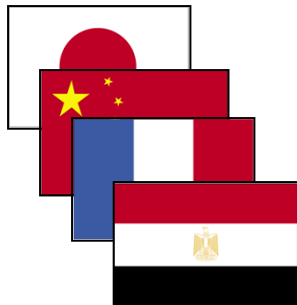


Fancier GMs: reinforcement learning

- Partially observed Markov decision processes (POMDP)

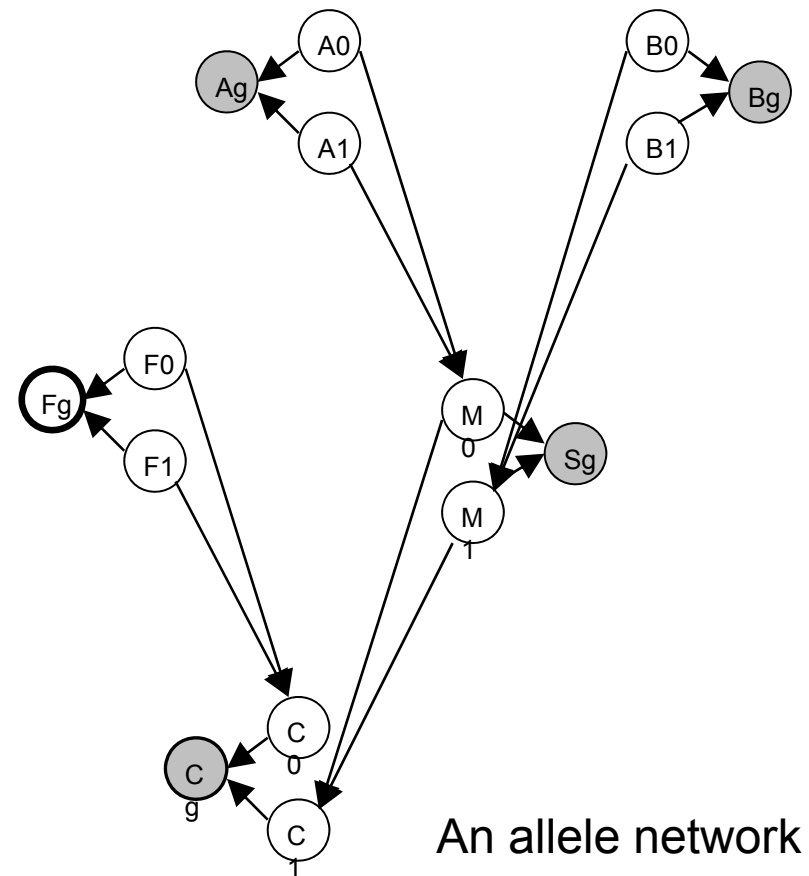
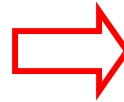
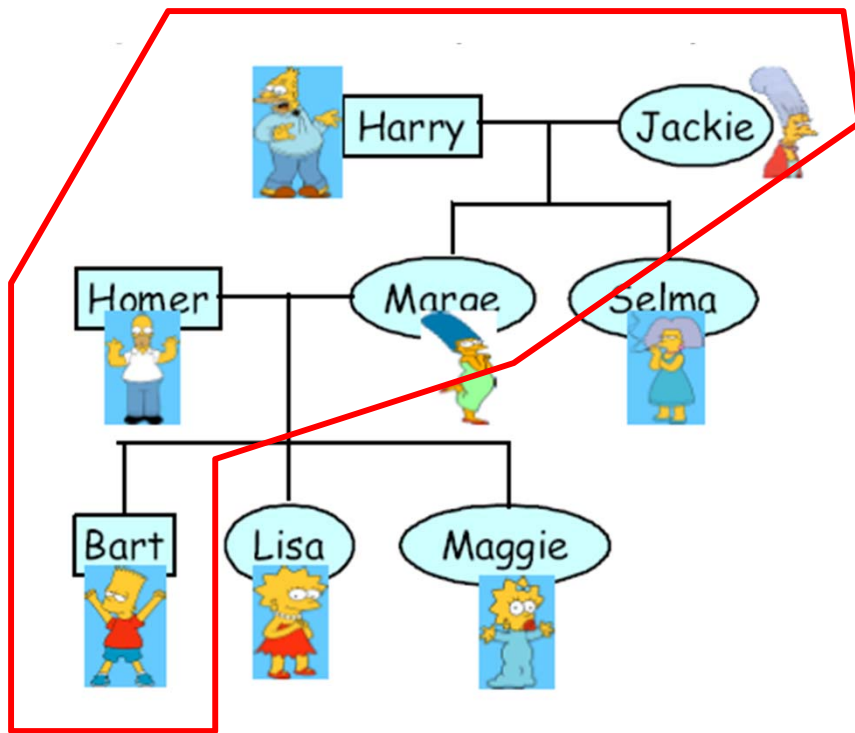


Fancier GMs: machine translation

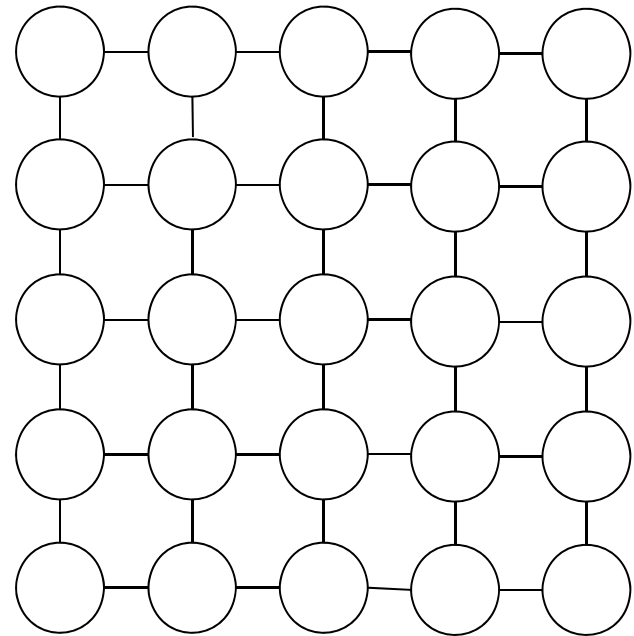
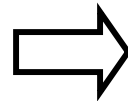
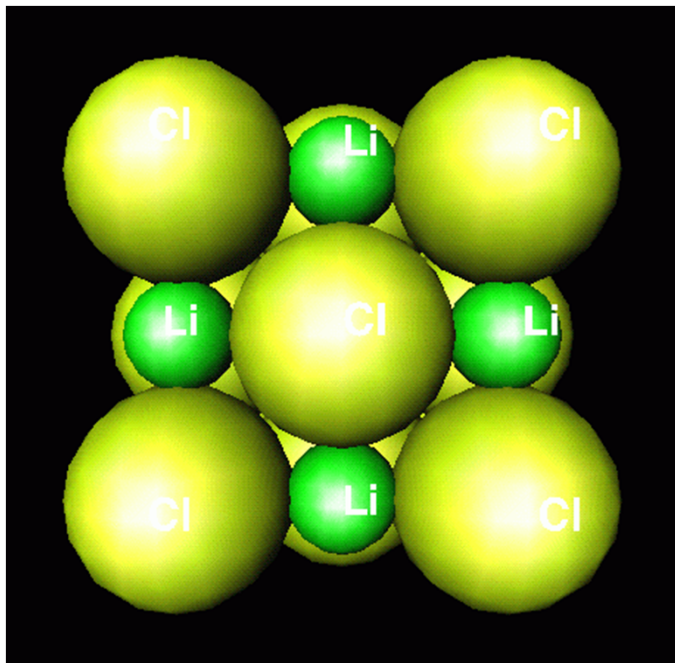


The HM-BiTAM model
(B. Zhao and E.P Xing,
ACL 2006)

Fancier GMs: genetic pedigree



Fancier GMs: solid state physics



Ising/Potts model

Application of GMs

- Machine Learning
- Computational statistics
- Computer vision and graphics
- Natural language processing
- Informational retrieval
- Robotic control
- Decision making under uncertainty
- Error-control codes
- Computational biology
- Genetics and medical diagnosis/prognosis
- Finance and economics
- Etc.

Why graphical models

- A language for communication
 - A language for computation
 - A language for development
-
- Origins:
 - Wright 1920's
 - Independently developed by Spiegelhalter and Lauritzen in statistics and Pearl in computer science in the late 1980's

Why graphical models

- **Probability theory** provides the **glue** whereby the parts are combined, ensuring that the system as a whole is consistent, and providing ways to interface models to data.
- The **graph theoretic** side of graphical models provides both an intuitively appealing interface by which humans can model highly-interacting sets of variables as well as a data structure that lends itself naturally to the design of efficient general-purpose algorithms.
- **Many of the classical multivariate probabilistic systems** studied in fields such as statistics, systems engineering, information theory, pattern recognition and statistical mechanics **are special cases of the general graphical model formalism**
- The graphical model framework provides a way to view all of these systems as instances of a **common underlying formalism**.

A few myths about graphical models

- They require a localist semantics for the nodes ✓
- They require a causal semantics for the edges ✗
- They are necessarily Bayesian ✗
- They are intractable ✓✗