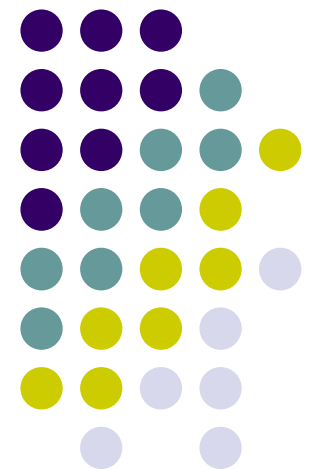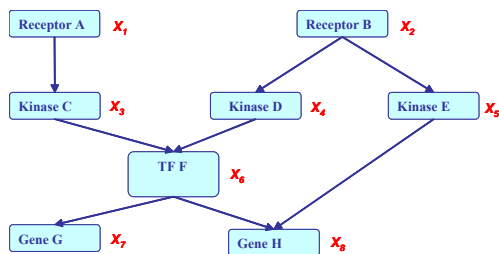# Probabilistic Graphical Models

## Introduction to GM

**Eric Xing**

**Lecture 1, January 18, 2017**

**Reading: see class homepage**

1

# Logistics

- Class webpage:
    - http://www.cs.cmu.edu/~epxing/Class/10708-17/

# Logistics

- Text books:
  - Daphne Koller and Nir Friedman, **Probabilistic Graphical Models**
  - M. I. Jordan, **An Introduction to Probabilistic Graphical Models**
- Mailing Lists:
  - To contact the instructors: 10708-instructor@cs.cmu.edu
  - Class announcements list: 10708-students@cs.cmu.edu.
- TA:
  - Maruan Al-Shedivat, GHC 8223, Office Hour: Wednesday, 4:30 - 5:30pm
  - Haohan Wang, GHC 5507, Office Hour: Friday, 6:00pm - 7:00pm
  - David Dai, GHC 8116, Office hours: TBA

- Lecturers: Eric Xing
- Assistant Instructor: Sarah Schultz
- Class Assistant:
  - Amy Protos, GHC 8221
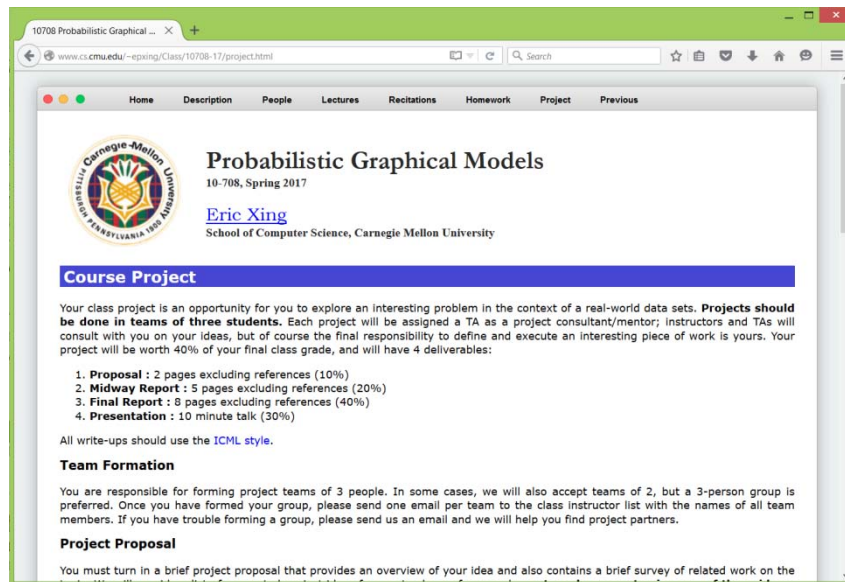- Instruction aids: Piazza

# Logistics

- 4 homework assignments: 40% of grade
  - Theory exercises, Implementation exercises

- Scribe duties: 10% (~once to twice for the whole semester)

- Short reading summary: 10%  (due at the beginning of every lecture)

- Final project: 40% of grade
  - Applying PGM to the development of a real, substantial ML system
    - Design and Implement a (record-breaking) distributed Logistic Regression, Gradient Boosted Tree, Deep Network, or  Topic model on Petuum and apply to ImageNet, Wikipedia, and/or other data
    - Build a web-scale topic or story line tracking system for news media, or a paper recommendation system for conference review matching
    - An online car or people or event detector for web-images and webcam
    - An automatic "what's up here?" or "photo album" service on iPhone
  - Theoretical and/or algorithmic work
    - a more efficient approximate inference or optimization algorithm, e.g., based on stochastic approximation, proximal average, or other new techniques
    - a distributed sampling scheme with convergence guarantee
  - 3-member team to be formed in the first three weeks, proposal, mid-way report, oral presentation & demo, final report, peer review  → possibly conference submission !

# Past projects:



- **Award Winning Projects:**

  J. Yang, Y. Liu, E. P. Xing and A. Hauptmann, **Harmonium-Based Models for Semantic Video Representation and Classification** , *Proceedings of The Seventh SIAM International Conference on Data Mining* **(SDM 2007 best paper)**
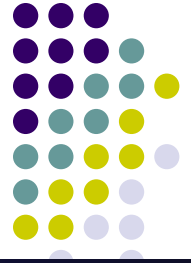
  Manaal Faruqui, Jesse Dodge, Sujay Kumar Jauhar, Chris Dyer, Eduard Hovy, Noah A. Smith, **Retrofitting Word Vectors to Semantic Lexicons**, NAACL 2015 best paper

  Others … such as KDD 2014 best paper

- **Other projects:**

  Andreas Krause, Jure Leskovec and Carlos Guestrin, **Data Association for Topic Intensity Tracking,** *23rd International Conference on Machine Learning* **(ICML 2006).**

  M. Sachan, A. Dubey, S. Srivastava, E. P. Xing and Eduard Hovy, **Spatial Compactness meets Topical Consistency: Jointly modeling Links and Content for Community Detection** , *Proceedings of The 7th ACM International Conference on Web Search and Data Mining* **(WSDM 2014).**
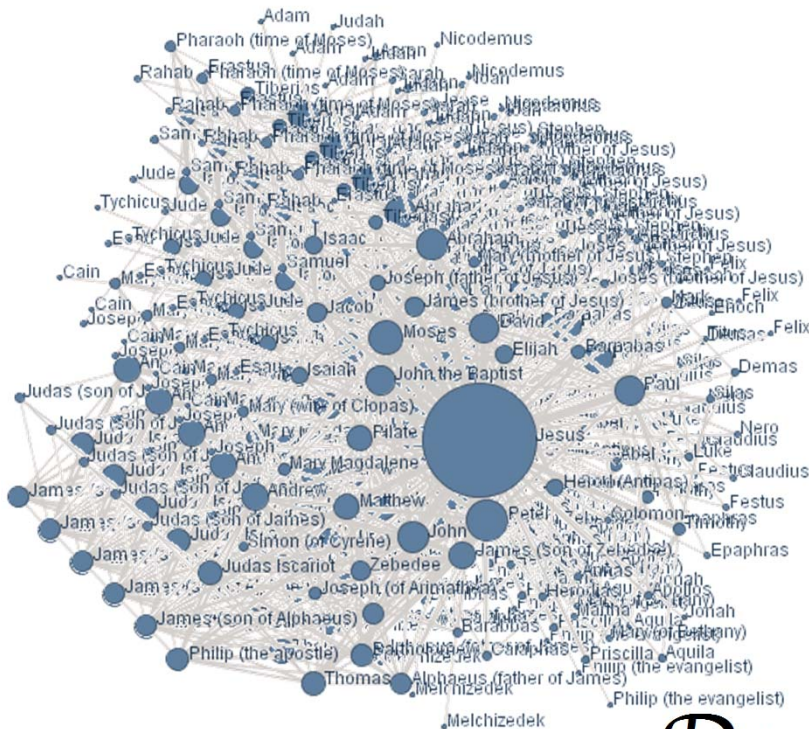
- ## We will have a prize for the best project(s) …

# What Are Graphical Models?
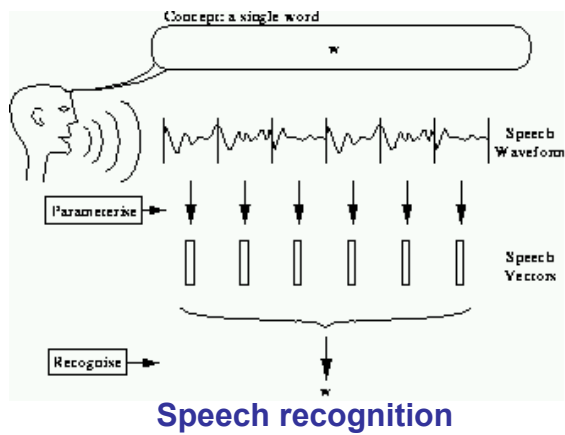
**Graph**

**Model**

$$\mathcal{M}_G$$

**Data**

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

# Reasoning under uncertainty!



Speech recognition



Information retrieval
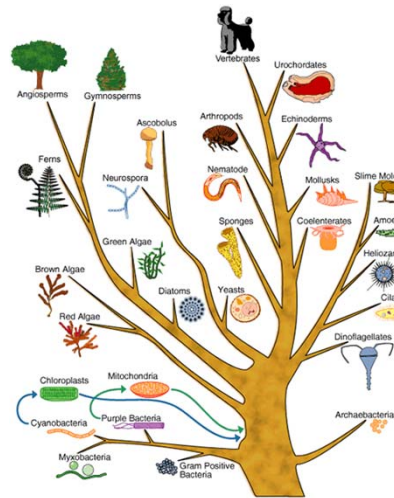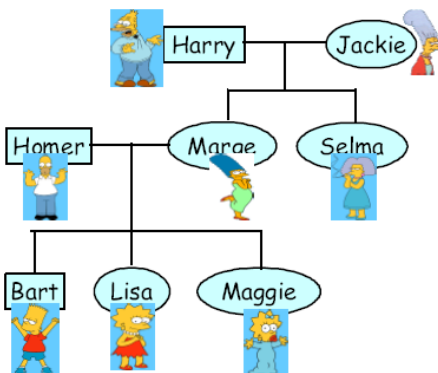


Computer vision



Pedigree



Evolution



Games



Robotic control



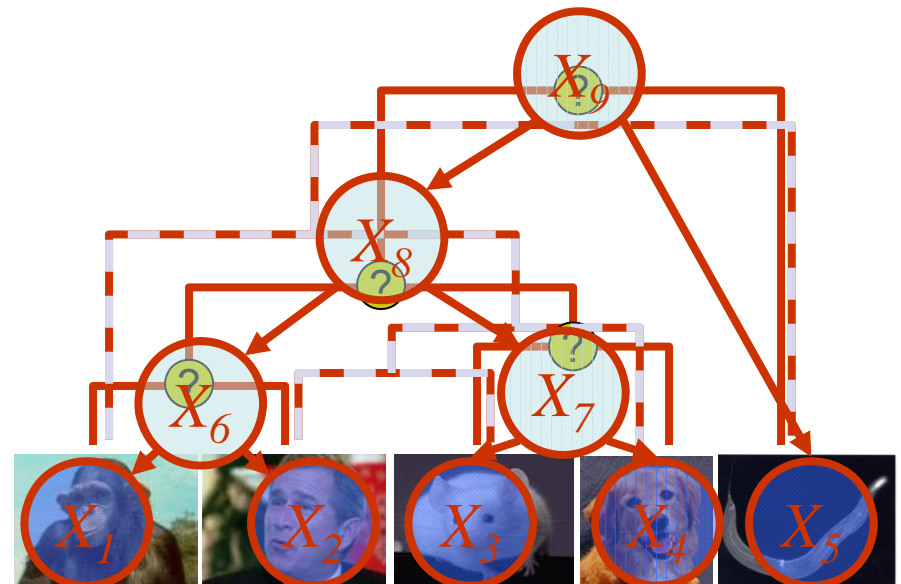Planning

# 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?

  $$\text{e.g.:} \quad P(X_i \mid \mathbf{D})$$

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

  $$\text{e.g.:} \quad \mathcal{M} = \arg\max_{\mathcal{M} \in M} F(\mathbf{D}; \mathcal{M})$$
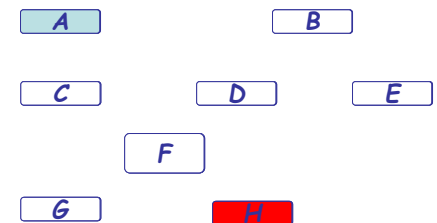
# Recap of Basic Prob. Concepts

- 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?**

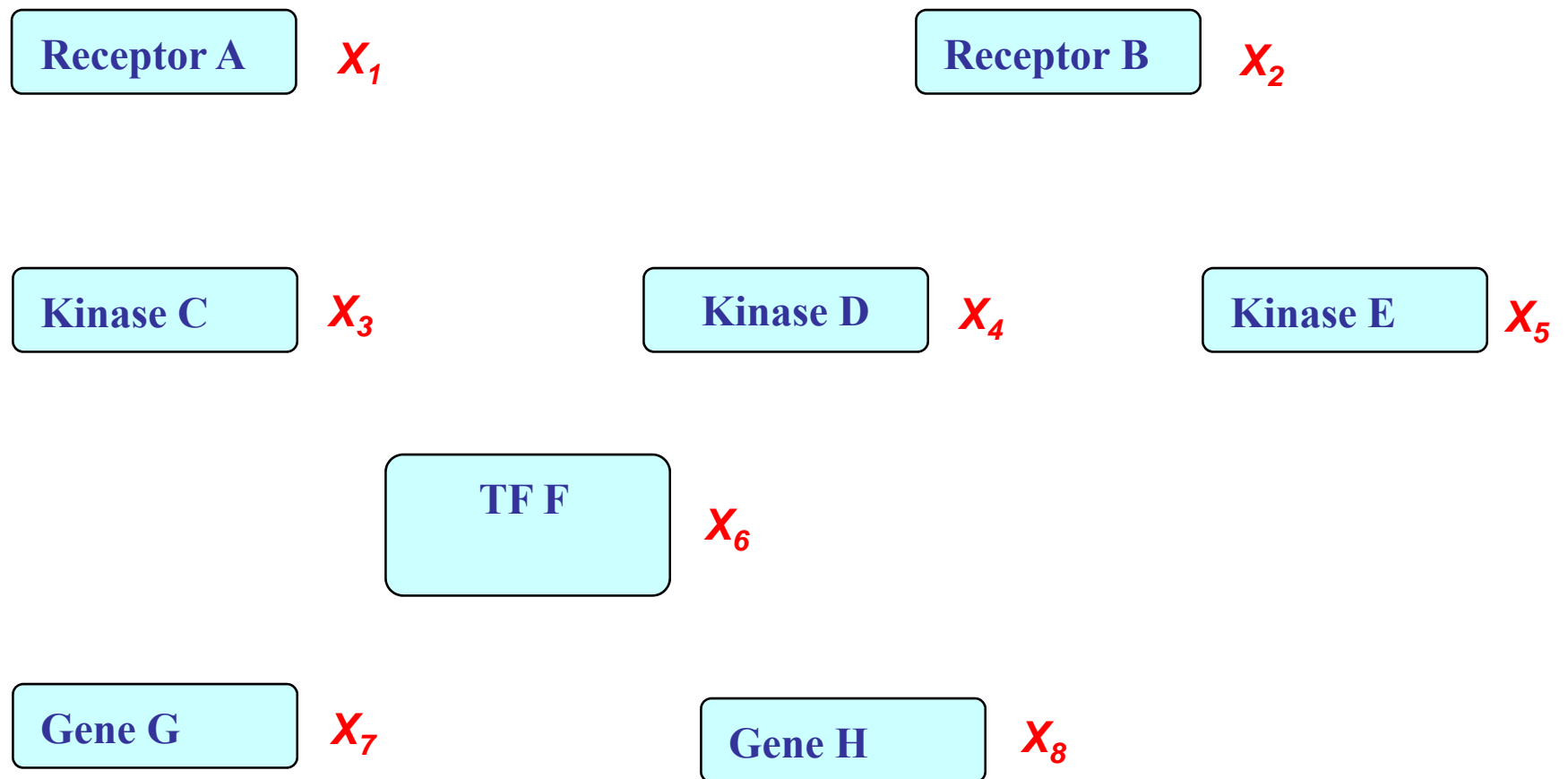| A |  | B |
|---|---|---|
| C | D | E |
|   | F |   |
| G | H |   |

- 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?
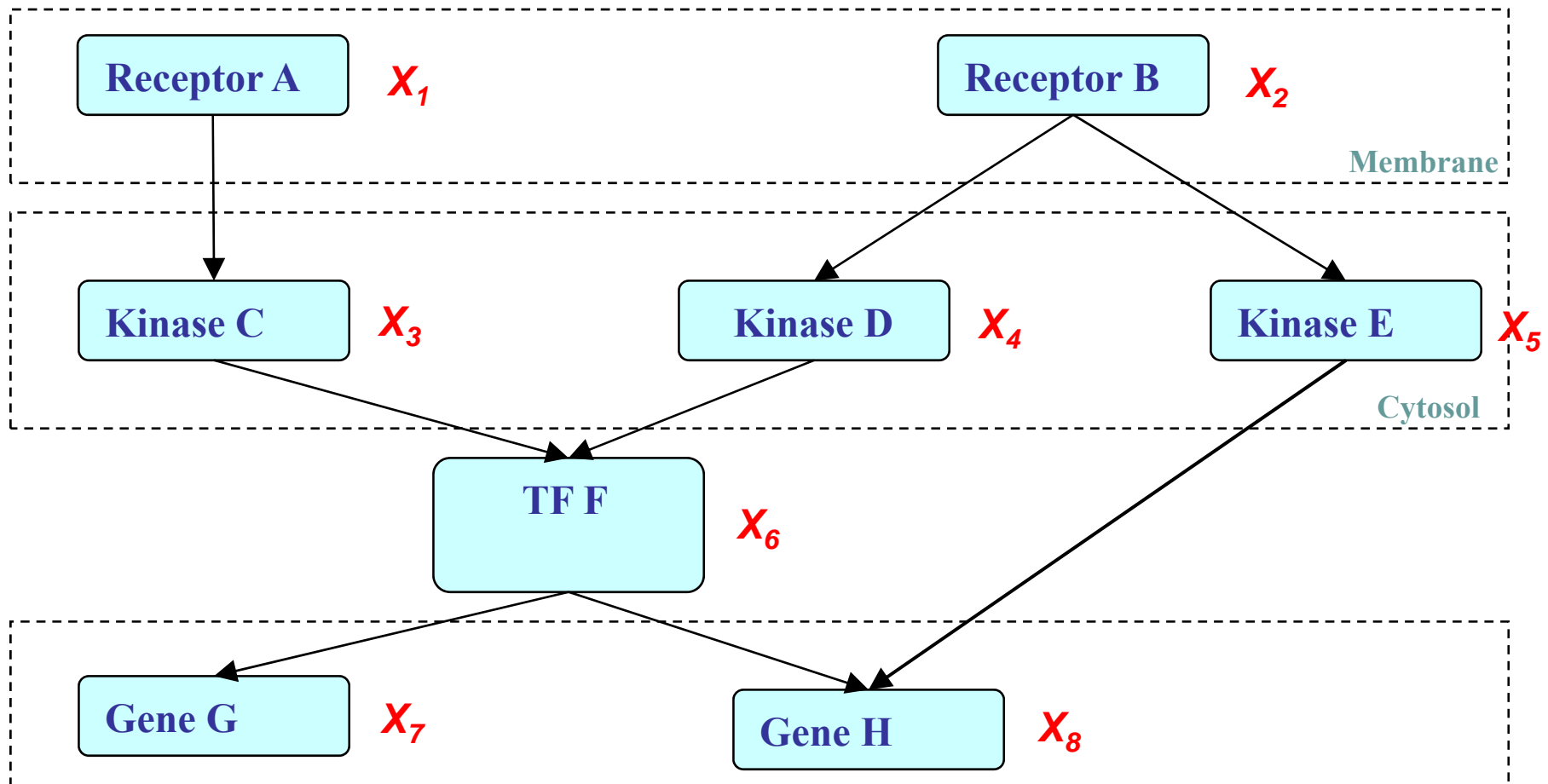## --- Multivariate Distribution in High-D Space

- A possible world for cellular signal transduction:

Receptor A $X_1$

Receptor B $X_2$

Kinase C $X_3$

Kinase D $X_4$

Kinase E $X_5$

TF F $X_6$

Gene G $X_7$

Gene H $X_8$

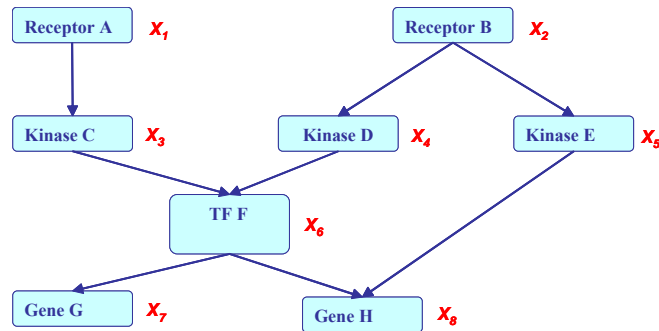# GM: Structure Simplifies Representation

- 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.,



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

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

**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 + Transcritome + 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.,
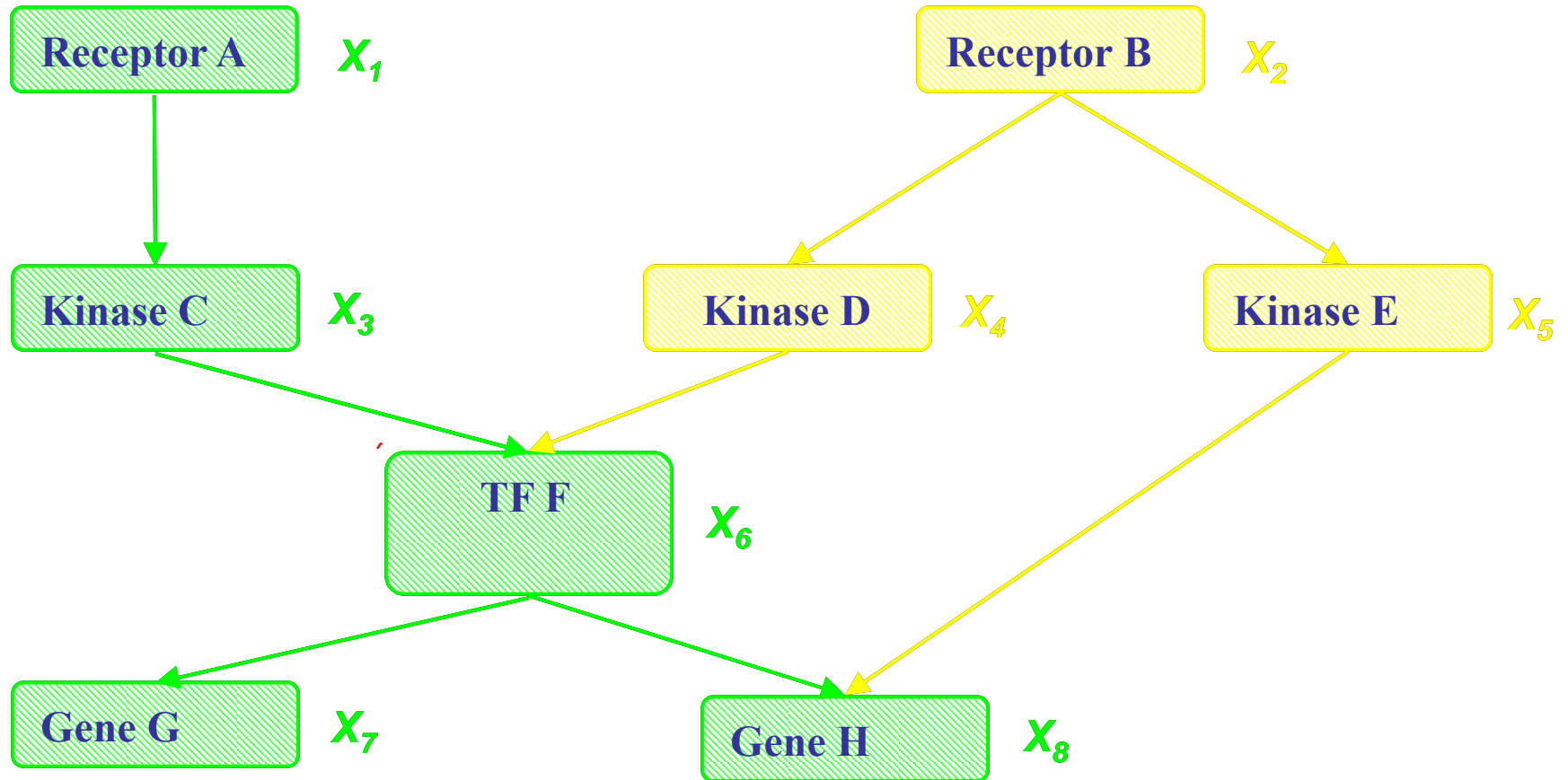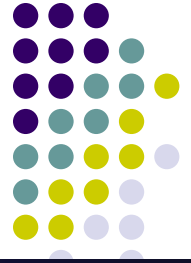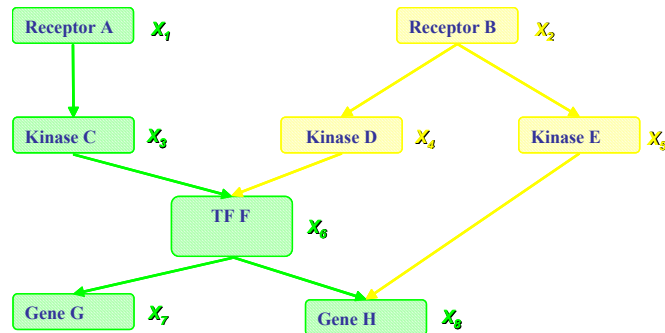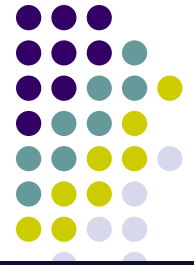


$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$
$$= P(X_2) P(X_4/X_2) P(X_5/X_2) P(X_1) P(X_3/X_1)$$
$$P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)$$

❑ 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
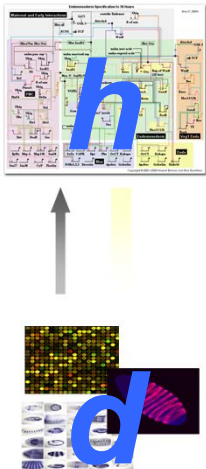
# Rational Statistical Inference

**The Bayes Theorem:**
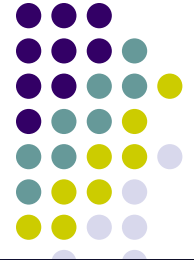
Posterior probability

Likelihood

Prior probability

$$p(h \mid d) = \frac{p(d \mid h)\, p(h)}{\sum_{h' \in H} p(d \mid h')\, p(h')}$$
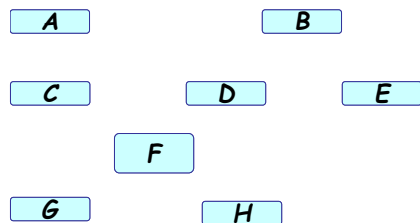
Sum over space of hypotheses

*h*

*d*

- 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 $\Theta$ is conditioned on the values of the observed variables $\mathbf{A}_{obs}$ and prior $p(\,|\chi)$
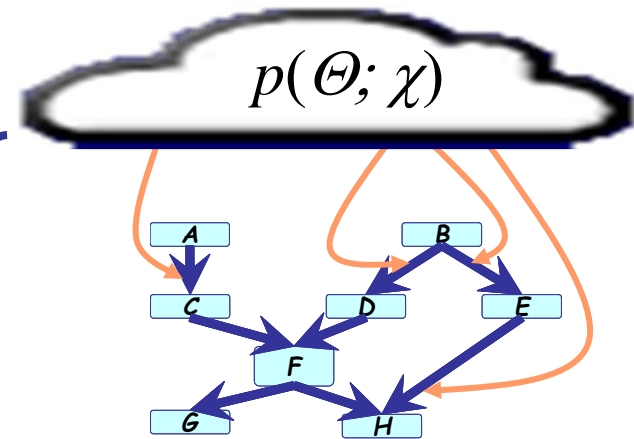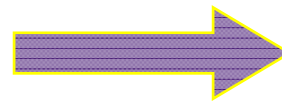
$$p(\Theta;\chi)$$

$$(A,B,C,D,E,\ldots)=(T,F,F,T,F,\ldots)$$
$$\mathbf{A}=(A,B,C,D,E,\ldots)=(T,F,T,T,F,\ldots)$$
$$\ldots\ldots$$
$$(A,B,C,D,E,\ldots)=(F,T,T,T,F,\ldots)$$

| C | D | P(F | C,D) | |
|---|---|---|---|
| c | d | 0.9 | 0.1 |
| c | d̄ | 0.2 | 0.8 |
| c̄ | d | 0.9 | 0.1 |
| c̄ | d̄ | 0.01 | 0.99 |

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

$$p(\Theta\mid\mathbf{A};\chi)\propto p(\mathbf{A}\mid\Theta)\,p(\Theta;\chi)$$

posterior     likelihood     prior

# 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.,



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

❑ 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

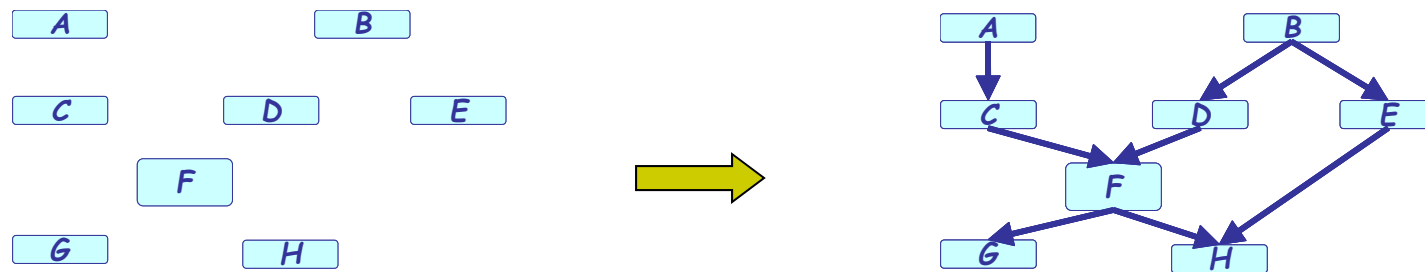# So What Is a PGM After All?

In a nutshell:

PGM   =   Multivariate Statistics + Structure

GM   =   Multivariate Obj. Func. + Structure

# So What Is a PGM After All?

- 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*



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

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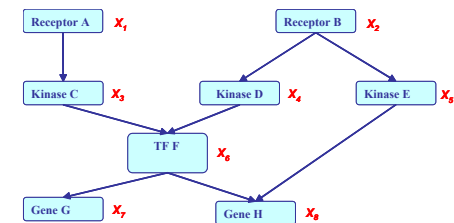$$P(X_6 \mid X_3, X_4)P(X_7 \mid X_6)P(X_8 \mid X_5, X_6)$$

- 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
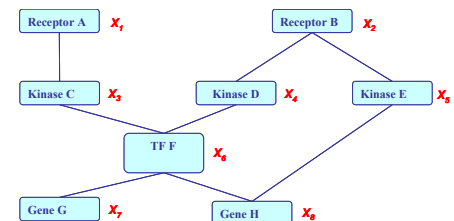
# Two types of GMs

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

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

$$= P(X_1)\, P(X_2)\, P(X_3/X_1)\, P(X_4/X_2)\, P(X_5/X_2)$$
$$P(X_6/X_3, X_4)\, P(X_7/X_6)\, P(X_8/X_5, X_6)$$
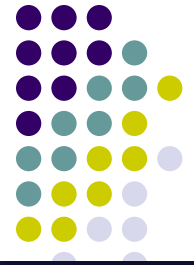
- Undirected edges simply give correlations between variables (**Markov Random Field** or **Undirected Graphical model**):

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

$$= 1/Z \, \exp\{E(X_1)+E(X_2)+E(X_3, X_1)+E(X_4, X_2)+E(X_5, X_2)$$
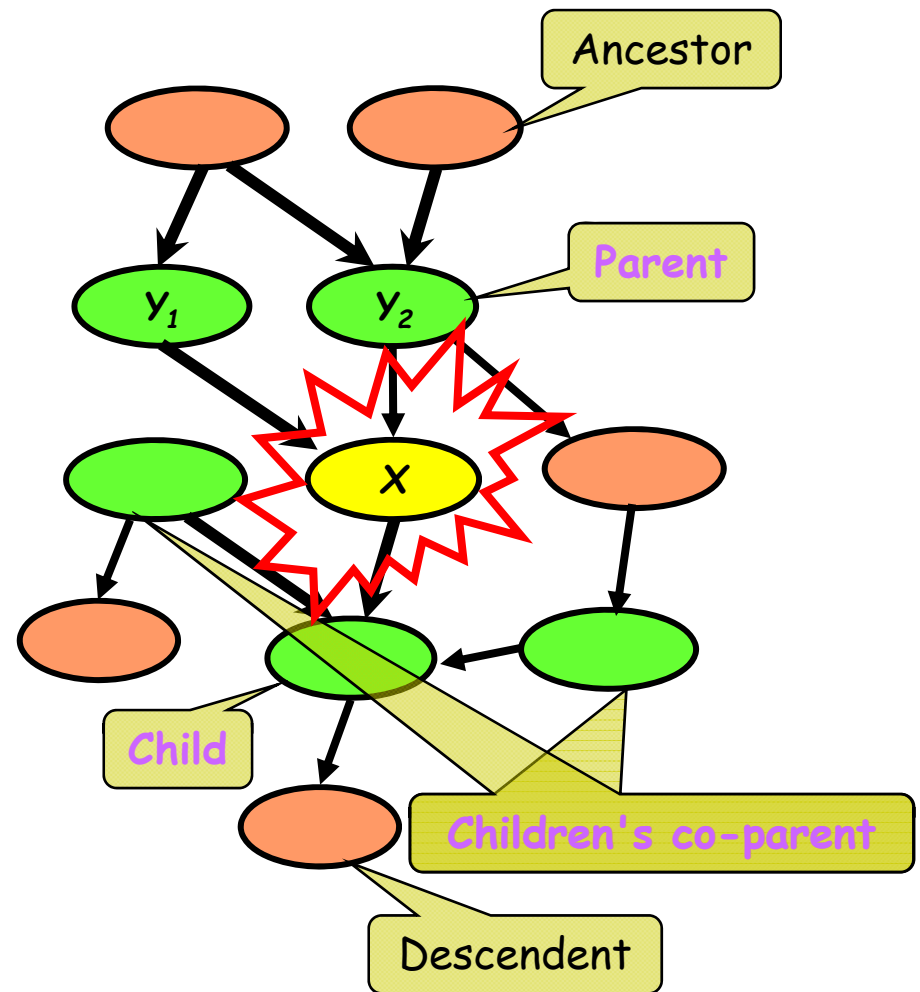$$+ E(X_6, X_3, X_4)+E(X_7, X_6)+E(X_8, X_5, X_6)\}$$

# 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.

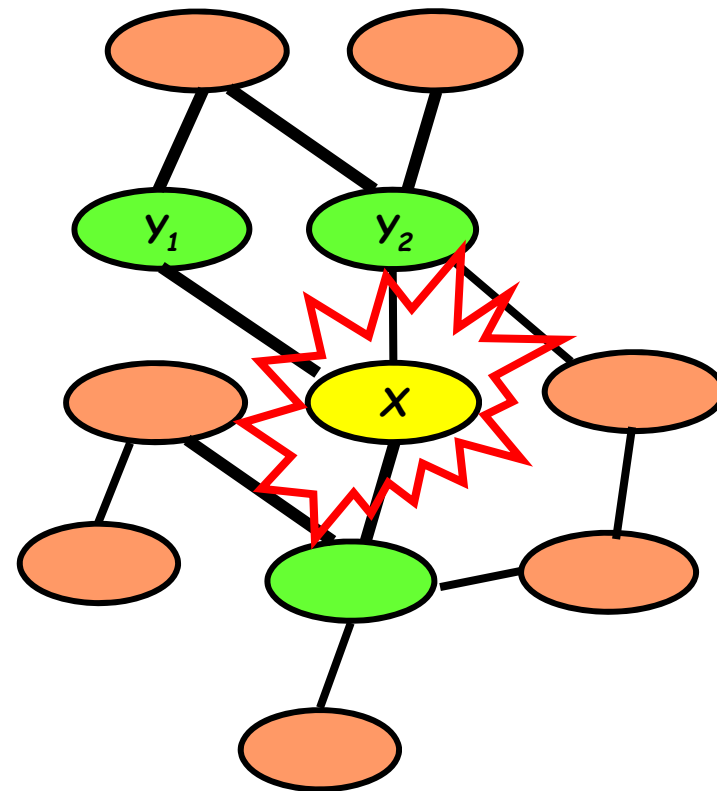- Give **causality** relationships, and facilitate a **generative** process



Ancestor

Parent

$Y_1$  $Y_2$

X

Child

Children's co-parent

Descendent

# 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),

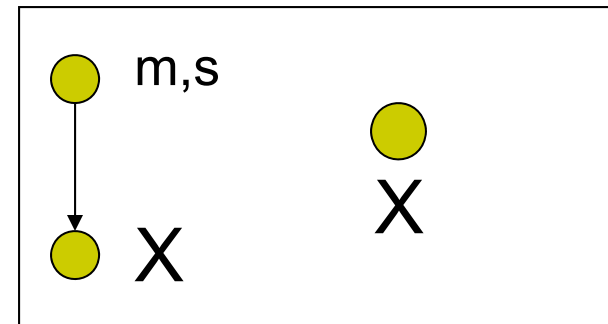  Let $\mathcal{D}_2$ denote the family of all distributions that factor according to G,

  Then $\mathcal{D}_1 \equiv \mathcal{D}_2$.
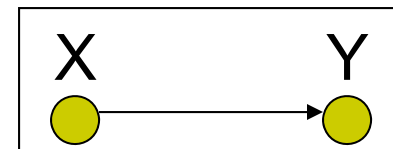
# GMs are your old friends

## Density estimation

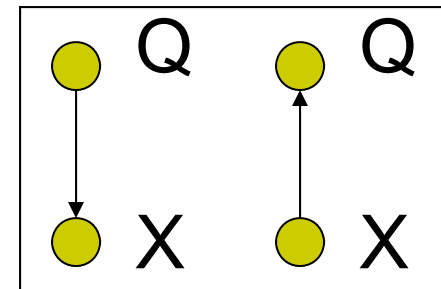Parametric and nonparametric methods

## Regression

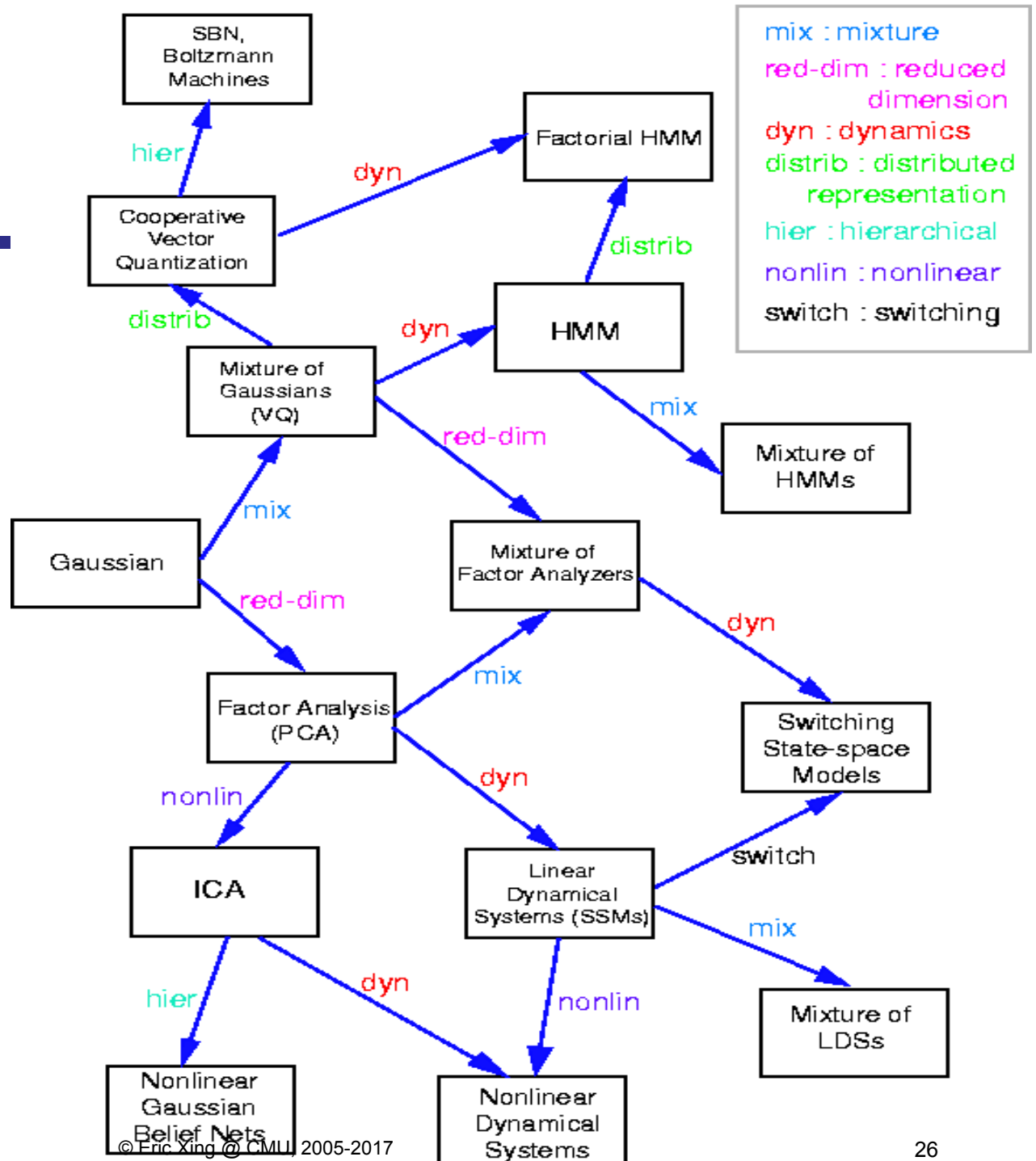Linear, conditional mixture, nonparametric

## Classification
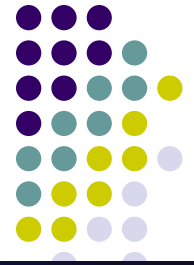
Generative and discriminative approach
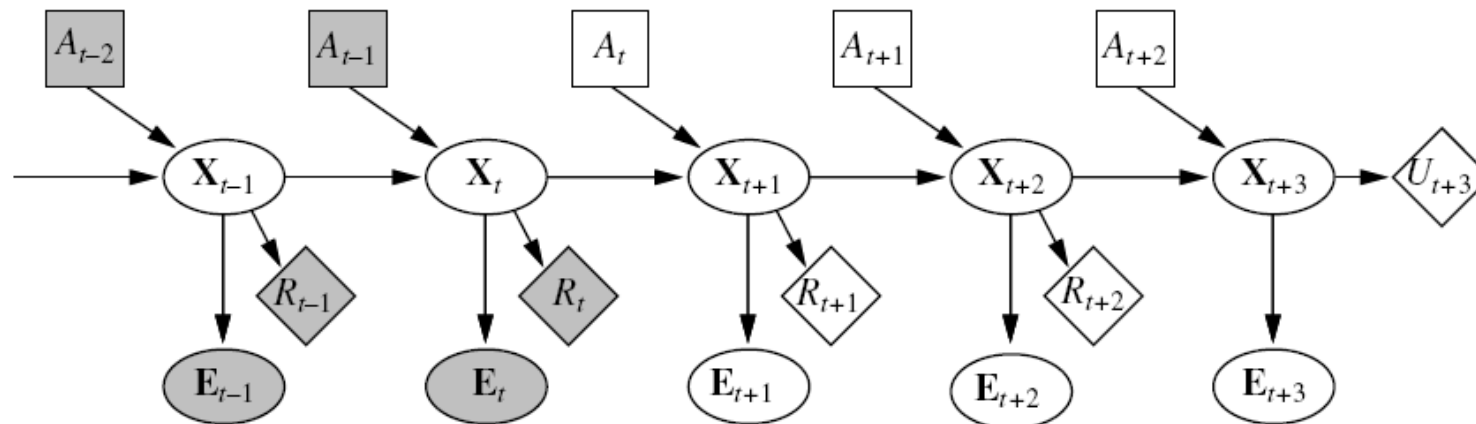
## Clustering

# An (incomplete) genealogy of graphical models



(Picture by Zoubin Ghahramani and Sam Roweis)

26

# Fancier GMs: reinforcement learning

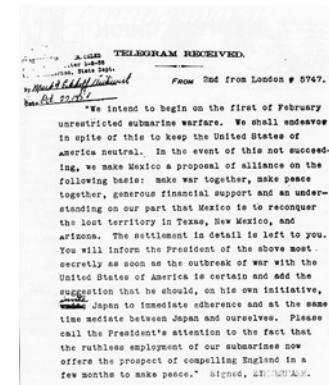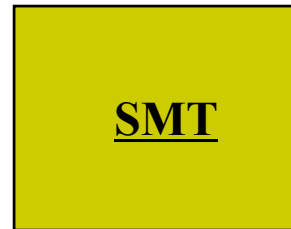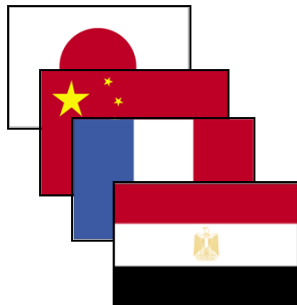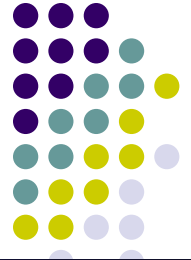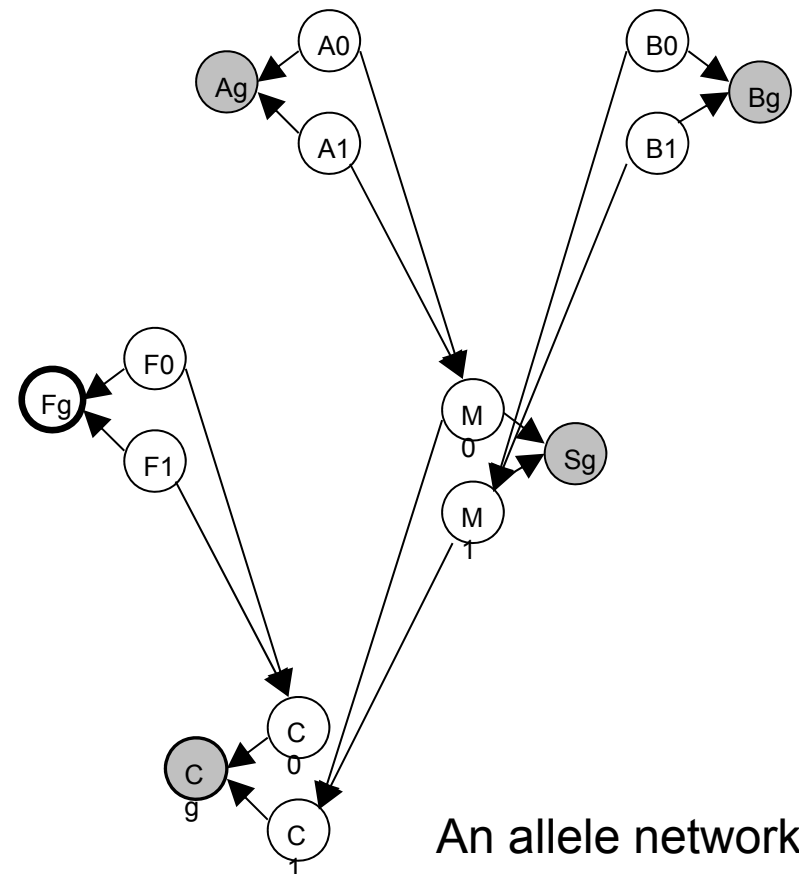- Partially observed Markov decision processes (POMDP)

# Fancier GMs: machine translation



SMT

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

# Fancier GMs: genetic pedigree



An allele network

# Fancier GMs:
# solid state physics



Ising/Potts model

# Deep Neural Networks



input
1@32x32

**Layer 1**
**6@28x28**

**Layer 2**
**6@14x14**

**Layer 3**
**12@10x10**

**Layer 4**
**12@5x5**

**Layer 5**
**100@1x1**

**Layer 6: 10**

10

**5x5
convolution**

**2x2
pooling/
subsampling**

**5x5
convolution**

**2x2
pooling/
subsampling**

**5x5
convolution**



Layer 3
Layer 2
Layer 1



High-level
linguistic representations



$O_i^0$  $w_{ij}^1$  $O_j^1$  $w_{jk}^2$  $O_k^2$

$x_1^p$  $O_1^0$  $w_{11}^1$  $w_{21}^1$  $O_1^1$  $w_{11}^2$

$w_{13}^1$  $w_{22}^1$  $O_2^1$  $w_{21}^2$

$w_{23}^1$  $O_3^1$  $w_{31}^2$  $O_1^2$  $t_1^p$

$x_2^p$  $w_{15}^1$  $w_{24}^1$  $O_4^1$  $w_{41}^2$

$O_2^0$  $w_{25}^1$  $O_5^1$  $w_{51}^2$

Layer 0
(Input)

Layer 1
(Hidden)

Layer 2
(Output)

# What makes it work? Why?

# An MLer's View of the World



Loss functions
(likelihood, reconstruction, margin, …)

Structures
(Graphical, group, chain, tree, iid, …)

Constraints
(normality, sparsity, label, prior, KL, sum, …)

Algorithms
MC (MCMC, Importance), Opt (gradient, IP), …

Stopping criteria
Change in objective, change in update …

Empirical Performances?

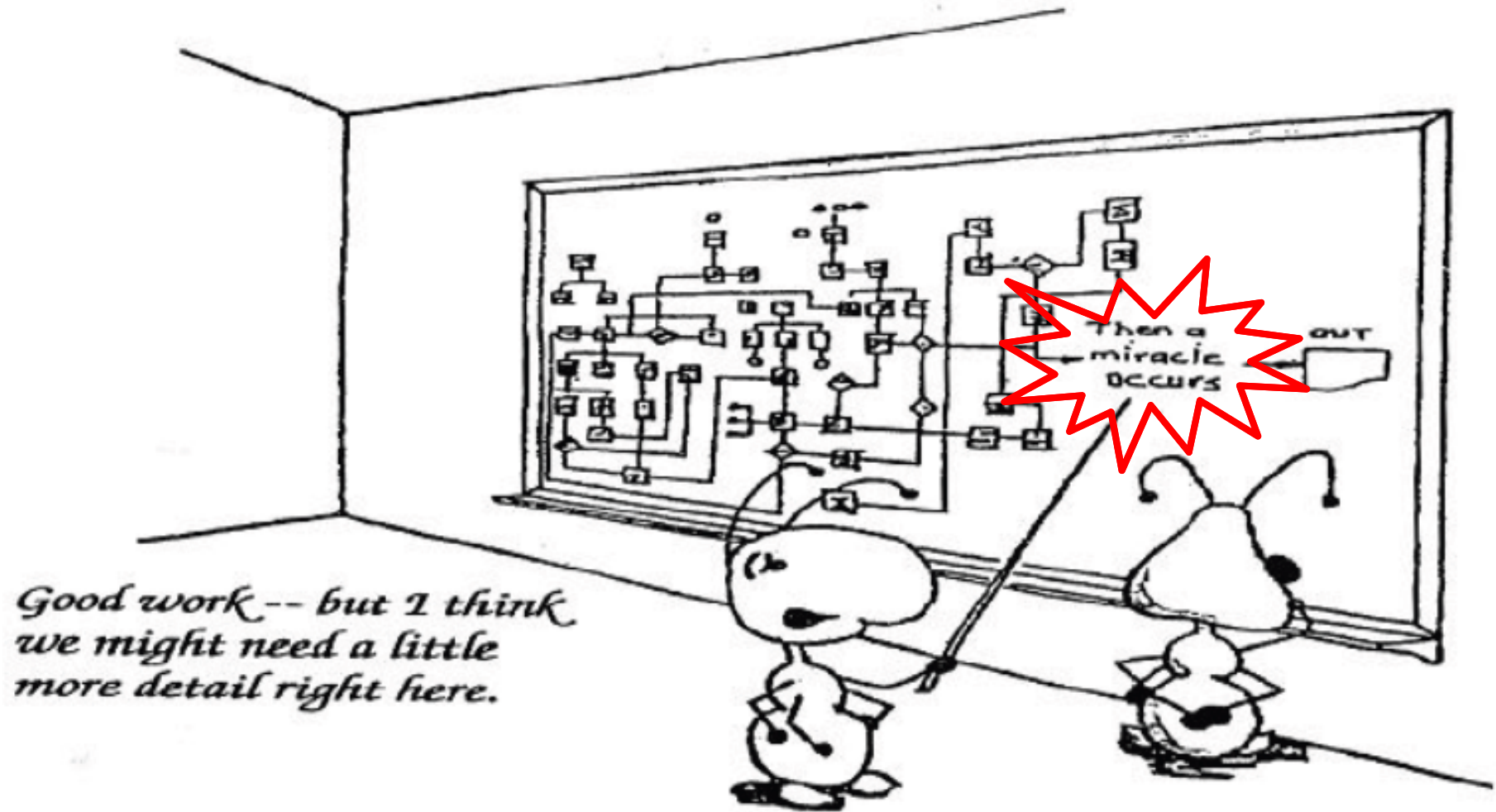| | DL | ⋛ ? ML (e.g., GM) |
|---|---|---|
| Empirical goal: | e.g., classification, feature learning, generating samples … | e.g., supervised/unsupervised learning, transfer learning, latent variable inference |
| Structure: | Graphical | Graphical |
| Objective: | Something aggregated from local functions | Something aggregated from local functions |
| | | |
| Vocabulary: | Neuron, activation/gate function … | Variables, potential function |
| Algorithm: | A single, unchallenged, inference algorithm – BP | A major focus of open research, many algorithms, and more to come |
| Evaluation: | On a black-box score -- end performance | On almost every intermediate quantity |
| Implementation: | Many untold-tricks | More or less standardized |
| Experiments: | Massive, real data (GT unknown) | Modest, often simulated data (GT known) |

# 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**.

**--- M. Jordan**

# Plan for the Class

- Fundamentals of Graphical Models:
    - Bayesian Network and Markov Random Fields
    - Discrete, Continuous and Hybrid models, exponential family, GLIM
    - Basic representation, inference, and learning
    - …

- Advanced topics and latest developments
    - Approximate inference
        - Monte Carlo algorithms
        - Vatiational methods and theories
    - "Infinite" GMs: nonparametric Bayesian models
    - Optimization-theoretic formulations for GMs, e.g., Structured sparsity
    - Nonparametric and spectral graphical models, where GM meets kernels and matrix algebra
    - Alternative GM learning paradigms,
        - e.g., Margin-based learning of GMs (where GM meets SVM)
        - e.g., Regularized Bayes: where GM meets SVM, and meets Bayesian, and meets NB …

- Case studies: popular GMs and applications
    - Multivariate Gaussian Models
    - Conditional random fields
    - Mixed-membership, aka, Topic models

# Questions ?