# ECE 448: Artificial Intelligence

## Lecture 16: Bayesian Networks

Prof. Hongwei Wang hongweiwang@intl.zju.edu.cn

Prof. Mark Hasegawa-Johnson jhasegaw@illinois.edu

Spring 2023

- A general scenario:
  - Query *variables:* **X**
  - *Evidence* (*observed*) variables and their values: **E = e**
- **Inference problem**: answer questions about the query variables given the evidence variables
- This can be done using the posterior distribution P(**X | E = e**)
- Example of a useful question: **Which X is true?**
- More formally: what value of **X** has the least probability of being wrong?
- Answer: **MPE = MAP** (argmin P(error) = argmax P(X=x|E=e))

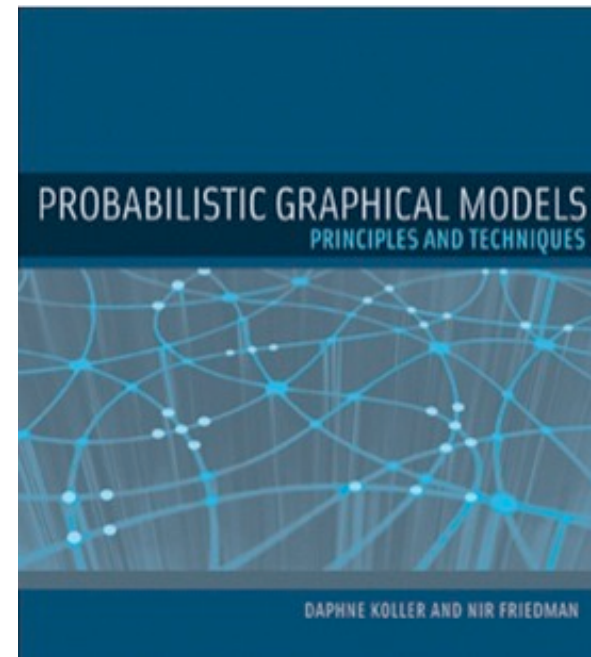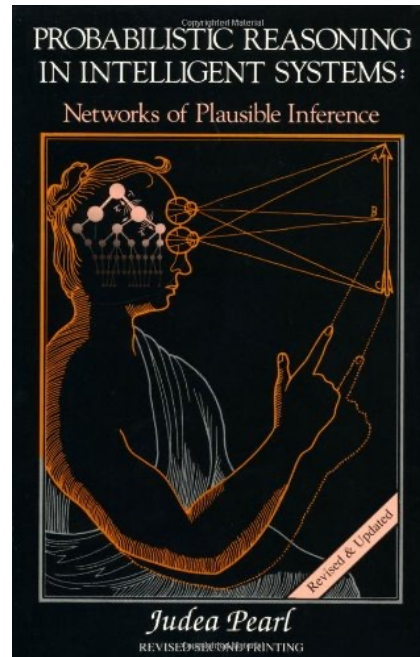- Very, very common problem: P(X,E) is complicated because both X and E depend on some hidden variable Y

- SOLUTION:
  - Draw a bunch of circles and arrows that represent the dependence
  - When your algorithm performs inference, make sure it does so in the order of the graph

- FORMALISM: Bayesian Network

- A general scenario:
  - Query *variables:* **X**
  - *Evidence* (*observed*) variables and their values: **E** = **e**
  - *Unobserved* variables: **Y**
- **Inference problem**: answer questions about the query variables given the evidence variables
  - This can be done using the posterior distribution P(**X** | **E** = **e**)
  - In turn, the posterior needs to be derived from the full joint P(**X**, **E**, **Y**)

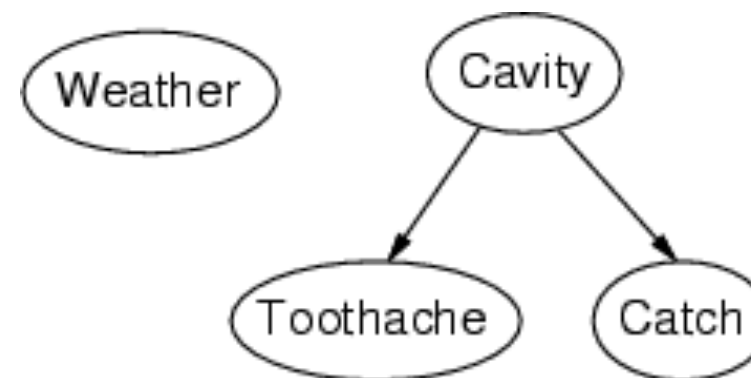$$P(X \mid E = e) = \frac{P(X, e)}{P(e)} \propto \sum_y P(X, e, y)$$

- Bayesian networks are a tool for representing joint probability distributions efficiently

- More commonly called *graphical models*
- A way to depict conditional independence relationships between random variables
- A compact specification of full joint distributions

ZJUI

1.   Review: Bayesian inference

2.   Bayesian network: graph semantics

3.   The Los Angeles burglar alarm example

4.   Conditional independence ≠ Independence

5.   Constructing a Bayesian network: Structure learning

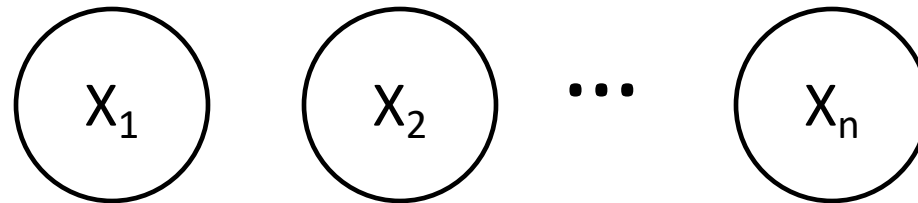6.   Constructing a Bayesian network: Hire an expert

1. **Review: Bayesian inference**

2. **Bayesian network: graph semantics**

3. **The Los Angeles burglar alarm example**

4. **Conditional independence ≠ Independence**

5. **Constructing a Bayesian network: Structure learning**

6. **Constructing a Bayesian network: Hire an expert**
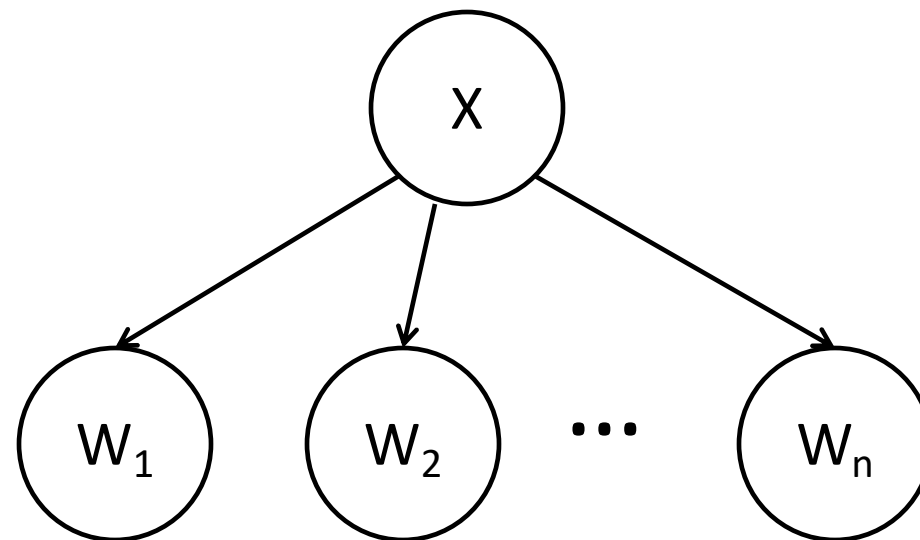
ZJUI



- **Nodes:** random variables

- **Arcs:** interactions
  - An arrow from one variable to another indicates direct influence
  - Must form a directed, *acyclic* graph
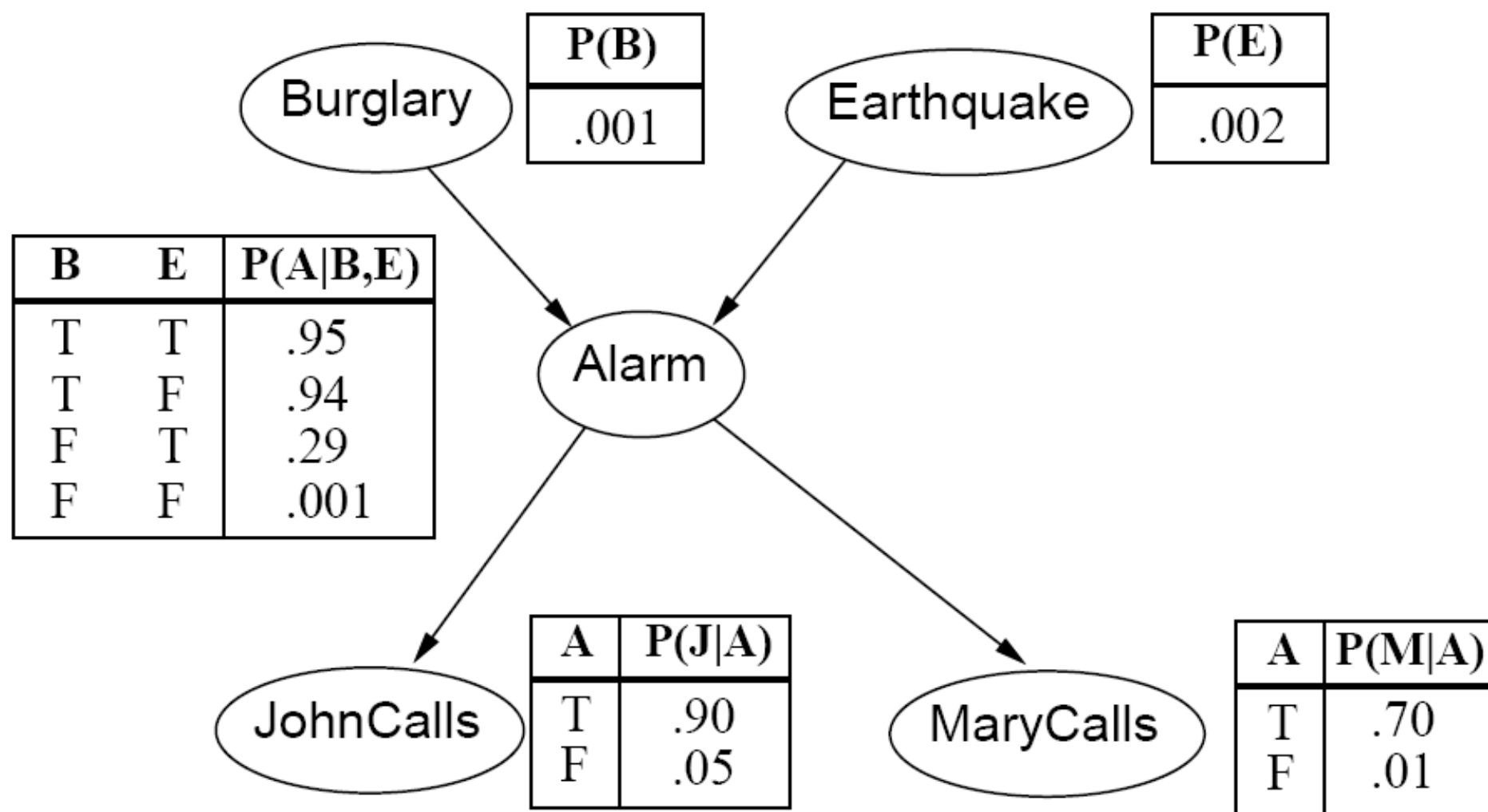
ZJUI

- Complete independence: no interactions

$$X_1 \quad X_2 \quad \cdots \quad X_n$$

- Random variables:
  - X: document class
  - $W_1, ..., W_n$: words in the document

1. **Review: Bayesian inference**

2. **Bayesian network: graph semantics**

3. **The Los Angeles burglar alarm example**

4. **Conditional independence ≠ Independence**

5. **Constructing a Bayesian network: Structure learning**

6. **Constructing a Bayesian network: Hire an expert**

- I have a burglar alarm that is sometimes set off by minor earthquakes. My two neighbors, John and Mary, promised to call me at work if they hear the alarm
  - Example inference task: suppose Mary calls and John doesn't call. What is the probability of a burglary?

- What are the random variables?
  - *Burglary, Earthquake, Alarm, JohnCalls, MaryCalls*

- What are the direct influence relationships?
  - A burglar can set the alarm off
  - An earthquake can set the alarm off
  - The alarm can cause Mary to call
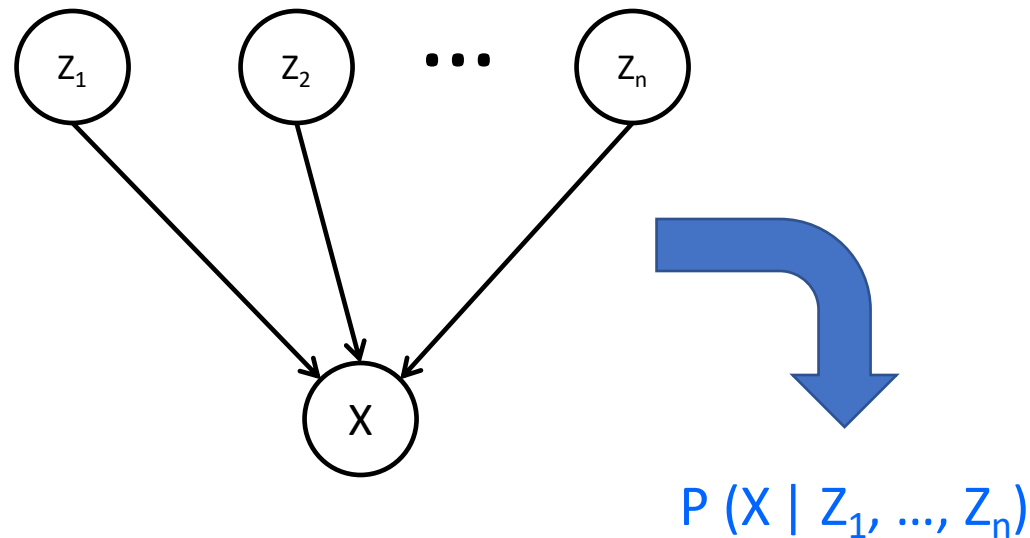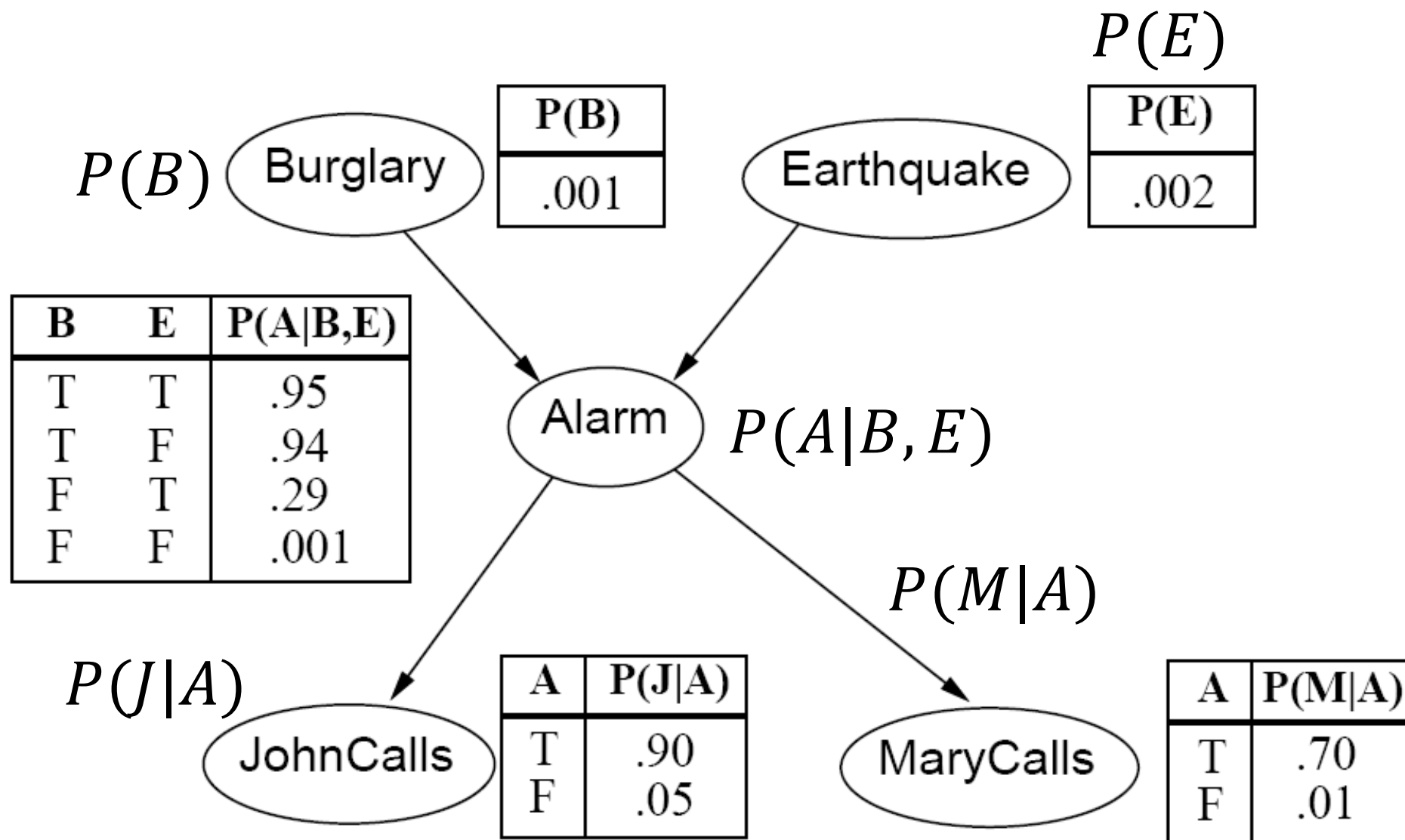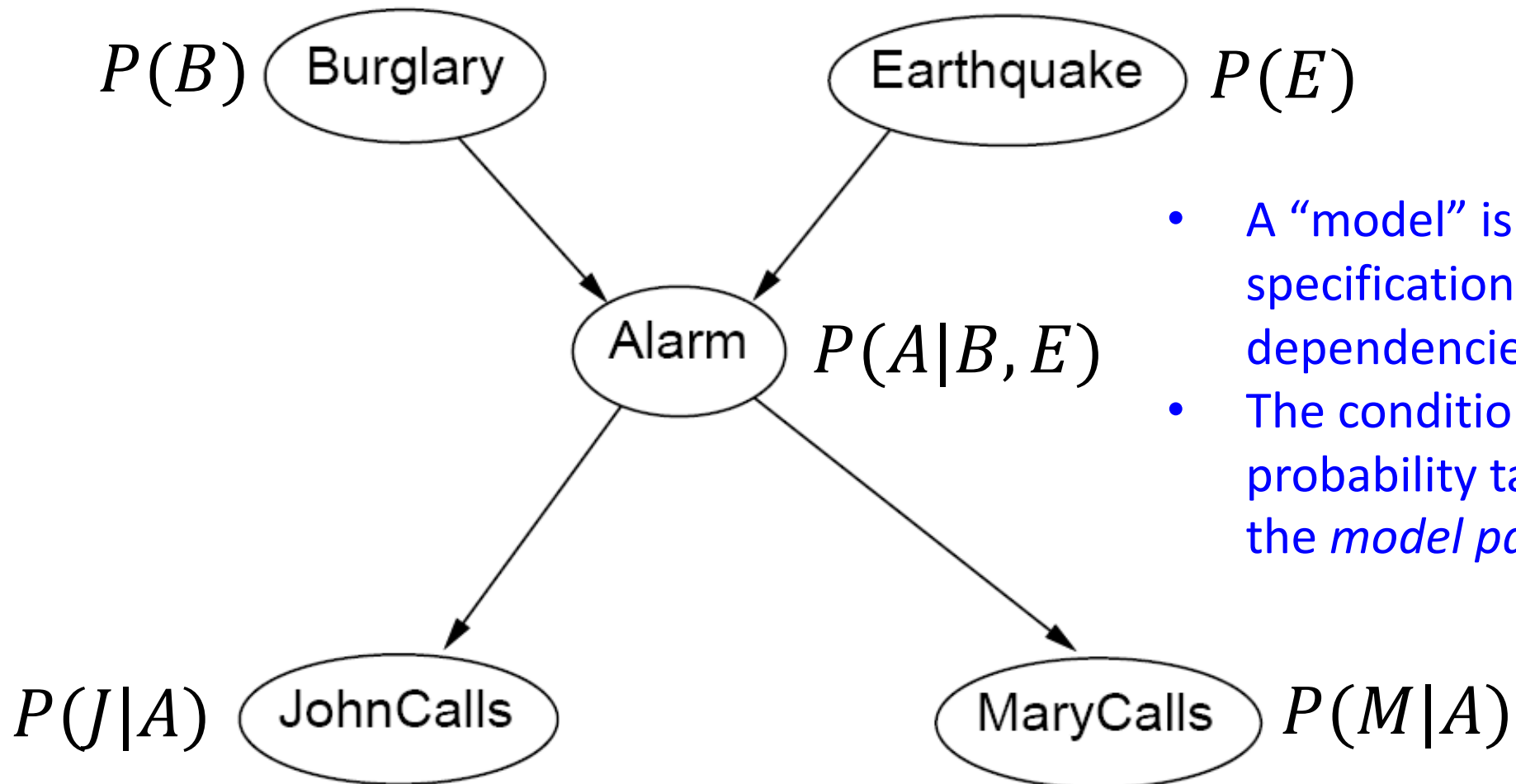  - The alarm can cause John to call

- Key property: each node is conditionally independent of its *non-descendants* given its *parents*

- Suppose the nodes $X_1, ..., X_n$ are sorted in topological order

- To get the joint distribution $P(X_1, ..., X_n)$, use chain rule:

$$P(X_1, ..., X_n) = \prod_{i=1}^{n} P(X_i \mid X_1, ..., X_{i-1})$$

$$= \prod_{i=1}^{n} P(X_i \mid Parents(X_i))$$

- To specify the full joint distribution, we need to specify a *conditional* distribution for each node given its parents:
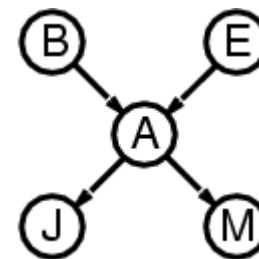  P (X | Parents(X))



$P(X \mid Z_1, \ldots, Z_n)$

$P(E)$

| | P(E) |
|---|---|
| | .002 |

| | P(B) |
|---|---|
| | .001 |

$P(B)$   Burglary

Earthquake

| B | E | P(A\|B,E) |
|---|---|---|
| T | T | .95 |
| T | F | .94 |
| F | T | .29 |
| F | F | .001 |

Alarm   $P(A|B,E)$

$P(M|A)$

$P(J|A)$

JohnCalls

| A | P(J\|A) |
|---|---|
| T | .90 |
| F | .05 |

MaryCalls

| A | P(M\|A) |
|---|---|
| T | .70 |
| F | .01 |

$P(B)$ Burglary          Earthquake $P(E)$

Alarm $P(A|B,E)$

- A "model" is a complete specification of the dependencies.
- The conditional probability tables are the *model parameters.*

$P(J|A)$ JohnCalls          MaryCalls $P(M|A)$

ZJUI

$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P(X_i \mid Parents(X_i))$$



For example,

P(j, m, a, ¬b, ¬e) = P(¬b) P(¬e) P(a|¬b,¬e) P(j|a) P(m|a)

- By saying that $X_i$ and $X_j$ are independent, we mean that
$$\mathrm{P}(X_j, X_i) = \mathrm{P}(X_i)\mathrm{P}(X_j)$$

- $X_i$ and $X_j$ are independent if and only if they have no common ancestors

- Example: *independent coin flips*
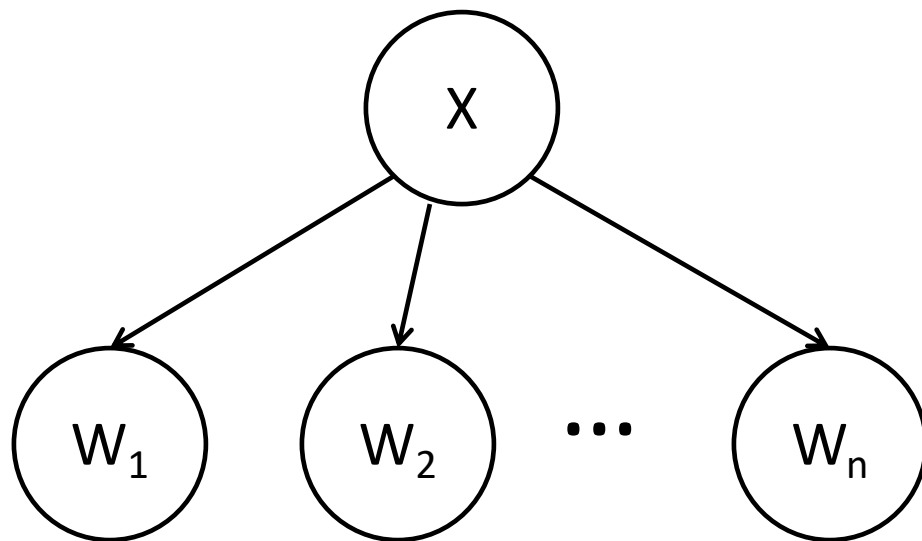


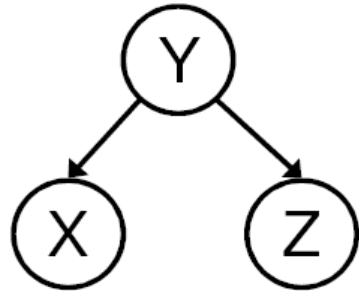- Another example: Weather is independent of all other variables in this model.

- By saying that $W_i$ and $W_j$ are conditionally independent given $X$, we mean that

$$\mathrm{P}(W_i, W_j | X) = \mathrm{P}(W_i | X)\mathrm{P}(W_j | X)$$

- $W_i$ and $W_j$ are conditionally independent given $X$ if and only if they have no common ancestors other than the ancestors of $X$.

- Example: *naïve Bayes model:*

Common cause: Conditionally
Independent



Y: Project due

X: Newsgroup
busy

Z: Lab full

Common effect: Independent



X: Raining
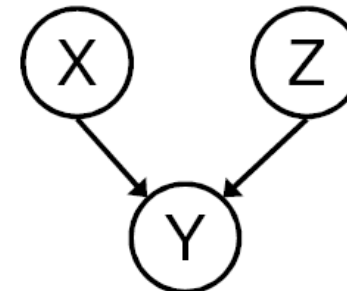
Z: Ballgame

Y: Traffic

Are X and Z independent? **No**

$$P(Z,X) = \sum_Y P(Z|Y)P(X|Y)P(Y)$$

$$P(Z)P(X) = \left(\sum_Y P(Z|Y)P(Y)\right)\left(\sum_Y P(X|Y)P(Y)\right)$$

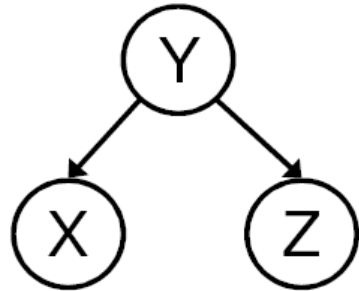Are they conditionally independent given Y? **Yes**

$$P(Z,X|Y) = P(Z|Y)P(X|Y)$$

Are X and Z independent? **Yes**

$$P(X,Z) = P(X)P(Z)$$

Are they conditionally independent given Y? **No**

$$P(Z,X|Y) = \frac{P(Y|X,Z)P(X)P(Z)}{P(Y)}$$

$$\neq P(Z|Y)P(X|Y)$$

ZJUI

## Common cause: Conditionally Independent



Y: Project due

X: Newsgroup busy

Z: Lab full

## Common effect: Independent
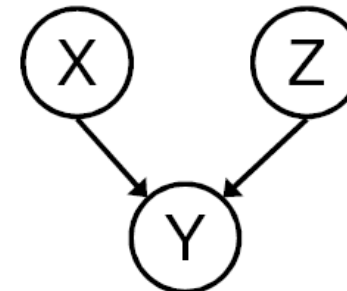


X: Raining

Z: Ballgame

Y: Traffic

Are X and Z independent? **No**

Knowing X tells you about Y, which tells you about Z.

Are they conditionally independent given Y? **Yes**

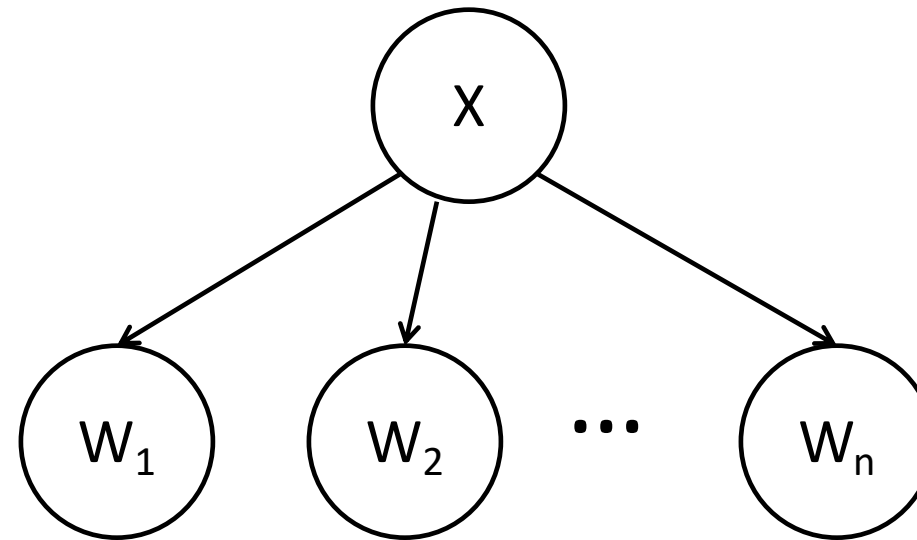If you already know Y, then X gives you no useful information about Z.

Are X and Z independent? **Yes**

**Knowing X tells you nothing about Z.**

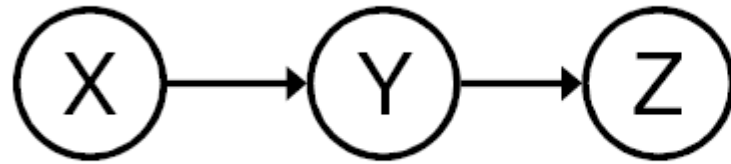Are they conditionally independent given Y? **No**

If Y is true, then either X or Z must be true.
Knowing that X is false means Z must be true.
We say that X "explains away" Z.

Being conditionally independent given X does NOT mean that $W_i$ and $W_j$ are independent.  Quite the opposite.  For example:

- The document topic, X, can be either "sports" or "pets", equally probable.
- $W_1=1$ if the document contains the word "food," otherwise $W_1=0$.
- $W_2=1$ if the document contains the word "dog," otherwise $W_2=0$.
- Suppose you don't know X, but you know that $W_2=1$ (the document has the word "dog").  Does that change your estimate of $p(W_1=1)$?

Another example: *causal chain*



X: Low pressure

Y: Rain

Z: Traffic

- X and Z are conditionally independent given Y, because they have no common ancestors other than the ancestors of Y.

- Being conditionally independent given Y does NOT mean that X and Z are independent. Quite the opposite. For example, suppose $P(X) = 0.5$, $P(Y|X) = 0.8$, $P(Y|\neg X) = 0.1$, $P(Z|Y) = 0.7$, and $P(Z|\neg Y) = 0.4$. Then we can calculate that $P(Z|X) = 0.64$, but $P(Z) = 0.535$

1.   **Review: Bayesian inference**

2.   **Bayesian network: graph semantics**

3.   **The Los Angeles burglar alarm example**

4.   **Conditional independence ≠ Independence**

5.   **Constructing a Bayesian network: Structure learning**

6.   **Constructing a Bayesian network: Hire an expert**

1. "Structure Learning" a.k.a. "Analysis of Causality:"
    1. Suppose you know the variables, but you don't know which variables depend on which others. You can learn this from data.
    2. This is an exciting new area of research in statistics, where it goes by the name of "analysis of causality."
    3. … but it's almost always harder than method #2. You should know how to do this in very simple examples (like the Los Angeles burglar alarm), but you don't need to know how to do this in the general case.

2. "Hire an Expert:"
    1. Find somebody who knows how to solve the problem.
    2. Get her to tell you what are the important variables, and which variables depend on which others.
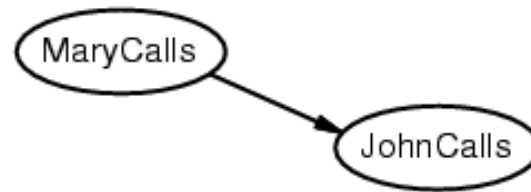    3. THIS IS ALMOST ALWAYS THE BEST WAY.

1. Choose an ordering of variables $X_1, \ldots , X_n$

2. For i = 1 to n
   - add $X_i$ to the network
   - Check your training data. If there is any variable $X_1, \ldots ,X_{i-1}$ that CHANGES the probability of $X_i$=1, then add that variable to the set Parents($X_i$) such that
     $P(X_i \mid Parents(X_i)) = P(X_i \mid X_1, \ldots X_{i-1})$

3. Repeat the above steps for every possible ordering (complexity: n!).

4. Choose the graph that has the smallest number of edges.

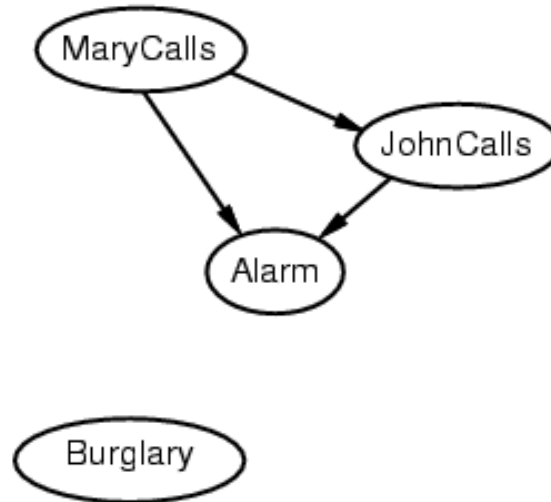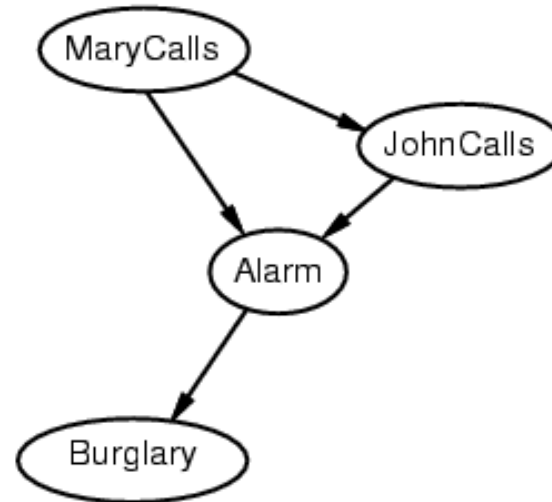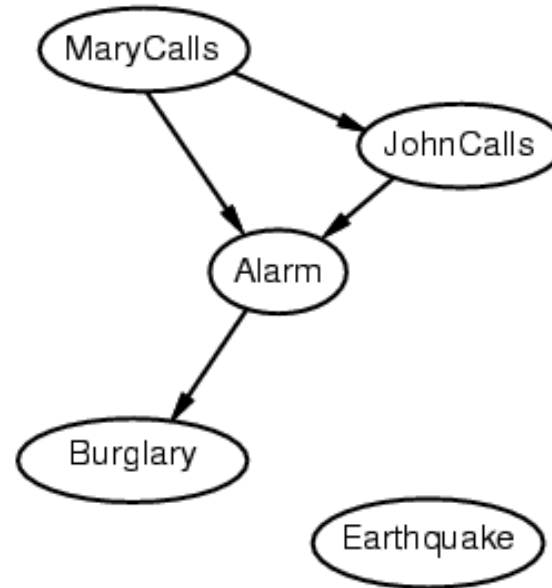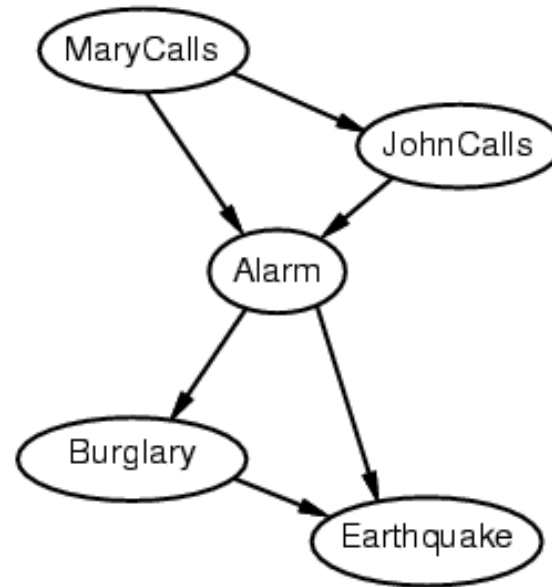- Suppose we choose the ordering M, J, A, B, E

MaryCalls

JohnCalls

- Suppose we choose the ordering M, J, A, B, E

- Suppose we choose the ordering M, J, A, B, E

- Suppose we choose the ordering M, J, A, B, E

ZJUI

- Suppose we choose the ordering M, J, A, B, E
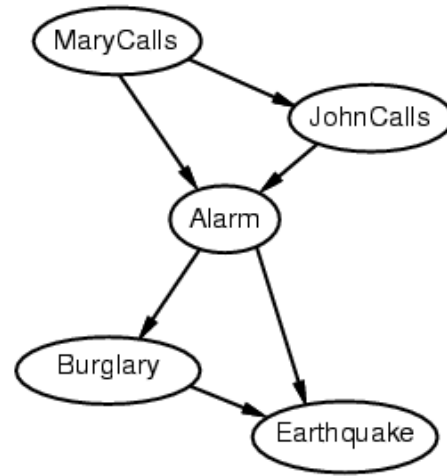
- Suppose we choose the ordering M, J, A, B, E

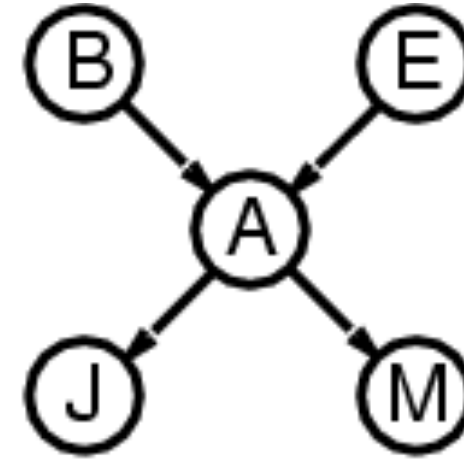- Suppose we choose the ordering M, J, A, B, E

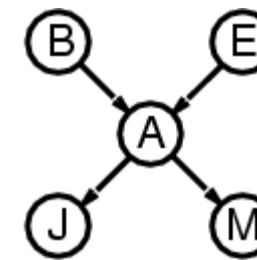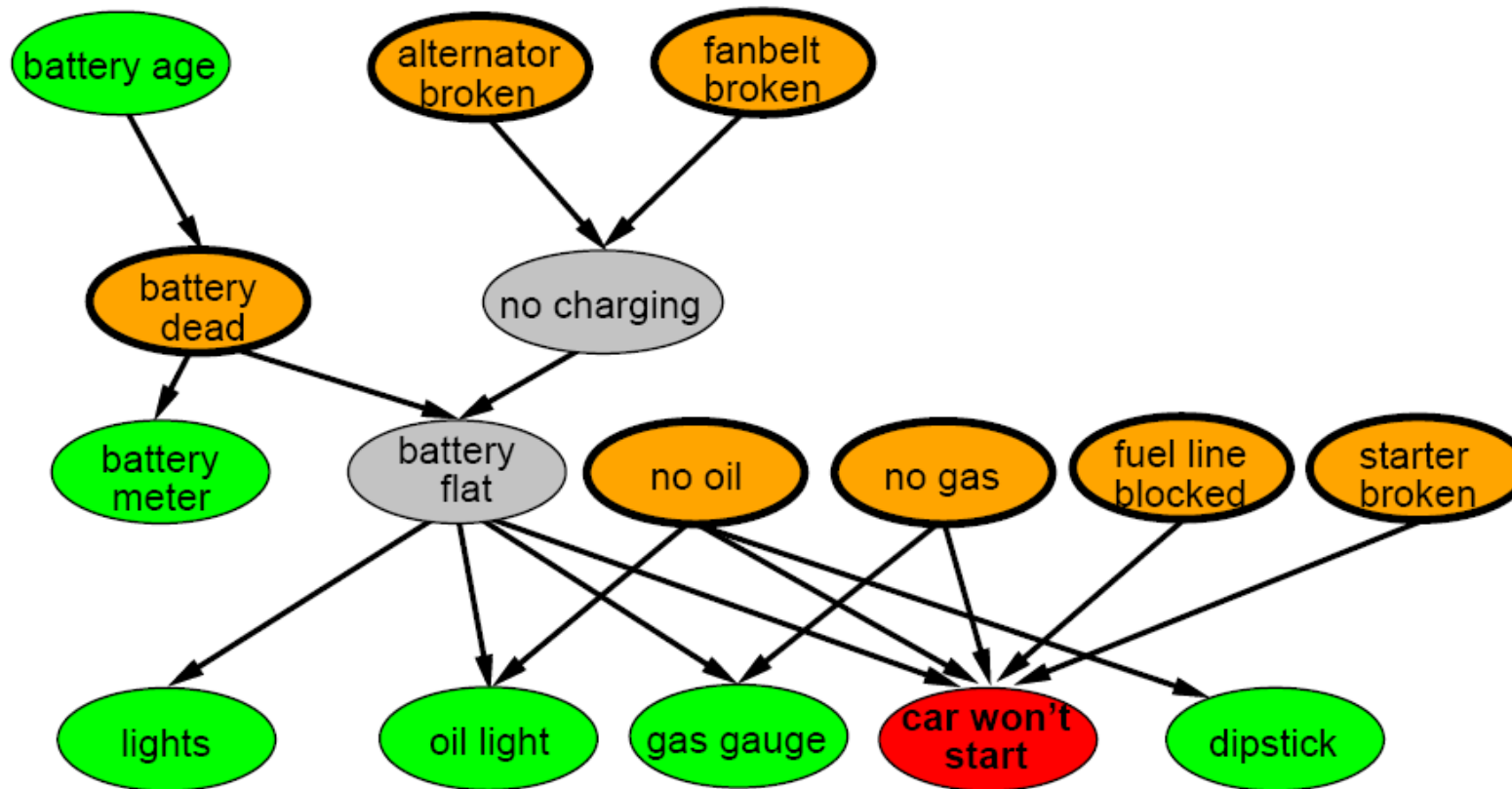- Suppose we choose the ordering M, J, A, B, E

versus

- Deciding conditional independence is hard in noncausal directions
  - The causal direction seems much more natural
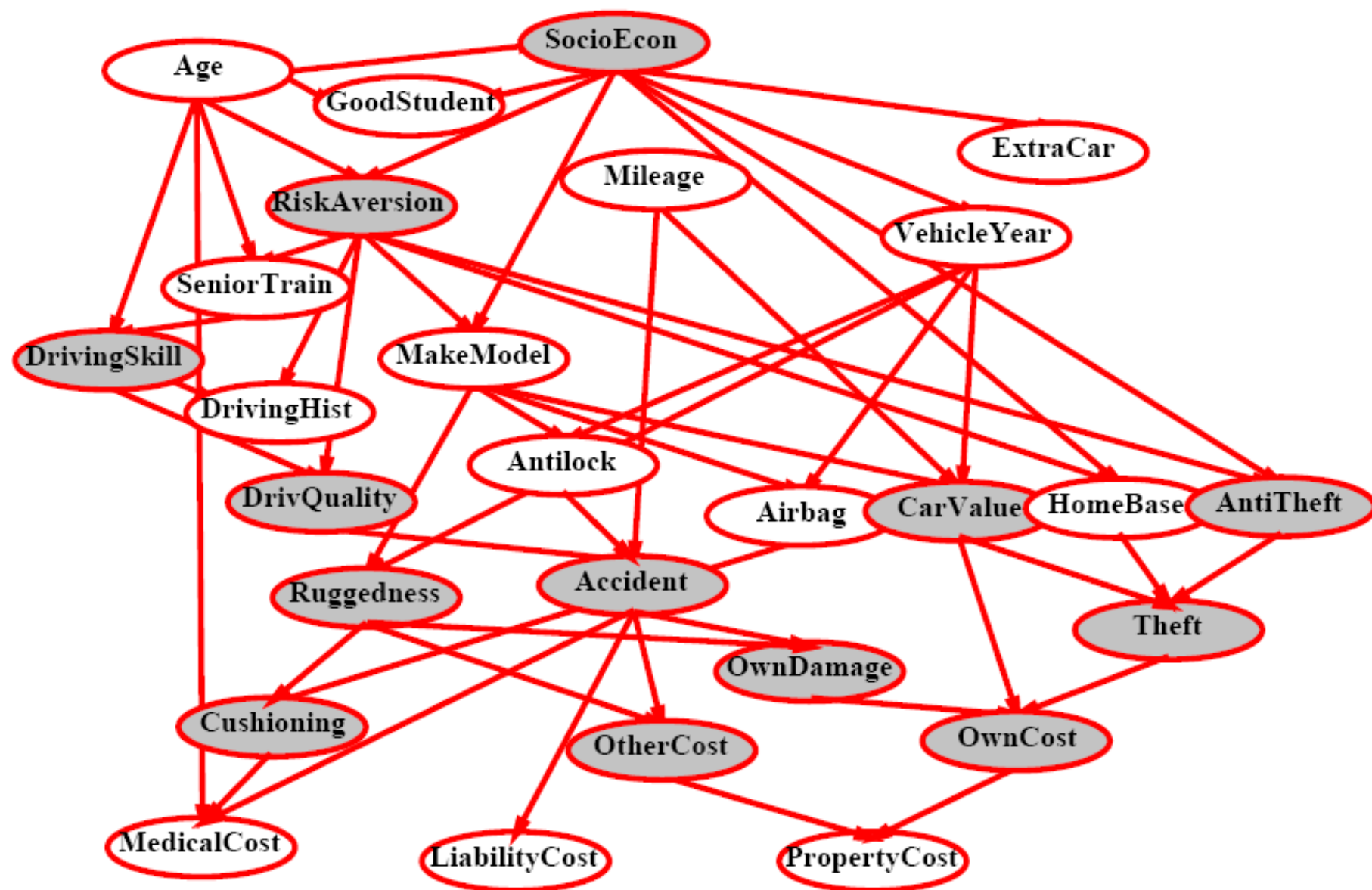- Network is less compact: 1 + 2 + 4 + 2 + 4 = 13 numbers needed (vs. 1+1+4+2+2=10 for the causal ordering)

# Why store it in causal order? A: Saves memory

ZJUI

- Suppose we have a Boolean variable $X_i$ with k Boolean parents. How many rows does its conditional probability table have?
  - $2^k$ rows for all the combinations of parent values
  - Each row requires one number for P($X_i$ = true | parent values)

- If each variable has no more than k parents, how many numbers does the complete network require?
  - $O(n \cdot 2^k)$ numbers – vs. $O(2^n)$ for the full joint distribution

- How many nodes for the burglary network?
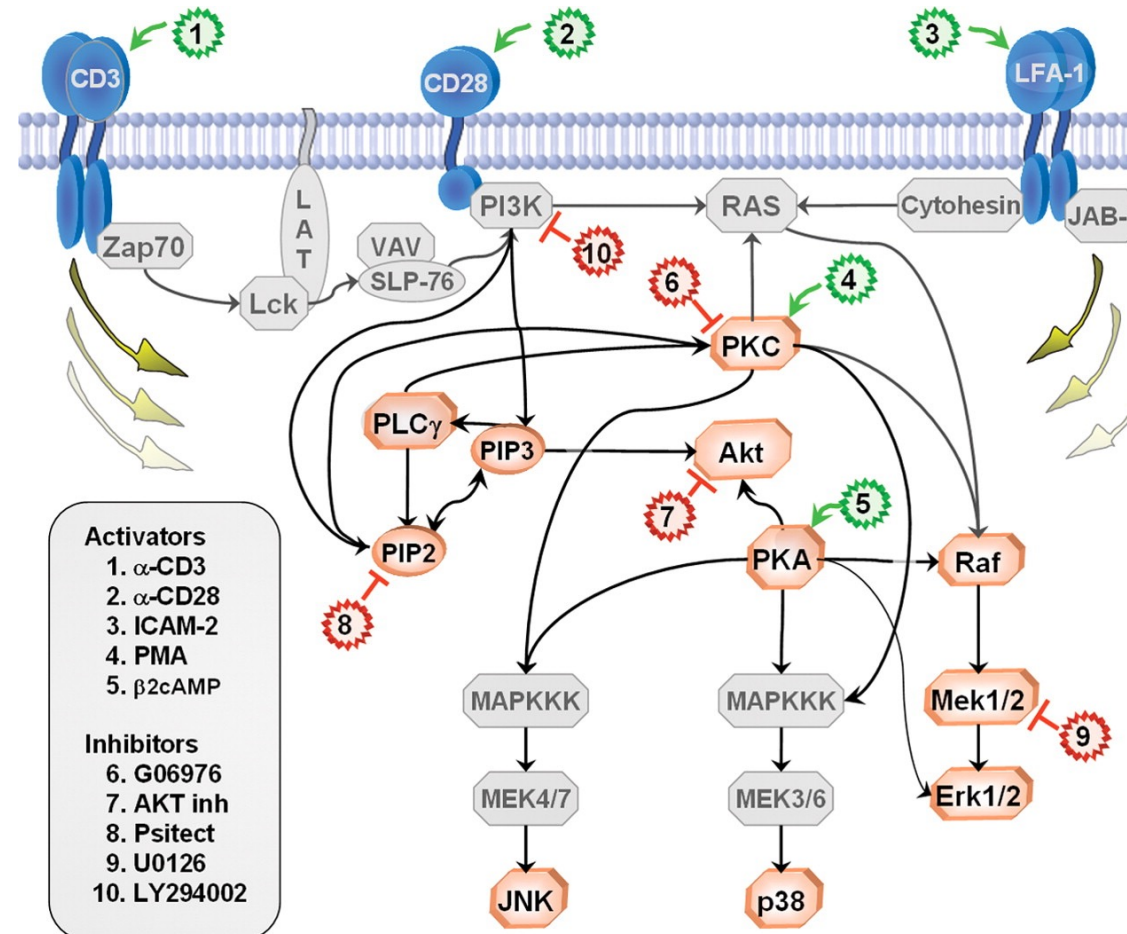  - 1 + 1 + 4 + 2 + 2 = 10 numbers (vs. $2^5$-1 = 31)

Car

ZJUI

# A more realistic Bayes Network:  Car diagnosis

- Initial observation: car won't start
- Orange: "broken, so fix it" nodes
- Green: testable evidence
- Gray: "hidden variables" to ensure sparse structure, reduce parameters

**Causal Protein-Signaling Networks Derived from Multiparameter Single-Cell Data**

Karen Sachs, Omar Perez, Dana Pe'er, Douglas A. Lauffenburger, and Garry P. Nolan

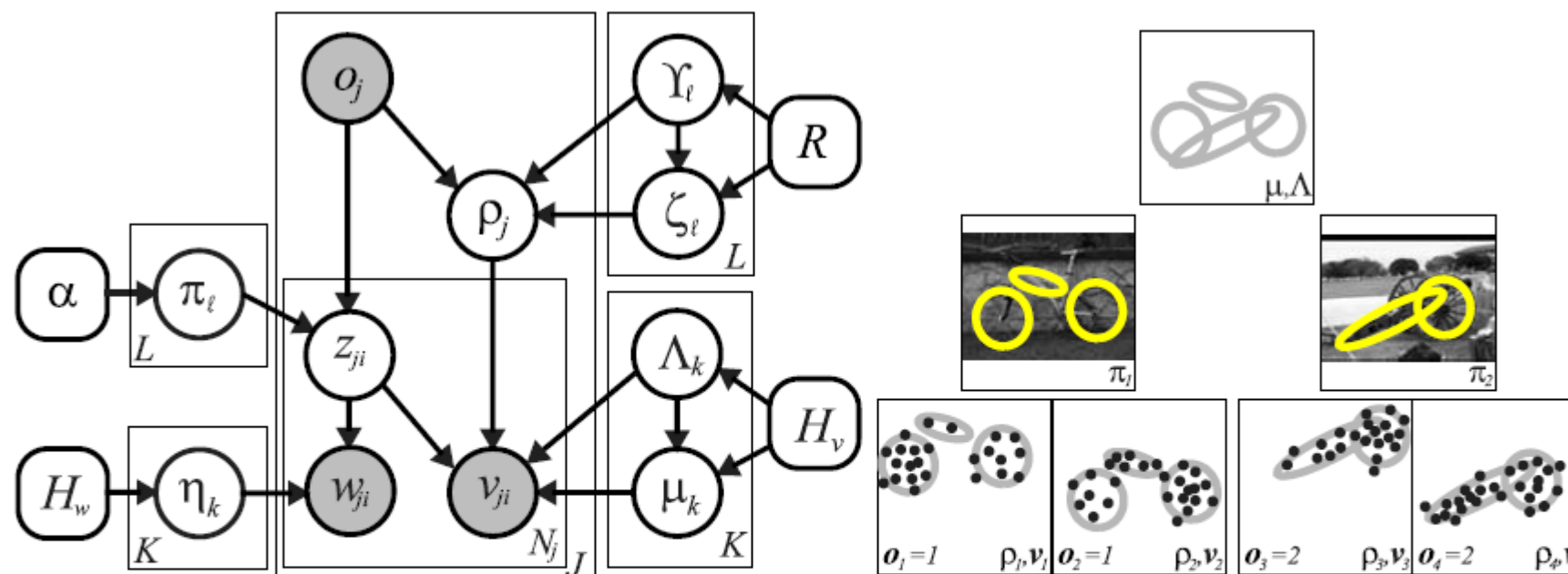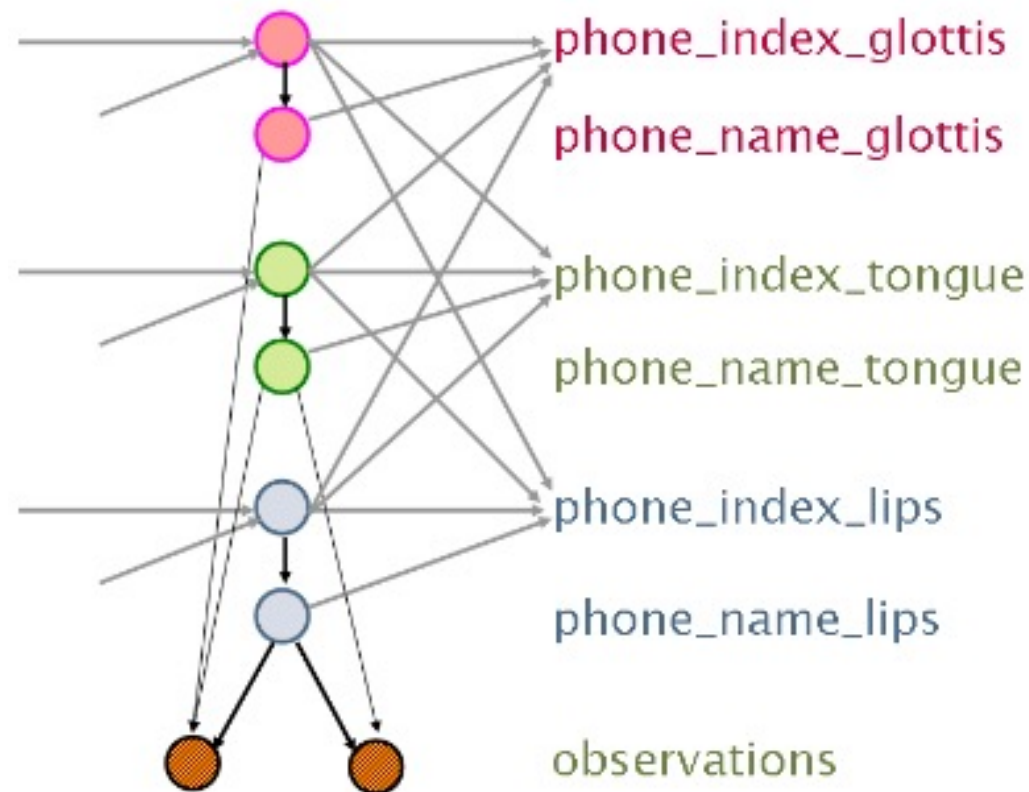(22 April 2005) *Science* **308** (5721), 523.

Fig. 3 A parametric, fixed-order model which describes the visual appearance of $L$ object categories via a common set of $K$ shared parts. The $j^{th}$ image depicts an instance of object category $o_j$, whose position is determined by the reference transformation $\rho_j$. The appearance $w_{ji}$ and position $v_{ji}$, relative to $\rho_j$, of visual features are determined by assignments $z_{ji} \sim \pi_{o_j}$ to latent parts. The cartoon example illustrates how a wheel part might be shared among two categories, *bicycle* and *cannon*. We show feature positions (but not appearance) for two hypothetical samples from each category

**Describing Visual Scenes Using Transformed Objects and Parts**

E. Sudderth, A. Torralba, W. T. Freeman, and A. Willsky.

*International Journal of Computer Vision*, No. 1-3, May 2008, pp. 291-330.

**Audiovisual Speech Recognition with Articulator Positions as Hidden Variables**

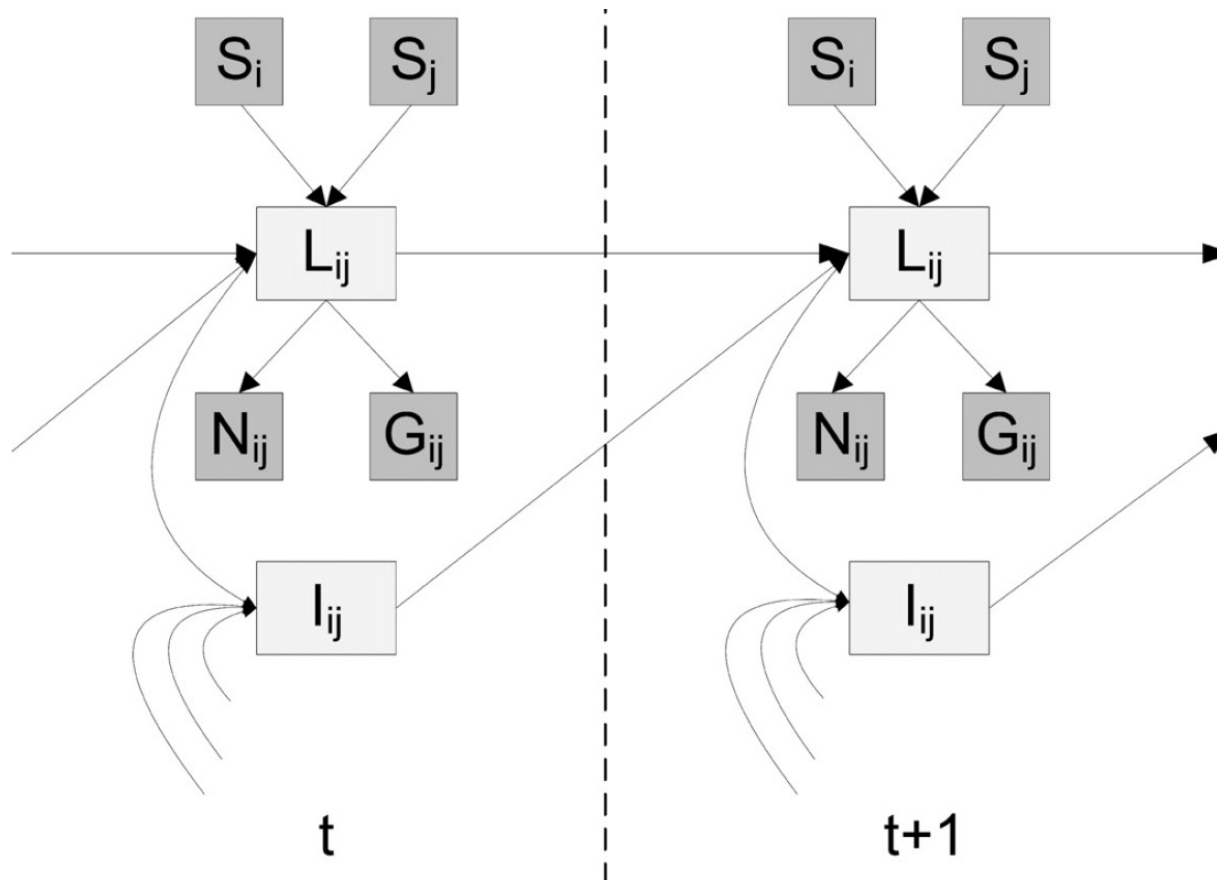Mark Hasegawa-Johnson, Karen Livescu, Partha Lal and Kate Saenko

*International Congress on Phonetic Sciences* 1719:299-302, 2007

**Detecting interaction links in a collaborating group using manually annotated data**

S. Mathur, M.S. Poole, F. Pena-Mora, M. Hasegawa-Johnson, N. Contractor

*Social Networks* 10.1016/j.socnet.2012.04.002

- **Link:** $L_{ij} = 1$ if #i is listening to #j.
- **Indirect:** $I_{ij} = 1$ if #i and #j are both listening to the same person.
- **Speaking**: $S_i = 1$ if the i'th person is speaking.
- **Gaze:** $G_{ij} = 1$ if #i is looking at #j.
- **Neighborhood:** $N_{ij} = 1$ if they're near one another

**Detecting interaction links in a collaborating group using manually annotated data**

- Bayesian networks provide a natural representation for (causally induced) conditional independence

- Topology + conditional probability tables

- Generally easy for domain experts to construct