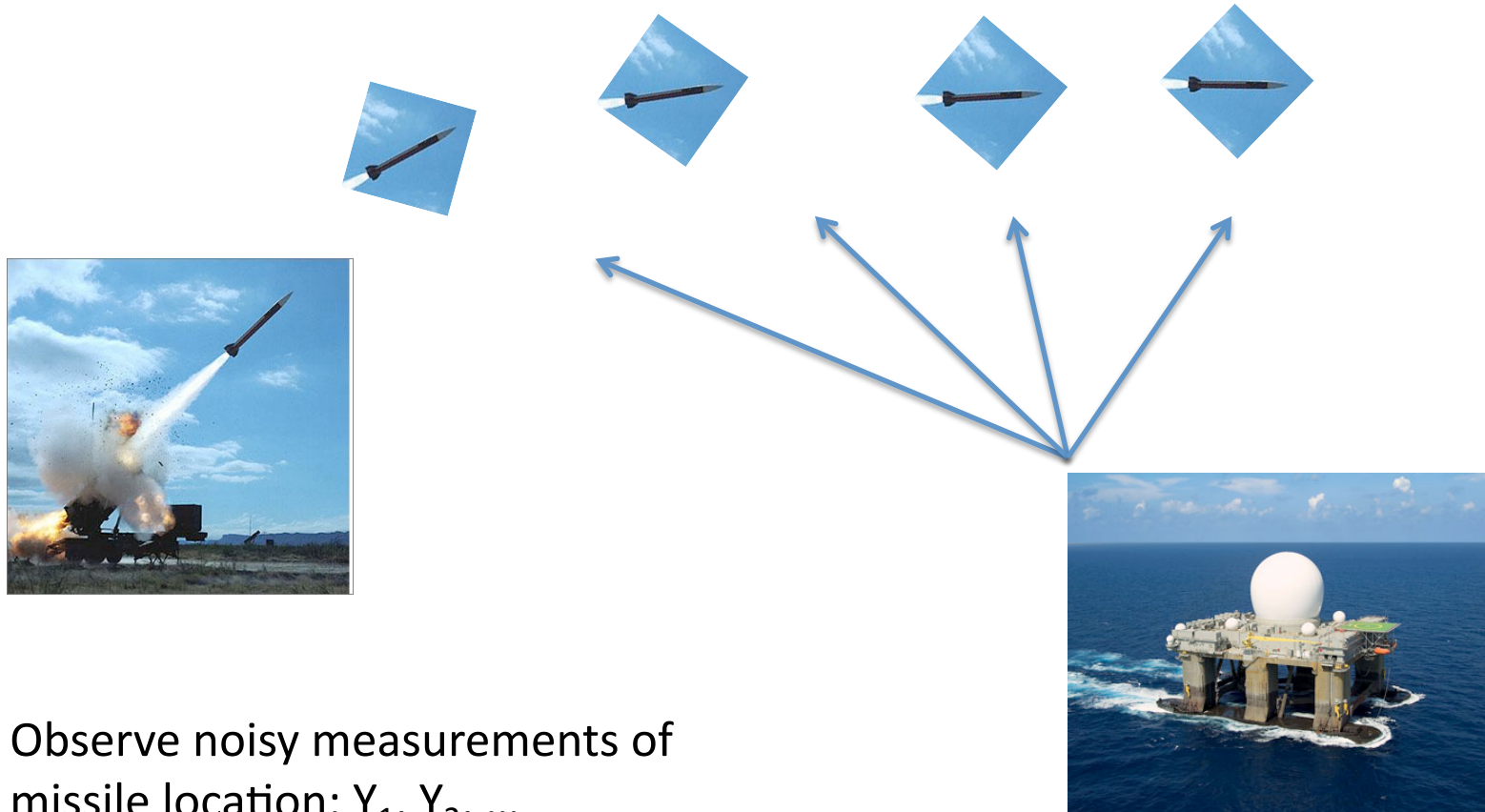# Hidden Markov models
# Lecture 23

David Sontag

New York University

# Example application: Tracking



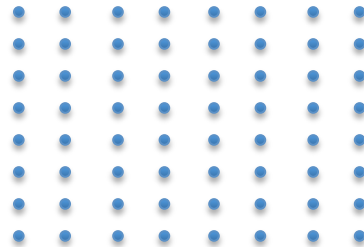Observe noisy measurements of missile location: $Y_1$, $Y_2$, ...

Radar

Where is the missile **now**? Where will it be in 10 seconds?
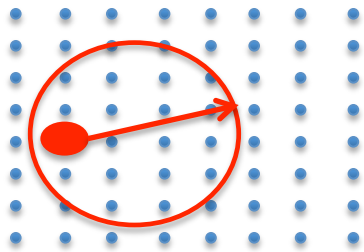
# Probabilistic approach

- Our measurements of the missile location were $Y_1, Y_2, ..., Y_n$

- Let $X_t$ be the *true* <missile location, velocity> at time t

- To keep this simple, suppose that everything is discrete, i.e. $X_t$ takes the values 1, ..., k

Grid the space:

# Probabilistic approach

- First, we specify the *conditional* distribution $\Pr(X_t \mid X_{t-1})$:



From basic physics, we can bound the distance that the missile can have traveled

- Then, we specify $\Pr(Y_t \mid X_t = <(10,20),\ 200\ \text{mph}$ toward the northeast$>)$:

With probability ½, $Y_t = X_t$ (ignoring the velocity). Otherwise, $Y_t$ is a uniformly chosen grid location

# Hidden Markov models

- Assume that the **joint** distribution on $X_1, X_2, \ldots, X_n$ and $Y_1, Y_2, \ldots, Y_n$ factors as follows:

$$\Pr(x_1, \ldots x_n, y_1, \ldots, y_n) = \Pr(x_1)\Pr(y_1 \mid x_1)\prod_{t=2}^{n}\Pr(x_t \mid x_{t-1})\Pr(y_t \mid x_t)$$

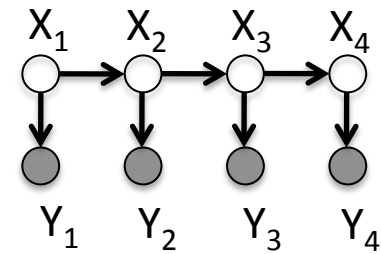- To find out where the missile is *now*, we do **marginal inference**:

$$\Pr(x_n \mid y_1, \ldots, y_n)$$

- To find the most likely *trajectory*, we do **MAP (maximum a posteriori) inference**:

$$\arg\max_{\mathbf{x}}\Pr(x_1, \ldots, x_n \mid y_1, \ldots, y_n)$$

# Inference

- Recall, to find out where the missile is now, we do marginal inference: $\Pr(x_n \mid y_1, \ldots, y_n)$



- How does one **compute** this?

- Applying rule of conditional probability, we have:

$$\Pr(x_n \mid y_1, \ldots, y_n) = \frac{\Pr(x_n, y_1, \ldots, y_n)}{\Pr(y_1, \ldots, y_n)}$$

- Naively, would seem to require $k^{n-1}$ summations,

$$\Pr(x_n, y_1, \ldots, y_n) = \sum_{x_1, \ldots, x_{n-1}} \Pr(x_1, \ldots, x_n, y_1, \ldots, y_n)$$

Is there a more efficient algorithm?

# Marginal inference in HMMs

- Use **dynamic programming**

$$\Pr(A = a) = \sum_b \Pr(B = b, A = a)$$

$$\Pr(x_n, y_1, \ldots, y_n) = \sum_{x_{n-1}} \Pr(x_{n-1}, x_n, y_1, \ldots, y_n)$$

$$\Pr(\vec{A} = \vec{a}, \vec{B} = \vec{b}) = \Pr(\vec{A} = \vec{a}) \Pr(\vec{B} = \vec{b} \mid \vec{A} = \vec{a})$$

$$= \sum_{x_{n-1}} \Pr(x_{n-1}, y_1, \ldots, y_{n-1}) \Pr(x_n, y_n \mid x_{n-1}, y_1, \ldots, y_{n-1})$$

<span style="color:red">Conditional independence in HMMs</span>

$$= \sum_{x_{n-1}} \Pr(x_{n-1}, y_1, \ldots, y_{n-1}) \Pr(x_n, y_n \mid x_{n-1})$$

$$\Pr(A = a, B = b) = \Pr(A = a) \Pr(B = b \mid A = a)$$

$$= \sum_{x_{n-1}} \Pr(x_{n-1}, y_1, \ldots, y_{n-1}) \Pr(x_n \mid x_{n-1}) \Pr(y_n \mid x_n, x_{n-1})$$

<span style="color:red">Conditional independence in HMMs</span>

$$= \sum_{x_{n-1}} \Pr(x_{n-1}, y_1, \ldots, y_{n-1}) \Pr(x_n \mid x_{n-1}) \Pr(y_n \mid x_n)$$
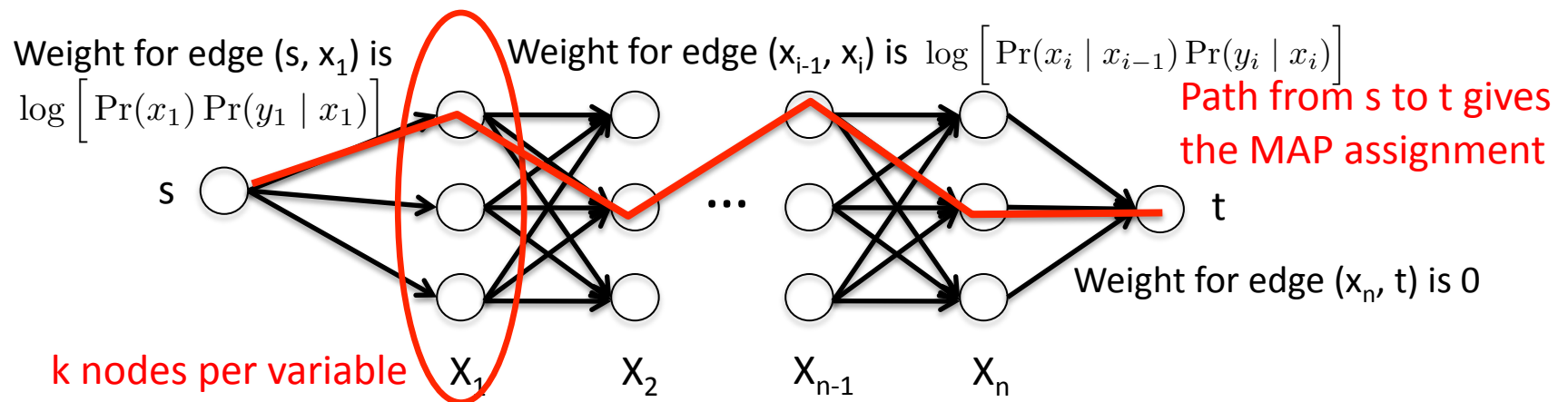
- For n=1, initialize  $\Pr(x_1, y_1) = \Pr(x_1) \Pr(y_1 \mid x_1)$

- Total running time is O(nk) – linear time!  <span style="color:red">Easy to do **filtering**</span>

# MAP inference in HMMs

- MAP inference in HMMs can *also* be solved in linear time!

$$\arg\max_{\mathbf{x}} \Pr(x_1, \ldots x_n \mid y_1, \ldots, y_n) = \arg\max_{\mathbf{x}} \Pr(x_1, \ldots x_n, y_1, \ldots, y_n)$$

$$= \arg\max_{\mathbf{x}} \log \Pr(x_1, \ldots x_n, y_1, \ldots, y_n)$$

$$= \arg\max_{\mathbf{x}} \ \log\Big[\Pr(x_1)\Pr(y_1 \mid x_1)\Big] + \sum_{i=2}^{n} \log\Big[\Pr(x_i \mid x_{i-1})\Pr(y_i \mid x_i)\Big]$$

- Formulate as a shortest paths problem

Weight for edge (s, $x_1$) is
$\log\Big[\Pr(x_1)\Pr(y_1 \mid x_1)\Big]$

Weight for edge ($x_{i-1}$, $x_i$) is $\log\Big[\Pr(x_i \mid x_{i-1})\Pr(y_i \mid x_i)\Big]$



Path from s to t gives
the MAP assignment

Weight for edge ($x_n$, t) is 0

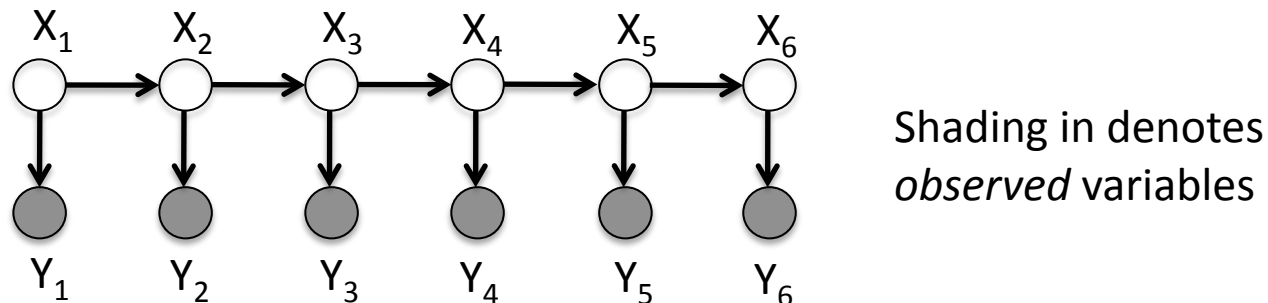k nodes per variable $X_1$   $X_2$   $X_{n-1}$   $X_n$

Called the Viterbi algorithm

# Applications of HMMs

- Speech recognition
  - Predict phonemes from the sounds forming words (i.e., the actual signals)

- Natural language processing
  - Predict parts of speech (verb, noun, determiner, etc.) from the words in a sentence

- Computational biology
  - Predict intron/exon regions from DNA
  - Predict protein structure from DNA (locally)

- And many many more!

# Hidden Markov models

- We can represent a hidden Markov model with a graph:



Shading in denotes *observed* variables

$$\Pr(x_1, \ldots x_n, y_1, \ldots, y_n) = \Pr(x_1) \Pr(y_1 \mid x_1) \prod_{t=2}^{n} \Pr(x_t \mid x_{t-1}) \Pr(y_t \mid x_t)$$

- There is a 1-1 mapping between the graph structure and the factorization of the joint distribution

- More generally, a **Bayesian network** is defined by a graph *G=(V,E)* with one node per variable, and a distribution for each variable conditioned on its parents' values:

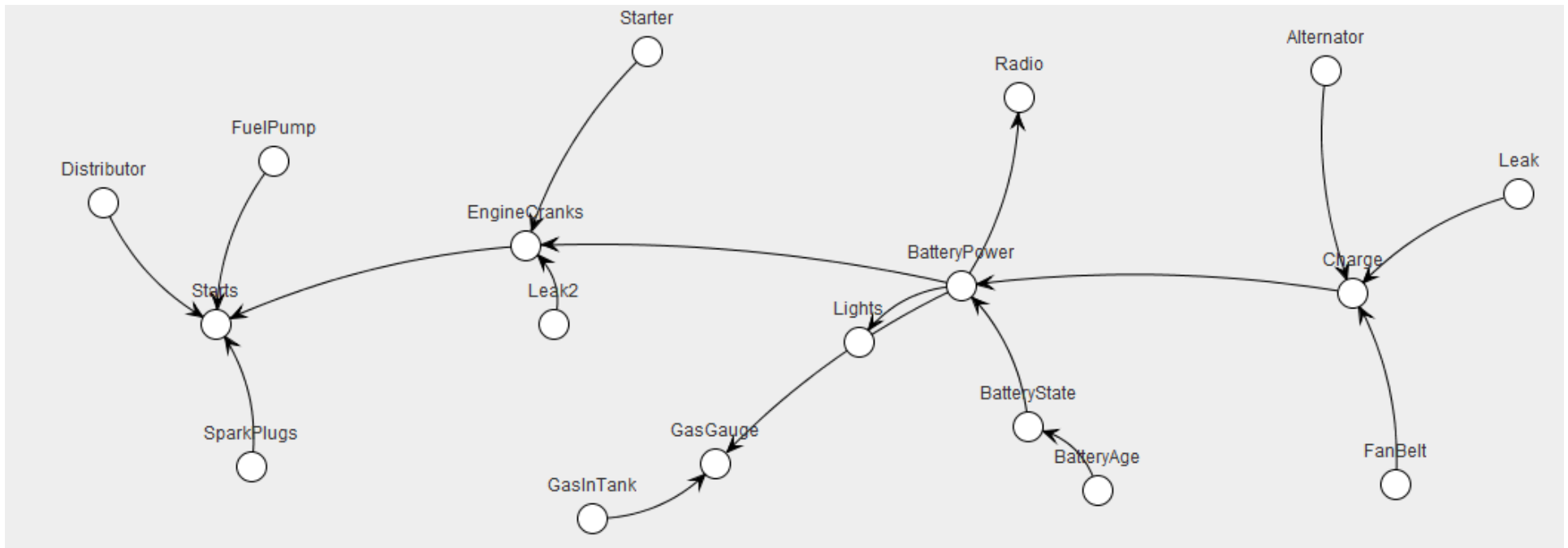$$\Pr(\mathbf{v}) = \prod_{i \in V} \Pr(v_i \mid \mathbf{v}_{pa(i)})$$

pa(i) denotes the parents of variable i

# Bayesian networks

$$\Pr(\mathbf{v}) = \prod_{i \in V} \Pr(v_i \mid \mathbf{v}_{pa(i)})$$

## Will your car start this morning?



Heckerman *et al.*, Decision-Theoretic Troubleshooting, 1995

# Bayesian networks

$$\Pr(\mathbf{v}) = \prod_{i \in V} \Pr(v_i \mid \mathbf{v}_{pa(i)})$$

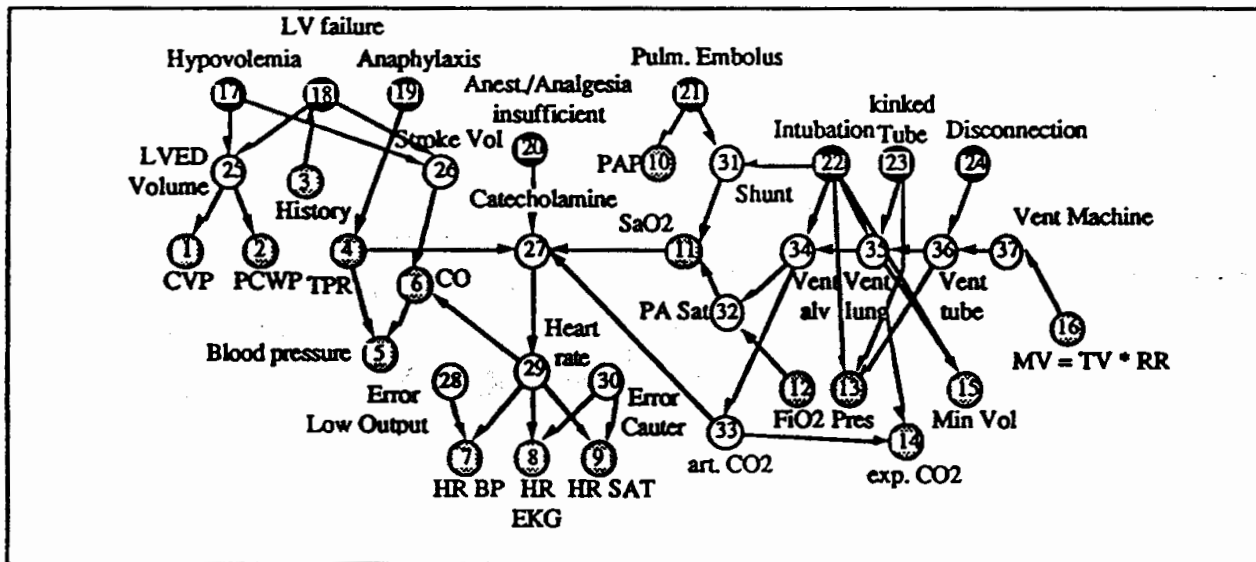What is the differential diagnosis?



Fig. 1 The ALARM network representing causal relationships is shown with diagnostic (●), intermediate (○) and measurement (◉) nodes. CO: cardiac output, CVP: central venous pressure, LVED volume: left ventricular end-diastolic volume, LV failure: left ventricular failure, MV: minute ventilation, PA Sat: pulmonary artery oxygen saturation, PAP: pulmonary artery pressure, PCWP: pulmonary capillary wedge pressure, Pres: breathing pressure, RR: respiratory rate, TPR: total peripheral resistance, TV: tidal volume

Beinlich *et al.*, The ALARM Monitoring System, 1989

acm

MORE ACM AWARDS

A.M. TURING AWARD

Search TYPE HERE

A.M. TURING AWARD WINNERS BY...

| ALPHABETICAL LISTING | YEAR OF THE AWARD | RESEARCH SUBJECT |
|---|---|---|

📖 **Photo-Essay**

**BIRTH:**

September 4, 1936, Tel Aviv.

**EDUCATION:**

B.S., Electrical Engineering (Technion, 1960); M.S., Electronics (Newark College of Engineering, 1961); M.S., Physics (Rutgers University, 1965); Ph.D., Electrical Engineering (Polytechnic Institute of Brooklyn, 1965).

**EXPERIENCE:**

Research Engineer, New York University Medical School (1960–1961); Instructor,

# JUDEA PEARL

United States – 2011

## CITATION

For fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning.

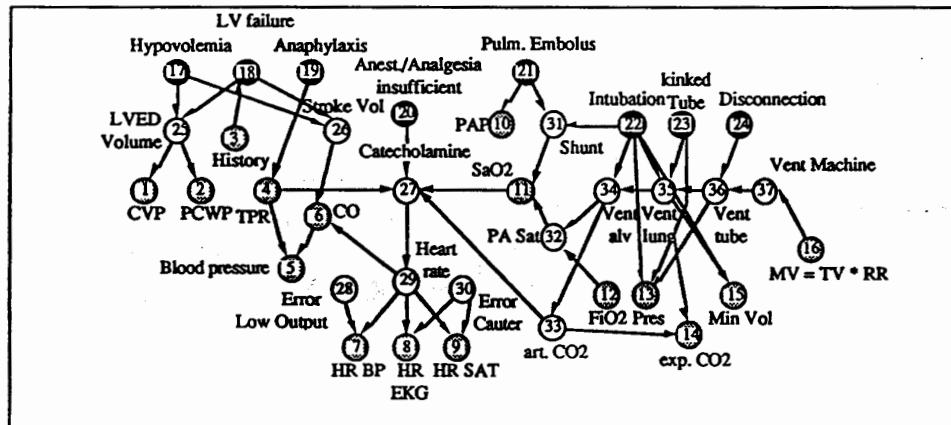| SHORT ANNOTATED BIBLIOGRAPHY | ACM DL AUTHOR PROFILE | ACM TURING AWARD LECTURE VIDEO | RESEARCH SUBJECTS | ADDITIONAL MATERIALS |
|---|---|---|---|---|

Judea Pearl created the representational and computational foundation for the processing of information under uncertainty.

He is credited with the invention of *Bayesian networks*, a mathematical formalism for defining complex probability models, as well as the principal algorithms used for inference in these models. This work not only revolutionized the field of artificial intelligence but also became an important tool for many other branches of engineering and the natural sciences. He later created a mathematical framework for causal inference that has had significant impact in the social sciences.

Judea Pearl was born on September 4, 1936, in Tel Aviv, which was at that time administered under the British Mandate for Palestine. He grew up in Bnei Brak, a Biblical town his grandfather went to reestablish in 1924. In 1956, after serving in the Israeli army and joining a Kibbutz, Judea decided to study engineering. He attended the Technion, where he met his wife, Ruth, and received a B.S. degree in Electrical Engineering in 1960. Recalling the Technion faculty members in a 2012 interview in the *Technion Magazine*, he emphasized the thrill of discovery:

# Inference in Bayesian networks

- Computing marginal probabilities in **tree** structured Bayesian networks is easy
    - The algorithm called "belief propagation" generalizes what we showed on the previous slides to arbitrary trees

- Wait... this isn't a tree! What can we do?

# Inference in Bayesian networks

- In some cases (such as this) we can *transform* this into what is called a "junction tree", and then run belief propagation
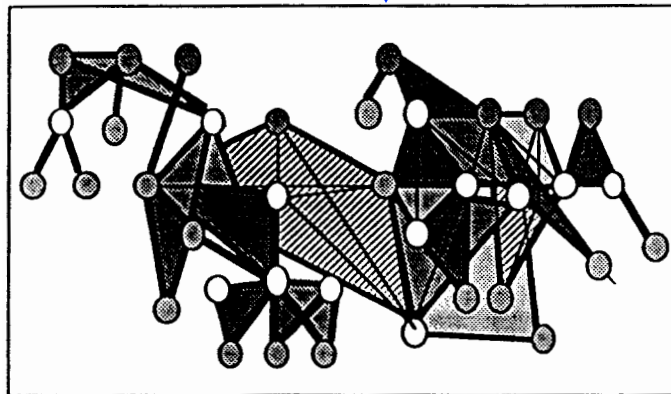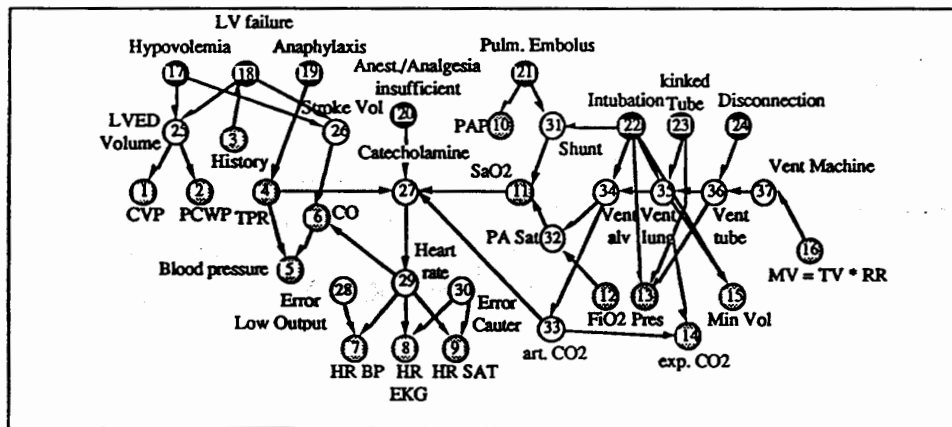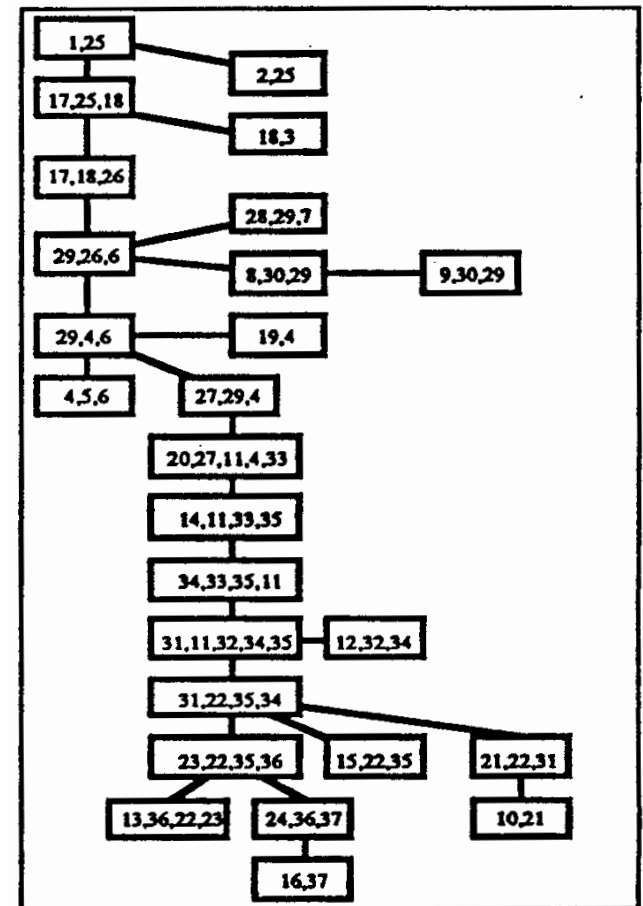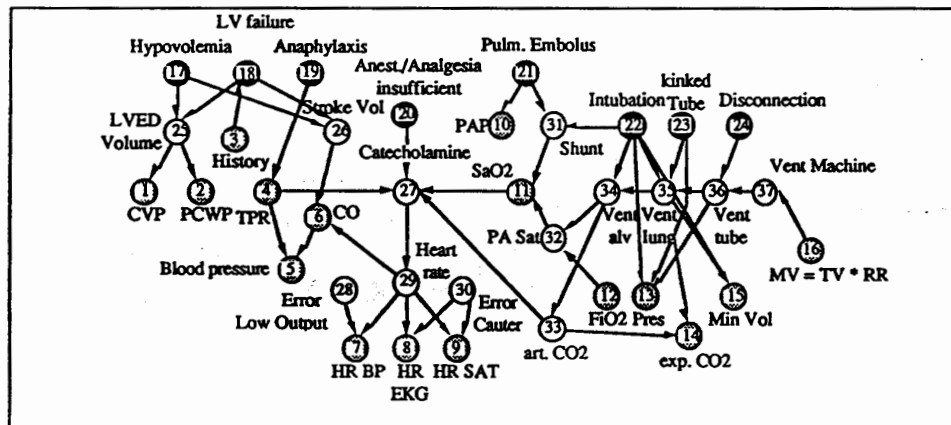




Fig. 7

Spiegelhalter's algorithm re-arranges the ALARM network by triangulation and clique formation. The cliques are shaded differently to make them visible.

# Approximate inference

- There is also a wealth of **approximate** inference algorithms that can be applied to Bayesian networks such as these



- Markov chain Monte Carlo algorithms repeatedly sample assignments for estimating marginals
- Variational inference algorithms (which are deterministic) attempt to fit a simpler distribution to the complex distribution, and then computes marginals for the simpler distribution