# Probabilistic Graph Models:
# from Bayesian to Factor Graphs

ECE/CS 498 DS U/G

Lecture 13

Ravi K. Iyer

Dept. of Electrical and Computer Engineering

University of Illinois at Urbana Champaign

# Announcements

- Today:
  - Monday March 4
  - Hidden Markov Model
  - Factor Graphs
- Wednesday (March 6):
  - Solve last year's midterm
- HW3 has been released; due on Friday, March 8
- MP2 Checkpoint 2 due on Friday, March 8
- Midterm on Monday (March 11) (**Note: Room change**)
  - ECE room 2013 (NetID starting with **a-m**), room 3013 (NetID starting with **n-z**)
  - Will start on time at 12:30 – 1:50pm
  - Bring one 8x11 sheet of notes
  - Closed book, no calculators or electronic devices
  - Bring your Univ IDs

# Announcements

# Overview of PGM Data Analytics/Modeling Process



**DATA COLLECTION**　　　**DATA ANALYSIS/Training**　　　**MODELING**

ECE ILLINOIS　　　　　　　　　　　　　I ILLINOIS

# Measurements from NCSA@Illinois: Five minute Snap Shot

- ## Goals:
  - Provide a system-level characterization of incidents and evaluate the intricacies of real-time diagnosis
  - Design protection strategies to reduce missed incidents and false positives
  - Experimentally Demonstrate new techniques in a sandbox

- ## Challenges



Five-Minute
Snapshot
of In-and-Out
Traffic
at NCSA

ECE ILLINOIS

ILLINOIS

# Five-Minute Snapshot of In-and-Out Traffic within NCSA@Illinois



(a)

(b)

# An Application in Security Data Analytics
## Individual components of an attack as attack progresses



**Attack stages for the credential stealing attack**

# Annotation and extracting patterns in past attacks

**Naïve Bayes**

a) Annotated events and attack stages in a pair of attacks

Database of 200 attacks

Attack 1
events $\epsilon_1$ $\epsilon_2$ $\epsilon_3$ $\epsilon_3$ $\epsilon_4$  time
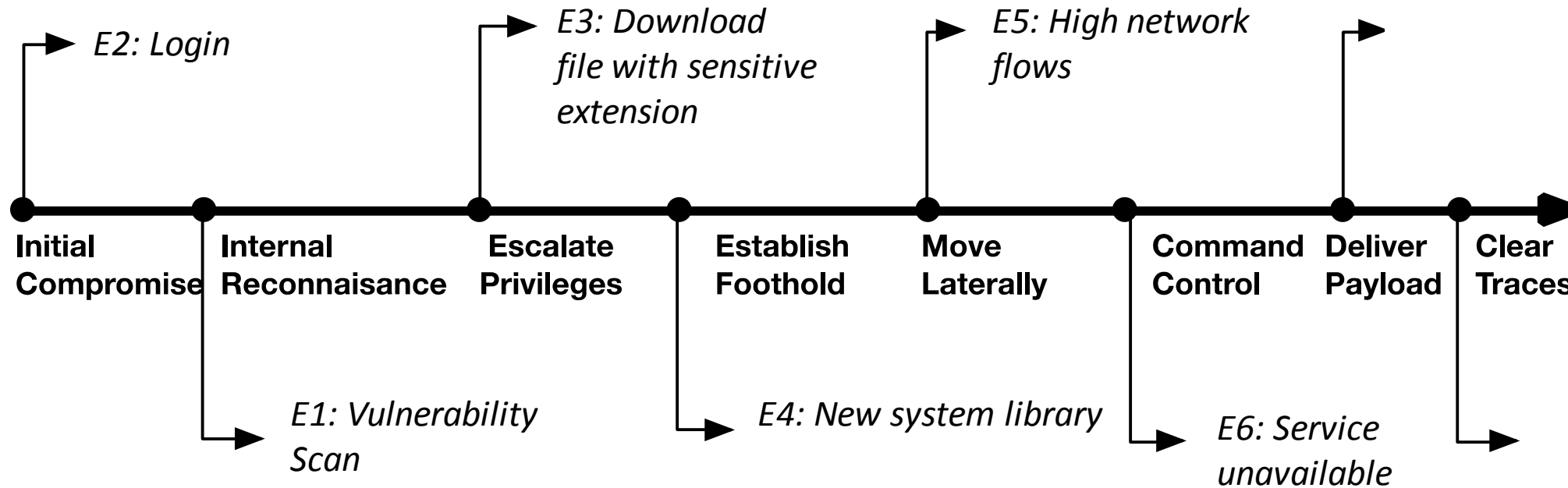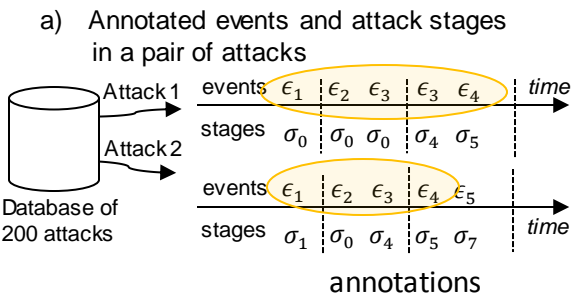stages $\sigma_0$ $\sigma_0$ $\sigma_0$ $\sigma_4$ $\sigma_5$

Attack 2
events $\epsilon_1$ $\epsilon_2$ $\epsilon_3$ $\epsilon_4$ $\epsilon_5$  time
stages $\sigma_1$ $\sigma_0$ $\sigma_4$ $\sigma_5$ $\sigma_7$

annotations

c) Example patterns, stages, probabilities, and significance learned from the attack pair

| Pattern | Attack stages | Probability in past attacks | Significance (p-value) |
|---|---|---|---|
| $[\epsilon_1, \epsilon_3, \epsilon_4]$ | $[\sigma_1, \sigma_4, \sigma_5]$ | $q_a$ | $p_a$ |
| $[\epsilon_1]$ | $[\sigma_0 \mid \sigma_1]$ | $q_b$ | $p_b$ |

…

**Bayesian Network**

b) Event-stage annotation table for the attack pair (Attack 1 and Attack 2)

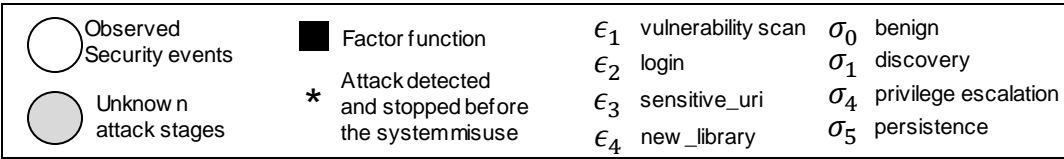| Event | Attack stage |
|---|---|
| $\{\epsilon_1\}$ | $\{\sigma_0 \mid \sigma_1\}$ |
| $\{\epsilon_2\}$ | $\{\sigma_0\}$ |
| $\{\epsilon_3\}$ | $\{\sigma_4\}$ |
| $\{\epsilon_4\}$ | $\{\sigma_5\}$ |
| $\{\epsilon_5\}$ | $\{\sigma_7\}$ |

**Dynamic Bayesian Network**

**Hidden Markov Model**

**Factor Graphs**

**OFFLINE ANNOTATION ON PAST ATTACKS**

**OFFLINE LEARNING OF PATTERNS**

**PROBABILISTIC GRAPHICAL MODELS**

Note: $\epsilon_i$ is the corresponding value of an event $E_t$

- ◯ Observed Security events
- ⬤ Unknown attack stages
- ■ Factor function
- * Attack detected and stopped before the system misuse
- $\epsilon_1$ vulnerability scan
- $\epsilon_2$ login
- $\epsilon_3$ sensitive_uri
- $\epsilon_4$ new_library
- $\sigma_0$ benign
- $\sigma_1$ discovery
- $\sigma_4$ privilege escalation
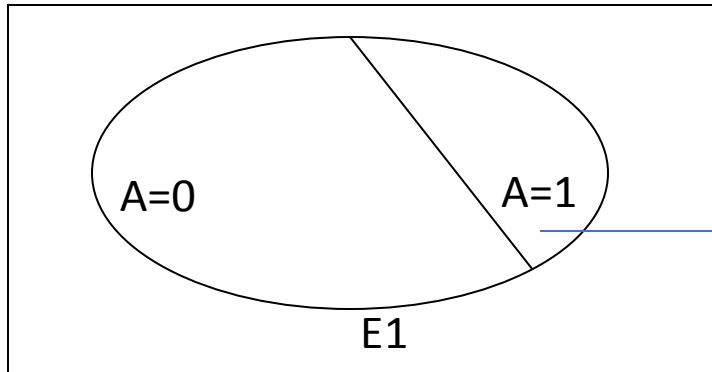- $\sigma_5$ persistence

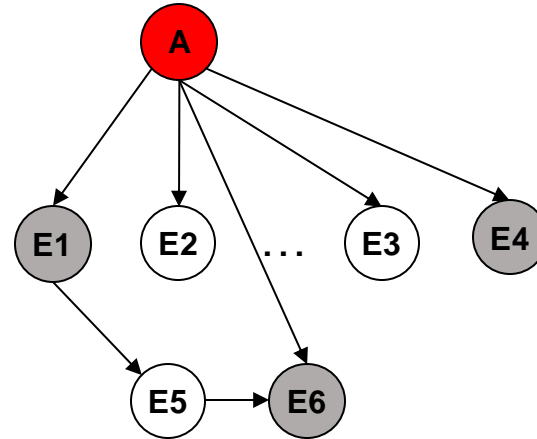# Modeling the credential stealing attack using Naïve Bayes vs. Bayesian Network

**Naïve Bayes**

$$P(A|E_1, E_2, ..., E_4) = P(A) \prod_i P(E_i|A)$$

Is $(E1, E2, ..., E4)$ represents Benign activity?
$[P(E_1|A = \text{Benign}) .... P(E_4|A = \text{Benign})]P(A = \text{Benign}) > [P(E_1|A = Attack) ... P(E_4|A = Attack)]P(A = Attack)$

A=0       A=1

E1

$P(E_1|A = 1)$

**Bayesian Network**

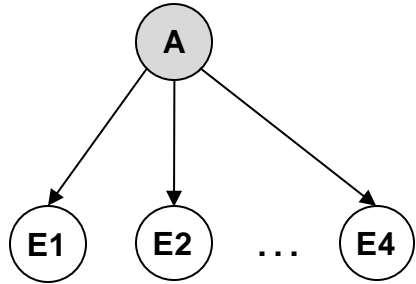Joint Distribution: $P(E_1, E_2, ..., E_n, A) = P(A) \prod_{i=1}^{n} P(E_i|parents(E_i))$

Hypothesis:

$$P(A = attack|E_1, E_4\ E_6) = ?$$

$$P(A = benign|E_1, E_4\ E_6) = ?$$

| ID | Description |
|----|-------------|
| A | Attack |
| E1 | Vulnerability scan |
| E2 | Login |
| E3 | Download file with sensitive extension |
| E4 | New system library |
| E5 | High network flows |
| E6 | Service unavailable |

ECE ILLINOIS                    I ILLINOIS

# Modeling the credential stealing attack using Naïve Bayes vs. Bayesian Network



**Naïve Bayes**

**Bayesian Network**

Possible causal loop

| ID | Description |
|----|-------------|
| A | Attack |
| E1 | Vulnerability scan |
| E2 | Login |
| E3 | Download file with sensitive extension |
| E4 | New system library |
| E5 | High network flows |
| E6 | Service unavailable |

**Model assumptions**
1. All events share the same parent variable
2. All events are conditionally independent

**Advantage:**
Simplify calculation of posterior probability on A

**Model assumptions**
1. An event can be preceded (causal) by another event
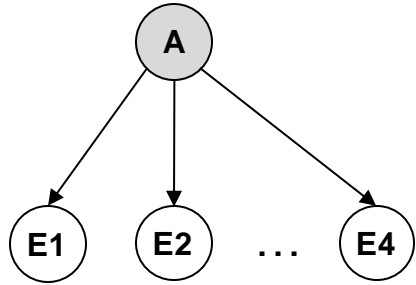2. There is no cycle in the network

**Disadvantage**
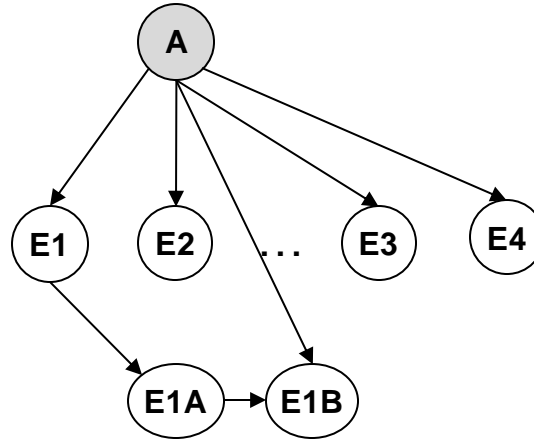Explicitly assume causal relationships
(Causality may not be clear from the data)
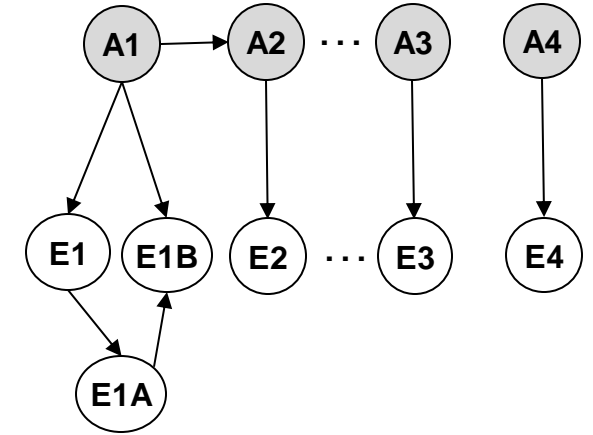For complicated attacks, causal loops may form and render the BN invalid

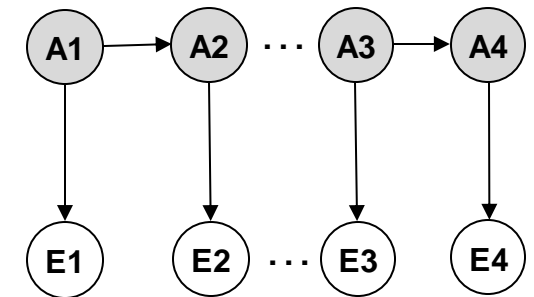# Modeling the credential stealing attack using Naïve Bayes vs. Bayesian Network



**Naïve Bayes**

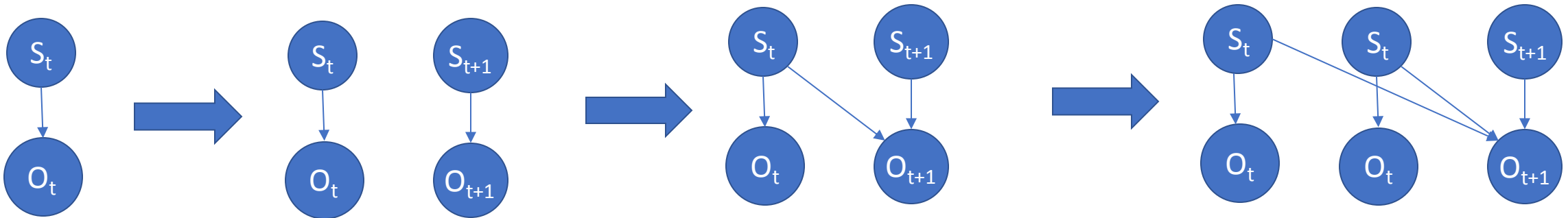**Bayesian Network**

**Dynamic Bayesian Network**

- When we consider the time evolution of the BN each variable in each timestep together, e.g., *t* and *t+1,* we have a Dynamic Bayesian Network that captures the first-order dependency --> referred to as the Markov Property

- This concept can be extended to higher order dependencies e.g on , t-2, t-3, … and is called a higher-order Markov property, e.g., 2nd or 3rd Markov property.

$$P(A_1, E_1, …, A_n, E_n) = P(A_1)P(E_1|A_1) … P(E_{t+1}|A_{t+1})P(A_{t+1}|A_t)$$

**Hidden Markov Model**

# Dynamic Bayesian Networks

- We have considered BNs with a static set of random variables, e.g., two variables: only one measurement variable and one state variable of the system.

- In reality, data is often time series in which each time step $t$ has one measurement variable $O_t$ and one state variable $S_t$. Thus, the number of random variables is proportional with the number of timesteps.

- Without correlating the random variables in each timestep, we have T disconnected BNs

- When we correlate each variable in each timestep together, e.g., $t$ and $t+1$, we have a Dynamic Bayesian Network that captures the first-order Markov property.

- This concept can be extended for t, t+1, t+2, … and is called a higher-order Markov property, e.g., 2nd or 3rd



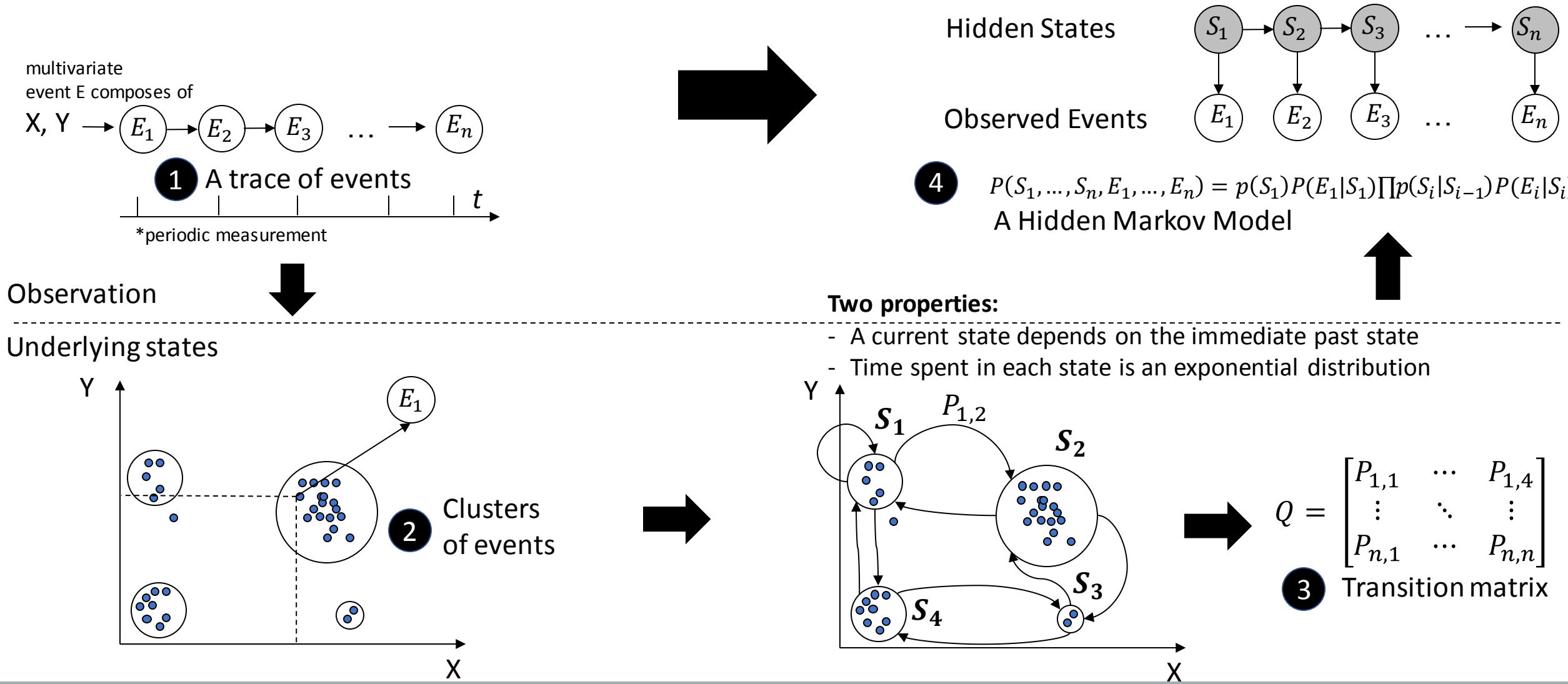$$P(S_t, O_t) = P(S_t)P(O_t|S_t)$$

$$P(S_t, O_t) = P(S_t)P(O_t|S_t)$$

$$P(S_{t+1}, O_{t+1}) = P(S_{t+1})P(O_{t+1}|S_{t+1})$$

$$P(S_t, S_{t+1}, O_t, O_{t+1}) = P(S_t)P(O_t|S_t)P(O_{t+1}|S_t, S_{t+1})P(S_{t+1})$$

# From a trace of events to a Hidden Markov Model



multivariate event E composes of X, Y

**1** A trace of events

*periodic measurement

Observation

Underlying states

**2** Clusters of events

Hidden States

Observed Events

**4** $P(S_1, \ldots, S_n, E_1, \ldots, E_n) = p(S_1)P(E_1|S_1)\prod p(S_i|S_{i-1})P(E_i|S_i)$

A Hidden Markov Model

**Two properties:**
- A current state depends on the immediate past state
- Time spent in each state is an exponential distribution

$$Q = \begin{bmatrix} P_{1,1} & \cdots & P_{1,4} \\ \vdots & \ddots & \vdots \\ P_{n,1} & \cdots & P_{n,n} \end{bmatrix}$$

**3** Transition matrix

# Hidden Markov Models

**Model assumptions**

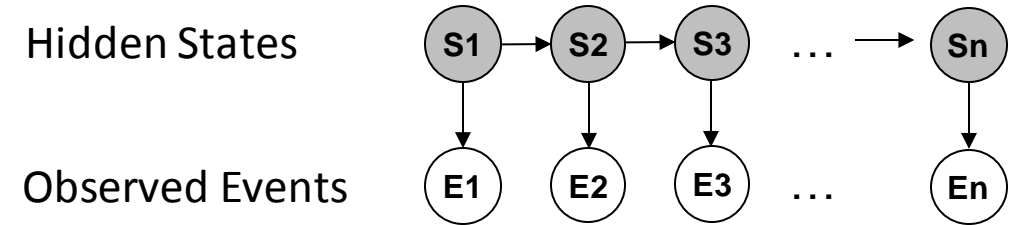An observation depends on its hidden state

A state variable only depends on the immediate previous state ( Markov assumption)

The future observations and the past observations are conditionally independent given the current hidden state

**Advantages:**

HMM can model sequential nature of input data (future depends on the past)

HMM has a linear-chain structure that clearly separates system state and observed events.

Hidden States

Observed Events

$$P(S_1, \ldots, S_n, E_1, \ldots, E_n) = p(S_1)P(E_1|S_1)\prod p(S_i|S_{i-1})P(E_i|S_i)$$

**A Hidden Markov model on observed events and system states**

# Markov Model

- Consider a system which can occupy one of *N* discrete *states* or *categories*

$$x_t \in \{1, 2, \ldots, N\} \longrightarrow \text{state at time } t$$
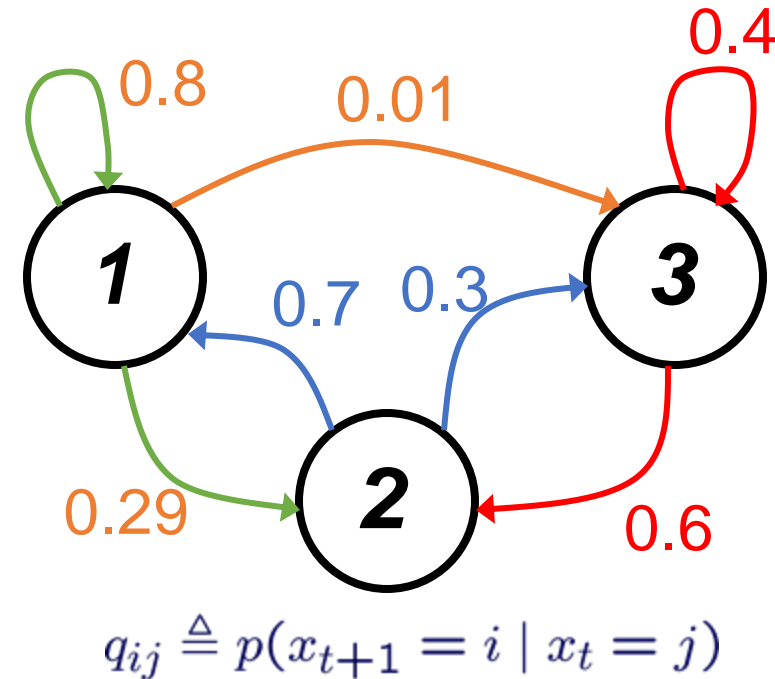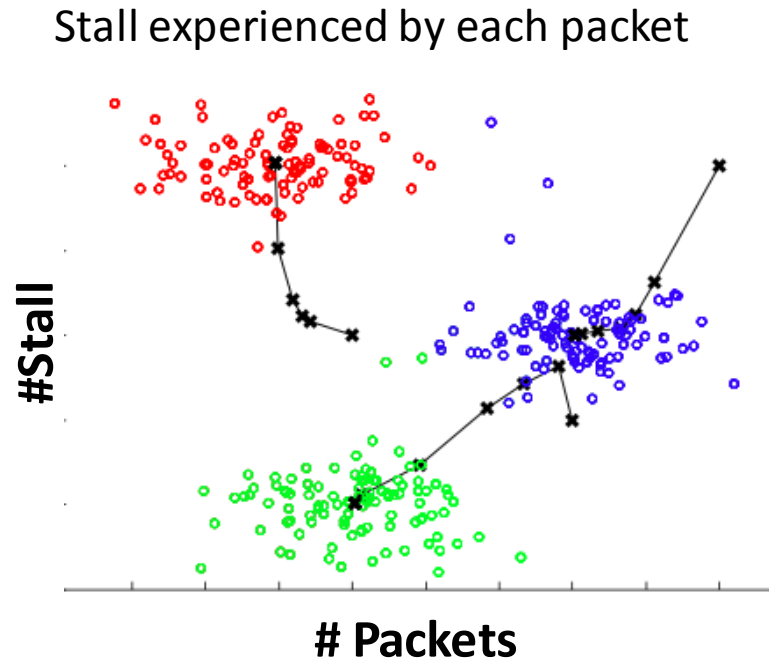
- We are interested in *stochastic* systems, in which state evolution is random

- Any *joint* distribution can be factored into a series of *conditional* distributions:

$$p(x_0, x_1, \ldots, x_T) = p(x_0) \prod_{t=1}^{T} p(x_t \mid x_0, \ldots, x_{t-1})$$

- For a *Markov* process, the next state depends only on the current state:

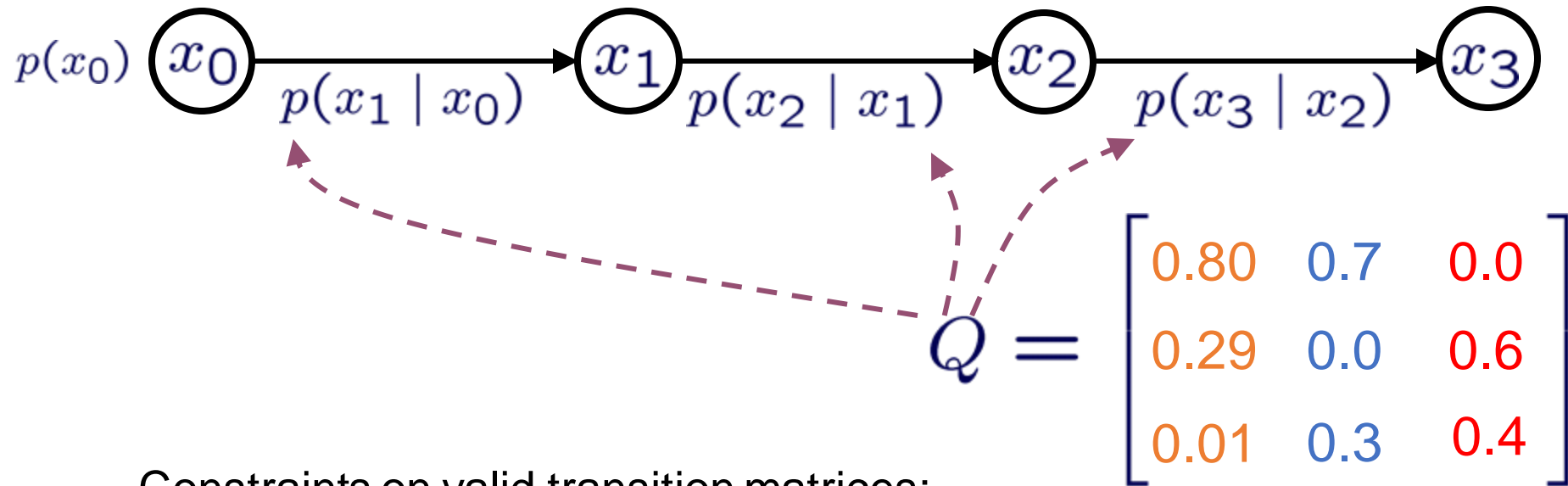$$p(x_{t+1} \mid x_0, \ldots, x_t) = p(x_{t+1} \mid x_t)$$

# State Transition Diagrams

Stall experienced by each packet



**# Packets**



$$q_{ij} \triangleq p(x_{t+1} = i \mid x_t = j)$$

- Think of a particle randomly following an arrow at each discrete time step
- Most useful when *N* small, and *Q* *sparse*

# Markov Chains: Graphical Models

$$p(x_0, x_1, \ldots, x_T) = p(x_0) \prod_{t=1}^{T} p(x_t \mid x_{t-1})$$
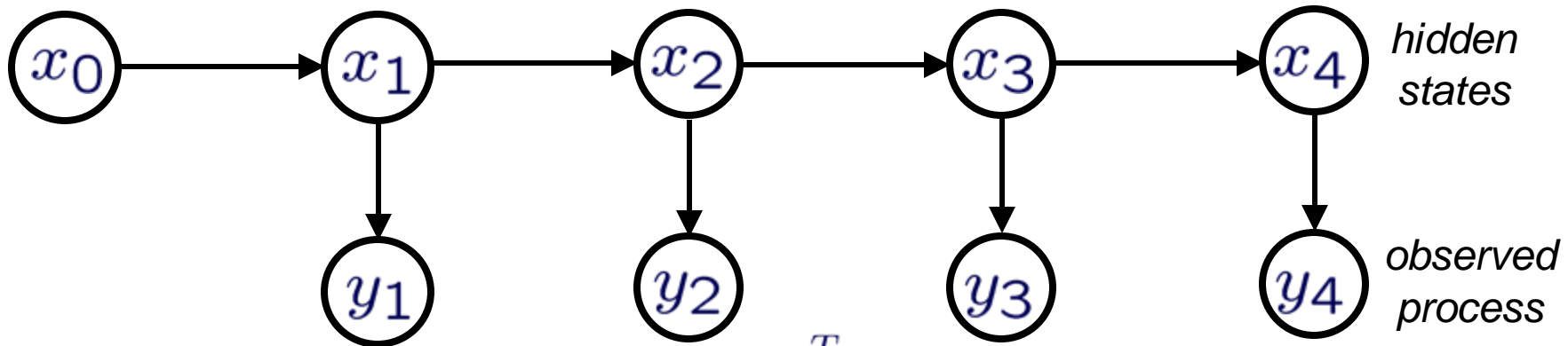


$p(x_0)$   $x_0$    $p(x_1 \mid x_0)$    $x_1$    $p(x_2 \mid x_1)$    $x_2$    $p(x_3 \mid x_2)$    $x_3$

$$Q = \begin{bmatrix} 0.80 & 0.7 & 0.0 \\ 0.29 & 0.0 & 0.6 \\ 0.01 & 0.3 & 0.4 \end{bmatrix}$$

Constraints on valid transition matrices:

$$q_{ij} \geq 0 \, , \quad \sum_{i=1}^{N} q_{ij} = 1 \quad \text{for all } j \qquad q_{ij} \triangleq p(x_{t+1} = i \mid x_t = j)$$

ECE ILLINOIS      I ILLINOIS

# Hidden Markov Models

- Stall exists due to congestion
- Not directly measurable at runtime (hidden)
- Motivates *hidden Markov models (HMM):*



hidden states

observed process

$$p(x_0, x_1, \ldots, x_T) = p(x_0) \prod_{t=1}^{T} p(x_t \mid x_{t-1}) p(y_t \mid x_t)$$

Given $x_t$, previous observations impact future observations

$$p(y_t, y_{t+1}, \ldots \mid x_t, y_{t-1}, y_{t-2}, \ldots) = p(y_t, y_{t+1}, \ldots \mid x_t)$$

# State Transition Matrices

- A *stationary* <span style="color:red">*stationary*</span> Markov chain with *N* states is described by an *NxN* *transition matrix:*

$$Q = \begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix}$$

$$q_{ij} \triangleq p(x_{t+1} = i \mid x_t = j)$$

- Constraints on valid transition matrices:

$$q_{ij} \geq 0 \qquad \sum_{i=1}^{N} q_{ij} = 1 \quad \text{for all } j$$

# State Transition Diagrams(Another Example)

$$q_{ij} \triangleq p(x_{t+1} = i \mid x_t = j)$$

$$Q = \begin{bmatrix} 0.5 & 0.1 & 0.0 \\ 0.3 & 0.0 & 0.4 \\ 0.2 & 0.9 & 0.6 \end{bmatrix}$$



- Think of a particle randomly following an arrow at each discrete time step
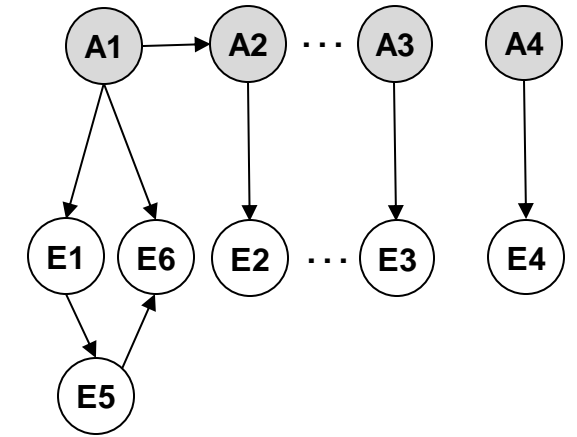- Most interesting when $Q$ *sparse*

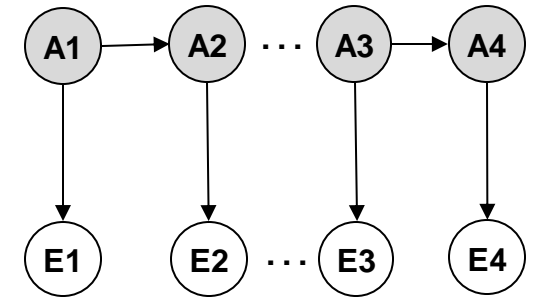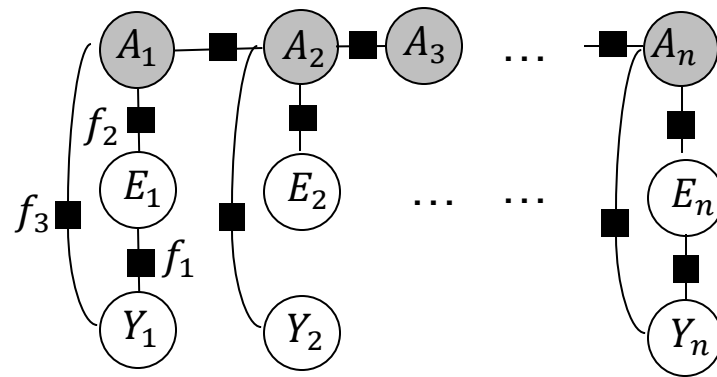# Modeling the credential stealing attack using Naïve Bayes vs. Bayesian Network



**Naïve Bayes**

**Bayesian Network**

**Dynamic Bayesian Network**

**Factor Graphs**

**Hidden Markov Model**
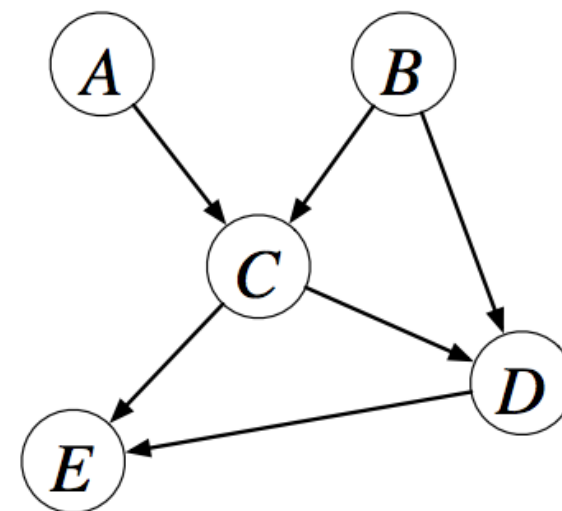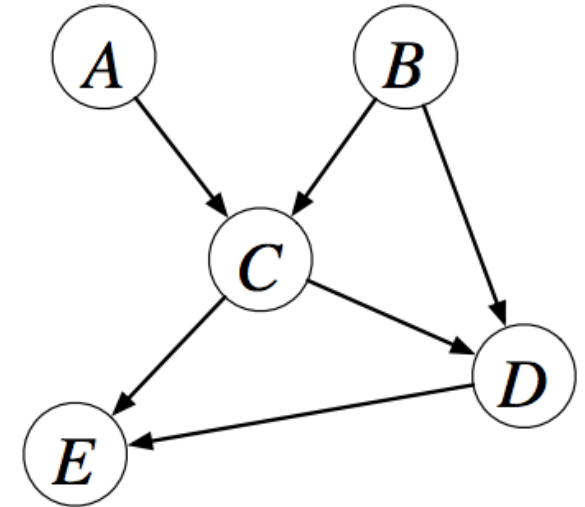
# Representing knowledge through graphical models

- A PGM encodes structural aspects of a joint probability distribution
  - G = {V,E}
- A node corresponds to a random variable
- An edge represent a dependencies between the variables
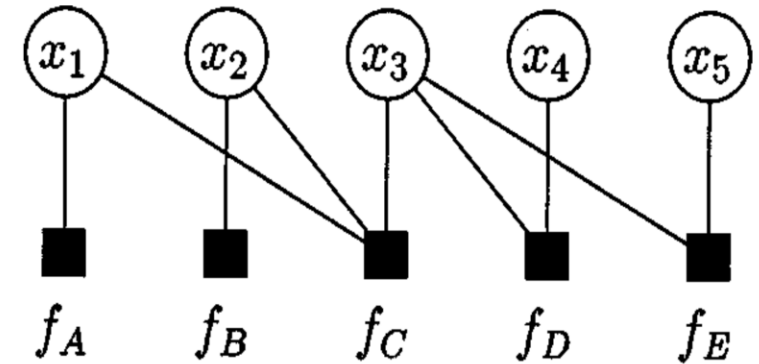
# Why do we need graphical models?

- Graphs are an intuitive way of visualizing relationship among variables
- A graph shows the conditional independence between variables via edges
- Effective inference algorithms can be run on graphs such as belief propagation to infer marginal probabilities of variables
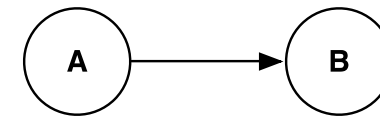
# Definition of a Factor Graph

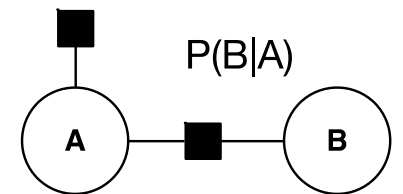A factor graph is a **bipartite, undirected graph** of **random variables and factor functions.** [Frey et. al. 01]

A factor function is a mathematical definition of *prior beliefs* or expert knowledge. *FG can represent both causal and non-causal relations.*
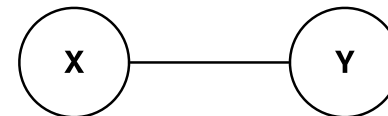


A factor graph for the product $f_A(x_1) f_B(x_2) f_C(x_1, x_2, x_3)$ $\cdot f_D(x_3, x_4) f_E(x_3, x_5)$.
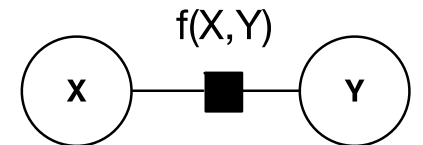


Bayesian Network (BN)

Factor Graph equivalent of BN



Markov Random Fields (MRF)

Factor Graph equivalent of MRF

# Applications of Probabilistic Graphs in Security Domain

**Problem statement.** Given a set of security events, infer whether an attack is in progress?

**Modeling Approach.**

Each security event is a known variable **e**, each takes value from a discrete set of events **E**.

An attack happens in a chain of exploits, thus we have a sequence of events in time dimension.

Each event is associated with a corresponding attack state **s,** which is unknown. The simplest approach is to classify **s** as a binary {0,1}. However, when we can infer **s** it is often too late (the attacker is already in the system)

Thus, we want to discretize **s** to smaller attack stages and provide update on such stages as soon as an event is observed.

# Applications in the Security Domain (cont.)

**Problem statement.** Given a set of security events, infer whether an attack is in progress?

Formally, the problem becomes

1. Define a joint probability distribution function (joint pdf)

$P(e_1,e_2,..,e_n,s_1,s_2,...,s_n)$

2. Derive a conditional probability

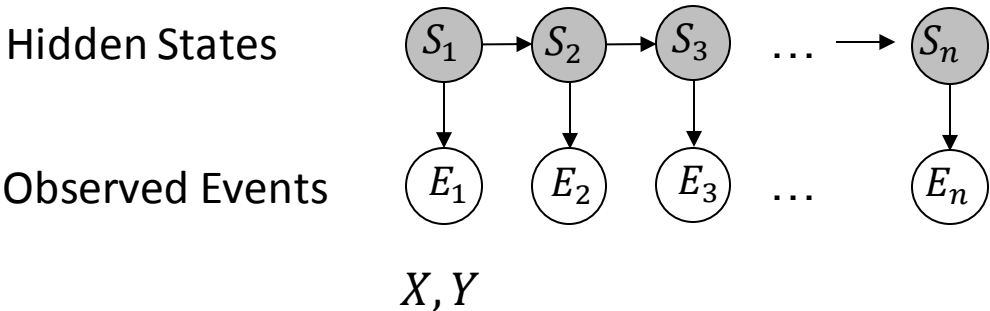$P(e_1,e_2,..,e_n|s_1,s_2,...,s_n)$

However, the search space is exponentially large (by the order of the number of observed stages and events) and the joint pdf is sophisticated.

We want to break the joint pdf into smaller components that are easier to compute, i.e., factorize the joint pdf.

# Underlying representation of a Hidden Markov Model and conversion to a Factor Graph
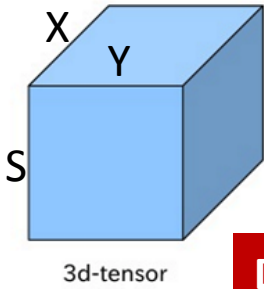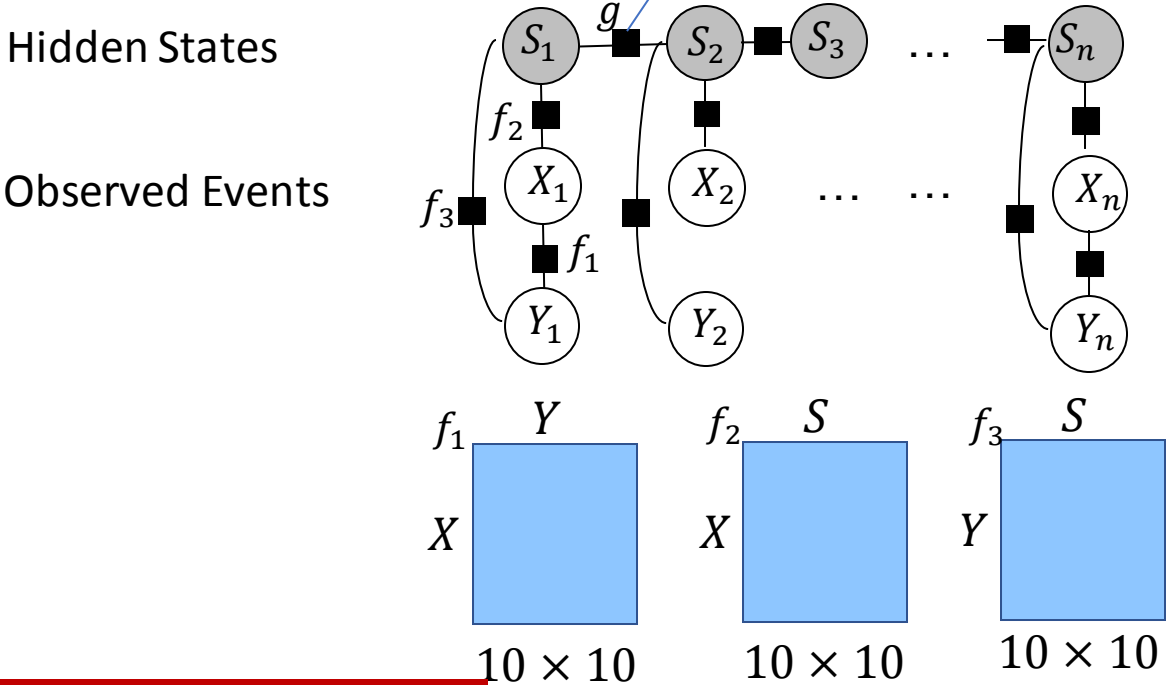
**Hidden Markov Model**

**Factor Graph of the HMM**

Hidden States

Observed Events

$X, Y$

Example
$|S| = 10$
$|X| = 10$
$|Y| = 10$

X
Y
S

3d-tensor

**Domain knowledge:** *"variables are pair-wise related"* **reduces dimensionality**

*size of tensor*
$10 \times 10 \times 10 = 1000$

$1000^n \gg 400 \times n$

$S_{t+1}$
$S_t$

$g$

$10 \times 10$

Hidden States

Observed Events

$f_1$ $Y$
$X$
$10 \times 10$

$f_2$ $S$
$X$
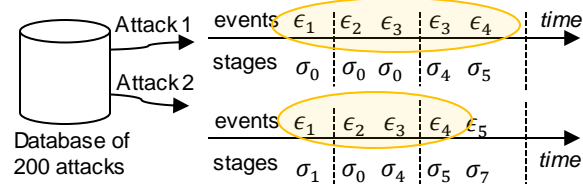$10 \times 10$

$f_3$ $S$
$Y$
$10 \times 10$

*size of three matrices + one transition*
$10 \times 10 + 10 \times 10 + 10 \times 10 + 10x10 = 400$

# Modeling the credential stealing attack using Factor Graphs

## OFFLINE ANNOTATION ON PAST ATTACKS

a) Annotated events and attack stages in a pair of attacks



Database of 200 attacks

b) Event-stage annotation table for the attack pair (Attack 1 and Attack 2)

| Event | Attack stage |
|---|---|
| $\{\epsilon_1\}$ | $\{\sigma_0|\sigma_1\}$ |
| $\{\epsilon_2\}$ | $\{\sigma_0\}$ |
| $\{\epsilon_3\}$ | $\{\sigma_4\}$ |
| $\{\epsilon_4\}$ | $\{\sigma_5\}$ |
| $\{\epsilon_5\}$ | $\{\sigma_7\}$ |

## OFFLINE LEARNING OF PATTERNS

c) Example patterns, stages, probabilities, and significance learned from the attack pair

| Pattern | Attack stages | Probability in past attacks | Significance (p-value) |
|---|---|---|---|
| $[\epsilon_1, \epsilon_3, \epsilon_4]$ | $[\sigma_1, \sigma_4, \sigma_5]$ | $q_a$ | $p_a$ |
| $[\epsilon_1]$ | $[\sigma_0|\sigma_1]$ | $q_b$ | $p_b$ |

...

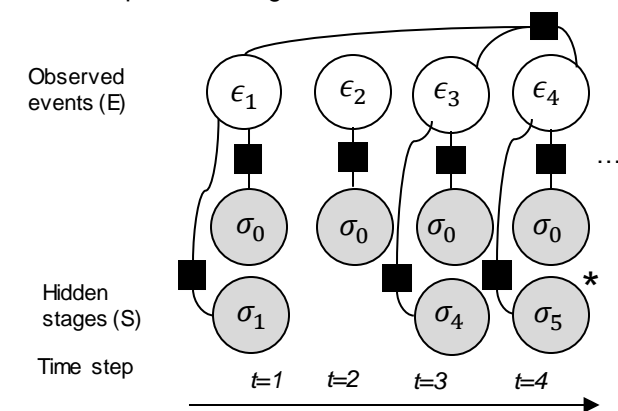$$f(E) = \exp\{q_E(1 - p_E)\}$$

A factor function defined on the learned pattern, stages, and its significance

## Model assumptions

1. There are multivariate relationships among the events
2. Such relationships are represented by factor functions
3. There is no restriction on order of the relationships like causal in Bayesian Network

More suitable for modeling highly complex attacks, where the causal relations among the events are not immediately clear.
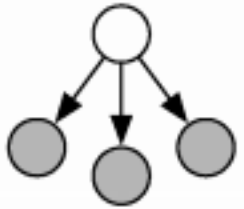
## RUNTIME DETECTION OF UNSEEN ATTACKS

d) An evolution of the Factor Graph for the port knocking attack at run-time
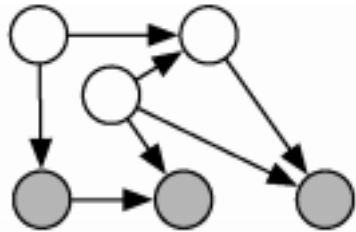


Observed events (E)

Hidden stages (S)

Time step    t=1    t=2    t=3    t=4

| | Observed Security events | ■ | Factor function | $\epsilon_1$ vulnerability scan | $\sigma_0$ benign |
|---|---|---|---|---|---|
| | Unknown attack stages | * | Attack detected and stopped before the system misuse | $\epsilon_2$ login | $\sigma_1$ discovery |
| | | | | $\epsilon_3$ sensitive_uri | $\sigma_4$ privilege escalation |
| | | | | $\epsilon_4$ new_library | $\sigma_5$ persistence |

ECE ILLINOIS          ILLINOIS
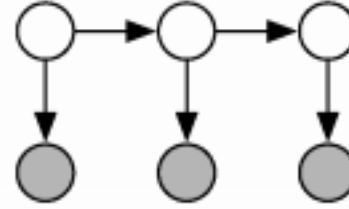
# Taxonomy of graphical models
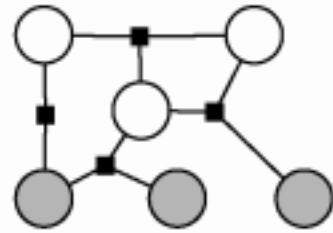


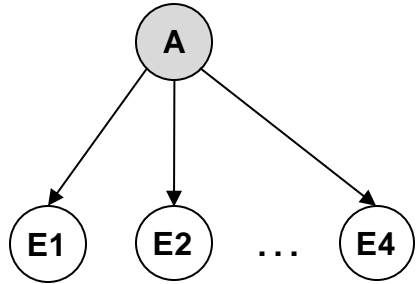Naïve Bayes        Bayesian Network        Hidden Markov Model        Factor Graph

Conditional probabilities and statistical dependencies can be represented by a general type of graph: Factor Graph
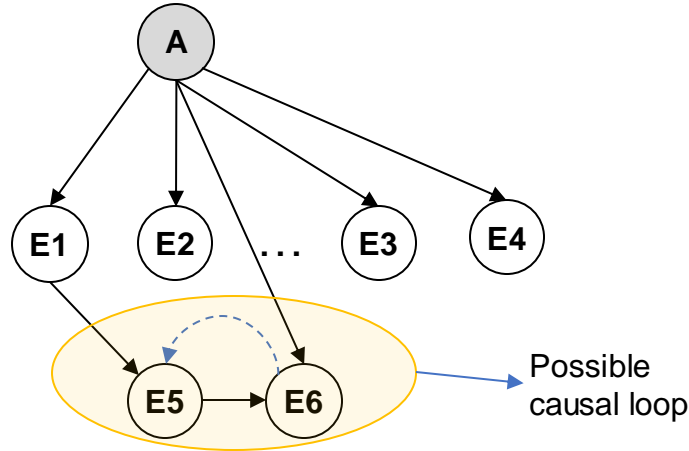
# Model structure and inference in PGMs

| | | Naïve Bayes | Bayesian Network | Hidden Markov Model | Factor Graphs |
|---|---|---|---|---|---|
| *MODEL STRUCTURE* | *Graph type* | Directed | Directed | Directed | **Undirected** |
| | *Graph structure* | Parent-child | Hierarchical parent-child | Sequential | **Arbitrary structure** |
| | *Variable of interest* | Attack (0 or 1) | Attack (0 or 1) | Sequence of system states | **Sequence of attack stages** |
| | *Relationship* | Conditional independence | Prior Conditional independence | State transitions Emission of event | **Temporal relationships (patterns of events)** **Statistical relationships (severity or repetitiveness of events)** |
| *INFERENCE* | *Algorithm* | Multiplication of conditional probabilities | Multiplication of conditional probabilities and priors | Dynamic Programming | **Belief Propagation Sampling** |

# Modeling the credential stealing attack using Naïve Bayes vs. Bayesian Network



**Naïve Bayes**

**Bayesian Network**

Possible causal loop

| ID | Description |
|----|-------------|
| A | Attack |
| E1 | Vulnerability scan |
| E2 | Login |
| E3 | Download file with sensitive extension |
| E4 | New system library |
| E5 | High network flows |
| E6 | Service unavailable |

**Model assumptions**
1. All events share the same parent variable
2. All events are conditionally independent

**Advantage:**
Simplify calculation of posterior probability on A

**Model assumptions**
1. An event can be preceded (causal) by another event
2. There is no cycle in the network

**Disadvantage**
Explicitly assume causal relationships
(Causality may not be clear from the data)
For complicated attacks, causal loops may form and render the BN invalid

**ECE ILLINOIS**          **I ILLINOIS**

# Bayesian Networks vs. Markov Random Fields vs. Factor Graphs
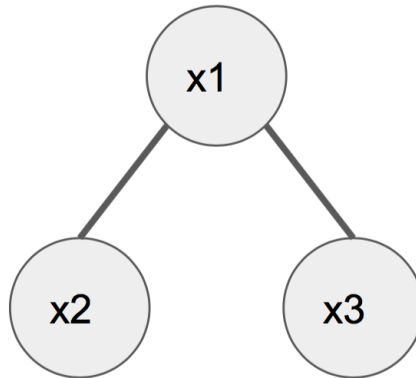


**Bayesian networks**

$$p(x_1)p(x_2|x_1)p(x_3|x_1)$$

Product of conditional probabilities
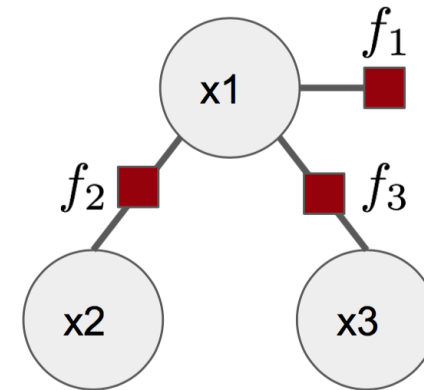
Causal relationships

**Markov random fields**

$$\frac{1}{Z}\psi_1(x_1, x_2)\psi_2(x_1, x_3)$$

Product of dependencies among variable cliques

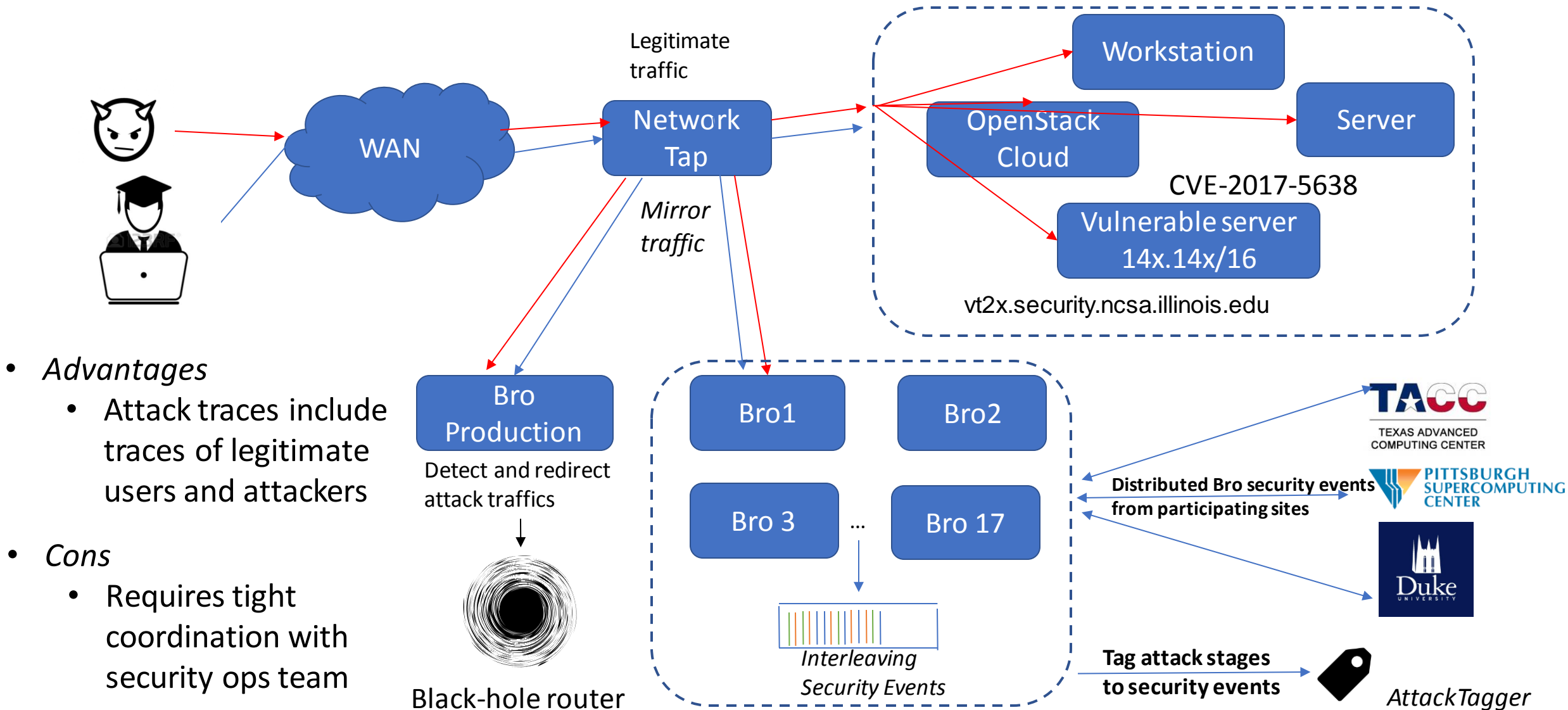Statistical dependencies

**Factor graph**

$$\frac{1}{Z}f_1(x_1)f_2(x_2, x_1)f_3(x_1, x_3)$$

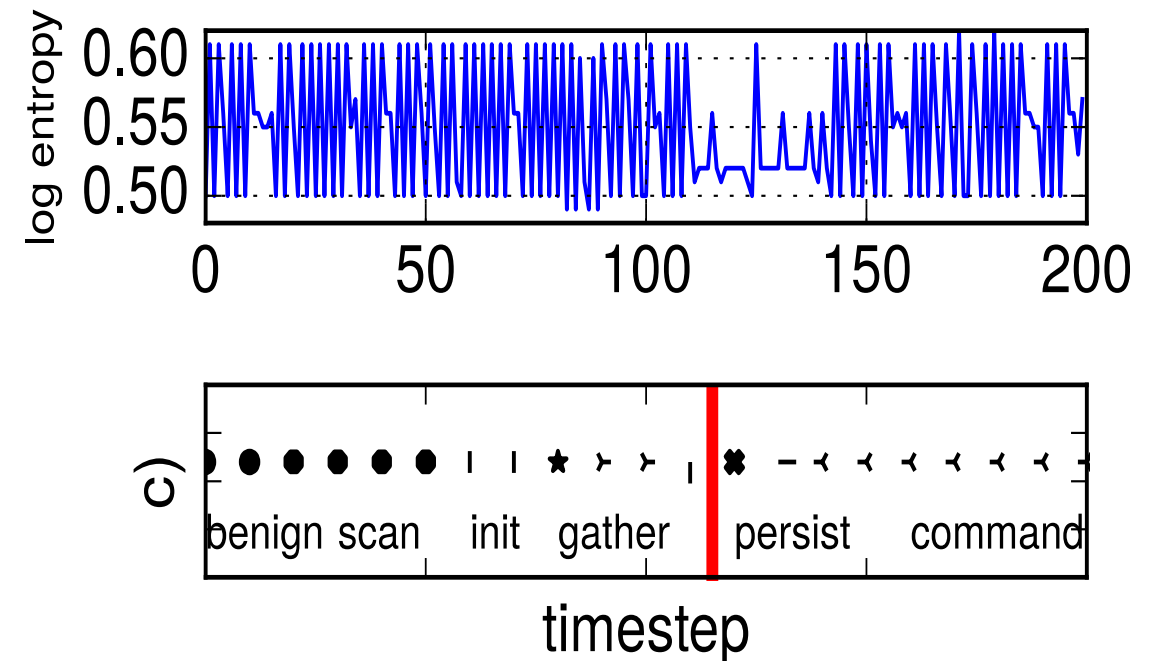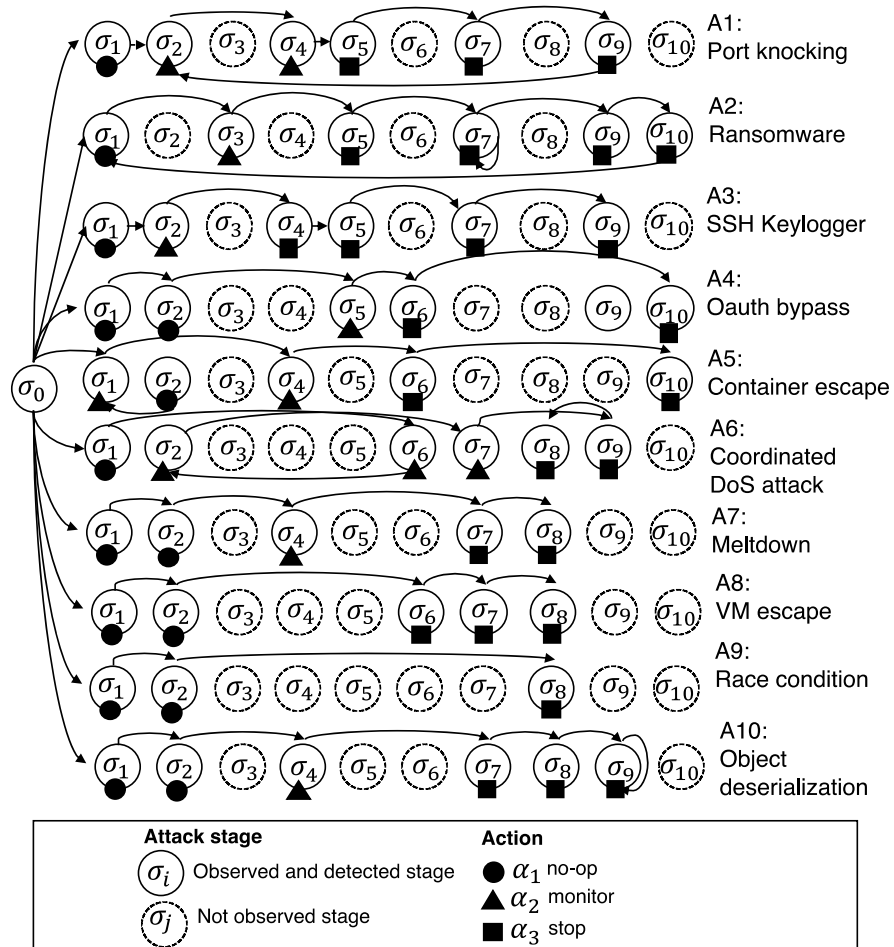Product of dependencies using univariate, bivariate, or multivariate functions

Both types of relations (including prior on a variable)

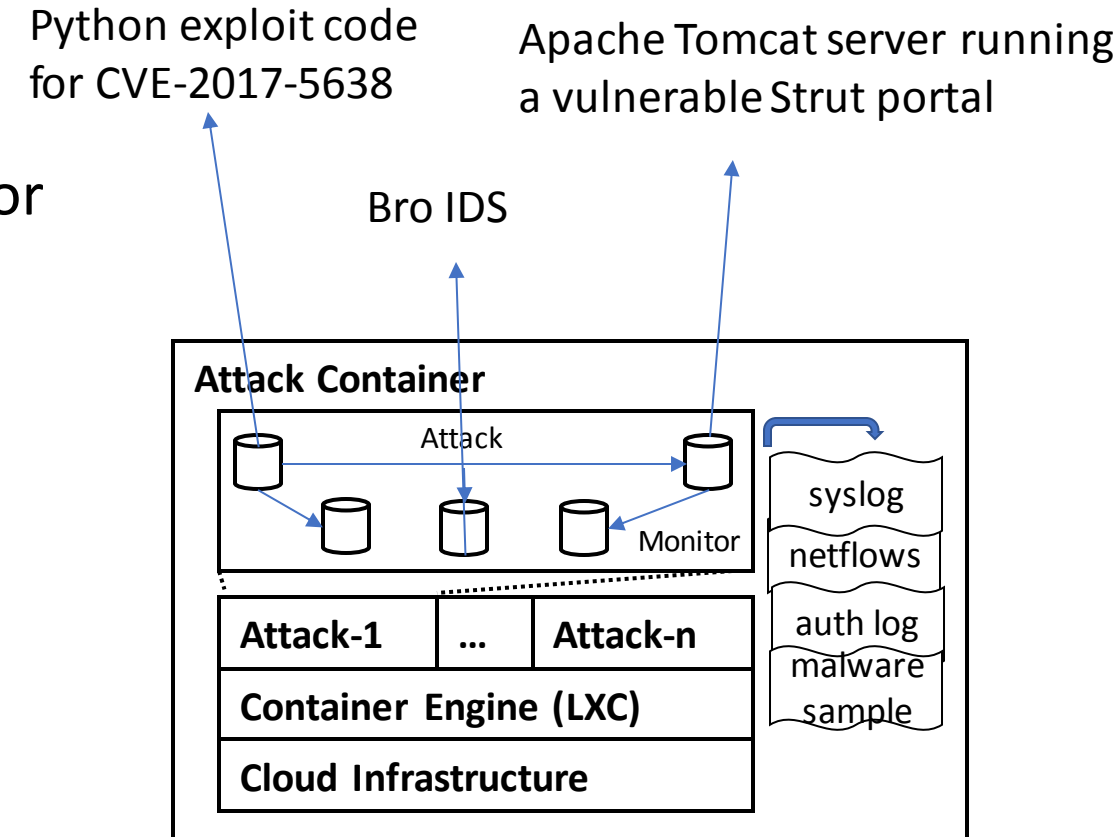# An attack testbed in real production traffic – an experiment at NCSA



- *Advantages*
  - Attack traces include traces of legitimate users and attackers

- *Cons*
  - Requires tight coordination with security ops team

Legitimate traffic

WAN

Network Tap

*Mirror traffic*

Workstation

OpenStack Cloud

Server

CVE-2017-5638

Vulnerable server 14x.14x/16

vt2x.security.ncsa.illinois.edu

Bro Production

Detect and redirect attack traffics

Black-hole router

Bro1   Bro2

Bro 3   …   Bro 17

*Interleaving Security Events*

TACC
TEXAS ADVANCED COMPUTING CENTER

**Distributed Bro security events from participating sites**

PITTSBURGH SUPERCOMPUTING CENTER

Duke UNIVERSITY

**Tag attack stages to security events**

*AttackTagger*

# Evaluation Result

# Stage transition of a multi-stage attack that exploits CVE-2017-5638
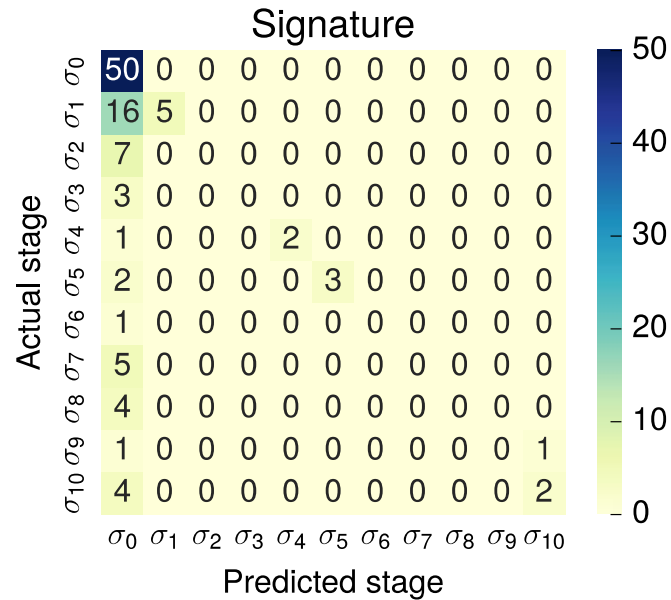
# Emulating CVE-2017-5638 in a container-based environment

- *Advantages*
  - We were able to create an exact environment for the vulnerable Strut application
  - Monitors are in place to collect attack traces
  - Network policies are implemented to isolate potential outbreak of the attack
- *Limitations*
  - Containers are not exposed to a real network thus are not visible to attackers
  - Traces only include attack activities

Python exploit code for CVE-2017-5638

Apache Tomcat server running a vulnerable Strut portal

Bro IDS

**Attack Container**

Attack

Monitor

syslog

netflows

auth log

malware sample

| **Attack-1** | **...** | **Attack-n** |
|---|---|---|

**Container Engine (LXC)**

**Cloud Infrastructure**

# Evaluation Results



Signature / Anomaly / Factor Graph confusion matrices (Actual stage vs Predicted stage)

# Concluding Remarks

**1. Probabilistic Graphical Models appear to be the way to integrate disparate issues on failure and attack pre-emption**

**2. Continuous and dynamic monitoring and adaptive abstraction offered by the factor graph based learning is critical**

**3. Going forward: Factor graphs could combine both security logs and error logs for diagnosis**

ECE ILLINOIS        ILLINOIS

**6000+**
users

**5+ millions**
connections

**34M+**
log events

**4.5+ GB**
Compressed final log

Heterogeneous host and network logs
- Syslog
- Netflows
- IDS alerts
- Human-written reports

200+ incidents in the past years (2008-2017)

Brute-force attacks

Credential compromise

Abusing computing infrastructure
- Send spam
- Launch Denial of

Service attacks.

5-minute snapshot of network traffic in and out of NCSA

# Why attack injection?

- Vulnerabilities are discovered on a daily basis, however, is a target system immune from such vulnerabilities?

- **Our goals are to:**

  - Evaluate ability of security monitoring systems in capturing attack-related security events

  - Run live, integration tests on applied security patches

  - Provide a dynamic blueprint of an attack (in terms of attack stages) as the attack unfolds across a production network

ECE ILLINOIS

I ILLINOIS

# What is a Linux Container (LXC)?



**Virtual Machine (VM)** is an efficient, isolated duplicate of a real computer machine.

| Features | Virtual Machine | Linux Container |
|---|---|---|
| Emulation | A real machine | **A Linux system** |
| Guest OS | Almost any OS | **Only Linux system** |
| Isolation and Resource management | Fully virtualized | **Kernel namespace and control groups** |
| File system | Separated file system for each VM | **Layered filesystem (AUFS)** |
| Disk and Memory | GBs | **MBs** |
| Startup time | Minutes | **Seconds** |



**Linux Container (LXC)** is a virtualization technology for running multiple isolated Linux systems (containers).

# How does AttackTagger work?