# Multimodal Machine Learning

## Lecture 7.2: Generative Models

**Louis-Philippe Morency**

\* Original version co-developed with Tadas Baltrusaitis

# Outline

- Probabilistic graphical models
  - Joint probabilistic distribution
  - Example: creating a graphical model
- Bayesian networks
  - Conditional probability distribution
  - Dynamic Bayesian Network
- Variational Auto Encoder
- Generative Adversarial Network
  - cGAN, infoGAN, cycleGAN

# Probabilistic Graphical Models

# Probabilistic Graphical Model

**Definition:** A probabilistic graphical model (PGM) is a graph formalism for compactly modeling joint probability distributions and dependence structures over a set of random variables.

- Random variables: $X_1,\ldots,X_n$
- P is a joint distribution over $X_1,\ldots,X_n$

Why do we want to learn the joint distribution?

# Inference for Known Joint Probability Distribution

When we know the joint probability distribution :

$$P(A, B, C, D, E)$$ ➡ If A, B C, D and E are discrete variables, then P(A,B,C,D, E) will be a 5-D tensor (matrix)

Two main forms of inference:

①   Joint probability for a particular assignment

$$P(A = 1, B = 'car', C = 2, D = 'banana', E = 10)$$

➡ A specific entry in the 5-D tensor

# Inference for Known Joint Probability Distribution

**②** Probability of a subset of variables (query) given known assignments of other variables (evidences)

$$P(A, D|C = 3)$$  ➡️  Use the product rule to *marginalize* the other variables B and E

$$P(A, D|C = 3) = \sum_{\forall b \in B, e \in E} P(A, D, b, e|C = 3)$$

➡️  Use the inverse of product rule $P(X|Y) = P(X, Y)/P(Y)$

$$P(A, D|C = 3) = \frac{1}{P(C)} \sum_{\forall b \in B, e \in E} P(A, D, b, e|C = 3)$$

# Inference for Known Joint Probability Distribution

**(2)** Probability of a subset of variables (query) given known assignments of other variables (evidences)

$$P(x|y) = \alpha \sum_{\forall z \in Z} P(x, y, z)$$

where $x$ is the subset of query variables

$y$ is the subset of evidence assignments

$Z$ is the set of all other variables (not in $x$ or $y$)

Can we represent P more compactly?

- Key: Exploit independence properties

# Independent Random Variables

- Two variables X and Y are independent if
  - $P(X=x|Y=y) = P(X=x)$ for all values x,y
  - Equivalently, knowing Y does not change predictions of X
- If X and Y are independent then:
  - $P(X, Y) = P(X|Y)P(Y) = P(X)P(Y)$
- If $X_1,\ldots,X_n$ are independent then:
  - $P(X_1,\ldots,X_n) = P(X_1)\ldots P(X_n)$

X        Y
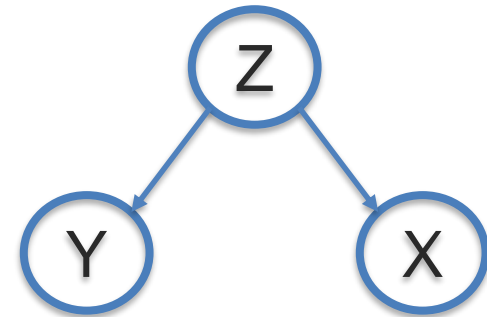
# Conditional Independence

- X and Y are conditionally independent given Z if
    - P(X=x|Y=y, Z=z) = P(X=x|Z=z) for all values x, y, z
    - Equivalently, if we know Z, then knowing Y does not change predictions of X
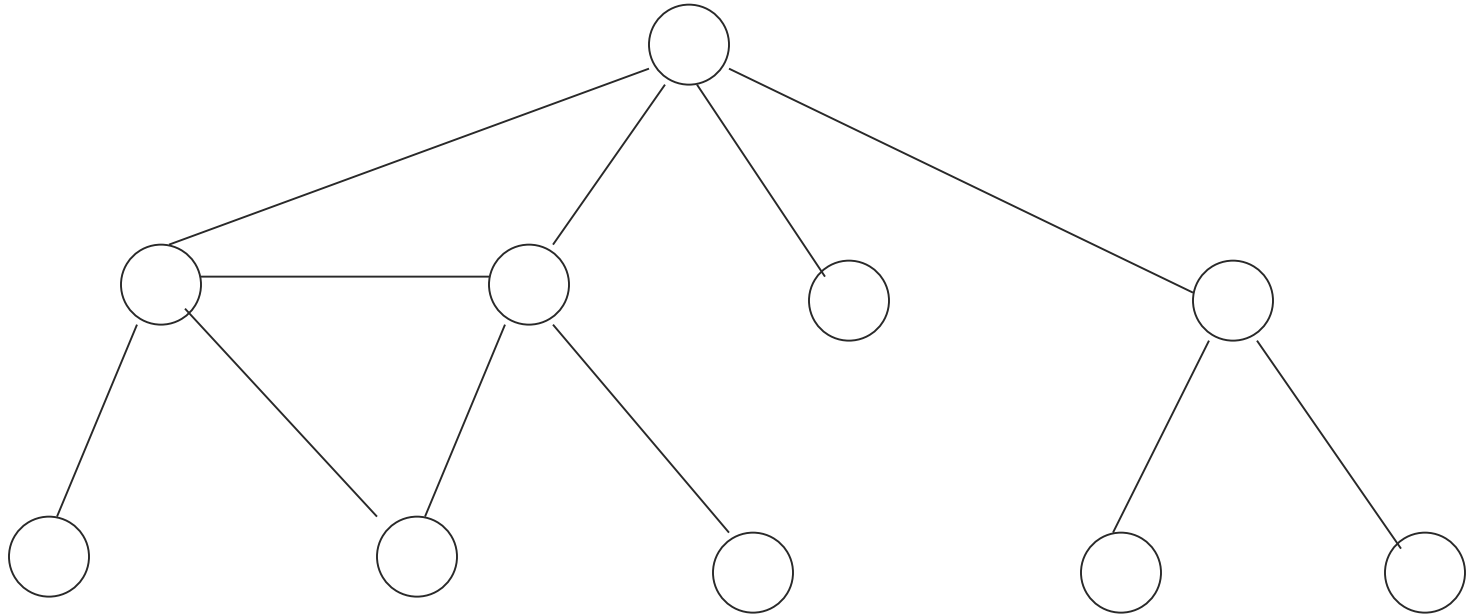
# Graphical Model

- A tool that visually illustrate <u>conditional independence</u> among variables in a given problem.

- Consisting of nodes (Random variables or States) and edges (Connecting two nodes, directed or undirected).

- The lack of edge represents conditional independence between variables.
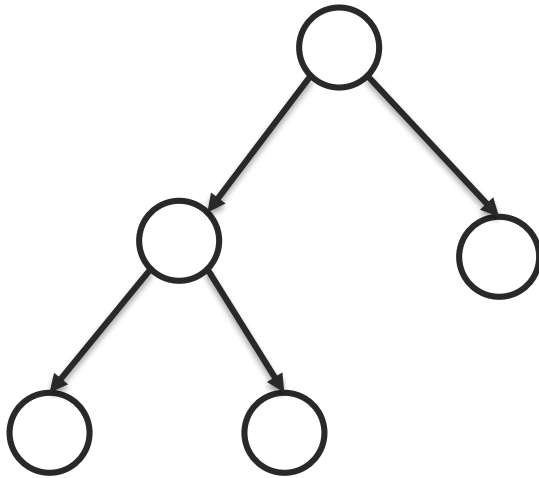
# Graphical Model



Different types of graphical models:

- ■ Chain, Path, Cycle, Directed Acyclic Graph (DAG), Parents and Children

# Two Main Types of Graphical Models

## Bayesian networks

## Markov Models (in 2 weeks)

- Directed acyclic graph
- Conditional dependencies

- Undirected graphical model
- Cyclic dependencies

# Creating a Graphical Model

# Example: Inferring Emotion from Interaction Logs

[Sabourin et al., 2011]

**Student**

**Tutoring System**

Student Traits

Anxious
Bored
Confused
Curious
Excited
Focused
Frustrated

**Emotion?**

Logs

Language Technologies Institute

Carnegie Mellon University

# Example: Bayesian Network Representation

[Sabourin et al., 2011]

**Outcome** (non-observable)

Emotion

- Anxious
- Bored
- Confused
- Curious
- Excited
- Focused
- Frustrated

**Evidences** (observable)

| # book views | # correct ans. |
| # notes taken | # incorrect ans. |
| # poster views | Total goals |

**Observable environment variables**

| Openness | Mastery avoidance |
| Agreeableness | |
| Conscientious | Mastery approach |

**Survey-based personality variables**

Language Technologies Institute

Carnegie Mellon University

# Example: Naïve Bayes Approach

[Sabourin et al., 2011]



Outcome (non-observable)

Evidences (observable)

Emotion

# book views
# notes taken
# poster views
# correct ans.
# incorrect ans.
Total goals

Openness
Agreeableness
Conscientious
Mastery avoidance
Mastery approach

**Observable environment variables**

**Survey-based personality variables**

# Appraisal Theory of Emotion

World Events

Metal State (beliefs, goals)

Body

Expression

Action tendency

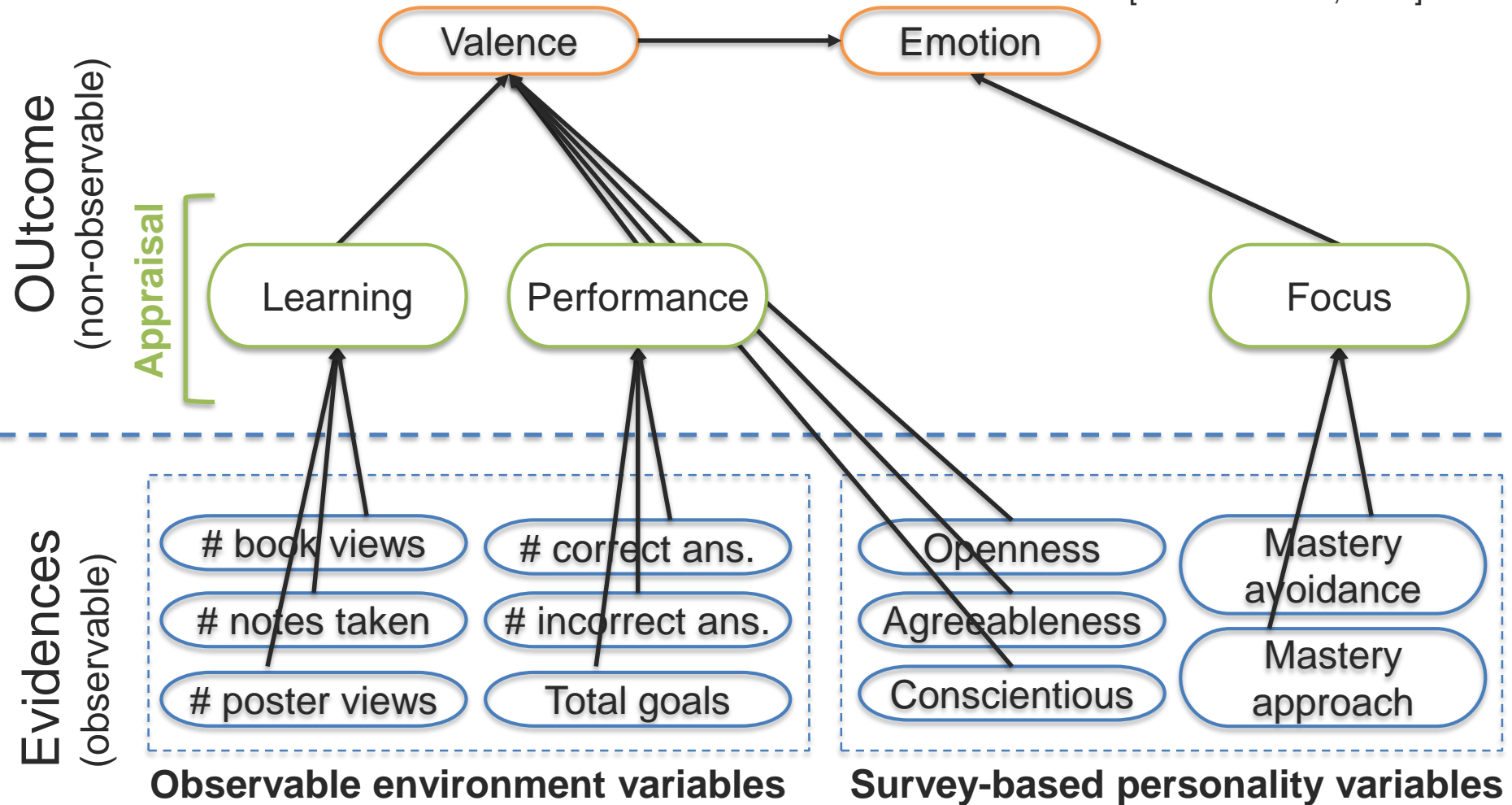Physiological response

Argues for importance of three interrelated concepts
- World events
- Mental state
- Emotional Response

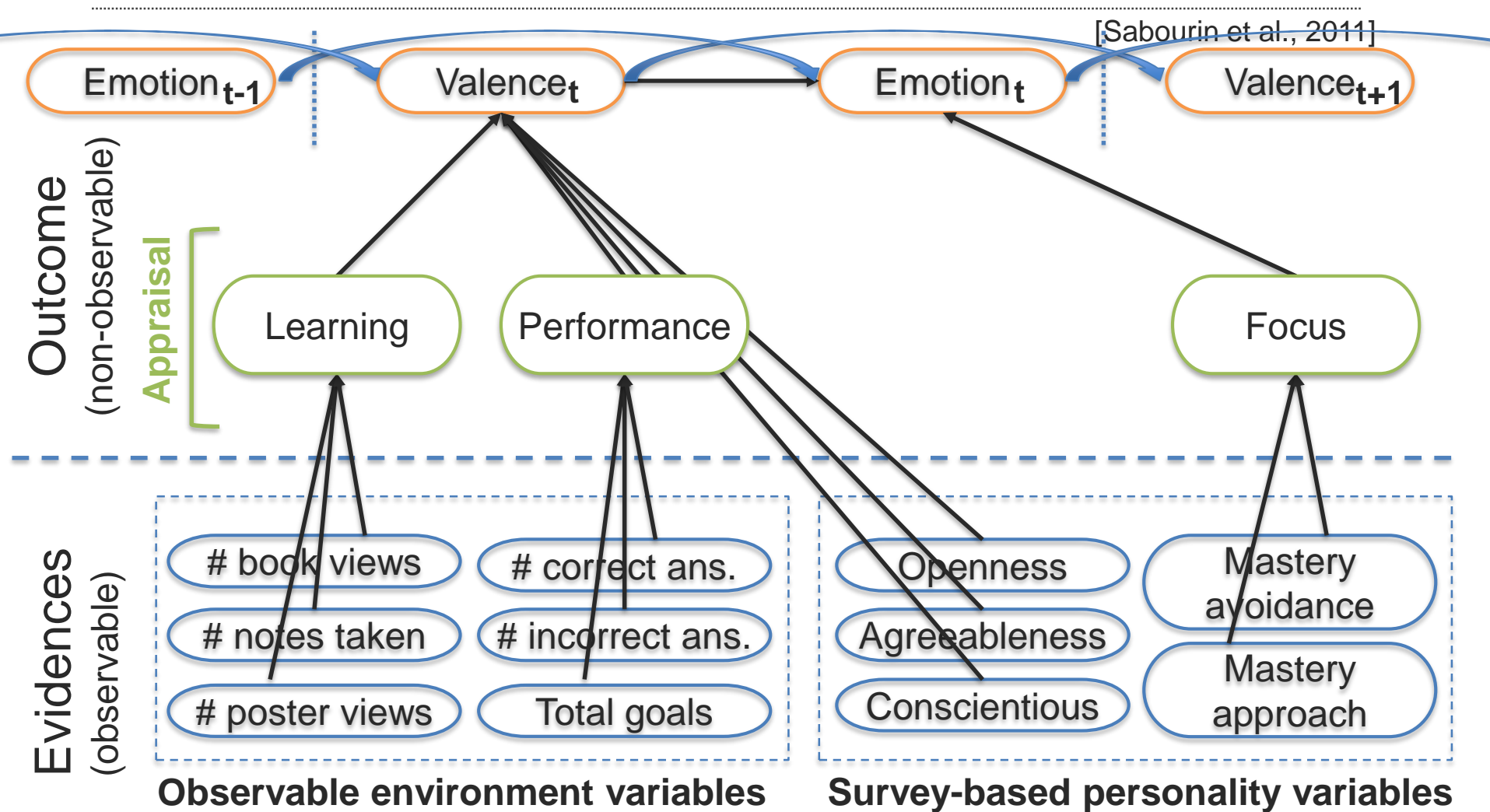If we know two of these variables, we can make predictions about the third

Response= f(Env., Mind)

Language Technologies Institute

Carnegie Mellon University

# Example: Bayesian Network Approach

[Sabourin et al., 2011]

OUtcome (non-observable)

Appraisal

Evidences (observable)

Valence → Emotion

Learning   Performance   Focus

# book views
# notes taken
# poster views

# correct ans.
# incorrect ans.
Total goals

Openness
Agreeableness
Conscientious

Mastery avoidance
Mastery approach

**Observable environment variables**    **Survey-based personality variables**

Language Technologies Institute

Carnegie Mellon University

# Example: Dynamic Bayesian Network Approach



[Sabourin et al., 2011]

Emotion$_{t-1}$  Valence$_t$  Emotion$_t$  Valence$_{t+1}$

**Outcome** (non-observable)

**Appraisal**

Learning    Performance    Focus

**Evidences** (observable)

# book views    # correct ans.    Openness    Mastery avoidance

# notes taken    # incorrect ans.    Agreeableness    Mastery approach

# poster views    Total goals    Conscientious

**Observable environment variables**    **Survey-based personality variables**

Language Technologies Institute

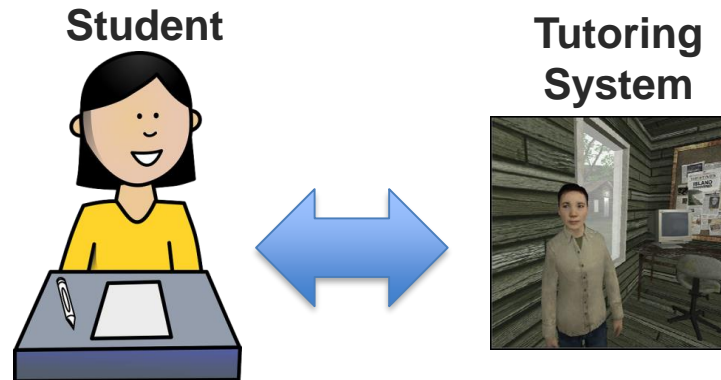Carnegie Mellon University

# Example: Dynamic Bayesian Network Approach

[Sabourin et al., 2011]

# Example: Inferring Emotion from Interaction Logs

[Sabourin et al., 2011]

**Student**

**Tutoring System**

|  | Emotion Accuracy | Valence Accuracy |
|---|---|---|
| **Baseline** | 22.4% | 54.5% |
| **Naïve Bayes** | 18.1% | 51.2% |
| **Bayes Net** | 25.5% | 66.8% |
| **Dynamic BN** | 32.6% | 72.6% |

Carnegie Mellon University

# Bayesian Networks
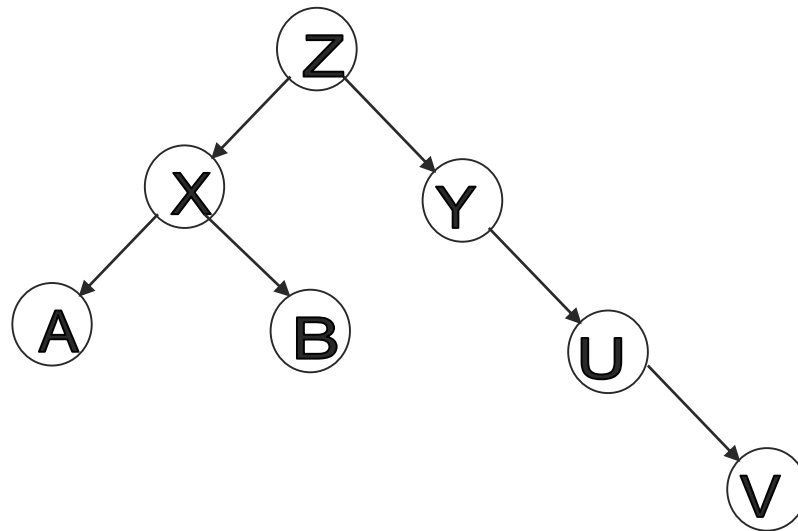
# Bayesian networks

- A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions

- Syntax:
    - a set of nodes, one per variable
    - a directed, acyclic graph (link ≈ "directly influences")
    - a conditional distribution for each node given its parents:
        $$\mathbf{P}\,(X_i \mid \text{Parents}\,(X_i))$$

- In the simplest case, conditional distribution represented as a conditional probability distribution (CPD) giving the distribution over $X_i$ for each combination of parent values

# Bayesian Network (BN)

- A specific type of graphical model that is represented as a Directed Acyclic Graph.
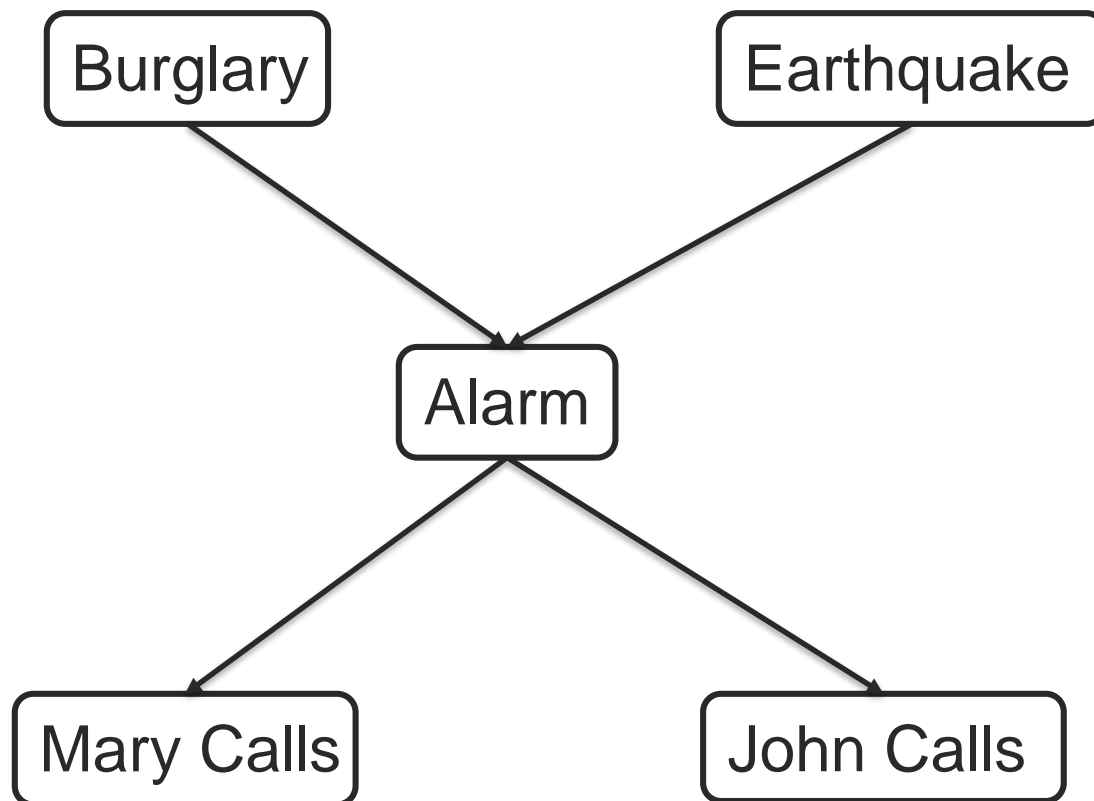
# Example

*"I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?"*

- Variables?
  - *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

- "Causal" knowledge?
  - 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

# Example – Network Topology

Language Technologies Institute

Carnegie Mellon University

# Joint Probability in Graphical Models

With chain-rule, the joint probability can be restated:

$$P(A, B, C, D, E) = P(A|B, C, D, E)P(B, C, D, E)$$

$$= P(A|B, C, D, E)P(B|C, D, E)P(C|D, E)$$

$$= P(A|B, C, D, E)P(B|C, D, E)P(C, D, E)$$

$$= P(A|B, C, D, E)P(B|C, D, E)P(C|D, E)P(D, E)$$

$$= P(A|B, C, D, E)P(B|C, D, E)P(C|D, E)P(D|E)P(E)$$

⟹  The order in applying the chain-rule is arbitrary.

How can we simplify the joint probability even more, given the graphical model?
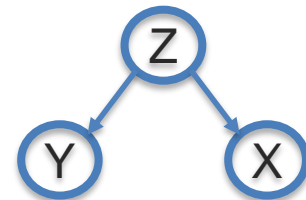
# Joint Probability in Graphical Models

With chain-rule, the joint probability can be reshaped:

$$P(A, B, C, D, E) = P(A|B, C, D, E)P(B|C, D, E)P(C|D, E)P(D|E)P(E)$$

➡ Remember these concepts:

$X$  $Y$

Independent variables

$Z$
$Y$  $X$

conditionally independent

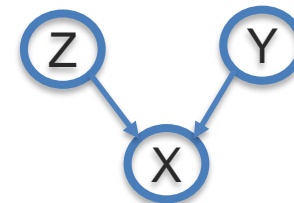➡ In a Bayesian network, each conditional probability for a specific variable X only depends on its parents:

$$P(X|\ all\ variables) = P(X|parents(X))$$

Conditional Probability Distribution (CPD)

# Conditional Probability Distribution (CPD)

Given a variable X and its parents (Y and Z):

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

**Definition:** probability distribution of X when the assignment of it parents is known (Y and Z)

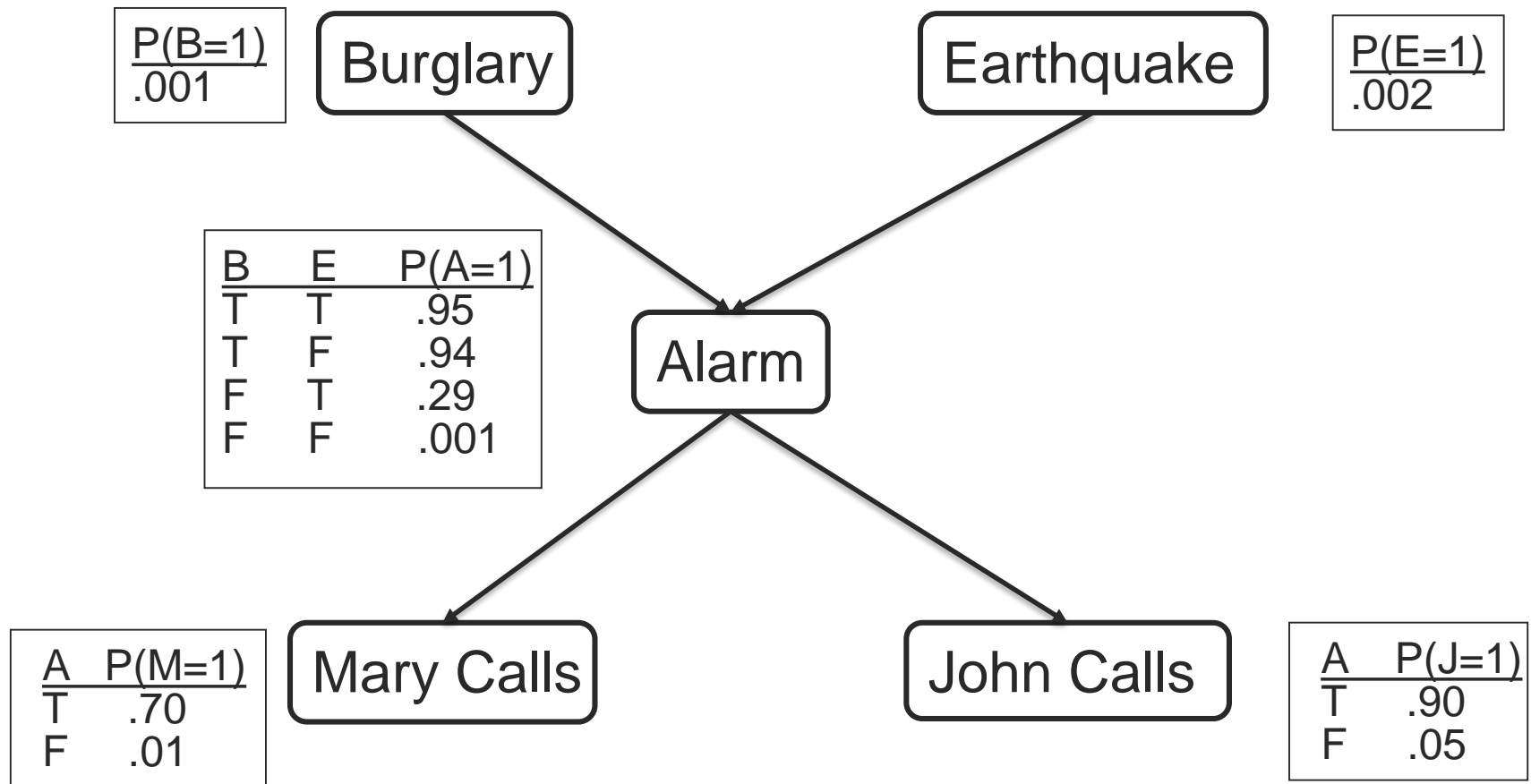❑ For **categorical variable**: expressed as a conditional probability table

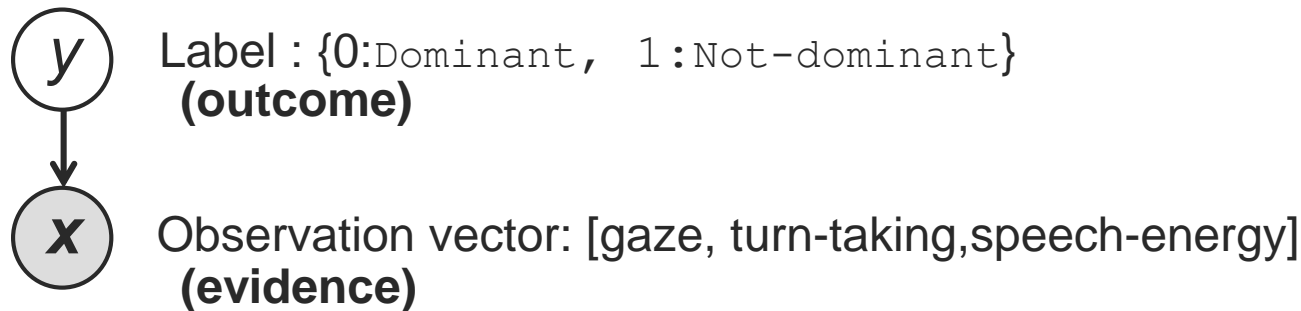|          | Y=0 | Y=1 |
|----------|-----|-----|
| P(X=0|Y) | 4/6 | 1/3 |
| P(X=1|Y) | 2/6 | 2/3 |

❑ For **continuous variable**: expressed as a conditional density function

- For example, multivariate normal density function or Gaussian linear regression (used by Bayes RegressionLinear Model)

Language Technologies Institute

Carnegie Mellon University

# Example – Conditional Probability Distributions



P(B=1)
.001

Burglary

Earthquake

P(E=1)
.002

| B | E | P(A=1) |
|---|---|--------|
| T | T | .95 |
| T | F | .94 |
| F | T | .29 |
| F | F | .001 |

Alarm

Mary Calls

John Calls

| A | P(M=1) |
|---|--------|
| T | .70 |
| F | .01 |

| A | P(J=1) |
|---|--------|
| T | .90 |
| F | .05 |

Language Technologies Institute

Carnegie Mellon University

# Generative Model: Naïve Bayes Classifier

$y$  Label : {0:Dominant, 1:Not-dominant}
   **(outcome)**

$x$  Observation vector: [gaze, turn-taking,speech-energy]
   **(evidence)**

**Score function:** $P(y = a | \boldsymbol{x_i})$

Likelihood                    Prior          Chain rule

Bayes' theorem: $P(y|\boldsymbol{x}) = \dfrac{P(\boldsymbol{x}|y)P(y)}{P(\boldsymbol{x})} \approx P(\boldsymbol{x}|y)P(y) = P(\boldsymbol{x}, y)$

Posterior

Marginal likelihood
(partition)   $P(\boldsymbol{x}) = \sum_{y} P(\boldsymbol{x}|y)P(y)$
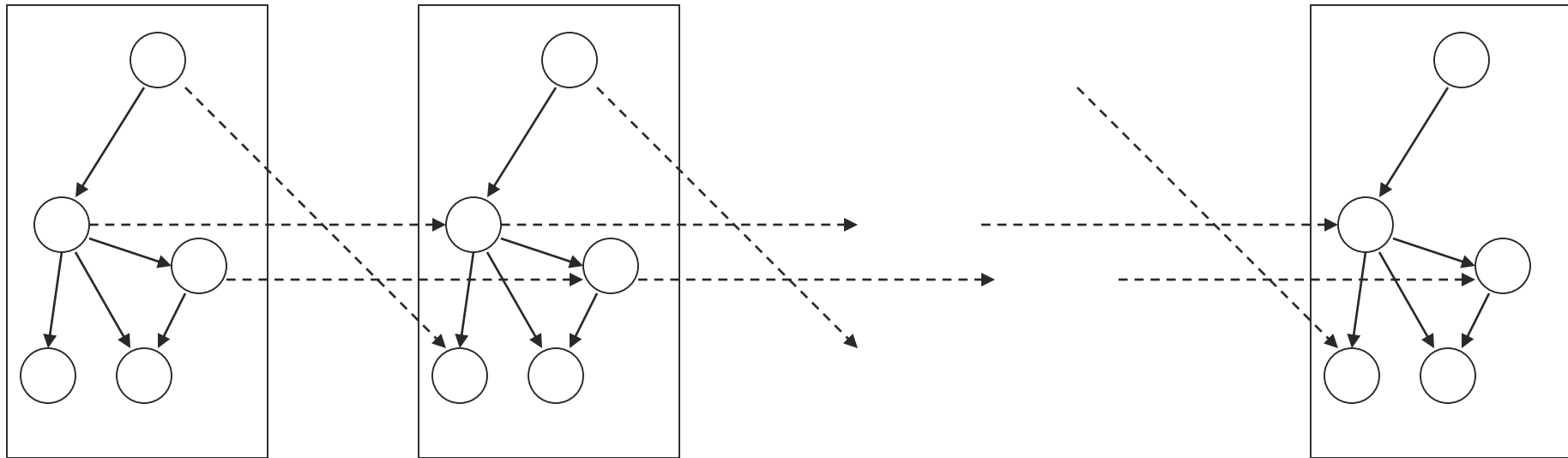
# Dynamic Bayesian Network

# Dynamic Bayesian Network (DBN)

- Bayesian network with time-series to represent temporal dependencies.

- Dynamically changing or evolving over time.

- Directed graphical model of stochastic processes.

- Especially aiming at time series modeling.


- Satisfying the Markovian condition:

  *The state of a system at time t depends only on its immediate past state at time t-1.*
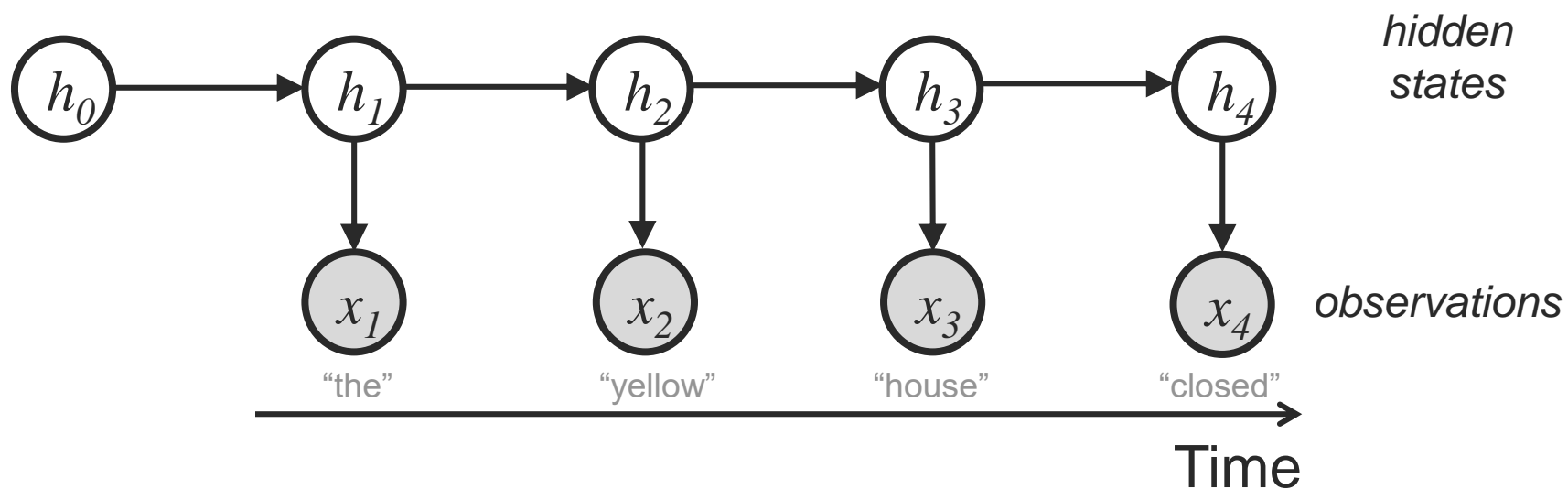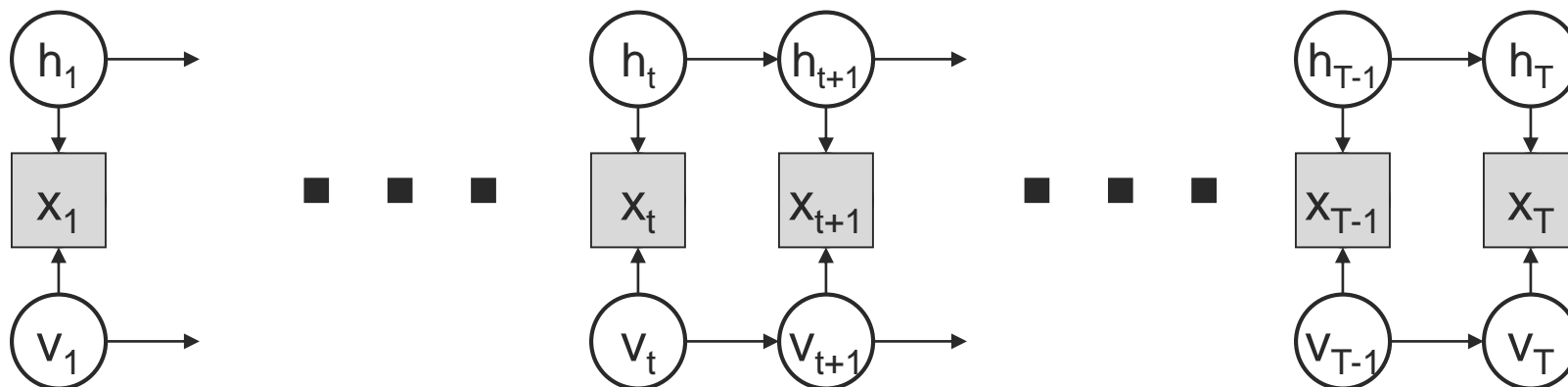
# Dynamic Bayesian Network (DBN)
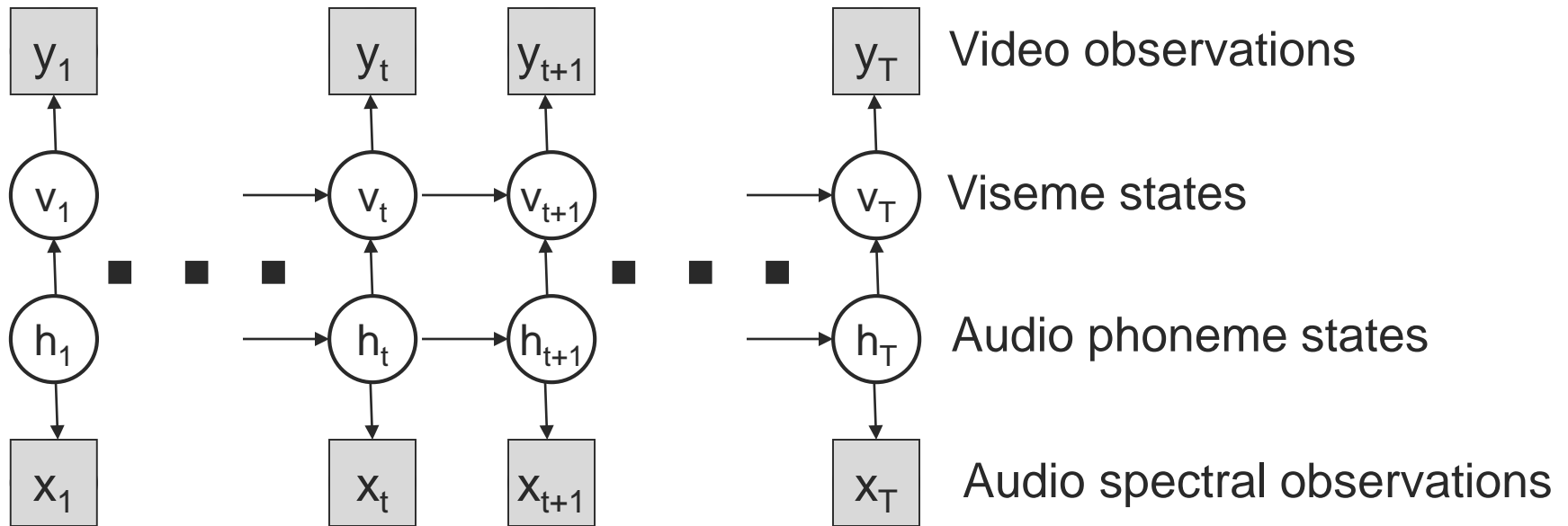
# Hidden Markov Models

# Factorial HMM



- Factorial HMM:
    - $h_t$ and $v_t$ represent two different types of background information, each with its own history
    - Observations $x_t$ depend on both hidden processes
- Model parameters:
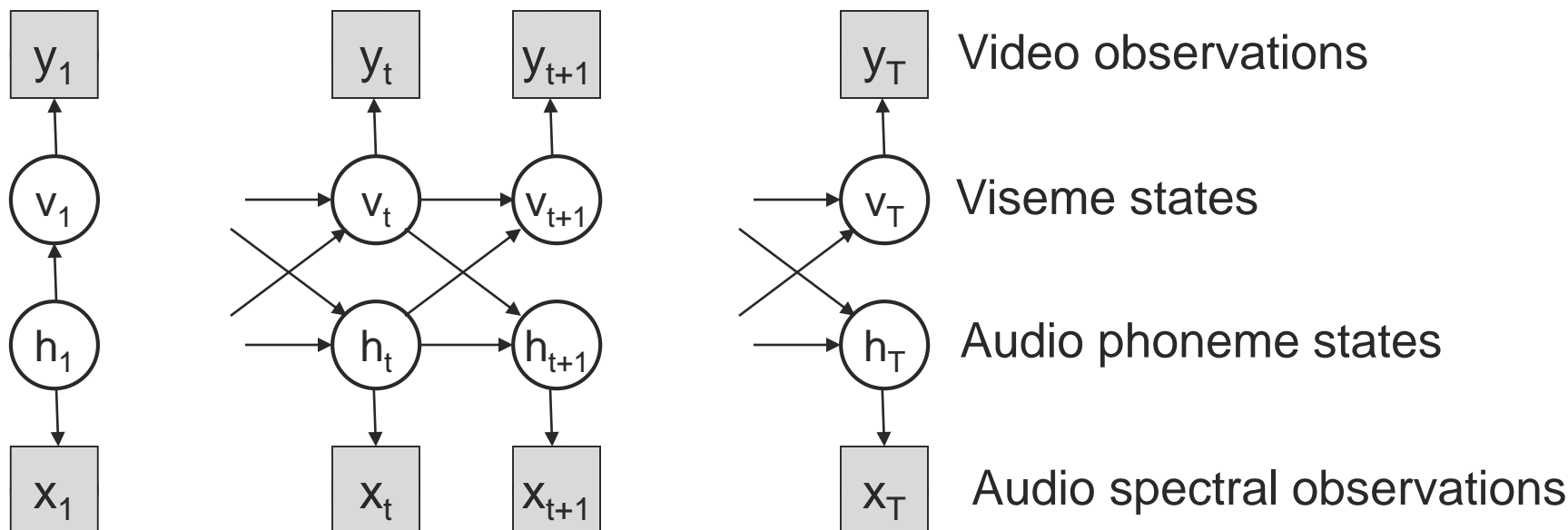    - $p(v_{t+1}|v_t)$, $p(h_{t+1}|h_t)$, $p(x_t|h_t,v_t)$

# The Boltzmann Zipper



$y_1$ $y_t$ $y_{t+1}$ $y_T$ — Video observations

$v_1$ $v_t$ $v_{t+1}$ $v_T$ — Viseme states

$h_1$ $h_t$ $h_{t+1}$ $h_T$ — Audio phoneme states

$x_1$ $x_t$ $x_{t+1}$ $x_T$ — Audio spectral observations

- Both streams have a "memory" ($h_t$ and $v_t$)

- Model parameters:
  - $p(h_{t+1}|h_t)$, $p(x_t|h_t)$
  - $p(v_{t+1}|v_t,h_{t+1})$, $p(y_t|h_t)$

# The Coupled HMM

$y_1$    $y_t$    $y_{t+1}$    $y_T$    Video observations

$v_1$    $v_t$    $v_{t+1}$    $v_T$    Viseme states

$h_1$    $h_t$    $h_{t+1}$    $h_T$    Audio phoneme states

$x_1$    $x_t$    $x_{t+1}$    $x_T$    Audio spectral observations

- Advantage over Boltzmann Zipper: More flexible, because neither vision nor sound is "privileged" over the other.
  - $p(h_{t+1}|v_t,h_t)$, $p(x_t|h_t)$
  - $p(v_{t+1}|v_t,h_t)$, $p(y_t|h_t)$

# Learning (Dynamic) Bayesian Networks

- Multiple techniques exist to learn the model parameters based on data
    - Maximum likelihood estimator
    - Bayesian estimator, which allows to include prior information
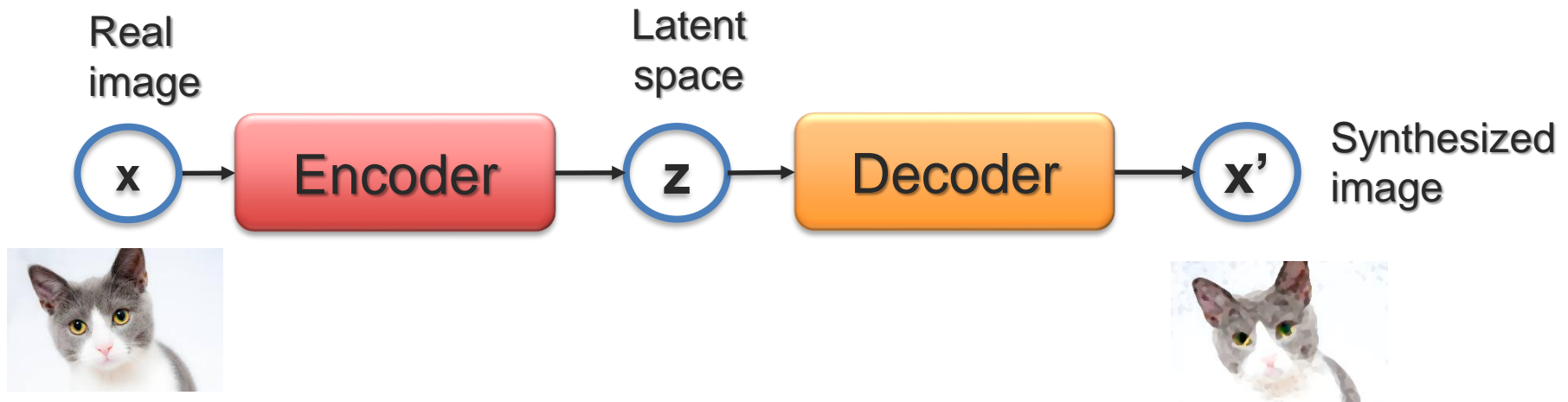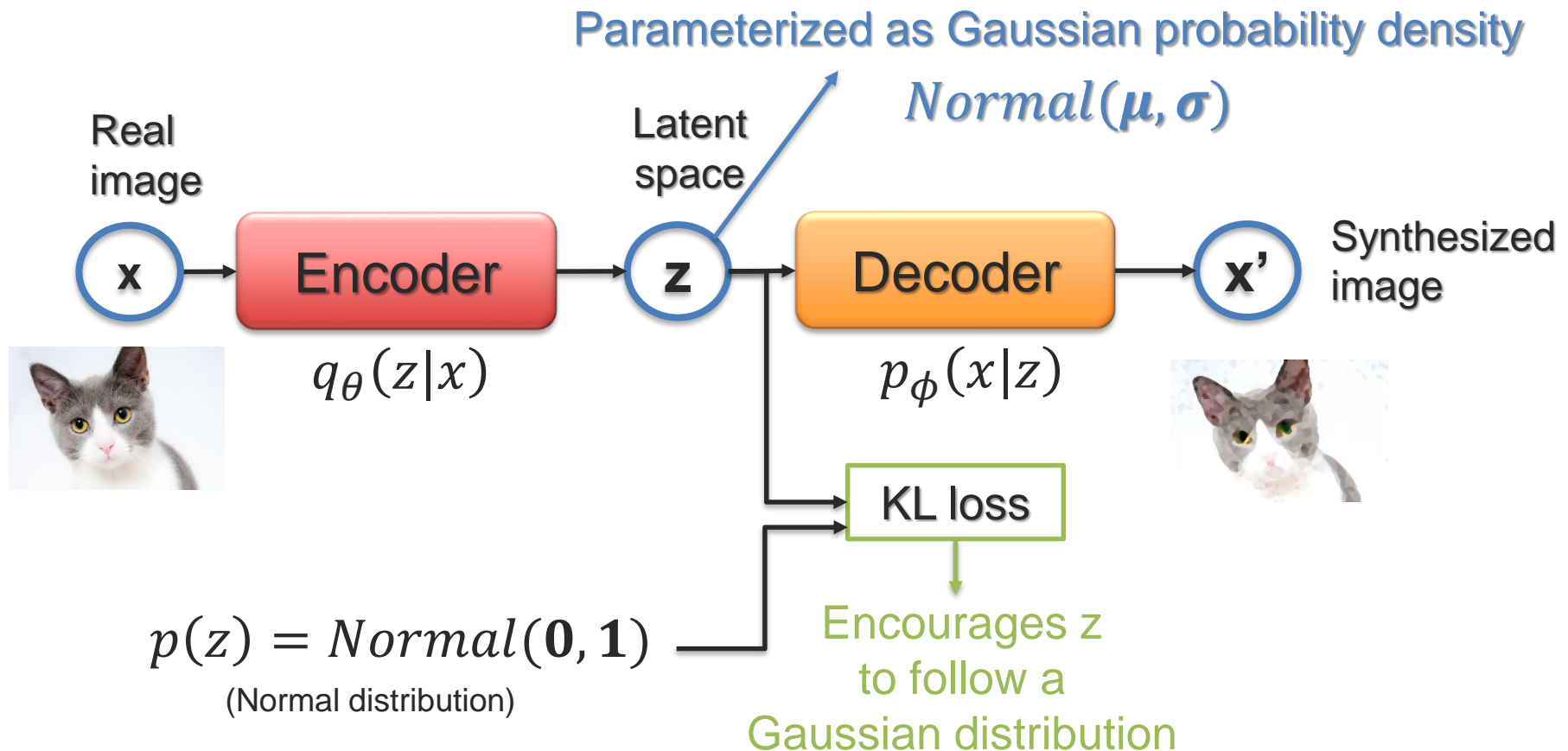- Python libraries:
    - http://pgmpy.org/
    - http://www.bayespy.org
    - https://pomegranate.readthedocs.io/en/latest/

# Variational AutoEncoder

# Auto-encoder



Real image → **x** → Encoder → Latent space **z** → Decoder → **x'** → Synthesized image

After learning this autoencoder,
can I input any z vector in the decoder?

# Variational Autoencoder

Parameterized as Gaussian probability density

$$Normal(\boldsymbol{\mu}, \boldsymbol{\sigma})$$

Real image

$x$ → Encoder → $z$

Latent space

$q_\theta(z|x)$

Decoder → $x'$

Synthesized image

$p_\phi(x|z)$

KL loss

$$p(z) = Normal(\boldsymbol{0}, \boldsymbol{1})$$
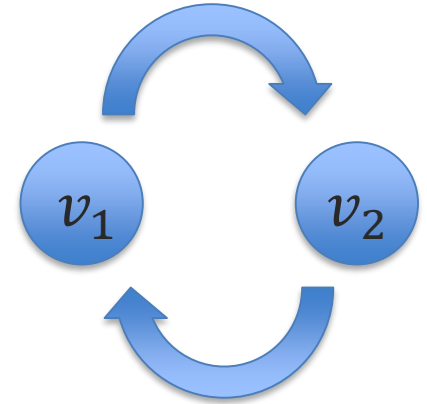
(Normal distribution)

Encourages z to follow a Gaussian distribution

# Variational Inference

- When inference is not possible
  - Either relax the problem
  - Or use variational methods
- Variational inference:
  - Unroll through time (MCMC, Gibbs) – RBM
  - Mean-field Approximation (Fully Connected CRF)

$v_1$ $v_2$

# Variational Auto-encoder
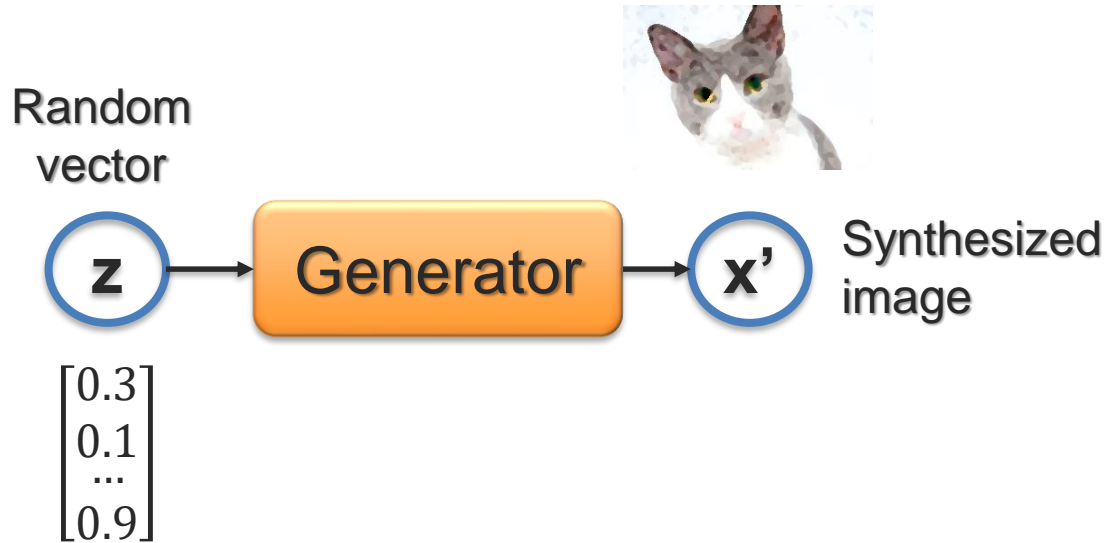
- The normal distribution has nice properties.



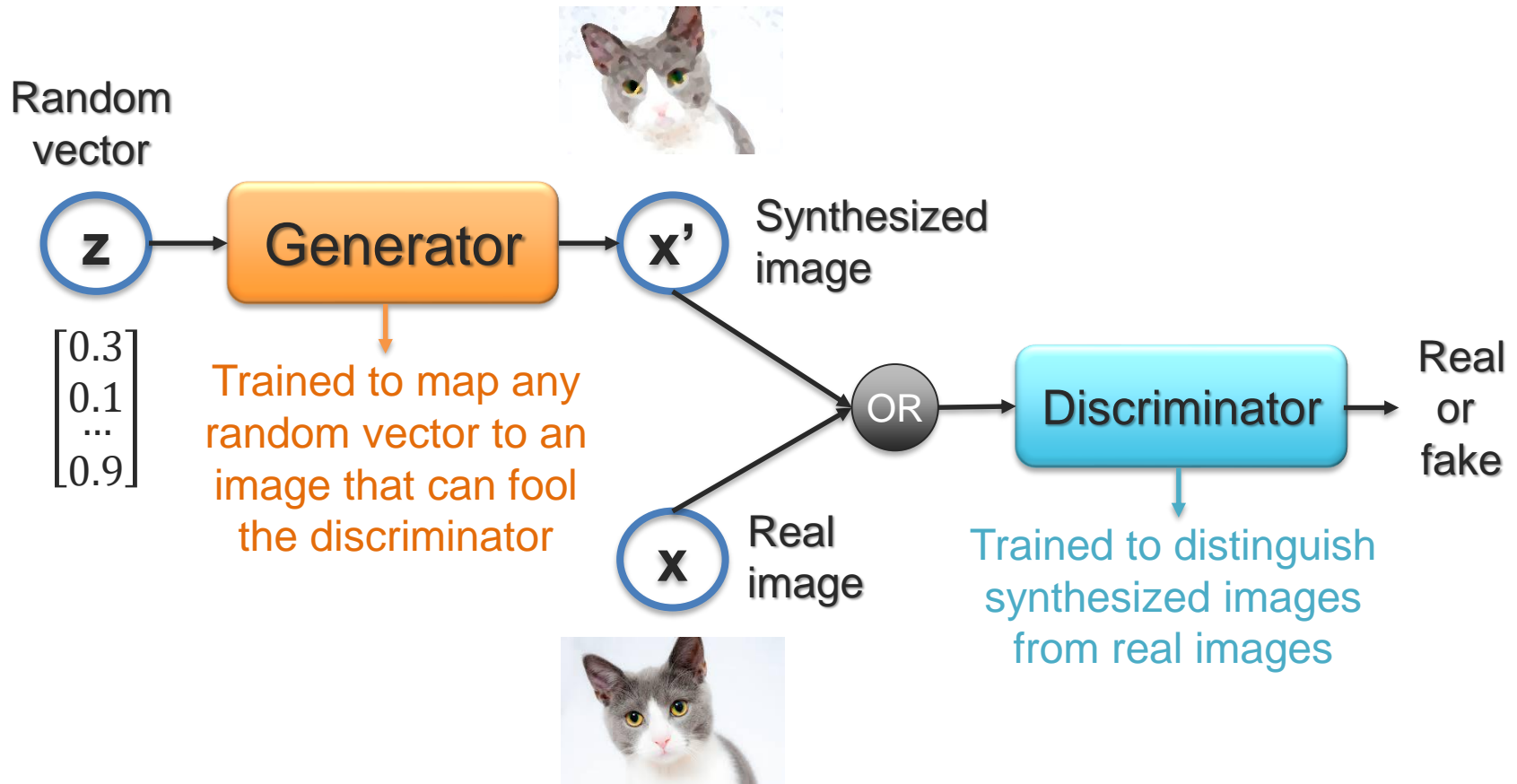But these images are not as realistic looking…

# Generative Adversarial Networks

Language Technologies Institute

Carnegie Mellon University

# Generative Network

Random
vector

$\mathbf{z}$ → Generator → $\mathbf{x'}$    Synthesized
image

$$\begin{bmatrix} 0.3 \\ 0.1 \\ ... \\ 0.9 \end{bmatrix}$$

How to train the generator to synthesize realistic images?

# Generative Adversarial Network (GAN)



**Random vector**

$$\begin{bmatrix} 0.3 \\ 0.1 \\ \cdots \\ 0.9 \end{bmatrix}$$

**z** → **Generator** → **x'**   Synthesized image

Trained to map any random vector to an image that can fool the discriminator

**x**   Real image

OR → **Discriminator** → Real or fake

Trained to distinguish synthesized images from real images

## How to train both the generator and the discriminator?

# GAN Training

$$\max_{\mathcal{D}} \min_{\mathcal{G}} V(\mathcal{G}, \mathcal{D})$$

How do we optimize
this objective function?

$$V(\mathcal{G}, \mathcal{D}) = \mathbb{E}_{p_{data}(\mathbf{x})} \log \mathcal{D}(\mathbf{x}) + \mathbb{E}_{p_g(\mathbf{x})} \log(1 - \mathcal{D}(\mathbf{x}))$$

Random
vector

$\mathbf{z}$ → Generator → $\mathbf{x}'$  Synthesized image

$$\begin{bmatrix} 0.3 \\ 0.1 \\ \cdots \\ 0.9 \end{bmatrix}$$

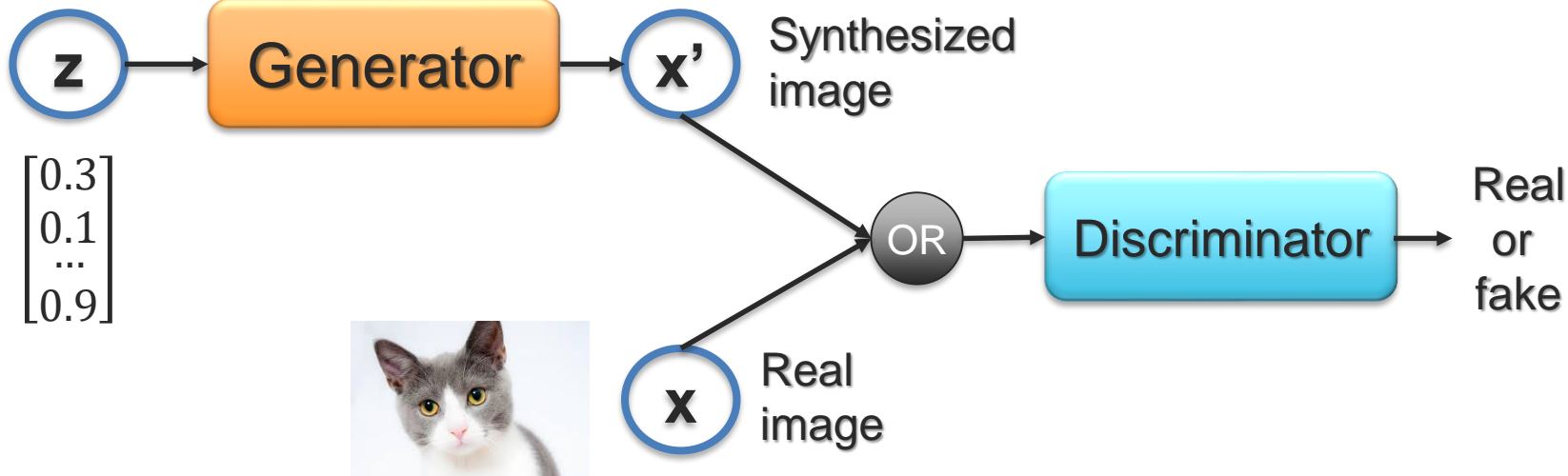$\mathbf{x}$  Real image

OR → Discriminator → Real or fake

# GAN Training

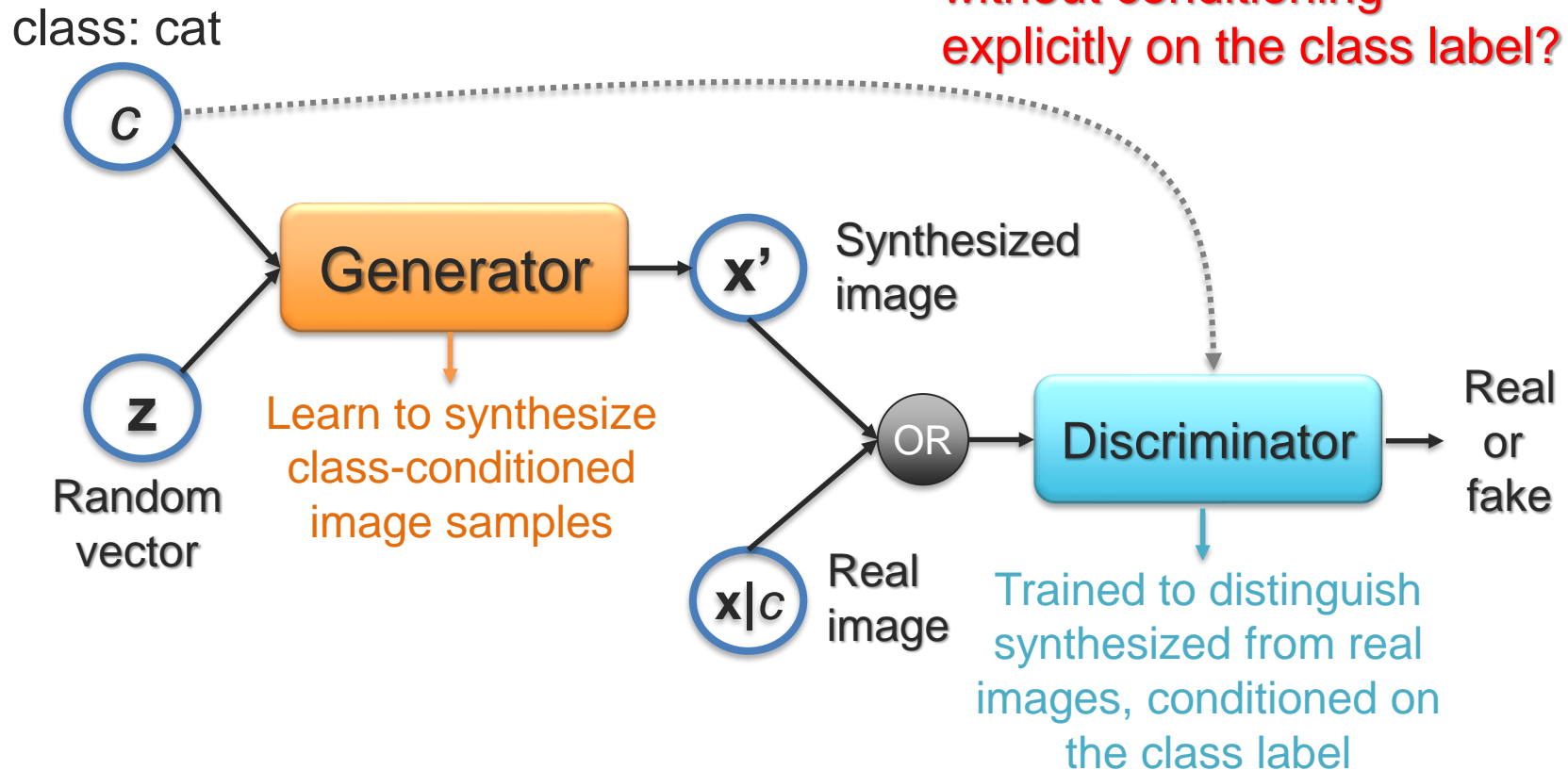$$\max_{\mathcal{D}} \min_{\mathcal{G}} V(\mathcal{G}, \mathcal{D})$$

Optimization:

(1) Fix generator, and update discriminator

(2) Fix discriminator, and update generator
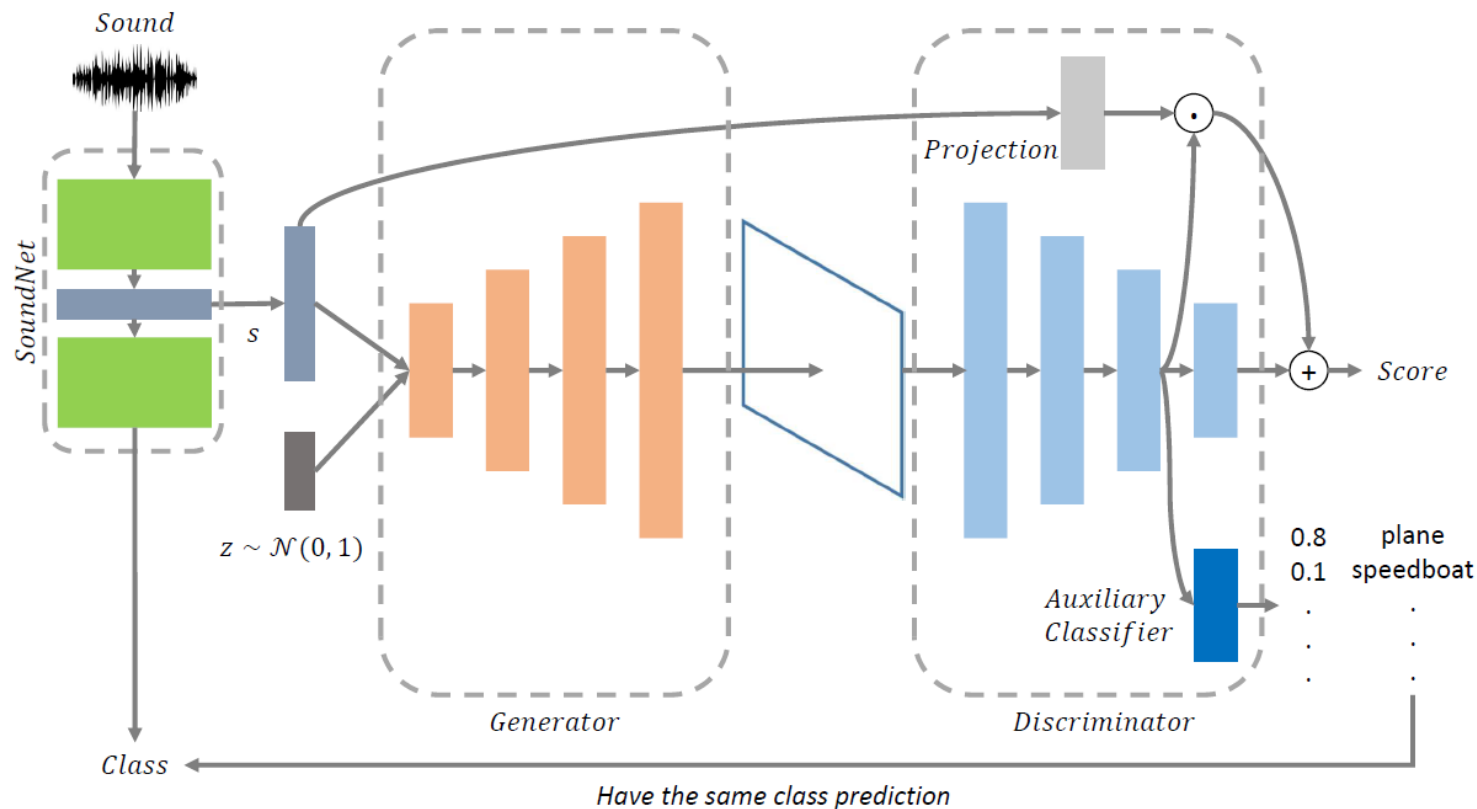
Random vector

$$\begin{bmatrix} 0.3 \\ 0.1 \\ \cdots \\ 0.9 \end{bmatrix}$$

z → Generator → x' Synthesized image

OR → Discriminator → Real or fake

x Real image

# Conditional GAN

class: cat



How to train discriminator without conditioning explicitly on the class label?

$c$

Generator

Learn to synthesize class-conditioned image samples

$z$

Random vector

$x'$ Synthesized image

OR

$x|c$ Real image

Discriminator

Trained to distinguish synthesized from real images, conditioned on the class label

Real or fake

# Audio to Scene

Language Technologies Institute

Carnegie Mellon University

# Audio to Scene



https://wjohn1483.github.io/audio_to_scene/index.html

# Talking Head

Language Technologies Institute

Carnegie Mellon University

# Info GAN

What if we would like to also learn an encoder from x to z?

class: cat

$c$

Generator

Learn to synthesize class-conditioned image samples

$z$

Random vector

x'

Synthesized image

OR

$x|c$

Real image

Discriminator

Outputs both an authenticity decision and a class label prediction

$\hat{c}$

Real or fake

# Talking Head

Language Technologies Institute

Carnegie Mellon University

# Bidirectional GAN

Random vector

Synthesized image

z → Generator → x'

Selects either
(z',x) or (z,x')

OR → Discriminator → Real or fake

z' ← Encoder ← x

Encoding of real image

Learn to map image to latent space

Real image

What if we would like to also learn an encoder from x to z?

# Paired and Unpaired Data



Paired

$x_i$ $y_i$

Unpaired

$X$ $Y$

Language Technologies Institute

Carnegie Mellon University

# Cycle GAN

Language Technologies Institute

Carnegie Mellon University

# cAE-GAN



Real image → x → Encoder → z (Latent space) → Generator → x' (Synthesized image)

x OR x' → Discriminator → Real or fake

# cVAE-GAN

Real
image

Latent
space

Synthesized
image

x → Encoder → z → Generator → x'

KL loss

$N(z)$

OR → Discriminator → Real
or
fake

Language Technologies Institute

Carnegie Mellon University

# BiCycle GAN



**Let's put everything in one model!!**

Legend:
- Input Image
- Ground truth output
- Network output
- Loss
- Deep network
- Target latent distribution
- Sample from distribution

(c) Training cVAE-GAN

(d) Training cLR-GAN

| Input | Ground truth | Generated samples |
|---|---|---|