



# 10-708 Probabilistic Graphical Models

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University



## Course Overview

Matt Gormley  
Lecture 1  
Feb. 01, 2021

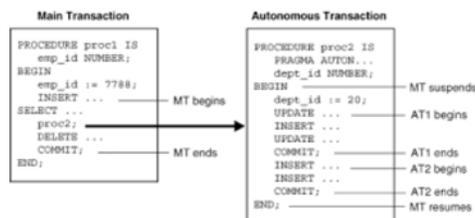
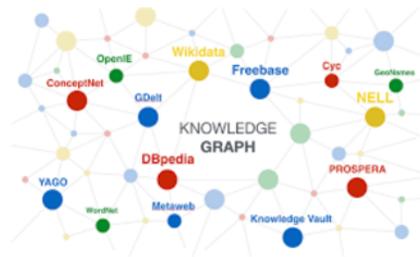
How to define a structured prediction problem

# **STRUCTURED PREDICTION**

# Structured vs. Unstructured Data

## Structured Data Examples

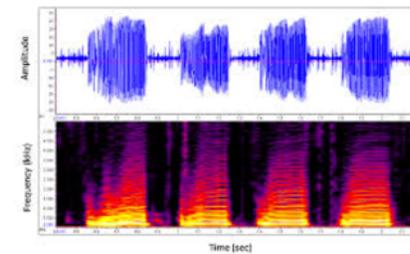
- database entries
- transactional information
- wikipedia infobox
- knowledge graphs
- hierarchies



## Unstructured Data Examples

- written text
- images
- videos
- spoken language
- music
- sensor data

مساء الخير! مرحبا بكم في الدرجة



# Structured Prediction

- The focus of most Intro ML courses is **classification**
  - Given observations:  $\mathbf{x} = (x_1, x_2, \dots, x_K)$
  - Predict a (binary) **label**:  $y$
- Many real-world problems require **structured prediction**
  - Given observations:  $\mathbf{x} = (x_1, x_2, \dots, x_K)$
  - Predict a **structure**:  $\mathbf{y} = (y_1, y_2, \dots, y_J)$
- Some *classification* problems benefit from **latent structure**

# Structured Prediction

## Classification / Regression

1. Input can be semi-structured data
2. Output is a **single number (integer / real)**
3. In linear models, features can be arbitrary combinations of [input, output] pair
4. Output space is **small**
5. Inference **is trivial**

## Structured Prediction

1. Input can be semi-structured data
2. Output is a **sequence of numbers representing a structure**
3. In linear models, features can be arbitrary combinations of [input, output] pair
4. Output space **may be exponentially large in the input space**
5. Inference **problems are NP-hard or #P-hard in general and often require approximations**

# Structured Prediction Examples

- **Examples of structured prediction**
  - Part-of-speech (POS) tagging
  - Handwriting recognition
  - Speech recognition
  - Object detection
  - Scene understanding
  - Machine translation
  - Protein sequencing

# Part-of-Speech (POS) Tagging

|           |       |       |      |       |       |
|-----------|-------|-------|------|-------|-------|
| Sample 1: | n     | v     | p    | d     | n     |
|           | time  | flies | like | an    | arrow |
| Sample 2: | n     | n     | v    | d     | n     |
|           | time  | flies | like | an    | arrow |
| Sample 3: | n     | v     | p    | n     | n     |
|           | flies | fly   | with | their | wings |
| Sample 4: | p     | n     | n    | v     | v     |
|           | with  | time  | you  | will  | see   |

# Dataset for Supervised Part-of-Speech (POS) Tagging

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

|           |   |   |   |   |   |   |           |
|-----------|---|---|---|---|---|---|-----------|
| Sample 1: |    |    |    |    |    |    | $y^{(1)}$ |
|           |    |    |    |    |    |    | $x^{(1)}$ |
| Sample 2: |    |    |    |    |    |    | $y^{(2)}$ |
|           |    |    |    |    |    |    | $x^{(2)}$ |
| Sample 3: |    |    |    |    |    |    | $y^{(3)}$ |
|           |   |   |   |   |   |   | $x^{(3)}$ |
| Sample 4: |  |  |  |  |  |  | $y^{(4)}$ |
|           |  |  |  |  |  |  | $x^{(4)}$ |

# Handwriting Recognition

Sample 1:

u n e x p e c t e d



Sample 2:

v o l c a n i c



Sample 2:

e m b r a c e s



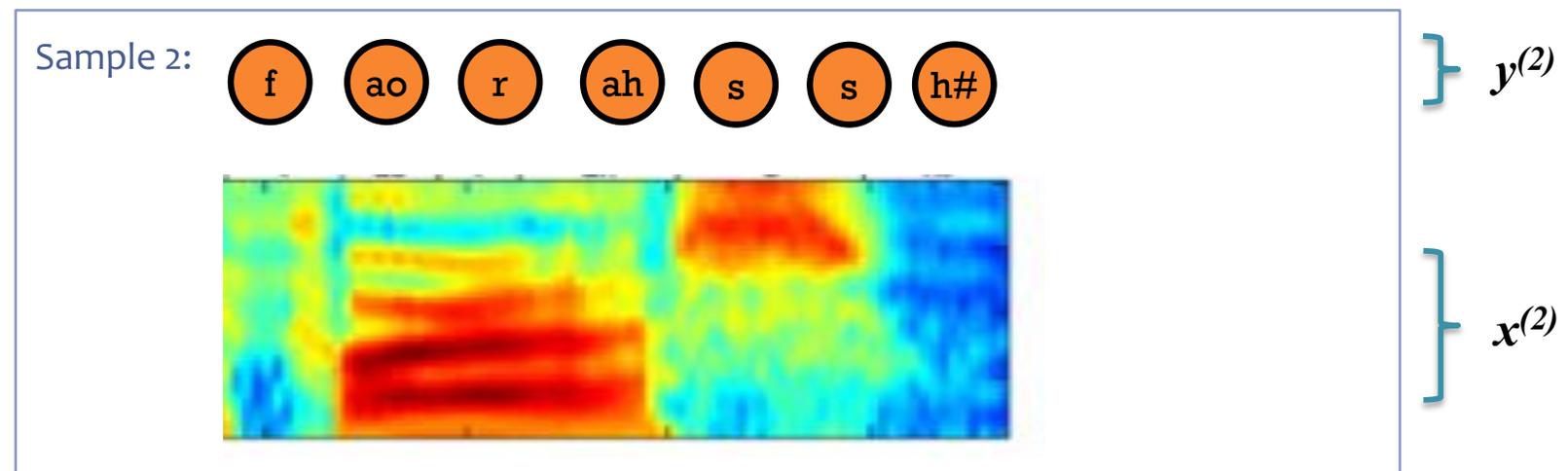
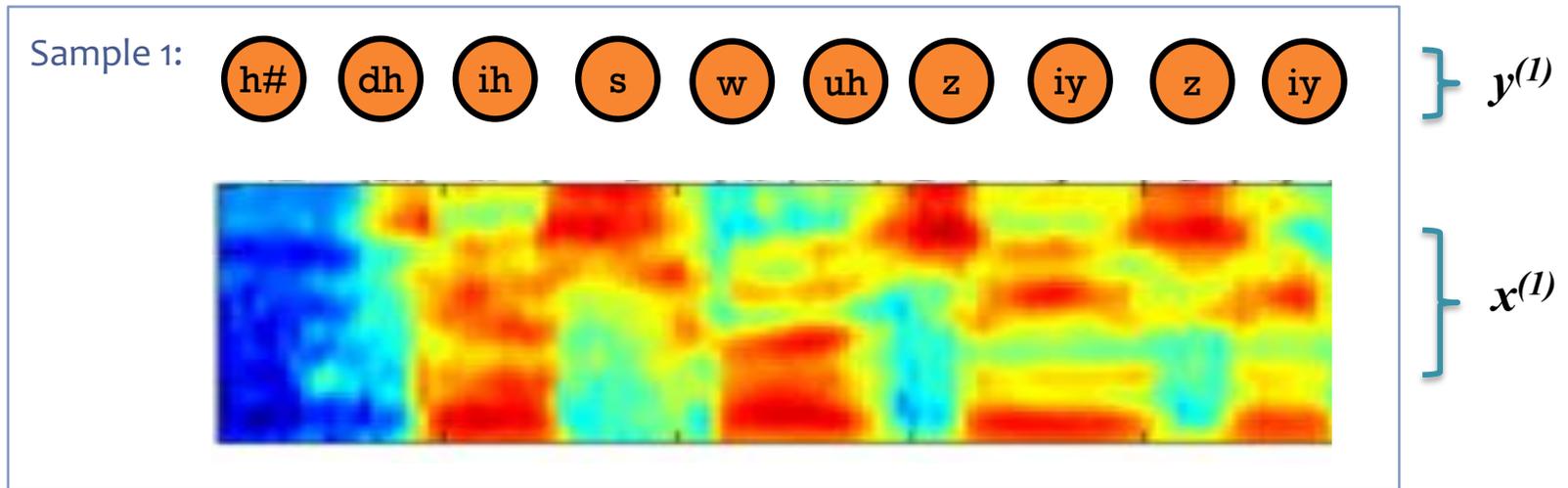
# Dataset for Supervised Handwriting Recognition

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$



# Dataset for Supervised Phoneme (Speech) Recognition

Data:  $\mathcal{D} = \{ \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \}_{n=1}^N$



# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .



pigeon

$x^{(1)}$

$y^{(1)}$



rhinoceros

$x^{(2)}$

$y^{(2)}$



leopard

$x^{(3)}$

$y^{(3)}$



llama

$x^{(4)}$

$y^{(4)}$

# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time

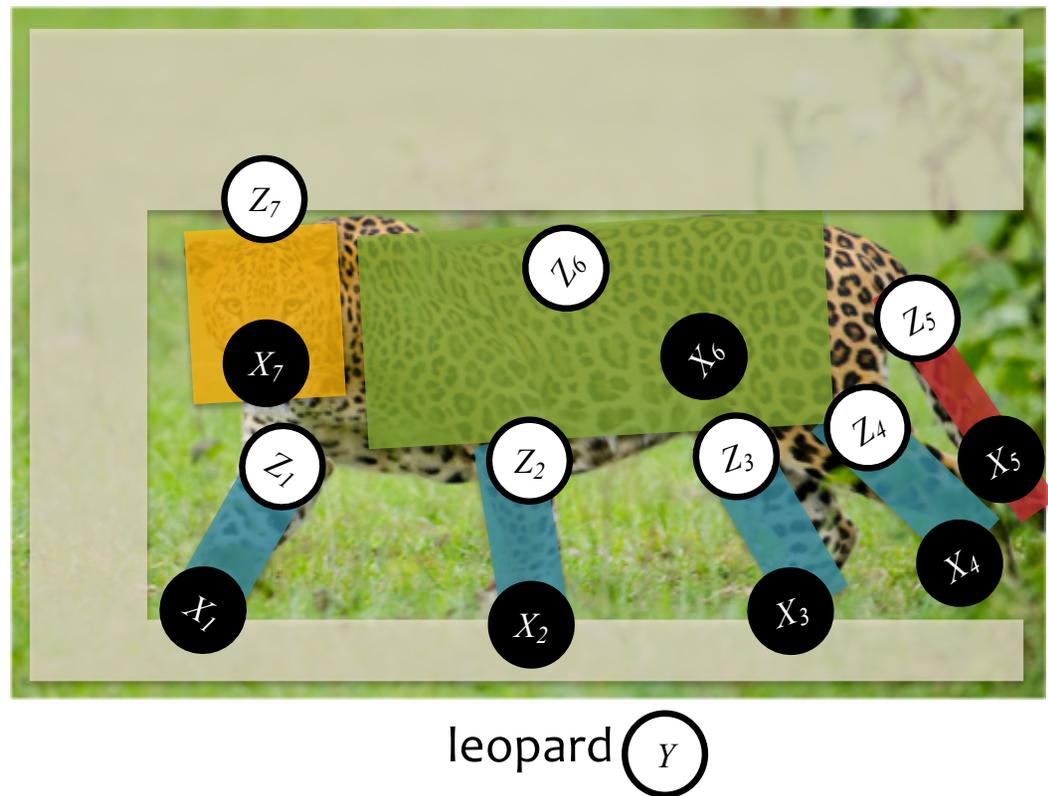


leopard

# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

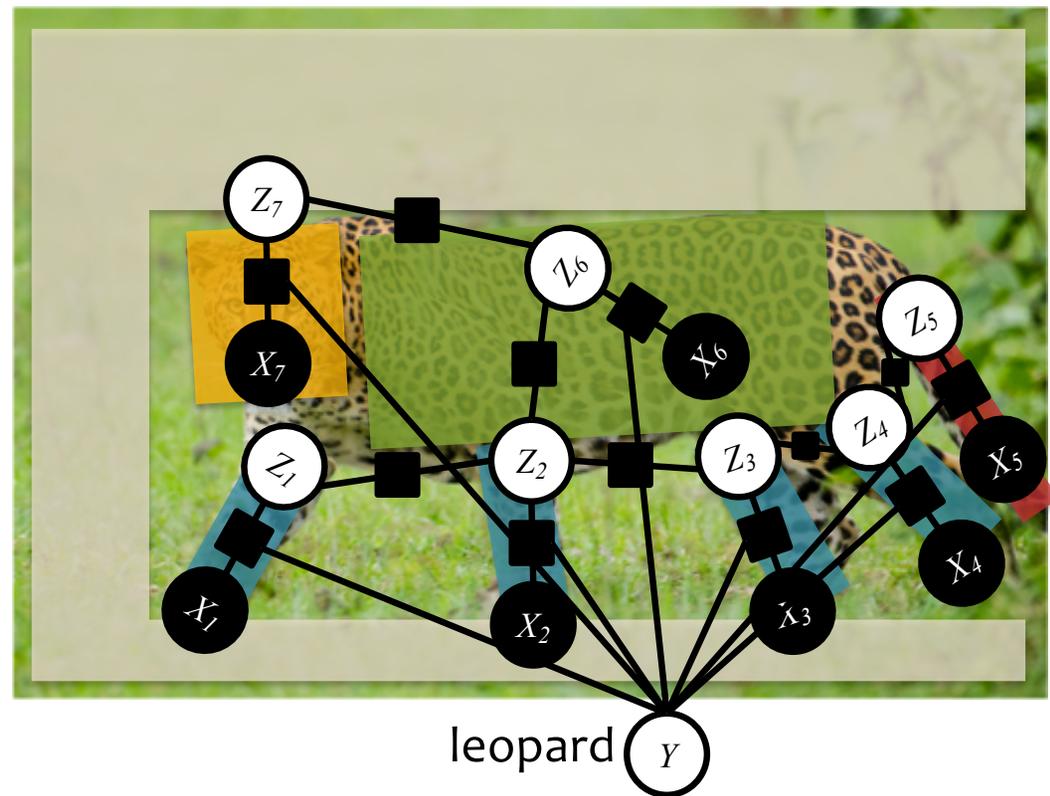
- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time



# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time



# Structured Prediction

## Preview of challenges to come...

- Consider the task of finding the **most probable assignment** to the output

Classification

$$\hat{y} = \operatorname{argmax}_y p(y|\mathbf{x})$$

where  $y \in \{+1, -1\}$

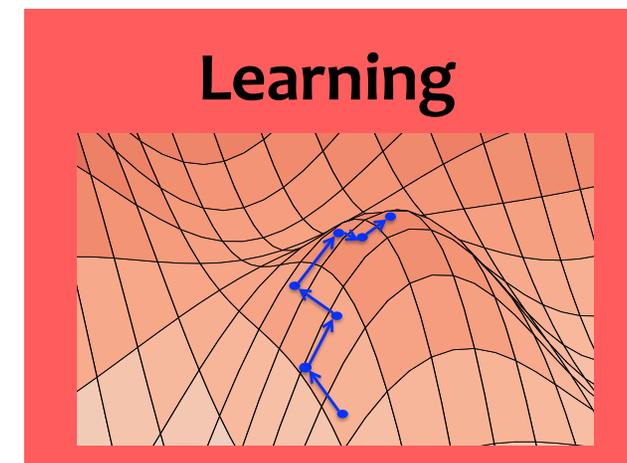
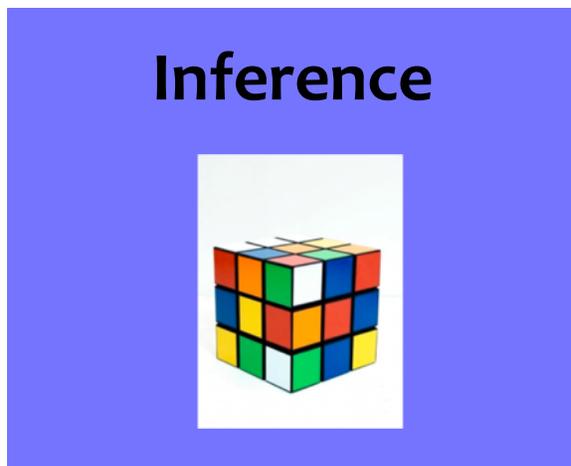
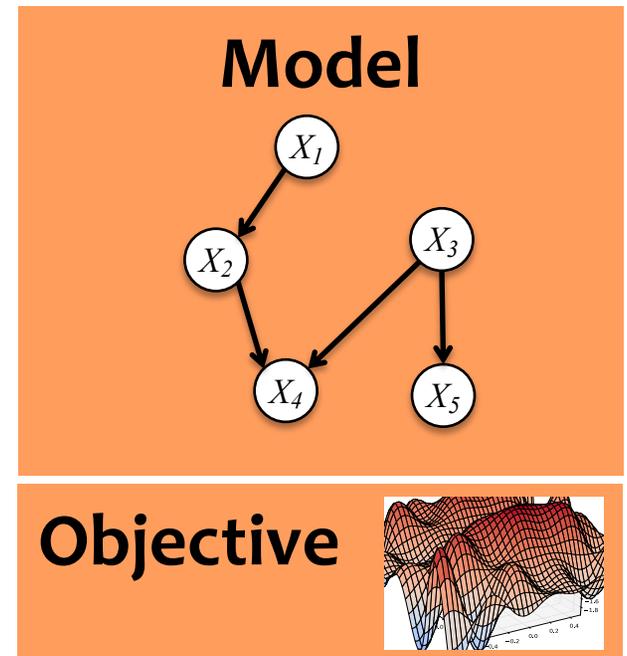
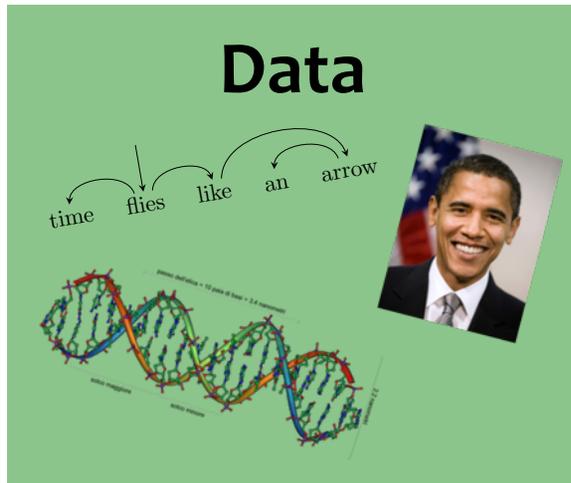
Structured Prediction

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$$

where  $\mathbf{y} \in \mathcal{Y}$

and  $|\mathcal{Y}|$  is very large

# Structured Prediction



(Inference is usually called as a subroutine in learning)

# Structured Prediction

The **data** inspires the structures we want to predict



Our **model** defines a score for each structure

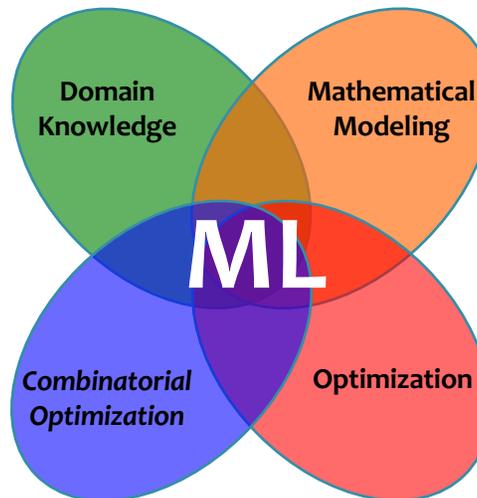
It also tells us what to optimize



**Learning** tunes the parameters of the model

**Inference** finds {best structure, marginals, partition function} for a new observation

(**Inference** is usually called as a subroutine in learning)



# Decomposing a Structure into Parts

- Why divide a **structure** into its **pieces**?
  - amenable to **efficient inference**
  - enable natural **parameter sharing** during learning
  - easier definition of fine-grained **loss functions**
  - clearer depiction of **model's uncertainty**
  - easier specification of **interactions** between the parts
  - (may) lead to natural definition of a **search problem**
- A key step in **formulating a task as a structured prediction**

# Scene Understanding

- **Variables:**
  - boundaries of image regions
  - tags of regions
- **Interactions:**
  - semantic plausibility of nearby tags
  - continuity of tags across visually similar regions (i.e. patches)

Labels **with** top-down information

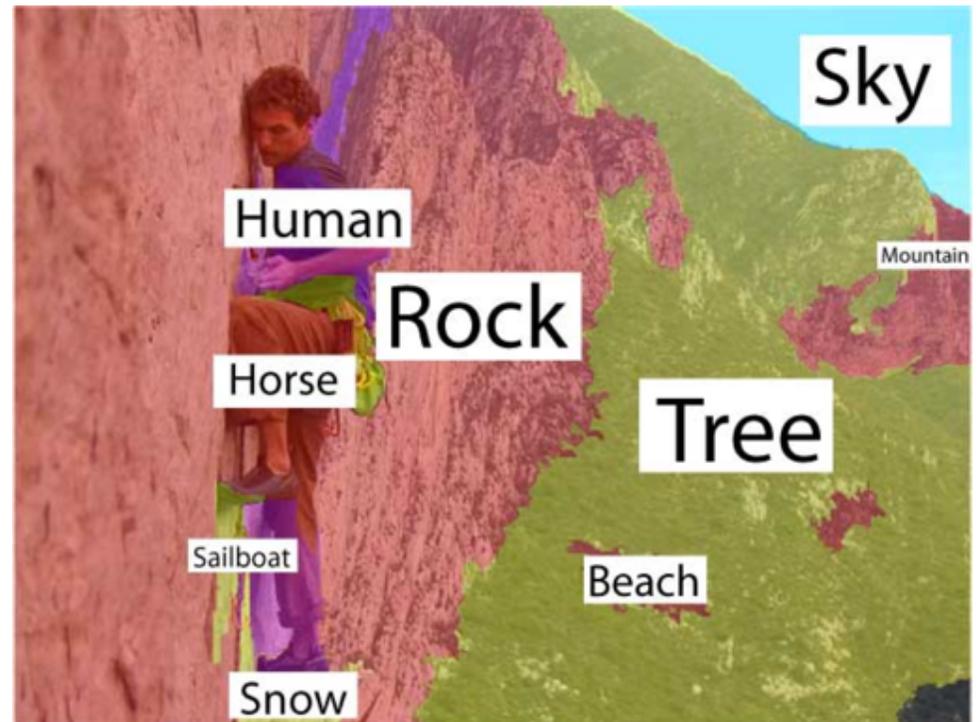


(Li et al., 2009)

# Scene Understanding

- **Variables:**
  - boundaries of image regions
  - tags of regions
- **Interactions:**
  - semantic plausibility of nearby tags
  - continuity of tags across visually similar regions (i.e. patches)

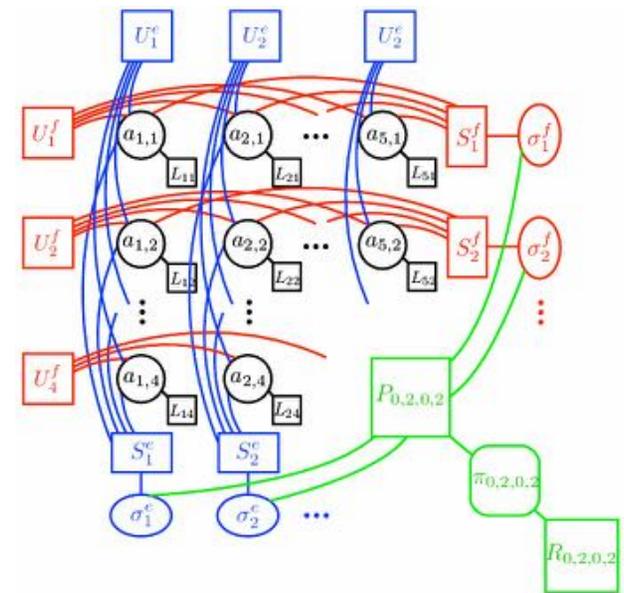
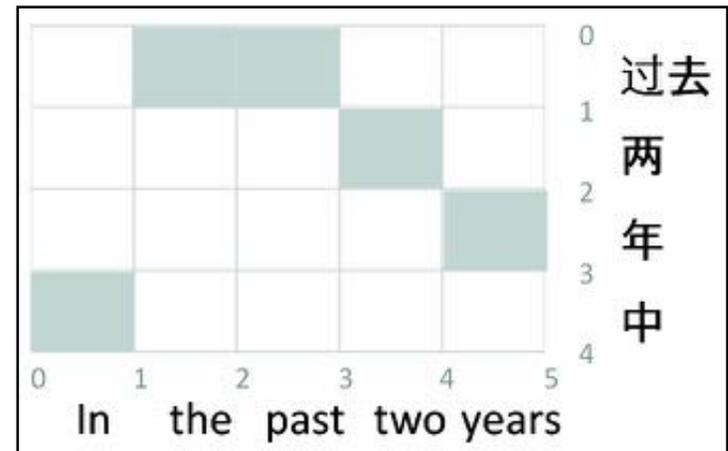
Labels **without** top-down information



(Li et al., 2009)

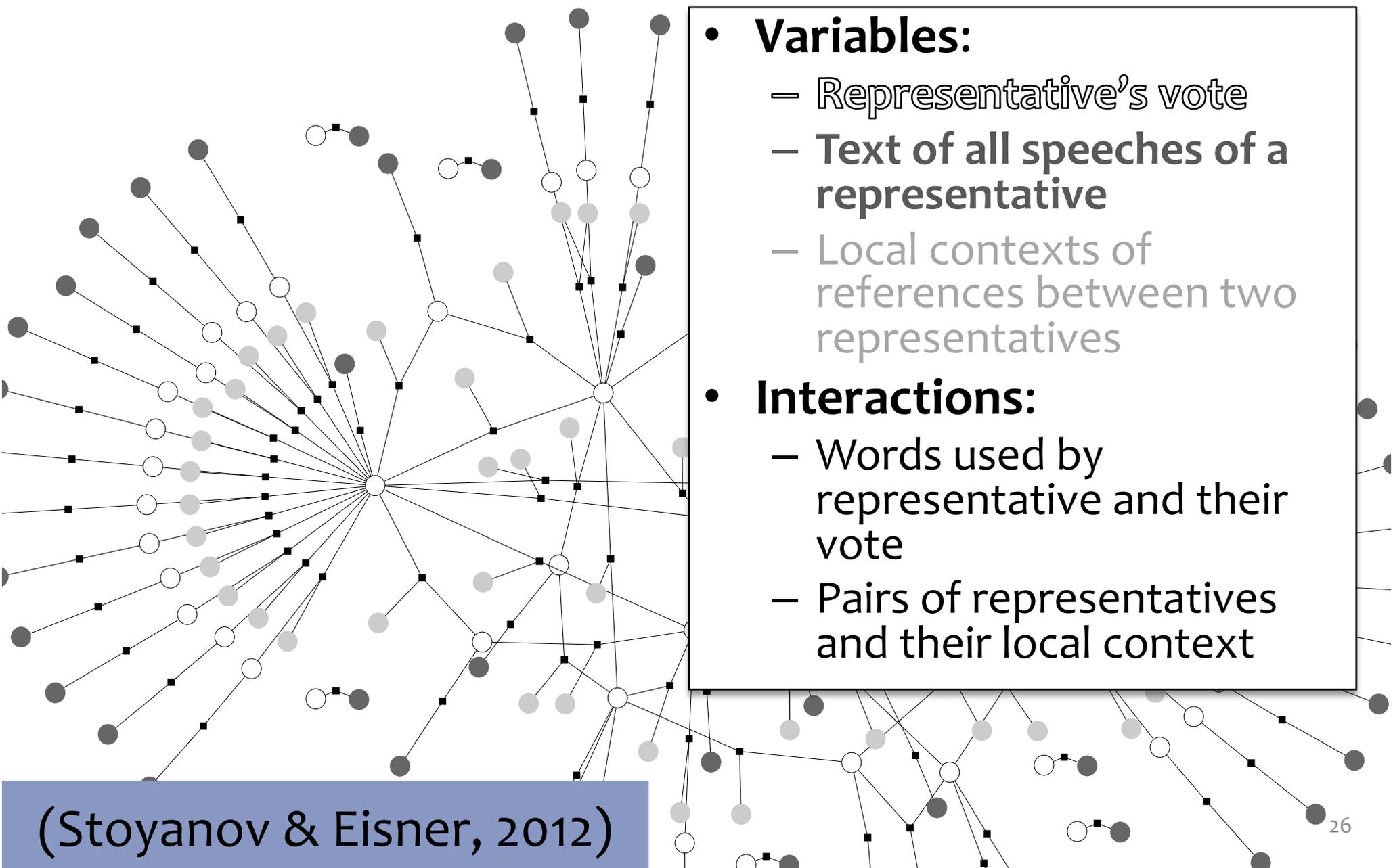
# Word Alignment / Phrase Extraction

- **Variables (boolean):**
  - For each (Chinese phrase, English phrase) pair, are they linked?
- **Interactions:**
  - Word fertilities
  - Few “jumps” (discontinuities)
  - Syntactic reorderings
  - “ITG constraint” on alignment
  - Phrases are disjoint (?)



(Burkett & Klein, 2012)

# Congressional Voting



- **Variables:**
  - Representative's vote
  - Text of all speeches of a representative
  - Local contexts of references between two representatives
- **Interactions:**
  - Words used by representative and their vote
  - Pairs of representatives and their local context

(Stoyanov & Eisner, 2012)

# Medical Diagnosis

The screenshot shows a web-based medical diagnosis form. At the top, there is a navigation bar with links for 'Application', 'CMI', 'Eoster', 'List', 'Dashboard', 'ED', 'Order', 'Responses', 'Reports', 'User', 'Feedback', and 'Help'. Below this is a toolbar with various icons. The main form area is divided into several sections:

- Assessment**: Includes tabs for 'Physical', 'Investigations', and 'Discharge'. It contains fields for 'Temp (\*C)', 'P', 'BP(L)', 'Resp', 'O2 Sat(%)', and 'BP(R)'. There are also checkboxes for 'Emerg. Phys.' and 'Resident'.
- Assessment**: Contains a dropdown for 'Onset of Pain', a text field for 'age', a dropdown for 'Duration', and a dropdown for 'Severity'. It also has a 'Describe Pain' section with checkboxes for 'Aching', 'Pressure', 'Spurting', 'Burning', and 'Stabbing', and a 'Radiating' section with checkboxes for 'L Arm' and 'R Arm'. There is an 'Other:' text field.
- Pain worse with**: Includes checkboxes for 'Activity', 'Eating', 'Movement', 'Deep Breathing', 'Lying', and 'Sitting'. There is an 'Other:' text field.
- Pain relieved with**: Includes checkboxes for 'Deep breathing', 'NTG', 'Eating', and 'Rest'. There is an 'Other:' text field and a 'Specify:' dropdown.
- Cardiac Risk Factors**: Includes checkboxes for 'Hx MI', 'Diabetes', 'Hx IHD', 'Hypertension', 'CABG', 'Increased Cholesterol', 'CHF', 'Family Hx IHD age < 60 years', 'PCA/Stenting', 'Smoker', and 'Other:'.
- Associated Symptoms**: Includes checkboxes for 'Nausea', 'Presyncope', 'Rumbling', 'Syncope', 'Diaphoresis', 'Cough', 'Shortness of Breath', 'Peripheral Edema(New)', 'Palpitations', and 'Orthopnea'.
- Family History**: Includes a text field for 'ETOH'.

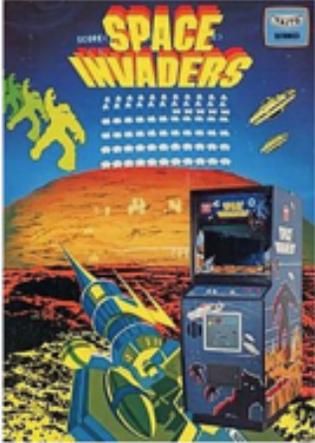
- **Variables:**
  - content of text field
  - checkmark
  - dropdown menu
- **Interactions:**
  - groups of related symptoms (e.g. that are predictive of a disease)
  - social history (e.g. smoker) and symptoms
  - risk factors (e.g. infant) and lab results

# Wikipedia Infoboxes

| Gryan Miers   |  |                      |
|---|--|----------------------|
| Personal information  |  |                      |
| <b>Date of birth</b>  | 30 March 1999 (age 20)                             |                      |
| <b>Original team(s)</b>   | Grovedale (GFL)                                    |                      |
| <b>Draft</b>  | No. 57, 2017 national draft                        |                      |
| <b>Debut</b>  | Round 1, 2019, Geelong vs. Collingwood, at the MCG |                      |
| <b>Height</b>   | 178 cm (5 ft 10 in)                                |                      |
| <b>Weight</b>   | 78 kg (172 lb)                                     |                      |
| <b>Position(s)</b>  | Small forward                                      |                      |
| Club information  |  |                      |
| <b>Current club</b>   | Geelong  |                      |
| <b>Number</b>   | 32   |                      |
| Playing career <sup>1</sup>   |  |                      |
| <b>Years</b>  | <b>Club</b>  | <b>Games (Goals)</b> |
| 2019–   | Geelong  | 19 (19)              |
| <sup>1</sup> Playing statistics correct to the end of round 17, 2019.           |  |                      |
| Career highlights   |  |                      |
| <ul style="list-style-type: none"> <li>AFL Rising Star nominee: 2019</li> </ul> |  |                      |
| Sources: AFL Tables, AustralianFootball.com                                     |  |                      |

| Pather Panchali   |  |
|---|--|
|  |  |
| A poster of Pather Panchali   |  |
| <b>Directed by</b>  | Satyajit Ray   |
| <b>Screenplay by</b>  | Satyajit Ray   |
| <b>Based on</b>   | Pather Panchali by Bibhutibhusan Bandyopadhyay   |
| <b>Starring</b>   | Subir Banerjee<br>Kanu Banerjee<br>Karuna Banerjee<br>Uma Dasgupta<br>Chunibala Devi<br>Tulsi Chakrabarti                              |
| <b>Music by</b>   | Ravi Shankar   |
| <b>Cinematography</b>   | Subrata Mitra  |
| <b>Edited by</b>  | Dulal Dutta  |
| <b>Production company</b>   | Government of West Bengal  |
| <b>Distributed by</b>   | Aurora Film Corporation (1955)<br>Edward Harrison (1958)<br>Merchant Ivory Productions<br>Sony Pictures Classics (1995) <sup>[a]</sup> |
| <b>Release date</b>   | 26 August 1955 (India)   |

| Changsha Kingdom   |   |
|--|---|
| 長沙國  |   |
| 203/202 BC–AD 33   |   |
|         |   |
| Silk map unearthed from Mawangdui, showing Changsha and the neighboring kingdom of Nanyue. |   |
| <b>Capital</b>   | Linxiang (present-day Changsha)   |
| <b>Government</b>  | Monarchy  |
| <b>History</b>   | <ul style="list-style-type: none"> <li>Established</li> <li>Extinction of the Wu family line</li> <li>Reestablishment under the Liu family</li> <li>Dissolution under Wang Mang</li> <li>Restoration</li> <li>Disestablished</li> </ul> |
|  | <ul style="list-style-type: none"> <li>203/202 BC</li> <li>157 BC</li> <li>155 BC</li> <li>AD 9</li> <li>AD 26</li> <li>AD 33</li> </ul>  |

| Space Invaders  |  |
|---|--|
|  |  |
| Promotional flyer   |  |
| <b>Developer(s)</b>   | Taito  |
| <b>Publisher(s)</b>   | JP: Taito<br>NA: Midway<br>EU: Midway <sup>[1]</sup><br>AU: Leisure & Allied Industries <sup>[2]</sup><br>Atari, Inc. (home) |
| <b>Designer(s)</b>  | Tomohiro Nishikado   |
| <b>Platform(s)</b>  | Arcade, Atari 2600, Atari 5200, Atari 8-bit, MSX   |
| <b>Release</b>  | JP: June 1978 <sup>[3]</sup><br>NA: July 1978  |
| <b>Genre(s)</b>   | Fixed shooter  |
| <b>Mode(s)</b>  | Single-player, 2 players alternating   |
| <b>Cabinet</b>  | Upright, cocktail <sup>[4]</sup>   |
| <b>Arcade system</b>  | Taito 8080 <sup>[5]</sup>  |
| <b>CPU</b>  | 8080 @ 2 MHz <sup>[5]</sup>  |
| <b>Sound</b>  | SN76477 @ 1.9968 MHz   |
| <b>Display</b>  | Fujitsu MB14241, <sup>[6]</sup> monochrome raster, vertical orientation, 224×256 resolution <sup>[5]</sup>                   |

# Exercise: Wikipedia Infoboxes

## Question:

Suppose you want to populate missing infobox fields. What model would you pick for the job, and why?

- A. Multiclass classifier
- B. RNN
- C. Graphical model

## Answer:

### Central Park Conservancy

From Wikipedia, the free encyclopedia

Coordinates: 40.76424°N 73.97169°W﻿ / ﻿40.76424°N 73.97169°W﻿ / 40.76424; -73.97169

The **Central Park Conservancy** is a private, nonprofit park conservancy that manages Central Park under a contract with the City of New York and NYC Parks. The conservancy employs most maintenance and operations staff in the park. It effectively oversees the work of both the private and public employees under the authority of the publicly appointed Central Park administrator, who reports to the parks commissioner and the conservancy's president.<sup>[1]</sup>

The Central Park Conservancy was founded in 1980 in the aftermath of Central Park's decline in the 1960s and 1970s.<sup>[2]</sup> Initially devoted to fundraising for projects to restore and improve the park, it took over the park's management duties in 1998.<sup>[3]</sup> The organization has invested more than \$800 million toward the restoration and enhancement of Central Park since its founding.<sup>[4]</sup> With an endowment of over \$200 million, consisting of contributions from residents, corporations, and foundations,<sup>[5]</sup> the Conservancy provides 75 percent of the Park's \$65 million annual operating budget and is responsible for all basic care of the park.<sup>[6]</sup> The Conservancy also provides maintenance support and staff training programs for other public parks in New York City, and has assisted with the development of new parks, such as the High Line and Brooklyn Bridge Park.<sup>[7]</sup><sup>[45–46]</sup>

Central Park Conservancy



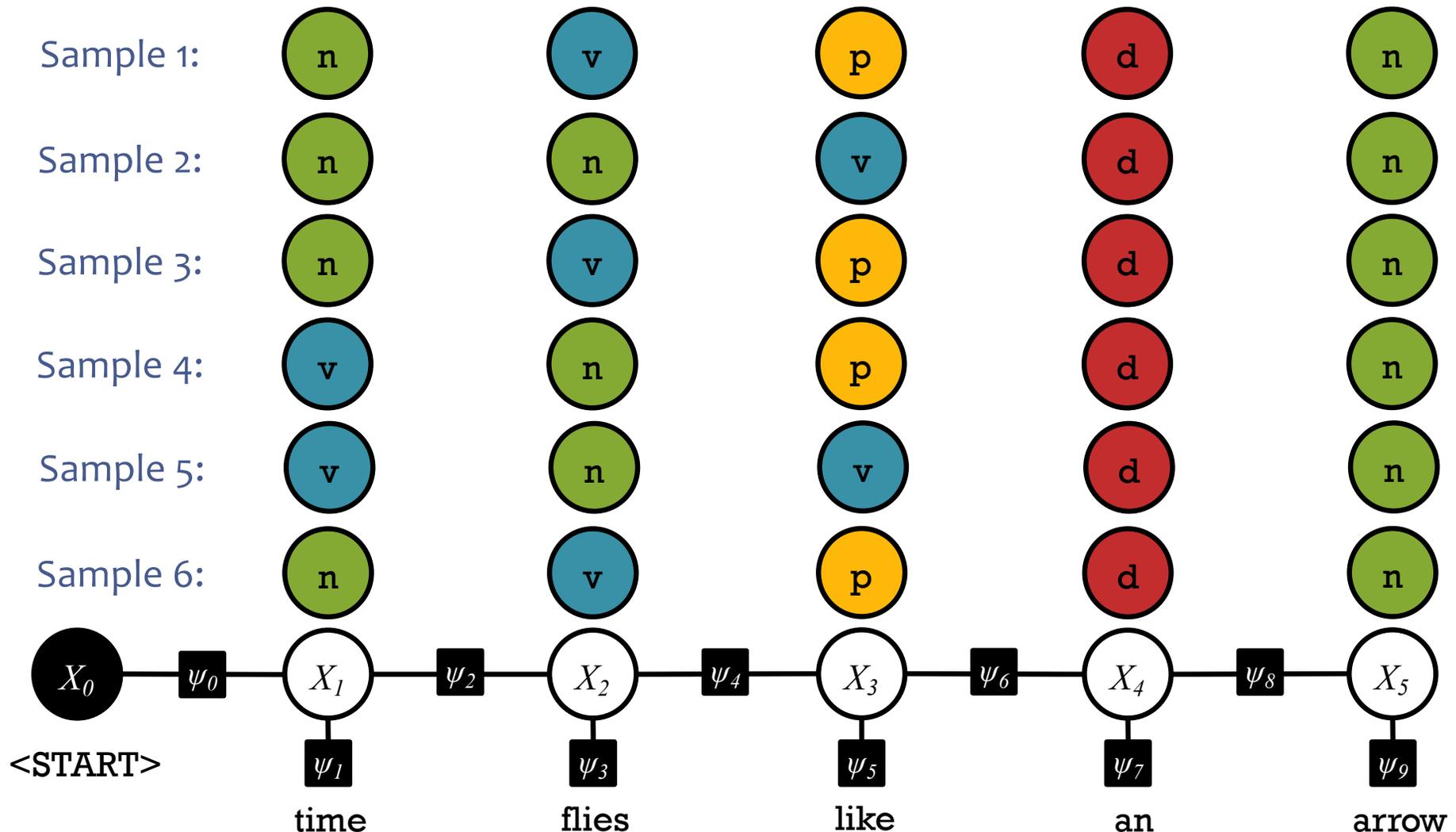
|                    |  |
|--------------------|--|
| <b>Predecessor</b> | Central Park Task Force, Central Park Community Fund   |
| <b>Founded</b>     | 1980   |
| <b>Founders</b>    | Elizabeth Barlow Rogers, William Sperry Beinecke, Ed Koch, Gordon Davis, Andrew Stein ( <i>ex officio</i> )  |
| <b>Tax ID no.</b>  | 13-3022855   |
| <b>Location</b>    | New York City, U.S.  |
| <b>Coordinates</b> | <span><span><span><span><span>40.76424°N</span> <span>73.97169°W</span></span></span><span><span>﻿</span> / <span>﻿</span></span><span><span>40.76424°N 73.97169°W</span><span><span>﻿</span> / <span>40.76424; -73.97169</span></span></span></span></span> |
| <b>Area served</b> | Central Park   |
| <b>Key people</b>  | Elizabeth W. Smith (President & CEO)   |
| <b>Website</b>     | <a href="http://centralparknyc.org">centralparknyc.org</a>   |

A Visual Language for Variables and Interactions

# **GRAPHICAL MODELS**

# Sampling from a Joint Distribution

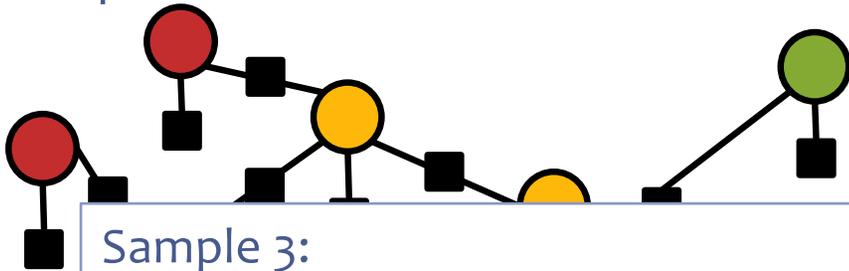
A **joint distribution** defines a probability  $p(x)$  for each assignment of values  $x$  to variables  $X$ . This gives the **proportion** of samples that will equal  $x$ .



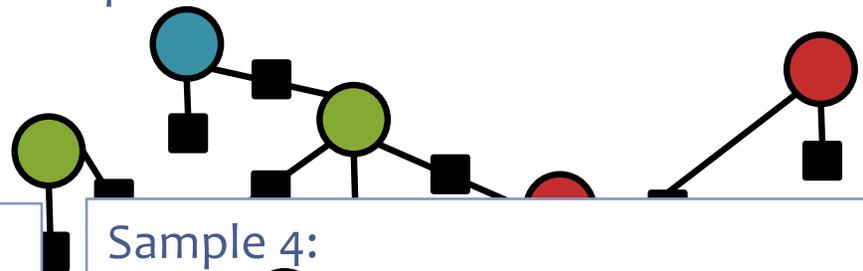
# Sampling from a Joint Distribution

A **joint distribution** defines a probability  $p(x)$  for each assignment of values  $x$  to variables  $X$ . This gives the **proportion** of samples that will equal  $x$ .

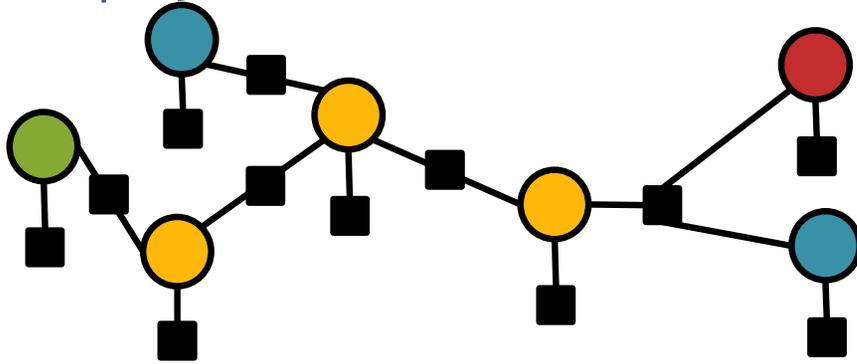
Sample 1:



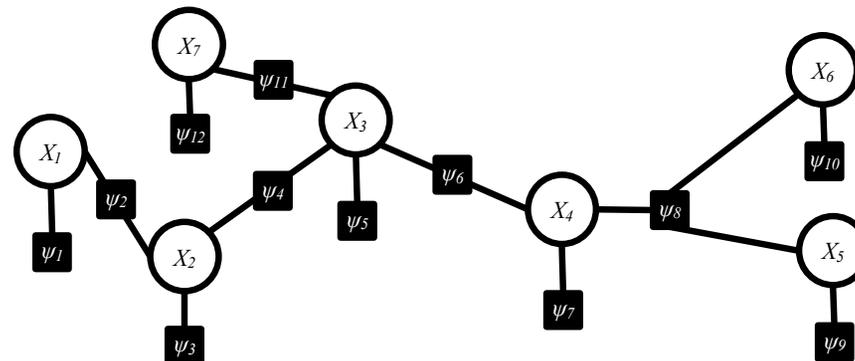
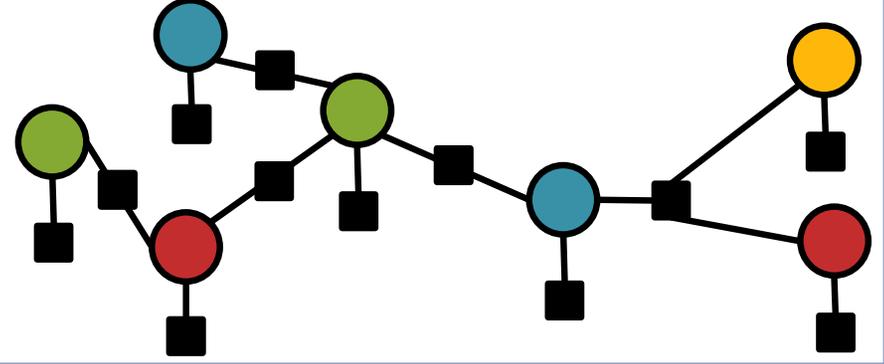
Sample 2:



Sample 3:

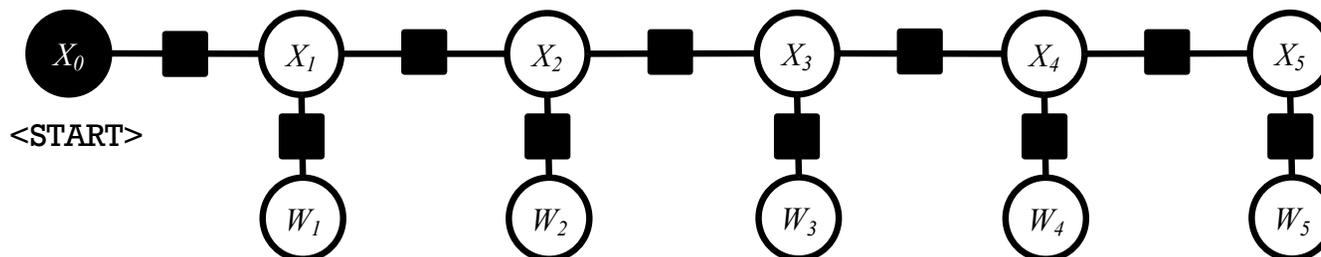
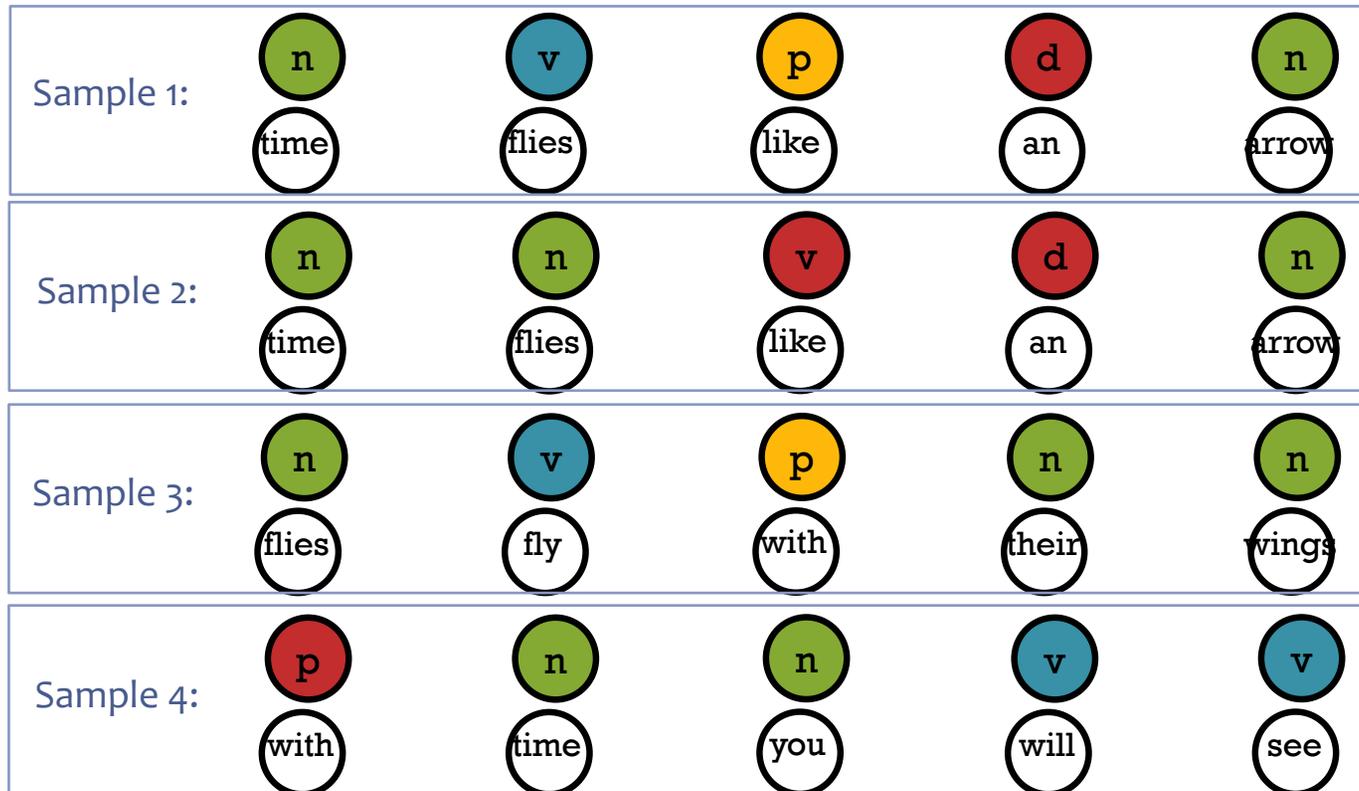


Sample 4:



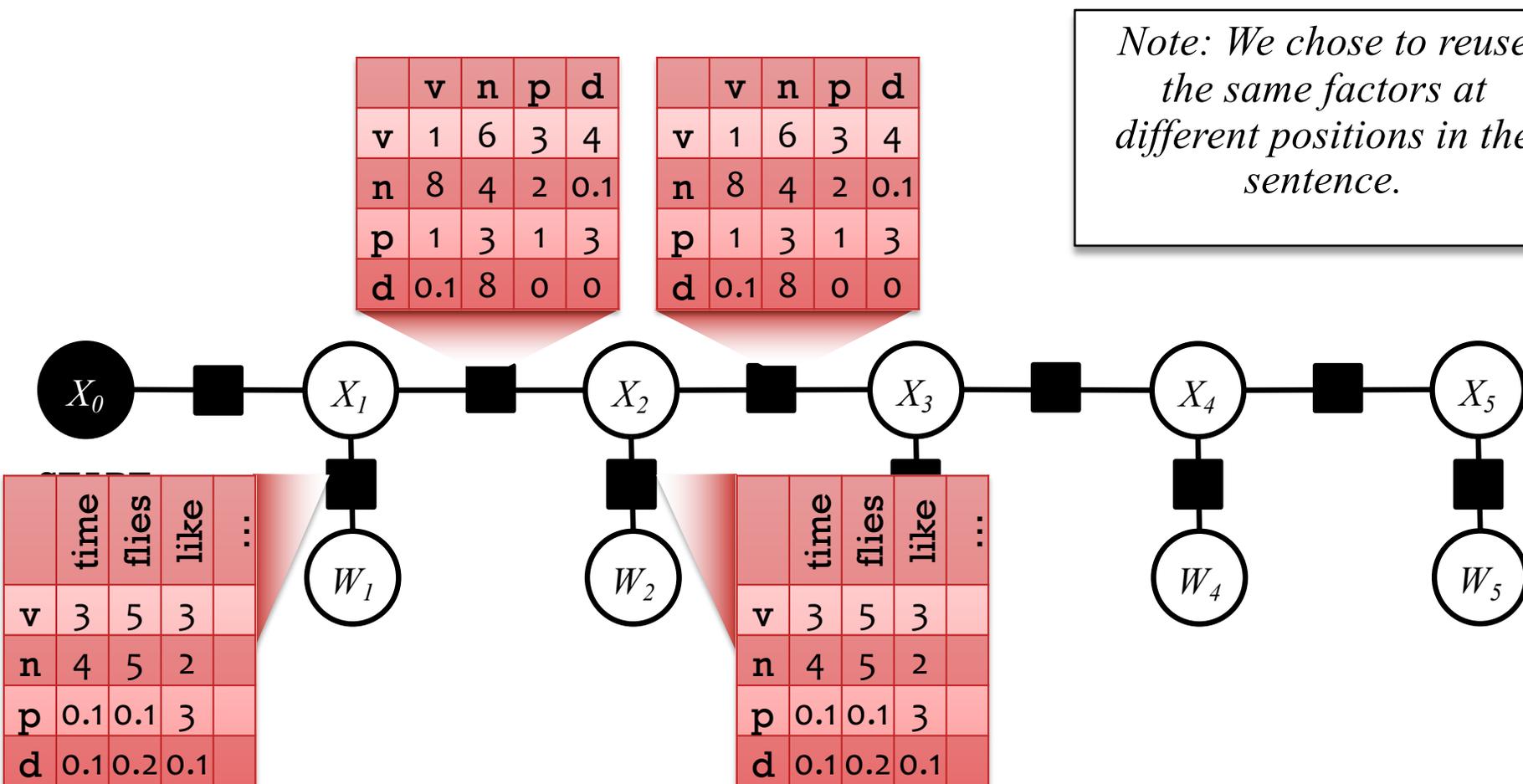
# Sampling from a Joint Distribution

A **joint distribution** defines a probability  $p(x)$  for each assignment of values  $x$  to variables  $X$ . This gives the **proportion** of samples that will equal  $x$ .



# Factors have local opinions ( $\geq 0$ )

Each black box looks at some of the tags  $X_i$  and words  $W_i$

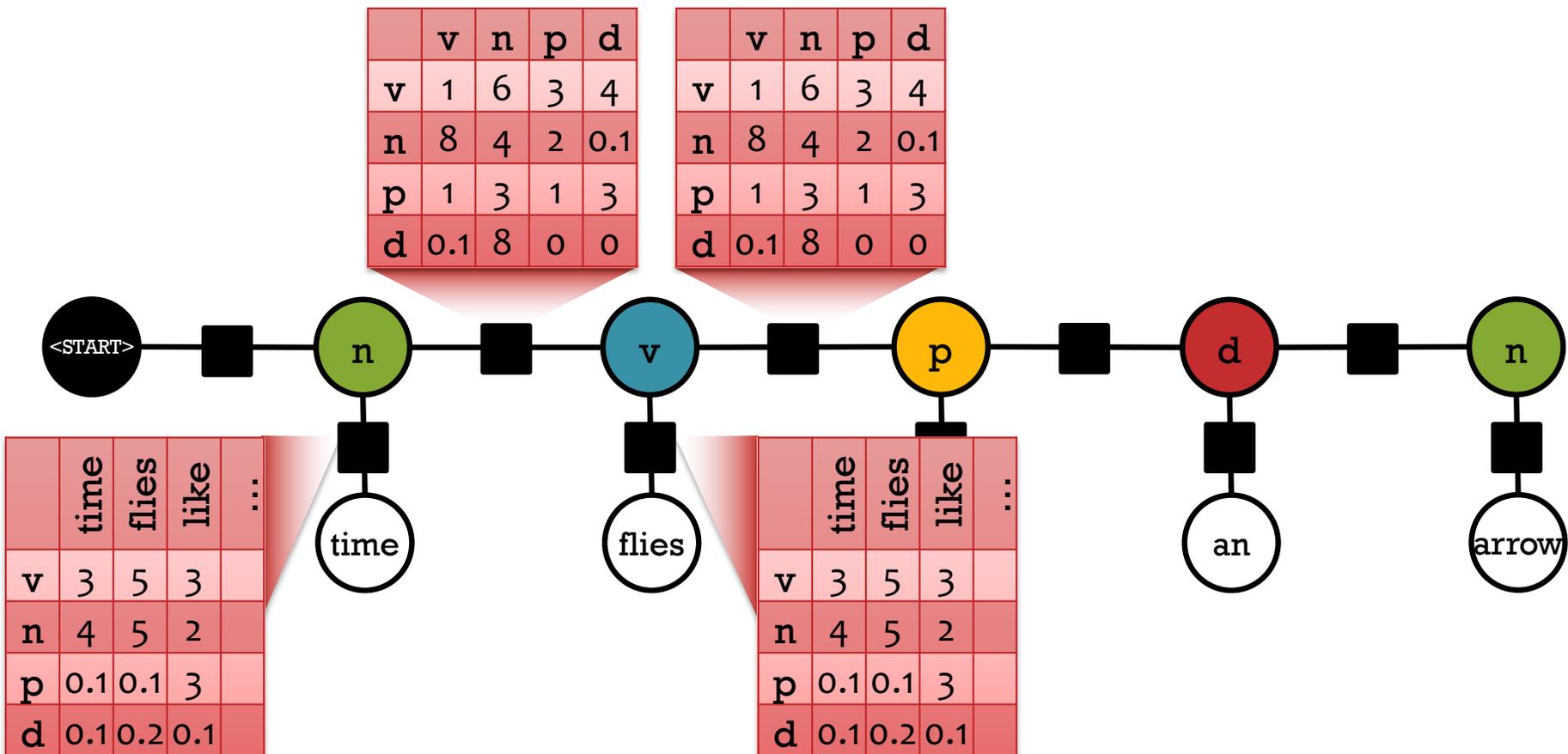


*Note: We chose to reuse the same factors at different positions in the sentence.*

# Factors have local opinions ( $\geq 0$ )

Each black box looks at some of the tags  $X_i$  and words  $W_i$

$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = ?$$



# Global probability = product of local opinions

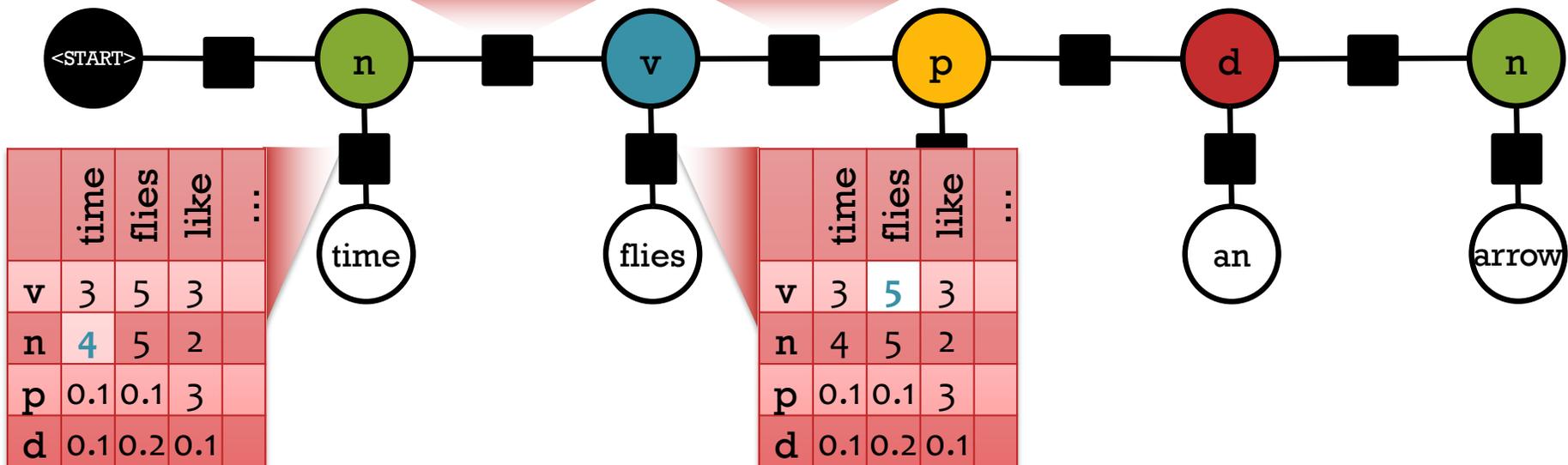
Each black box looks at some of the tags  $X_i$  and words  $W_i$

$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$

|   | v   | n | p | d   |
|---|-----|---|---|-----|
| v | 1   | 6 | 3 | 4   |
| n | 8   | 4 | 2 | 0.1 |
| p | 1   | 3 | 1 | 3   |
| d | 0.1 | 8 | 0 | 0   |

|   | v   | n | p | d   |
|---|-----|---|---|-----|
| v | 1   | 6 | 3 | 4   |
| n | 8   | 4 | 2 | 0.1 |
| p | 1   | 3 | 1 | 3   |
| d | 0.1 | 8 | 0 | 0   |

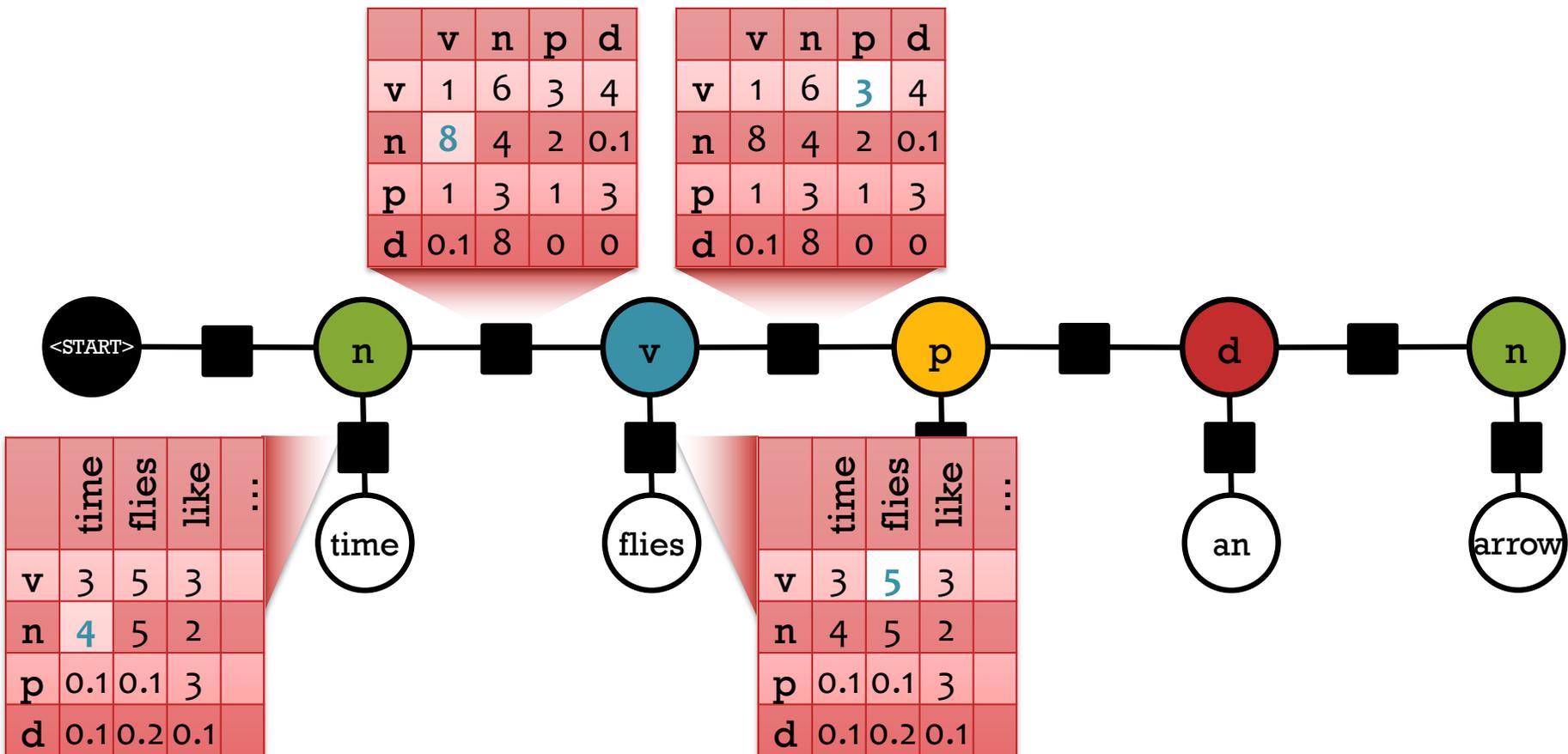
*Uh-oh! The probabilities of the various assignments sum up to  $Z > 1$ .  
So divide them all by  $Z$ .*



# Markov Random Field (MRF)

Joint distribution over tags  $X_i$  and words  $W_i$   
 The individual factors aren't necessarily probabilities.

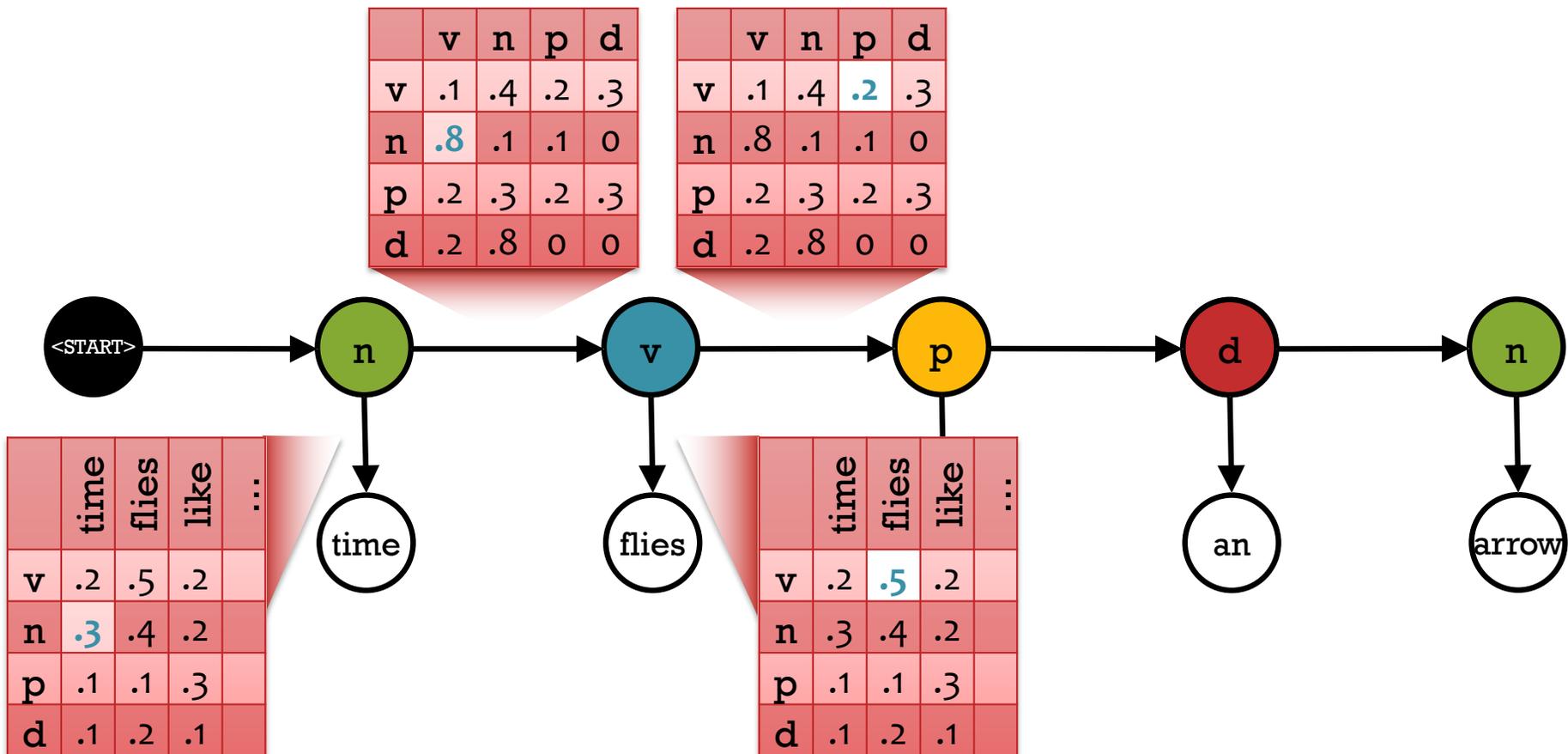
$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



# Hidden Markov Model

But sometimes we *choose* to make them probabilities.  
 Constrain each row of a factor to sum to one. Now  $Z = 1$ .

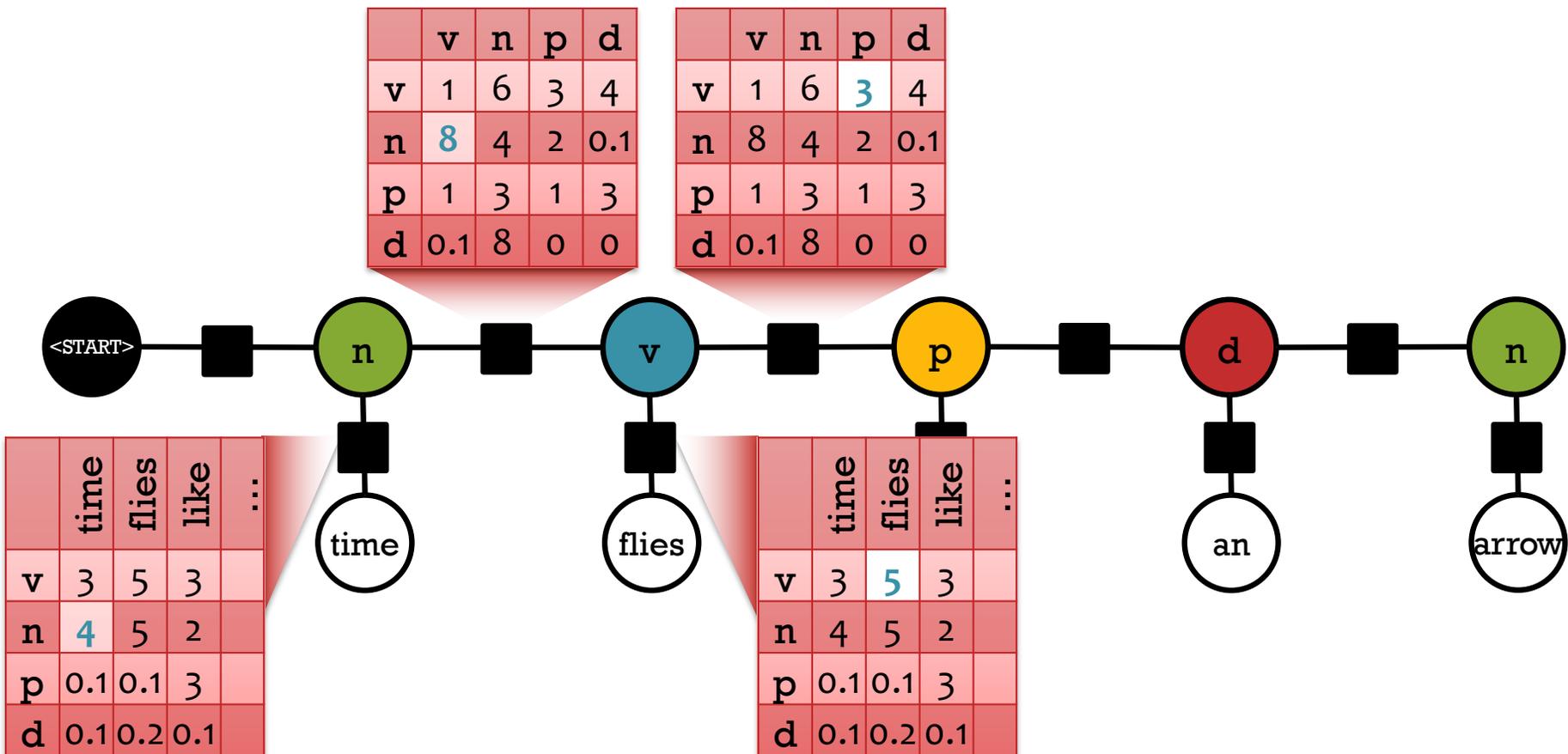
$$p(n, v, p, d, n, \text{time}, \text{flies}, \text{like}, \text{an}, \text{arrow}) = \frac{1}{Z} (.3 * .8 * .2 * .5 * \dots)$$



# Markov Random Field (MRF)

Joint distribution over tags  $X_i$  and words  $W_i$

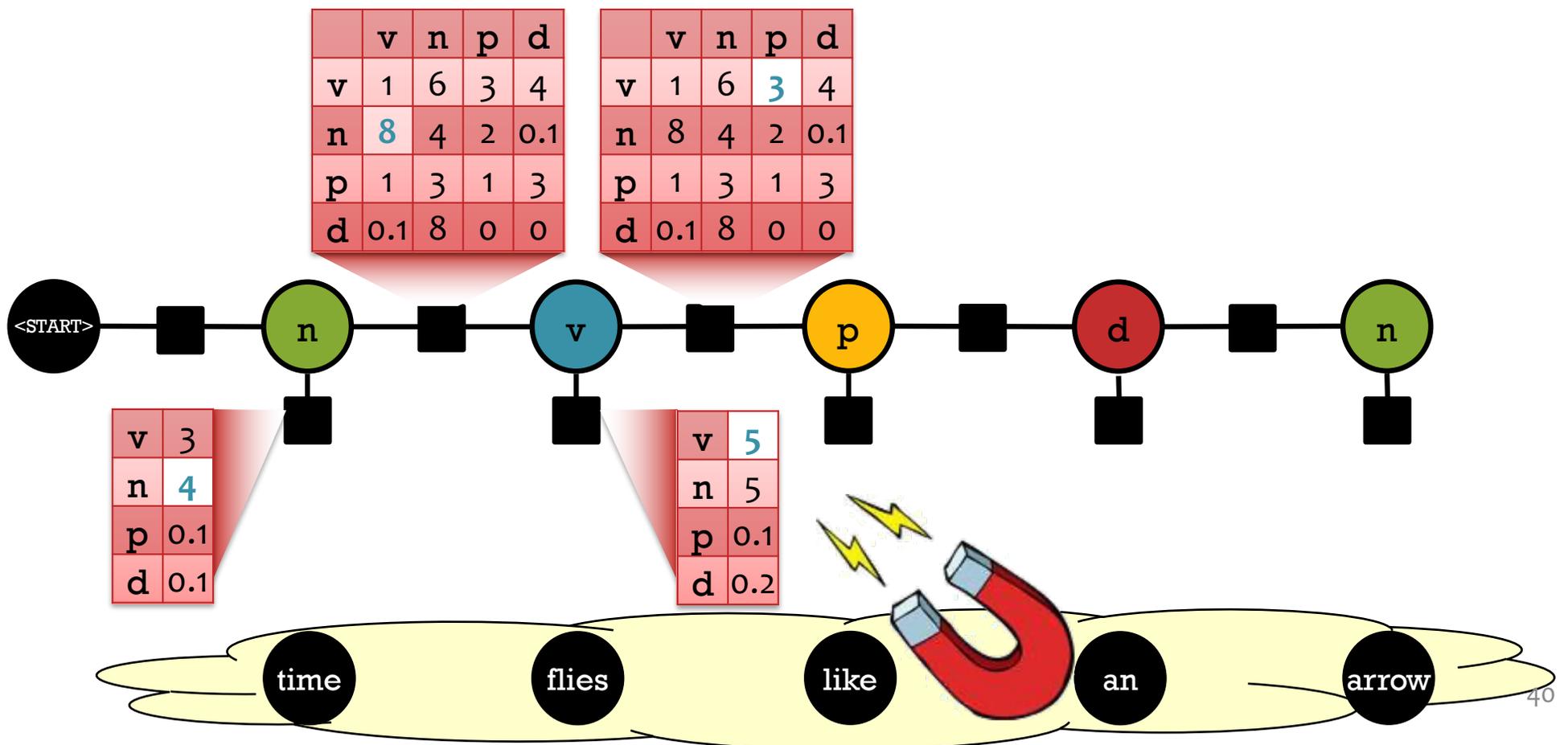
$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



# Conditional Random Field (CRF)

Conditional distribution over tags  $X_i$  given words  $w_i$ .  
The factors and  $Z$  are now specific to the sentence  $w$ .

$$p(n, v, p, d, n \mid \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



# Exercise: Wikipedia Infoboxes

## Question:

Suppose you want to populate missing infobox fields.

1. What are the variables?
2. What are the interactions?

## Answer:

### Central Park Conservancy

From Wikipedia, the free encyclopedia

Coordinates: 40.76424°N 73.97169°W﻿ / ﻿

The **Central Park Conservancy** is a private, nonprofit park conservancy that manages Central Park under a contract with the City of New York and NYC Parks. The conservancy employs most maintenance and operations staff in the park. It effectively oversees the work of both the private and public employees under the authority of the publicly appointed Central Park administrator, who reports to the parks commissioner and the conservancy's president.<sup>[1]</sup>

The Central Park Conservancy was founded in 1980 in the aftermath of Central Park's decline in the 1960s and 1970s.<sup>[2]</sup> Initially devoted to fundraising for projects to restore and improve the park, it took over the park's management duties in 1998.<sup>[3]</sup> The organization has invested more than \$800 million toward the restoration and enhancement of Central Park since its founding.<sup>[4]</sup> With an endowment of over \$200 million, consisting of contributions from residents, corporations, and foundations,<sup>[5]</sup> the Conservancy provides 75 percent of the Park's \$65 million annual operating budget and is responsible for all basic care of the park.<sup>[6]</sup> The Conservancy also provides maintenance support and staff training programs for other public parks in New York City, and has assisted with the development of new parks, such as the High Line and Brooklyn Bridge Park.<sup>[7]</sup><sup>[45–46]</sup>

#### Central Park Conservancy

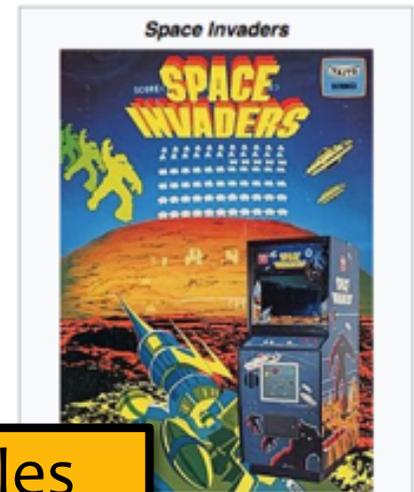


|                    |   |
|--------------------|---|
| <b>Predecessor</b> | Central Park Task Force, Central Park Community Fund  |
| <b>Founded</b>     | 1980  |
| <b>Founders</b>    | Elizabeth Barlow Rogers, William Sperry Beinecke, Ed Koch, Gordon Davis, Andrew Stein ( <i>ex officio</i> )   |
| <b>Tax ID no.</b>  | 13-3022855  |
| <b>Location</b>    | New York City, U.S.   |
| <b>Coordinates</b> | <span><span><span><span><span>40.76424°N</span> <span>73.97169°W</span></span></span><span><span>﻿</span> / <span>﻿</span></span><span><span></span><span><span></span></span></span></span></span> |
| <b>Area served</b> | Central Park  |
| <b>Key people</b>  | Elizabeth W. Smith (President & CEO)  |
| <b>Website</b>     | <span>centralparknyc.org</span> <sup>[?]</sup>  |

# Exercise: Wikipedia Infoboxes

## Populating Wikipedia Infoboxes

| Gryan Miers                 |  |               |
|-----------------------------|--|---------------|
| Personal Information        |  |               |
| Date of birth               | 30 March 1999 (age 20)                             |               |
| Original team(s)            | Grovedale (GFL)                                    |               |
| Draft                       | No. 57, 2017 national draft                        |               |
| Debut                       | Round 1, 2019, Geelong vs. Collingwood, at the MCG |               |
| Height                      | 178 cm (5 ft 10 in)                                |               |
| Weight                      | 78 kg (172 lb)                                     |               |
| Position(s)                 | Small forward                                      |               |
| Club Information            |  |               |
| Current club                | Geelong  |               |
| Number                      | 32   |               |
| Playing career <sup>1</sup> |  |               |
| Years                       | Club   | Games (Goals) |
| 2019                        | Geelong  | 1 (0)         |



Promotional flyer

Taito

JP: Taito  
NA: Midway  
EU: Midway<sup>[1]</sup>  
AU: Leisure & Allied Industries<sup>[2]</sup>  
Atari, Inc. (home)  
Tomohiro Nishikado  
Arcade, Atari 2600, Atari 5200, Atari 8-bit, MSX  
JP: June 1978<sup>[3]</sup>  
NA: July 1978  
Fixed shooter  
Single-player, 2 players alternating  
Upright, cocktail<sup>[4]</sup>  
Taito 8080<sup>[5]</sup>  
8080 @ 2 MHz<sup>[6]</sup>  
SN76477 @ 1.9968 MHz  
Fujitsu MB14241,<sup>[8]</sup>  
monochrome raster, vertical orientation, 224x256 resolution<sup>[9]</sup>

Q: Why do interactions appear between variables in this example?

A: Consider the test time setting:

- Author writes a new article (vector  $x$ )
- Infobox is empty
- ML system must populate all fields (vector  $y$ ) at once
- Interactions that were seen (i.e. vector  $y$ ) at training time are unobserved at test time – so we wish to model them

# **INFERENCE PROBLEMS**

# Inference

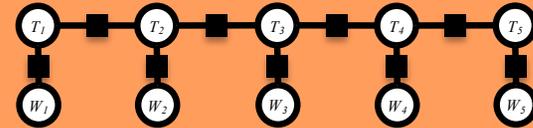
## 1. Data

$$\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$$

|           |            |            |           |           |            |
|-----------|------------|------------|-----------|-----------|------------|
| Sample 1: | n<br>time  | v<br>flies | p<br>like | d<br>an   | n<br>frown |
| Sample 2: | n<br>time  | n<br>flies | v<br>like | d<br>an   | n<br>frown |
| Sample 3: | n<br>flies | v<br>fly   | p<br>with | n<br>heir | n<br>ring  |
| Sample 4: | p<br>with  | n<br>time  | n<br>you  | v<br>will | v<br>see   |

## 2. Model

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$



## 3. Objective

$$\ell(\boldsymbol{\theta}; \mathcal{D}) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)} | \boldsymbol{\theta})$$

## 5. Inference

### 1. Marginal Inference

$$p(\mathbf{x}_C) = \sum_{\mathbf{x}': \mathbf{x}'_C = \mathbf{x}_C} p(\mathbf{x}' | \boldsymbol{\theta})$$

### 2. Partition Function

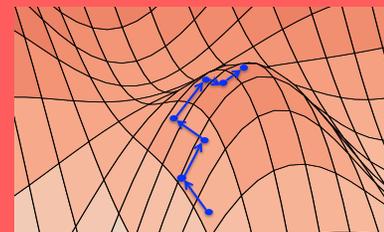
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

### 3. MAP Inference

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta})$$

## 4. Learning

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



# 5. Inference

## 1. Marginal Inference (#P-Hard)

Compute marginals of variables and cliques

$$p(x_i) = \sum_{\mathbf{x}' : x'_i = x_i} p(\mathbf{x}' | \boldsymbol{\theta}) \quad \Bigg| \quad p(\mathbf{x}_C) = \sum_{\mathbf{x}' : \mathbf{x}'_C = \mathbf{x}_C} p(\mathbf{x}' | \boldsymbol{\theta})$$

## 2. Partition Function (#P-Hard)

Compute the normalization constant

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

## 3. MAP Inference (NP-Hard)

Compute variable assignment with highest probability

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta})$$

## 4. Sampling (cf. convergence, variance)

Draw a sample variable assignment

$$\mathbf{x} \sim p(\cdot | \boldsymbol{\theta})$$

# Q&A

**Q:** But in **deep learning** we don't need to solve these inference problems, right?

**A:** Wrong... it's not that we don't *need* to solve them, it's that we can't!

**X** Questions you *could* ask your RNN-LM or seq2seq model:

- X** 1. What is the probability of the 7<sup>th</sup> token being 'zebra' (marginal inference)
- X** 2. For an unnormalized model, what is the normalization constant? (partition function)
- X** 3. What is the most probable output sequence? (MAP inference)
- ✓** 4. Give me 10 samples from the distribution.

# **SYLLABUS HIGHLIGHTS**

# Syllabus Highlights

The syllabus is located on the course webpage:

<http://708.mlcourse.org>  ...cs.cmu.edu...

The **course policies** are **required** reading.

# Syllabus Highlights

- **Grading:** 45% homework, 35% project, 15% quizzes, 5% participation
- **Quizzes:** ~3 quizzes, during lecture/recitation time
- **Homework:** ~5 assignments
  - 6 grace days for homework assignments
  - Late submissions: 80% day 1, 60% day 2, 40% day 3, 20% day 4
  - No submissions accepted after 4 days w/o extension
  - Extension requests: see syllabus
- **Recitations:** Fridays, same time/place as lecture (optional, interactive sessions)
- **Readings:** required, online PDFs, recommended for after lecture
- **Technologies:**
  - Piazza (discussion),
  - Gradescope (homework),
  - Google Forms (polls),
  - Gather.Town (office hours),
  - Zoom (livestream),
  - Panopto (video recordings)
- **Academic Integrity:**
  - Collaboration encouraged, but must be documented
  - Solutions must always be written independently
  - No re-use of found code / past assignments
  - Severe penalties (i.e.. failure)
- **Office Hours:** posted on Google Calendar on “People” page

# Lectures

- You should ask lots of questions
  - Interrupting (by raising a hand, turning on your video, and waiting to be called on) to ask your question is strongly encouraged
  - Use the chat to ask questions in real time (TAs will be monitoring the chat and will either answer or interrupt the instructor)
  - Asking questions later on Piazza is also great
- When I ask a question...
  - I want you to answer
  - Even if you don't answer, think it through as though I'm about to call on you
- Interaction improves learning (both in-class and at my office hours)

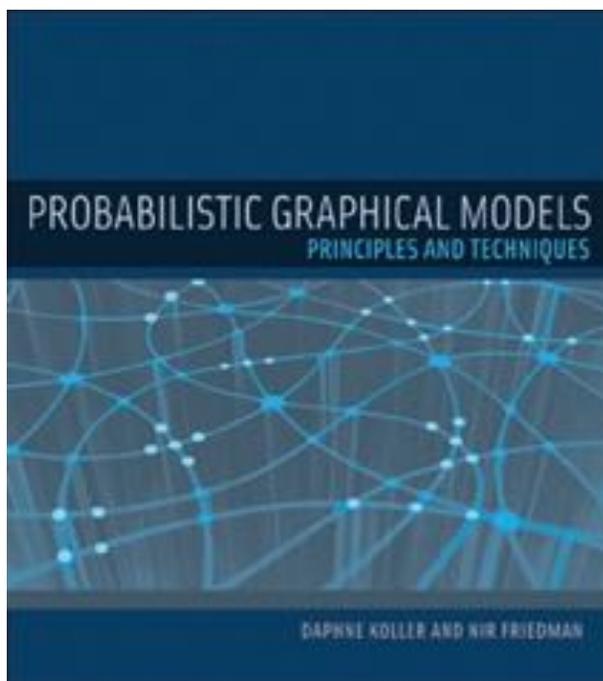
# Homework

There will be 5 homework assignments during the semester. The assignments will consist of both conceptual and programming problems.

|     | Main Topic                       | Implementation                  | Application Area            | Type                  |
|-----|----------------------------------|---------------------------------|-----------------------------|-----------------------|
| HW1 | Properties of graphical models   | NA                              | NA                          | written               |
| HW2 | Marginal inference and MLE       | RNN + Tree CRF                  | natural language processing | written + programming |
| HW3 | MAP inference and structured SVM | CNN + Ising model               | computer vision             | written + programming |
| HW4 | MCMC                             | word embeddings + Gibbs sampler | topic modeling              | written + programming |
| HW5 | Variational Inference            | variational inference           | TBD                         | written + programming |

# Textbooks

You are not *required* to read a textbook, but Koller & Friedman is a thorough reference text that includes a lot of the topics we cover.

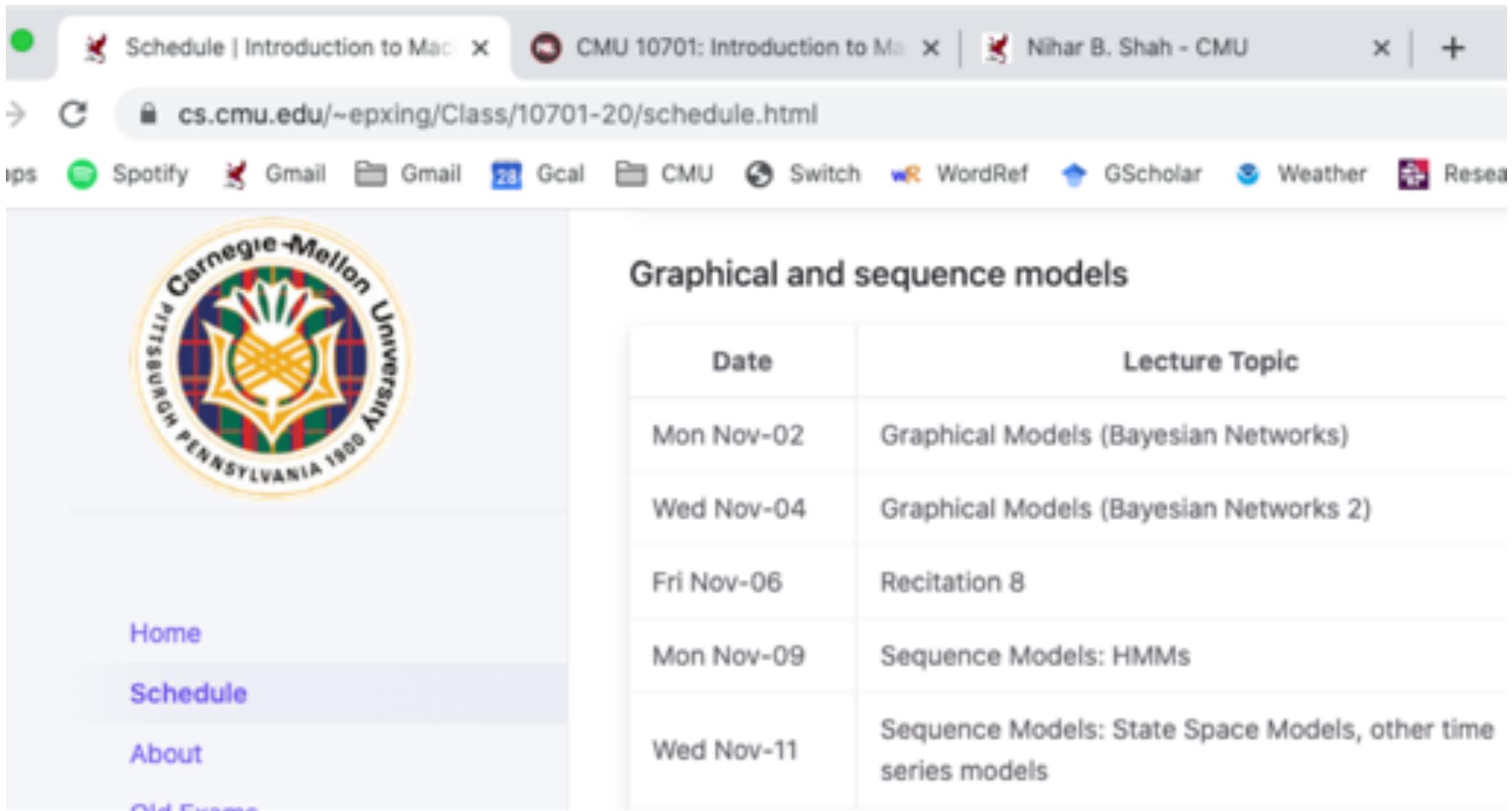


# Prerequisites

## What they are:

1. Introductory machine learning.  
(i.e. 10-701, 10-715)
2. Significant experience programming in a general programming language.
  - Some homework may require you to use Python, so you will need to at least be **proficient in the basics of Python.**
3. College-level probability, calculus, linear algebra, and discrete mathematics.

# Prerequisites



Home

Schedule

About

Old Exams

## Graphical and sequence models

| Date       | Lecture Topic   |
|------------|---|
| Mon Nov-02 | Graphical Models (Bayesian Networks)                          |
| Wed Nov-04 | Graphical Models (Bayesian Networks 2)                        |
| Fri Nov-06 | Recitation 8  |
| Mon Nov-09 | Sequence Models: HMMs   |
| Wed Nov-11 | Sequence Models: State Space Models, other time series models |

# Prerequisites

|     |        |    |                  |  |
|-----|--------|----|------------------|--|
| Wed | Apr-01 | 21 | Graphical Models | Graphical Models 1: Representing joint probability distributions       |
|     |        |    |                  |  |
|     |        |    |                  |  |
|     |        |    |                  |  |
| Fri | Apr-03 |    | Poster Session   | Poster Session (Online)  |
| Mon | Apr-06 | 22 | Graphical Models | Graphical Models 2: Inference and Supervised learning                  |
| Wed | Apr-08 | 23 | Graphical Models | Graphical Models 3: Learning Bayes nets, EM, semi-supervised learning. |
| Fri | Apr-10 |    |                  | No recitation  |
| Mon | Apr-13 | 24 | Graphical Models | Graphical Models 4: Mixture Model Clustering, D-Separation             |

# Prerequisites

|        |  |                 |
|--------|--|-----------------|
| Nov 9  | Boosting   | SB C            |
| Nov 11 | Online learning  | SB C            |
| Nov 16 | Semi-supervised learning, Active learning, Multi-armed bandits | Trans<br>Multi- |
| Nov 18 | Reinforcement learning   | Surve           |
| Nov 23 | Graphical models   | Graph           |
| Nov 25 | No class (Thanksgiving break)                                  |                 |

# Project

- Goals:
  - Explore an interesting problem in your domain of interest
  - Employ graphical models for a structured prediction task or latent variable modeling
  - For example:
    - compare models under the same inference technique
    - compare inference methods on the same model
    - compare learning methods on the same model
  - Deeper understanding of methods in real-world application
  - Work in teams of 3 – 4 students
- Milestones: (*due in 2<sup>nd</sup> half of semester*)
  1. Team Formation
  2. Proposal
  3. Midway Report
  4. Final Report
  5. Poster Presentation

Q&A