



**NUS**  
National University  
of Singapore

| **Computing**

# **CS4248: Natural Language Processing**

## **Lecture 7 — Sequences**

# Announcements



# Outline

- **Overview: Sequence Tasks**
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- Hidden Markov Models (HMM)
  - Basic setup and components
  - Core HMM tasks
    - Model Learning
    - Likelihood computation
    - Viterbi decoding

# Motivation

- So far: Bag-of-Words (BoW) models
  - Bag = whole document (e.g., Naive Bayes, Vector Space Model)
  - Bag = context of a word (e.g., PPMI and Word2Vec embeddings)
- Natural language: word order matters! (can vary greatly between languages, though)

*Bob kills mosquitoes using the book of Hamlet*

vs.

*Hamlet kills Bob using the book of mosquitoes*

*The food tastes good and does not look bad*

vs.

*The food tastes bad and does not look good*

Same words, very different meanings!

# Motivation — Example: English

- Fundamental rules word order
  - Subject—Verb—Object (SVO)
  - Adjectives only (immediately) before nouns
  - ...and many more
- "Informal" rules — e.g.: order of adjectives
  - Rule: opinion→size→physical quality→shape→age→color→origin→material→type→purpose

*I saw a beautiful, old, blue, German car.*

**vs.**

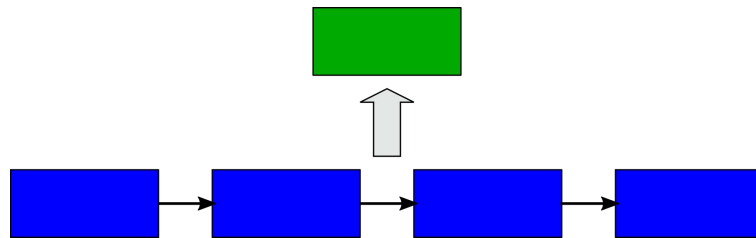
*I saw a German, blue, beautiful, old car.*

# Types of Sequence Tasks

- Sequence classification

(Many-to-One,  $N \rightarrow 1$ )

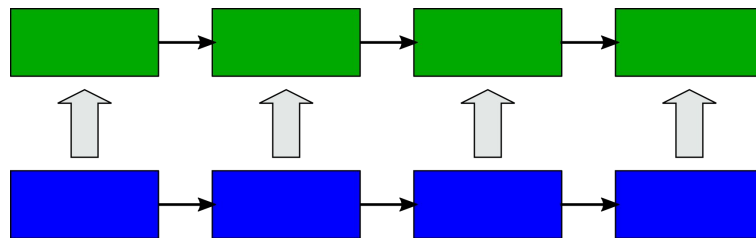
- Sentiment analysis
- Document categorization



- Sequence labeling/tagging

(Many-to-Many,  $N \rightarrow N$ )

- Part-of-Speech Tagging
- Named Entity Recognition

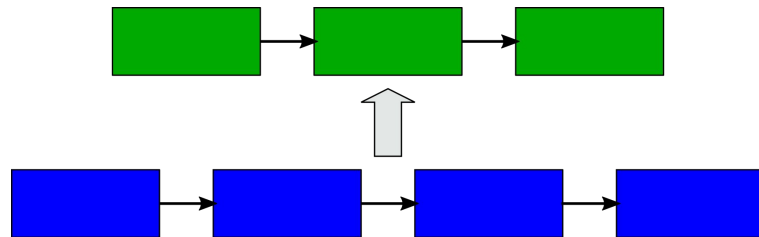


# Types of Sequence Tasks

- Sequence translation

(Many-to-Many,  $N \rightarrow M$ )

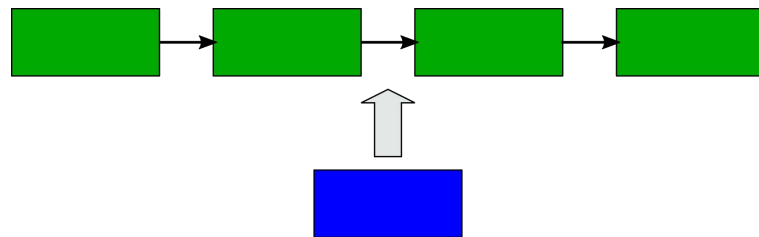
- Machine translation
- Sentence simplification
- Text summarization



- Sequence generation

(One-to-Many,  $1 \rightarrow N$ )

- Image captioning



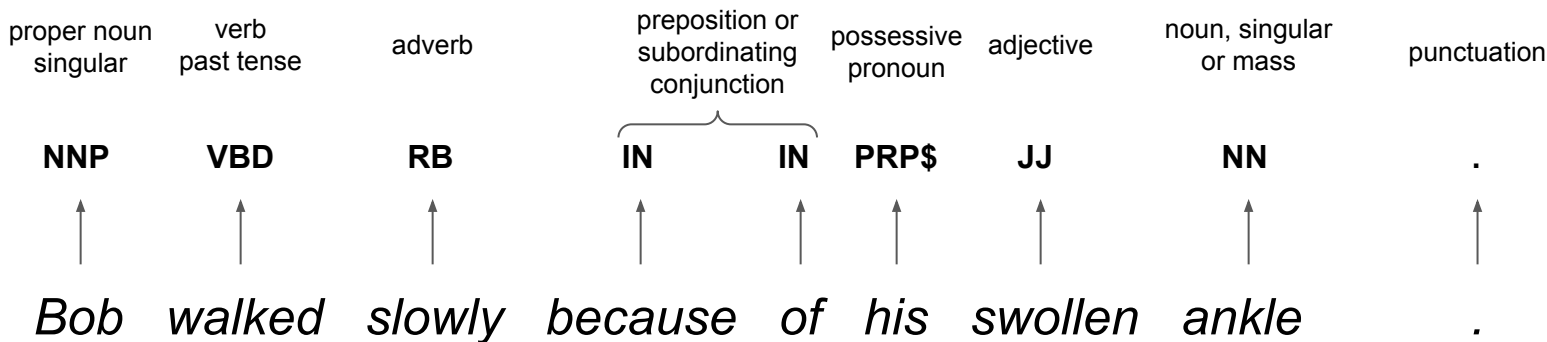
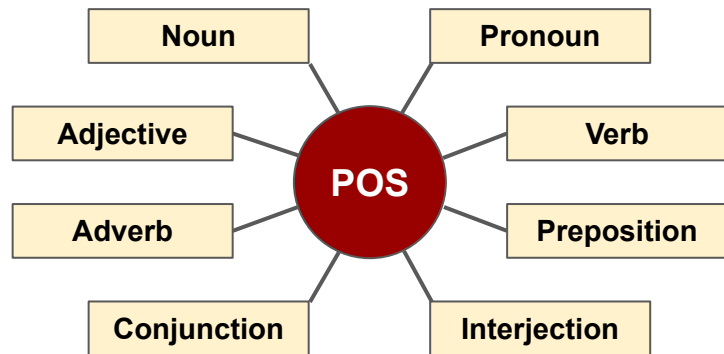
# Outline

- Overview: Sequence Tasks
- **POS Tagging**
  - **What are Parts of Speech?**
  - Why is this task important and challenging?
- Hidden Markov Models (HMM)
  - Basic setup and components
  - Core HMM tasks
    - Model Learning
    - Likelihood computation
    - Viterbi decoding



# Part-of-Speech Tagging

- **Part of Speech** (also: word class or syntactic category)
  - Each word belongs one or more of these classes
  - English: 8 main parts of speech / word classes (many additional classes and subclasses considered in practice)
- **Part-of-Speech (POS) tagging**
  - Assign each word in a text a part of speech (duh!)



# Penn Treebank Tag-Set

## Base set: 36 POS tags

|      |                                   |                                 |
|------|-----------------------------------|---------------------------------|
| CC   | Coordinating conjunction          | <i>and, or</i>                  |
| CD   | Cardinal number                   | <i>1, 2, 3, one, two, three</i> |
| DT   | Determiner                        | <i>the, a, an, any, some</i>    |
| EX   | Existential there                 |                                 |
| FW   | Foreign word                      |                                 |
| IN   | Preposition / subord. conjunction | <i>in, into, whether, if</i>    |
| JJ   | Adjective                         | <i>cleaner, nice</i>            |
| JJR  | Adjective (comparative)           | <i>cleaner, nicer</i>           |
| JJS  | Adjective (superlative)           | <i>cleanes, nicest</i>          |
| LS   | List item marker                  |                                 |
| MD   | Modal                             | <i>can, could, may</i>          |
| NN   | Noun (singular or mass)           | <i>machine, computer, air</i>   |
| NNS  | Noun (plural)                     | <i>machines, computers</i>      |
| NNP  | Proper noun (singular)            | <i>Clementi Mall</i>            |
| NNPS | Proper noun (plural)              | <i>Americas</i>                 |
| PDT  | Predeterminer                     | <i>all, both, half</i>          |
| POS  | Possessive ending                 | <i>'s</i>                       |
| PRP  | Personal pronoun                  | <i>him. himself, we</i>         |

|      |                                    |                                 |
|------|------------------------------------|---------------------------------|
| PP\$ | Possessive pronoun                 | <i>her, our, ours</i>           |
| RB   | Adverb                             | <i>quickly, swiftly</i>         |
| RBR  | Adverb (comparative)               | <i>further, greater, more</i>   |
| RBS  | Adverb (superlative)               | <i>furthest, greatest, most</i> |
| RP   | Particle                           | <i>across, up</i>               |
| SYM  | Symbol                             | <i>=, +, &amp;</i>              |
| TO   | to                                 | <i>to</i>                       |
| UH   | Interjection                       | <i>shucks, heck, oops</i>       |
| VB   | Verb (base form)                   | <i>be, assign, run</i>          |
| VBD  | Verb (past tense)                  | <i>was, assigned, ran</i>       |
| VBG  | Verb (gerund / present participle) | <i>being, assigning</i>         |
| VBN  | Verb (past participle)             | <i>been, assigned</i>           |
| VBP  | Verb (non-3rd pers. sing. present) | <i>am, are</i>                  |
| VBZ  | Verb (3rd pers. sing. present)     | <i>is</i>                       |
| WDT  | wh-determiner                      | <i>that, which, what</i>        |
| WP   | wh-pronoun                         | <i>that, which, whom</i>        |
| WP\$ | Possessive wh-pronoun              | <i>whose</i>                    |
| WRB  | wh-adverb                          | <i>how, however, why</i>        |

# Penn Treebank Tag-Set

## Extended set: 12 tags for punctuations and special symbols

|    |                             |             |
|----|-----------------------------|-------------|
| #  | Pound sign                  | #           |
| \$ | Dollar sign                 | \$          |
| .  | Sentence-final punctuation  | . ? !       |
| :  | Sentence-middle punctuation | : ; ... - — |
| ,  | Comma                       | ,           |
| (  | Left bracket character      | ( [ { <     |
| )  | Right bracket character     | ) ] } >     |
| "  | Straight double quote       |             |
| `  | Left open single quote      |             |
| `` | Left open double quote      |             |
| '  | Right close single quote    |             |
| "  | Right close double quote    |             |

# Part of Speech — Two Broad Categories

- **Closed class words**

- Small, fixed membership — reasonably easy to enumerate
- Generally, short function words that “structure” sentences
- Examples: prepositions, pronouns, participles, determiners, conjunctions, etc.

- **Open class words**

- Impossible to completely enumerate
- New words continuously being invented, borrowed, etc.
- For most languages: nouns, verbs, adjectives, adverbs

# Outline

- Overview: Sequence Tasks
- **POS Tagging**
  - What are Parts of Speech?
  - **Why is this task important and challenging?**
- Hidden Markov Models (HMM)
  - Basic setup and components
  - Core HMM tasks
    - Model Learning
    - Likelihood computation
    - Viterbi decoding

# POS Tagging — Why is it Important?

- Very useful or crucial for many NLP downstream tasks
  - Named Entity recognition (typically comprised of nouns and proper nouns)
  - Information extractions (e.g., verbs indicate relations between entities)
  - Parsing (information of word classes useful before creating parse trees)
  - Speech synthesis/recognition (e.g., noun "DIScount" vs. verb "disCOUNT")
  - Authorship Attribution (e.g., relative frequencies of nouns, verbs, adjectives, etc.)
  - Machine Translation (e.g., reordering of adjectives and nouns)

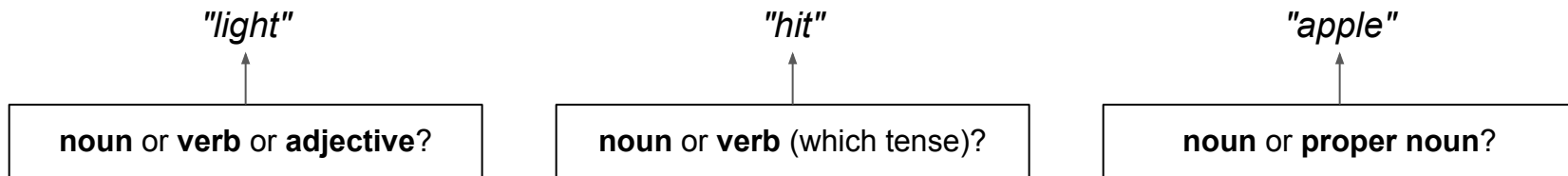
→ POS tagging: important low-level NLP task

# Quick Quiz



# POS Tagging — Why is it Difficult? And How Difficult?

- Our common problem: **Ambiguity**
  - Many common words have multiple meanings → multiple POS



- Often ambiguous even with additional context  
(even humans can often no agree on the correct labeling!)





# POS Tagging — Why is it Difficult? And How Difficult?

- POS tagging in English

- ~85% of word types are unambiguous (e.g., "quickly" is always an adverb, "Alice" is always a proper noun)
- ~15% of word types are ambiguous — but those are quite common!

- ➔ 55-65% of word tokens are ambiguous

- Ambiguous = 2 or more possible POS tags
- Results depend on text corpus

# Quick Quiz



# POS Tagging — Baseline Algorithm

- Most straightforward approach

- Label each word with its most frequent POS tag
- Label unknown words as nouns (most common open world class)

→ **Result: ~92% accuracy** (vs. ~97-98% accuracy for SOTA methods)

- Doesn't sound so bad right?
- 2 main problems:

(1) **Imbalanced errors**

- High accuracy due to common/frequent unambiguous words (e.g., "the", "a/an", "and", "or")
- Many of these words also often not that interesting for downstream NLP tasks

(2) **Downstream error propagation**

- POS tagging as low-level NLP task → errors quickly propagate up

# POS Tagging — Unsupervised Algorithms

- Basic intuition

- Utilize words with unambiguous POS tags → **anchor words**
- Observe patterns to group words into clusters of the same word class
- Use anchor words to assign clusters (and each containing word) to a POS tag

- Practical considerations

- No need for hand-labeled text corpora (only lexicon of anchor words required)
- Poorer performance compared to supervised methods

# POS Tagging — Supervised Methods

- Require hand-labeled text corpus

- Used as input training data for supervised models
- Challenging for low-resource languages  
(i.e., languages lacking in large, annotated datasets)

- Popular models (all yielding quite similar SOTA results)

- Hidden Markov Models (HMM)
- Conditional Random Fields (CRF)
- Neural sequence models (RNNs, Transformers)
- Large language models (e.g., BERT)

- Accuracies have reached "human ceiling"  
(i.e., POS taggers as good as human annotators)
- POS tagging considered a solved task  
for high-resource languages (e.g. English)
- Limitations: low-resource languages  
and special application domains

# Quick Quiz



# In-Lecture Activity (5 min)



# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- **Hidden Markov Models (HMM)**
  - **Basic setup and components**
  - Core HMM tasks
    - Model Learning
    - Likelihood computation
    - Viterbi decoding



# Markov Chains

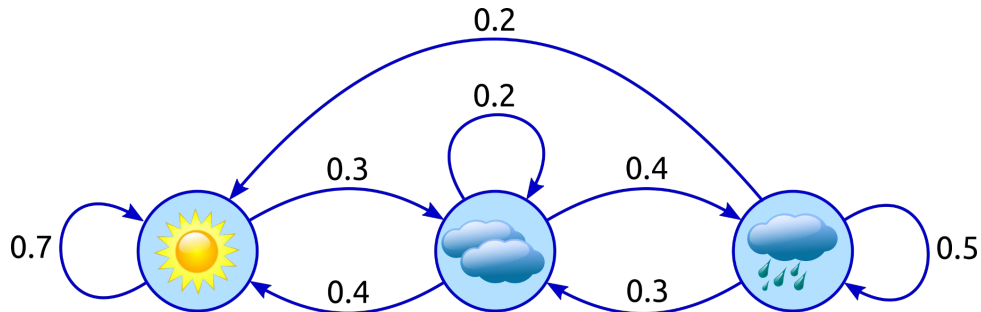
- Markov Chain

- Models transitions between a set of states using transition probabilities (captured by a transition matrix A)
- Transition only depends on current state (Markov assumption)
- Sequence: series of transitions

- Example: "Daily Weather"

- 3 states: *sunny*, *cloudy*, *rainy*

Example question: "What is the probability of getting a 5 sunny days in a row?"



$$A = \begin{matrix} & \text{☀️} & \text{☁️} & \text{☔️} \\ \text{☀️} & \begin{bmatrix} 0.7 & 0.3 & 0.0 \end{bmatrix} \\ \text{☁️} & \begin{bmatrix} 0.4 & 0.2 & 0.4 \end{bmatrix} \\ \text{☔️} & \begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix} \end{matrix}$$

# Markov Chain → Hidden Markov Models

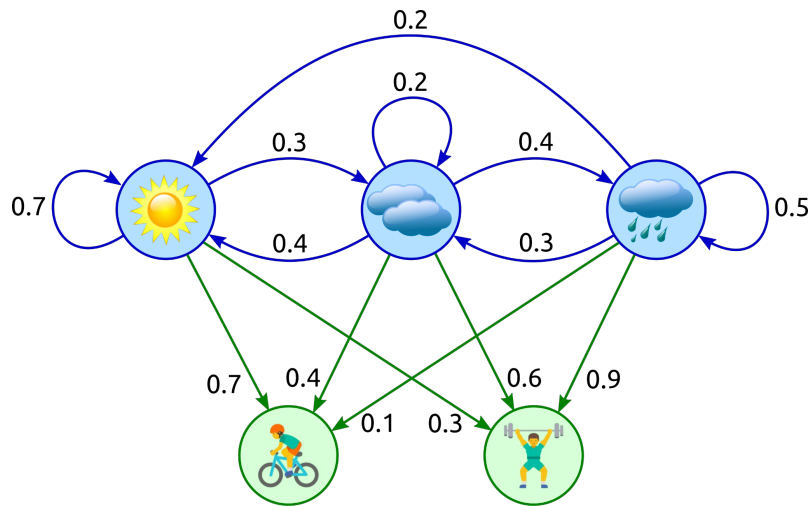
- Hidden Markov Models (HMM)

- States are hidden (i.e., not directly observable)
- Observable variables that depend on the states

- Example: "Exercising Routine"

- 3 hidden(!) states: *sunny*, *cloudy*, *rainy*
- 2 observed activities: *biking*, *lifting*  
(with the activity depending on the weather)

Example question: "Given that Chris went first 3 days lifting and then 3 days biking, what was the most likely weather over the last 6 days?"



# HMM — Components

Finite Set of **states**  $S = \{s_1, s_2, \dots, s_N\}$

Sequence of states

$Q = q_1, q_2, q_3, \dots, q_T$ , with  $q_t \in S$

Finite set of **symbols**  $V = \{v_1, v_2, \dots, v_M\}$

Sequence of observations

$O = o_1, o_2, o_3, \dots, o_T$ , with  $o_t \in V$

**Transition probability matrix**  $A$

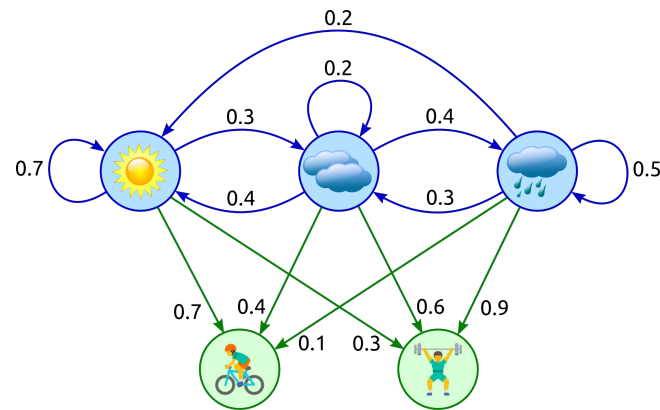
$A = \{a_{ij}\}$ ,  $a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$

**Observation / emission probability matrix**  $B$

$B = \{b_i(o_k)\}$ ,  $b_i(o_k) = P(o_t = v_k | q_t = s_i)$

**Initial state distribution**  $\pi$

$\pi = \{\pi_i\}$ ,  $\pi_i = P(q_1 = s_i)$



$$A = \{a_{ij}\} = \begin{matrix} \begin{matrix} \text{Sun} & \text{Cloud} & \text{Rain} \end{matrix} \\ \begin{matrix} \text{Sun} \\ \text{Cloud} \\ \text{Rain} \end{matrix} \end{matrix} \begin{bmatrix} 0.7 & 0.3 & 0.0 \\ 0.4 & 0.2 & 0.4 \\ 0.2 & 0.3 & 0.5 \end{bmatrix}$$

$$B = \{b_i(o_k)\} = \begin{matrix} \begin{matrix} \text{Sun} & \text{Person on a bicycle} \end{matrix} \\ \begin{matrix} \text{Cloud} \\ \text{Rain} \end{matrix} \end{matrix} \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \\ 0.1 & 0.9 \end{bmatrix}$$

# HMM — Probabilities (annotated)

## Transition probability matrix $A$

$$A = \{a_{ij}\}, \quad a_{ij} = P(q_{t+1} = s_j \mid q_t = s_i)$$

Probability of transitioning from state  $s_i$  to  $s_j$  at any time  $t$

$$\sum_j^N a_{ij} = 1 \quad \forall i$$

## Observation / emission probability matrix $B$

$$B = \{b_i(o_k)\}, \quad b_i(o_k) = P(o_t = v_k \mid q_t = s_i)$$

Probability of state  $s_i$  generating output  $v_k$  at any time  $t$

$$\sum_k^M b_i(o_k) = 1, \quad \forall i$$

## Initial state distribution $\pi$

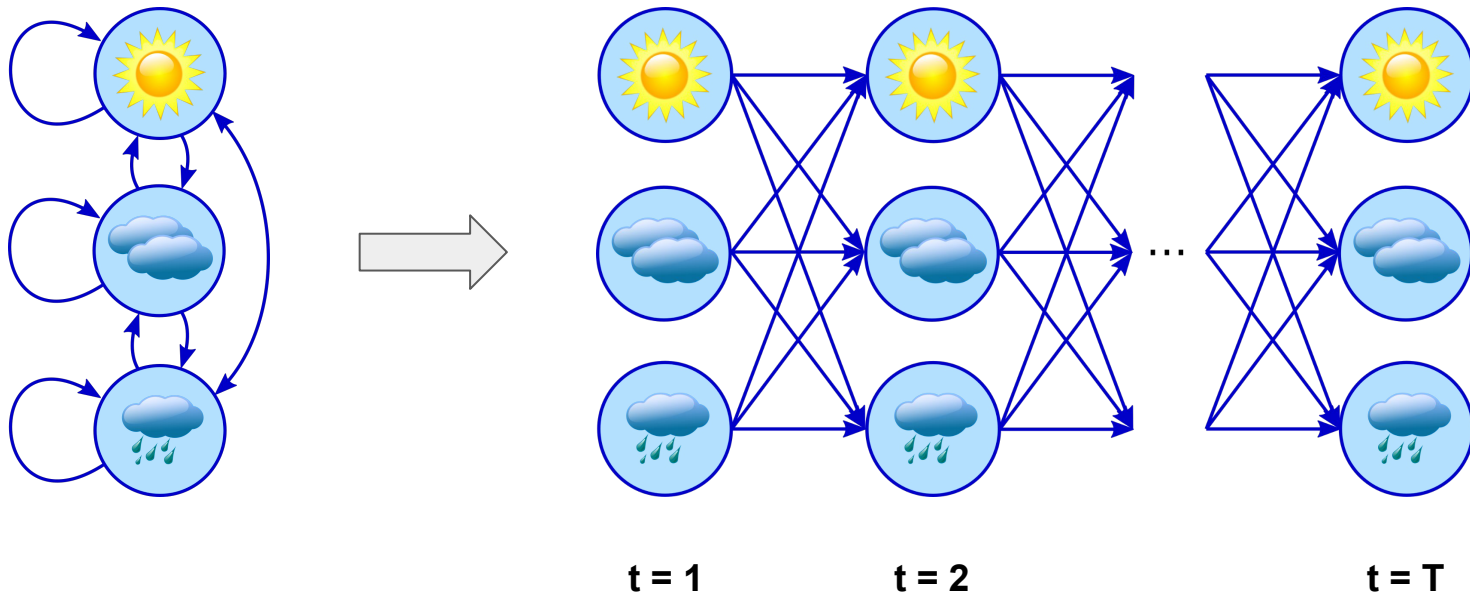
$$\pi = \{\pi_i\}, \quad \pi_i = P(q_1 = s_i)$$

Probability of sequence starting in state  $s_i$

$$\sum_i^N \pi_i = 1$$

# HMM — Unrolled Representation

**Trellis diagram** — graph representation of all possible states and transitions over time



# Quick Quiz



# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- **Hidden Markov Models (HMM)**
  - Basic setup and components
  - **Core HMM tasks**
    - Model Learning
    - Likelihood computation
    - Viterbi decoding

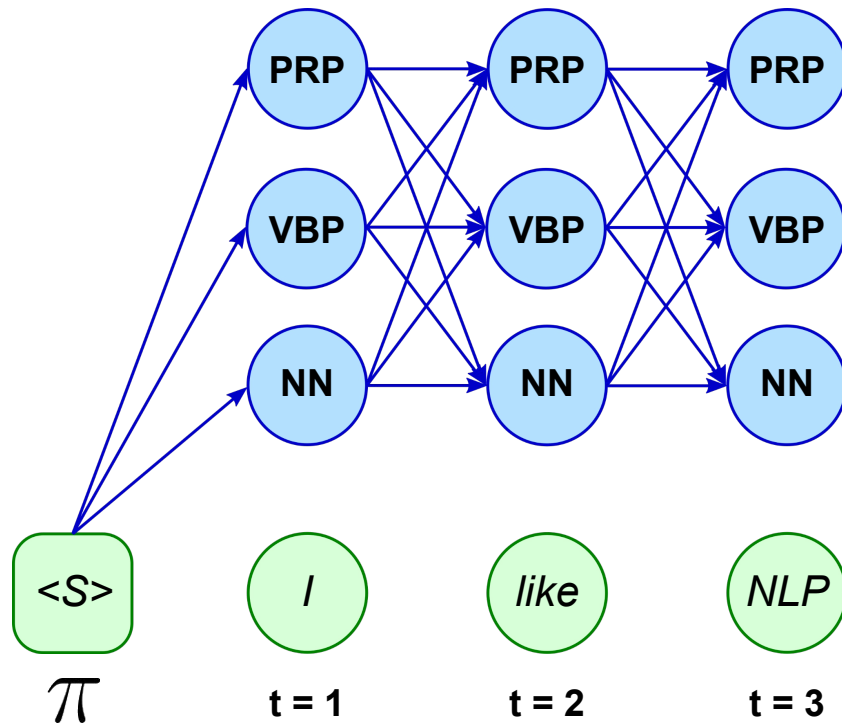
# POS Tagging with HMMs

- Basic setup

- Hidden states  $\rightarrow$  POS tags
- Observations  $\rightarrow$  words

- Example

- 3 states: {PRP, VBP, NN}





# HMM — Core Tasks

## (1) Model Learning

Given corresponding state and observation sequences  $Q$  and  $O$

→ Learn all model parameters, i.e., probabilities  $A$ ,  $B$  and  $\pi$

Training using an annotated dataset

## (2) Likelihood

Given an HMM  $\theta = (A, B, \pi)$

+ an state sequence  $Q$

+ an observation sequence  $O$

→ Calculate the probability  $P(Q, O|\theta)$

Given 2 POS tag sequences for a sentence, compare which is more likely

## (3) Decoding

Given an HMM  $\theta = (A, B, \pi)$  + an observation sequence  $O$

→ Find the most likely sequence of states

Given a sentence, find the most likely POS tags

# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- **Hidden Markov Models (HMM)**
  - Basic setup and components
  - **Core HMM tasks**
    - **Model Learning**
    - Likelihood computation
    - Viterbi decoding

# HMM — Model Learning

**Quick Quiz:** Spotting a familiar issue? How can we address it?

- Calculating probabilities using Maximum Likelihood Estimates

$$\pi_i = P(q_1 = s_i) = \frac{\text{Count}(\langle S \rangle s_i) \leftarrow \text{\#sentences starting with state } s_i}{\text{Count}(\langle S \rangle) \leftarrow \text{\#sentences}}$$

$$a_{ij} = P(q_{t+1} = s_j | q_t = s_i) = \frac{\text{Count}(s_i s_j) \leftarrow \text{\#occurences of state } s_i \text{ followed by state } s_j}{\text{Count}(s_i) \leftarrow \text{\#occurences of state } s_i}$$

$$b_i(o_k) = P(o_t = v_k | q_t = s_i) = \frac{\text{Count}(v_k, s_i) \leftarrow \text{\#occurences of observation } v_k \text{ in state } s_i}{\text{Count}(s_i) \leftarrow \text{\#occurences of state } s_i}$$

# HMM — Model Learning — Side Note

- POS tagging using HMM

- Full supervised task → corpus of words labeled with all the correct POS tags
- Fully "visible" Markov Model  
(we have the state and observation sequences)

"Direct" parameter learning using MLE  
(we just need simple counts)

- Often in other applications

- State sequences  $Q$  are not known
- Impossible to compute simple counts

# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- **Hidden Markov Models (HMM)**
  - Basic setup and components
  - **Core HMM tasks**
    - Model Learning
    - **Likelihood computation**
    - Viterbi decoding

# Likelihood

- Given: HMM  $\theta = (A, B, \pi)$  and

$$O = o_1, o_2, o_3, \dots, o_T$$

$$Q = q_1, q_2, q_3, \dots, q_T$$

- Calculate joint probability  $P(O, Q|\theta)$

$$P(O, Q|\theta) = P(O|Q) \cdot P(Q) = \prod_{i=1}^T P(o_i|q_i) \cdot P(q_i|q_{i-1})$$

emission  
probabilities

transition  
probabilities

with  $P(q_1|q_0) = P(q_1)$

initial state  
probability

# Likelihood — Example

$$P(O, Q | \theta) = P(O | Q) \cdot P(Q) = \prod_{i=1}^T P(o_i | q_i) \cdot P(q_i | q_{i-1})$$

- Given: sentence  $O$ , POS tags  $Q$

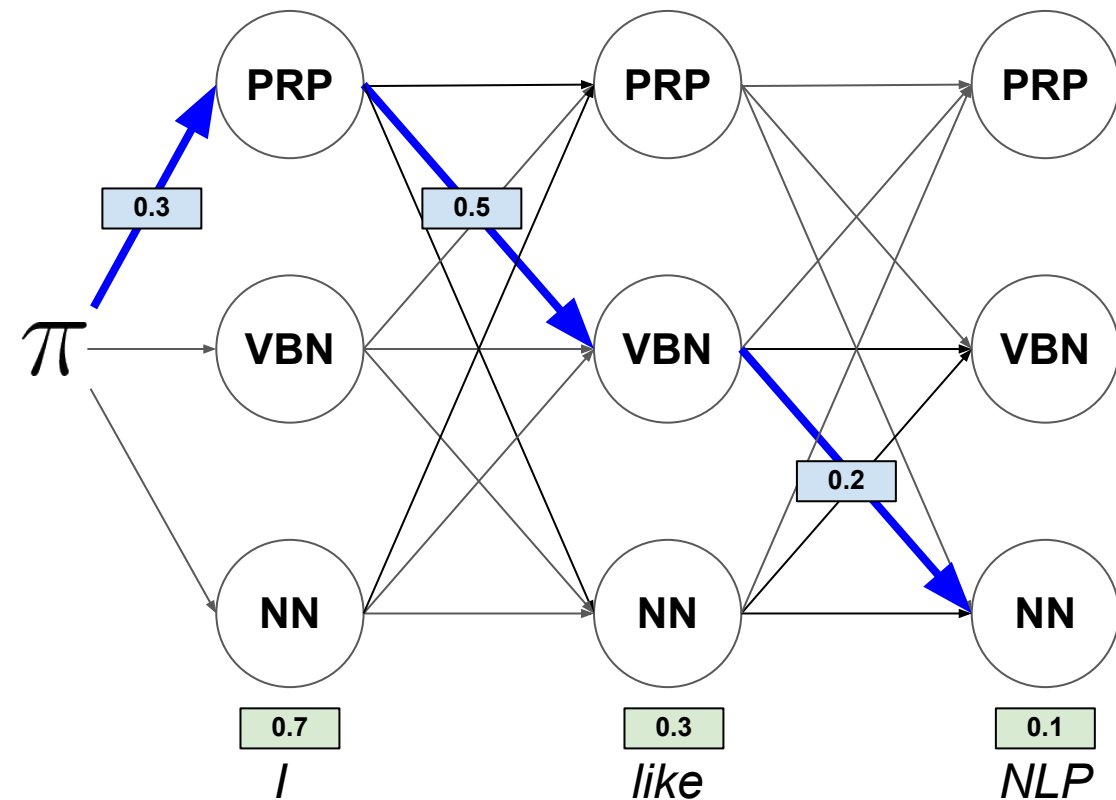
$$\left. \begin{array}{l} O = I, like, NLP \\ Q = PRP, VBN, NN \end{array} \right\} P("I, like, NLP" | PRP-VBN-NN) = ?$$

$$\begin{aligned} P("I, like, NLP" | PRP - VBN - NN) = & P(I | PRP) \cdot P(PR | \langle S \rangle) \cdot \\ & P(like | VBN) \cdot P(VBN | PRP) \cdot \\ & \underbrace{P(NLP | NN) \cdot P(NN | VBN)} \end{aligned}$$

All values can be directly taken from  $A$ ,  $B$ , and  $\pi$

# Likelihood — Example

Visualization using a Trellis diagram → just follow the path

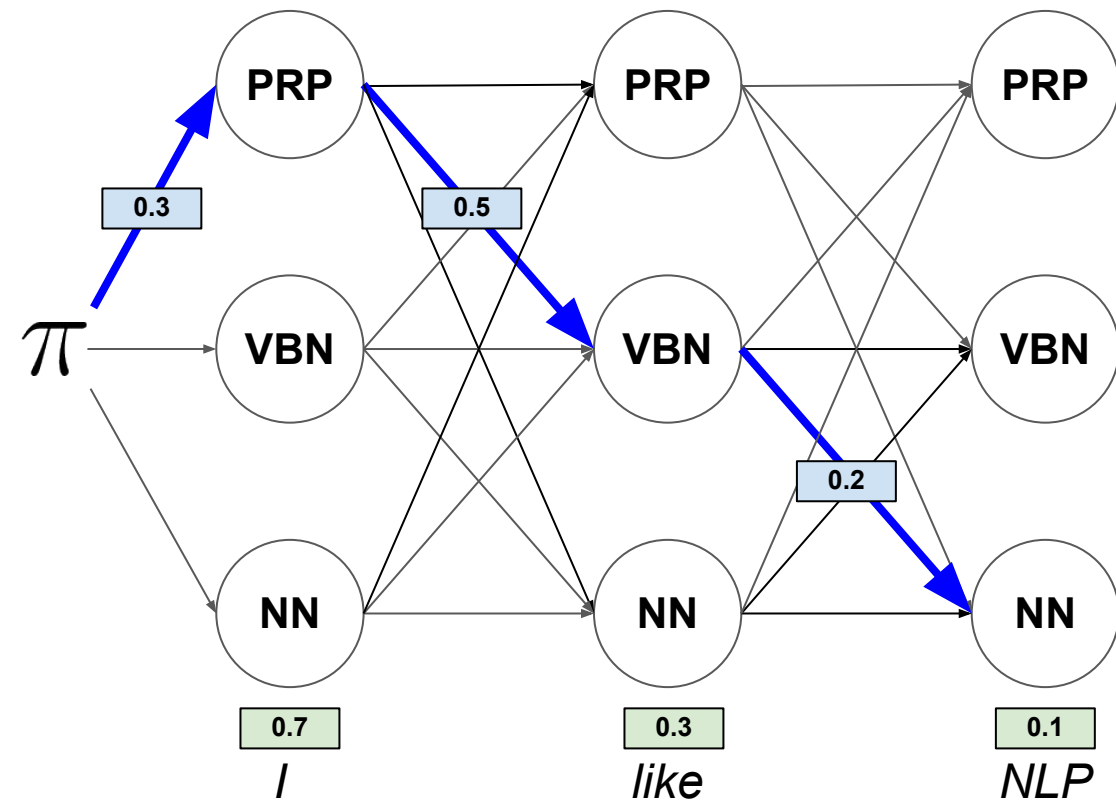


$$\begin{aligned} P("I, like, NLP" | PRP - VBN - NN) = \\ P(I | PRP) \cdot P(PRP | \langle S \rangle) \cdot \\ P(like | VBN) \cdot P(VBN | PRP) \cdot \\ P(NLP | NN) \cdot P(NN | VBN) \end{aligned}$$



# Likelihood — Example

Visualization using a Trellis diagram → just follow the path



$$\begin{aligned}
 P("I, like, NLP" | PRP - VBN - NN) &= \\
 &P(I | PRP) \cdot P(PRP | \langle S \rangle) \cdot \\
 &P(like | VBN) \cdot P(VBN | PRP) \cdot \\
 &P(NLP | NN) \cdot P(NN | VBN)
 \end{aligned}$$

$$\begin{aligned}
 P("I, like, NLP" | PRP - VBN - NN) &= \\
 &0.7 \cdot 0.3 \cdot \\
 &0.3 \cdot 0.5 \cdot \\
 &0.1 \cdot 0.2 \\
 &= 0.00063
 \end{aligned}$$

# Likelihood — Can we decode with it?

- Naive algorithm for decoding (for a given observation sequence  $O$ )
  - Enumerate all possible state sequences  $Q$
  - Compute all joint probabilities  $P(O, Q)$
  - Return state sequence  $Q$  with highest joint probability

→ What is the **runtime** of this algorithm?

# Quick Quiz



# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- **Hidden Markov Models (HMM)**
  - Basic setup and components
  - **Core HMM tasks**
    - Model Learning
    - Likelihood computation
    - **Viterbi decoding**

# Decoding

- Decoding task

- Given an HMM  $\theta = (A, B, \pi)$  + an observation sequence  $O$
- Find the most likely sequence of states  $Q$

$$Q = \operatorname{argmax}_{q_1 \dots q_t} \prod_{i=1}^T P(o_i | q_i) \cdot P(q_i | q_{i-1})$$

Diagram illustrating the components of the HMM decoding equation:

- The term  $P(o_i | q_i)$  is associated with the box labeled "emission probabilities".
- The term  $P(q_i | q_{i-1})$  is associated with the box labeled "transition probabilities".
- The term  $P(q_1 | q_0) = P(q_1)$  is associated with the box labeled "initial state probability".

→ **Dynamic Programming** to avoid checking all possible state sequences


# Viterbi Algorithm — Toy Example

- Oversimplified setup

- 3 POS tags: **DT** (determiner), **NN** (noun), **VB** (verb)
- Let's assume the following HMM

$$\pi = \begin{matrix} & \text{DT} & \text{NN} & \text{VB} \\ \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix} \end{matrix} \quad A = \begin{matrix} & \text{DT} & \text{NN} & \text{VB} \\ \begin{bmatrix} 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix} \begin{matrix} \text{DT} \\ \text{NN} \\ \text{VB} \end{matrix} \end{matrix}$$

**Note:** The rows in B do not up to 1 since B does not capture all words, only those we needs.


$$B = \begin{matrix} & the & fans & love & show \\ \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix} \begin{matrix} \text{DT} \\ \text{NN} \\ \text{VB} \end{matrix} \end{matrix}$$

- Task: Find the most likely sequence of state (i.e., POS tags) for:

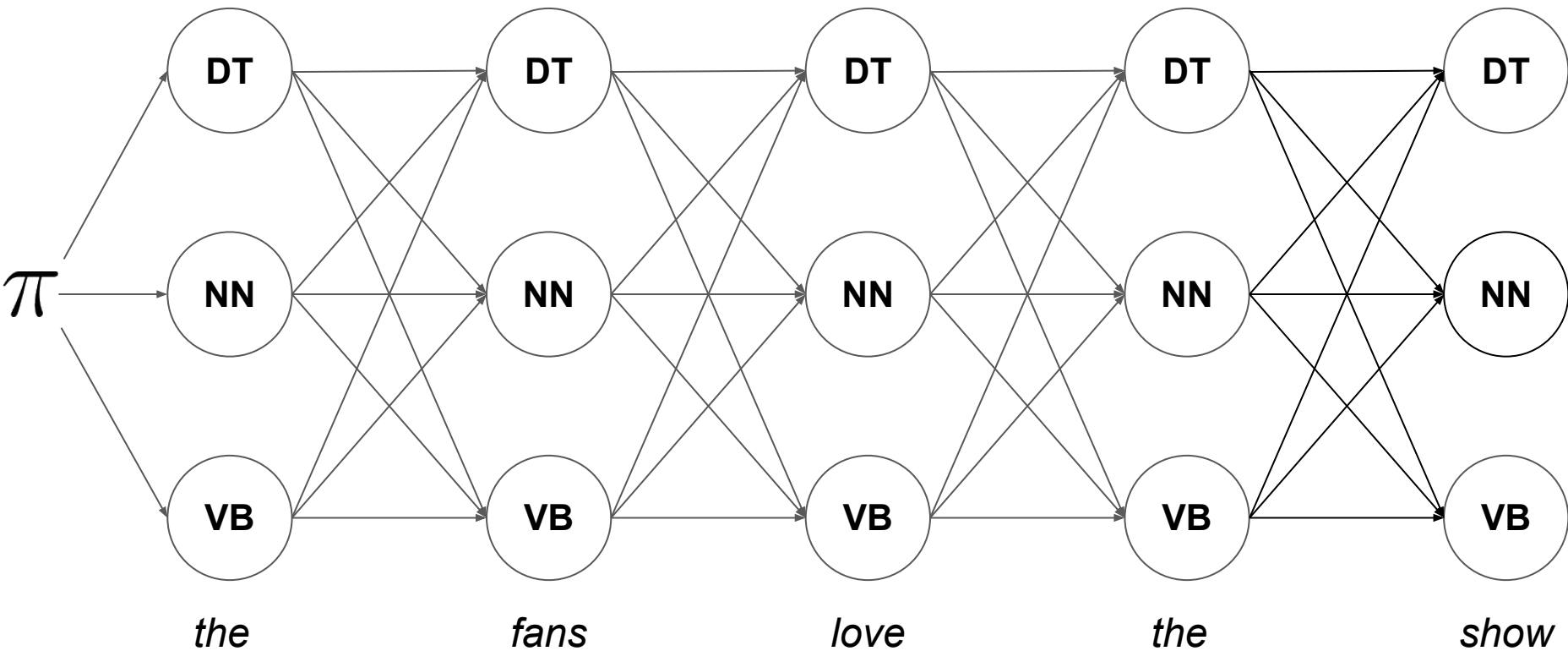
*"the fans love the show"*

## Example

$$\pi = \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix} \begin{matrix} \mathbf{DT} \\ \mathbf{NN} \\ \mathbf{VB} \end{matrix}$$

$$B = \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix} \begin{matrix} \mathbf{DT} \\ \mathbf{NN} \\ \mathbf{VB} \end{matrix}$$

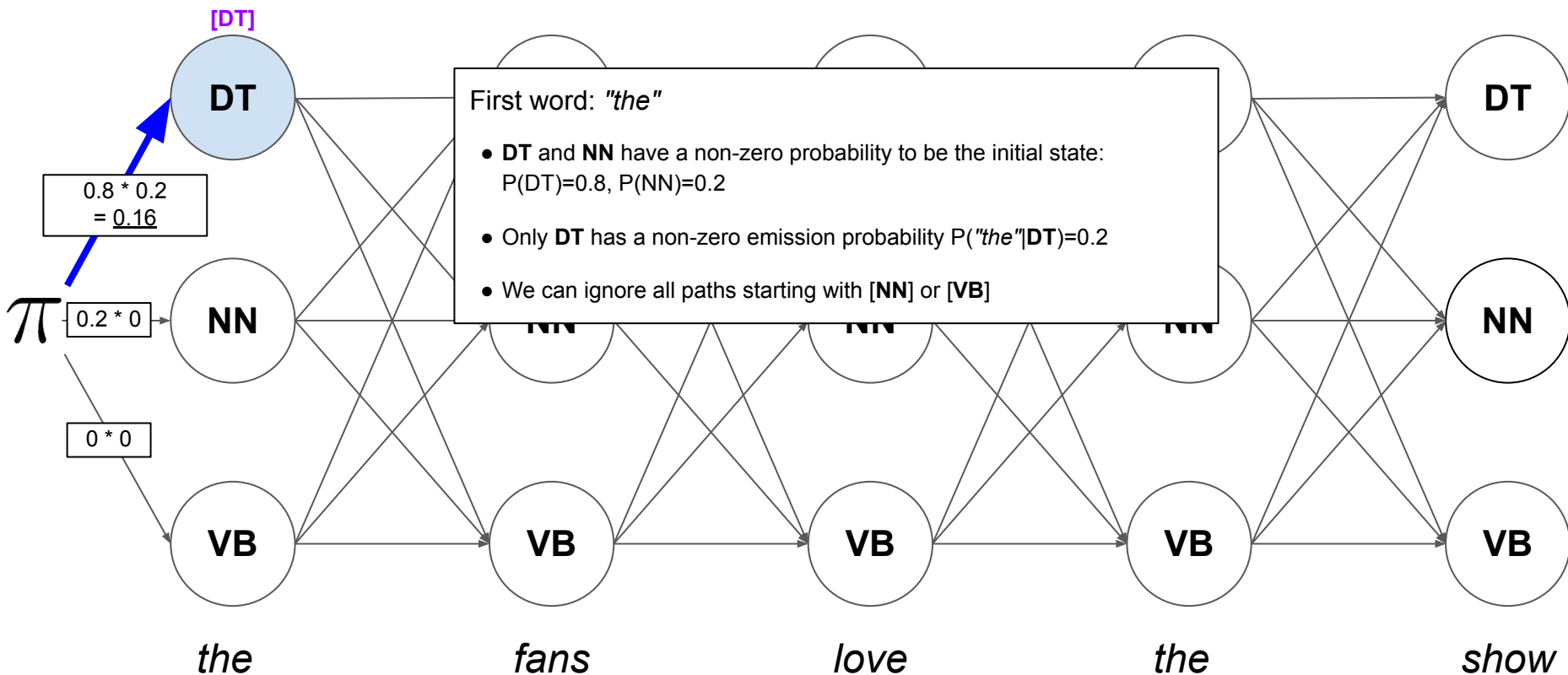


# Example

$$\pi = \begin{matrix} & \text{DT} & \text{NN} & \text{VB} \\ \begin{matrix} \text{DT} \\ \text{NN} \\ \text{VB} \end{matrix} & \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix} \end{matrix}$$

$$A = \begin{matrix} & \text{DT} & \text{NN} & \text{VB} \\ \begin{matrix} \text{DT} \\ \text{NN} \\ \text{VB} \end{matrix} & \begin{bmatrix} 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix} \end{matrix}$$

$$B = \begin{matrix} & \text{the} & \text{fans} & \text{love} & \text{show} \\ \begin{matrix} \text{DT} \\ \text{NN} \\ \text{VB} \end{matrix} & \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix} \end{matrix}$$





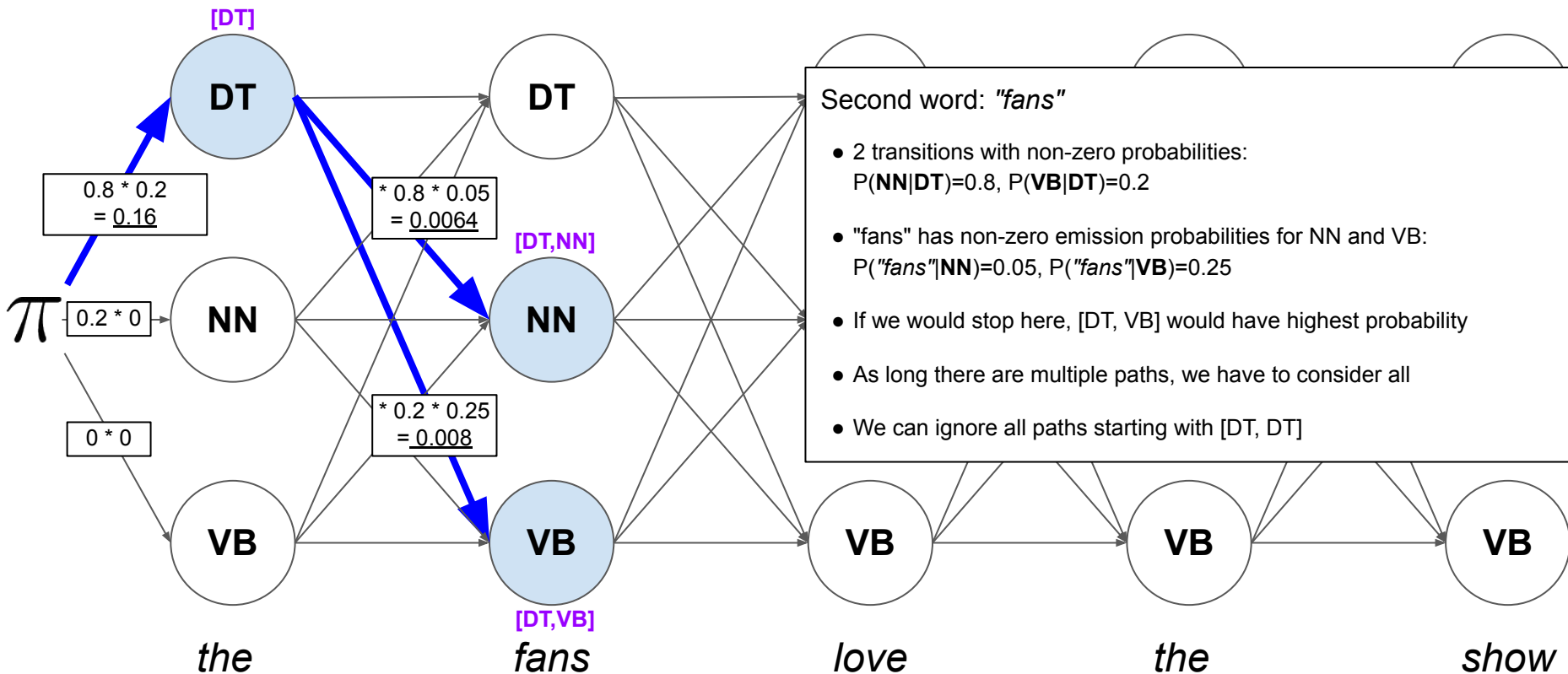
# Example

$$\pi = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.8 & 0.2 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix}$$

*the fans love show*

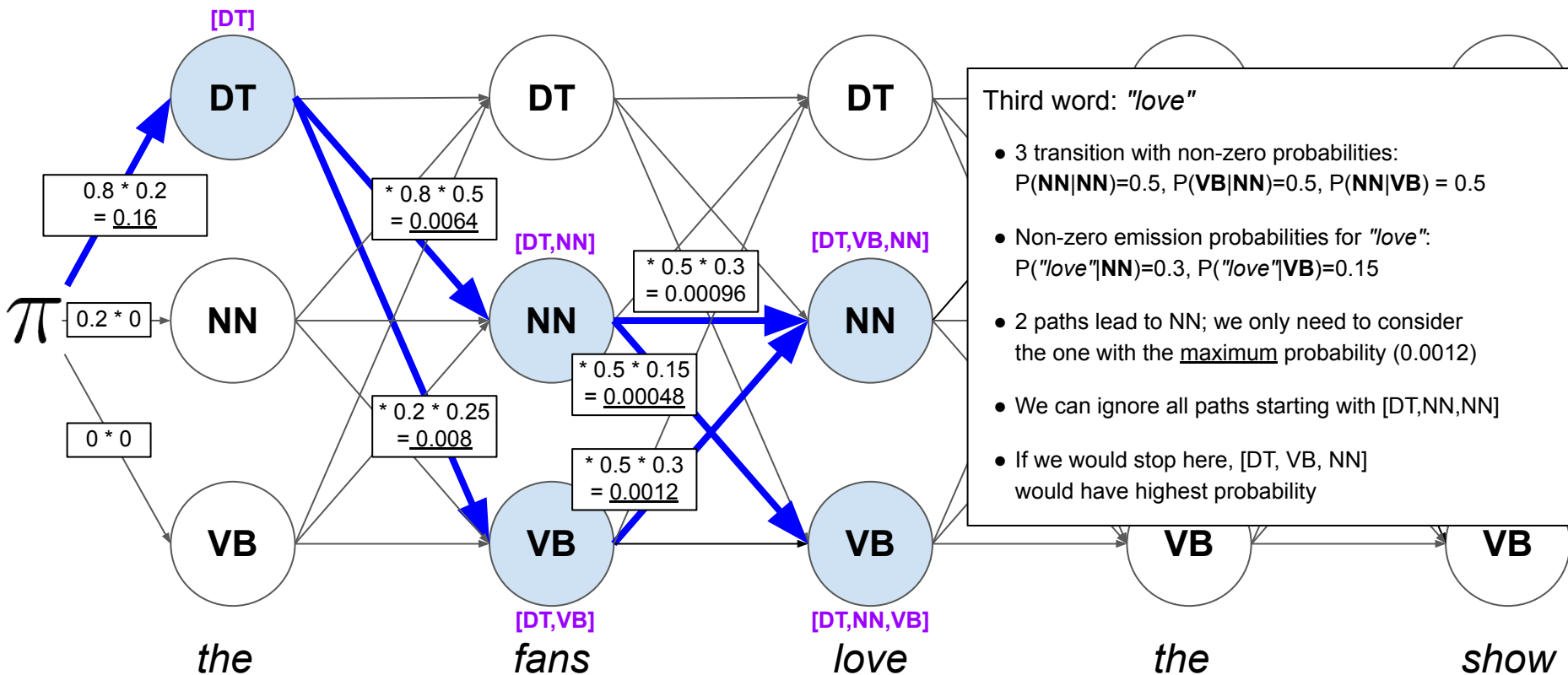


# Example

$$\pi = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.8 & 0.2 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix}$$



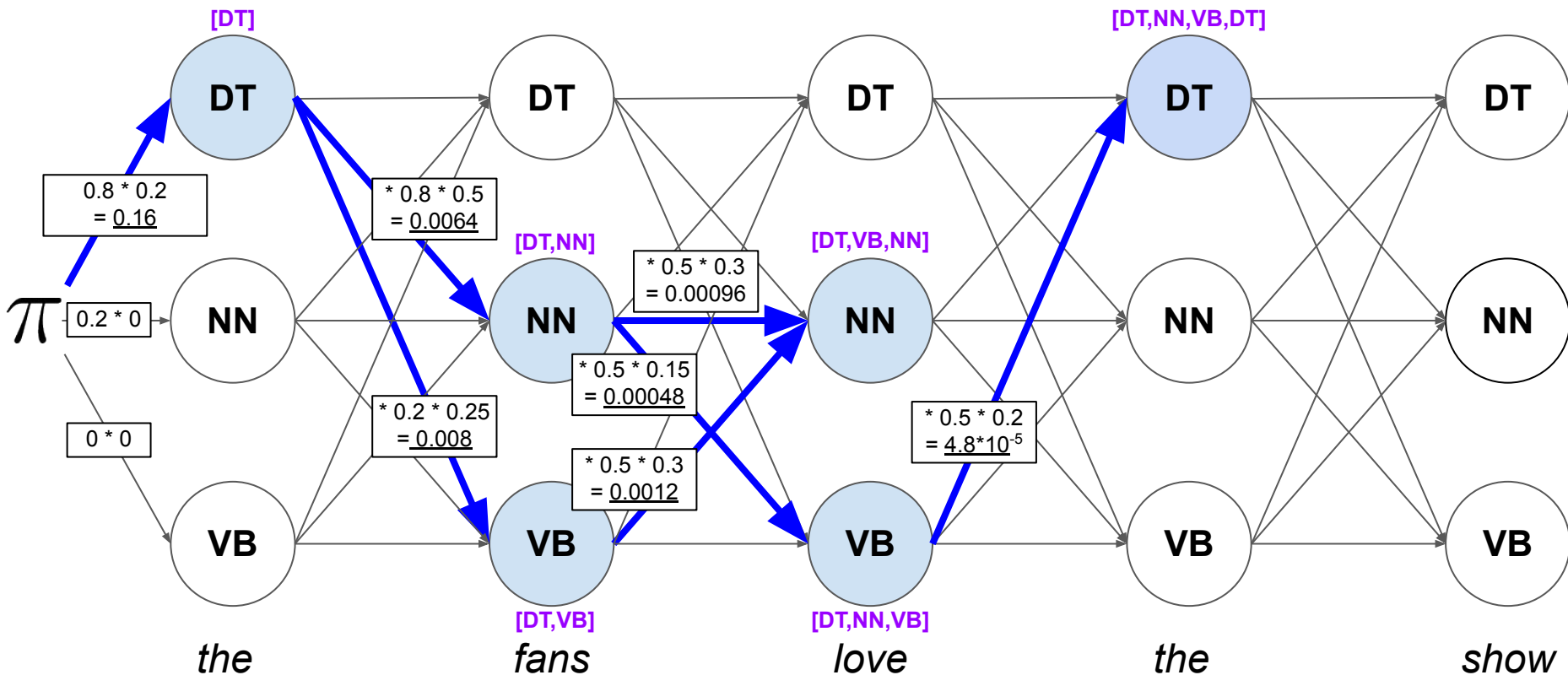
# Example

$$\pi = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.8 & 0.2 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix}$$

*the fans love show*



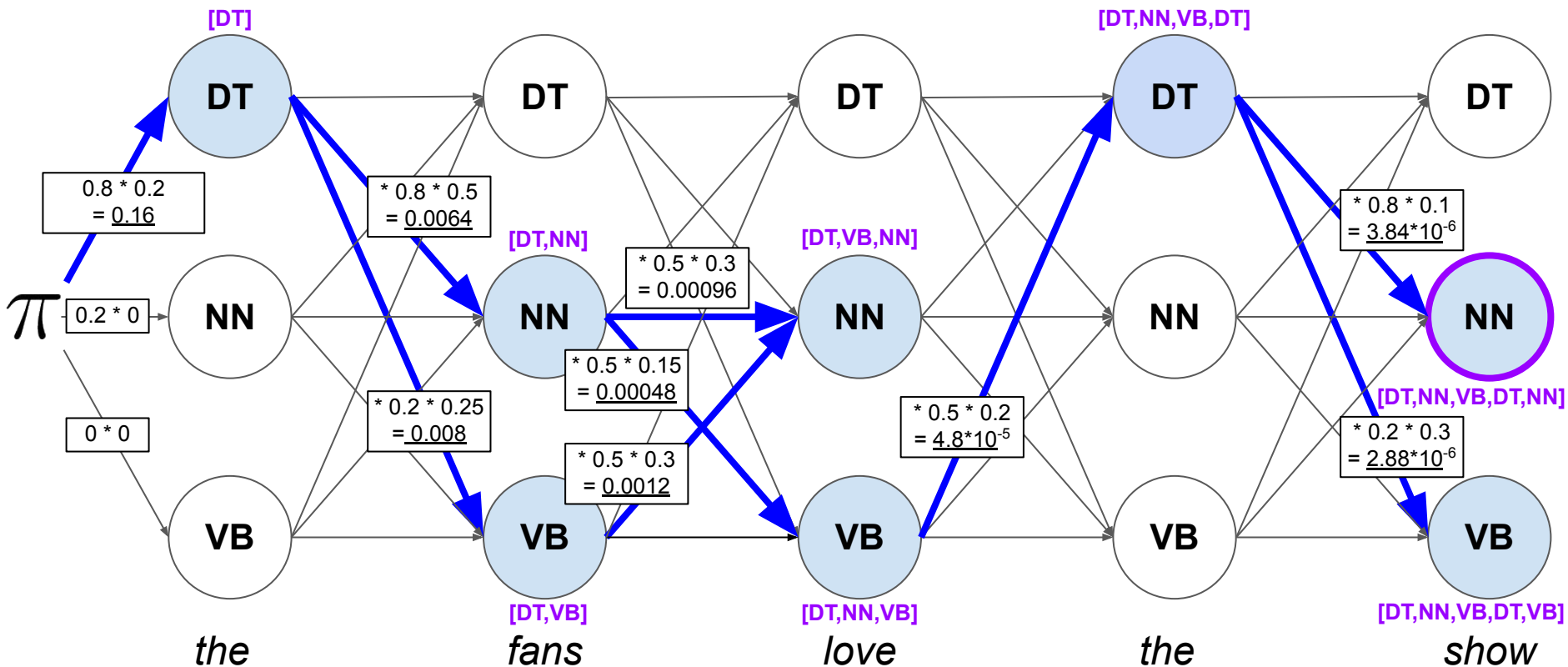
# Example

$$\pi = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.8 & 0.2 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} \text{DT} & \text{NN} & \text{VB} \\ 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix}$$

the fans love show

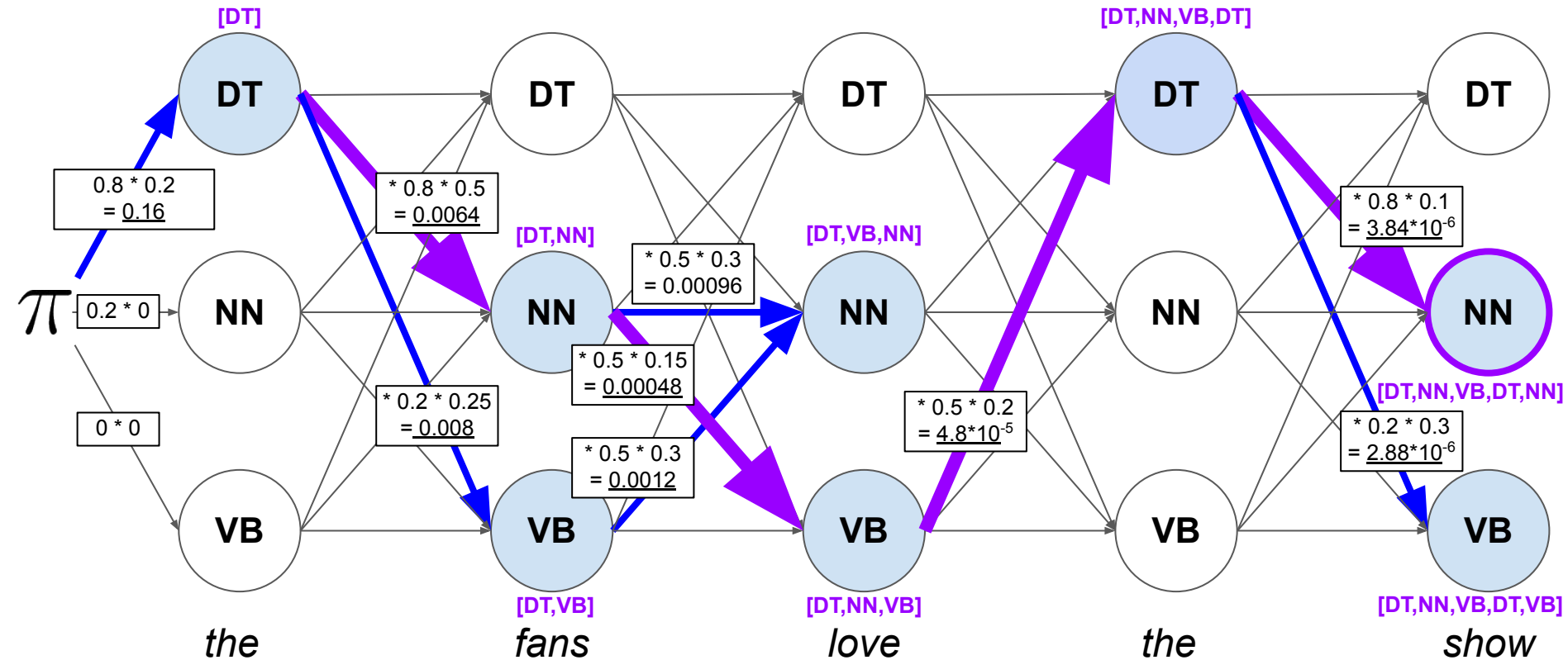


# Viterbi Algorithm

- 2 important question
  - How to get the final state sequence with the highest probability?
  - How exactly does the Viterbi algorithm reduces complexity?

# Backtracking

- During forward pass: remember input path with max probability
- Backtracking: follow paths with max probabilities back to beginning



# Quick Quiz

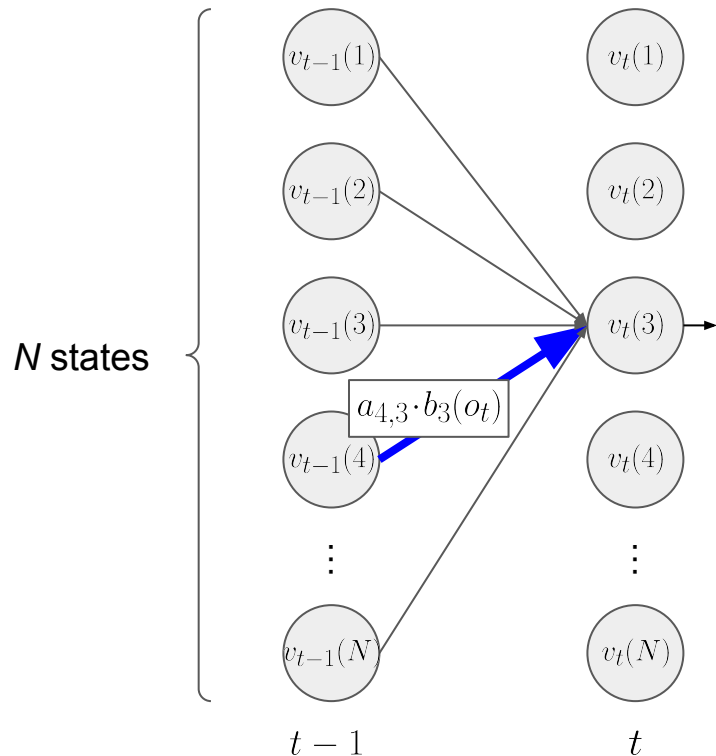


# Viterbi Algorithm — Complexity Analysis





# Viterbi Algorithm — The Basic Algorithm



## Initialization

$$v_1(t) = \pi_j \cdot b_j(o_1)$$

$$1 \leq j \leq N$$

$$bt_1(t) = 0$$

## Recursion

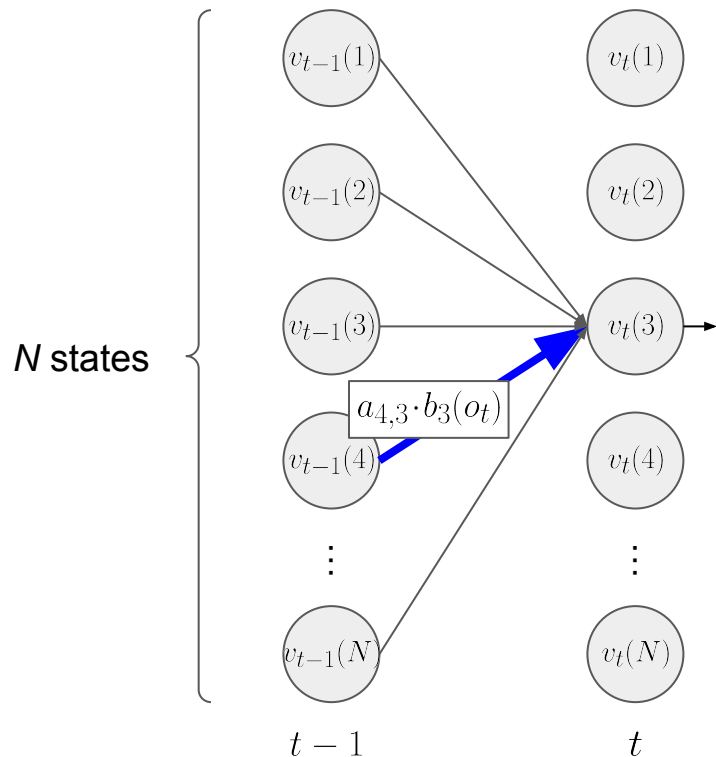
$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

$$1 \leq j \leq N, 1 < t \leq T$$

$$bt_t(j) = \operatorname{argmax}_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

Example for backtrace:  $bt_t(3) = 4$  since we get the highest probability for  $v_t(3)$  from the path coming from  $v_{t-1}(4)$

# Viterbi Algorithm — The Basic Algorithm



**Termination** (after computing all  $v_t(j)$  and  $bt_t(j)$ )

Probability of most likely path:  $P^* = \max_{i=1}^N v_T(i)$

Start of backtrace:  $q_T^* = \operatorname{argmax}_{i=1}^N v_T(i)$

# Viterbi Algorithm — Practical Consideration

- The "usual" problem: Risk of arithmetic underflow

$$\left. \begin{aligned} v_1(t) &= \pi_j \cdot b_j(o_1) \\ v_t(j) &= \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \end{aligned} \right\} \begin{array}{l} \text{Values for } v_t(j) \text{ become very small as we multiple} \\ \text{many (potentially very) small probability values} \end{array}$$

→ The "usual" solution: Logarithm

$$\begin{aligned} v_1(t) &= \log \pi_j + \log b_j(o_1) \\ v_t(j) &= \max_{i=1}^N v_{t-1}(i) + \log a_{ij} + \log b_j(o_t) \end{aligned}$$

# Viterbi Algorithm — Python/NumPy Implementation

```
def viterbi(tokens, A, B, PI):
    N, T = A.shape[0], len(tokens)
    M = np.zeros((N, T))          # Reflecting probabilities of trellis
    BT = np.zeros((N, T), dtype=np.int16) # For the Backtracking pointers

    # Initialization
    for s in range(N):
        M[s,0] = PI[s] * B[s, word2index[tokens[0]]]

    # Recursion (with dynamic programming)
    for t in range(1, T):
        for s in range(N):
            new_probs = M[:,t-1] * A[:,s] * B[s, word2index[tokens[t]]]
            max_idx = np.argmax(new_probs)
            M[s,t] = new_probs[max_idx]
            BT[s,t] = max_idx

    # Termination (start backtracking)
    state = np.argmax(M[:,-1])
    state_sequence = []
    for i in reversed(range(T)):
        state_sequence.append(state)
        state = BT[:,i][state]

    return [ index2tag[idx] for idx in reversed(state_sequence) ]
```

**Note:** This slide is only to show that it does not take much code to implement the Viterbi algorithm.

$$\left. \begin{array}{l} v_1(t) = \pi_j \cdot b_j(o_1) \\ bt_1(t) = 0 \end{array} \right\}$$

$$\left. \begin{array}{l} v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \\ bt_t(j) = \operatorname{argmax}_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \end{array} \right\}$$

$$\left. \begin{array}{l} q_T^* = \operatorname{argmax}_{i=1}^N v_T(i) \end{array} \right\}$$

# Viterbi Algorithm — Python/NumPy Implementation

- Using the HMM trained over 25k movie reviews

- 50 states (POS tags)
- 83k+ tokens (words, punctuation marks, etc.)

**Important:** I've cheated here by annotating the reviews using spaCy, not humans!

```
viterbi(['the', 'fans', 'love', 'the', 'show'], A, B, PI)  
['DT', 'NNS', 'VBP', 'DT', 'NN']
```

```
viterbi(['the', 'fans', 'like', 'the', 'show'], A, B, PI)  
['DT', 'NNS', 'IN', 'DT', 'NN']
```

```
viterbi(['funny', 'movies', 'are', 'the', 'best'], A, B, PI)  
['JJ', 'NNS', 'VBP', 'DT', 'JJS']
```

```
viterbi(['i', 'like', 'watching', 'comedies'], A, B, PI)  
['PRP', 'VBP', 'VBG', 'NNS']
```

# Outline

- Overview: Sequence Tasks
- POS Tagging
  - What are Parts of Speech?
  - Why is this task important and challenging?
- Hidden Markov Models (HMM)
  - Basic setup and components
  - Core HMM tasks
    - Model Learning
    - Likelihood computation
    - Viterbi decoding

# Summary

- Sequences

- A primary form of natural language data with many applications  
(a sentence is sequence of words; sequence captures meaning → BoW model intrinsically limited)
- Many sequence tasks in NLP

- Focus of this lecture: sequence labeling

- POS tagging as very fundamental sequence labeling task
- Different approaches, incl. Hidden Markov Models (HMM)

- Next lecture: encoder-decoder architecture

- Neural network-based architecture
- Applicable to all sequence tasks

# Pre-Lecture Activity for Next Week





# Solutions to Quick Quizzes

