

Natural Language Processing

Lecture 17

Dirk Hovy

dirk.hovy@unibocconi.it

 @dirk_hovy

Goals for Today

- Learn about **structured prediction** problems
- Understand the **structured perceptron**
- Understand the difference between the **greedy** and the **Viterbi** inference version

POS tagging

Grassfed highland Chianina beef with handcut fries and seasonal micro greens 29,—

Rich, tender, golden-brown beef with **crisp** fries and **tender** greens 18,—

Savory beef with **delicious** fries and **tasty** salad 12,—

ADJs = price?

Classification

Positive

My hovercraft is really bling

*CLASSIFICATION: 1 LABEL
FOR THE ENTIRE INPUT*

Sequence Classification

DET NOUN VERB ADP VERB

TOO MANY POSSIBLE LABELS!

DET NOUN VERB ADP NOUN

My hovercraft is really bling

*CLASSIFICATION: 1 LABEL
FOR THE ENTIRE INPUT*

Structured Prediction

WHAT IS THE BEST LABEL
SEQUENCE (Y)
FOR A GIVEN SENTENCE (X)?

DET NOUN VERB ADP ADJ

My hovercraft is really bling

STRUCTURED PREDICTION: 1 LABEL
PER SEQUENCE ELEMENT

Examples

NER

O

O

B-PERS

I-PERS

O

POS

PRON

VERB

PROPN

PROPN

PUNCT

|

|

|

|

|

I

admire

Rosa Parks

.

POS tagging

Open class words	Closed class words	Other
ADJ adjectives: <i>awesome, red</i> ADV adverbs: <i>quietly, where, never</i> INTJ interjections: <i>ouch, shhh</i> NOUN nouns: <i>book, war</i> PROPN proper nouns: <i>Rosa, Twitter</i> VERB full verbs: <i>(she) codes, (they) submitted</i>	ADP adpositions: <i>over, before</i> AUX auxiliary/modal verbs: <i>have (been), could (do), will (change)</i> CCONJ coordinating conjunctions: <i>and, or, but</i> DET determiners: <i>a, they, which</i> NUM numbers. Exactly what you would think it is... PART particles: <i>'s</i> PRON pronouns: <i>you, her, myself</i> SCONJ subordinating conjunctions: <i>since, if, that</i>	PUNCT punctuation marks: <i>!, ?, –</i> SYM symbols: <i>%, \$, :)</i> x other: <i>pfffrt</i>

Ambiguity & Context

can : {VERB, NOUN, AUX}

EXACT
LABEL
DEPENDS

PRON VERB

I can tuna

We can it for later

DET NOUN

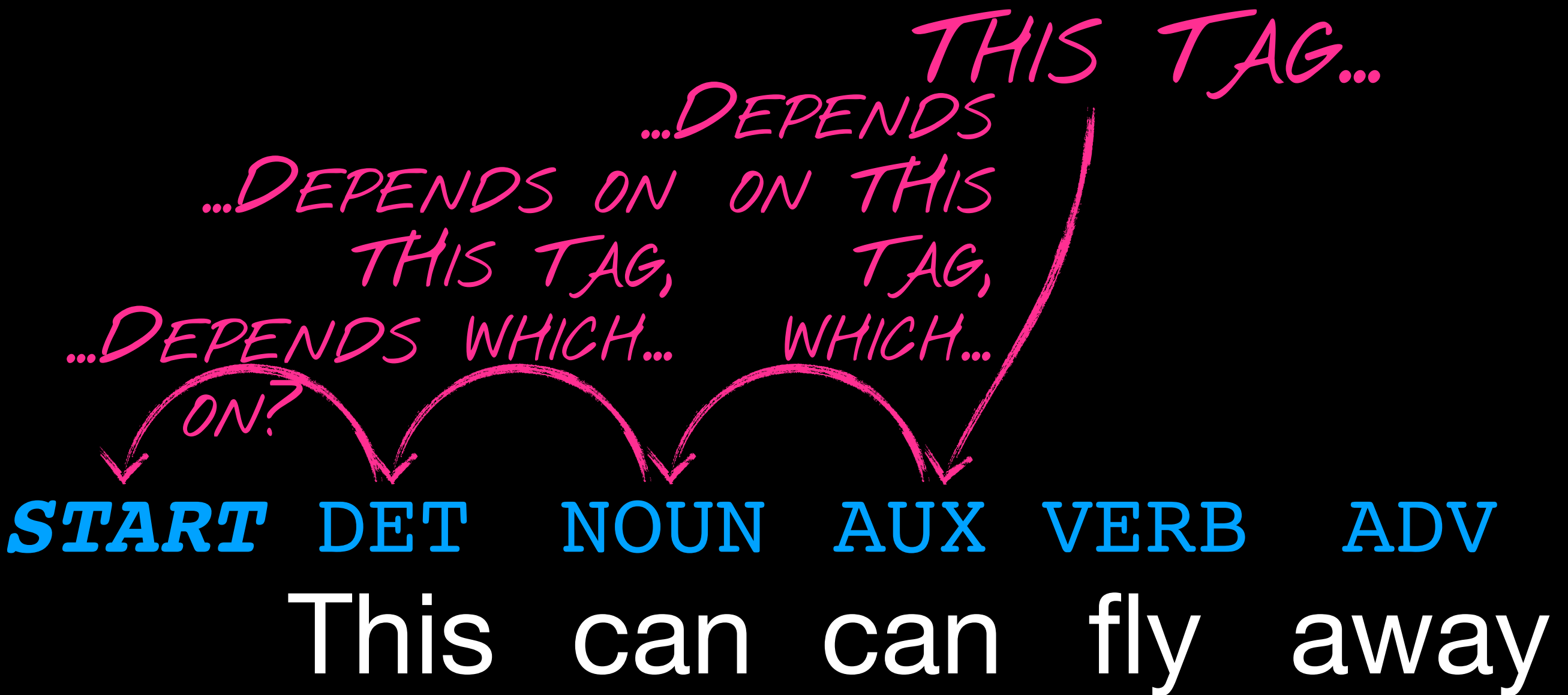
A can of tuna

The can is empty

PRON AUX

I can show you

A Recursive Problem?



SOME WORDS ARE UNAMBIGUOUS (= ONLY ONE POSSIBLE TAG)
THE START OF A SENTENCE IS UNAMBIGUOUS

The Idea: Structured Prediction

A Recursive Solution

guess **tag** for first word, given **previous tag** was *START*

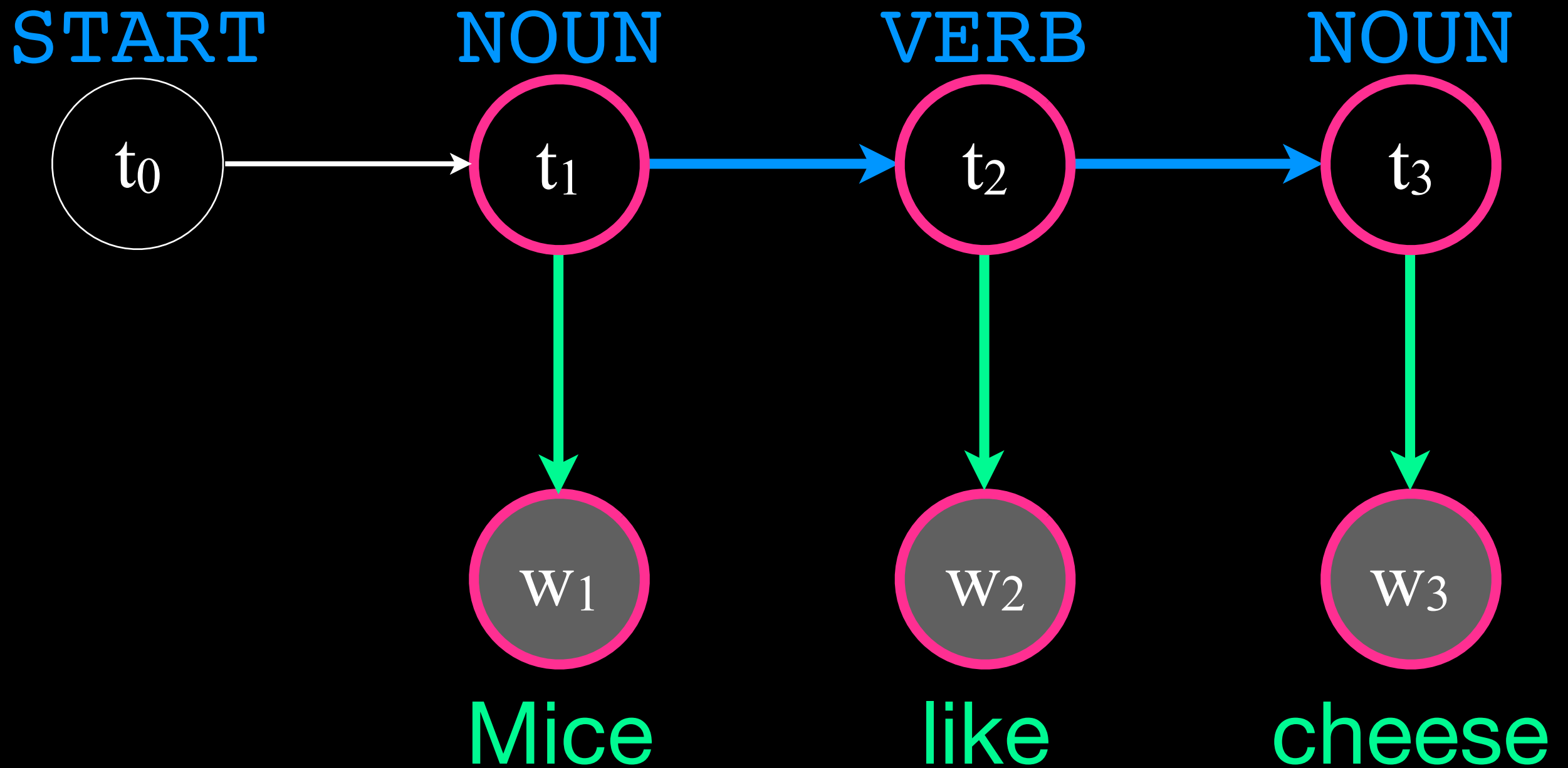
for each **word**:

guess **tag** based on current **word** and **previous tags**

The Structure: Hidden Markov Models

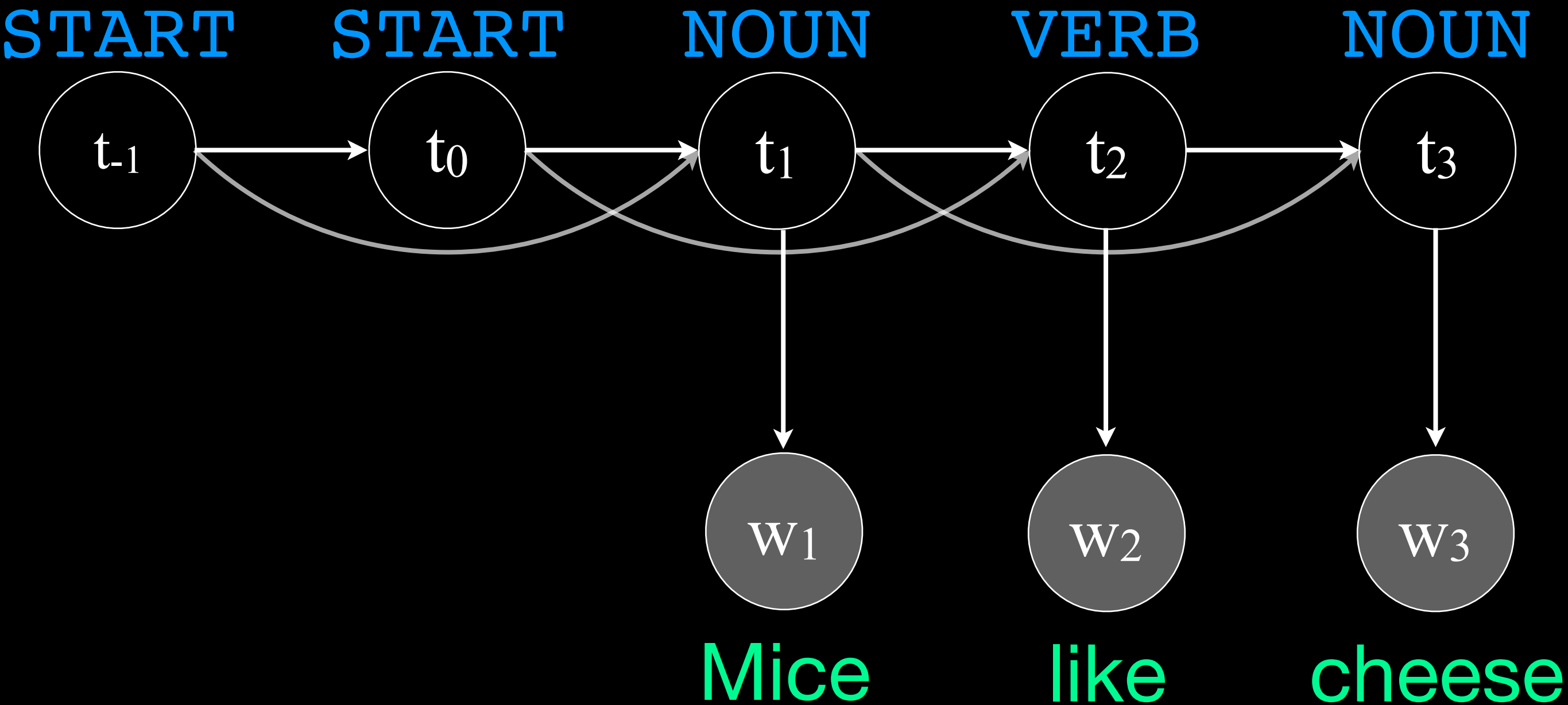
Structure

FIRST ORDER MARKOV CHAIN



Structure

SECOND ORDER MARKOV CHAIN



The Algorithm: Structured Perceptron

Structured Perceptron

for each *iteration*:

for *words*, *true_tags* in examples:

features = get_features(*words*, *true_tags*)

prediction = predict(*features*) *TWO*

if *prediction* != *true_tags*: *OPTIONS*

for *feature* in *features*:

UPDATE! weights[*feature*][*true_tags*] += 1

weights[*feature*][*prediction*] -= 1

The Works

INFERENC

$$T|W = \arg \max_{Y \in \mathcal{Y}} \sum_{d=1}^D \theta_d \cdot \Phi_d(W, Y)$$

DIMENSIONS

FEATURE

PARAMETERS

UPDATES

$$\theta_d = \theta_d + \Phi_d(W, T) - \Phi_d(W, Y)$$

TRUE FEATURES

PREDICTED FEATURES

FEATURES

$$\Phi_d(W, T) = \sum_{i=1}^N \phi_d(w_i, t_{i-2}, t_{i-1}, W, i)$$

TOKENS

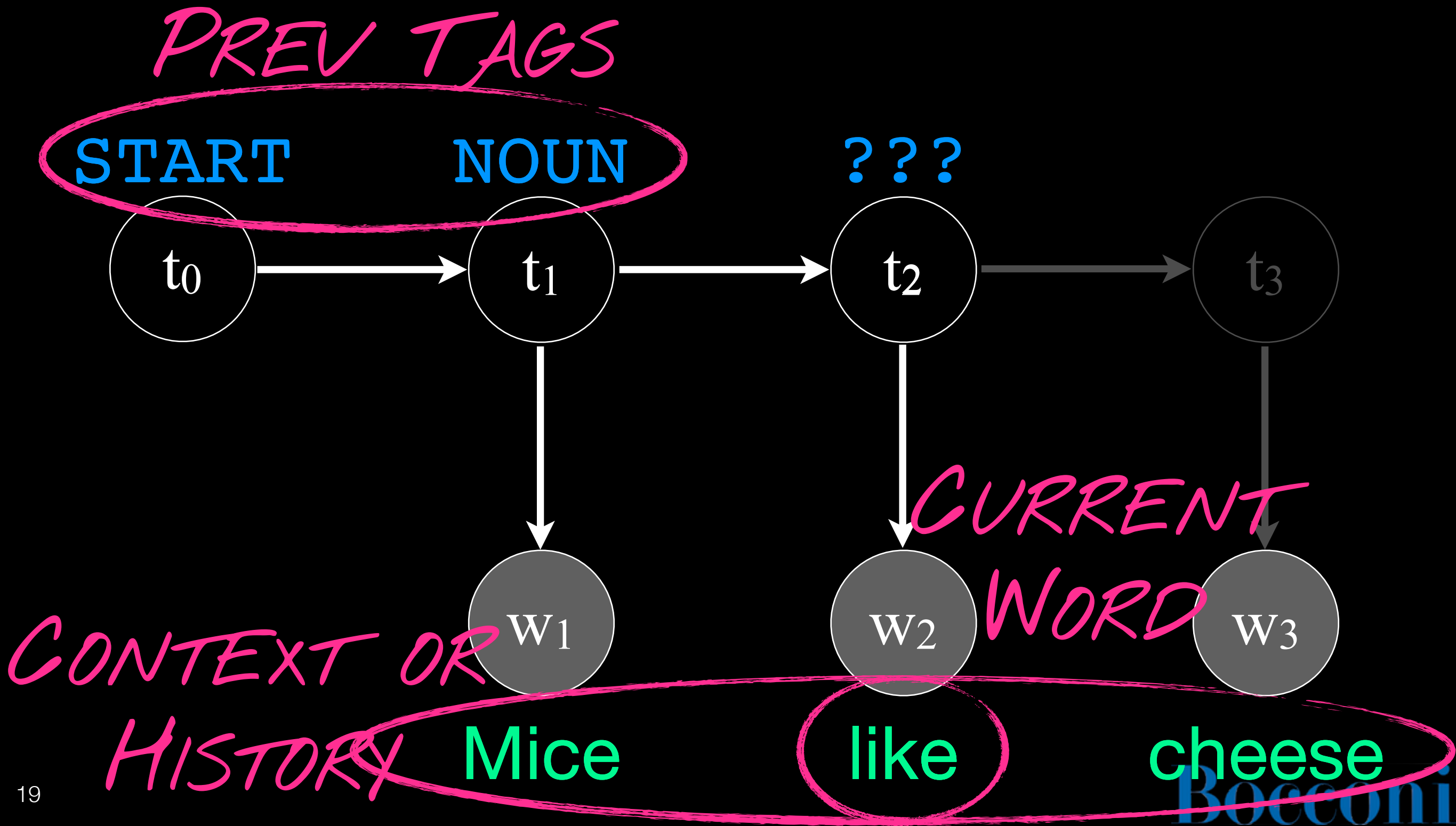
WORD

SENTENCE

TAG HISTORY

Local Features

$$\phi_2 = (\text{like}, \text{START}, \text{NOUN}, \text{Mice like cheese}, 2)$$



Local Feature Function

```
def features(word, prev_tag2, prev_tag, context, i):  
    return {'BIAS',  
            'WORD={}'.format(word),  
            'PREV_WORD={}'.format(context[i-1]),  
            'PREV_WORD+WORD={}+{}'.format(context[i-1], word),  
            'NEXT_WORD={}'.format(context[i+1]),  
            'WORD+NEXT_WORD={}'.format(word, context[i+1]),  
            'PREV_TAG={}'.format(prev_tag),  
            'PREV_TAG2={}'.format(prev_tag2),  
            'PREV_2TAGS={}+{}'.format(prev_tag2, prev_tag),  
            'PREV_TAG+WORD={}+{}'.format(prev_tag, word)  
    }
```

DERIVE FEATURES FROM THE

- CURRENT WORD

- PREVIOUS TAGS

- CONTEXT

Inference

Inference Tricks

*CONSTRAIN POSSIBLE TAGS
BASED ON OBSERVED OPTIONS*

NOUN	NOUN	NOUN	NOUN	NOUN
VERB	VERB	VERB	VERB	VERB
ADV	ADV	ADV	ADV	ADV
DET	DET	DET	DET	DET
This	can	can	fly	away

UNAMBIGUOUS!

Greedy Inference

*PREDICT MOST LIKELY LABEL IN
EACH POSITION*

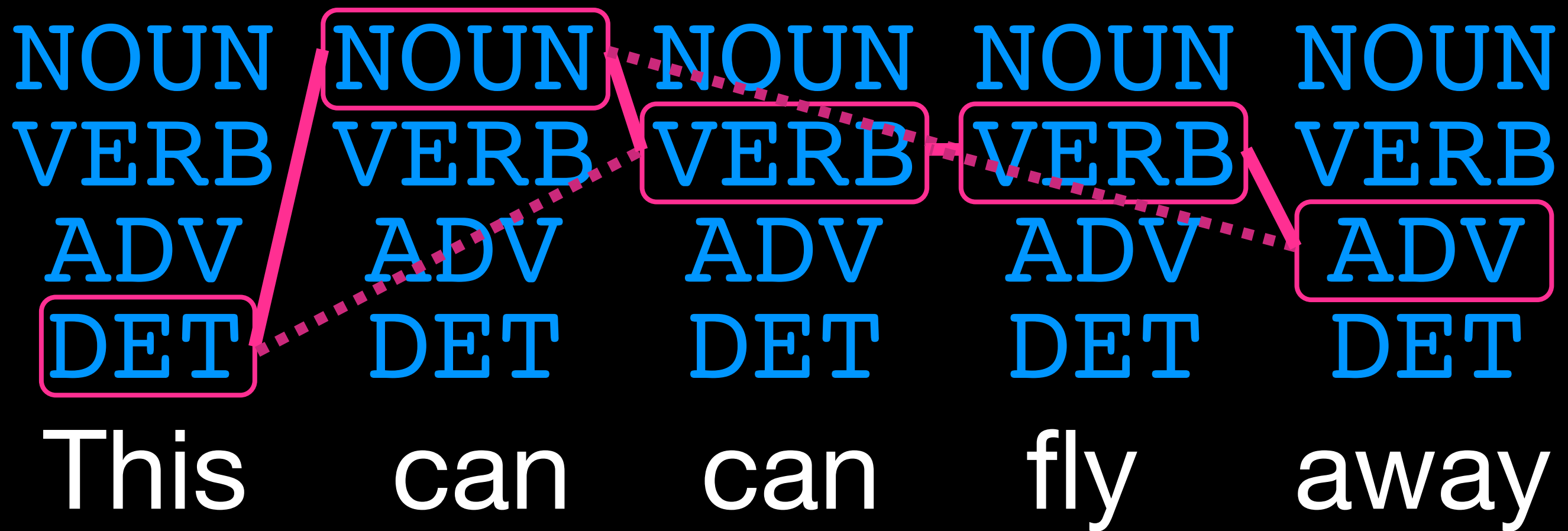
NOUN	NOUN	NOUN	NOUN	NOUN
VERB	VERB	VERB	VERB	VERB
ADV	ADV	ADV	ADV	ADV
DET	DET	DET	DET	DET
This	can	can	fly	away

Greedy Inference

```
for each words, tags:  
    predictions = []  
    for each word in words:  
        scores =  $\theta \phi(\text{word}, \text{tags}[-2], \text{tags}[-1], \text{words})$   
        tag = argmax(scores)  
        predictions(tag)
```


Viterbi Inference

*PREDICT EACH LABEL BASED ON
BEST PATH TO IT*



Viterbi Inference

```
# initialize scores
features = get_features(words[0], START, START, context, 1)
scores = get_scores(features)

for each allowed tag on word 1:
    Q[tag, word] = scores[tag]

# filling the lattice, for every position and every tag find Viterbi score Q
for each word i:
    for each allowed prev_tag on prev_word:
        best_score = float('-Inf')
        prev_score = Q[prev_tag, prev_word] # score of previous tag on previous word

        for each allowed prev2_tag:
            if i == 1:
                prev2_tag = START

            # get features of word i with the two previous tags
            features = get_features(word, prev2_tag, prev_tag, context, i)
            scores = get_scores(features)

            # update best score
            for each allowed tag on current word:
                tag_score = prev_score + scores[tag]

                if tag_score > best_score:
                    Q[tag, word] = tag_score
                    best_score = tag_score
                    backpointers[tag, word] = prev_tag
```

Viterbi Algorithm

This can can fly away

SCORE

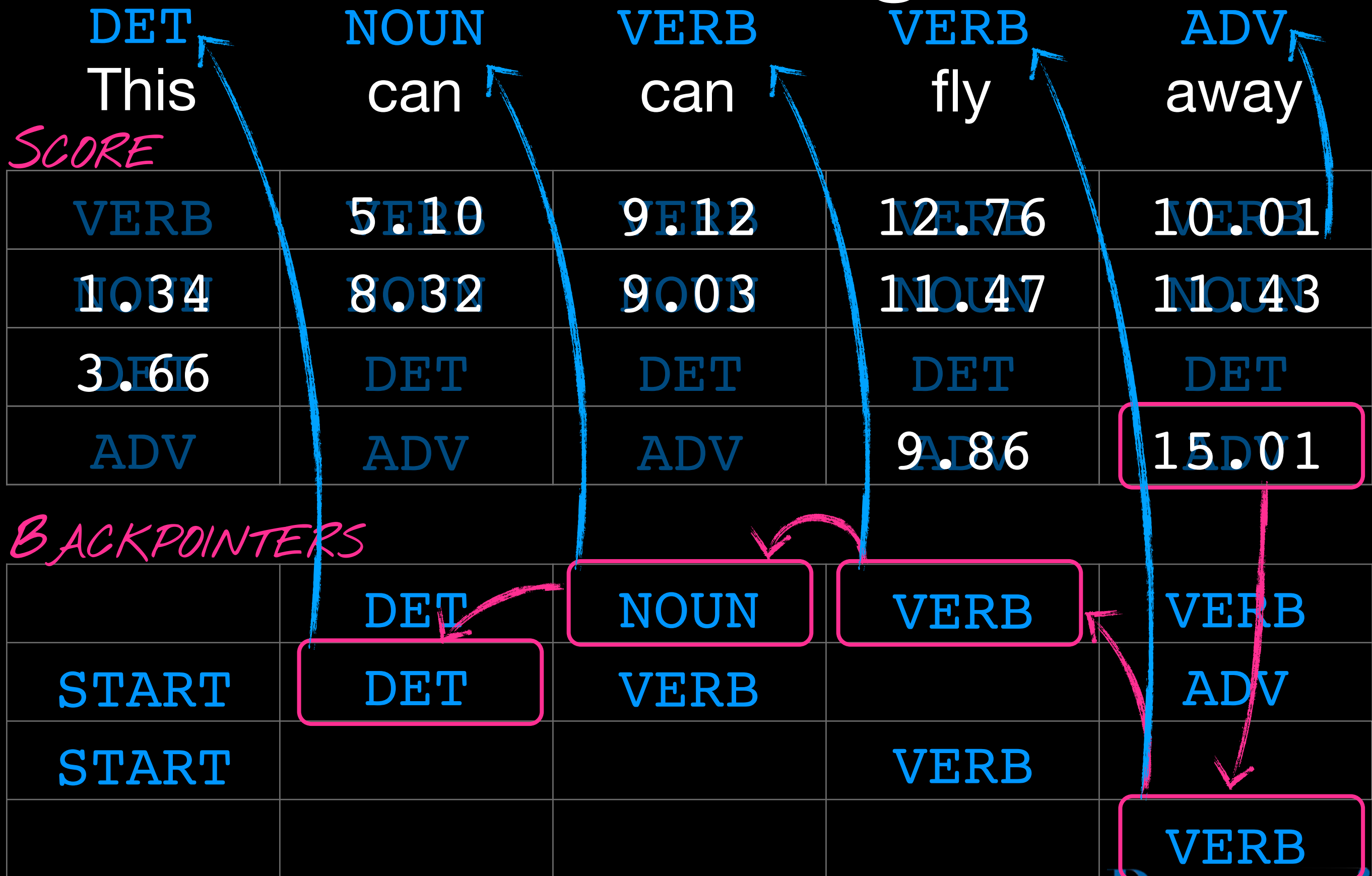
VERB 5.10	VERB 9.12	VERB 12.76	VERB 10.01
1.34	8.32	9.03	11.47
3.66	DET	DET	DET
ADV	ADV	ADV	9.86
			15.01

$$Q(TAG, 0) + \phi(\text{can}, \text{START}, TAG, \text{words}, 0) \cdot \theta$$

$$\phi(\text{This}, \text{START}, \text{START}, \text{words}, 0) \cdot \theta$$

27 *INITIALIZE FIRST ROWS*

Decoding

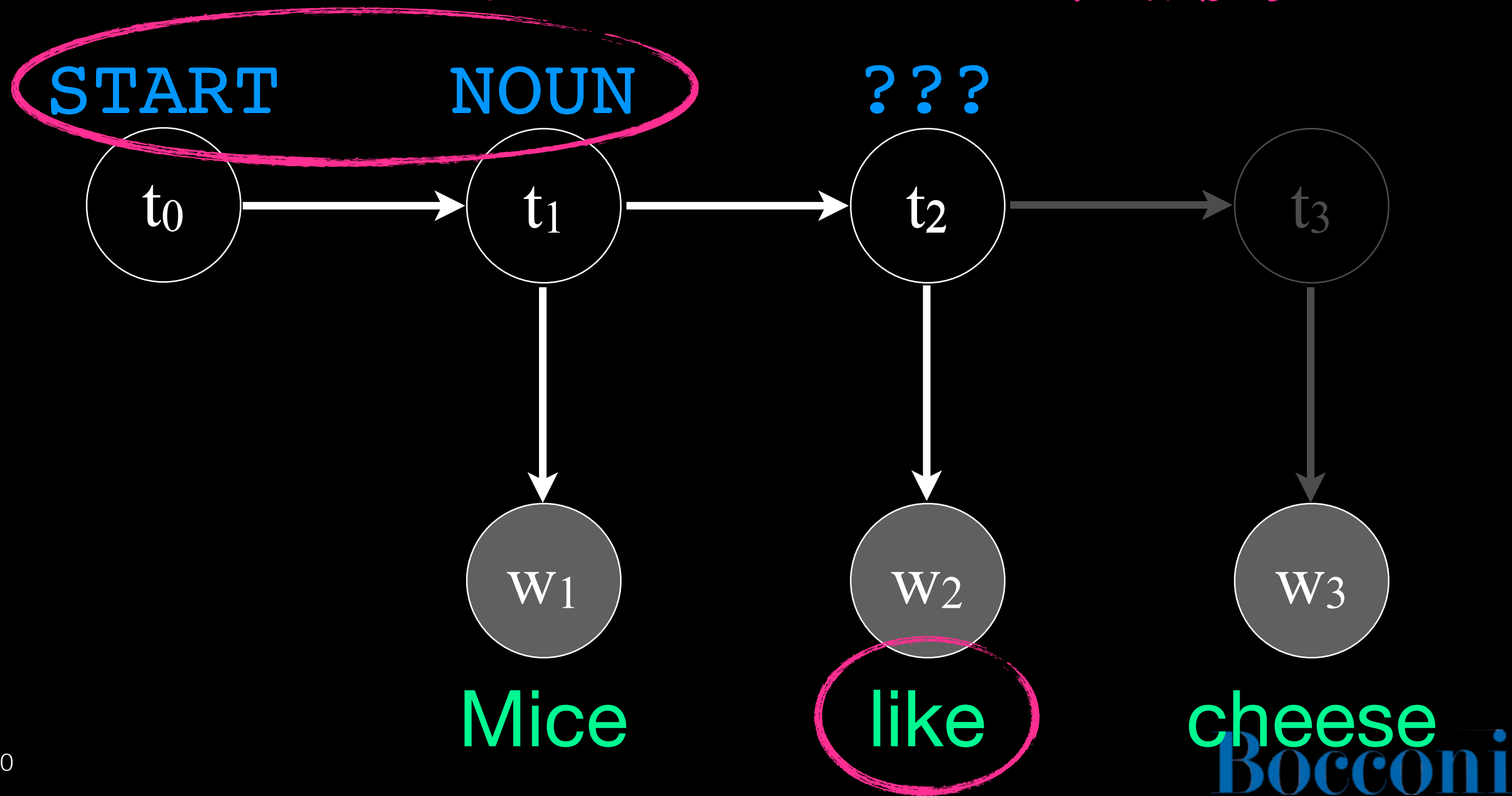


Hyper-Parameters

- Independent of inference, you need to tune
 - learning rate
 - number of iterations
 - features you generate

Using Embeddings

$$\phi_2 = (\text{PRETRAINED} \text{ FIXED})$$



Wrapping Up

Inference Comparison

	greedy	Viterbi
speed	fast	slower
accuracy	good	better on small data
sparsity	lower	higher
type	approximate	exact

Take-home points

- **Structured prediction** finds a tag for each token, based on context
- Assume invisible **START** tags at beginning
- Derive good **features** from context allows for unknown words
- **Inference** can be approximate (**greedy**) or exact (**Viterbi**)