

Sequence labeling

CMSC 723 / LING 723 / INST 725

Hal Daumé III [he/him]

15 Oct 2019

(many slides c/o Marine Carpuat or Graham Neubig)

Announcements, logistics

- Projects:
 - Google cloud \$\$\$ available if you need it (\$50/student)
 - Project pitches *one week* from today!
 - One slide (pdf, no animations), two minutes!
 - One page write-up
 - Due by NOON (so I have time to assemble them)
 - You will present in a (pseudo-)random order
- HW3 due Thursday (before class)

Last time

- Language modeling as a prediction problem
- Feed-forward neural language models (& embeddings)
- Recurrent neural language models

Today

- Sequence labeling as independent predictions
- Structured perceptron for sequence labeling
- Do we really need structured features?
- Recurrent neural network taggers

Sequence labeling is a workhorse of NLP

- NLP for NLP's sake:

- Part of speech tagging
- Word segmentation
- Morphological segmentation
- All-word word sense disambiguation
- ...even parsing!

- NLP for other's sake:

- Named entity recognition
- Grammatical error detection
- Spell correction
- ...

- General constraint: 1-1 mapping between input to

- Plus "BIO" trick for encoding subsequence problems

ກຸ່ມຊາດພັນໄທ-ກ
ບໍລິສຸດແລະເປັນຕົ້ນ
ສະດີດັງກ່າວຖືກຍື່ນ

Allergiantikropparna känner igen det ämne man är allergisk mot, till exempel pollen. När man andas in pollen sätts en allergisk reaktion igång och olika ämnen, bland annat histamin, frigörs. När histamin och andra ämnen frisätts vid den allergiska reaktionen startar en inflammation i ögonen och näsans slemhinnor. Det går inte att stoppa kroppens allergiska reaktioner helt och hållet, men mediciner kan dämpa besvären. Genom att använda läkemedel ska man kunna leva som vanligt och vistas utomhus trots att det finns pollen i luften. Man kan prova nässprej, ögondroppar eller tabletter mot allergi. Om man blir bättre av medicinerna är det troligt att pollenallergi är orsaken. Besvären kan också bero på en vanlig förkylning, och då hjälper inte medicinerna. Om man är osäker kan det vara bra att fråga om råd på ett apotek eller besöka en läkare. Ibland kan det hjälpa att bara skölja och rensa näsan från pollen.

uyuyorum	I am sleeping
uyuyorsun	you are sleeping
uyuyor	he/she/it is sleeping
uyuyoruz	we are sleeping
uyuyorsunuz	you are sleeping
uyuyorlar	they are sleeping
uyuduk	we slept
uyudukça	as long as (somebody) sleeps

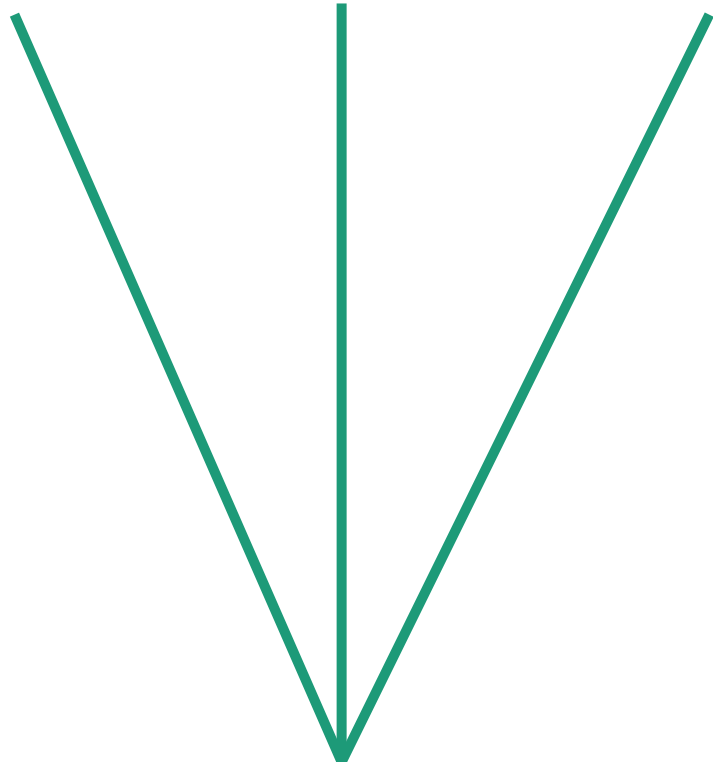
we must sleep
without sleeping
our sleeping
while (somebody) is sleeping
when (somebody) sleeps
to cause somebody to sleep
to cause (somebody) to
cause (another) to sleep
to cause (somebody) to
cause (some other) to
cause (yet another) to
sleep

uyutturtturmak

Sequence labeling as independent predictions

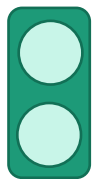
$x =$ 日本語 が あまり 話せません

$y =$ noun part adv verb



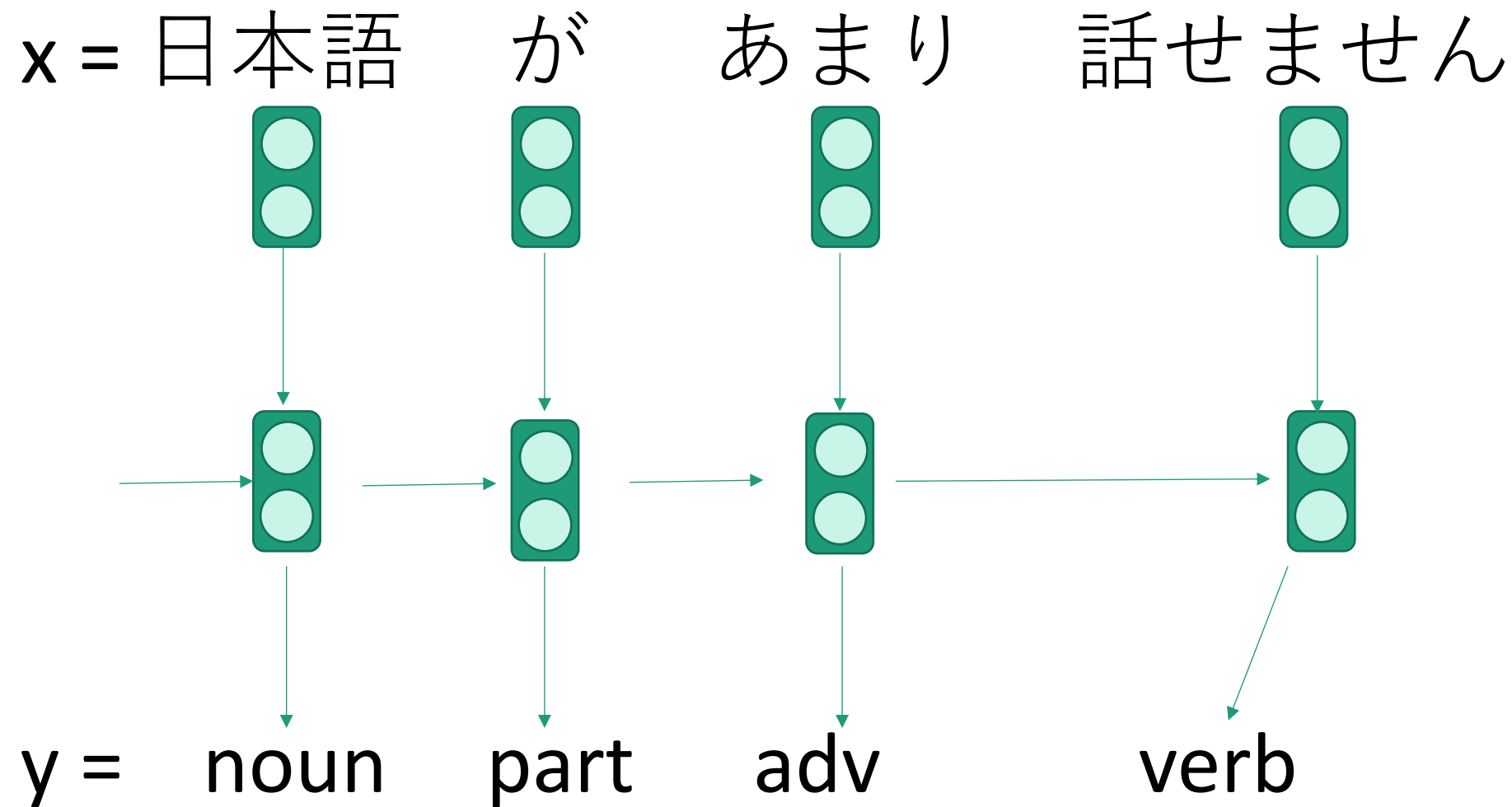
Sequence labeling as independent predictions

$x =$ 日本語 が あまり 話せません



$y =$ noun part adv verb

Sequence labeling as independent predictions



POS tagging

Sequence labeling with the perceptron

Sequence labeling problem

- Input:
 - sequence of tokens $x = [x_1 \dots x_L]$
 - Variable length L
- Output (aka label):
 - sequence of tags $y = [y_1 \dots y_L]$
 - # tags = K
 - Size of output space?

Structured Perceptron

- Perceptron algorithm can be used for sequence labeling
- But there are challenges
 - How to compute argmax efficiently?
 - What are appropriate features?
- Approach: leverage structure of output space

Reminder: multiclass perceptron

Algorithm 1 Perceptron learning algorithm

```
1: procedure PERCEPTRON( $\mathbf{x}_{1:N}, y_{1:N}$ )
2:   repeat
3:     Select an instance  $i$ 
4:      $\hat{y} \leftarrow \arg \max_y \boldsymbol{\theta}_t^\top \mathbf{f}(\mathbf{x}_i, y)$ 
5:     if  $\hat{y} \neq y_i$  then
6:        $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$ 
7:     else
8:       do nothing
9:   until tired
```

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$$

Feature
function
representation

Weights

Solving the argmax problem for sequences with dynamic programming

$x =$ “ monsters eat tasty bunnies ”

$y =$ noun verb adj noun

- Efficient algorithms possible if **the feature function decomposes over the input**
- This holds for unary and markov features used for POS tagging

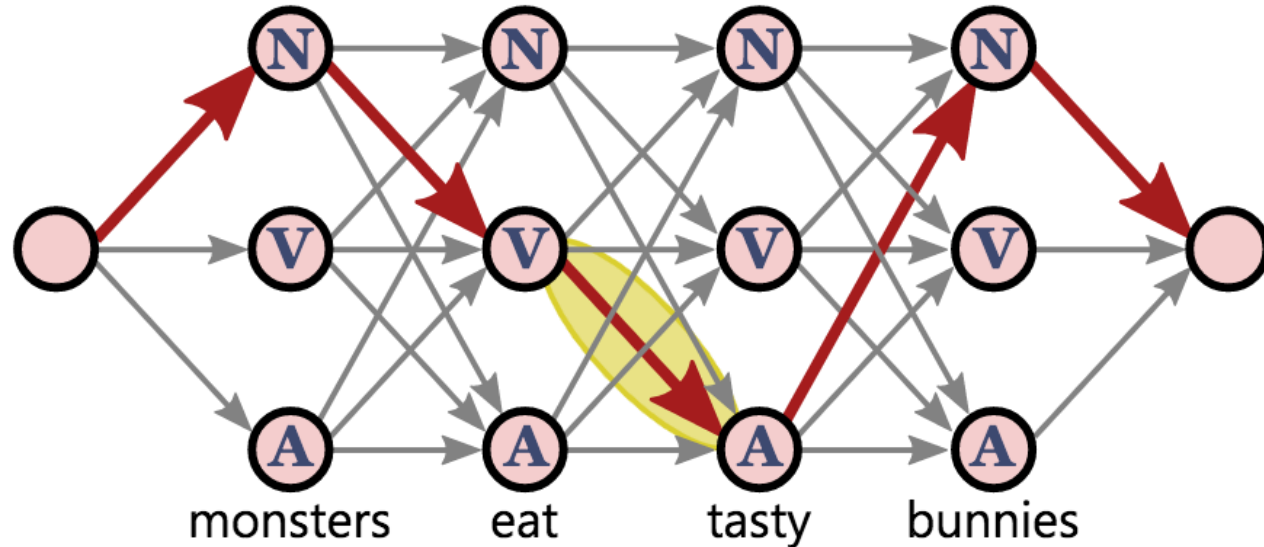
Feature functions for sequence labeling

$x =$ “ monsters eat tasty bunnies ”

$y =$ noun verb adj noun

- Standard features of POS tagging
 - **Unary features:** # times word w has been labeled with tag l for all words w and all tags l
 - **Markov features:** # times tag l is adjacent to tag l' in output for all tags l and l'
- Size of feature representation is constant wrt input length

Solving the argmax problem for sequences



- Trellis sequence labeling
 - Any path represents a labeling of input sentence
 - Gold standard path in red
 - Each edge receives a weight such that adding weights along the path corresponds to score for input/output configuration
- Any max-weight path algorithm can find the argmax
 - e.g. Viterbi algorithm $O(LK^2)$

Defining weights of edge in trellis

Unary features at position l together with Markov features that end at position l

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y) \quad \text{decomposition of structure} \quad (17.35)$$

$$= \sum_{l=1}^L w \cdot \phi_l(x, y) \quad \text{associative law} \quad (17.36)$$

- Weight of edge that goes from time $l-1$ to time l , and transitions from y to y'

$$w \cdot \phi_l(x, \dots \circ y \circ y')$$

Dynamic program

- Define: the score of best possible output prefix up to and including position l that labels the l -th word with label k

$$\alpha_{l,k} = \max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(x, \hat{y} \circ k)$$

- With decomposable features, alphas can be computed recursively

$$\alpha_{l+1,k} = \max_{k'} \left[\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right]$$

$$\alpha_{0,k} = 0 \quad \forall k \quad (17.41)$$

$$\zeta_{0,k} = \emptyset \quad \forall k \quad (17.42)$$

the score for any empty sequence is zero

$$\alpha_{l+1,k} = \max_{\hat{y}_{1:l}} w \cdot \phi_{1:l+1}(x, \hat{y} \circ k) \quad (17.43)$$

separate score of prefix from score of position $l+1$

$$= \max_{\hat{y}_{1:l}} w \cdot \left(\phi_{1:l}(x, \hat{y}) + \phi_{l+1}(x, \hat{y} \circ k) \right) \quad (17.44)$$

distributive law over dot products

$$= \max_{\hat{y}_{1:l}} \left[w \cdot \phi_{1:l}(x, \hat{y}) + w \cdot \phi_{l+1}(x, \hat{y} \circ k) \right] \quad (17.45)$$

separate out final label from prefix, call it k'

$$= \max_{\hat{y}_{1:l-1}} \max_{k'} \left[w \cdot \phi_{1:l}(x, \hat{y} \circ k') + w \cdot \phi_{l+1}(x, \hat{y} \circ k' \circ k) \right] \quad (17.46)$$

swap order of maxes, and last term doesn't depend on prefix

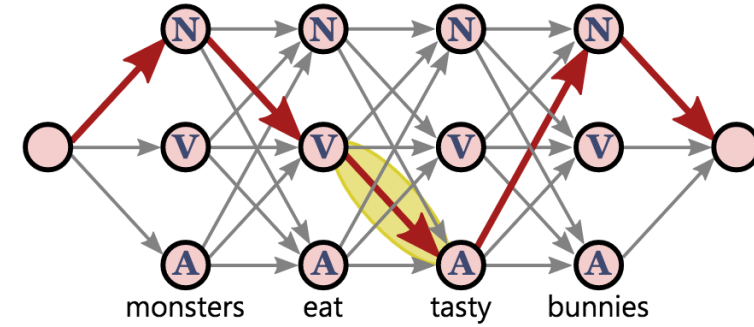
$$= \max_{k'} \left[\left[\max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(x, \hat{y} \circ k') \right] + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right] \quad (17.47)$$

apply recursive definition

$$= \max_{k'} \left[\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right] \quad (17.48)$$

Algorithm 42 ARGMAXFORSEQUENCES(x, w)

```
1:  $L \leftarrow \text{LEN}(x)$ 
2:  $\alpha_{l,k} \leftarrow 0, \quad \zeta_{k,l} \leftarrow 0, \quad \forall k = 1 \dots K, \quad \forall l = 0 \dots L$  // initialize variables
3: for  $l = 0 \dots L-1$  do
4:   for  $k = 1 \dots K$  do
5:      $\alpha_{l+1,k} \leftarrow \max_{k'} [\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle)]$  // recursion:
        // here,  $\phi_{l+1}(\dots k', k \dots)$  is the set of features associated with
        // output position  $l + 1$  and two adjacent labels  $k'$  and  $k$  at that position
6:      $\zeta_{l+1,k} \leftarrow$  the  $k'$  that achieves the maximum above // store backpointer
7:   end for
8: end for
9:  $y \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize predicted output to L-many zeros
10:  $y_L \leftarrow \text{argmax}_k \alpha_{L,k}$  // extract highest scoring final label
11: for  $l = L-1 \dots 1$  do
12:    $y_l \leftarrow \zeta_{l,y_{l+1}}$  // traceback  $\zeta$  based on  $y_{l+1}$ 
13: end for
14: return  $y$  // return predicted output
```



A more general approach for argmax Integer Linear Programming

- ILP: optimization problem of the form,
for a fixed vector a

$$\max_z \quad a \cdot z \quad \text{subj. to} \quad \text{linear constraints on } z$$

- With integer constraints
- Pro: can leverage well-engineered solvers (e.g., Gurobi)
- Con: not always most efficient

POS tagging as ILP

- Markov features as binary indicator variables

$$z_{l,k',k} = \mathbf{1}[\text{label } l \text{ is } k \text{ and label } l-1 \text{ is } k']$$

- Output sequence: $y(z)$ obtained by reading off variables z
- Define a such that $a \cdot z$ is equal to score

$$a_{l,k',k} = w \cdot \phi_l(x, \langle \dots, k', k \rangle)$$

- Enforcing constraints for well formed solutions

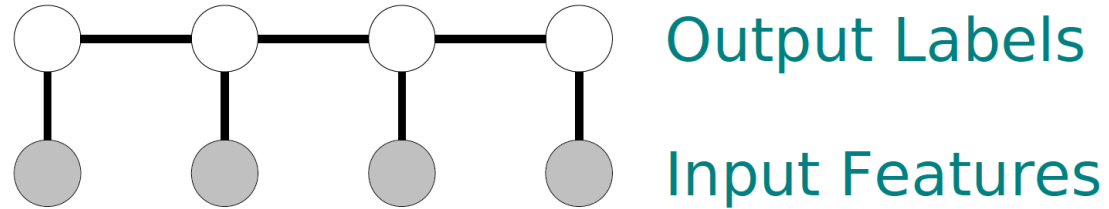
1. That all the z s are binary. That's easy: just say $z_{l,k',k} \in \{0, 1\}$, for all l, k', k .
2. That for a given position l , there is exactly one active z . We can do this with an equality constraint: $\sum_k \sum_{k'} z_{l,k',k} = 1$ for all l .
3. That the z s are internally consistent: if the label at position 5 is supposed to be "noun" then both $z_{5,..}$ and $z_{6,..}$ need to agree on this. We can do this as: $\sum_{k'} z_{l,k',k} = \sum_{k''} z_{l+1,k,k''}$ for all l, k . Effectively what this is saying is that $z_{5,?,\text{verb}} = z_{6,\text{verb},?}$ where the "?" means "sum over all possibilities."

Sequence labeling

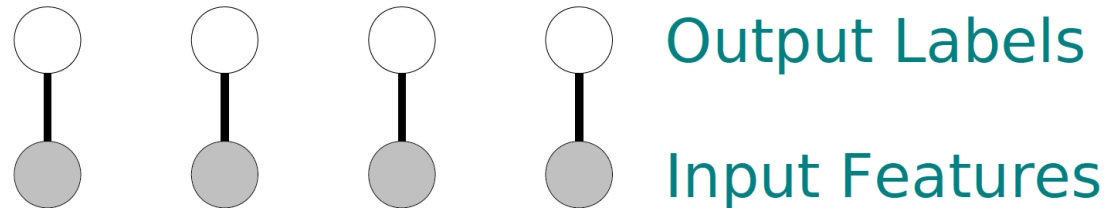
- Structured perceptron
 - A general algorithm for structured prediction problems such as sequence labeling
- The Argmax problem
 - Efficient argmax for sequences with Viterbi algorithm, given some assumptions on feature structure
 - A more general solution: Integer Linear Programming
- Loss-augmented argmax
 - Hamming Loss

How much does this structure matter?

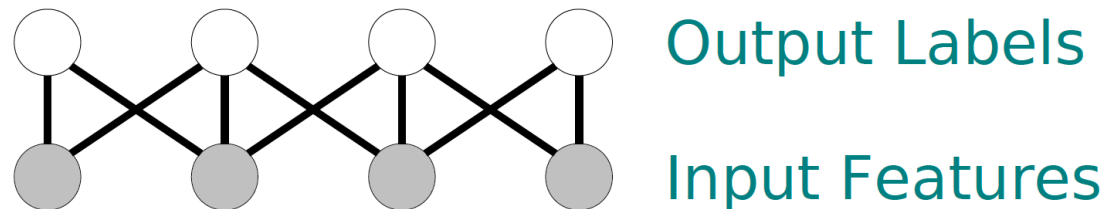
- **Structured models:** accurate but slow



- **Independent models:** less accurate but fast

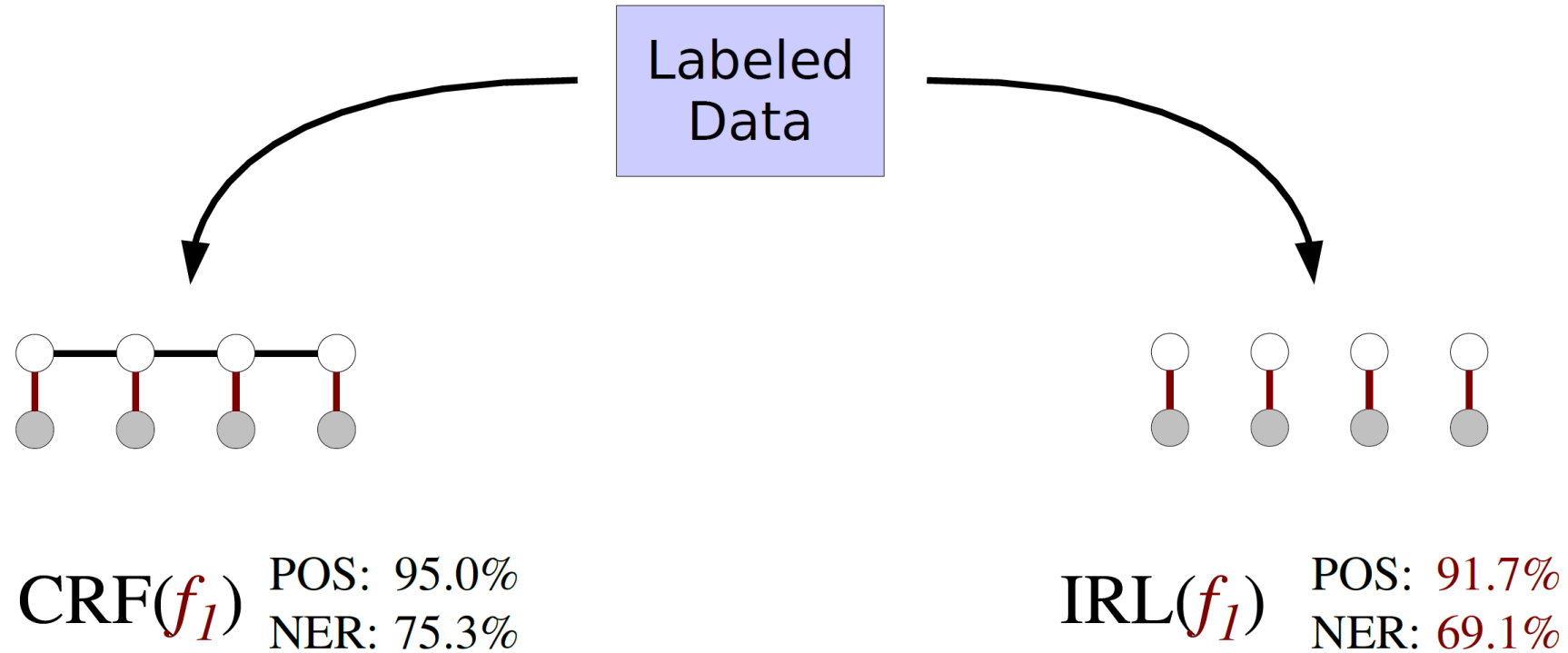


- **Goal:** transfer power to get fast+accurate



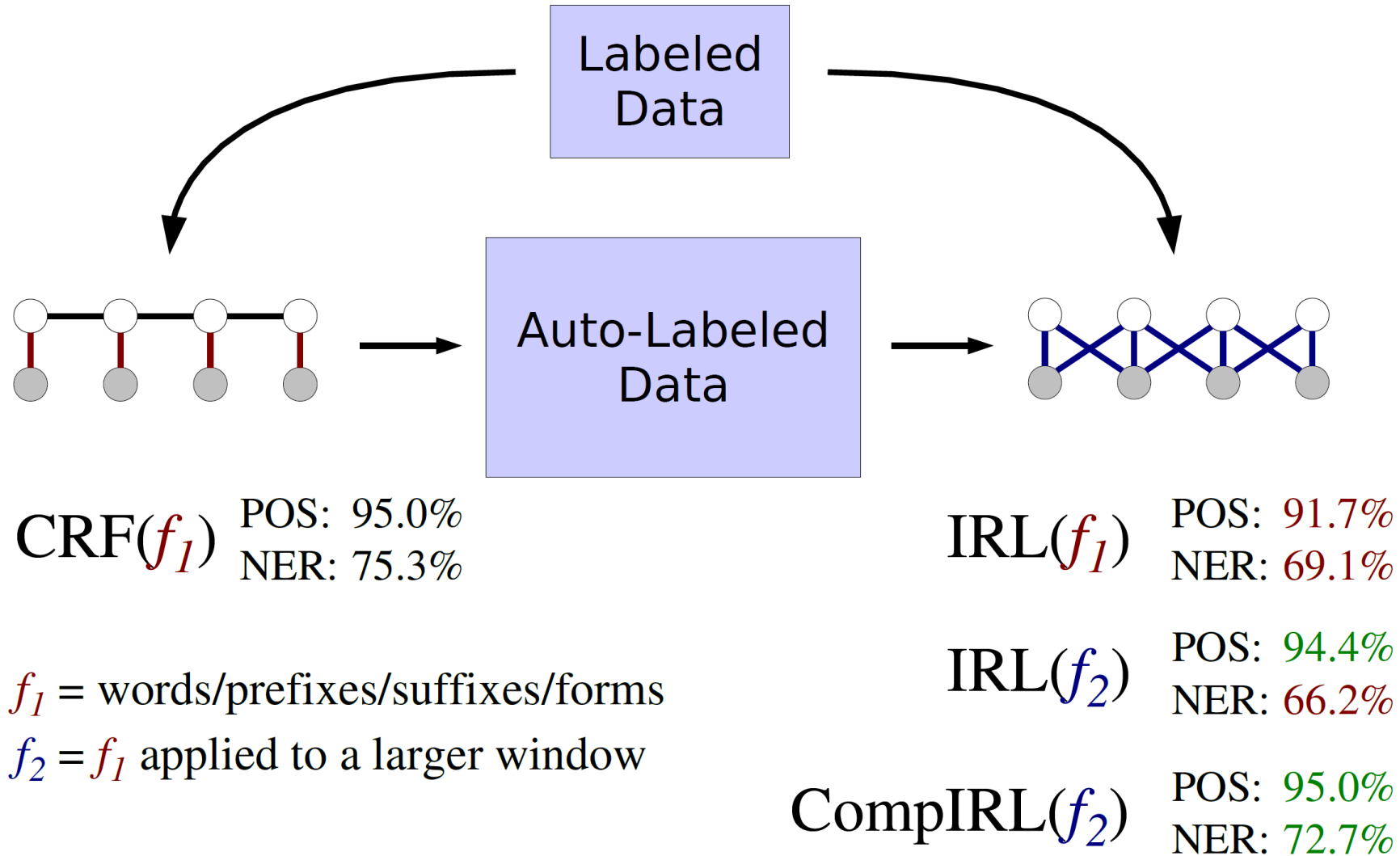
- **Questions:** are independent models...
 - ... expressive enough? (approximation error)
 - ... easy to learn? (estimation error)

“Compiling” structure away

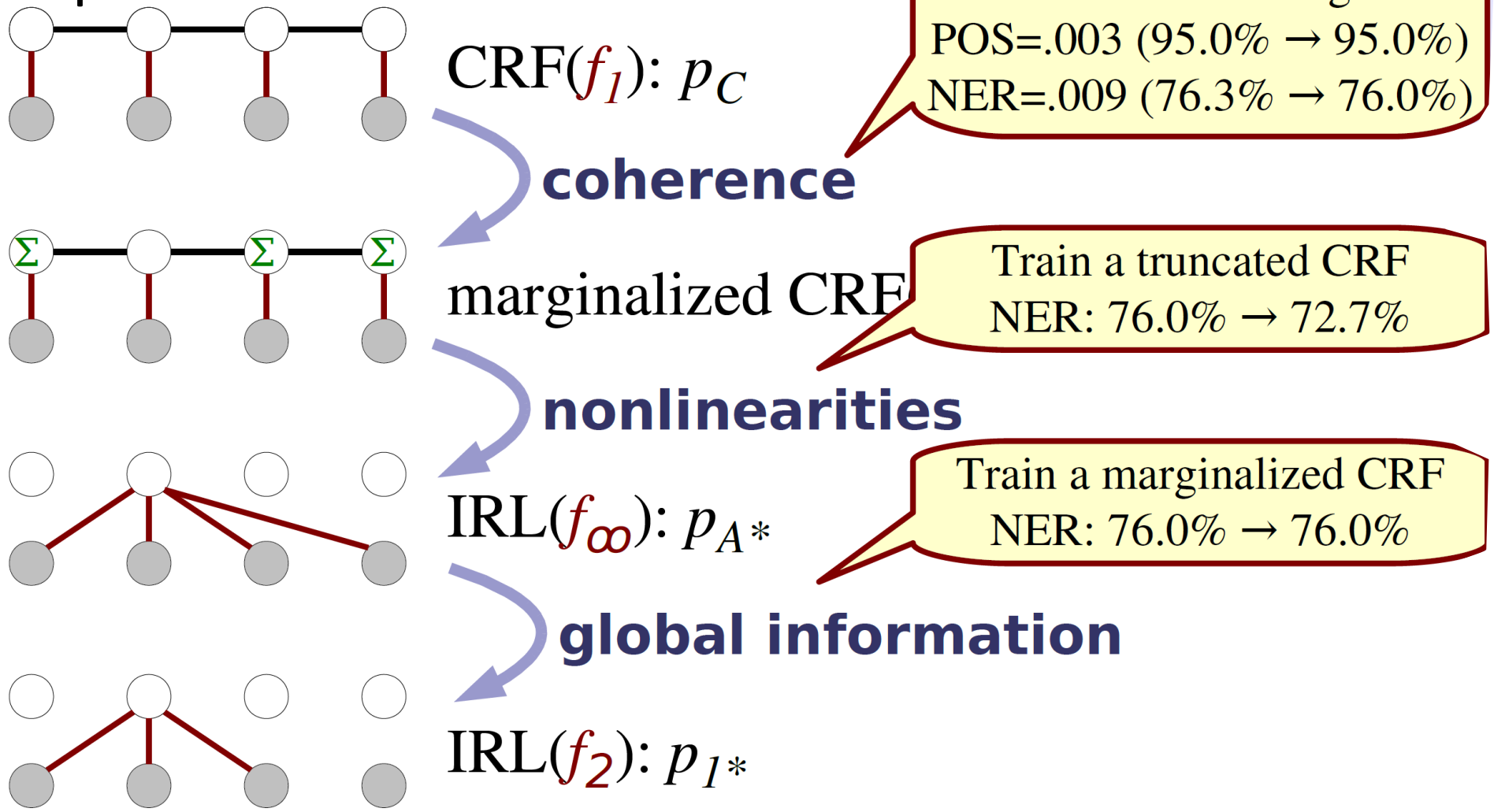


f_1 = words/prefixes/suffixes/forms

“Compiling” structure away



Decomposition of errors

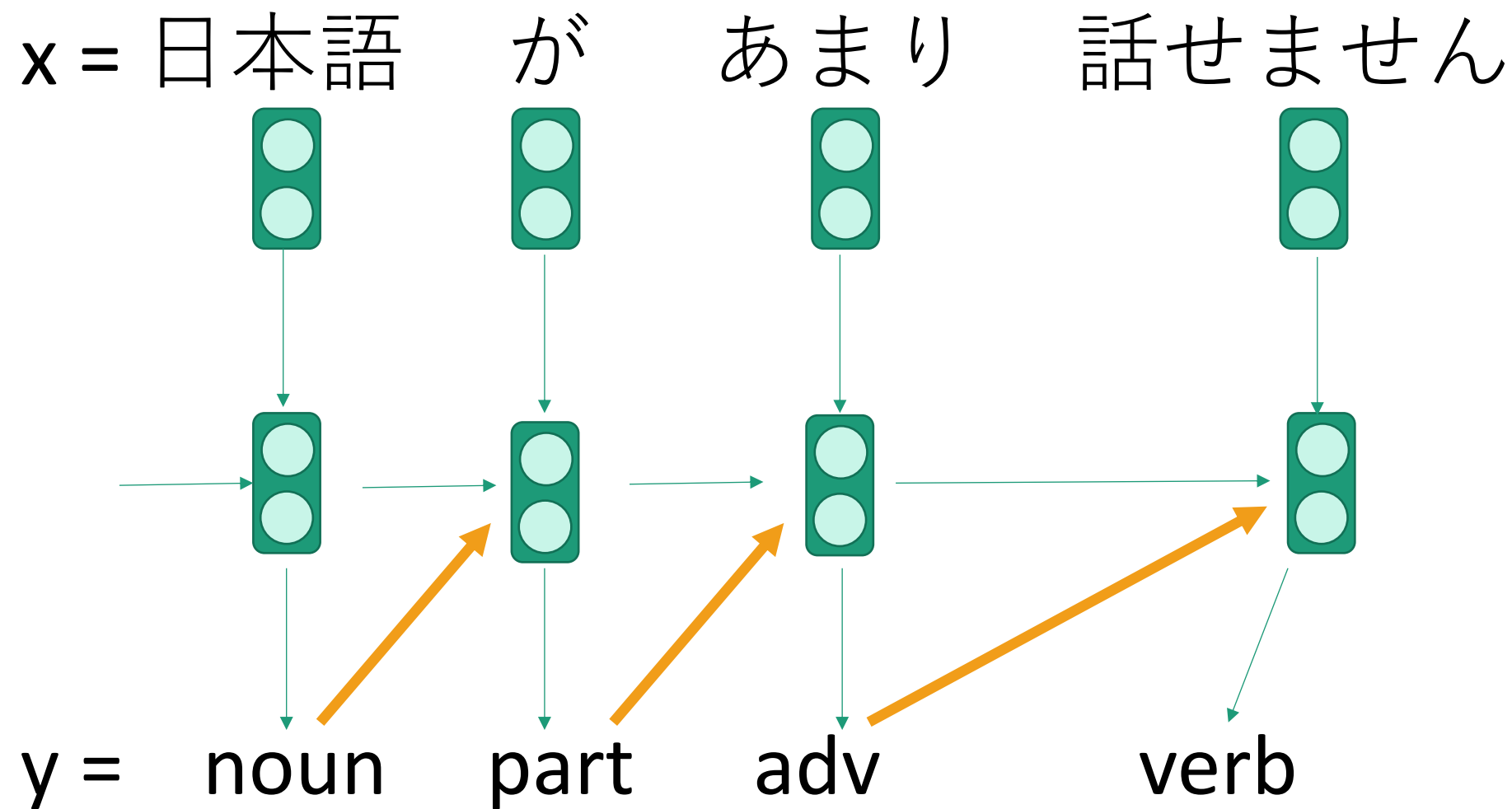


Theorem:

$$\text{KL}(p_C \parallel p_{I^*}) = \text{KL}(p_C \parallel p_{MC}) + \text{KL}(p_{MC} \parallel p_{A^*}) + \text{KL}(p_{A^*} \parallel p_{I^*})$$

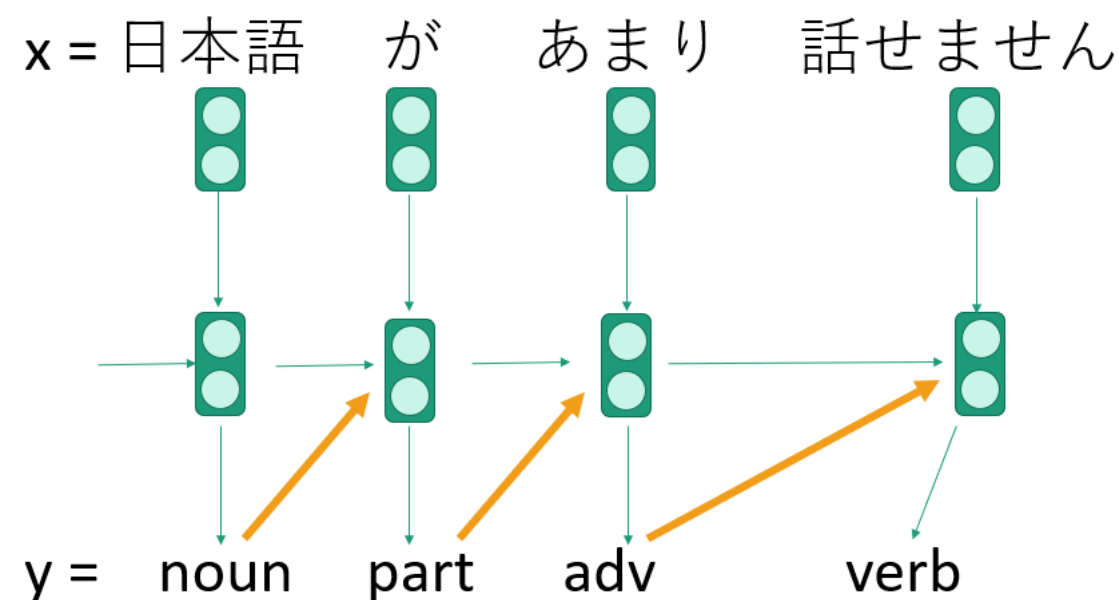
Non-linearities capture most information
... but structure still matters when not enough data

RNN sequence labeling



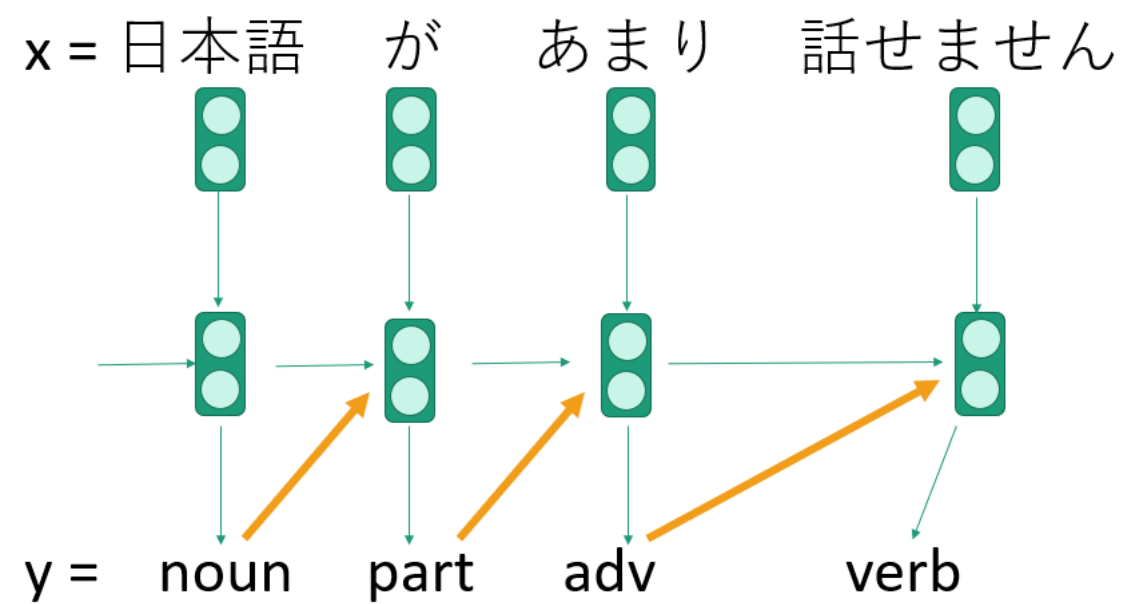
Prediction (aka “decoding”)

- Sampling:
- For $n = 1 \dots N$:
 - $y[n] \sim P(Y | X, y[1], \dots, y[n-1])$
- Return y



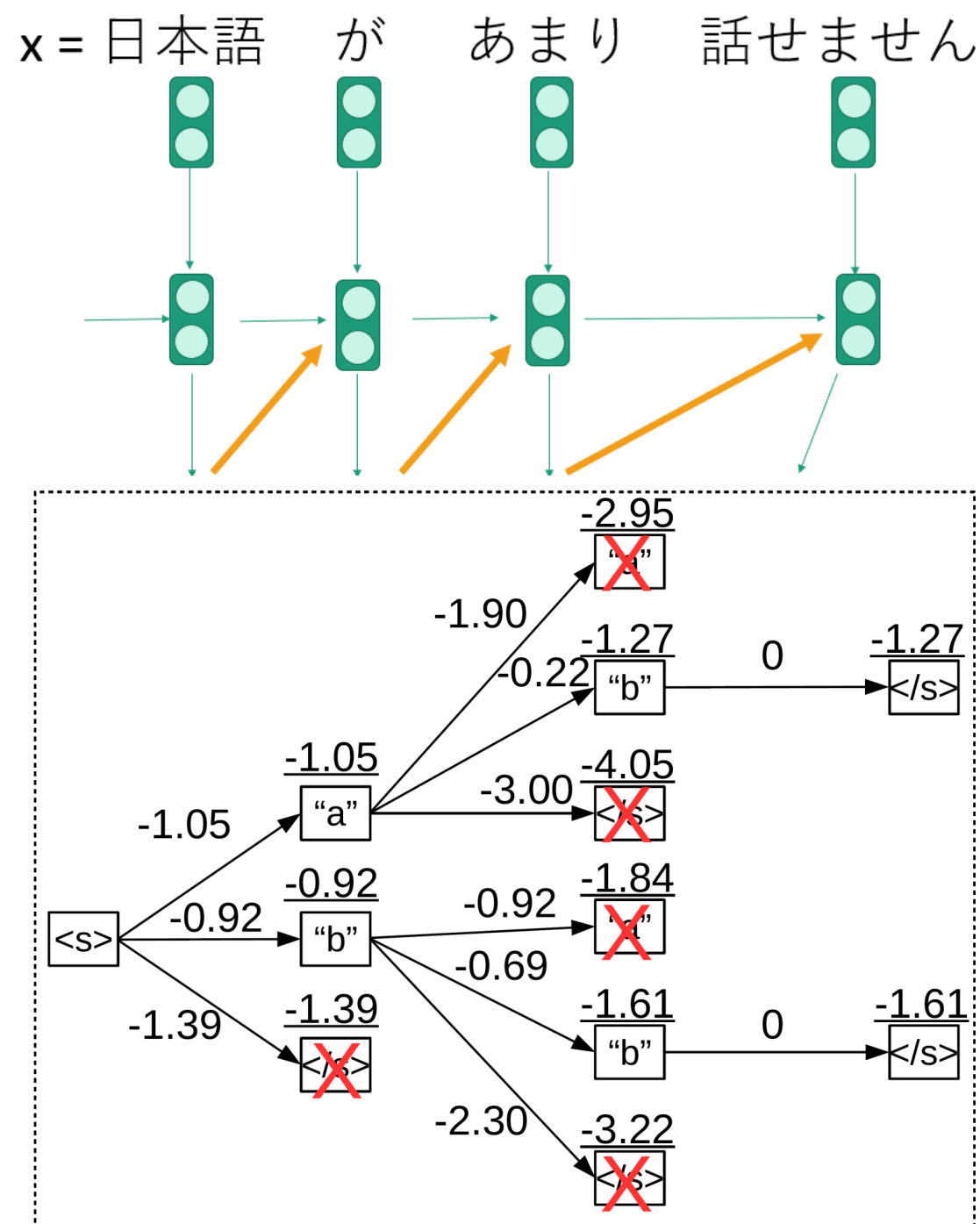
Prediction (aka “decoding”)

- Greedy search:
- For $n = 1 \dots N$:
 - $y[n] \leftarrow \operatorname{argmax}_y P(y \mid X, y[1], \dots, y[n-1])$
- Return y
- Issue: “label bias problem”



Prediction (aka “decoding”)

- Beam search:
- Consider b top hypothesis at each time step
- Example with beam size 2
 - (with log probs on edges)



Today

- Sequence labeling as independent predictions
- Structured perceptron for sequence labeling
- Do we really need structured features?
- Recurrent neural network taggers