

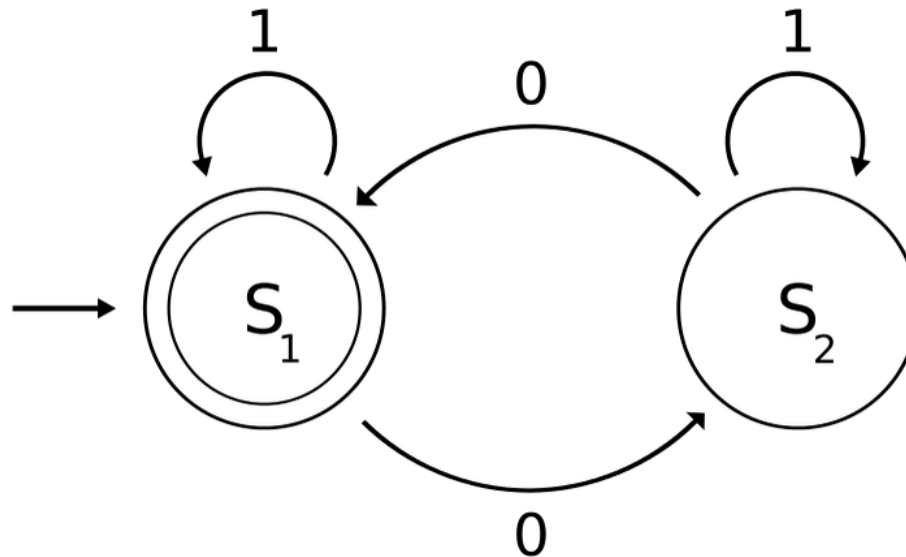
**Probabilistic Graphical Models:**  
**Hidden Markov Models & beyond**

Juan Miguel Cejuela

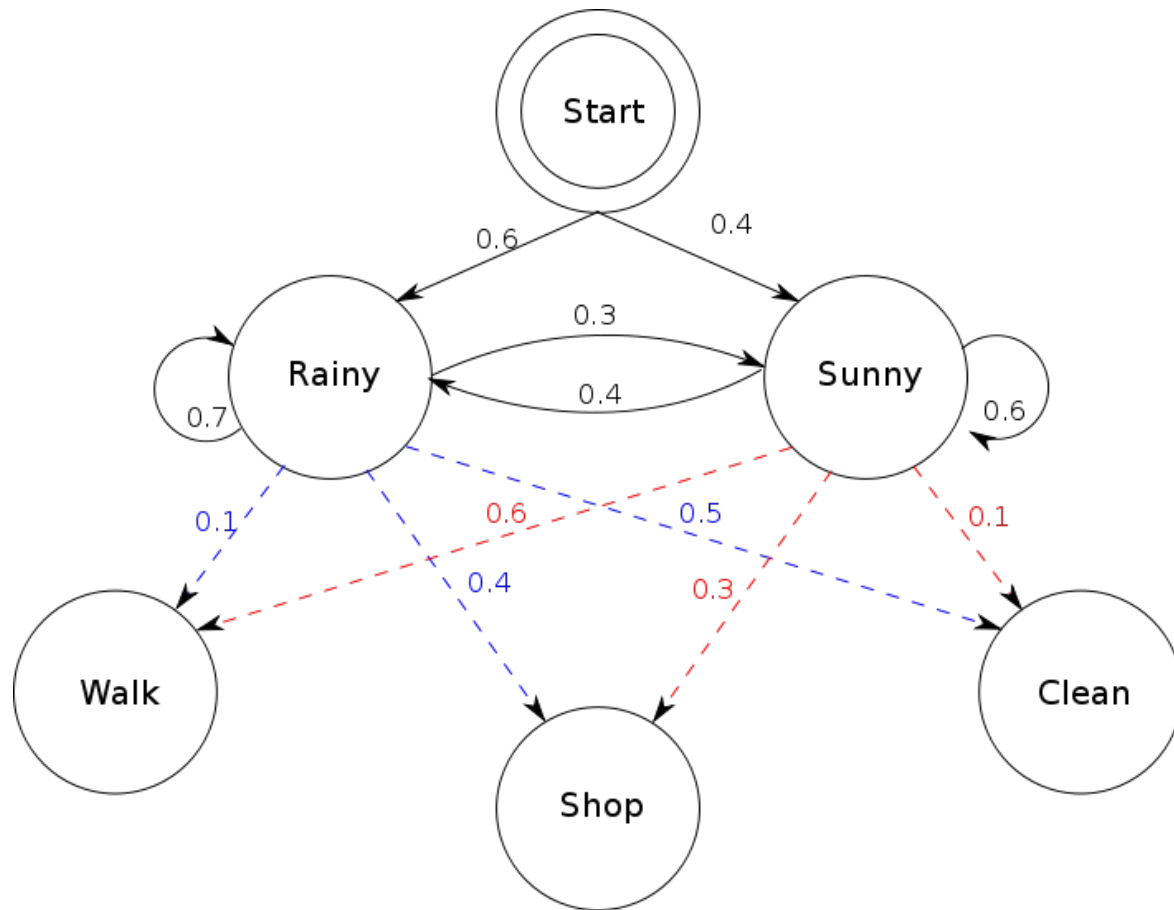
PhD Student in CS & Bioinformatics @ TUM



# (Probabilistic) Graphical Models



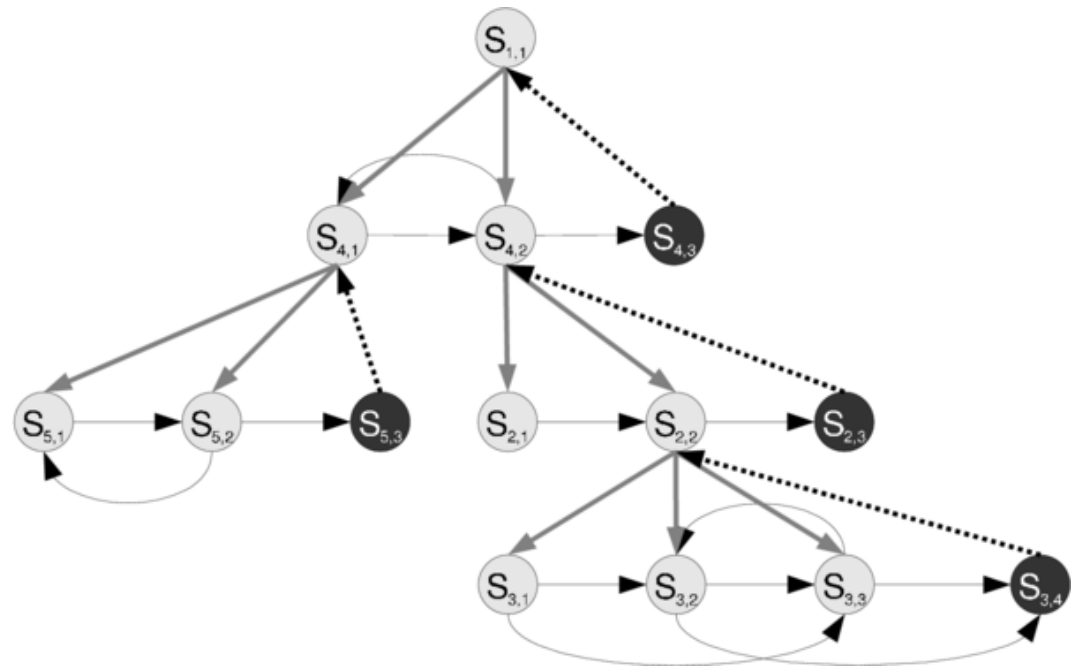
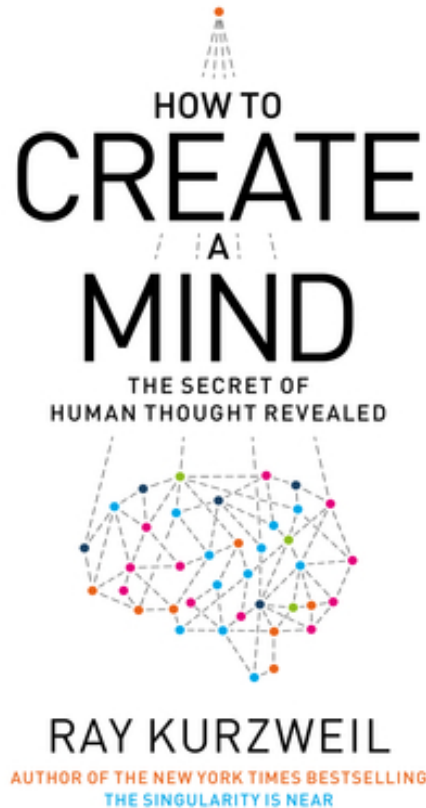
# Hidden Markov Models (HMM)



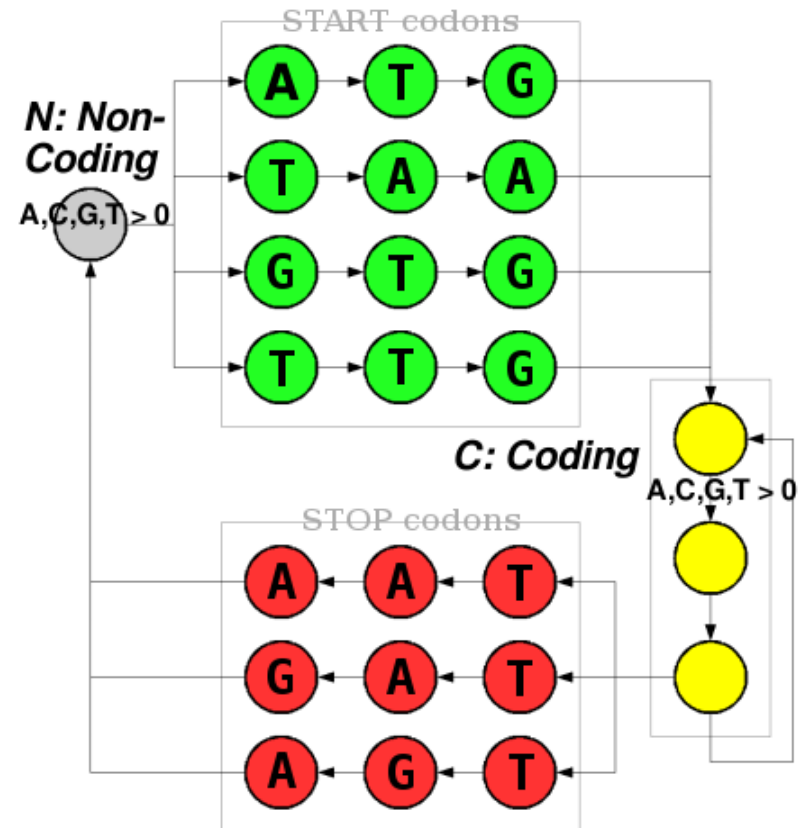
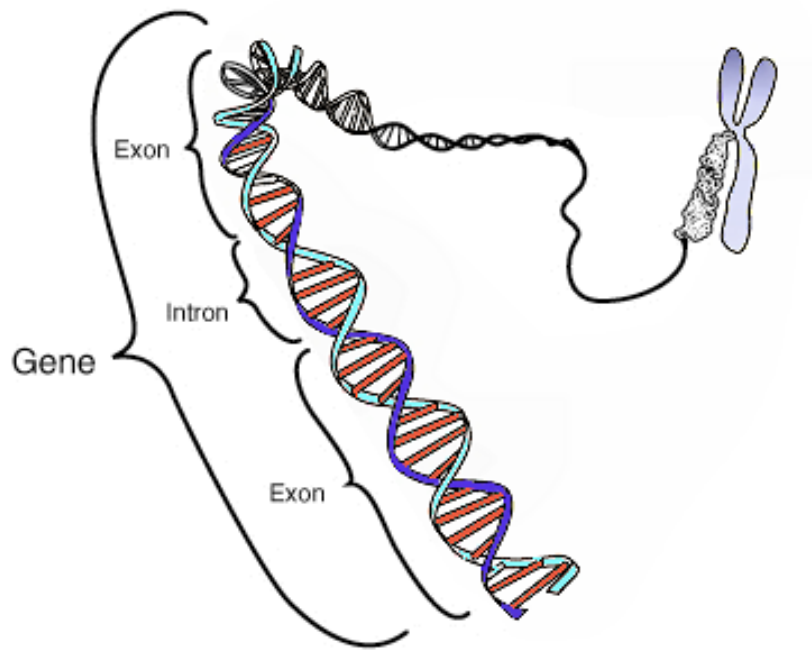
# History & Applications

- ~1960s, L. E. Baum & Ruslan L. Stratonovich
- Rabiner 1989 – A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition
- Sequence labeling
- Speech Recognition
- Gene Finding
- Machine Translation
- POS Tagging
- Protein Sequence Alignment
- ...

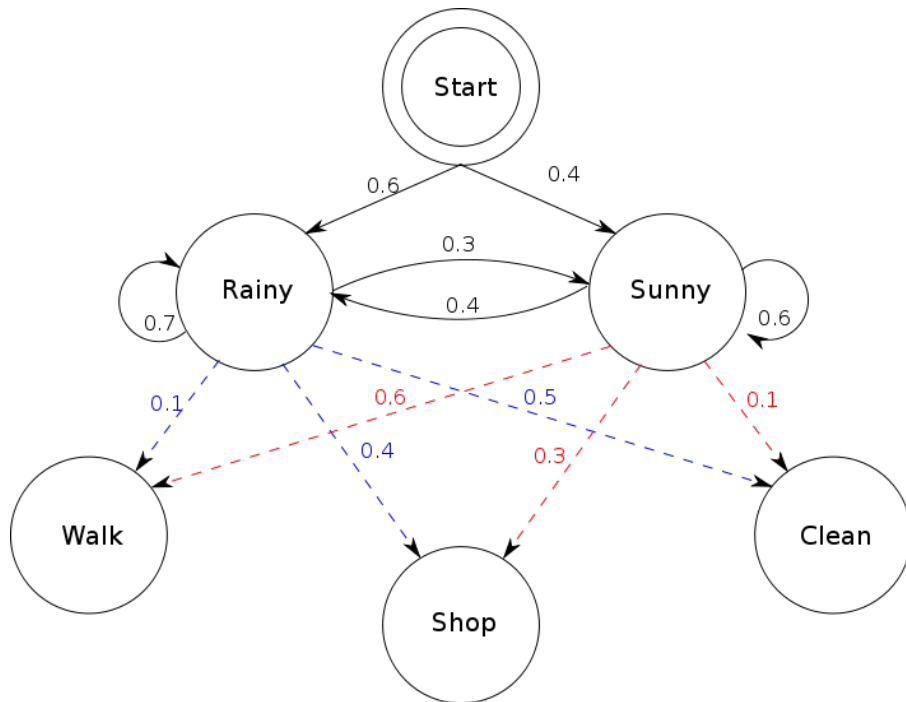
# Hierarchical Hidden Markov Models (HHMM)



# Application Case: Gene Finding



# HMM: Definition (I)



- Directed
- Generative
- Markov process

# HMM: Definition (II)

$$S = \{S_1, S_2, \dots S_N\}$$

$$O = \{O_1, O_2, \dots O_M\}$$

$$\pi = \{\pi_i\}$$

$$A = \{a_{ij}\}$$

$$B = \{b_j(k)\}$$

$$\lambda = \{\pi, A, B\}$$

```
states = ('Rainy', 'Sunny')

observations = ('walk', 'shop', 'clean')

start_probability = {'Rainy': 0.6, 'Sunny': 0.4}

transition_probability = {
    'Rainy' : {'Rainy': 0.7, 'Sunny': 0.3},
    'Sunny' : {'Rainy': 0.4, 'Sunny': 0.6},
}

emission_probability = {
    'Rainy' : {'walk': 0.1, 'shop': 0.4, 'clean': 0.5},
    'Sunny' : {'walk': 0.6, 'shop': 0.3, 'clean': 0.1},
}
```



# HMM: 3 problems

1. Scoring,  $P(X/\lambda)$
2. Decoding,  $\operatorname{argmax}_Y P(Y/X, \lambda)$
3. Training,  $\operatorname{argmax}_\lambda \prod P(X_k/\lambda)$

## Viterbi: decoding (predict state sequence)

- Find  $Y$ , path of states that best explains the given observation sequence  $X$
- Viterbi: path with highest probability

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1 \ q_2 \ \dots \ q_t = S_i, O_1 \ O_2 \ \dots \ O_t | \lambda]$$

## Initialisation

$$\begin{aligned}\delta_1(i) &= \pi_i b_i(O_1), \quad 1 \leq i \leq N \\ \psi_1(i) &= 0\end{aligned}$$

## Recursion

$$\begin{aligned}\delta_t(j) &= \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(O_t), \quad 2 \leq t \leq T \\ & \quad 1 \leq j \leq N\end{aligned}$$

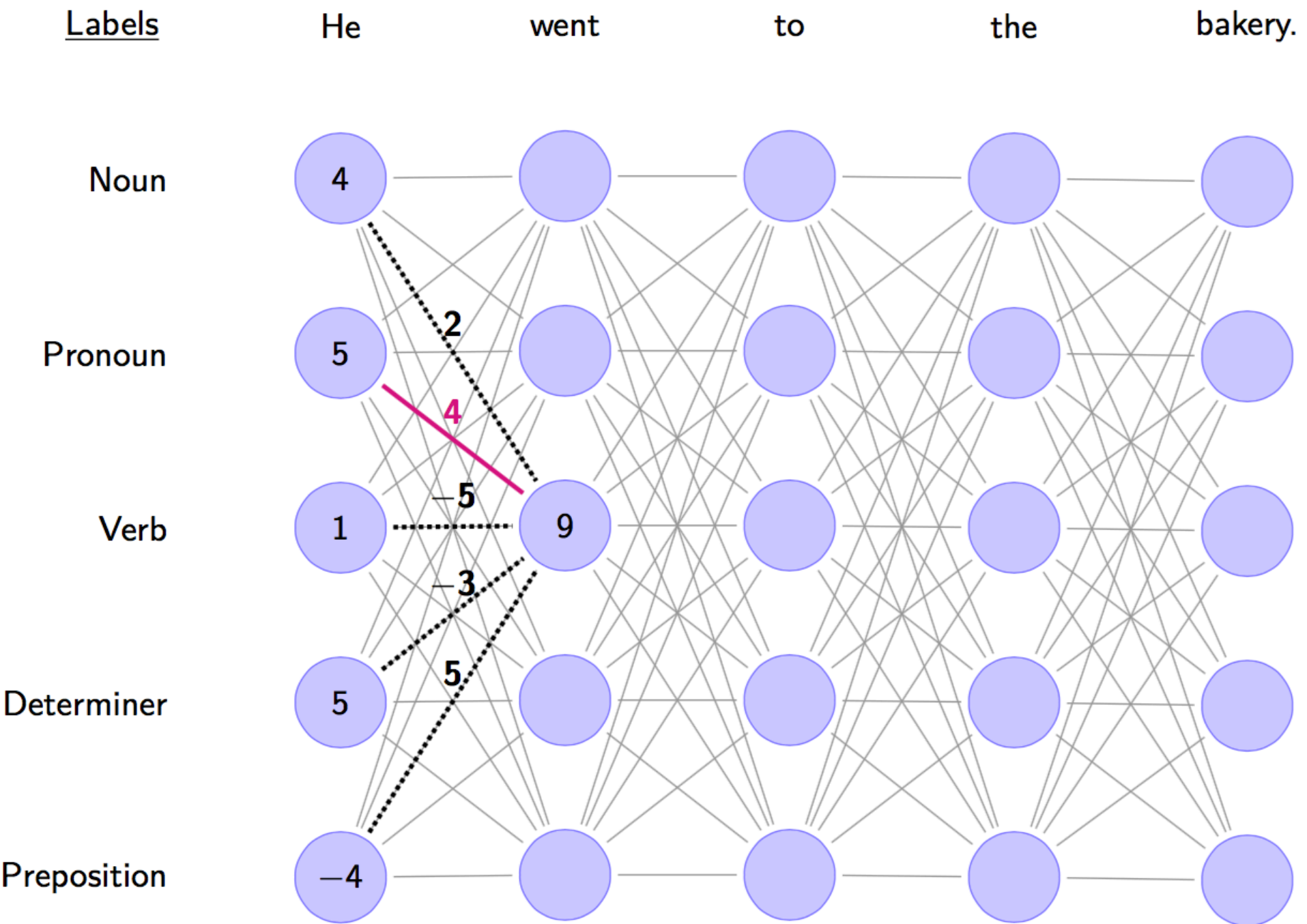
$$\begin{aligned}\psi_t(j) &= \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], \quad 2 \leq t \leq T \\ & \quad 1 \leq j \leq N\end{aligned}$$

## Termination

$$\begin{aligned}P^* &= \max_{1 \leq i \leq N} [\delta_T(i)] \\ q_T^* &= \arg \max_{1 \leq i \leq N} [\delta_T(i)]\end{aligned}$$

## Backtracking

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1$$



bakery.

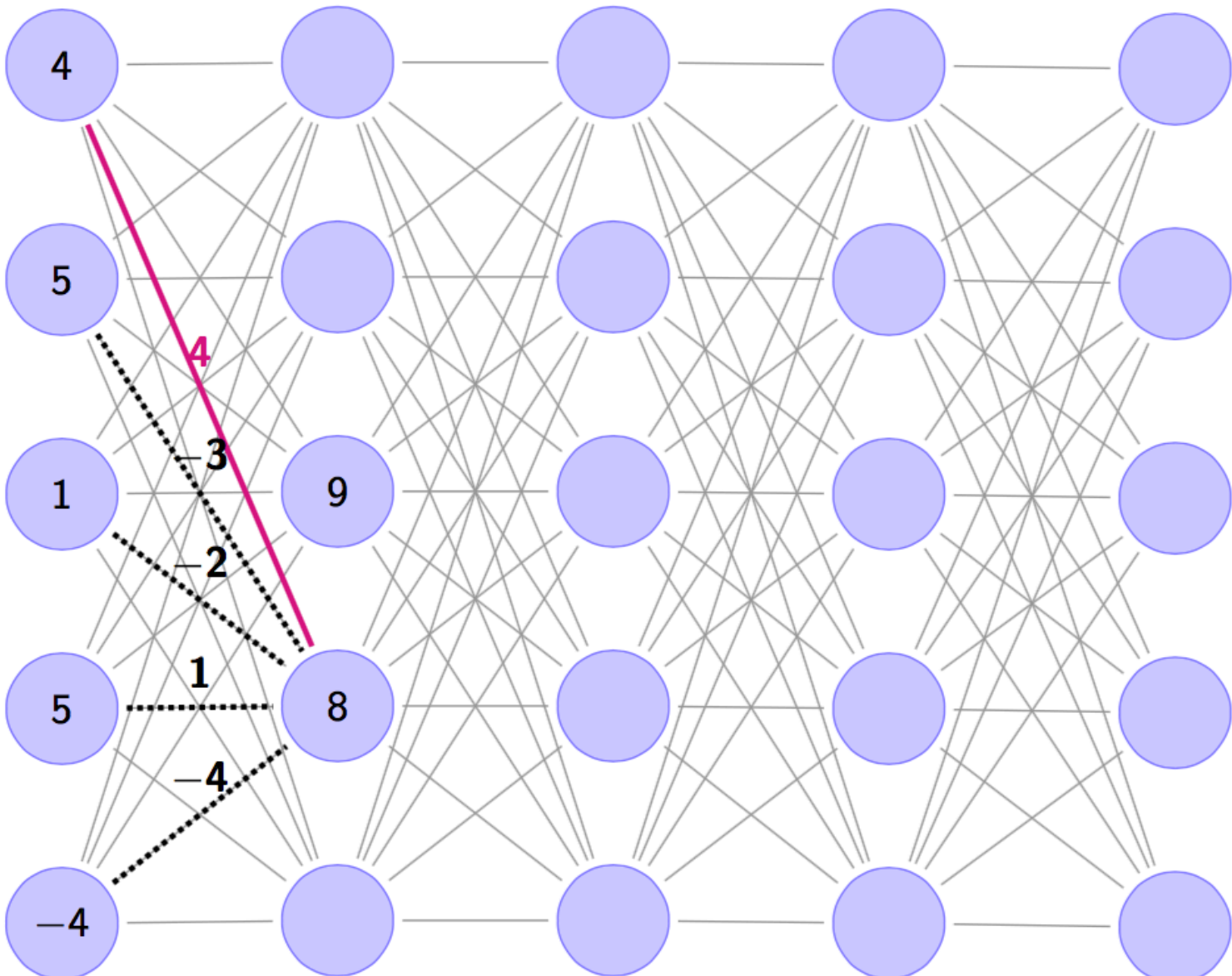
4

5

1

5

$-4$



bakery.

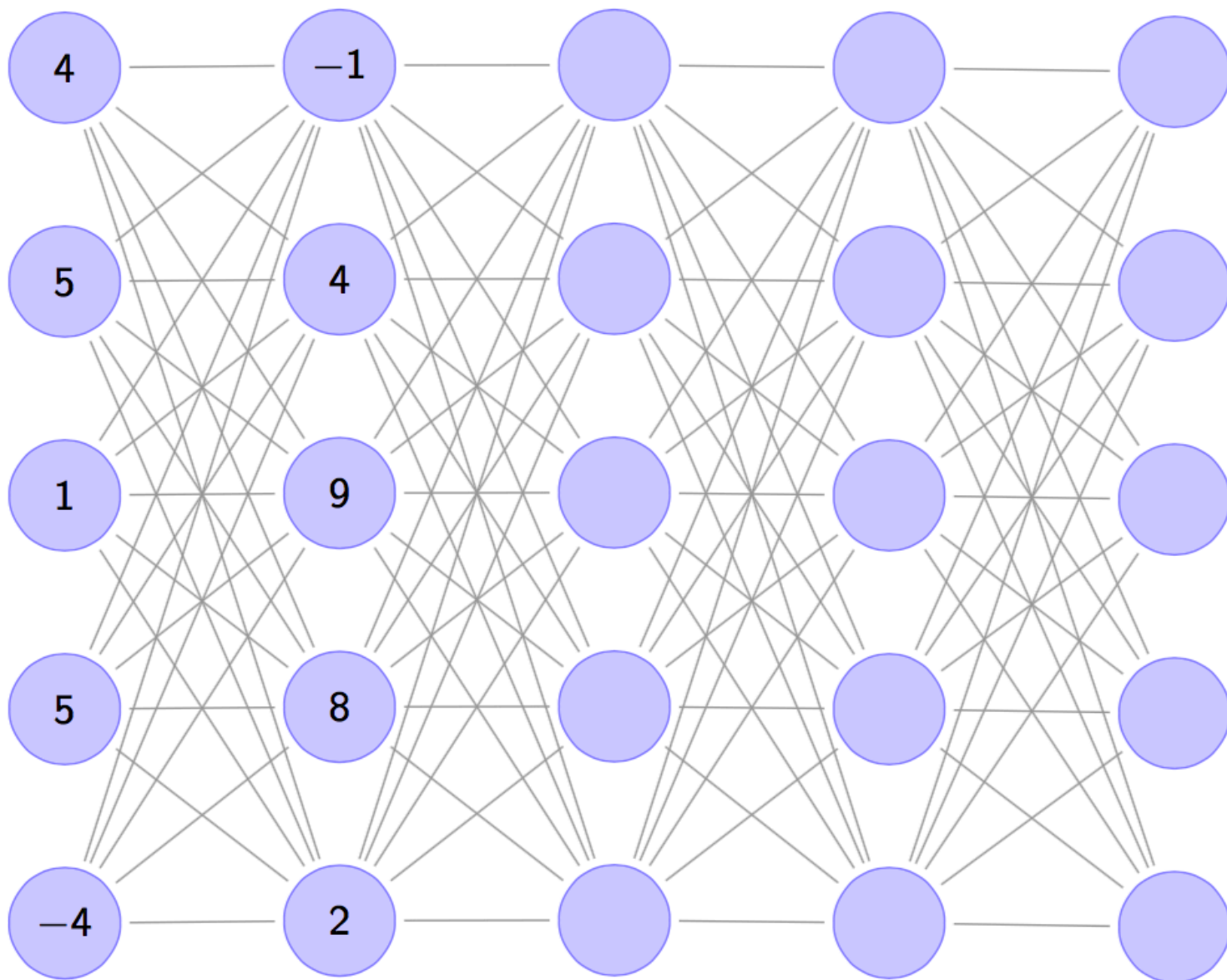
-1

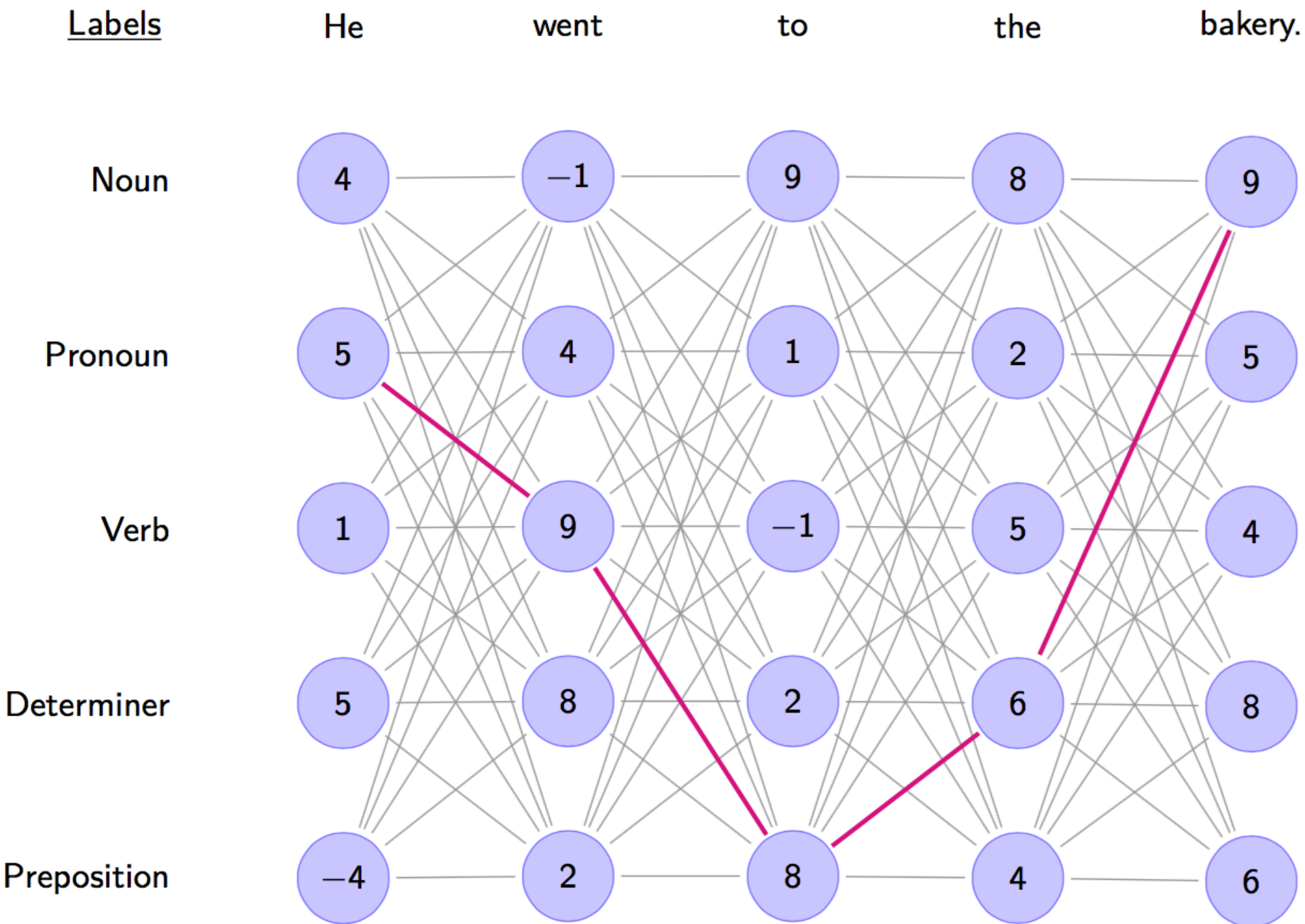
5

1

5

-4





# Training

- $\operatorname{argmax}_{\lambda} \prod P(X_k / \lambda)$
- Baum-Welch (forward & backward)  $\sim$  EM
- Supervised & semi-supervised learning

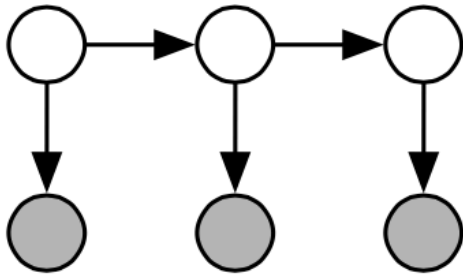


# Implementations

- [www.ghmm.org](http://www.ghmm.org) (C & Python)
- [github.com/jmcejuela/CL-HMM](https://github.com/jmcejuela/CL-HMM)

**log space!**

# Generative vs Discriminative

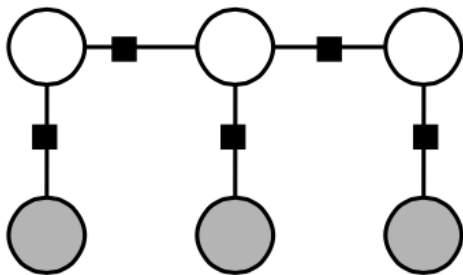


HMMs



$$P(X, Y) \text{ vs } P(Y/X)$$

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left\{ \sum_{k=1}^K \lambda_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}$$



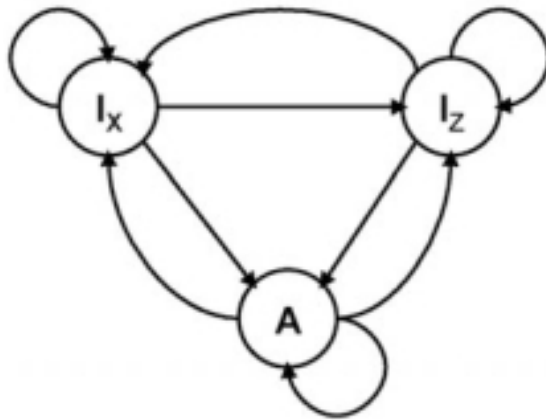
Linear-chain CRFs

# Pair Hidden Markov Models (PHMM)

- **Left & Right** observation alphabets
- $B = \{b_j(L, r)\}$
- $\operatorname{argmax}_Y P(Y/X, Z, \lambda)$
- forward & backward trickier in log space

# (Pairwise) Sequence Alignment

Pair HMM



$I_x$ : insertion in x (seq 1)

$I_z$ : insertion in z (seq 2)

A: aligned symbols in x and z

x (seq 1) : T T C C G - -

z (seq 2) : - - C C G T T

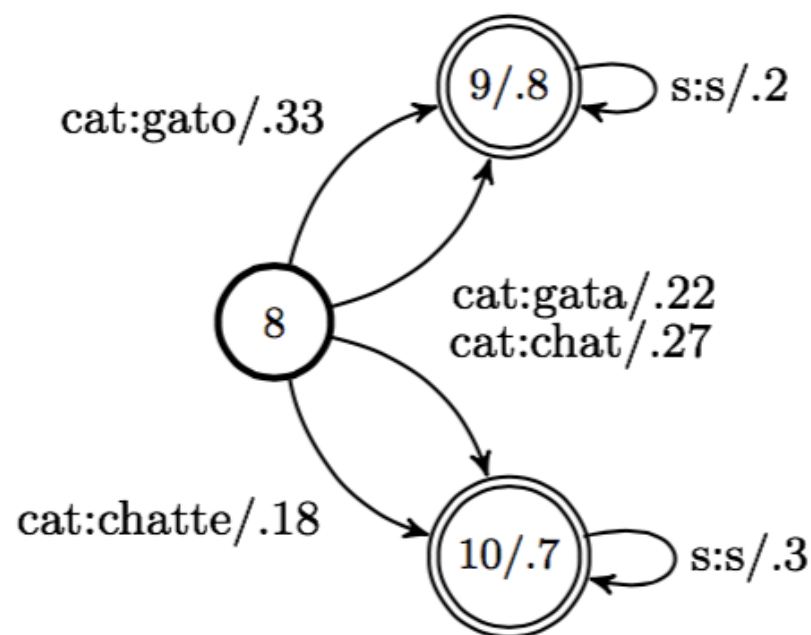
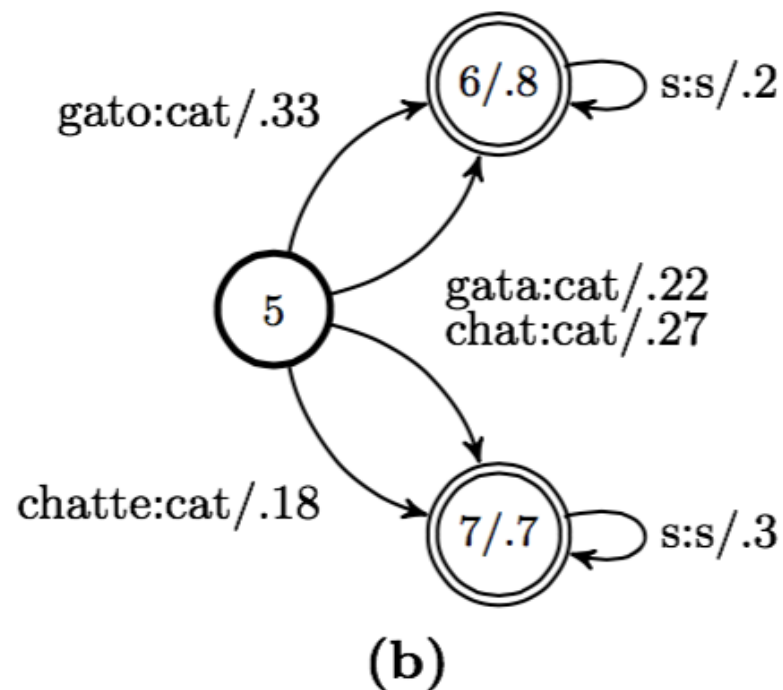
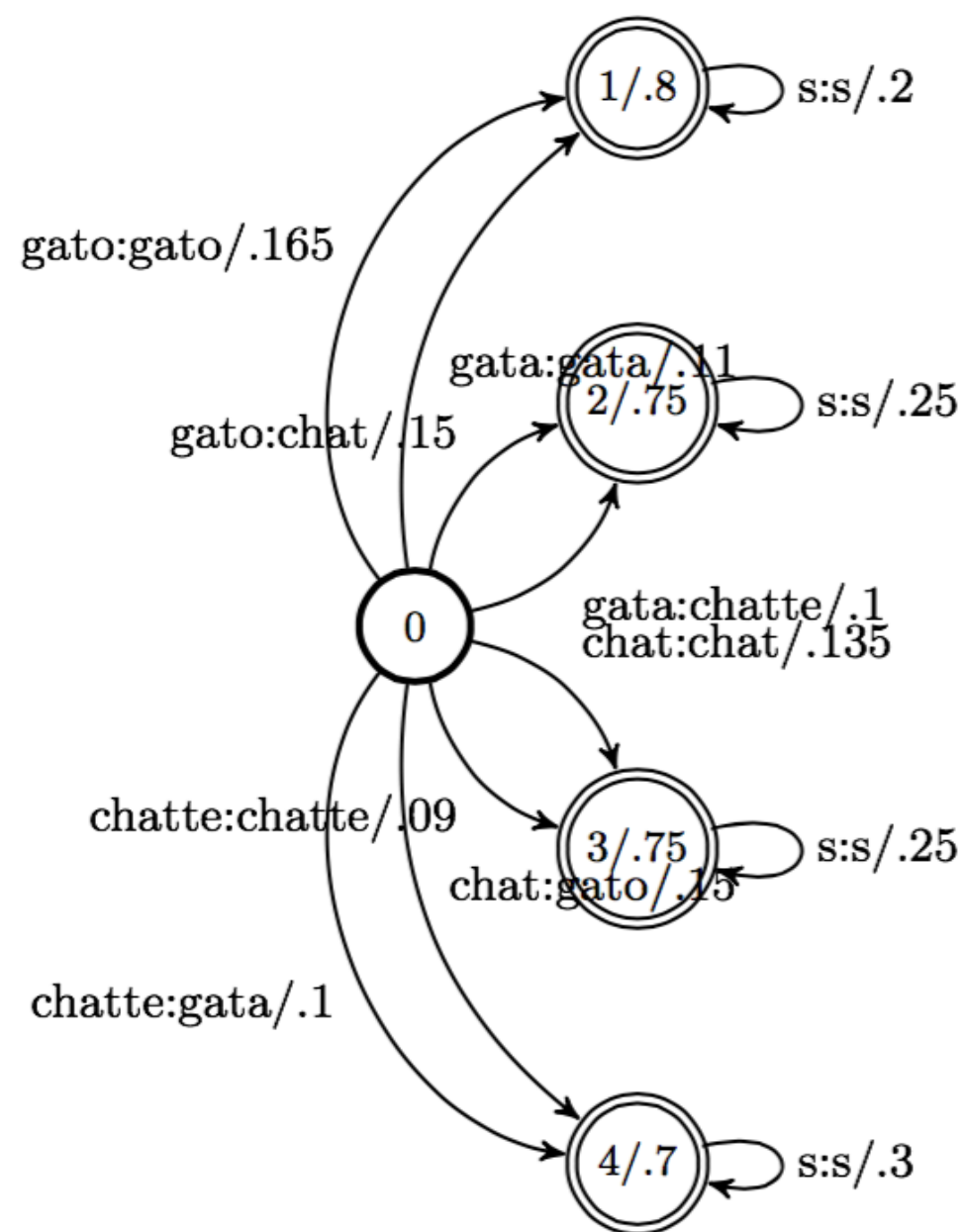
---

y (states) :  $I_x$   $I_x$  A A A  $I_z$   $I_z$

Levenshtein Distance, Needleman-Wunsch algorithm

# Machine Translation

- PHMMs as generalization of:
- **Weighted Finite State Transducers**

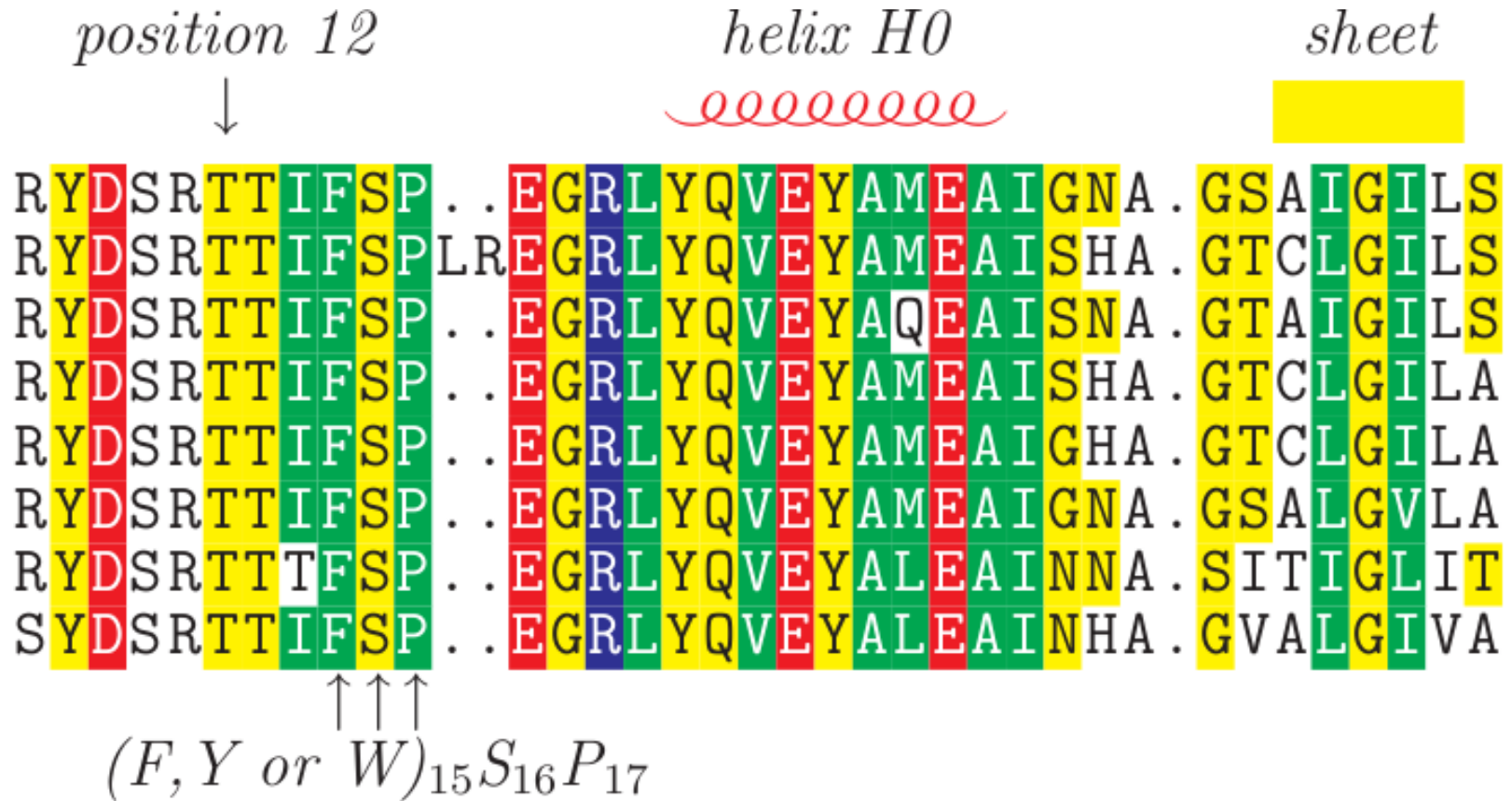


# WFST Implementation

- OpenFst Library (C++)

<http://openfst.org/twiki/bin/view/FST/WebHome>

# Profile-HMMs & Multiple Sequence Alignment





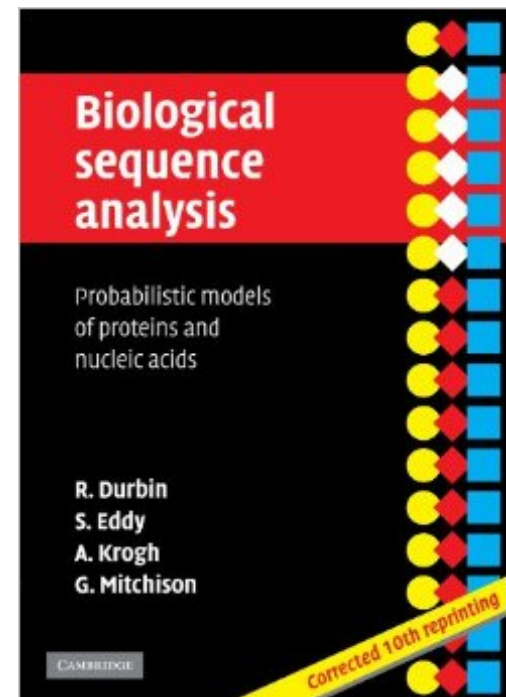
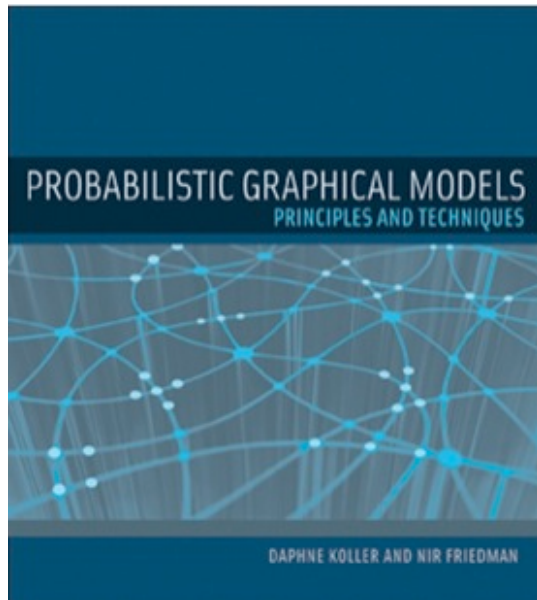
# More Complexity

- General Hidden Markov Models (GHMM)  
(linear vs fully connected models)
- Other orders: 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>, ...
- Discrete vs Continuous observations (speech recognition)
- Explicit State Duration Density

More...

Rabiner 1989

[coursera.org/course/pgm](https://coursera.org/course/pgm) by Daphne Koller



Thanks!

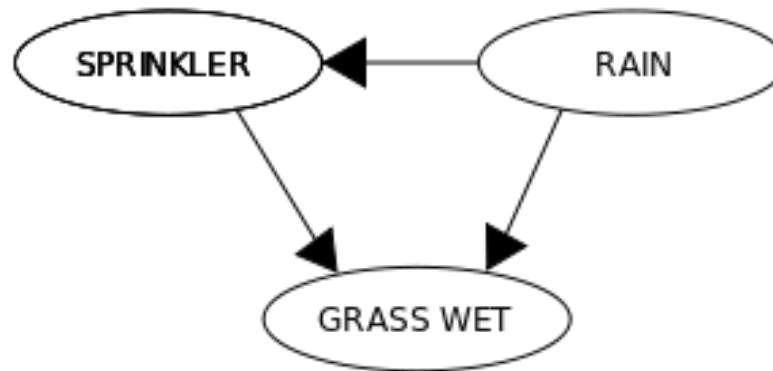
Juanmi @ tagtog.net

# Conditional Random Fields (CRF)

# Weighted Finite State Transducers (WFST)

# Bayesian Networks

RAIN	SPRINKLER	
	T	F
F	0.4	0.6
T	0.01	0.99



	RAIN	
	T	F
	0.2	0.8

		GRASS WET	
SPRINKLER	RAIN	T	F
F	F	0.0	1.0
F	T	0.8	0.2
T	F	0.9	0.1
T	T	0.99	0.01