



CS 7140: ADVANCED MACHINE LEARNING

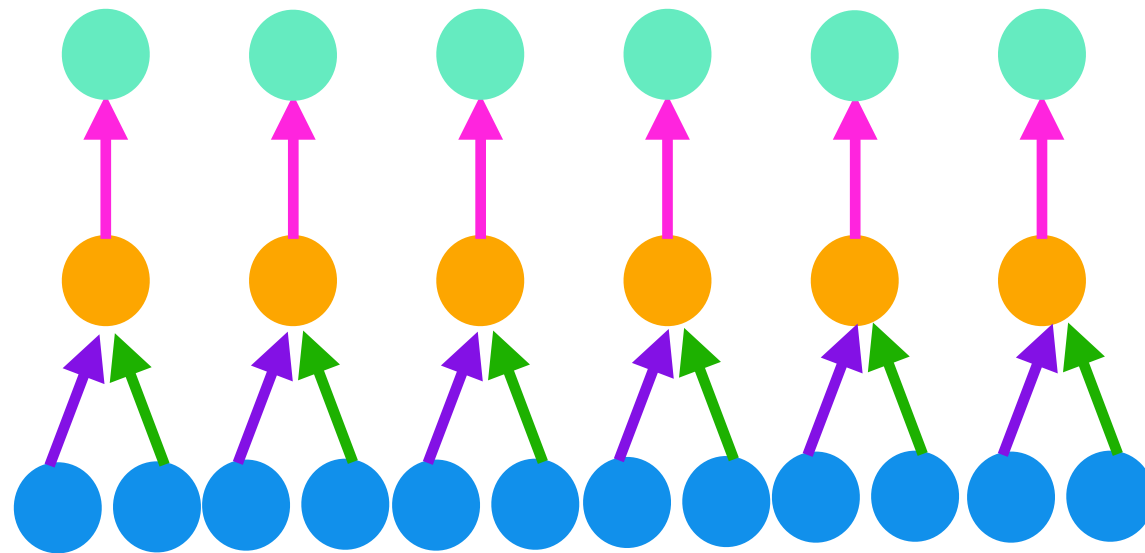
Recap: Structured Prediction

- A prediction function $f: \mathcal{X} \longrightarrow \mathcal{Y}$, where the output \mathcal{Y} domain is structured
 $y^\star = f(x) = \operatorname{argmax}_y g(x, y)$ MAP inference
- $g(x, y) = p(y | x)$ when the model is probabilistic
- Graph Cut, Local Search, Branch and Bound, Lagrangian Regularization.

Recap: CNN/RNN

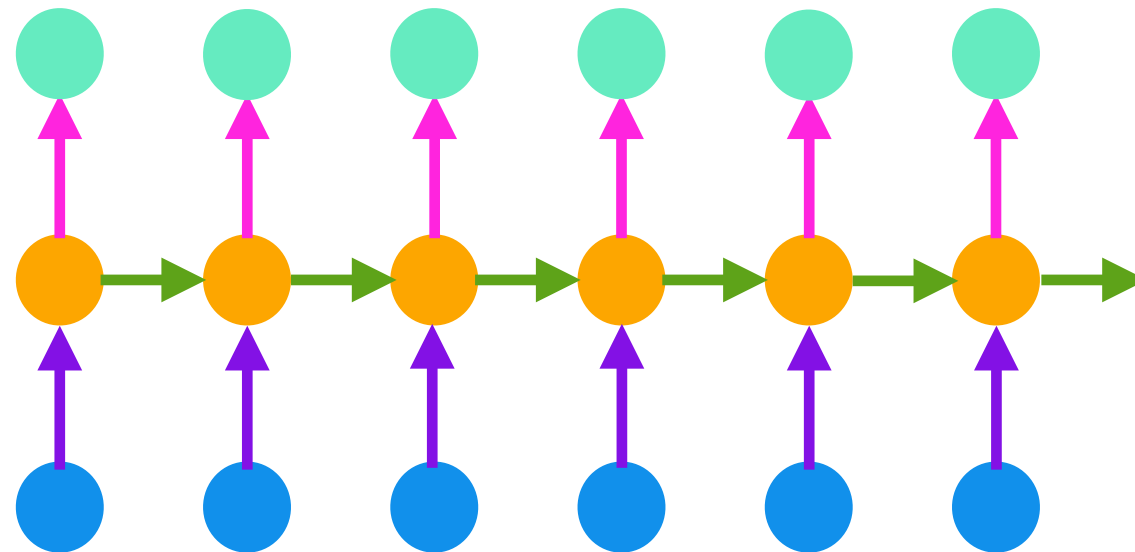
CNN

*sharing weights
across filters*



RNN

*sharing weights
across transitions*



INFERENCE AND HOPFIELD NETWORK

Exact Inference Problems

- The normalization problem:

$$P(\mathbf{x}) = \prod_{m=1}^M f(\mathbf{x}_m) \quad Z = \sum_{\mathbf{x}} \prod_{m=1}^M f(\mathbf{x}_m)$$

- The marginalization problem:

$$Z_n(x_n) = \sum_{\{x_{n'}\}, n' \neq n} P(\mathbf{x})$$

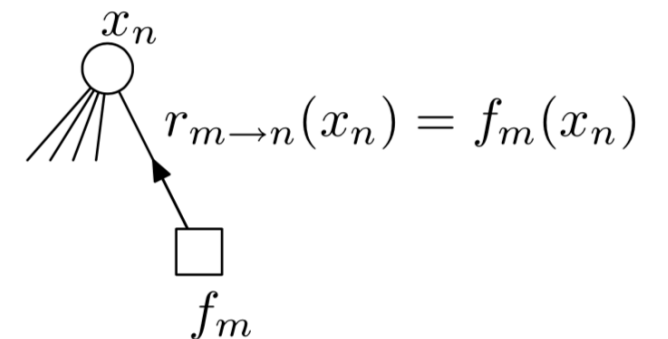
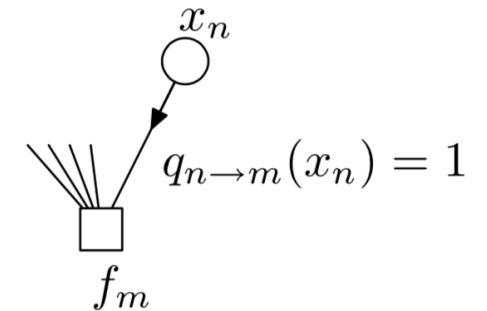
Message Passing

From variable to factor:

$$q_{n \rightarrow m}(x_n) = \prod_{m' \in \mathcal{M}(n) \setminus m} r_{m' \rightarrow n}(x_n). \quad (26.11)$$

From factor to variable:

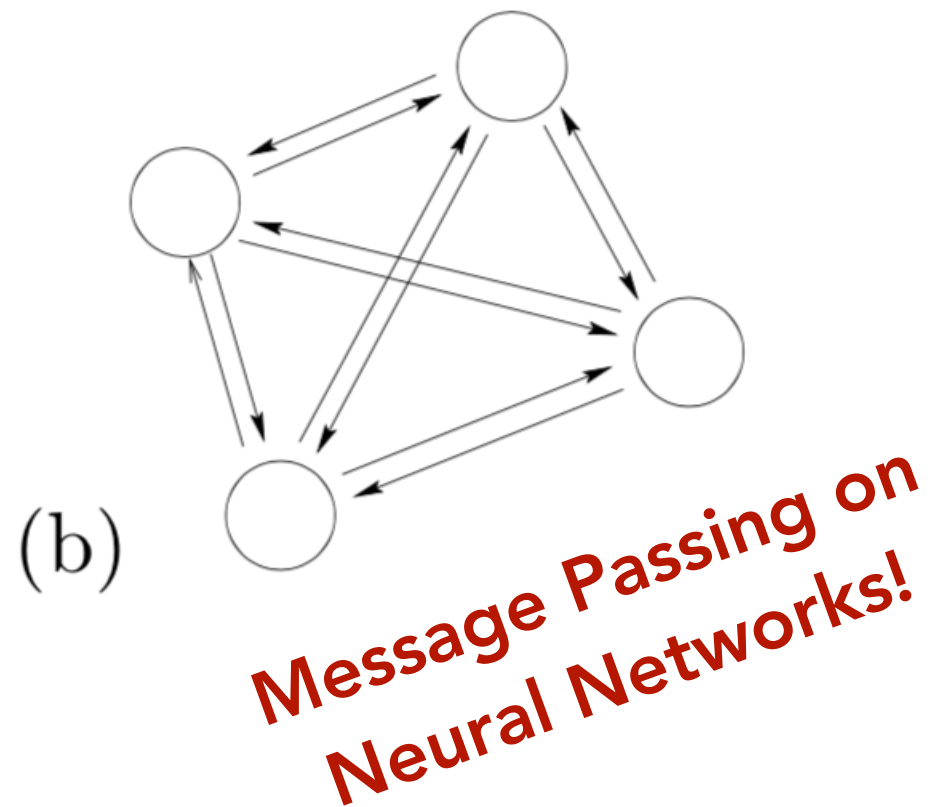
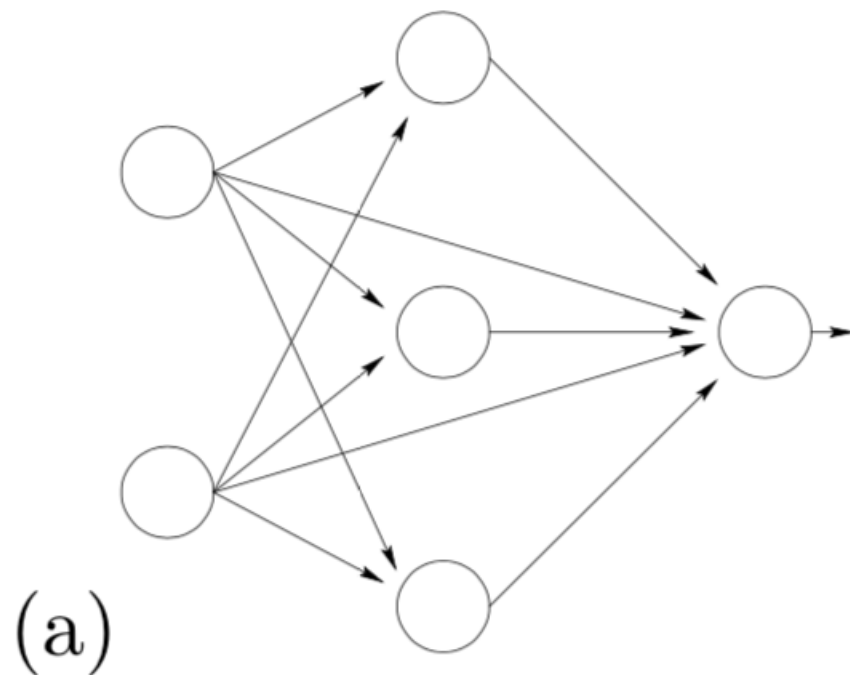
$$r_{m \rightarrow n}(x_n) = \sum_{\mathbf{x}_m \setminus n} \left(f_m(\mathbf{x}_m) \prod_{n' \in \mathcal{N}(m) \setminus n} q_{n' \rightarrow m}(x_{n'}) \right). \quad (26.12)$$



Two types of messages passing along the edges:

- messages $q_{n \rightarrow m}$ from variable nodes to factor nodes
- messages $r_{m \rightarrow n}$ from factor nodes to variable nodes.

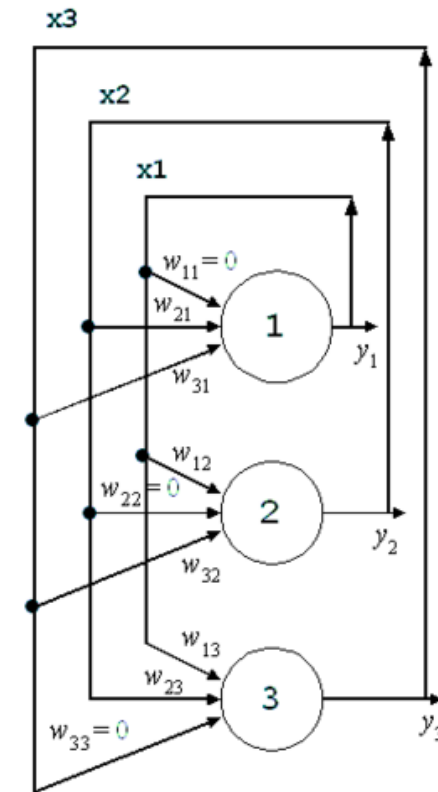
Feedforward/Feedback Networks



- **Feedforward:** all the connections are directed and form a directed acyclic graph
- **Feedback:** any network that is not a feedforward network
- **Applications:** Associated Memory, Optimization Problem

Binary Hopfield Network

- **Weights** w_{ij} denotes the connection from neuron i to neuron j
- **Architecture** symmetric, bidirectional connections $w_{ij} = w_{ji}$, no self-connections $w_{ii} = 0$
- **Activity rule** single neuron update
$$x(a) = \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases}$$
- **Updates** activations $a_i = \sum_j w_{ij}x_j$
- **Learning rule** make a set of desired memory $\{x^{(n)}\}$ be stable states, set weights $w_{ij} = \eta \sum x_i^{(n)} x_j^{(n)}$



Hopfield Network Convergence

- Binary probability distribution

$$P(x | \beta, J) = \frac{1}{Z} \exp[-\beta E(x; J)]$$

- Rewrite the expression

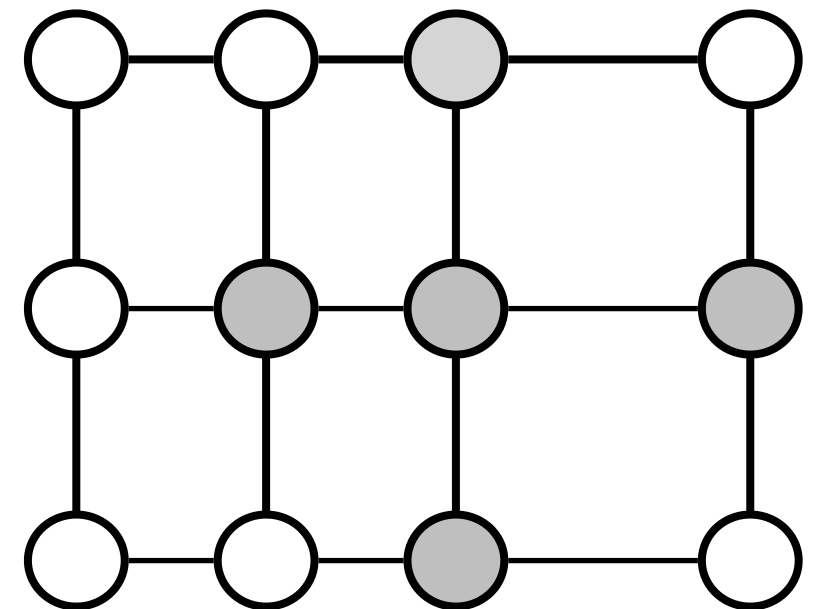
$$F(P^*, Q) = \beta \sum_x Q(x; \theta) E(x; J) - H_Q$$
$$\equiv \beta \mathbb{E}_Q[E(x; J)] - H_Q$$

mean energy

entropy

- Energy function

$$E(x; J) = -\frac{1}{2} \sum_{m,n} J_{mn} x_m x_n - \sum_n h_n x_n$$



Hopfield Network Convergence

- Variational Free Energy:

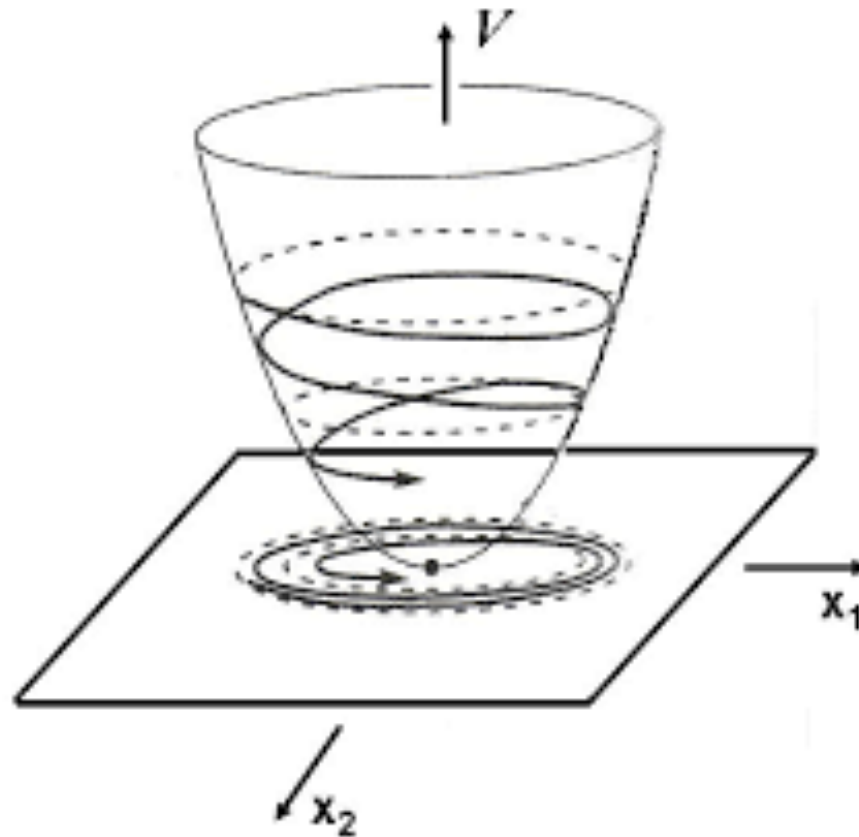
$$F(P^*, Q) = \beta \left(-\frac{1}{2} \sum_{mn} J_{mn} \mathbb{E}_Q[x_m] \mathbb{E}_Q[x_n] - \sum_n h_n \mathbb{E}_Q[x_n] \right) - H_Q$$

- Replace J by w , $\mathbb{E}_Q[x_n]$ by x_n , Hopfield network is identical to mean-field equations

$$F(P^*, Q) = -\beta \frac{1}{2} \mathbf{x}^\top \mathbf{W} \mathbf{x} - \sum_i H[(1 + x_n)/2] \quad \boxed{\text{Lyapunov function}}$$

- **Stability:** Hopfield network's dynamics will always converge to a stable fixed point

Lyapunov Function



- Generalized energy function for system, characterizes the dynamics of a system
- A function V is positive definite and $\dot{V}(x) \leq 0$ for all x

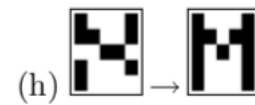
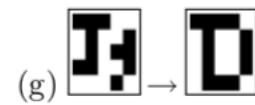
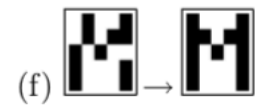
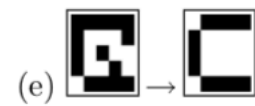
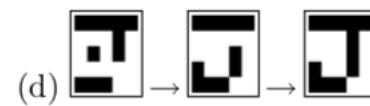
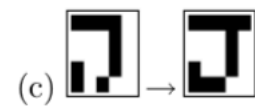
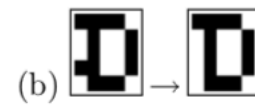
Associative Memory



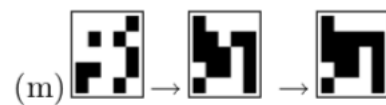
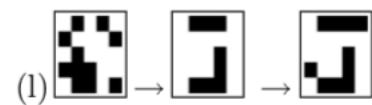
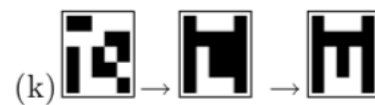
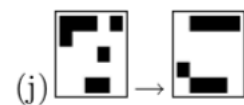
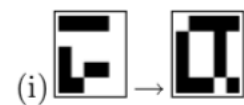
```

. 0 0 0 0 -2 2 -2 2 2 -2 0 0 0 2 0 0 -2 0 2 2 0 0 -2 -2
0 . 4 4 0 -2 -2 -2 -2 -2 -2 0 -4 0 -2 0 0 -2 0 -2 -2 4 4 2 -2
0 4 . 4 0 -2 -2 -2 -2 -2 -2 0 -4 0 -2 0 0 -2 0 -2 -2 4 4 2 -2
0 4 4 . 0 -2 -2 -2 -2 -2 -2 0 -4 0 -2 0 0 -2 0 -2 -2 4 4 2 -2
0 0 0 0 . 2 -2 -2 2 -2 2 -4 0 0 -2 4 -4 -2 0 -2 2 0 0 -2 2
-2 -2 -2 -2 2 . 0 0 0 0 4 -2 2 -2 0 2 -2 0 -2 0 0 -2 -2 0 4
2 -2 -2 -2 -2 0 . 0 0 4 0 2 2 -2 4 -2 2 0 -2 4 0 -2 -2 0 0
-2 -2 -2 -2 -2 0 0 . 0 0 0 2 2 2 0 -2 2 4 2 0 0 -2 -2 0 0
2 -2 -2 -2 2 0 0 0 . 0 0 -2 2 2 0 2 -2 0 2 0 4 -2 -2 -4 0
2 -2 -2 -2 -2 0 4 0 0 . 0 2 2 -2 4 -2 2 0 -2 4 0 -2 -2 0 0
-2 -2 -2 -2 2 4 0 0 0 0 . -2 2 -2 0 2 -2 0 -2 0 0 -2 -2 0 4
0 0 0 0 -4 -2 2 2 -2 2 -2 . 0 0 2 -4 4 2 0 2 -2 0 0 2 -2
0 -4 -4 -4 0 2 2 2 2 2 2 0 . 0 2 0 0 2 0 2 2 -4 -4 -2 2
0 0 0 0 0 -2 -2 2 2 -2 -2 0 0 . -2 0 0 2 4 -2 2 0 0 -2 -2
2 -2 -2 -2 -2 0 4 0 0 4 0 2 2 -2 . -2 2 0 -2 4 0 -2 -2 0 0
0 0 0 0 4 2 -2 -2 2 -2 2 -4 0 0 -2 . -4 -2 0 -2 2 0 0 -2 2
0 0 0 0 -4 -2 2 2 -2 2 -2 4 0 0 2 -4 . 2 0 2 -2 0 0 2 -2
-2 -2 -2 -2 -2 0 0 4 0 0 0 2 2 2 0 -2 2 . 2 0 0 -2 -2 0 0
0 0 0 0 0 -2 -2 2 2 -2 -2 0 0 4 -2 0 0 2 . -2 2 0 0 -2 -2
2 -2 -2 -2 -2 0 4 0 0 4 0 2 2 -2 4 -2 2 0 -2 . 0 -2 -2 0 0
2 -2 -2 -2 2 0 0 0 4 0 0 -2 2 2 0 2 -2 0 2 0 . -2 -2 -4 0
0 4 4 4 0 -2 -2 -2 -2 -2 -2 0 -4 0 -2 0 0 -2 0 -2 -2 . 4 2 -2
0 4 4 4 0 -2 -2 -2 -2 -2 -2 0 -4 0 -2 0 0 -2 0 -2 -2 4 . 2 -2
-2 2 2 2 -2 0 0 0 -4 0 0 2 -2 -2 0 -2 2 0 -2 0 -4 2 2 . 0
-2 -2 -2 -2 2 4 0 0 0 0 4 -2 2 -2 0 2 -2 0 -2 0 0 -2 -2 0 .

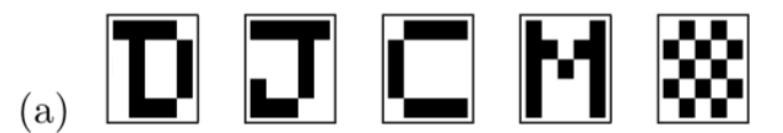
```



- 25 -unit binary Hopfield network
- (a) four patterns as 5x5 binary images
- (b)- (m) evolution of state of the network

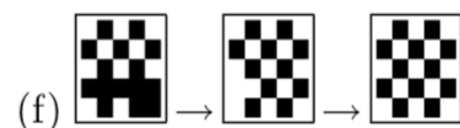
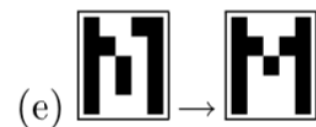
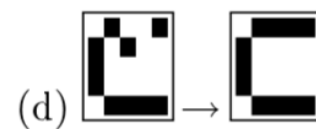
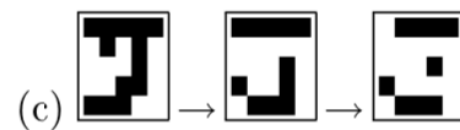
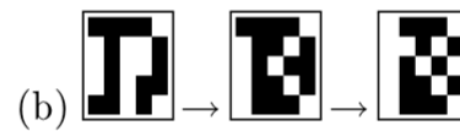


Robustness



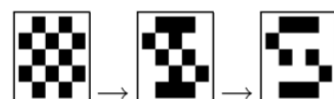
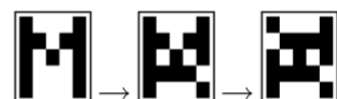
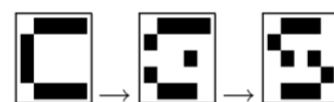
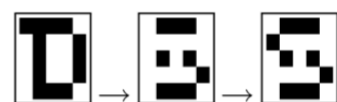
```

-1 1-1 1 x x-3 3 x x-1 1-1 x-1 1-3 x 1 3-1 1 x-1
-1 . 3 5-1-1-3-1-3-1-3 1 x 1-3 1-1-1-1-1-3 5 3 3-3
1 3 . 3 1-3-1 x-1-3-1-1 x-1-1-1 1-3 1-3-1 3 5 1-1
-1 5 3 .-1-1-3-1-3-1-3 1-5 1-3 1-1-1-1-1-3 5 x 3-3
1-1 1-1 . 1-1-3 x x 3-5 1-1-1 3 x-3 1-3 3-1 1-3 3
x-1-3-1 1 .-1 1-1 1 3-1 1-1-1 3-3 1 x 1 x-1-3 1 3
x-3-1-3-1-1 .-1 1 3 1 1 3-3 5-3 3-1-1 x 1-3-1-1 1
-3-1 x-1-3 1-1 .-1 1-1 3 1 x-1-1 1 5 1 1-1 x-3 1-1
3-3-1-3 x-1 1-1 .-1 1-3 3 1 1 1-1-1 3-1 5-3-1 x 1
x-1-3-1 x 1 3 1-1 .-1 3 1-1 3-1 x 1-3 5-1-1-3 1-1
x-3-1-3 3 3 1-1 1-1 .-3 3-3 1 1-1-1-1-1 1-3-1-1 5
-1 1-1 1-5-1 1 3-3 3-3 .-1 1 1-3 3 x-1 3-3 1-1 3-3
1 x x-5 1 1 3 1 3 1-1 .-1 3-1 1 1 1 1 3-5-3-3 3
-1 1-1 1-1-1-3 x 1-1-3 1-1 . x 1-1 3 3-1 1 1-1-1-3
x-3-1-3-1-1 5-1 1 3 1 1 3 x . x 3-1-1 3 1-3-1-1 1
-1 1-1 1 3 3-3-1 1-1 1-3-1 1 x .-5-1-1-1 1 1-1-1 1
1-1 1-1 x-3 3 1-1 x-1 3 1-1 3-5 . 1 1 1-1-1 1 1-1
-3-1-3-1-3 1-1 5-1 1-1 x 1 3-1-1 1 . 1 1-1-1-3 1-1
x-1 1-1 1 x-1 1 3-3-1-1 1 3-1-1 1 1 .-3 3-1 1-3-1
1-1-3-1-3 1 x 1-1 5-1 3 1-1 3-1 1 1-3 . x-1-3 1-1
3-3-1-3 3 x 1-1 5-1 1-3 3 1 1 1-1-1 3 x .-3-1-5 1
-1 5 3 5-1-1-3 x-3-1-3 1-5 1-3 1-1-1-1-3 . 3 x-3
1 3 5 x 1-3-1-3-1-3-1-1-3-1-1 1-3 1-3-1 3 . 1-1
x 3 1 3-3 1-1 1 x 1-1 3-3-1-1-1 1 1-3 1-5 x 1 .-1
-1-3-1-3 3 3 1-1 1-1 5-3 3-3 1 1-1-1-1-1 1-3-1-1 .
    
```



- randomly delete 26 weights (x)
- (a) five memories
- (b-f) initial states: some are restored, some converge to other states

Desired memories:

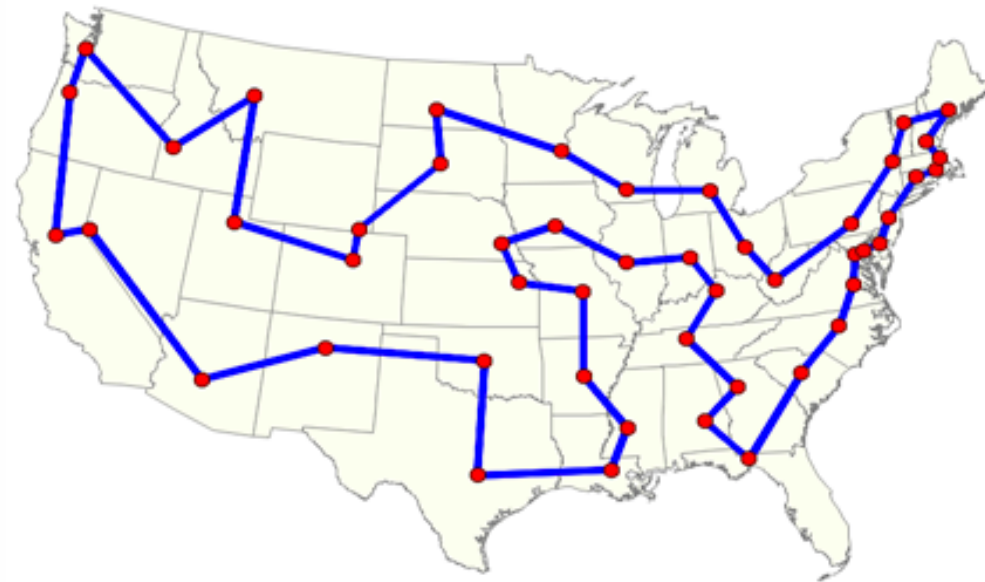


- Overloaded network fails catastrophically

Continuous Hopfield Network

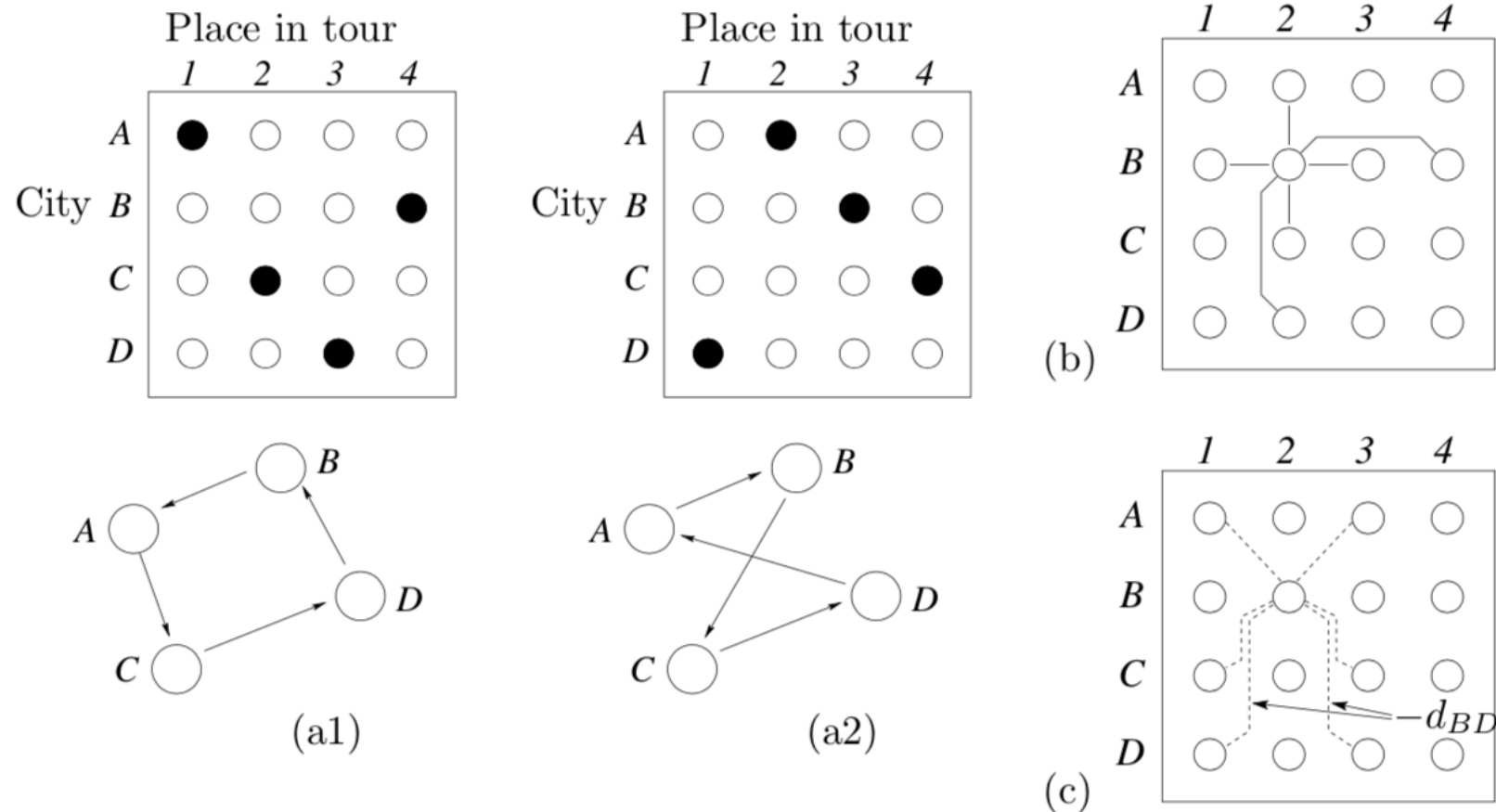
- Each neuron's activity x_i is a continuous function of time, activations $a_i(t) = \sum_j w_{ij}x_j(t)$
- Neuron's response mediated by the differential equation: $\frac{d}{dt}x_i(t) = -\frac{1}{\tau}(x_i(t) - f(a_i))$ where $f(a)$ is the activation function
- The weight matrix is symmetric \longrightarrow variational free energy is the Lyapunov function

Solving Optimization Problems



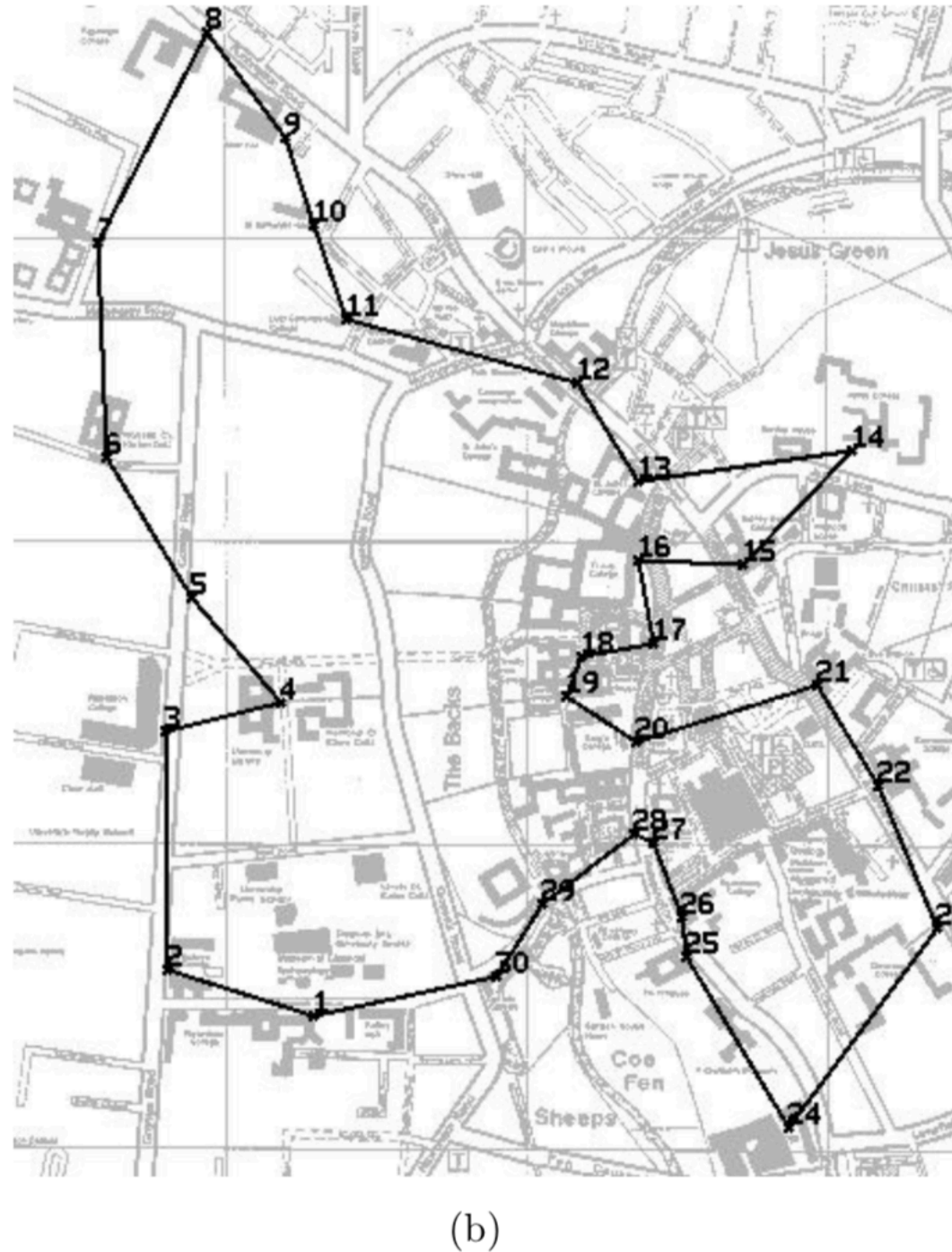
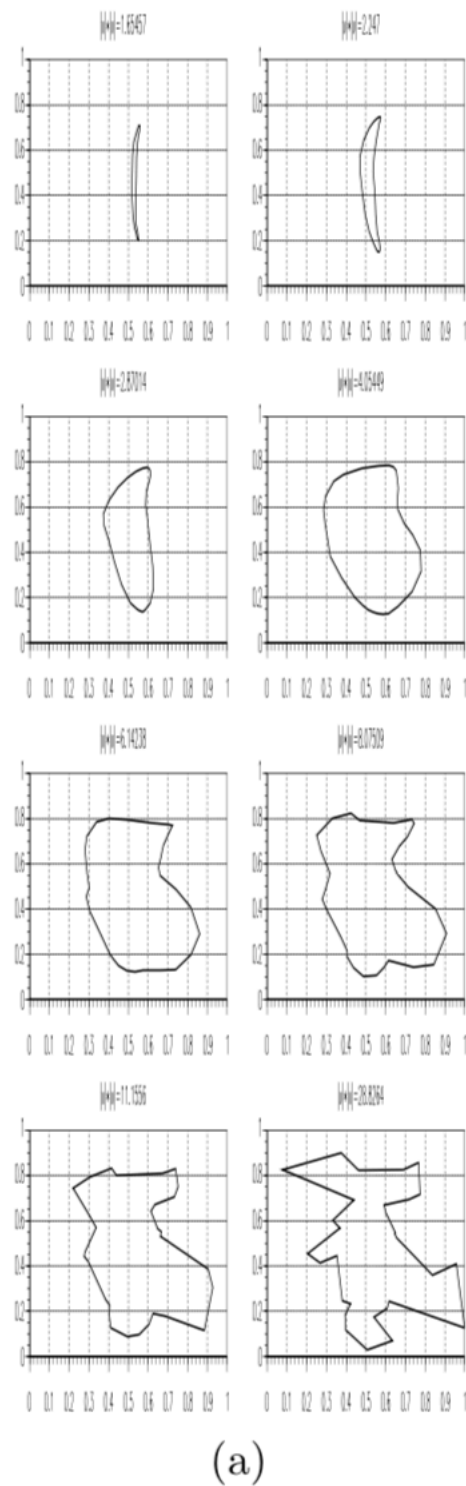
- Map interesting optimization problems onto Hopfield networks: *constrained satisfaction*
- e.g. Solving a traveling salesman problem with K cities
 - A set of K cities, a matrix of $K(K-1)/2$ distances between cities
 - Find a closed tour of the cities, visiting each city once, has the smallest total distance

TSP on Hopfield Network



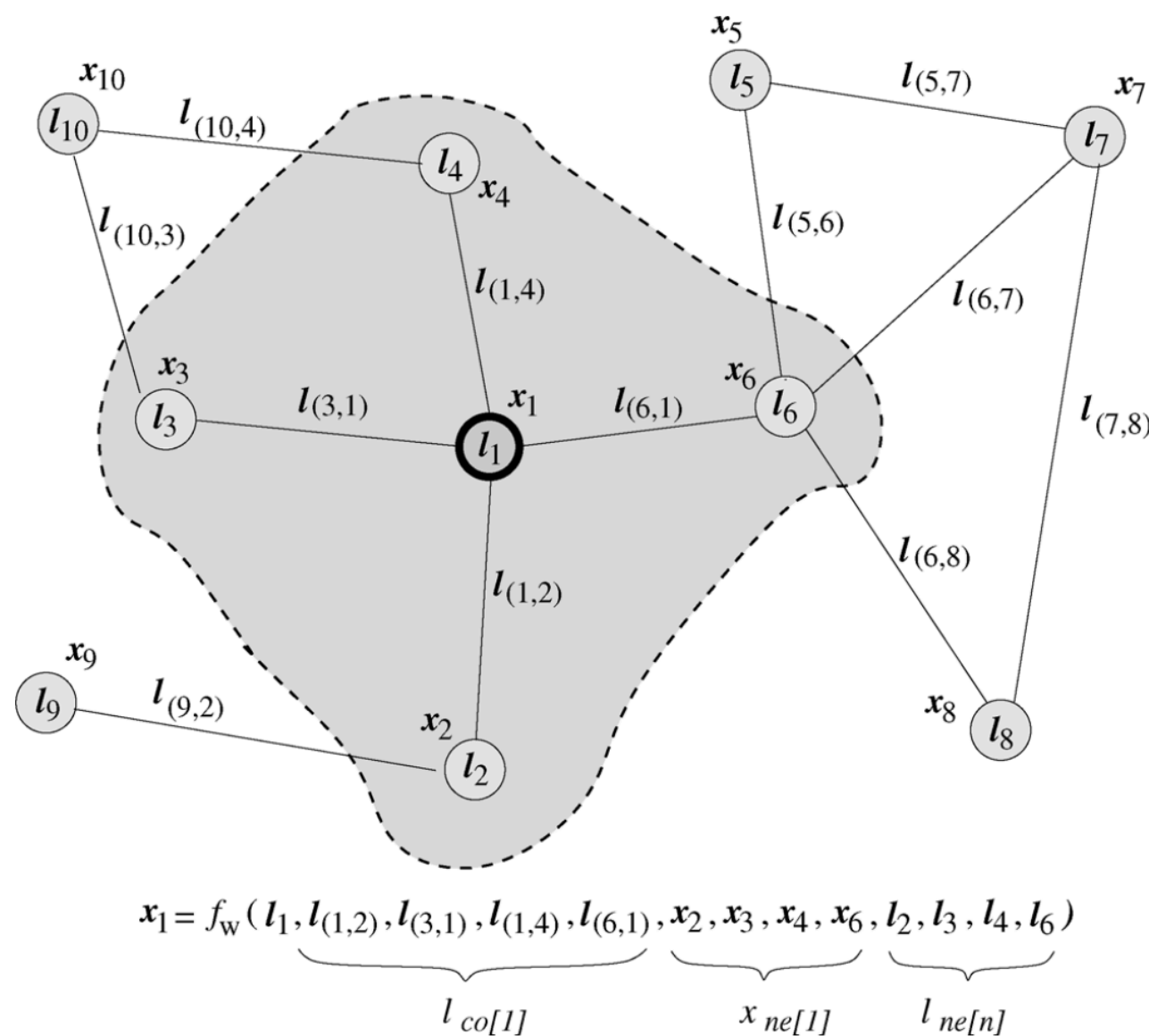
- **States have exactly one `1` in every row/column:** putting large negative weights between any pair of neurons in the same row/column
- **Weights encode the total distance:** putting negative weights proportional to the distance between the nodes

TSP on Hopfield Network



- Evolution of continuous Hopfield network
- Shortest path linking 27 Cambridge colleges

Graph Neural Networks



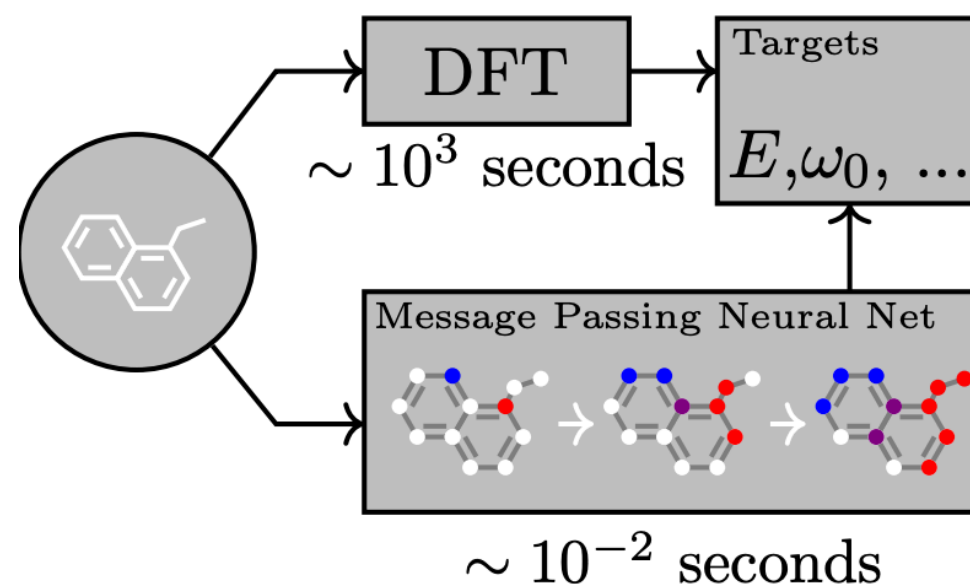
- State of a node depends on its neighbors

$$x_i = \sum_{j \in \mathcal{N}(i)} f(l_i, l_{i,j}, x_j, l_j)$$

where l is the label for nodes/edges

- Optimize by minimizing the energy

Neural Message Passing



- A graph G with node features x_v and edge features e_{vw}
- Message passing for T steps

$$m_v^{t+1} = \sum_{w \in \mathcal{N}(v)} M_t(h_v^t, h_w^t, e_{vw}) \quad h_v^{t+1} = U_t(h_v^t, m_v^{t+1}) \quad \hat{y} = R(\{h_v\})$$

Message function M_t Update function U_t and Readout function R are all neural networks