



UPPSALA
UNIVERSITET

Memristive Boltzmann Machine: A Hardware Accelerator for Combinatorial Optimization and Deep Learning*

Chris Kjellqvist

*Mahdi Nazm Bojnordi and Engin Ipek. Memristive boltzmann machine: A hardware accelerator for combinatorial optimization and deep learning. In Proceedings on the 22nd IEEE Symposium on High Performance Computer Architecture, Rochester, NY, March 2016.



UPPSALA
UNIVERSITET

The Boltzmann Machine, Problem Mapping, and Memristors

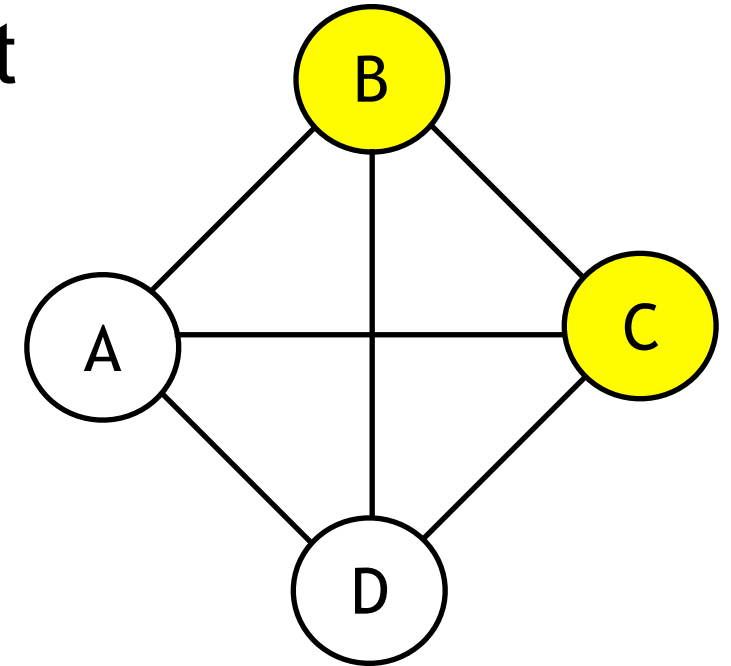
BACKGROUND



UPPSALA
UNIVERSITET

The Boltzmann Machine

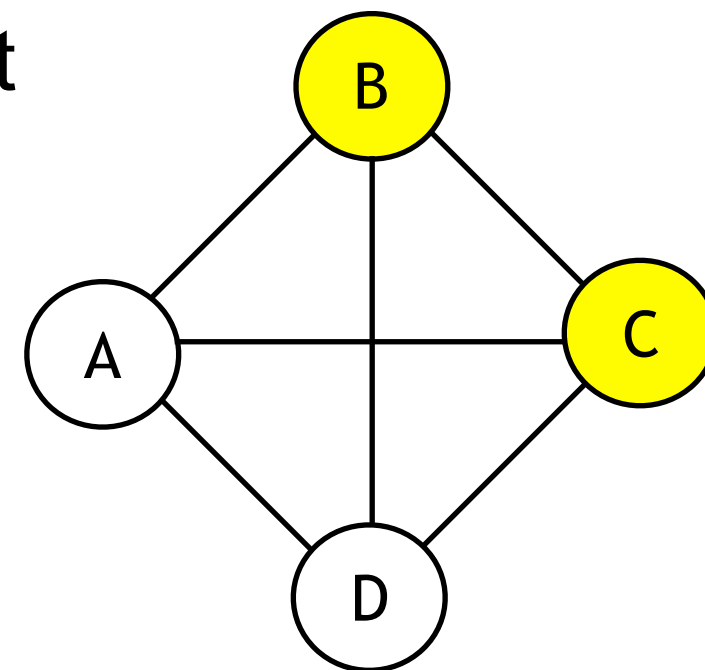
- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems





The Boltzmann Machine

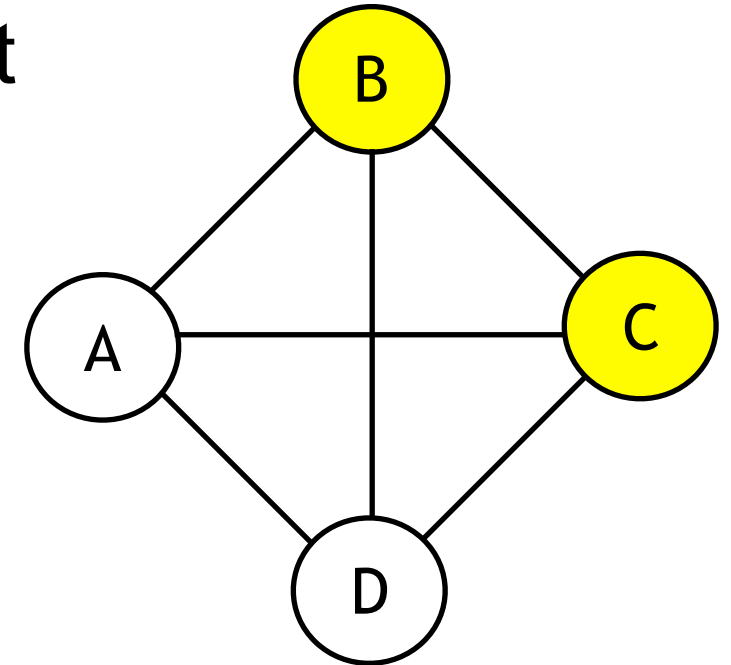
- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems
 - Traveling Salesman





The Boltzmann Machine

- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems
 - Traveling Salesman
 - Max-Cut

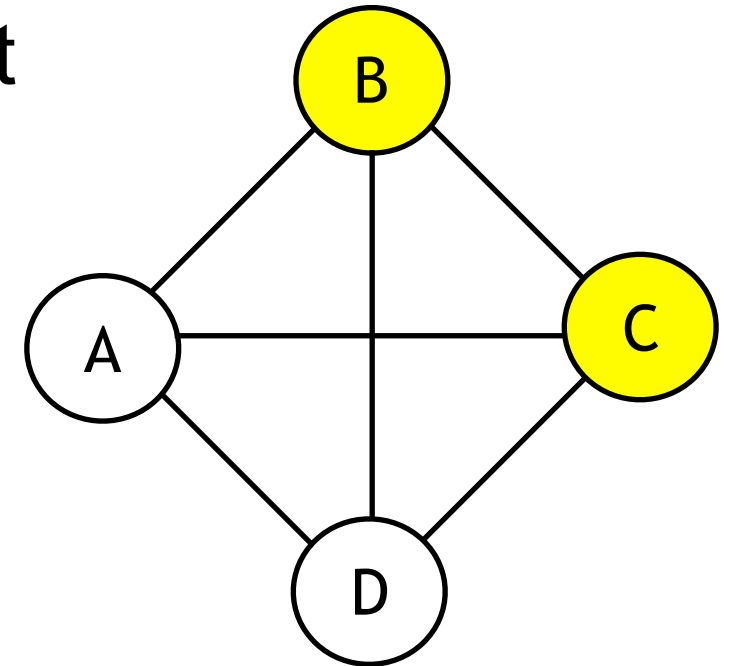




UPPSALA
UNIVERSITET

The Boltzmann Machine

- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems
 - Traveling Salesman
 - Max-Cut
 - Maximum Satisfiability

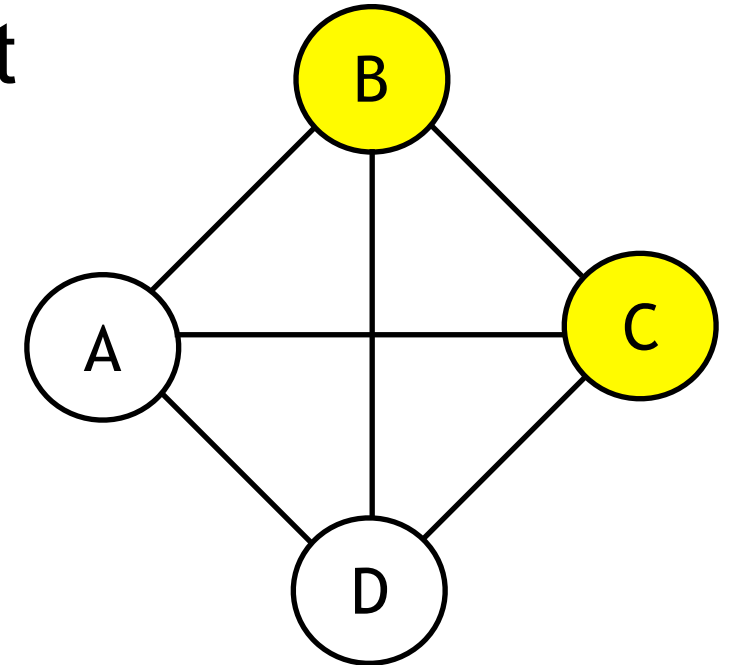




UPPSALA
UNIVERSITET

The Boltzmann Machine

- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems
 - Traveling Salesman
 - Max-Cut
 - Maximum Satisfiability
 - Deep Learning

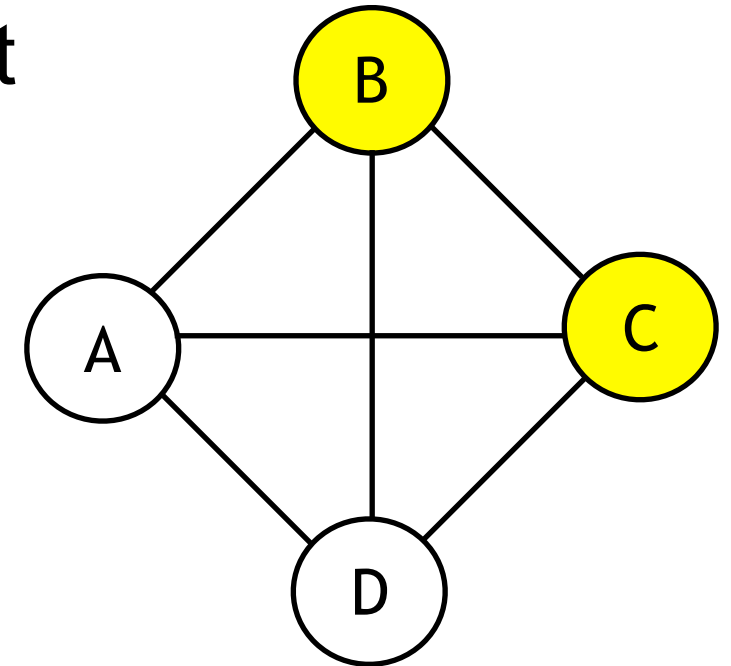


1) Ackley, David H; Hinton Geoffrey E; Sejnowski, Terrence J (1985), "A learning algorithm for Boltzmann machines", Cognitive Science, 9 (1): 147-169, doi:10.1207/s15516709cog0901_7



The Boltzmann Machine

- A Fully Connected Stochastic Recurrent Neural Network capable of solving optimization problems
- Nodes change state (on or off) as a function of their neighbor states and corresponding edge weights
- Goal to maximize “consensus” function





UPPSALA
UNIVERSITET

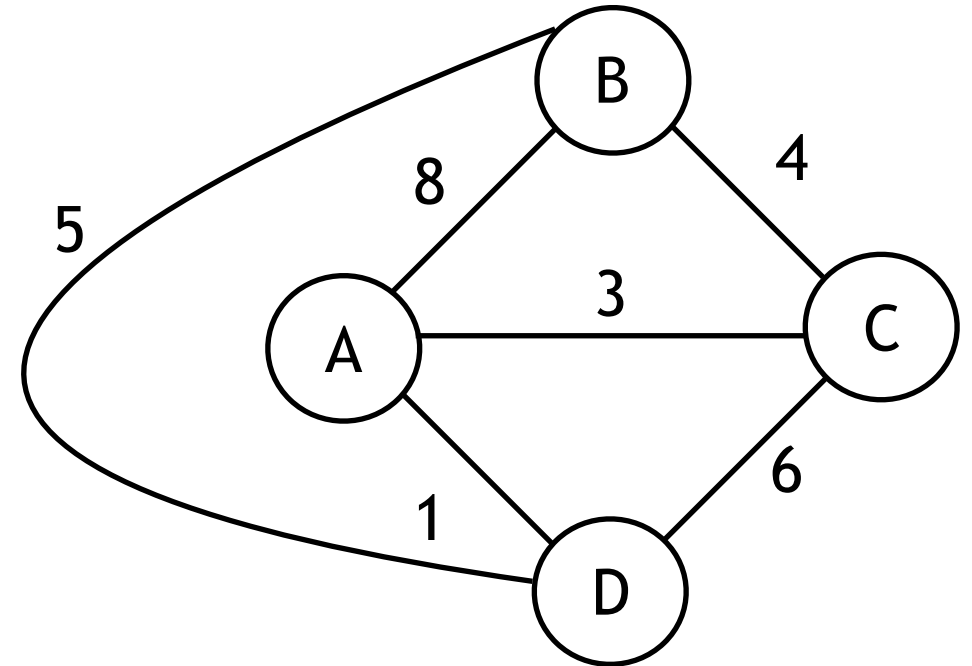
Example Problem: Max-Cut

- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized



Example Problem: Max-Cut

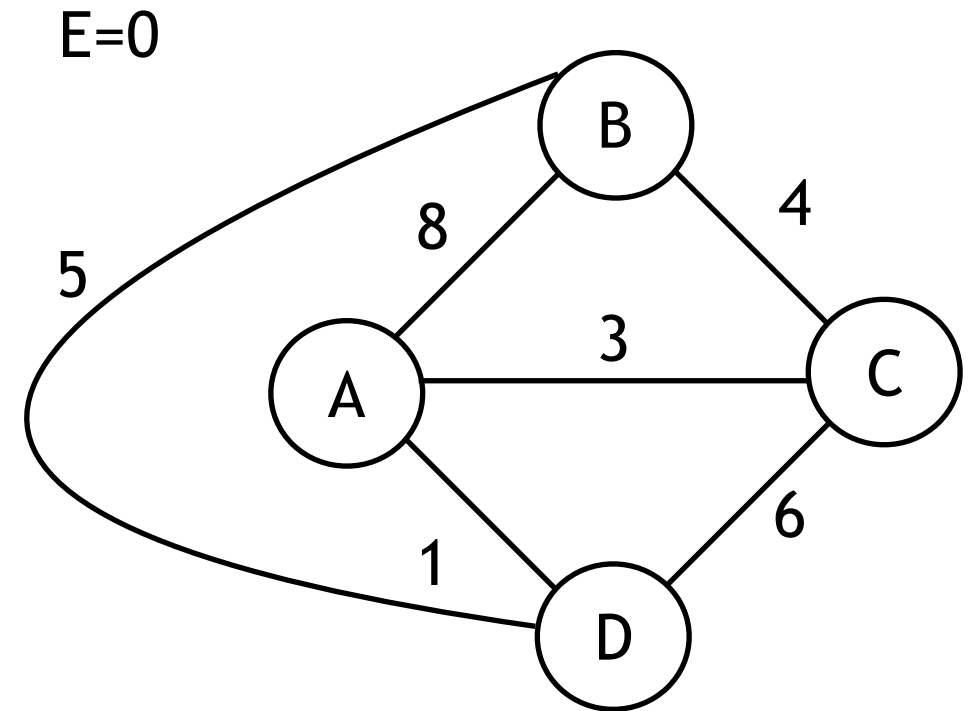
- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized





Example Problem: Max-Cut

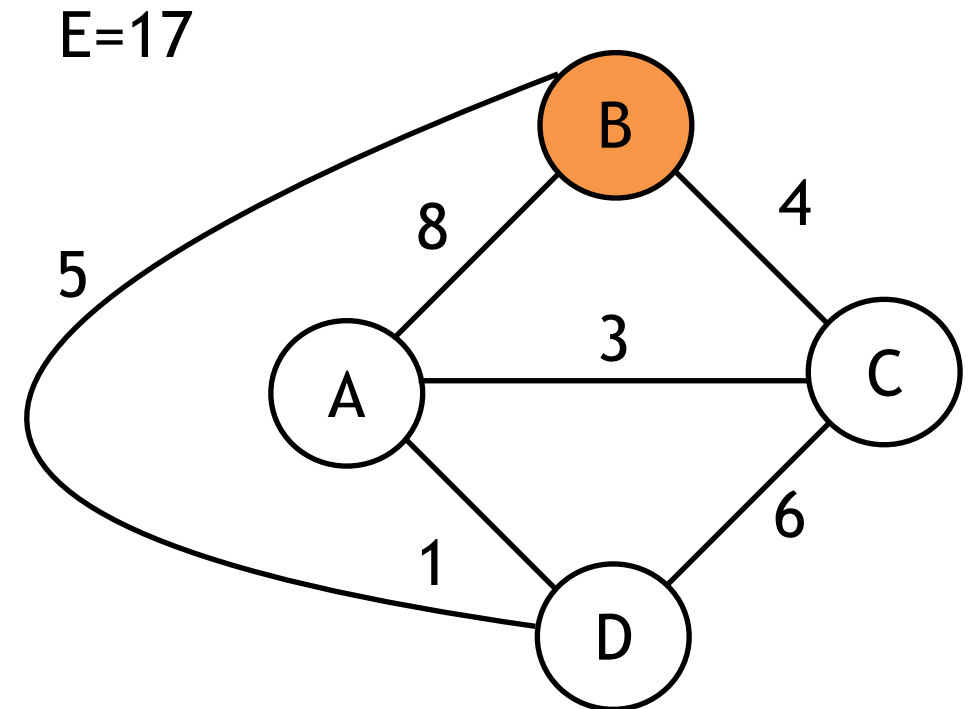
- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized





Example Problem: Max-Cut

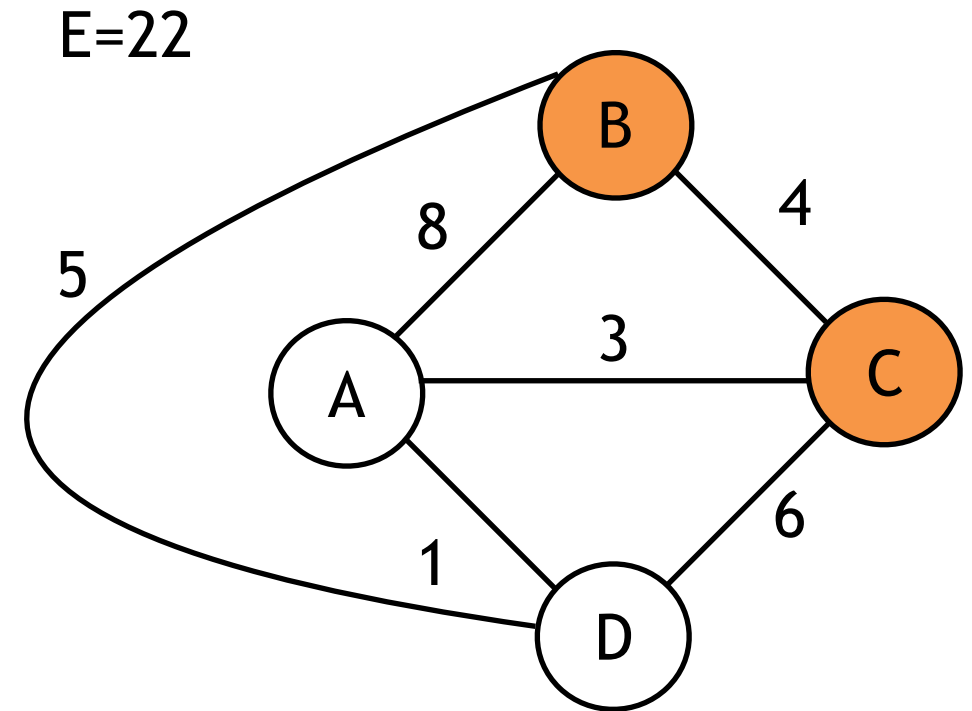
- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized





Example Problem: Max-Cut

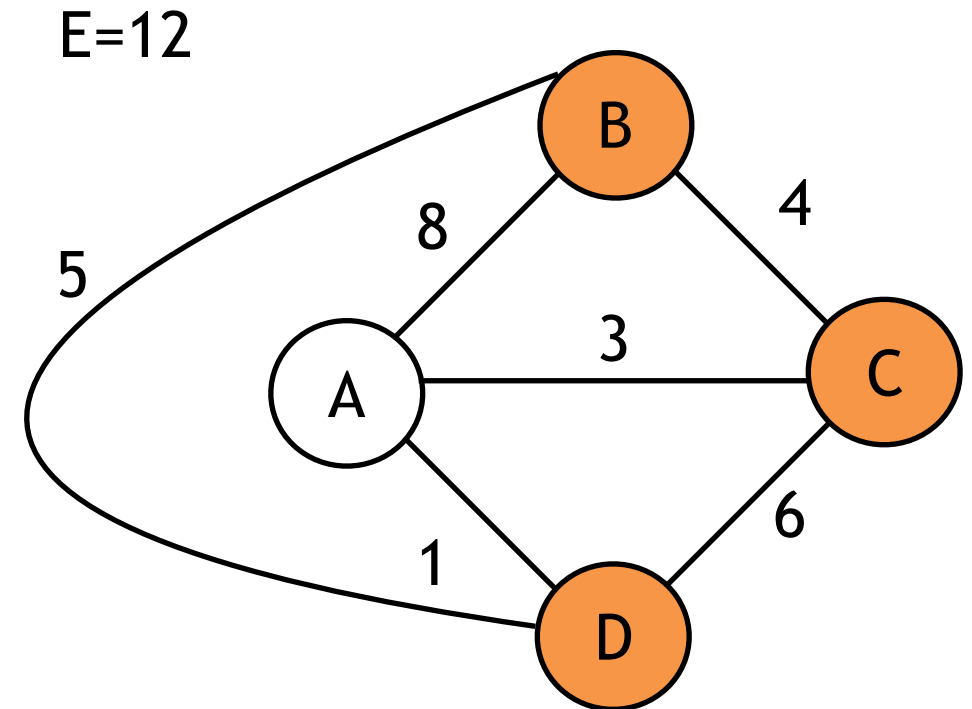
- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized





Example Problem: Max-Cut

- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized

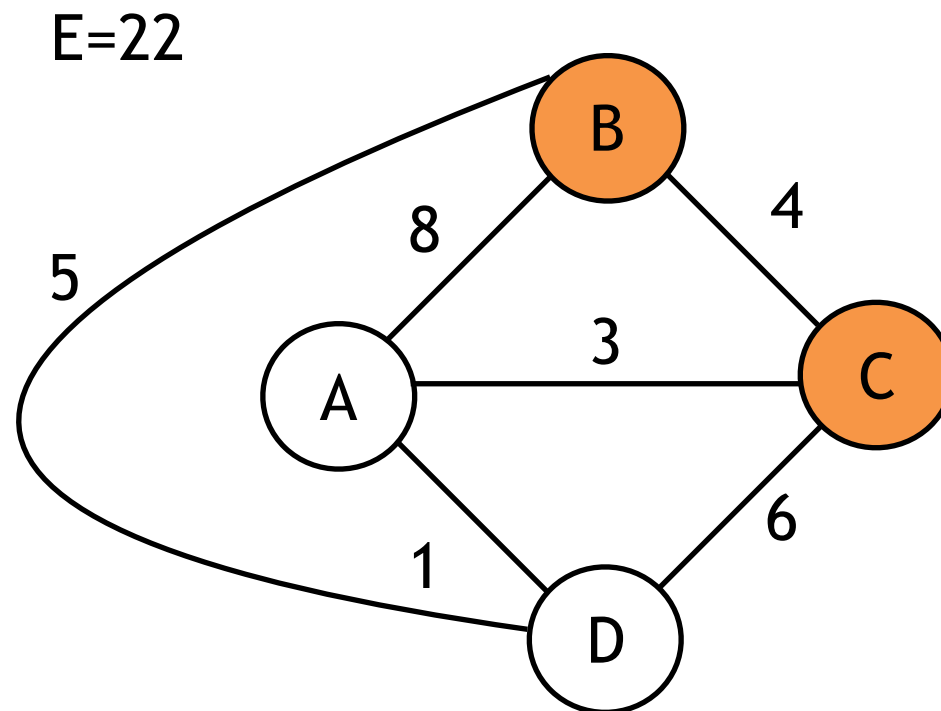




Example Problem: Max-Cut

- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized

$$S = \{B, C\}$$

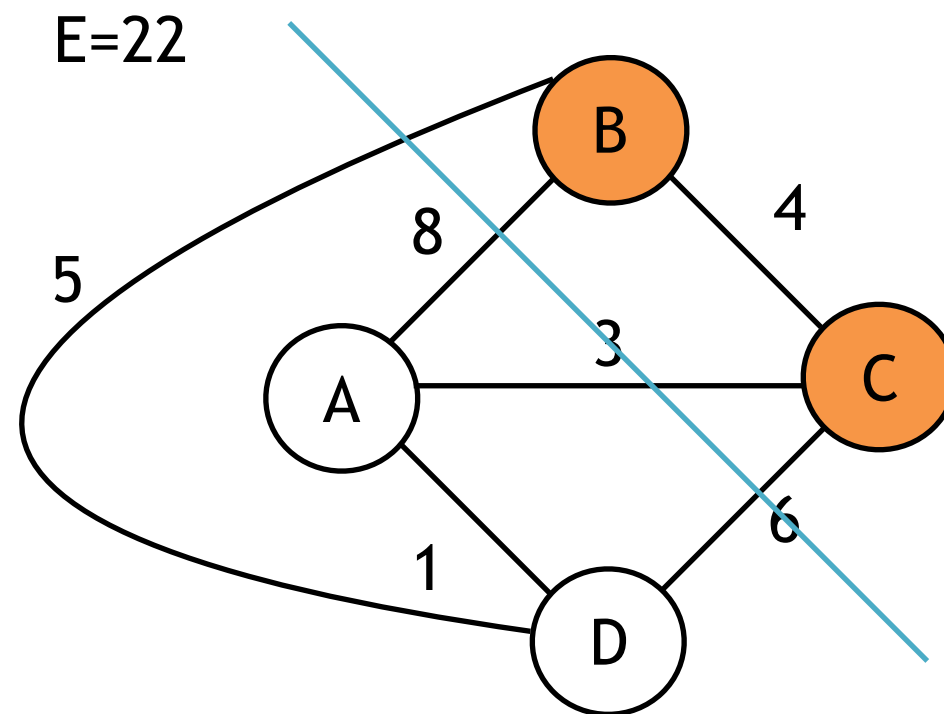




Example Problem: Max-Cut

- Given graph G , find a subset of graph G , S , such that the sum of the connection weights between the graphs are maximized

$$S = \{B, C\}$$





UPPSALA
UNIVERSITET

Mapping

- In order to map all graphs to the Boltzmann Machine, we add an incidence matrix I
- Change edge weights to allow proper training



UPPSALA
UNIVERSITET

Mapping

- In order to map all graphs to the Boltzmann Machine, we add an incidence matrix I
- Change edge weights to allow proper training

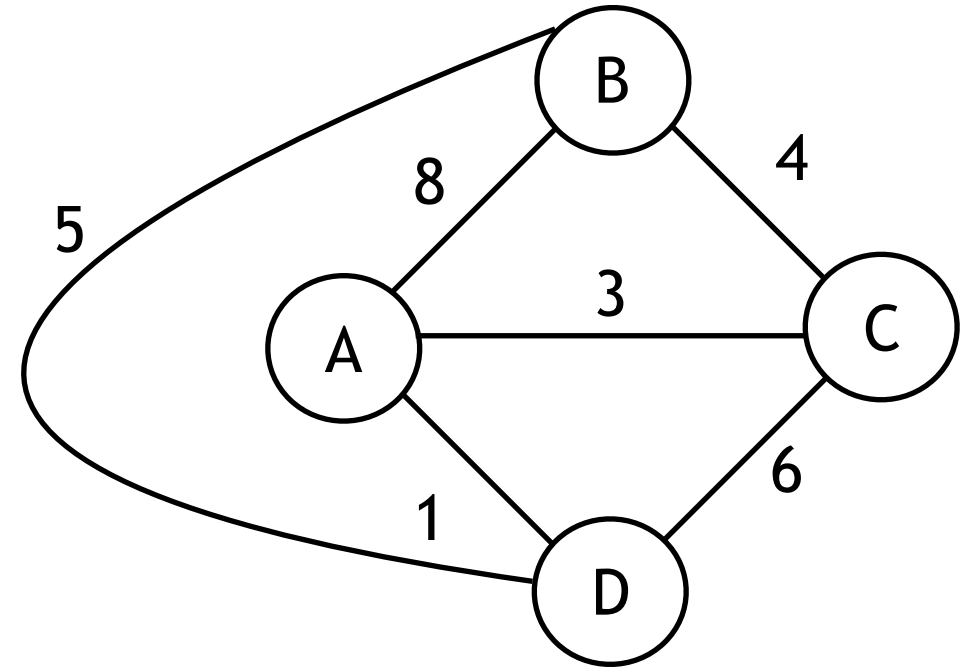
$$w_{jj} = \sum_i d_{ij} \quad w_{ji} = -2d_{ji}$$



Mapping

- In order to map all graphs to the Boltzmann Machine, we add an incidence matrix I
- Change edge weights to allow proper training

$$w_{jj} = \sum_i d_{ij} \quad w_{ji} = -2d_{ji}$$

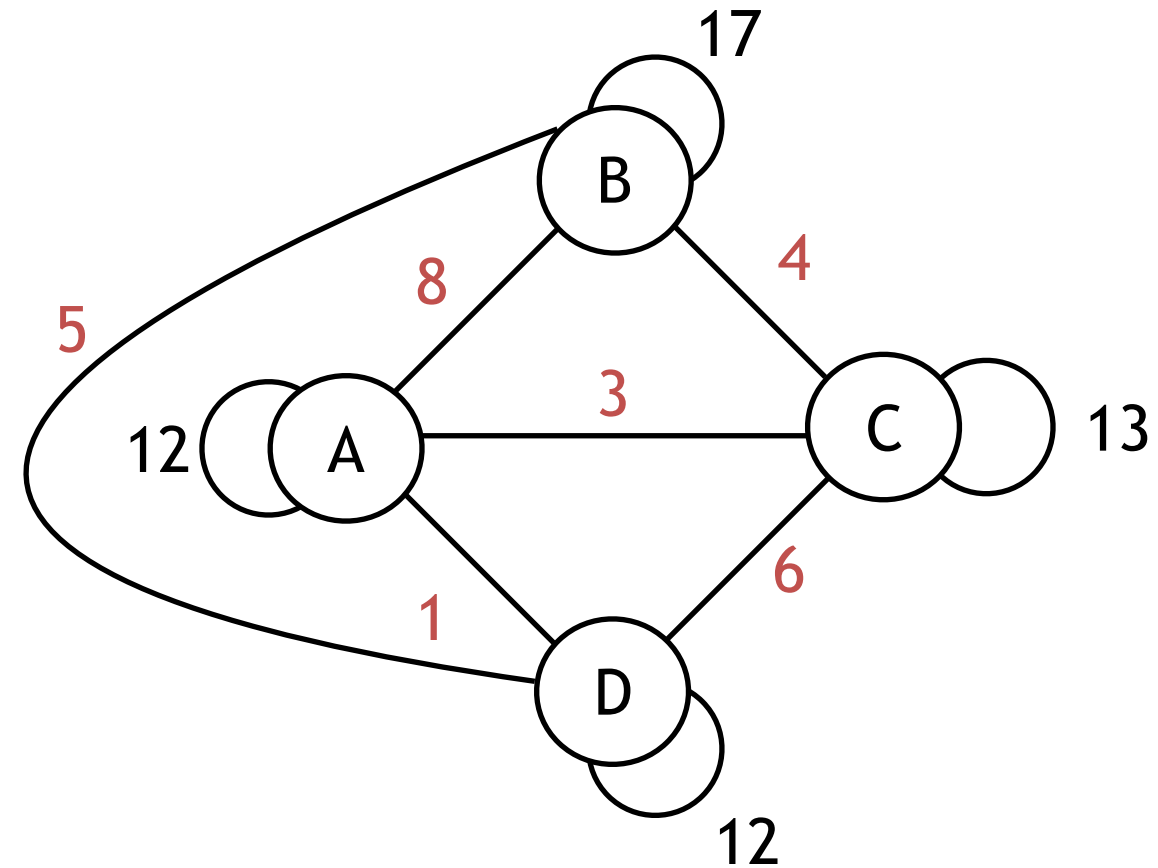




Mapping

- In order to map all graphs to the Boltzmann Machine, we add an incidence matrix I
- Change edge weights to allow proper training

$$w_{jj} = \sum_i d_{ij} \quad w_{ji} = -2d_{ji}$$



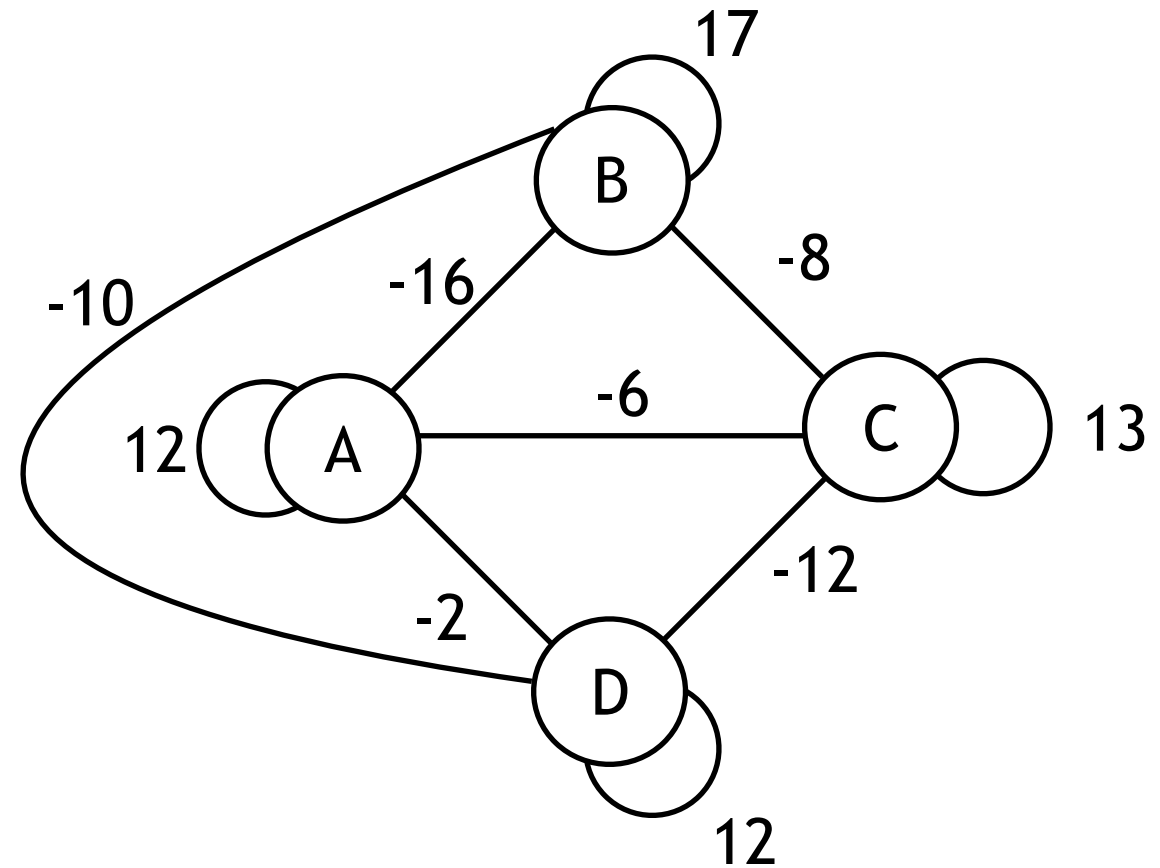


UPPSALA
UNIVERSITET

- In order to map all graphs to the Boltzmann Machine, we add an incidence matrix I
- Change edge weights to allow proper training

$$w_{jj} = \sum_i d_{ij} \quad w_{ji} = -2d_{ji}$$

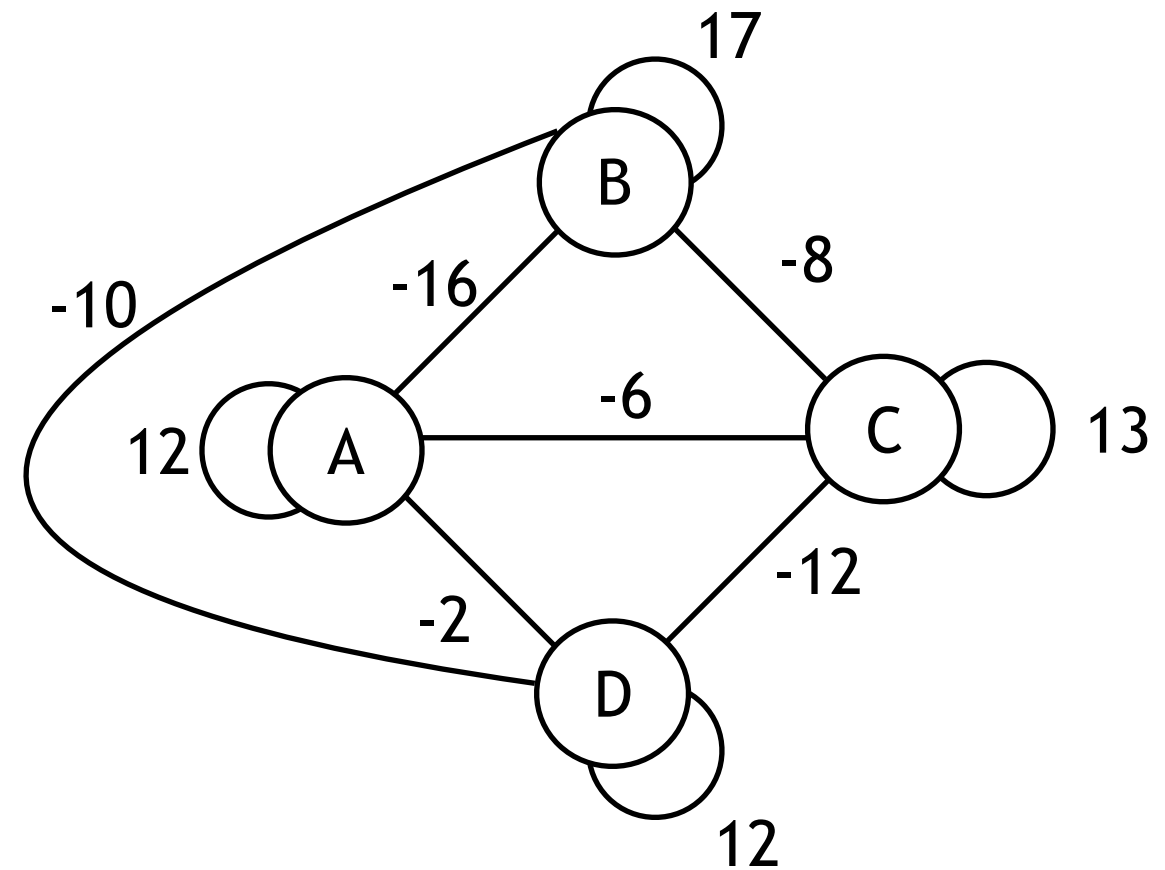
Mapping





UPPSALA
UNIVERSITET

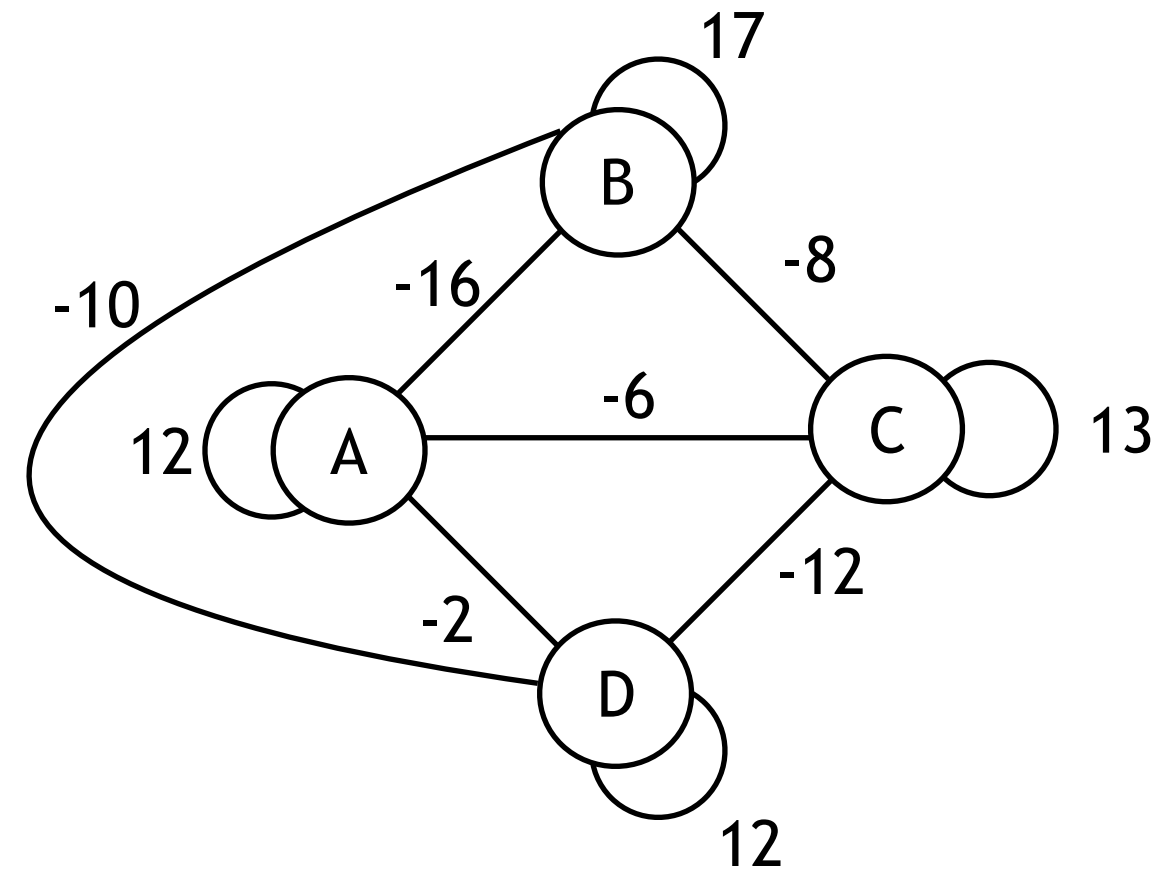
Mapping





Mapping

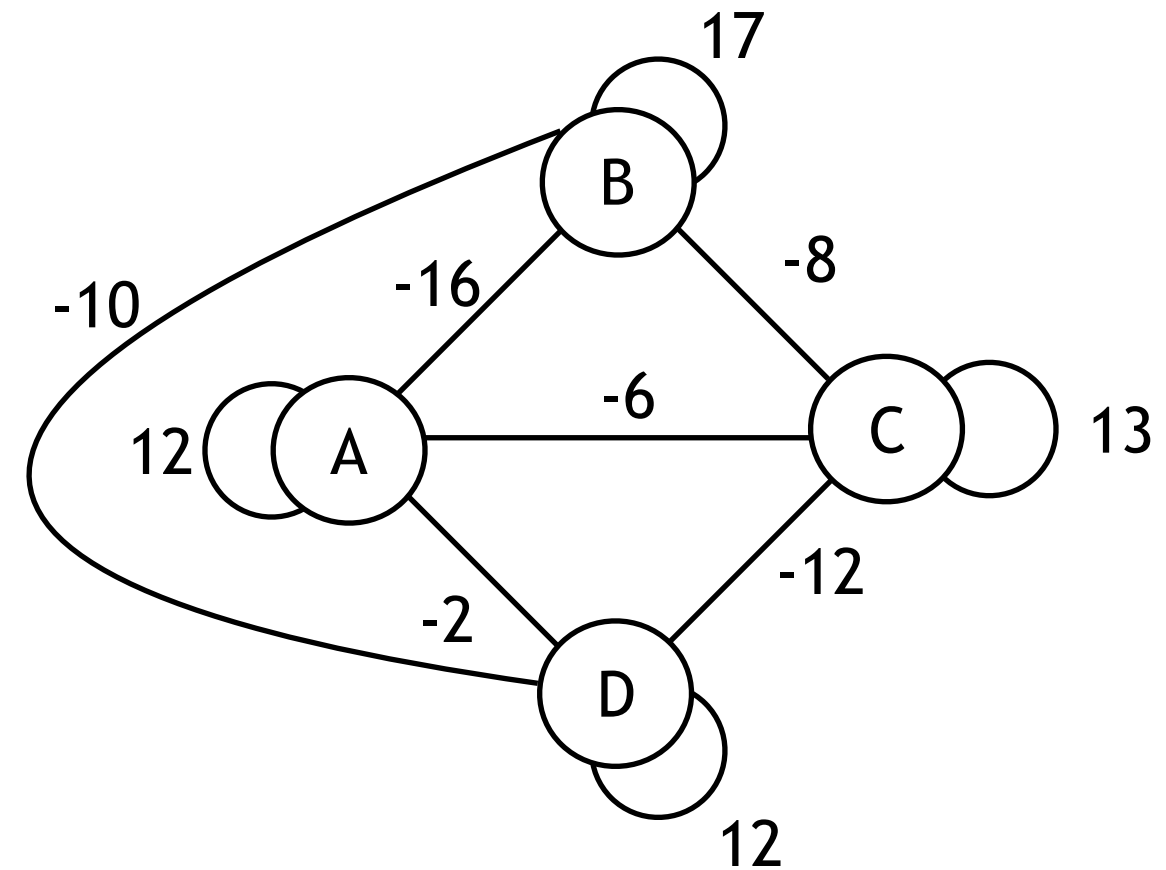
- The graph is mapped to a Boltzmann Machine but what's the task now?





Mapping

- The graph is mapped to a Boltzmann Machine but what's the task now?
- Minimize "Energy"

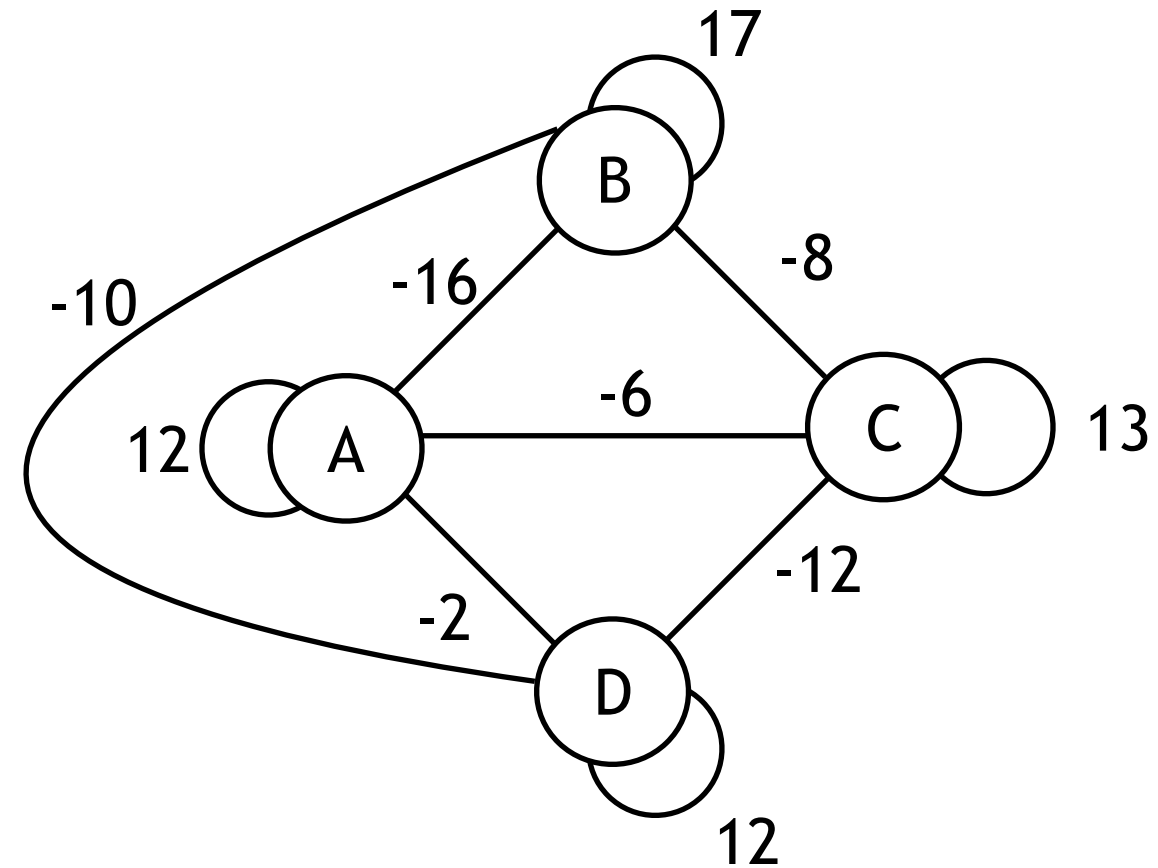




Mapping

- The graph is mapped to a Boltzmann Machine but what's the task now?
- Minimize “Energy”

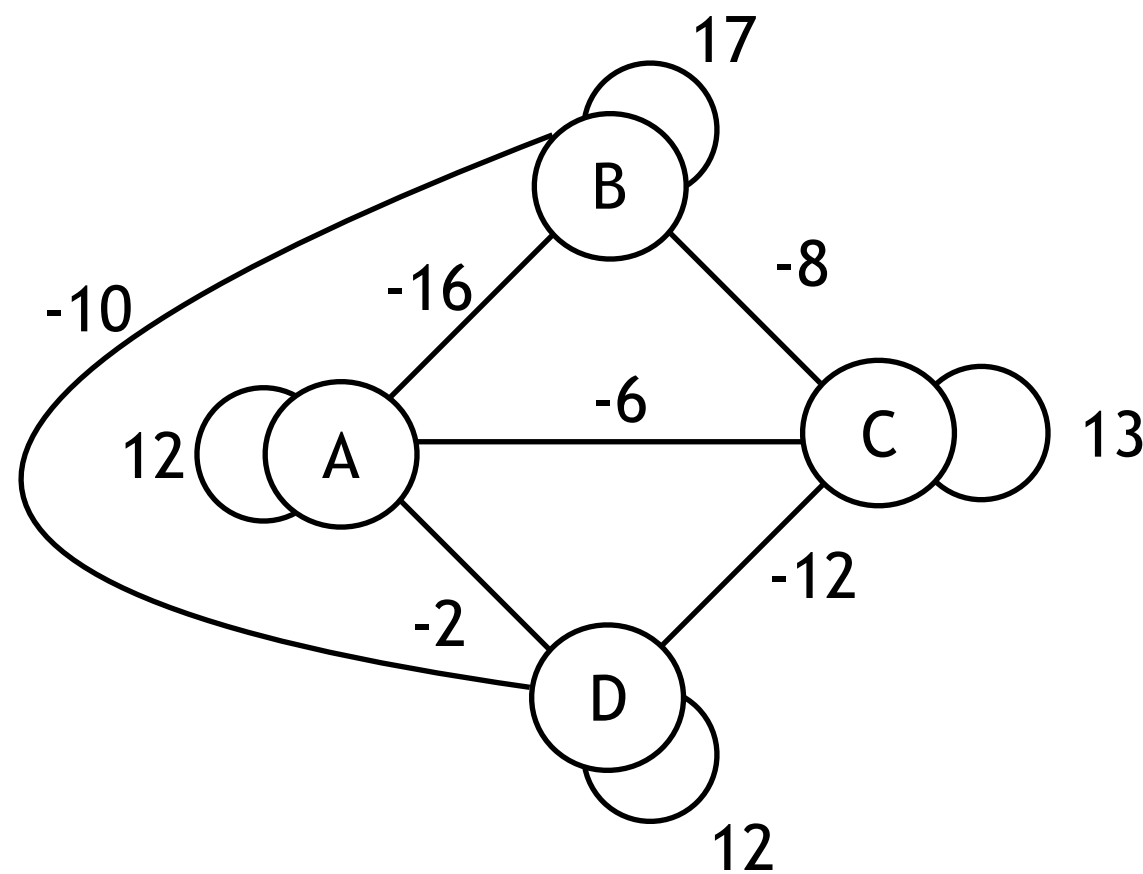
$$E(X) = -\frac{1}{2} \sum_j \sum_{i \neq j} x_j x_i w_{ij} - \sum_j x_j w_{jj}$$





Example

$$E(X) = -\frac{1}{2} \sum_j \sum_{i \neq j} x_j x_i w_{ij} - \sum_j x_j w_{jj}$$

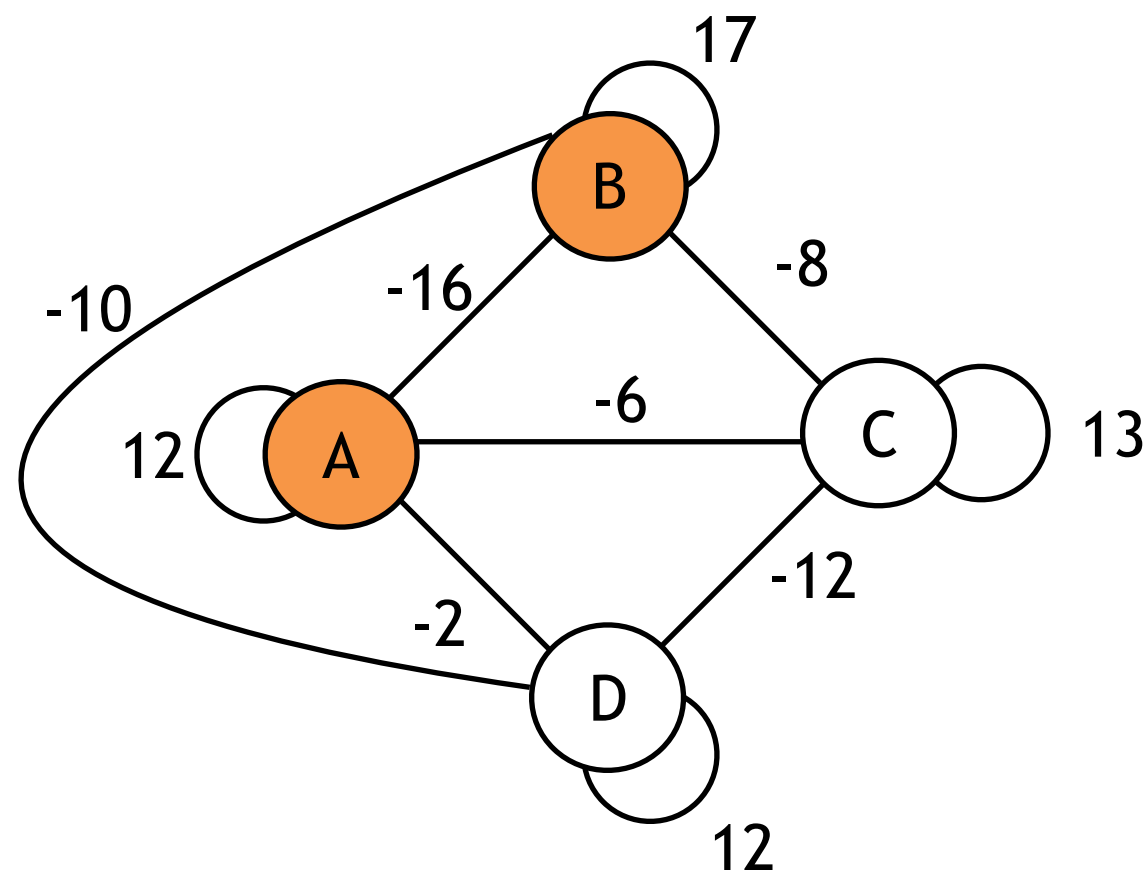




Example

$$E(X) = -\frac{1}{2} \sum_j \sum_{i \neq j} x_j x_i w_{ij} - \sum_j x_j w_{jj}$$

$$\begin{aligned} E(X) &= -\frac{1}{2} \cdot (-16 + -16) - (12 + 17) \\ &= 16 - 29 = -13 \end{aligned}$$





UPPSALA
UNIVERSITET

Simulated Annealing



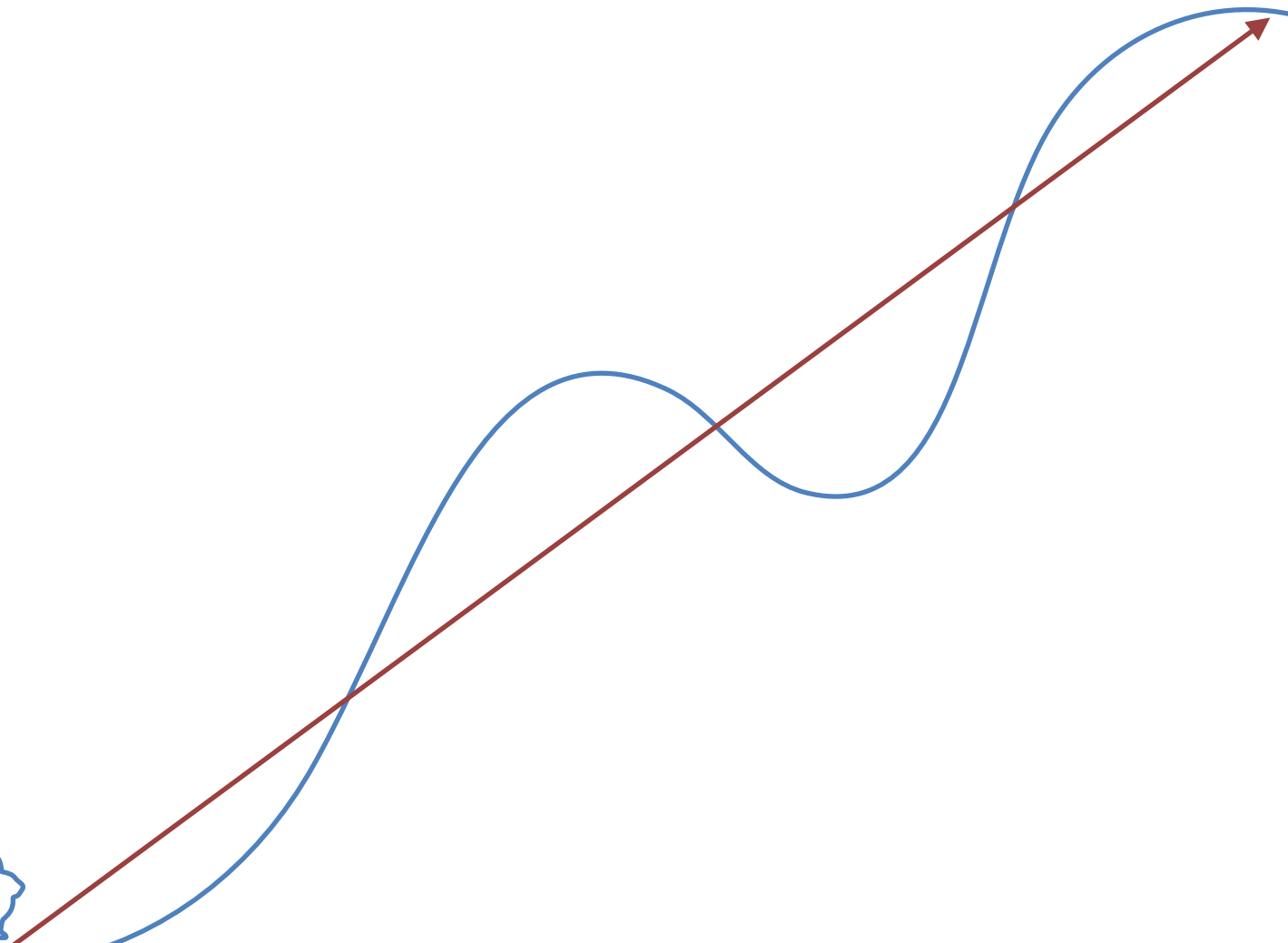


UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET



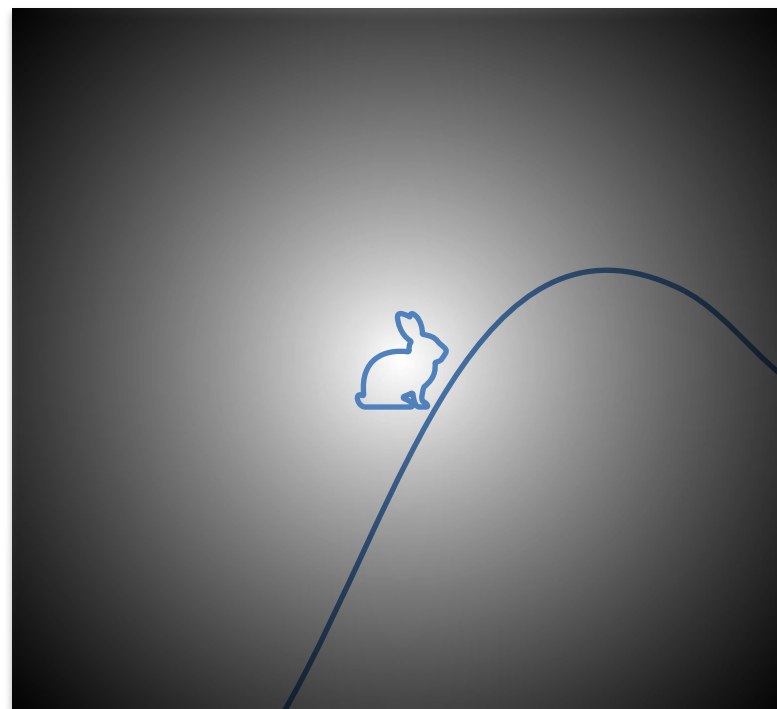


UPPSALA
UNIVERSITET



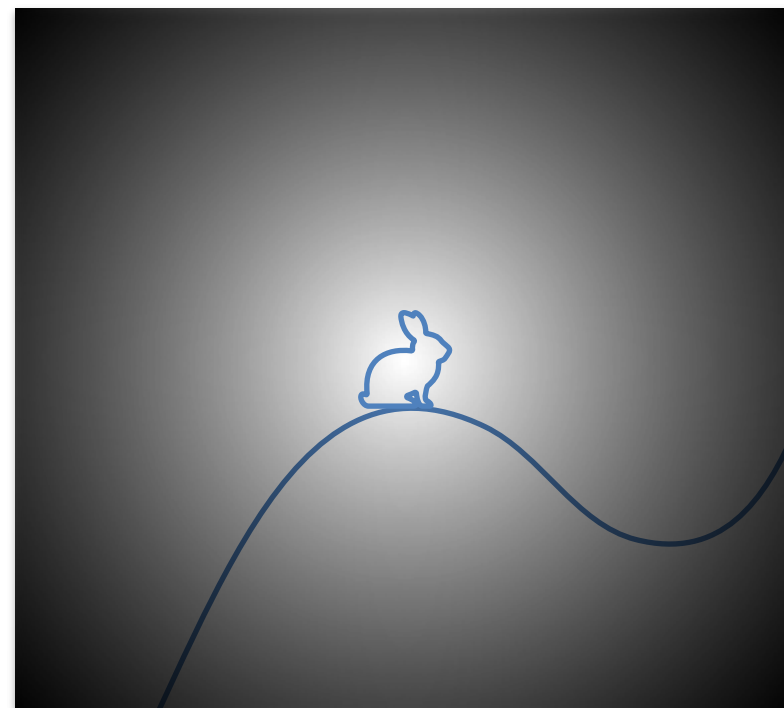


UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET

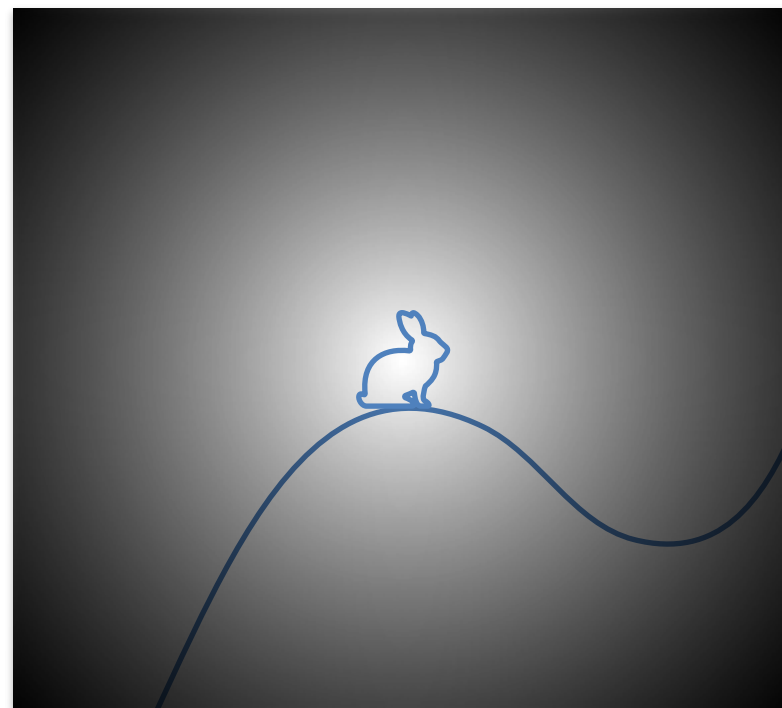


I am at the top!



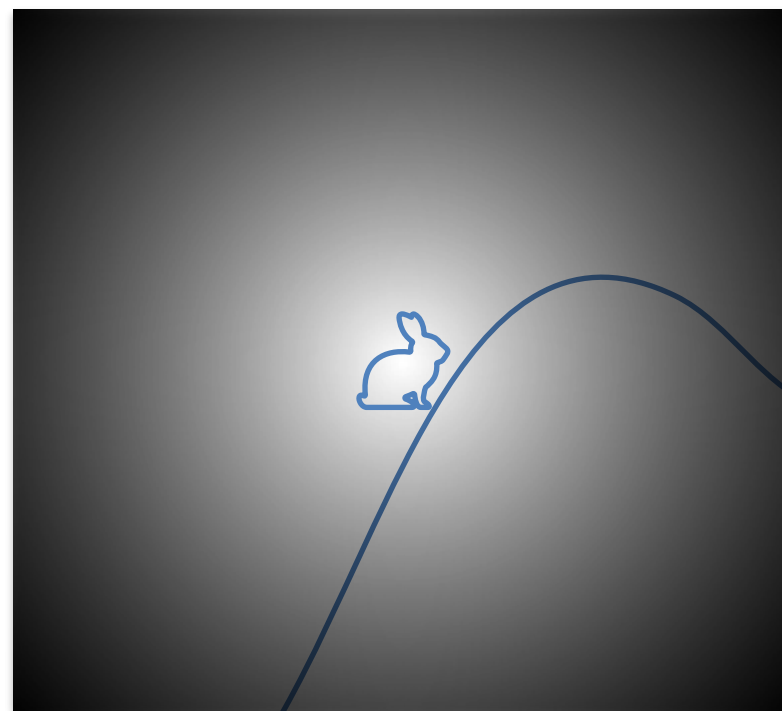
UPPSALA
UNIVERSITET

What if the rabbit is feeling risky?



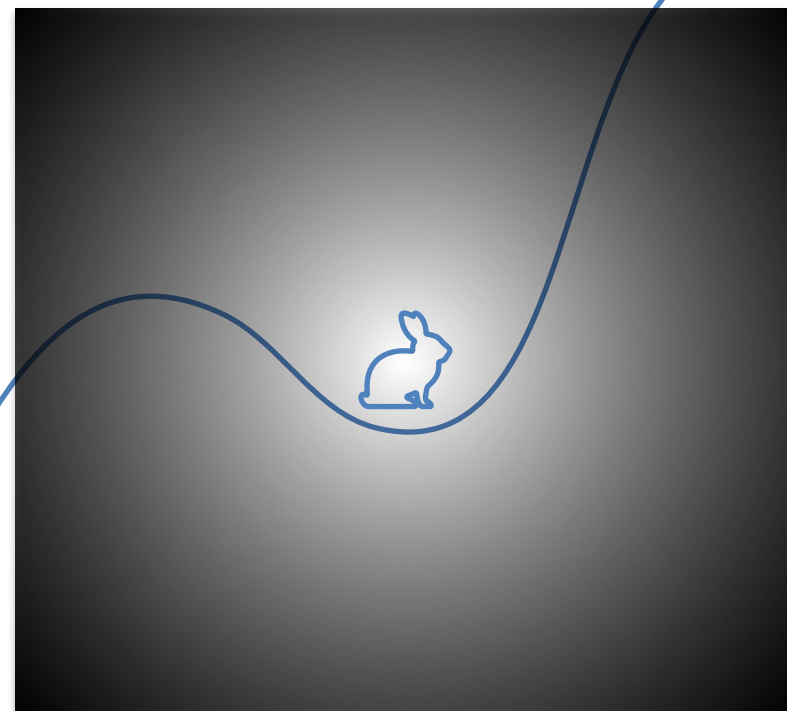


UPPSALA
UNIVERSITET



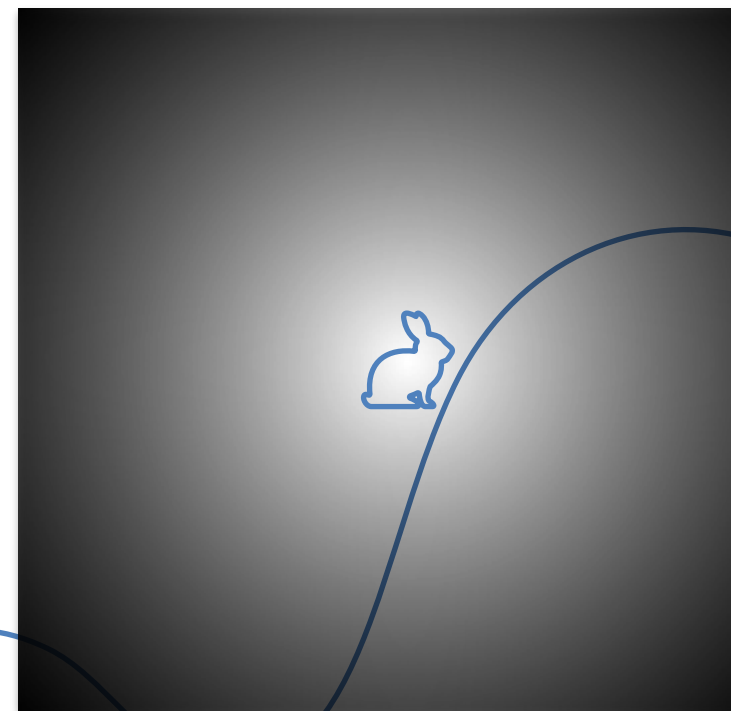


UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET





UPPSALA
UNIVERSITET

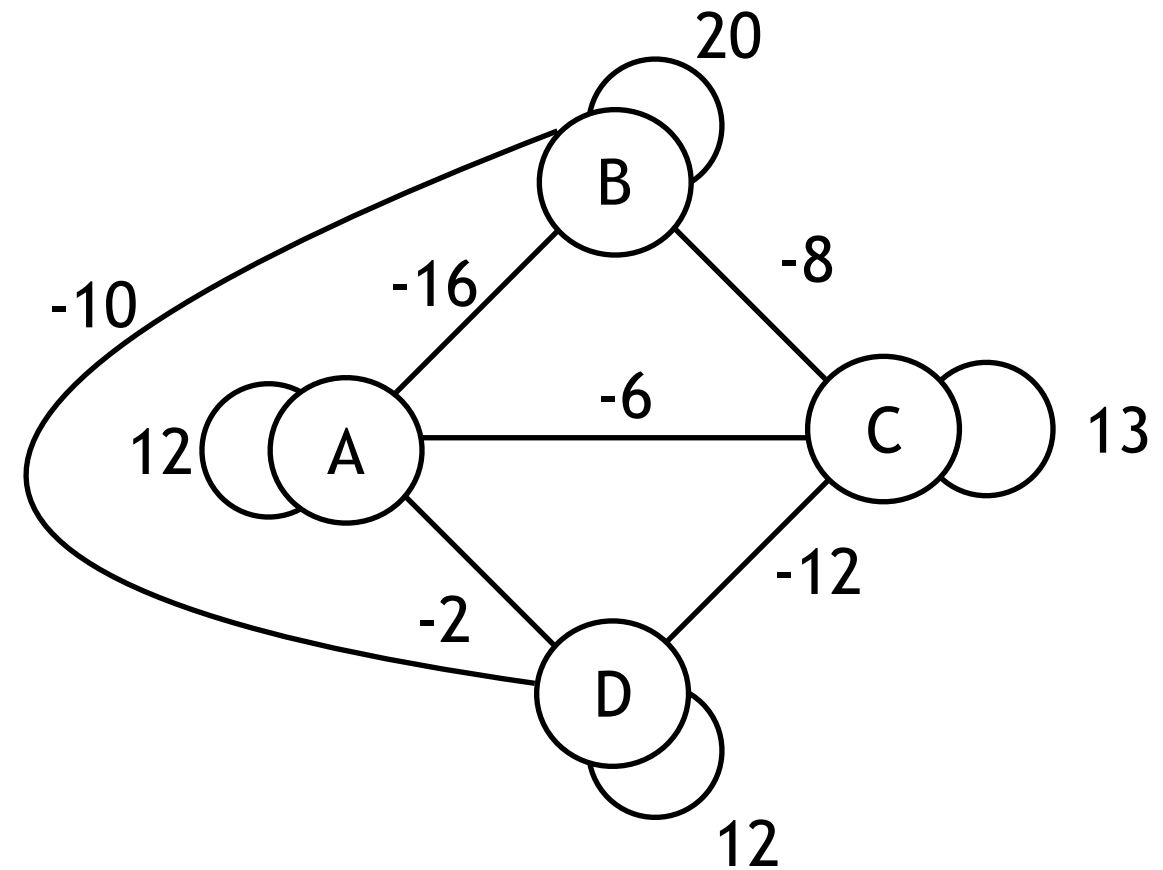




Learning the Solution

- Stochastically learn solution over time

$$\Delta E = (2x_j - 1) \left(\sum_{i \neq j} x_i w_{ji} + w_{jj} \right)$$





Learning the Solution

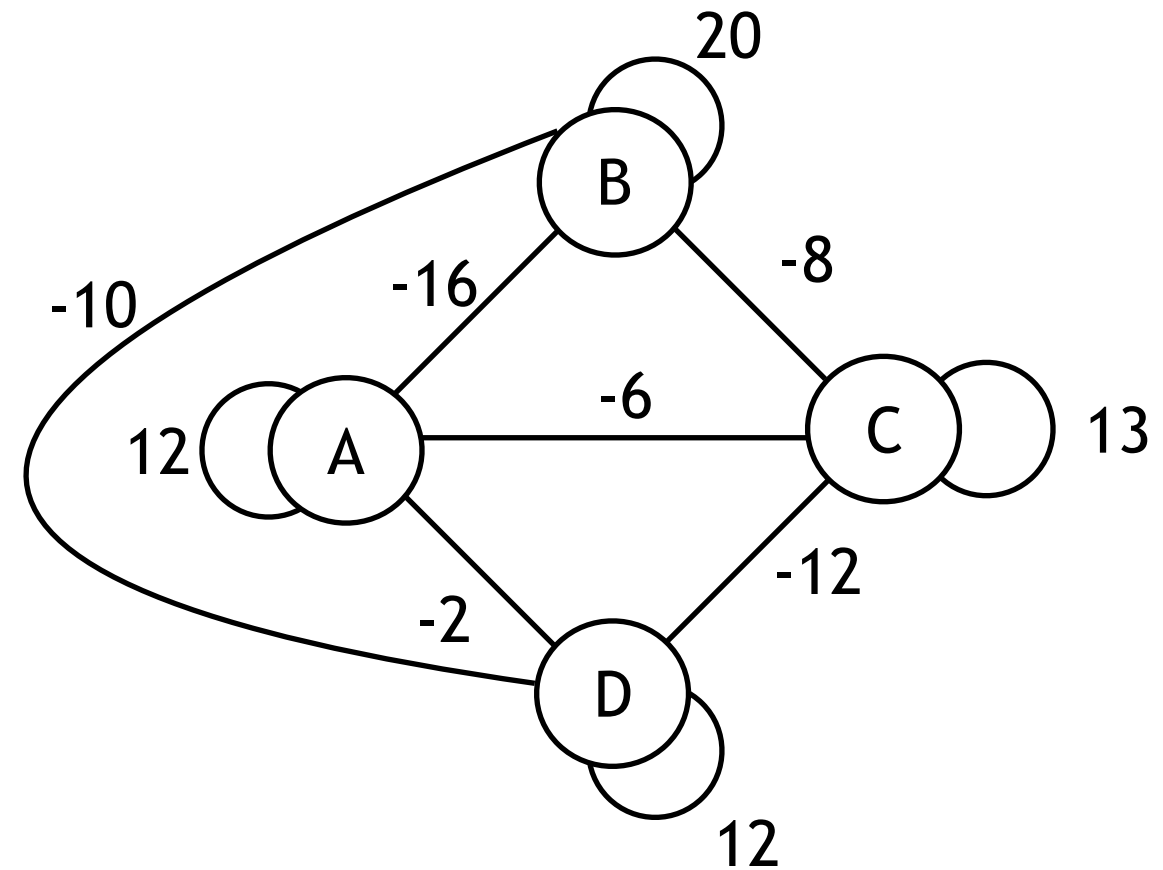
- Stochastically learn solution over time

$$\Delta E = (2x_j - 1) \left(\sum_{i \neq j} x_i w_{ji} + w_{jj} \right)$$

$$x^j = \langle x_0, x_1, \dots, \neg x_j, \dots, x_n \rangle$$

$$P(x^j | x) = \frac{1}{1 + e^{\frac{\Delta E}{c}}}$$

- C allows for probability control of suboptimal state changes

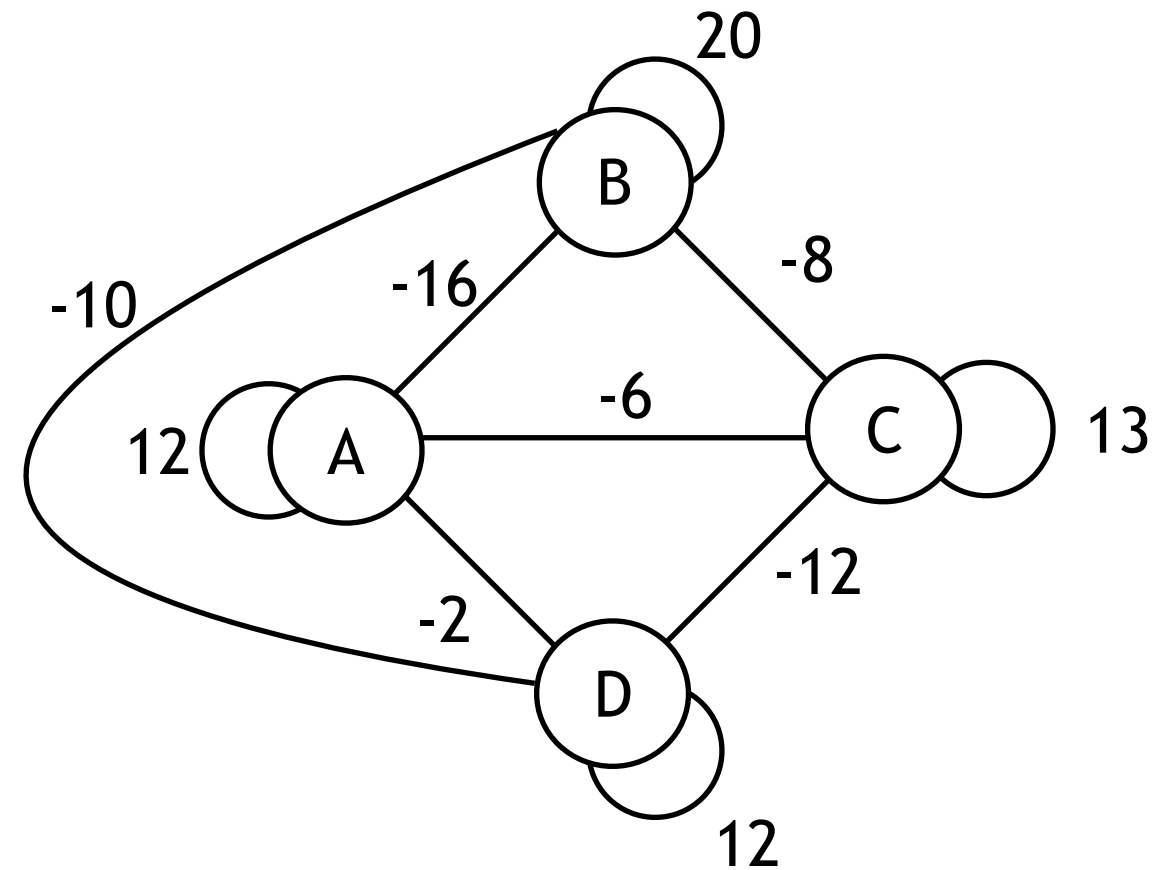




$$P(x^j | x) = \frac{1}{1 + e^{\frac{\Delta E}{c}}}$$

- The change in E can be continuously calculated *in situ* for each node
- C allows for simulated annealing, eventually *probably* resulting in a optimized graph

Learning the Solution

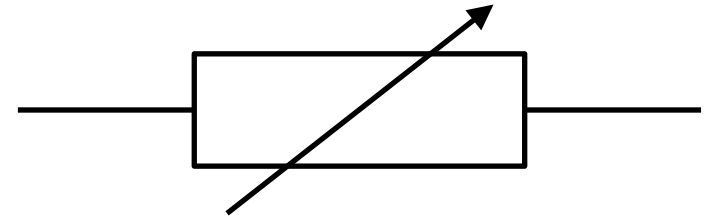




UPPSALA
UNIVERSITET

Memristors

- A component with variable conductance

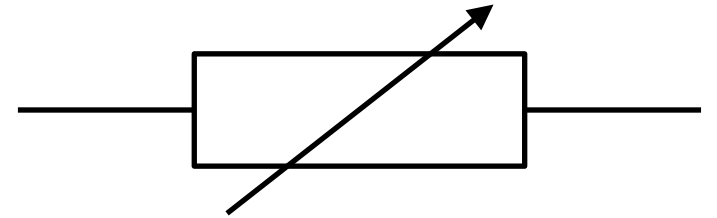




UPPSALA
UNIVERSITET

Memristors

- A component with variable conductance
- Recall $V = I \cdot R$

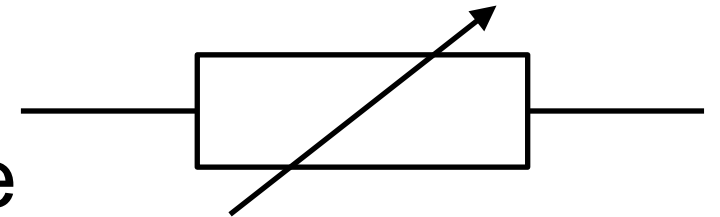




UPPSALA
UNIVERSITET

Memristors

- A component with variable conductance
- Recall $V = I \cdot R$
- Used to perform multiplications with the assistance of hardware to detect current

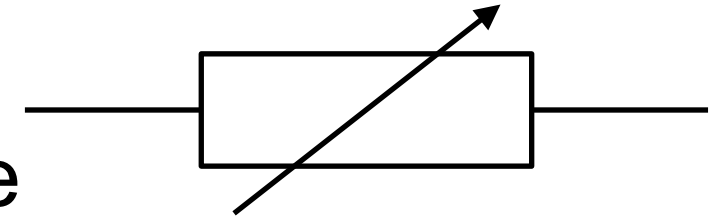




UPPSALA
UNIVERSITET

Memristors

- A component with variable conductance
- Recall $V = I \cdot R$
- Used to perform multiplications with the assistance of hardware to detect current
- Components of RRAM cells - memory with density of FLASH but DRAM speed

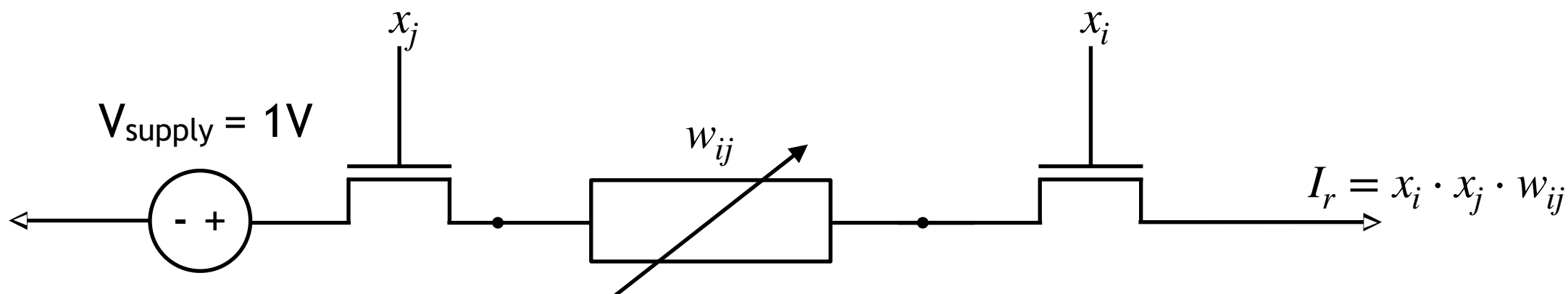




UPPSALA
UNIVERSITET

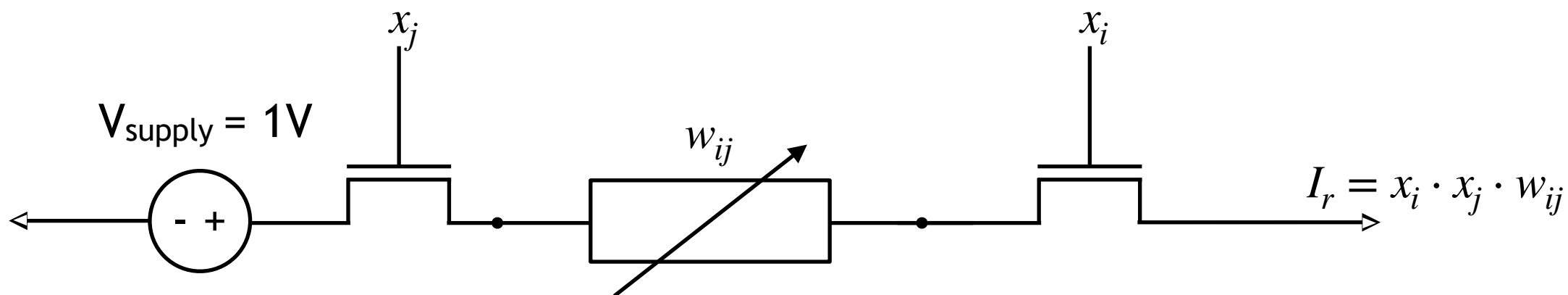
$$V = I \cdot R$$

Memristors



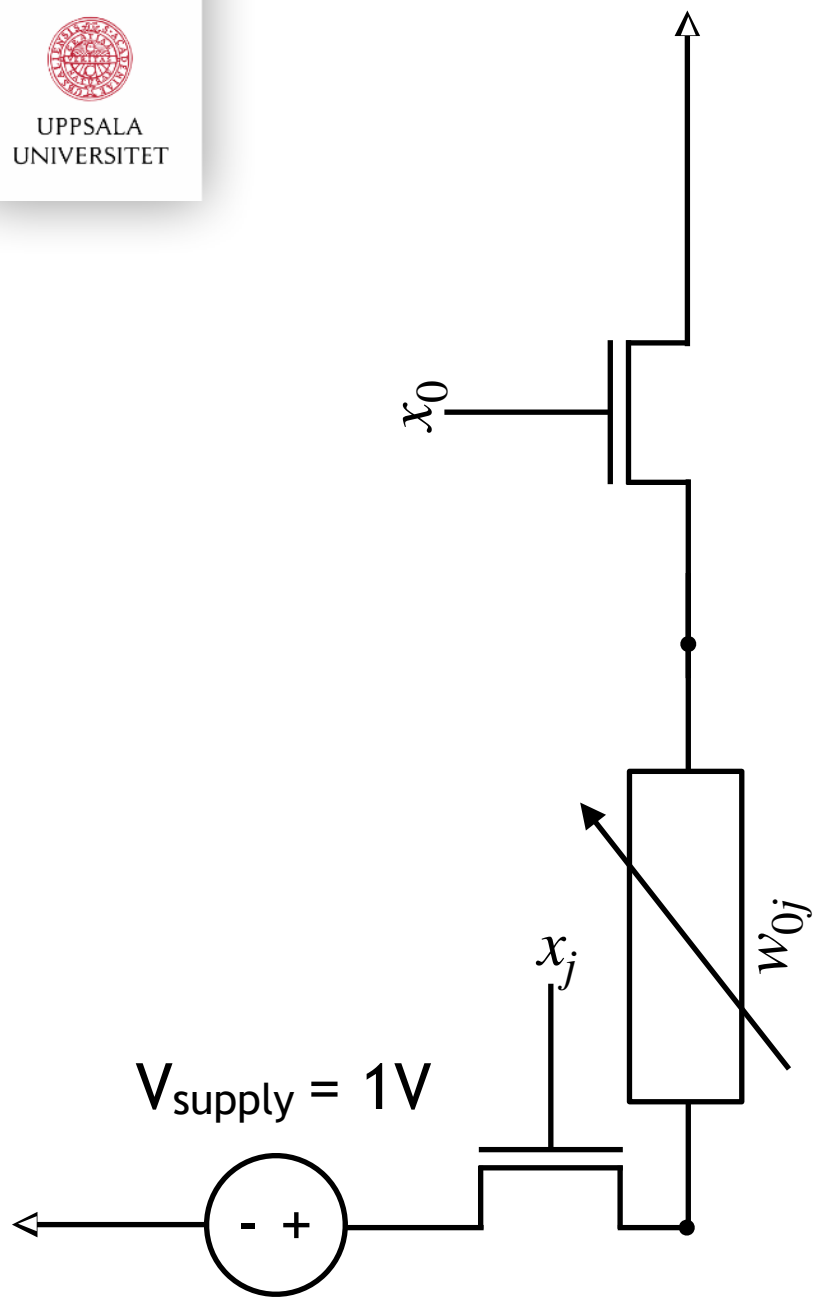


UPPSALA
UNIVERSITET



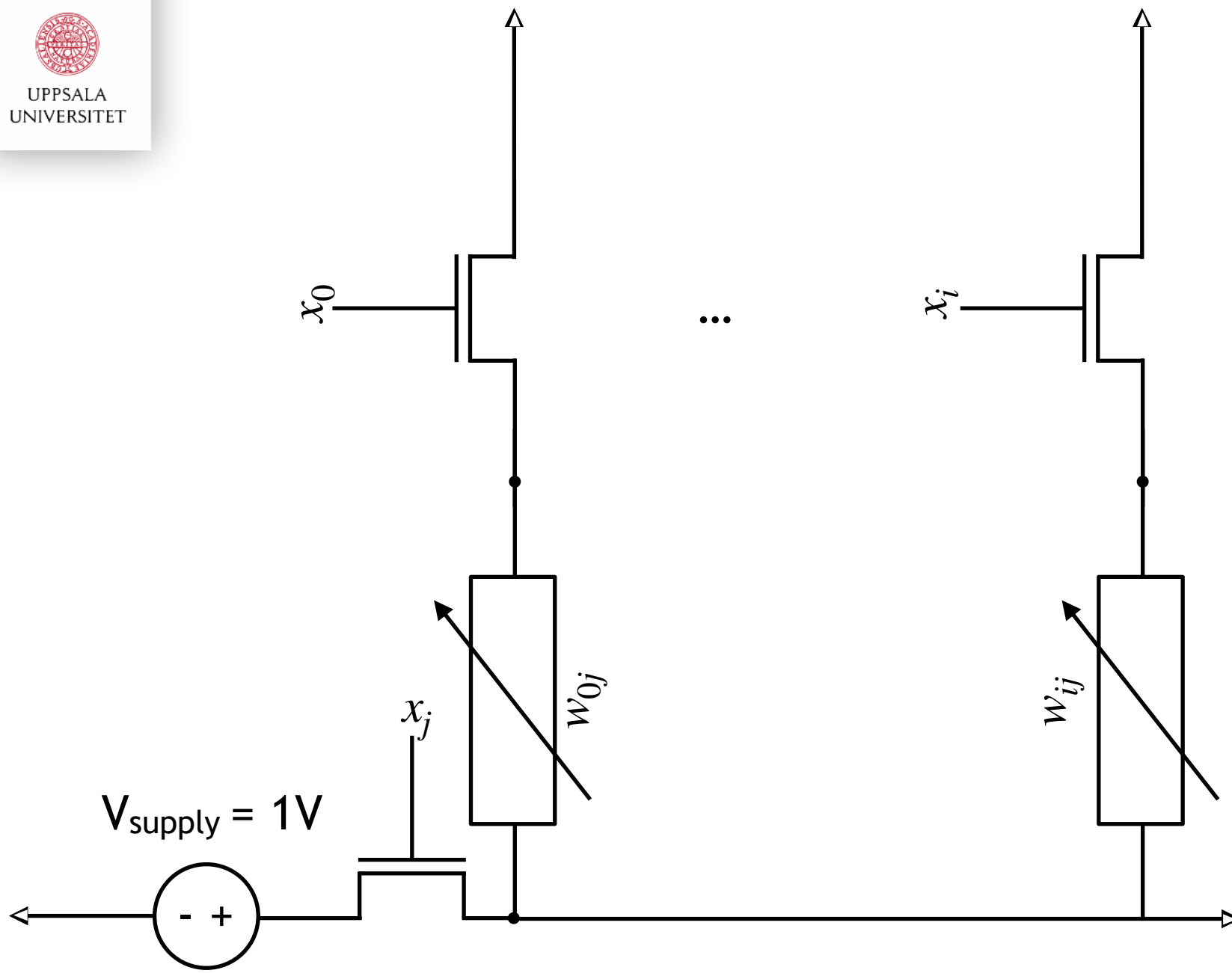


UPPSALA
UNIVERSITET



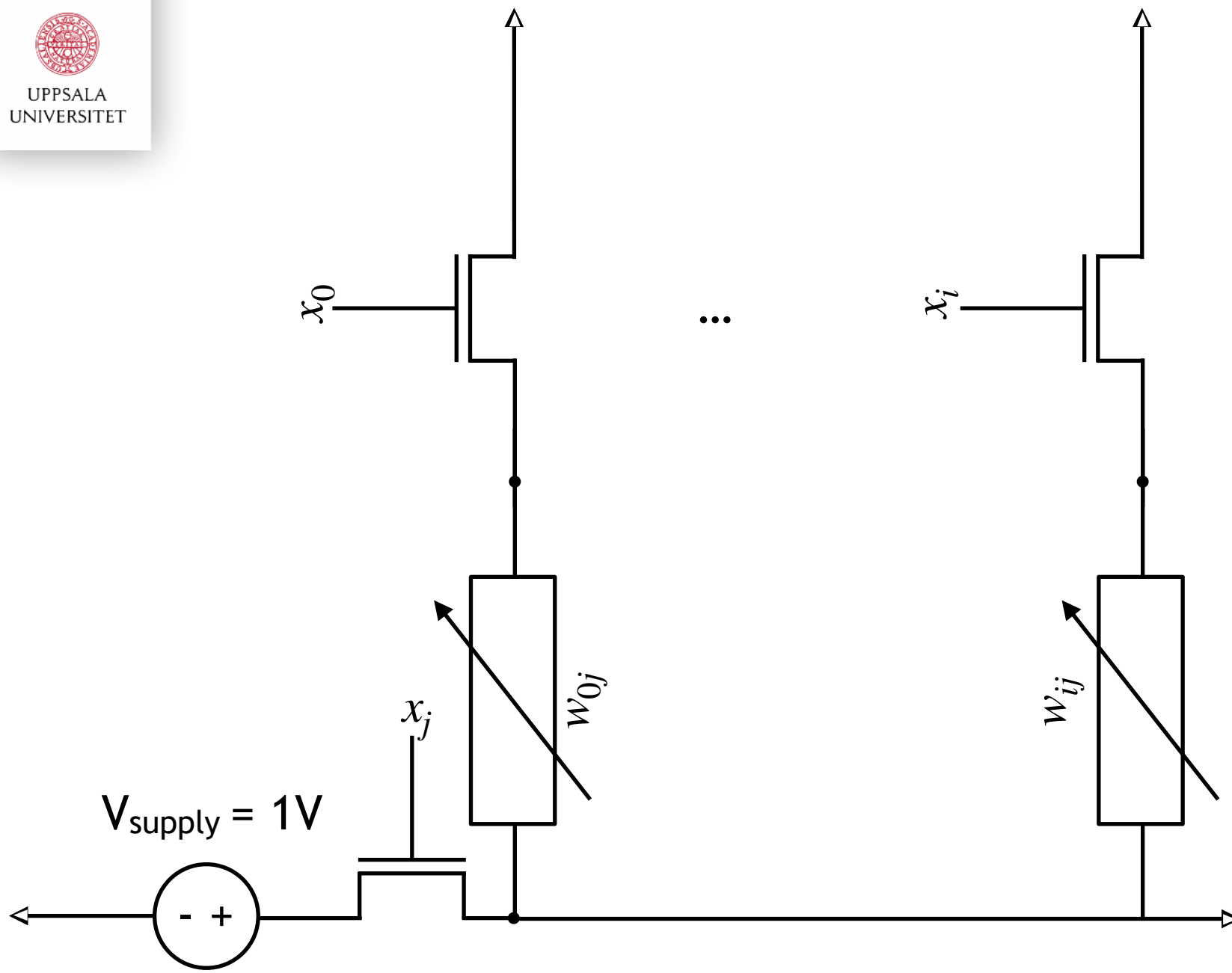


UPPSALA
UNIVERSITET





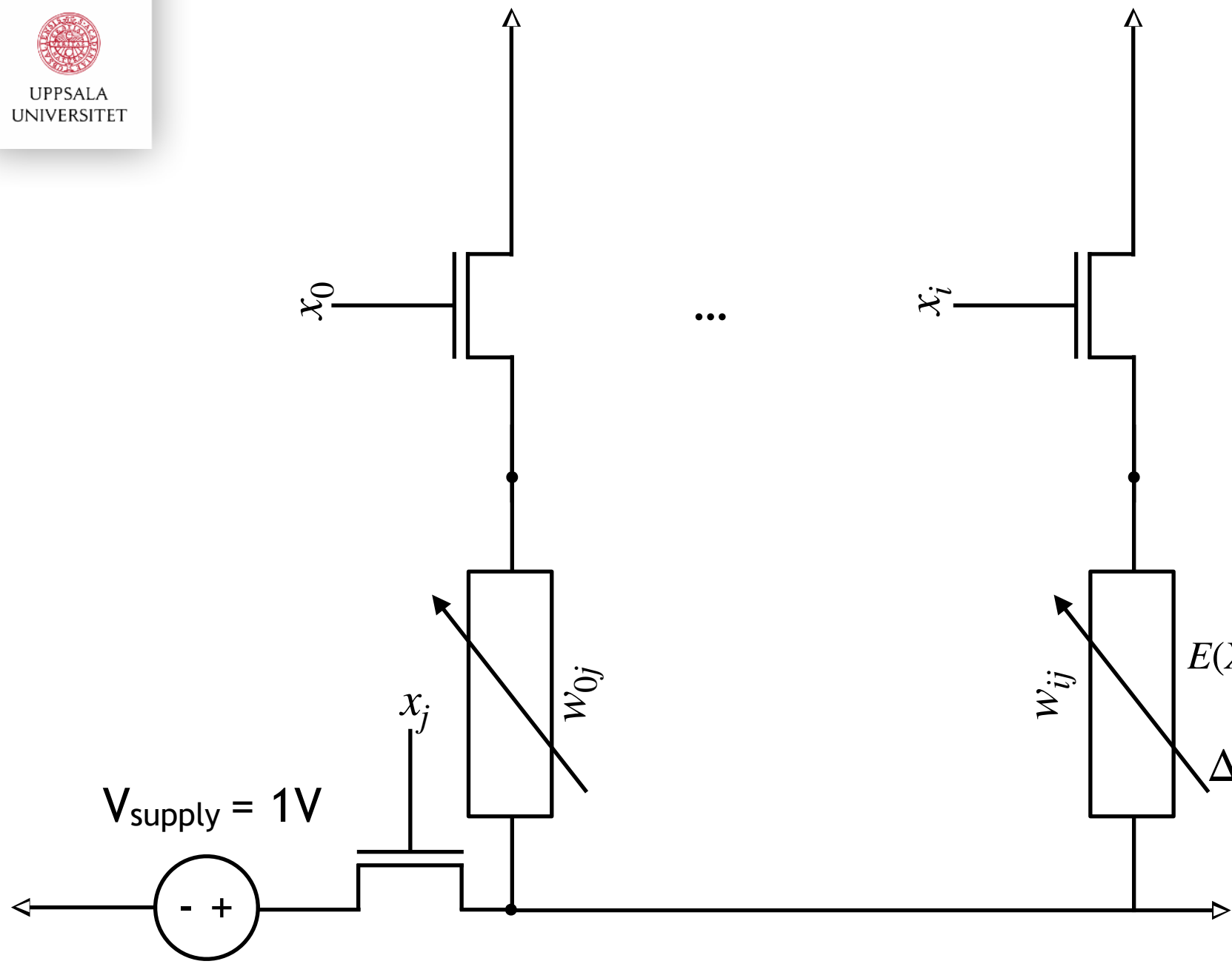
UPPSALA
UNIVERSITET



$$I_j = \sum_{i=0} I_{ji} = \sum_{i=0} x_j x_i w_{ij}$$



UPPSALA
UNIVERSITET



$$E(X) = -\frac{1}{2} \sum_j \sum_{i \neq j} x_j x_i w_{ij} - \sum_j x_j w_{jj}$$

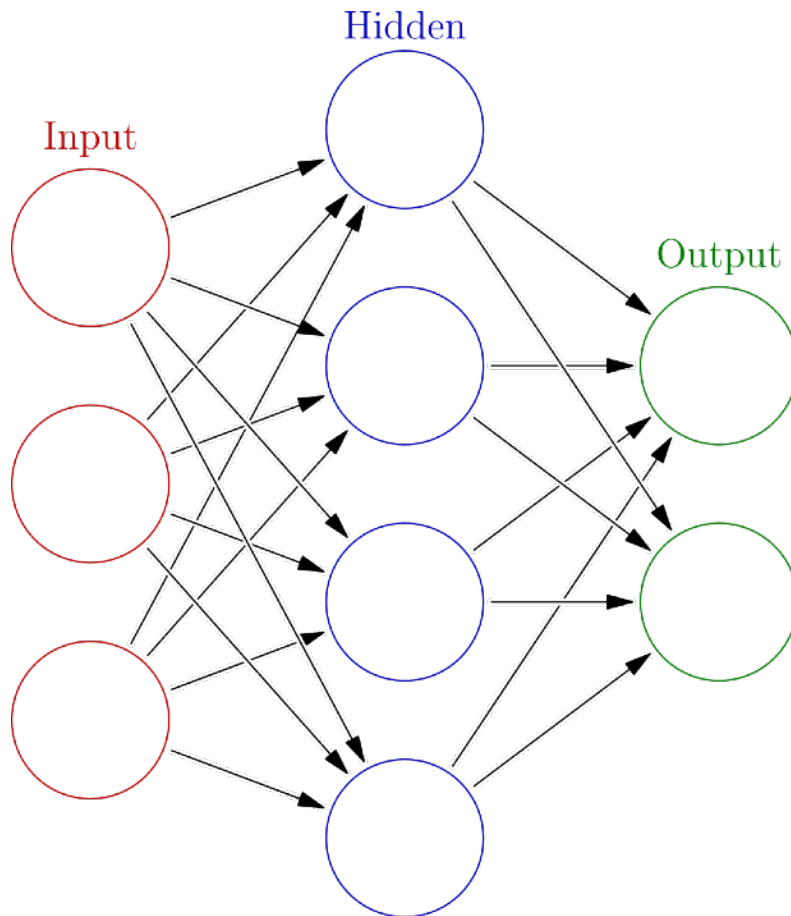
$$\Delta E = (2x_j - 1) \left(\sum_{i \neq j} x_i w_{ji} + w_{jj} \right)$$

$$I_j = \sum_{i=0} I_{ji} = \sum_{i=0}^{i \neq j} x_j x_i w_{ij}$$



UPPSALA
UNIVERSITET

Application to Other Problems



- Different mappings for different problems
- Deep learning uses “Restricted Boltzmann Machines”
 - Combination of layered approach of modern neural nets and Boltzmann Machines



UPPSALA
UNIVERSITET

System Organization and Training

ORGANIZATION



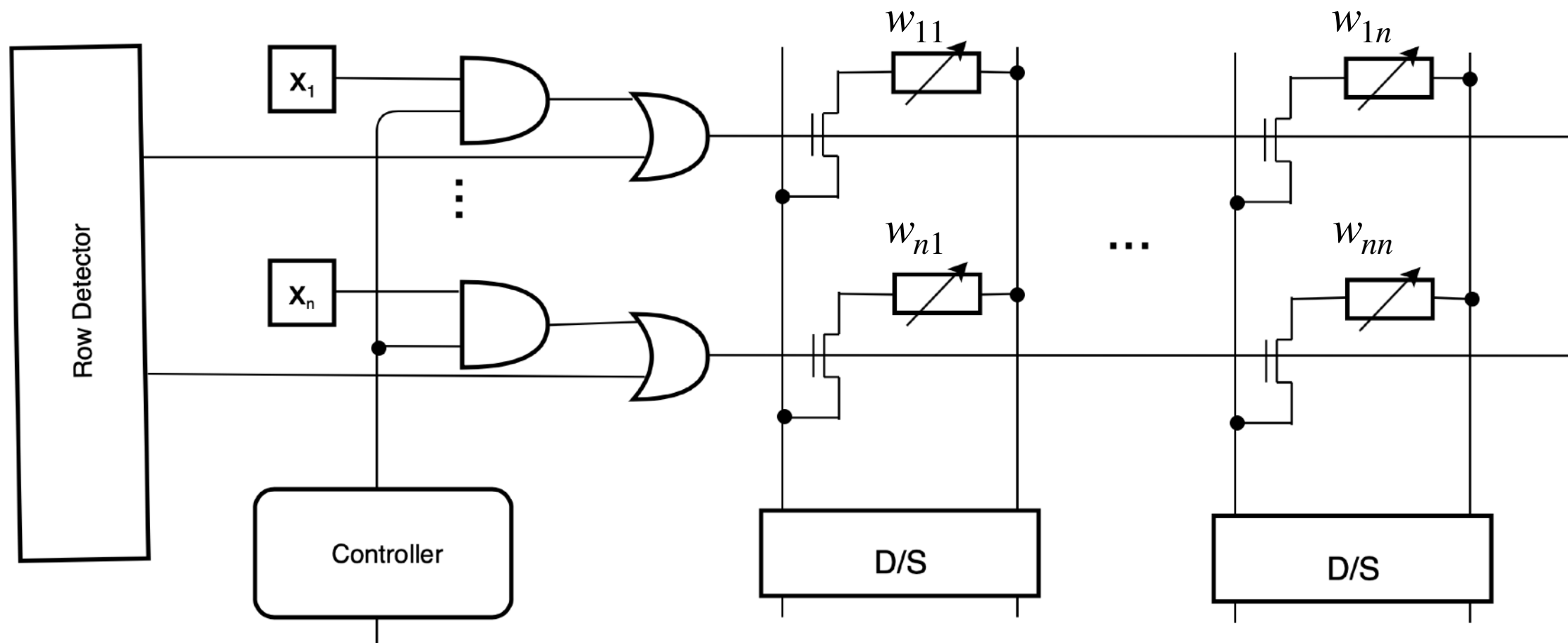
UPPSALA
UNIVERSITET

Interface

- Interfaces through regular DDRx
- Envisioned as a modular addition to any system with applicable workload



Organization





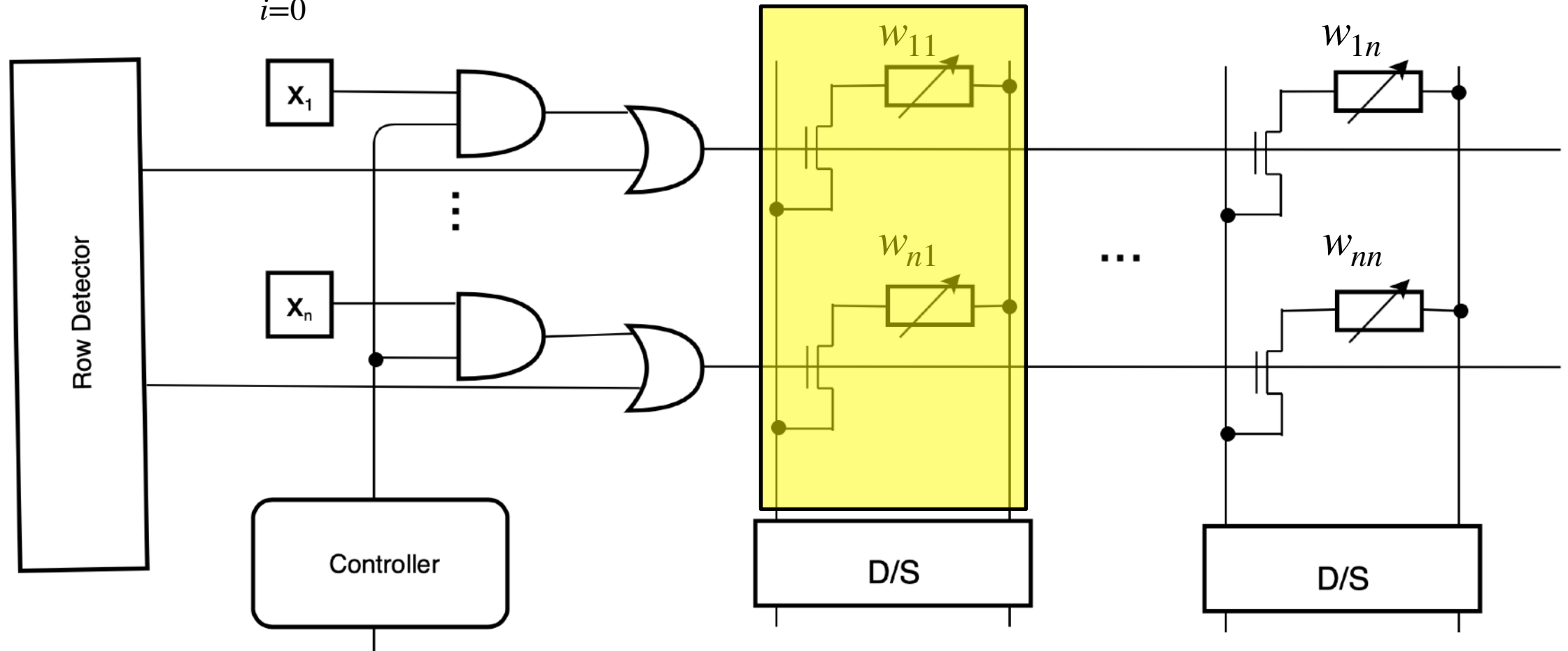
UPPSALA
UNIVERSITET

Current pulled by this column is

$$x_1 \cdot \sum_{i=0}^n x_i w_{ij}$$

Sensed by 'D/S' circuit

Organization





UPPSALA
UNIVERSITET

$$P(x^j | x) = \frac{1}{1 + e^{\frac{\Delta E}{c}}}$$

Consensus



UPPSALA
UNIVERSITET

$$P(x^j | x) = \frac{1}{1 + e^{\frac{\Delta E}{c}}}$$

Consensus

- Sigmoid function is hard to implement in hardware

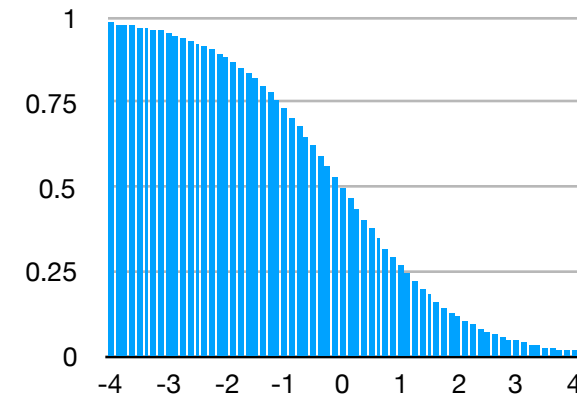


UPPSALA
UNIVERSITET

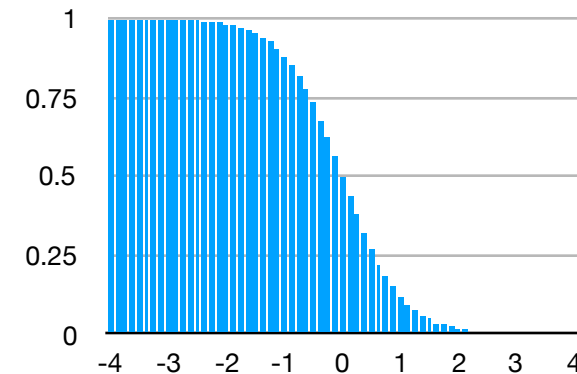
$$P(x^j | x) = \frac{1}{1 + e^{\frac{\Delta E}{C}}}$$

- Sigmoid function is hard to implement in hardware
- Use 64 evenly sampled points from expected range for values of ΔE [-4, 4]
- Different Cs are valid but paper uses C=1

Consensus



C=1



C=.5



UPPSALA
UNIVERSITET

Notes on Sensing Circuits

- Sense amplifiers sense one bit of the current
- Tree structures to sum them in quick fashion
- A large machine must tolerate a lot of inaccuracy and noise generated by current sensing

D/S



UPPSALA
UNIVERSITET

Performance and Energy

RESULTS



UPPSALA
UNIVERSITET

Optimization Problems

- 57.75x speedup over single threaded kernel
- 6.19x over PIM (process in memory)
- Many workloads achieved closer to 100x speedup over single threaded kernel



Optimization Problems

- 57.75x speedup over single threaded kernel
- 6.19x over PIM (process in memory)
- Many workloads achieved closer to 100x speedup over single threaded kernel
- 25x less energy used than single threaded kernel
- 5.2x less energy used than PIM



UPPSALA
UNIVERSITET

Deep Learning

- Even faster... 68.79x faster than single thread
- 6.89x faster than PIM



UPPSALA
UNIVERSITET

Deep Learning

- Even faster... 68.79x faster than single thread
- 6.89x faster than PIM
- 63x less energy than single thread
- 5.3x less energy than PIM