

# QIC 710: Introduction to Quantum Information Processing

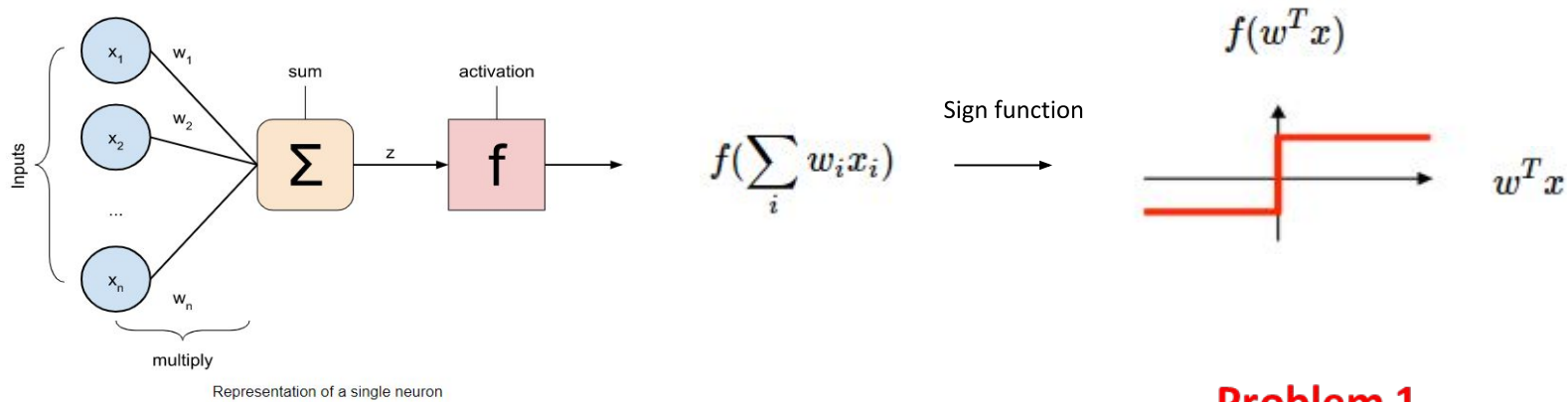
Instructor: Professor Richard Cleve

Prepared by: **Asad Raza**,  
Exchange Student at the  
Faculty of Mathematics,  
University of Waterloo.

# Quantum speed-up in global optimization of binary neural nets (arXiv: 1810.12948)

- Refresher on Classical Neural Networks and their training
  - Quantum Binary Neural Networks
  - Quantum Advantage in Training

# Classical Neuron



## Problem 1

NP-complete: No efficient algorithm to find the optimal  $\vec{w}$  exists.

$$C(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m, \vec{w}, \vec{y}) = \frac{1}{m} \sum_{i=1}^m (y_i - f(\vec{x}_i, \vec{w}))^2$$

↓ Number of training samples  $m$   
↑ True Labels  $y_i$

Spoiler Alert!

**Goal:** Find the *optimal*  $\vec{w}^*$  such that the cost,  $C(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m, \vec{w}, \vec{y})$ , is the lowest

# Training Classical Neural Networks Using Gradient Descent

## Problem 3

Gradient descent does not guarantee global convergence for non-convex functions

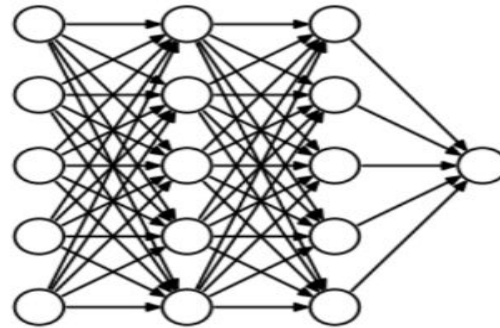
Weight of edge between neuron  $i$  in layer  $l - 1$  and neuron  $j$  in layer  $l$ , at  $(t - 1)$ -st training cycle

Small learning rate

$$w_{ij}^{(l-1,l)}(t) = w_{ij}^{(l-1,l)}(t-1) - \eta \frac{\partial C}{\partial w_{ij}^{(l-1,l)}}$$

But how to calculate  $\frac{\partial C}{\partial w_{ij}^{(l-1,l)}}$  ?

$$\frac{\partial C}{\partial w_{ij}^{(L-1,L)}} = \frac{\partial C}{\partial f(z_j^{(L)})} \frac{\partial f(z_j^{(L)})}{\partial w_{ij}^{(L-1,L)}} = \frac{\partial C}{\partial f(z_j^{(L)})} \frac{\partial f(z_j^{(L)})}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L-1,L)}}$$



$$f(z) = f\left(\sum_i w_i x_i\right) = f(w^T x)$$

## Problem 2

Negligible  $w_{ij}^{(l-1,l)}(t)$  updates  
 → No learning

The derivatives for most activation functions range between 0 and 1, which are multiplied  $L-1$  times to compute

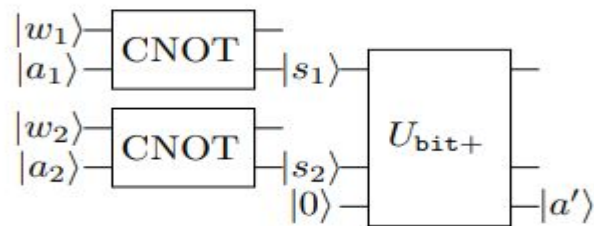
$$\frac{\partial C}{\partial w_{ij}^{(L-1,L)}} \rightarrow \frac{\partial C}{\partial w_{ij}^{(l-1,l)}} \approx 0$$

# Quantum Binary Neuron

Encode the weights, the inputs and the activations take values either +1 or -1. Dramatic Simplification (?)

Not as bad as it sounds. See *Bengio et al. (arXiv: 1602.02830)*

Scheme: Represent -1 as  $|1\rangle$  and 1 as  $|0\rangle$



$$V(|w_1\rangle_1|a_1\rangle_2|w_2\rangle_3|a_2\rangle_4|0\rangle_5) = |w_1\rangle_1|s_1\rangle_2|w_2\rangle_3|s_2\rangle_4|a'\rangle_5$$

Figure 5. *Functioning of QBN*: the input data is encoded into quantum states. The multiplication succeeds by the CNOT gate and outputs the states  $|s_1\rangle$ ,  $|s_2\rangle$ . Finally, the Toffoli gate, which has known decompositions into elementary gates, executes  $U_{\text{bit+}}$ , leading to the output state  $|a'\rangle$ .

Generalize to Quantum Binary Feedforward Neural Network (QBFFNN)

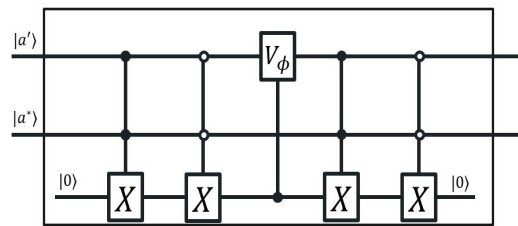
$$U(|a_1\rangle|a_2\rangle \dots |a_n\rangle|w_1\rangle|w_2\rangle \dots |w_N\rangle|0\rangle^{\otimes p}) = |\widetilde{a}_1, \widetilde{a}_2, \dots, \widetilde{a}_n, w_1, w_2, \dots, w_N, \text{rest}, a'_1, a'_2, \dots, a'_k\rangle$$

↓  
Saved values to achieve unitarity

QBFFNN outputs, encoded by ancillas

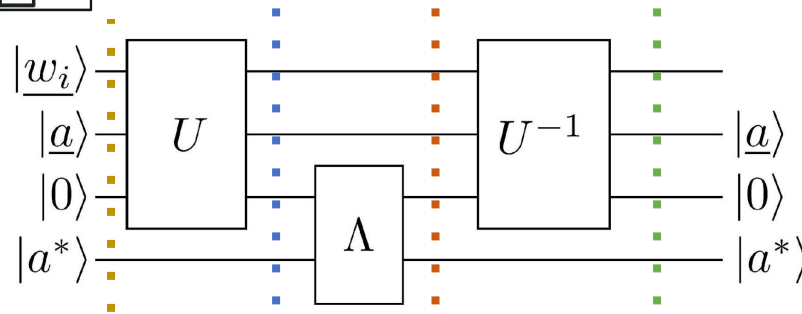
# Training QBFNN: Single Weight Configuration

Pick **one pair** of training data and call it  $(\vec{a}, a^*)$ , and **one possible** weight configuration  $\vec{w}_i$  (out of  $2^N$ )



3. After oracle  $\Lambda$  :  
 $\Lambda(|\underline{w}_i, \tilde{a}, a'_i\rangle|a^*\rangle) = e^{i\Delta\theta_{\vec{w}_i}} |\underline{w}_i, \tilde{a}, a'_i\rangle|a^*\rangle$

4. Uncomputation to disentangle:  
 $U^\dagger (e^{i\Delta\theta_{\vec{w}_i}} |\underline{w}_i, \tilde{a}, a'_i\rangle|a^*\rangle) = e^{i\Delta\theta_{\vec{w}_i}} |\underline{w}_i\rangle|\underline{a}\rangle|0\rangle|a^*\rangle$



1. Input  
 $|\underline{w}_i\rangle|\underline{a}\rangle|0\rangle|a^*\rangle$

2. After neuron/NN action:  
 $U(|\underline{w}_i\rangle|\underline{a}\rangle|0\rangle|a^*\rangle) = |\underline{w}_i, \tilde{a}, a'_i\rangle|a^*\rangle$

$$e^{i\Delta\theta_{\vec{w}_i}} |\underline{w}_i\rangle$$

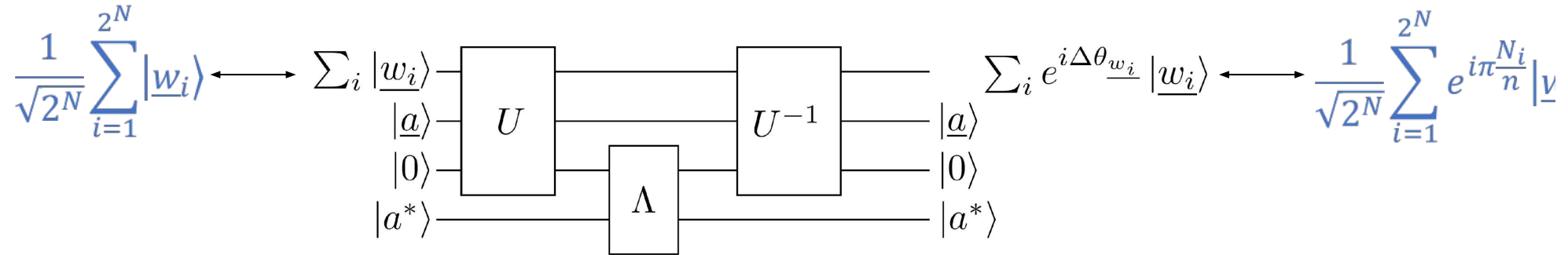
Reuse this weight state,  $e^{i\Delta\theta_{\vec{w}_i}} |\vec{w}_i\rangle$ , for the next training sample

After looping through  $n$  data samples, the final weight state will read:

$N_i \in \{0, 1, \dots, n\}$  is the number of correct predictions

$$e^{i\pi \frac{N_i}{n}} |\vec{w}_i\rangle$$

# Training QBFNN: Superposition of Weights



ONE loop instead of  $2^N$  loops for each weight configuration

Since we loop through all the  $n$  training samples coherently, we can represent  $n$  rounds of  $U^\dagger \Lambda U$  as **one gigantic unitary  $\tilde{U}$**

$$\left( \frac{1}{\sqrt{2^N}} \sum_{i=1}^{2^N} |\vec{w}_i\rangle \right) \otimes |\vec{a}_1\rangle |\vec{a}_2\rangle \dots |\vec{a}_n\rangle |0\rangle^{\otimes p} |a_1^*\rangle |a_2^*\rangle \dots |a_n^*\rangle$$

$$\xrightarrow{\tilde{U}} \left( \frac{1}{\sqrt{2^N}} \sum_{i=1}^{2^N} e^{i\pi \frac{N_i}{n}} |\vec{w}_i\rangle \right) \otimes |\vec{a}_1\rangle |\vec{a}_2\rangle \dots |\vec{a}_n\rangle |0\rangle^{\otimes p} |a_1^*\rangle |a_2^*\rangle \dots |a_n^*\rangle$$

# Phase Estimation

Observe that  $|\psi_i\rangle = |\vec{w}_i\rangle (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle$ ,  $i = 1, 2, \dots, 2^N$  are **eigenvectors** of  $\tilde{U}$  with **eigenvalues**  $e^{i\pi \frac{N_i}{n}}$  since  $\tilde{U}|\psi_i\rangle = e^{i2\pi \frac{N_i}{2n}} |\psi_i\rangle$

But what is a “good”  $N_i$ ?

Estimate  $\varphi^{(i)} = \frac{N_i}{2n}$  using phase estimation to get the t-bit binary encoding  $\varphi^{(i)} = |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle$

Fix a threshold count  $N_t \in \mathbb{N}^+$ . If  $N_i \geq N_t$ ,  $N_i$  classifies as a “good enough” configuration.

Since phase estimation works for a superposition of eigenvectors,  $\sum_i |\vec{w}_i\rangle \otimes (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle \otimes |0\rangle^{\otimes t}$   
 $\xrightarrow{PE} \sum_i |\vec{w}_i\rangle \otimes |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle$



# Getting the “quality” weight vector

Define another **oracle**  $O_{\pm 1}$ , which acts on  $|\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle$  as follows:

$$O_{\pm 1} |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle = \begin{cases} -|\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle, & \text{if } N_i \geq N_t \\ |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle, & \text{if } N_i < N_t \end{cases}$$

$$|\widehat{W}\rangle = \frac{1}{\sqrt{2^N}} \sum_i (-1)^{N_i \geq N_t} |\underline{w}_i\rangle$$

Started off with  $\sum_i |\vec{w}_i\rangle \otimes (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle \otimes |0\rangle^{\otimes t}$

$$\xrightarrow{PE} \sum_i |\vec{w}_i\rangle \otimes |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle$$

$$\xrightarrow{O_{\pm 1}} \sum_i (-1)^{N_i \geq N_t} |\vec{w}_i\rangle \otimes |\varphi_1^{(i)} \varphi_2^{(i)} \dots \varphi_t^{(i)}\rangle (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle \xrightarrow{PE^{-1}} (\sum_i (-1)^{N_i \geq N_t} |\vec{w}_i\rangle) (\otimes_{j=1}^n |\vec{a}_j\rangle \otimes |a_j^*\rangle) \otimes |0\rangle \otimes |0\rangle^{\otimes t}$$

**Problem: Entanglement between weights and phases!**

Solution:  $PE^\dagger$

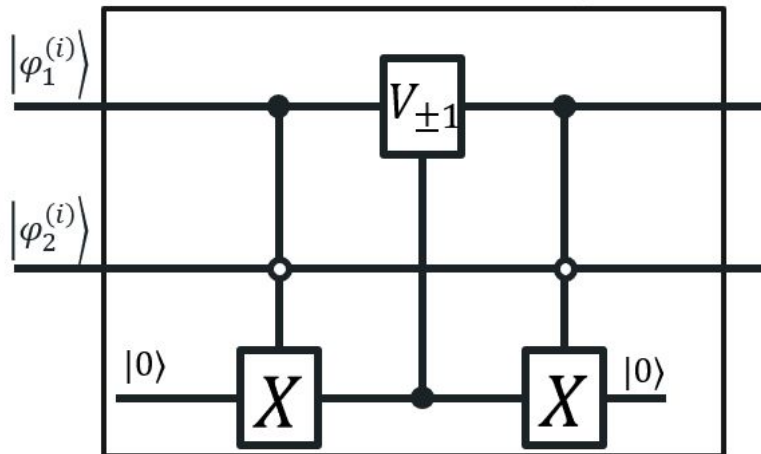
**Uncompute phase estimation!**

# Applying Grover's to find the "good" weights

Example:  $n = 2$

$$N_t = n \rightarrow \frac{N_i}{2n} = \frac{1}{2}$$

$\varphi^{(i)} = \frac{N_i}{2n}$ , the encoding  $|\varphi_1^{(i)} \varphi_2^{(i)}\rangle$  will be correspondingly  $|10\rangle$



After  $PE^\dagger$ , we had the weight superposition  $|\widehat{W}\rangle = \frac{1}{\sqrt{2^N}} \sum_i (-1)^{N_i \geq N_t} |w_i\rangle$



**A Grover state!**

The Grover oracle acted as

$$O|x\rangle = \begin{cases} -|x\rangle, & \text{if } x \text{ is solution to search} \\ |x\rangle, & \text{else} \end{cases}$$



$$O(\sum_i |x_i\rangle) = \sum_i (-1)^{f(x_i)} |x_i\rangle, \text{ where } f(x_i) = 1 \text{ if } x \text{ is solution and } 0 \text{ else}$$



Iterate for  $k^* = \sqrt{\frac{2^N}{M}} \frac{\pi}{4}$  times to get the optimal weight vector,  $\vec{w}^*$

# What have we achieved?

- To ensure globally optimal set of weight, need  $\approx 2^N$  steps ( $\approx$  Problem 1)
- Vanishing/Exploding gradient issues for gradient based trainings (Problem 2)
- Gradient descent almost never finds globally optimal solutions. Difference between local and global optimum can be huge (Problem 3)

- Find globally optimal set of weights in  $\approx \sqrt{2^N}$  steps
- No vanishing gradients (Problem 2 resolved). Even better, no gradients involved.
- Guaranteed to find globally optimal weight configuration (Problem 3 resolved)

Thank you for  
listening!  
Questions?