

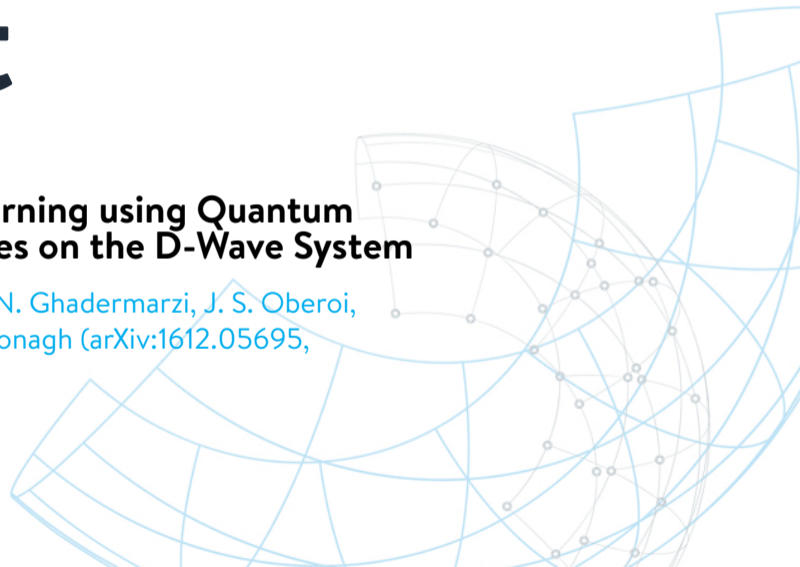


Reinforcement Learning using Quantum Boltzmann Machines on the D-Wave System

Joint work with A. Levit, N. Ghadermarzi, J. S. Oberoi, E. Zahendinejad, and P. Ronagh (arXiv:1612.05695, 1706.00074)

Daniel Crawford | 1QBit

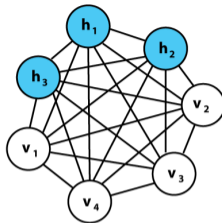
September 27, 2017



Outline

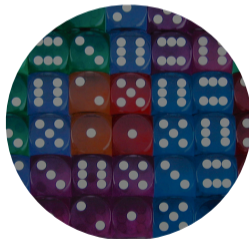
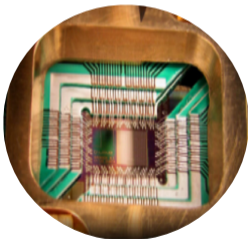
Reinforcement learning
Quantum Boltzmann machines
Quantum Monte Carlo simulations
Results

RL



GBM

QBM



QMC
Simulations

Markov Decision Process

- Markov decision process

states and actions: (finite) sets S and A ;

controlled Markov chain: defined by a transition kernel $\mathbb{P}(s' \in S | s \in S, a \in A)$;

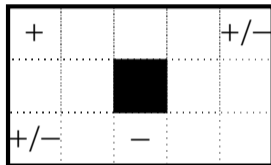
immediate reward structure: given via $r : S \times A \rightarrow \mathbb{R}$;

discount factor: a constant $\gamma \in [0, 1)$.

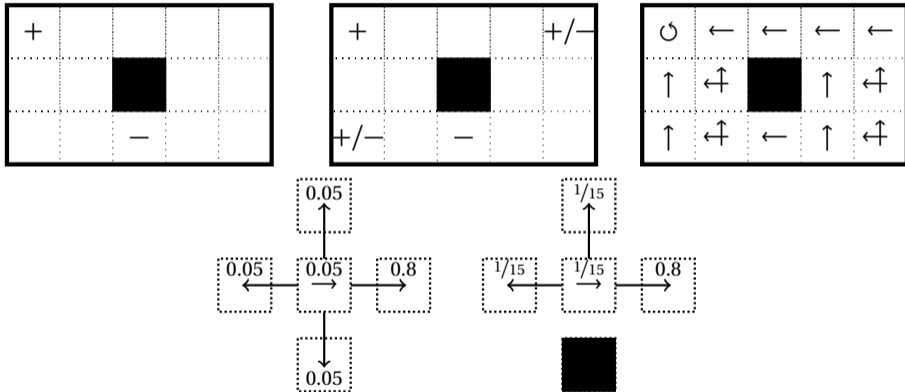
- The stationary policy

$$\pi : S \longrightarrow A$$

reduces the MDP into a time-homogeneous Markov chain, Π , with kernel $\mathbb{P}(s' | s, \pi(s))$.



Example: Grid-World Problem



Markov Decision Problem

- Goal: to solve the optimization problem

$$\pi^* = \operatorname{argmax}_{\pi} V(\pi, s).$$

- The discounted reward function

$$V(\pi, s) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r(\Pi_i^s, \pi(\Pi_i^s)) \right].$$

- Bellman recursion:

$$\begin{aligned} V(\pi, s) &= \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r(\Pi_i^s, \pi(\Pi_i^s)) \right] \\ &= \mathbb{E}[r(\Pi_0^s, \pi(\Pi_0^s))] \\ &\quad + \gamma \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r(\Pi_{i+1}^s, \pi(\Pi_{i+1}^s)) \right] \\ &= \mathbb{E}[r(s, \pi(s))] + \gamma \sum_{s' \in S} \mathbb{P}(s'|s, \pi(s)) V(\pi, s') \end{aligned}$$

Q-function

- Maps a triplet (π, s, a) to the expected value of the reward of the Markov chain that begins with taking action a at initial state s and continuing according to π :

$$Q(\pi, s, a) = \mathbb{E}[r(s, a)] + \mathbb{E}\left[\sum_{i=1}^{\infty} \gamma^i r(\Pi_i^s, \pi(\Pi_i^s))\right].$$

- Reconstructs π^* and $V^*(s) = V(\pi^*, s)$:

$$V^*(s) = \max_a Q^*(s, a),$$

for $Q^*(s, a) = \max_{\pi} Q(\pi, s, a)$ and

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a).$$

Parametrization of the Temporal Difference

- The goal is to minimize

$$Q_{n+1}(s, a) - Q_n(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} \mathbb{P}(s'|s, a) \max_a Q_n(s', a) - Q_n(s, a)$$

Parametrization of the Temporal Difference

- The goal is to minimize

$$Q_{n+1}(s, a) - Q_n(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} \mathbb{P}(s'|s, a) \max_a Q_n(s', a) - Q_n(s, a)$$

↑
temporal difference (E_{TD})

Parametrization of the Temporal Difference

- The goal is to minimize

$$Q_{n+1}(s, a) - Q_n(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} \mathbb{P}(s'|s, a) \max_a Q_n(s', a) - Q_n(s, a)$$

↑
temporal difference (E_{TD})

- Local parametrization

$$Q(s, a) = Q(s, a; \boldsymbol{\theta})$$

Parametrization of the Temporal Difference

- The goal is to minimize

$$Q_{n+1}(s, a) - Q_n(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} \mathbb{P}(s'|s, a) \max_a Q_n(s', a) - Q_n(s, a)$$

↑ temporal difference (E_{TD})

- Local parametrization

$$Q(s, a) = Q(s, a; \boldsymbol{\theta})$$

- Descent along

$$-\nabla_{\boldsymbol{\theta}} (E_{TD})^2 = -E_{TD} \nabla_{\boldsymbol{\theta}} E_{TD}$$
$$\xrightarrow{\text{SGD}} \left(r_n(s_n, a_n) + \gamma \max_{a_{n+1}} Q(s_{n+1}, a_{n+1}) - Q(s_n, a_n) \right) \frac{\partial}{\partial \boldsymbol{\theta}} Q(s_n, a_n)$$

Quantum Boltzmann Machines

GBMs, DBMs, and QBMs as function approximators

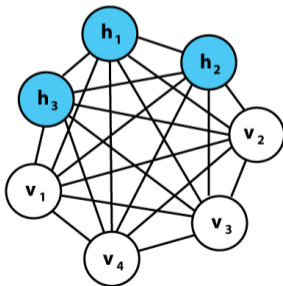


Quantum Boltzmann Machines

Amin et al., 2016

- Classical GBM

$$\mathcal{E}(v, h) = - \sum_{v \in V, h \in H} w^{vh} v h - \sum_{\{v, v'\} \subseteq V} w^{vv'} v v' - \sum_{\{h, h'\} \subseteq H} w^{hh'} h h'$$



Quantum Boltzmann Machines

- Classical GBM

$$\mathcal{E}(v, h) = - \sum_{v \in V, h \in H} w^{vh} v h - \sum_{\{v, v'\} \subseteq V} w^{vv'} v v' - \sum_{\{h, h'\} \subseteq H} w^{hh'} h h'$$

- Clamped GBM (fixed assignment v of the visible binary variables)

$$\mathcal{E}_v(h) = - \sum_{v \in V, h \in H} w^{vh} v h - \sum_{\{v, v'\} \subseteq V} w^{vv'} v v' - \sum_{\{h, h'\} \subseteq H} w^{hh'} h h'$$

- Clamped QBM

$$\mathcal{H}_v = - \sum_{v \in V, h \in H} w^{vh} v \sigma_h^z - \sum_{\{v, v'\} \subseteq V} w^{vv'} v v' - \sum_{\{h, h'\} \subseteq H} w^{hh'} \sigma_h^z \sigma_{h'}^z - \Gamma \sum_{h \in H} \sigma_h^x$$

Free Energy of a QBM

- Equilibrium free energy

$$F(\nu) := -\frac{1}{\beta} \ln Z_\nu = \langle \mathcal{H}_\nu \rangle + \frac{1}{\beta} \text{tr}(\rho_\nu \ln \rho_\nu).$$

$\frac{1}{K_B T}$ $\text{tr}(e^{-\beta \mathcal{H}_\nu})$ $\frac{1}{Z_\nu} e^{-\beta \mathcal{H}_\nu}$

- Entropy $-\text{tr}(\rho_\nu \ln \rho_\nu)$
- Gibbs measure

$$\langle \mathcal{H}_\nu \rangle = \frac{1}{Z_\nu} \text{tr}(\mathcal{H}_\nu e^{-\beta \mathcal{H}_\nu})$$

Free Energy as a Function Approximator

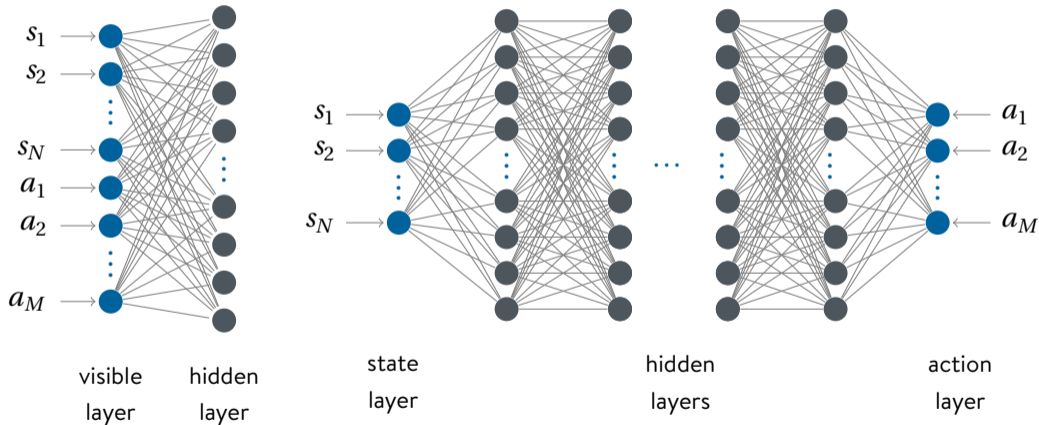
$$Q(s, a) \approx -F(s, a) = -F(s, a; \omega) \quad \leftarrow \text{QBM coupling strength}$$

$$\Delta\omega = -\varepsilon(r(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)) \frac{\partial F}{\partial \omega}$$

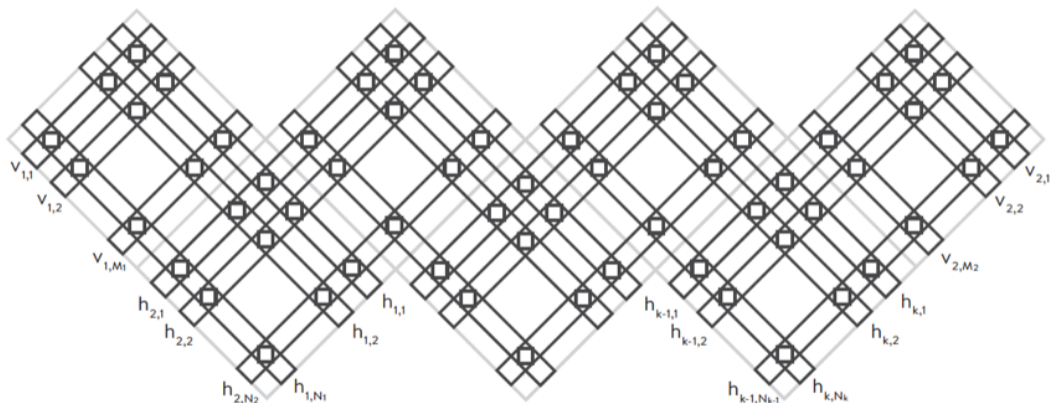
$$\begin{aligned} \frac{\partial F(s, a)}{\partial \omega} &= \frac{1}{Z_{s,a}} \frac{\partial}{\partial \omega} \text{tr}(e^{-\beta \mathcal{H}_{s,a}}) \\ &= -\frac{1}{Z_{s,a}} \text{tr}(\beta e^{-\beta \mathcal{H}_{s,a}} \frac{\partial}{\partial \omega} \mathcal{H}_{s,a}) \\ &= -\beta \left\langle \frac{\partial}{\partial \omega} \mathcal{H}_{s,a} \right\rangle \end{aligned}$$

$$\begin{aligned} \Delta\omega^{vh} &= \varepsilon(r(s, a) \\ &\quad - \gamma \min_{a'} F(s', a') + F(s, a)) \nu \langle \sigma_h^z \rangle \\ \Delta\omega^{hh'} &= \varepsilon(r(s, a) \\ &\quad - \gamma \min_{a'} F(s', a') + F(s, a)) \langle \sigma_h^z \sigma_{h'}^z \rangle \end{aligned}$$

RBM and DBM Layout



A DBM as a Layout of Superconducting Qubits



Quantum Monte Carlo Simulations

Approximating free energy and spin expectations

Suzuki–Trotter Expansion

The effective Hamiltonian of an Ising model with transverse field

Transverse field Ising Hamiltonian:

$$H = - \sum_{(i,j)} J_{ij} \sigma_i^z \sigma_j^z - \Gamma \sum_{i=1}^N \sigma_i^x.$$

Suzuki–Trotter Expansion

The effective Hamiltonian of an Ising model with transverse field

Transverse field Ising Hamiltonian:

$$H = - \sum_{(i,j)} J_{ij} \sigma_i^z \sigma_j^z - \Gamma \sum_{i=1}^N \sigma_i^x.$$

The partition function is approximated as

$$Z \approx \left(\frac{1}{2} \sinh \frac{2\beta\Gamma}{M} \right)^{\frac{nM}{2}} \text{tr} \left[e^{(-\beta \mathcal{H}_{\text{eff}})} \right]$$

with an *effective* Hamiltonian

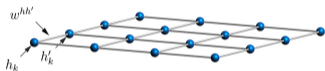
$$\mathcal{H}_{\text{eff}}(\sigma) = - \sum_{i,j,k} \frac{J_{ij}}{M} \sigma_{ik} \sigma_{jk} - \sum_{i,k} \frac{1}{2\beta} \ln \coth \left(\frac{\beta\Gamma}{M} \right) \sigma_{ik} \sigma_{ik+1}.$$

Approximating the Expected Values

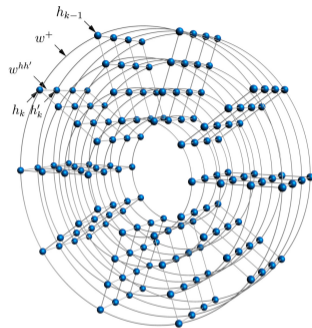
Theorem [Suzuki, 1976] Spin correlations $\langle \sigma_i^z \sigma_j^z \rangle$ are calculated as those of the corresponding $(d + 1)$ -dimensional Ising model:

$$\langle \sigma_i^z \sigma_j^z \rangle = \frac{1}{Z_{d+1}} \lim_{n \rightarrow \infty} \sum_{\sigma} \sigma_i \sigma_j e^{(-\beta \mathcal{H}_{\text{eff}}(\sigma))}.$$

Simulated Quantum Annealing



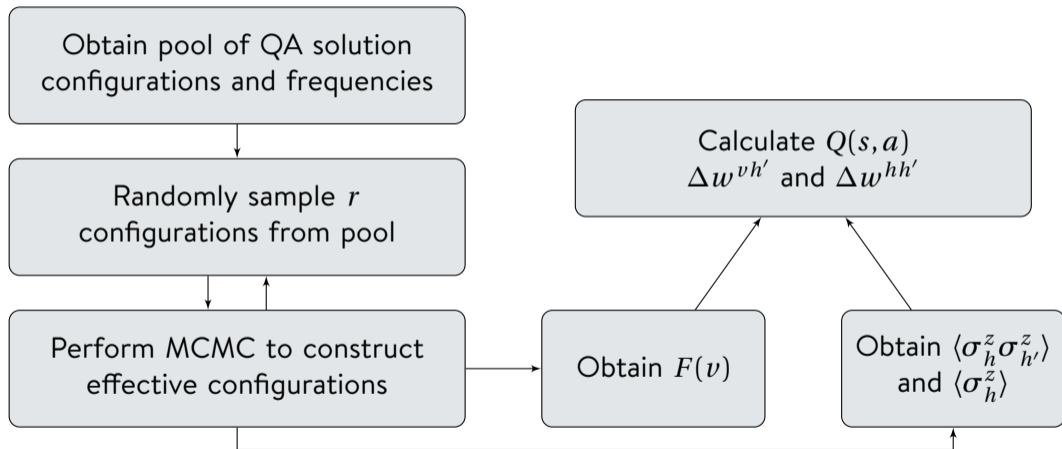
(a) Original transverse field Ising Hamiltonian



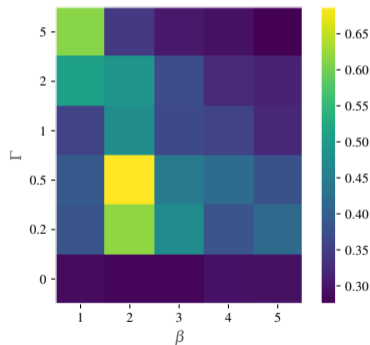
(b) Effective classical Hamiltonian in one dimension higher

The dynamics are governed by the MCMC method using Metropolis acceptance probabilities, as the transverse field strength Γ is slowly decreased to zero.

Using QA Samples to Build Effective Configurations



Issue: Effective Temperature and Transverse Field



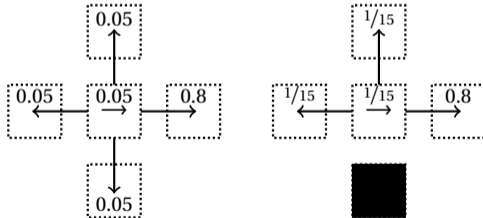
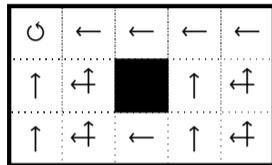
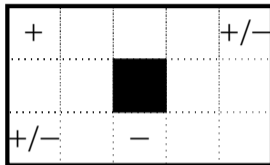
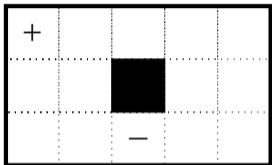
$$\mathcal{H}_{\text{eff}}(\sigma) = - \sum_{i,j,k} \frac{J_{ij}}{M} \sigma_{ik} \sigma_{jk} - \sum_{i,k} \frac{1}{2\beta} \ln \coth \left(\frac{\beta \Gamma}{M} \right) \sigma_{ik} \sigma_{ik+1}$$

Experimental Results

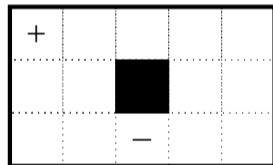
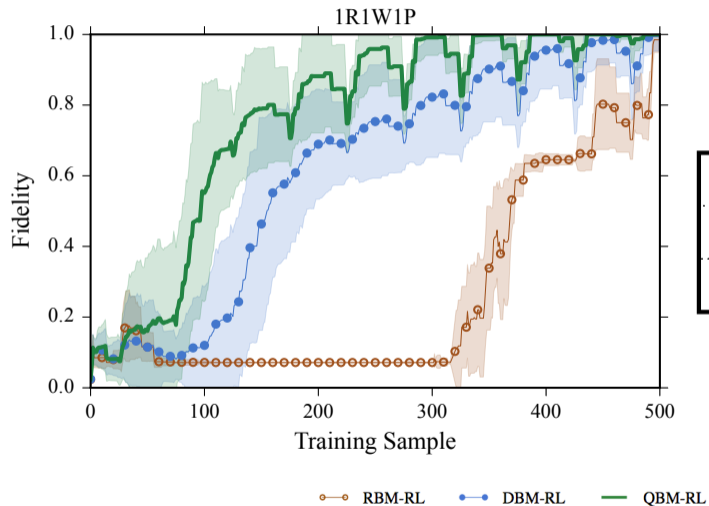
Reinforcement learning of a maze



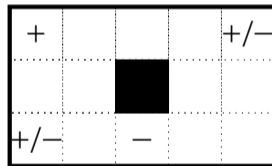
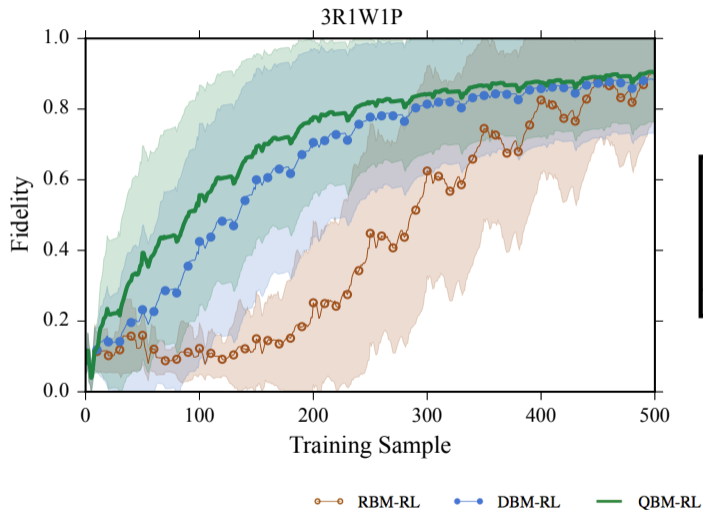
Grid-World Problem



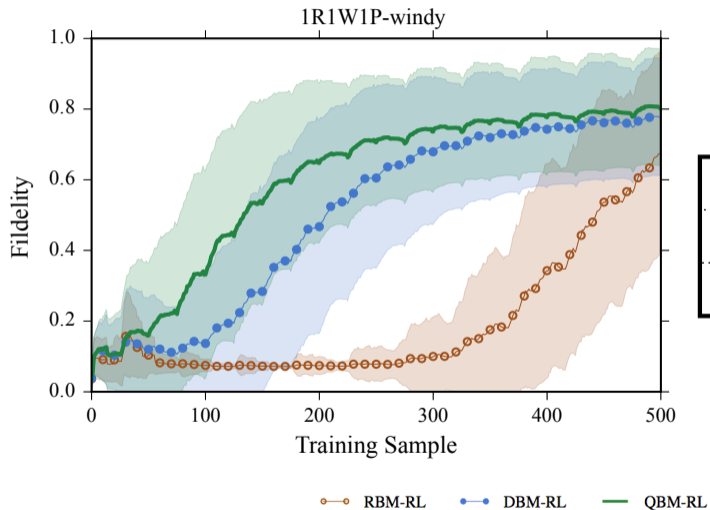
Results



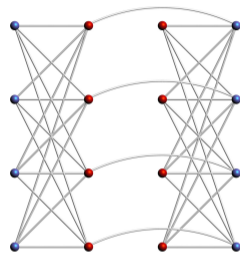
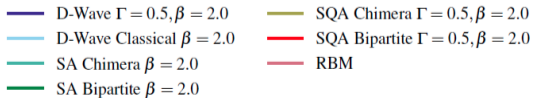
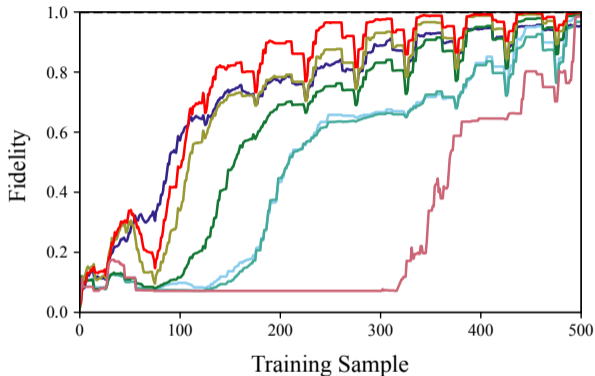
Results – Stochastic Rewards



Results – Stochastic Transitions



Results



Chimera Graph



Thank you!

Daniel.Crawford@1qbit.com