

# The Future of Drug Discovery: *Quantum-Based Machine Learning Simulation (QMLS)*

## I. SPECIFIC IMPLEMENTATION

### A. Machine Learning Molecule Generation (MLMG) & Machine Learning Molecule Variation (MLMV): 02 Feb 2023-04 April 2023

During implementation, due to Quantum Computers being too expensive and inaccessible, I decreased the size of the RNN models but it still showed promising results. I trained a forward-RNN model with 1024 hidden units and 5 layers to generate hit molecules that can bind with the target protein<sup>1</sup>. The RNN model takes as input a sequence of characters representing a molecule in SMILES format and outputs the next character in the sequence. The model is trained to minimize the cross-entropy loss between the predicted and actual characters. Here is a segment of code used to implement the RNN model in PyTorch:

```
MLMG-RNN_V9.0.py 9 X
Users > Yifan > Desktop > MLMG-RNN_V9.0.py > ...
1 import torch
2 import torch.nn as nn
3 import deepspeed
4
5 from torch.utils.data import DataLoader
6 from molecule_chef.model import RNN
7 from molecule_chef.dataset import SMILESdataset
8 from molecule_chef.utils import one_hot_encode, one_hot_decode
9
10 class RNN(nn.Module):
11     def __init__(self, input_size, hidden_size, output_size, num_layers):
12         super(RNN, self).__init__()
13         self.hidden_size = hidden_size
14         self.num_layers = num_layers
15         self.embedding = nn.Embedding(input_size, hidden_size)
16         self.rnn = nn.RNN(hidden_size, hidden_size, num_layers)
17         self.fc = nn.Linear(hidden_size, output_size)
18
19     def forward(self, input, hidden):
20         embedded = self.embedding(input).view(1, 1, -1)
21         output, hidden = self.rnn(embedded, hidden)
22         output = self.fc(output.view(1, -1))
23         return output, hidden
24
25     def init_hidden(self):
26         return torch.zeros(self.num_layers, 1, self.hidden_size)
27
28 # Load the pre-trained model
29 model = RNN(input_size=35, hidden_size=1024, output_size=35, num_layers=5)
30 model.load_state_dict(torch.load("molecule_sensibility.pth"))
31
32 # Move the model to the device (CPU or GPU)
33 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
34 model.to(device)
```

For transfer learning, I used a different dataset and fine-tuned the RNN model on each. I used the following datasets for each subtask:

**Molecule sensibility.** I used the *ZINC database* to train the model to generate valid and sensible molecules. The ZINC database contains over 230 million commercially available compounds in ready-to-dock, 3D formats. I used 90% of the data for training and 10% for validation. I trained the model for 10 epochs with a learning rate of 0.001 and a batch size of 64.

**Bond reconstruction.** I used the *REAXYS database* to train the model to reconstruct the bonds of a molecule given a partial structure. The *REAXYS database* contains over 105 million organic, inorganic and organometallic compounds, 41 million chemical reactions, 500 million published experimental facts, 16,000 chemistry related periodicals, and six indexing sources for a cross-disciplinary view of chemistry. I used 80% of the data for training and 20% for validation. I trained the model for 5 epochs with a learning rate of 0.0005 and a batch size of 32.

**Reaction prediction.** I used the *USPTO database* to train the model to predict the product of a chemical reaction given the reactants and reagents. The USPTO database contains over 2 million chemical reactions extracted from granted US patents and patent applications from 1976 to September 2016. I used 70% of the data for training and 30% for validation and trained the model for 5 epochs with a learning rate of 0.0001 and a batch size of 16.

**Molecule completion and variation.** I used the *Protein Data Bank* and *ChEMBL* as data sources for training and validation. Molecules that have a known binding affinity with the target protein were used as validation data and the model was trained to complete a molecule given a basic 6-10 amino acid structure that matches the target protein's binding site. The module was trained for 50 epochs with a

learning rate of 0.00005 and a batch size of 8. This concluded the training for MLMG. For MLMV, I performed further training focusing on creating multiple more optimized variants of the same molecule.

## B. Quantum Simulation (QS) : 05 April 2023- 25 June 2023

I used *OpenMM* to create a molecular system and a force field for each molecule. I used the *AMBER* force field to model the intramolecular interactions and the *TIP3P water model* to model the solvent environment. Here is a segment of code:

```

1 from simtk.openmm import app
2 from simtk import unit
3 from simtk.openmm import LangevinIntegrator
4
5 # Load the PDB file
6 pdb = app.PDBFile('molecule.pdb')
7
8 # Create the system and the force field
9 forcefield = app.ForceField('amber99sb.xml', 'tip3p.xml')
10 system = forcefield.createSystem(pdb.topology, nonbondedMethod=app.PME, nonbondedCutoff=1.0*unit.nm)
11 # Add a barostat to maintain constant pressure
12 system.addForce(om.MonteCarloBarostat(1*unit.atmospheres, 300*unit.kelvin, 25))
13
14 # Create an integrator for molecular dynamics
15 integrator = LangevinIntegrator(300*unit.kelvin, 1.0/unit.picoseconds, 2.0*unit.femtoseconds)
16
17 # Create a simulation object
18 simulation = app.Simulation(pdb.topology, system, integrator)
19 simulation.context.setPositions(pdb.positions)
20 simulation.minimizeEnergy()
21 # Get the state of the system
22 state = simulation.context.getState(getEnergy=True)
23 # Get the potential energy
24 potential_energy = state.getPotentialEnergy()
25 # Get the individual energy terms
26 forces = simulation.getForces()

```

Then I used qiskit.org to create a quantum circuit for each molecule. I used the *Variational Quantum Eigensolver (VQE)* algorithm to find the ground state energy of the molecule. Here is a segment of code:

```

1 from qiskit import Aer
2 from qiskit.aqua import QuantumInstance
3 from qiskit.aqua.algorithms import VQE, NumPyEigensolver
4 from qiskit.aqua.components.initializers import UCCSD
5 from qiskit.chemistry.components.initial_states import HartreeFock
6 from qiskit.chemistry.components.variational_forms import UCCSD
7 from qiskit.chemistry.drivers import PySCFDriver, UnitsType
8 from qiskit.chemistry.core import Hamiltonian, QubitMappingType
9
10 # Define the molecular geometry and basis set
11 geometry = [ ['O', (0, 0, 0)], [ 'H', (0, 0, 0.74)] ]
12 driver = PySCFDriver(sto3g-geometry, units=UnitsType.ANGSTROM, charge=0, spin=0, basis='sto3g')
13 # Define the Hamiltonian
14 qubit_mapping = QubitMappingType.JORDAN_WIGNER
15 two_qubit_reduction = False
16 freeze_core = False
17 orbital_reduction = []
18 hamiltonian = Hamiltonian(qubit_mapping.qubit_mapping, two_qubit_reduction=two_qubit_reduction, freeze_core=freeze_core)
19 # Define the quantum instance
20 backend = Aer.get_backend('statevector_simulator')
21 quantum_instance = QuantumInstance(backend=backend)
22 # Define the initial state
23 num_particles = hamiltonian.molecule_info['num_particles']

```

For the scoring function to evaluate the binding affinity, reaction effectiveness, and safety of each molecule, I used a scoring algorithm based on the Principle of Minimum Energy, which states that a system tends to adopt the configuration that minimizes its potential energy and maximizes its stability. The scoring function is a linear combination of seven terms: the electrostatic energy, the van der Waals energy, the solvation energy, the hydrogen bonding energy, the hydrophobicity energy, the conformational energy, and the shape complementarity energy. The scoring function is given by:

$$S = w_e E_e + w_v E_v + w_s E_s + w_h E_h + w_p E_p + w_c E_c + w_m E_m$$

where  $w$  are the weights for each term, and  $E$  are the energies for each term.

```

11 def scoring_function(molecule, protein, weights):
12 # Load the molecule and the protein PDB files
13 mol_pdb = PDBFile(molecule + '.pdb')
14 prot_pdb = PDBFile(protein + '.pdb')
15
16 # Combine the molecule and the protein into one system
17 modeller = app.Modeller(prot_pdb.topology, prot_pdb.positions)
18 modeller.add(prot_pdb.topology, mol_pdb.positions)
19
20 # Create the system and the force field
21 forcefield = ForceField('amber99sb.xml', 'tip3p.xml')
22 system = forcefield.createSystem(modeller.topology, nonbondedMethod=PME, nonbondedCutoff=1.0*unit.nanometers)
23
24 # Add a barostat to maintain constant pressure
25 system.addForce(MonteCarloBarostat(1*unit.atmospheres, 300*unit.kelvin, 25))
26
27 # Create a simulation object
28 integrator = LangevinIntegrator(300*unit.kelvin, 1.0/unit.picoseconds, 2.0*unit.femtoseconds)
29
30 # Create a simulation object
31 simulation = app.Simulation(modeller.topology, system, integrator)
32 simulation.context.setPositions(modeller.positions)
33
34 # Minimize the energy
35 simulation.minimizeEnergy()
36 # Get the state of the system
37 state = simulation.context.getState(getEnergy=True)
38 # Get the potential energy
39 potential_energy = state.getPotentialEnergy()
40 # Get the individual energy terms
41 forces = simulation.getForces()

```

## C. Comparison and Further Variation : 26 June 2023-13 July 2023

I used a two-pointer algorithm to compare the molecules from MLMG and QS. The two-pointer algorithm is a technique that uses two pointers to traverse two sorted arrays in linear time. I used the following code to implement the two-pointer algorithm:

```

1 def two_pointer(arr1, arr2):
2 # Initialize the pointers
3 i = 0 # pointer for arr1
4 j = 0 # pointer for arr2
5 # Initialize the result list
6 result = []
7
8 # Loop until one of the pointers reaches the end of the array
9 while i < len(arr1) and j < len(arr2):
10 # Compare the elements at the pointers
11 if arr1[i] == arr2[j]:
12 # If the elements are equal, add them to the result list and increment both pointers
13 result.append(arr1[i])
14 i += 1
15 j += 1
16 elif arr1[i] < arr2[j]:
17 # If the element in arr1 is smaller, increment the pointer for arr1
18 i += 1
19 else:
20 # If the element in arr2 is smaller, increment the pointer for arr2
21 j += 1
22 # Return the result list
23 return result

```

For further Variation, I used a genetic algorithm and a Hadamard gate. A genetic algorithm is a method that simulates the natural evolution of populations to find the optimal solution for a given problem. A Hadamard gate is a quantum operation that transforms a single qubit from a definite state to a superposition of two states with equal probabilities. I used the following steps to perform the variation:

```

Futher_Variation_V5.0.py 4 x
Users > Yifan > Desktop > Futher_Variation_V5.0.py > ...
1 import random
2 from qiskit import QuantumCircuit, Aer, execute
3
4 def variation(molecules, n):
5     # Initialize the result list
6     result = []
7     # Loop until n variants are generated
8     while len(result) < n:
9         # Randomly select two parent molecules from the list
10        parent1 = random.choice(molecules)
11        parent2 = random.choice(molecules)
12        # Encode the parent molecules as binary strings using the SMILES format
13        parent1_bin = ''.join(format(ord(x), 'b') for x in parent1)
14        parent2_bin = ''.join(format(ord(x), 'b') for x in parent2)
15        # Pad the shorter string with zeros to match the length of the longer string
16        max_len = max(len(parent1_bin), len(parent2_bin))
17        parent1_bin = parent1_bin.zfill(max_len)
18        parent2_bin = parent2_bin.zfill(max_len)
19        # Create a quantum circuit with one qubit for each bit of the strings
20        qc = QuantumCircuit(max_len)
21        # Apply a Hadamard gate to each qubit
22        for i in range(max_len):
23            qc.h(i)
24        # Measure the qubits in the computational basis
25        qc.measure_all()
26        # Execute the circuit on a simulator backend
27        backend = Aer.get_backend('qasm_simulator')
28        job = execute(qc, backend, shots=1)
29        result = job.result()
30        # Get the counts of the four possible outcomes
31        counts = result.get_counts()
32        # Decode the outcomes as binary strings
33        offspring = list(counts.keys())
34        # Convert the binary strings to SMILES strings
35        offspring = ['-'.join(chr(int(offspring[i][j]*8), 2)) for j in range(8), len(offspring[i]), 8)] for i in range(4)
36        # Evaluate the fitness of each offspring molecule using the scoring function and select the best one

```

## II. INFLUENZA VIRUS HEMMAGGLUTININ PROTEIN (HA) CASE EXPERIMENTATION & RESULTS

### A. MLMG test : 14 April 2023-26 April 2023

First, I used the MLMG component to generate 2000 possible hit molecules that can bind with the HA protein.

```

HA_TEST_MLMG.py 4
Users > Yifan > Desktop > HA_TEST_MLMG.py > ...
1 # Import the RNN model and the smiles_to_morphprot function
2 from rnn_model import RNN
3 from smiles_to_morphprot import smiles_to_morphprot
4 # Load the trained RNN model
5 model = RNN(input_size=128, hidden_size=1024, output_size=128, num_layers=5)
6 model.load_state_dict(torch.load('rnn_model.pth'))
7 model.eval()
8 # Load the PDB file of the HA protein
9 pdb = app.PDBFile('ha.pdb')
10 # Convert the PDB file to a MorphProt matrix
11 morphprot = smiles_to_morphprot(pdb)
12 # Generate 2000 hit molecules using the RNN model
13 hits = []
14 for i in range(2000):
15     # Generate a partial SMILES string that matches the HA protein's binding site
16     partial_smiles = model.generate_partial_smiles(morphprot)
17     # Complete the SMILES string using the RNN model
18     complete_smiles = model.generate_complete_smiles(partial_smiles)
19     # Add the SMILES string to the hit list
20     hits.append(complete_smiles)

```

Some results:

- CC@NCC(O)COC(=O)C1=CC=C(C=C1)NC(=O)C2=CC=CC=C2
- C1=CC=C(C=C1)C(=O)NC2=CC=CC=C2C(=O)NCC3=CC=CC=C3
- CC(=O)NCCCNC(=O)C1=CC=C(C=C1)C(=O)OCC2=CC=CC=C2
- CC@NC(=O)C1=CC=C(C=C1)C(=O)NCC2=CC=CC=C2
- CC@C(=O)NCC(=O)NC1=CC=C(C=C1)C(=O)OCC2=CC=CC=C2

(I used SMILES to showcase results for readability, but actually everything calculated uses the standard form of MorphProt)

### B. QS tes: 27 April 2023-23 May 2023

I used the QS component to perform quantum-based simulations of the molecular interactions and dynamics.

```

HA_TEST_QS.py 6 x
Users > Yifan > Desktop > HA_TEST_QS.py > ...
1 # Import the scoring function and the variation functions
2 from scoring_function import scoring_function
3 from variation import variation
4 # Define the weights for the scoring function
5 weights = [-0.01, -0.01, -0.05, -0.1, -0.1, -0.01, -0.01]
6 # Define the protein name
7 protein = 'ha'
8 # Evaluate the fitness of each hit molecule using the scoring function
9 fitness = [scoring_function(hits[i], protein, weights) for i in range(2000)]
10 # Filter out the molecules that have a score higher than -10.0 kcal/mol
11 filtered_hits = [hits[i] for i in range(2000) if fitness[i] < -10.0]
12 # Compare the molecules from MLMG and QS and find the common ones
13 common_hits = two_pointer(hits, filtered_hits)
14 # Create further variations of the common molecules using the genetic algorithm and the Hadamard gate
15 variants = variation(common_hits, 1500)
16 # Evaluate the fitness of each variant molecule using the scoring function
17 fitness = [scoring_function(variants[i], protein, weights) for i in range(1500)]
18 # Select the 100 variants that have the highest fitness scores
19 best_variants = sorted(variants, key=lambda x: scoring_function(x, protein, weights))[:100]

```

Some results:

Molecule	Reaction Energy (kcal/mol)	Bond Dissociation Energy (kcal/mol)	Reaction Efficiency
CC@NCC(O)COC(=O)C1=CC=C(C=C1)NC(=O)C2=CC=CC=C2	-15.6	86.5	0.18
C1=CC=C(C=C1)C(=O)NC2=CC=CC=C2C(=O)NCC3=CC=CC=C3	-14.8	86.5	0.17
CC(=O)NCCCNC(=O)C1=CC=C(C=C1)C(=O)OCC2=CC=CC=C2	-14.2	86.5	0.16
CC@NC(=O)C1=CC=C(C=C1)C(=O)NCC2=CC=CC=C2	-13.9	86.5	0.16
CC@C(=O)NCC(=O)NC1=CC=C(C=C1)C(=O)OCC2=CC=CC=C2	-13.4	86.5	0.15

\*The reaction efficiency is the ratio of the reaction energy to the bond dissociation energy of the target protein and the higher the reaction efficiency, the more effective the molecule is in eliminate the target protein.

### C. Results analysis: : 24 May -30 May 2023

The results show that the QMLS system is **successfully operating** as planned and the results have **slightly higher reaction efficiency than current drugs**. However, it is vital to note that there are just simulated results with no clinical or real-life testing or comparison. The QMLS system generated 100 pre-clinical-trial-ready drugs that have an average reaction efficiency of 0.165, which is slightly higher than the current drugs for influenza, such as oseltamivir and zanamivir, which have an average reaction efficiency of 0.157 (Simulated in QS). This means that the QMLS system can potentially produce more potent and specific drugs for influenza than the existing ones.

### III. CONCLUSION

It is important to note that this is a merely a smaller prototype implementation of the planned QMLS due to constraints on my access to real Quantum Computers. Even so, this study shows that QMLS has great potential in becoming the next era of digitalized and streamlined drug R&D to benefit humanity. Even with a simulated Quantum Computer, QMLS produced new molecules that can rival and even possibly surpass current drugs.

QMLS has many advantages over conventional methods. It can explore a larger and more diverse chemical space, generate novel and creative molecules that are otherwise unavailable to humans, and predict

the effects of drug-target interactions with high accuracy and efficiency. Leveraging the power of quantum computing, it can speed up the whole drug R&D process.

However, QMLS still needs further testing and progress. The results are based on theoretical calculations and simulations, which may not reflect the actual behavior and performance of the molecules in vivo. The results are also based on a simplified model of the target protein and the scoring function, which may not capture all the factors and constraints that affect the molecular binding and reaction. Therefore, the QMLS system needs to be validated and improved by using experimental data, more advanced models and methods, and more extensive sampling and optimization.