

# Roadmap for a Postdoc in Computational Drug Discovery

## Table of Contents

- [Introduction](#)
- [Phase 1 – Foundational Knowledge and Hands-On Skills](#)
- [Week 1-2: Programming and ML Foundations](#)
- [Week 3-4: Cheminformatics & Molecular Modeling Basics](#)
- [Week 5-6: Deep Learning for Molecular Data](#)
- [Week 7-8: Quantum Chemistry Fundamentals](#)
- [Week 9-10: Quantum Computing Fundamentals](#)
- [Phase 2 – Mini Project and Hypothesis Generation](#)
- [Selecting a Mini Project](#)
- [Implementing the Mini Project \(4-week plan\)](#)
- [Phase 3 – Publication-Level Project Planning](#)
- [Literature Review and Gap Analysis](#)
- [Formulating a Research Proposal](#)
- [Planning and Next Steps](#)
- [Additional Resources and Tips](#)

## Introduction

Preparing for a postdoctoral position in **computational drug discovery** requires gaining expertise across multiple domains – from machine learning (ML) and data science to quantum chemistry, quantum computing, and molecular modeling. Modern drug discovery is increasingly interdisciplinary: for example, ML models can rapidly predict molecular properties and drug-target interactions, accelerating discovery <sup>1</sup>. However, purely data-driven models often ignore quantum-mechanical details of molecules <sup>2</sup>, and emerging techniques like **quantum computing** promise to enhance simulations for drug design <sup>3</sup>. This roadmap provides a comprehensive, phased learning plan to build these skills in a practical way, with **4 hours of focused learning per day**.

**How to use this guide:** The roadmap is divided into three phases, each with specific goals: - **Phase 1** establishes foundational knowledge in each area (with hands-on coding practice).

- **Phase 2** involves undertaking a small independent project to apply and integrate the skills, helping you generate research hypotheses.

- **Phase 3** guides you in formulating a publication-level research plan.

Each phase is broken into weekly sub-goals (based on a typical learning curve and ~20 hours/week effort). We include curated resources — tutorials, libraries, GitHub projects, and datasets — to facilitate a *code-first learning* approach. You are encouraged to use **GitHub Copilot** in VS Code as you work through coding exercises and projects, allowing you to iterate faster by auto-generating boilerplate code and discovering

functions (always double-check Copilot's suggestions for accuracy!). Throughout the guide, references are provided for deeper reading on concepts or tools (in the format `[tsource]` ).

By following this structured plan, you will incrementally build proficiency in: **(1)** core machine learning techniques and their application to molecular data, **(2)** classical molecular modeling methods (cheminformatics, molecular simulations), **(3)** quantum chemistry fundamentals, **(4)** the basics of quantum computing and its emerging role in drug discovery. The ultimate aim is to help you confidently discuss and pursue research at the intersection of these areas, design your own hypotheses, and hit the ground running in a postdoc interview or early research phase.

## Phase 1 – Foundational Knowledge and Hands-On Skills

**Goal:** Develop a strong foundation in the essential domains (programming/ML, cheminformatics & molecular modeling, quantum chemistry, and quantum computing) and gain hands-on experience with key tools. By the end of Phase 1, you should be comfortable writing code (preferably Python) to solve basic problems in each domain and understand how these areas connect to real-world drug discovery challenges.

Each sub-section below corresponds to roughly **2 weeks (10 weekdays, ~40 hours)** of work. The timeline can be adjusted based on your background; for example, if you already have ML experience, you can shorten the ML basics section. Aim to spend about 4 hours a day following the suggested topics and exercises. Short, frequent coding exercises will reinforce learning — try to integrate Copilot to assist with syntax and boilerplate, but make sure you understand the code you write.

### Week 1-2: Programming and ML Foundations

**Objective:** Ensure you have the programming and basic machine learning skills needed for more advanced topics. If you are already proficient in Python and familiar with basic ML, you can skim this section or use it as a quick refresher.

- **Python Refresh:** Dedicate the first few days to solidify your Python skills, especially if your background is primarily chemistry/biology. Focus on writing and running scripts in VS Code. Key topics:
  - Python syntax and best practices (functions, classes, using libraries).
  - Scientific Python stack: **NumPy** (numerical arrays), **pandas** (dataframes), and **matplotlib/seaborn** (for basic plotting).
  - If needed, follow a quick tutorial or **Jupyter Notebook** introduction on using these libraries for data manipulation. (Use Copilot to suggest code snippets, e.g., ask it to show how to load a CSV and plot data.)
- *Resource: "Python Data Science Handbook" by J. VanderPlas (freely available) has chapters on NumPy, pandas, matplotlib – skim relevant sections for any gaps.*
- **Machine Learning Basics:** In the remainder of Week 1 and into Week 2, learn core ML concepts and practice on simple datasets. Even if your ultimate focus is deep learning, understanding fundamentals (train/test split, evaluation metrics, overfitting, etc.) is crucial. Key tasks:

- Learn the basics of **scikit-learn** (a Python ML library): try a simple regression (predict molecule solubility from a property) or classification task with scikit-learn's built-in datasets (e.g., iris or a toy dataset) to familiarize yourself with model training/prediction workflow.
- Concepts to cover: linear regression, logistic regression, basic clustering (k-means), and cross-validation. Focus on understanding how to evaluate models (accuracy, RMSE, ROC-AUC, etc.).
- **Resource Links:** The scikit-learn official tutorials and examples are excellent for quick hands-on learning. For instance, the scikit-learn tutorial on model training and evaluation can guide you through a basic project (see scikit-learn documentation [\[1\]](#)). If you prefer a guided course, Coursera's *Machine Learning* by Andrew Ng or *fast.ai*'s free course can be done in parallel with a focus on concepts (but don't get bogged down in theory – prioritize practical understanding).
- **Coding with Copilot:** As an exercise, write a small script to train a simple neural network on dummy data using **PyTorch** or **TensorFlow** (to anticipate deep learning usage later). Use Copilot to help with the boilerplate (e.g., network architecture definition). This will get you comfortable with the syntax of a deep learning library.

By the end of Week 2, you should be able to load data, train a basic ML model, and understand results. This foundation will be essential when working with molecular data in subsequent weeks.

## Week 3-4: Cheminformatics & Molecular Modeling Basics

**Objective:** Learn how to represent and manipulate molecular data on a computer (cheminformatics), and grasp the basics of classical molecular modeling techniques used in drug discovery (such as docking and molecular dynamics). You will practice using key libraries and tools, especially **RDKit**, a fundamental cheminformatics toolkit.

- **Chemical Data Representation:** Understand common ways to represent molecules digitally: **SMILES** strings, molecular graphs, and 3D structures (PDB, SDF files). Learn how to go from a chemical identifier to a machine-readable format and vice versa.
- Install **RDKit** in your Python environment (it's easiest via Anaconda). RDKit is a widely used open-source toolkit for cheminformatics (the intersection of chemistry and computer science) [\[4\]](#). It provides functionality to read/write molecular files, compute molecular descriptors/fingerprints, and perform substructure searches. RDKit also has basic machine-learning tools and integrates with other libraries [\[5\]](#) [\[6\]](#).
- **Hands-On:** Use RDKit (with Copilot's help for function hints) to perform simple tasks: parse a SMILES string into a molecule object, compute some properties (molecular weight, rotatable bonds, etc.), and generate a 2D depiction of a molecule. For instance, you can follow an RDKit beginner tutorial like "**RDKit for Beginners**" [\[4\]](#) which introduces fundamental operations.
- Learn about **molecular descriptors** and **fingerprints** (e.g., RDKit's Morgan fingerprint). These are essential for ML on molecular data. Try computing a fingerprint for a set of molecules and understand how it can be used as input features for ML models [\[6\]](#).
- **Classical Molecular Modeling Concepts:** In parallel with cheminformatics, get acquainted with how drug discovery uses physics-based modeling:
- **Molecular Docking:** Understand the idea of docking – predicting how a small molecule (ligand) binds to a protein (receptor) and scoring the fit. While doing a full docking workflow requires a receptor structure and tools like AutoDock Vina, you can conceptually learn it this week. *Suggested*

*activity:* read a tutorial or watch a video on how docking software works (for example, the AutoDock Vina documentation and basic tutorials <sup>7</sup> for command-line docking). You don't need to master running docking now, but know the inputs/outputs (protein PDB, ligand, scoring function). This will help later when considering ML improvements to scoring or when discussing structure-based design.

- **Molecular Dynamics (MD):** Learn what MD simulations are – essentially, using physics (force fields) to simulate the motion of molecules over time. MD is used to study protein-ligand interactions, conformational changes, etc. For practicality, try a small simulation using an accessible tool:

**OpenMM** is a Python-based, high-performance toolkit for molecular simulation <sup>8</sup>. Install OpenMM and follow a basic tutorial (OpenMM's official documentation or examples). For instance, OpenMM's user guide has a simple example of simulating a few water molecules or a small protein.

- *Hands-On:* Using Copilot, write a Python script with OpenMM to simulate a simple system (even just a box of water or a small molecule in vacuum). This will teach you how to set up a system, apply a force field, and run a short simulation. Focus on the workflow (define system, integrator, run simulation loop) rather than the complex theory.
- Understand force fields at a high level (e.g., MMFF or AMBER force fields: parameters that approximate molecular forces). No need to dive too deep into parameterization, but know that MD relies on classical approximations, whereas quantum chemistry (next section) deals with first-principles calculations.

After Week 4, you should be comfortable handling molecular structures in code, computing basic chemical features, and you'll have a conceptual grasp of docking and molecular dynamics. You'll see how ML can complement these: e.g. ML models can predict binding affinities faster than docking in some cases, or learn from MD simulation data to predict protein flexibility.

## Week 5-6: Deep Learning for Molecular Data

**Objective:** Apply machine learning skills specifically to molecular problems. Focus on deep learning approaches (neural networks) which have become state-of-the-art in many drug discovery tasks <sup>9</sup>. Gain experience with models and libraries specialized for molecular data.

- **Learning from Molecular Datasets:** Become familiar with public datasets relevant to drug discovery and how to load them for ML. A great resource is **MoleculeNet**, a large benchmark suite of datasets for molecular machine learning <sup>10</sup>. MoleculeNet includes various tasks such as quantum chemistry datasets (like QM9 for molecular energies), physical chemistry (ESOL for solubility), biophysics (PDBbind for protein-ligand binding), and physiology (Tox21 for toxicity) <sup>11</sup>.
- *Hands-On:* Pick one dataset from MoleculeNet (for example, **ESOL** – a small dataset of molecules with known aqueous solubility, or **HIV** dataset for a bioactivity classification task). Using **DeepChem** (an open-source library that provides ready access to MoleculeNet datasets <sup>10</sup>), load the dataset. DeepChem is a convenient toolkit that “democratizes deep learning in drug discovery, materials science, quantum chemistry, and biology” <sup>12</sup>. It has functions to fetch datasets and split them into train/test.
- Once loaded, inspect the data (with pandas or directly via DeepChem) to understand features and labels. If you have molecular SMILES, consider how to featurize them (DeepChem can provide featurizers like ConvMolFeaturizer for graph conv nets, or you can use RDKit descriptors as features).

- **Deep Learning Models for Molecules:** Study common neural network architectures used in molecular modeling: **Graph Neural Networks (GNNs)** and **sequence-based models**. In recent years, GNNs (like graph convolutional networks) have shown outstanding performance for predicting molecular properties <sup>9</sup>, because molecules can be naturally represented as graphs (atoms as nodes, bonds as edges). Sequence models (like transformers) can be applied to SMILES strings treating them as text sequences.
- **Resources:** Read an overview or introductory article on GNNs in chemistry – for example, the MoleculeNet paper or a blog summary of it notes that graph convolutional networks outperformed many other methods on molecular benchmarks <sup>9</sup>. Also, explore DeepChem’s examples of graph models or the documentation of **Chemprop** (a popular GNN-based model for molecular property prediction).
- **Hands-On (Code):** Using either DeepChem or **PyTorch Geometric** (a library for GNNs), train a simple graph neural network on the dataset you picked. DeepChem has high-level APIs (e.g., `GraphConvModel`) that can train a graph conv net for regression or classification. Alternatively, if you want to be more hands-on, use PyTorch Geometric: it has tutorials on building a GNN for a molecular dataset (you might need to convert SMILES to graph input; RDKit can help generate the adjacency and feature matrices).
  - Use Copilot to speed up writing model definitions or data processing code. For instance, ask Copilot to help create a PyTorch dataset class for molecules. Ensure you evaluate the model on a test set and get a feel for performance metrics (e.g., if regression, what MAE/RMSE did you get for ESOL solubility?).
- **Alternative/Additional Task:** Try a simple **generative modeling** task. For example, use the **Jupyter Notebook tutorials** in DeepChem (or standalone projects on GitHub) that demonstrate training a variational autoencoder (VAE) or GAN to generate new SMILES molecules. This will expose you to the idea of de novo drug design with deep learning <sup>13</sup> (one of the key tasks in AI-driven drug discovery is generating novel chemical structures with desired properties).
- **Model Interpretation and Validation:** As you work with ML models, consider how to validate them in a chemistry context. Techniques like **cross-validation** on scaffold splits (to ensure your model is truly generalizing to novel scaffolds, not just similar compounds) are important in drug discovery. Read about common pitfalls like data leakage in screening datasets <sup>14</sup> <sup>13</sup>. At minimum, be aware that good performance on test sets is just a start – models should be assessed on diverse compounds and ideally prospective experimental data.

After Week 6, you should have concrete experience applying deep learning to molecular problems. You will have used at least one specialized library (DeepChem, PyTorch Geometric, or similar) and handled a real dataset from drug discovery. You’ve also seen how ML ties back to drug discovery tasks: for example, predicting molecular properties (solubility, bioactivity) or generating new molecules <sup>13</sup>. This prepares you for deeper topics and for creating your own projects in Phase 2.

## Week 7-8: Quantum Chemistry Fundamentals

**Objective:** Build a basic understanding of quantum chemistry methods and why they matter in drug discovery. Get hands-on experience with quantum chemistry calculations on small molecules, and learn how these first-principles techniques can complement ML.

- **Quantum Chemistry Concepts:** Learn the foundations of electronic structure theory, focusing on what a computational drug discovery researcher should know:
  - Key concepts: **Schrödinger equation** for molecules, what is a molecular **wavefunction**, the idea of **potential energy surface**.
  - Methods hierarchy: **Hartree-Fock (HF)** as a starting point, and **Density Functional Theory (DFT)** as a workhorse for calculating molecular properties (e.g., optimized geometries, energies) from first principles. Also be aware of post-HF methods like MP2, Coupled Cluster, etc., though you need not master them for now.
- Practical understanding: quantum chemistry methods provide accurate calculations of molecular properties (like binding energies, orbital energies, etc.) but are computationally expensive. This is why ML can be useful to approximate some of these results more quickly <sup>1</sup> <sup>2</sup>. Conversely, quantum chemistry can generate high-quality data for ML model training (e.g., quantum-calculated energies of molecules used to train a ML model, as in the QM9 dataset).
- **Hands-On with Quantum Chemistry Software:** Use a beginner-friendly quantum chemistry tool to perform a simple calculation: for example, optimize the geometry of a small molecule and calculate its energy. Two accessible routes:
  - **Psi4** – an open-source quantum chemistry package with a Python API. Psi4 can perform HF, DFT, etc. You can write a Python script (or use a Jupyter notebook) to run a DFT calculation for a molecule (e.g., water or ethanol) and get properties.
  - **Psi4NumPy** – a project that combines Psi4 and NumPy for educational purposes. It provides reference implementations of quantum algorithms in Python <sup>15</sup> <sup>16</sup>. Notably, it offers Jupyter notebook tutorials that intermix theory and code, bridging the gap between equations and implementation <sup>17</sup>.
  - *Hands-On Activity:* Follow one of the Psi4/Psi4NumPy tutorials. For instance, a tutorial on implementing Hartree-Fock from scratch in Python (with Psi4 providing integrals) will deepen your understanding of how quantum chemistry methods work under the hood <sup>17</sup>. By coding parts of the HF procedure (with guidance from Copilot or the provided tutorial steps), you'll reinforce concepts like basis sets, self-consistent field iterations, etc. This is time-consuming, so even if you don't complete the entire implementation, go through enough to appreciate the workflow.
- Alternatively, use Psi4 at a higher level: write a Python input script to run a simple DFT energy calculation on a molecule. Psi4's documentation has examples for this. Use Copilot to help with calling the Psi4 API (e.g., `psi4.geometry()` to define molecules, `psi4.optimize()` to optimize geometry with a chosen method/basis).
- **Quantum Chemistry in Drug Discovery Context:** Read about how quantum chemistry is used in practice for drug discovery. For example, computing accurate energies for enzyme reaction mechanisms, or polarization effects in binding, etc. A current trend is integrating quantum-chemical insights into ML models to improve their accuracy <sup>2</sup> <sup>18</sup>. Recent research even suggests new

molecular representations that include quantum chemical interactions to boost ML predictions <sup>1</sup> <sup>19</sup>. By the end of this week, make sure you can answer: *Why is quantum chemistry relevant to drug discovery, and what are its limits?* (For instance, you might note that while DFT can give accurate energies, it's too slow for screening thousands of compounds – hence the need for ML approximations and potentially quantum computing in the future.)

After Week 8, you will have demystified a bit of quantum chemistry. You've run a quantum calculation yourself and understood the basics of how molecules are treated quantum-mechanically. This will also allow you to converse with computational chemists and understand literature that mentions things like “DFT-optimized conformations” or “QM/MM calculations” etc. Moreover, this sets the stage for understanding the promise of **quantum computing** in chemistry.

## Week 9-10: Quantum Computing Fundamentals

**Objective:** Learn the basics of quantum computing, with an emphasis on its applications in chemistry and drug discovery. Get hands-on experience with quantum programming frameworks (like Qiskit) to understand how quantum algorithms can be used for molecular problems. Recognize the current limitations (NISQ era issues) and the potential of hybrid quantum-classical approaches.

- **Quantum Computing 101:** Devote a few days to understanding what a quantum computer is and how it differs from classical computers. Key concepts: qubits (quantum bits), superposition, entanglement, quantum gates, and circuits. You don't need to become a quantum algorithms expert, but grasp the *computational* aspect: how qubits can represent and manipulate information.
- **Learning Resource:** IBM's **Qiskit Textbook** (available online for free) provides an excellent introduction to quantum computing with practical examples. Read the introductory chapters that explain qubits and simple quantum circuits. When you reach the section on algorithms, focus on those relevant to chemistry, like the **Variational Quantum Eigensolver (VQE)**. VQE is a quantum algorithm to find minimum eigenvalues (e.g., molecular ground state energy) and is seen as a promising approach for quantum chemistry on near-term quantum hardware.
- Understand that we are currently in the **NISQ era** (Noisy Intermediate-Scale Quantum): today's quantum computers have limited qubits and are noisy <sup>20</sup> <sup>21</sup>. This means most applications, including chemistry, require hybrid approaches (quantum + classical) to be practical <sup>22</sup> <sup>23</sup>. Keep this in mind as a talking point in interviews – it shows you are realistic about quantum tech: it's promising but not magic (yet).
- **Quantum Programming with Qiskit:** Install **Qiskit**, a popular Python framework for quantum computing. In VS Code, set up a notebook or script and use Copilot to help with basic tasks (it can suggest Qiskit code structure). Try the following:
  - Write a simple Qiskit program to create and run a quantum circuit (for example, a Bell state or a basic 2-qubit circuit). This is just to get familiar with Qiskit's syntax for defining circuits, running simulations, and getting results.
  - Move to a chemistry-related quantum example: Qiskit has an application module called **Qiskit Nature** for problems in chemistry. Follow a tutorial from the Qiskit Nature documentation <sup>24</sup>, such as computing the ground state energy of a small molecule (like H<sub>2</sub>) using VQE. Qiskit Nature provides convenient tools to map a molecular Hamiltonian to a qubit Hamiltonian. The tutorial will walk you

through setting up the molecule, choosing a quantum algorithm (VQE with a particular ansatz, e.g., UCCSD), and running it on a simulator.

- *Note:* If Qiskit setup or VQE feels too heavy, an alternative is to explore **PennyLane** (another library for quantum ML) which has some quantum chemistry demos, or read through a high-level description of the VQE process without coding it fully. However, attempting the Qiskit tutorial is valuable hands-on experience, even if you only use a simulator backend (which is fine).
- **Quantum Computing for Drug Discovery – Context:** Spend some time researching or reading recent developments in applying quantum computing to drug discovery. It's an evolving field; being knowledgeable here will set you apart. Key points to understand or mention in interviews:
  - So far, quantum computing in drug discovery has mostly been **proof-of-concept** – e.g., estimating energies of small molecules, or small optimization problems, rather than directly finding new drugs <sup>3</sup>.
  - There was even a 2023 challenge (QCDDC'23) where teams used quantum algorithms to compute molecular ground state energies (OH<sup>+</sup> molecule) with integration of machine learning to cope with hardware noise <sup>25</sup> <sup>26</sup>. The outcome highlighted the importance of hybrid classical-quantum methods and error mitigation <sup>20</sup> <sup>23</sup>.
  - Companies and research labs are actively exploring quantum algorithms for tasks like protein folding, generative chemistry, and precise simulation of reaction mechanisms <sup>27</sup> <sup>28</sup>. One study even developed a hybrid quantum computing pipeline for drug design problems (like calculating prodrug activation free energies) to test quantum approaches on realistic tasks <sup>29</sup> <sup>30</sup>.
  - **Bottom line:** Quantum computing is a frontier area. Emphasize that you understand both its potential (speed-ups for certain problems, more accurate quantum-mechanical simulations) and its current limitations (noisy hardware, limited qubits, need for classical assistance) <sup>22</sup> <sup>31</sup>. Showing a balanced view will demonstrate critical thinking.

By the end of Week 10, you'll have written and executed basic quantum computing code, and you'll know how quantum computing ties into the computational chemistry sphere. You should feel confident explaining, for example, what VQE does and why it's relevant, or discussing a "hybrid classical-quantum workflow" for drug discovery. This knowledge, combined with your ML and modeling skills, completes the foundational toolkit.

## Phase 2 – Mini Project and Hypothesis Generation

**Goal:** Transition from guided learning to **independent exploration**. In this phase, you will undertake a mini research project that combines two or more skills from Phase 1. The project should be manageable (aim to complete in ~4 weeks of part-time effort), but complex enough to simulate the process of formulating and testing a research hypothesis. By doing this, you'll gain confidence in your ability to tackle novel problems, which is crucial for interviews and for hitting the ground running in a postdoc.



## Selecting a Mini Project

Choose a project that genuinely interests you and lies at the intersection of the domains you've learned. Here are a few **example project ideas** to illustrate the possibilities (feel free to come up with your own or modify these):

- **Project Idea 1: ML-Enhanced Molecular Property Prediction** – *Using quantum chemistry data to train an ML model.* For example, take a dataset of small molecules with experimentally measured properties (solubility, logP, etc.) or quantum-calculated properties (like the QM9 dataset for molecular energy levels). Train a **graph neural network** (or other ML model) to predict one of these properties. Hypothesis example: *"A graph neural network can learn to predict DFT-calculated dipole moments of molecules with  $\langle X \rangle$  accuracy, potentially much faster than doing new DFT calculations."*
- **Hands-on:** You could use the QM9 dataset (available via MoleculeNet/DeepChem) which contains DFT-calculated properties <sup>10</sup>. Train a model (e.g., using DeepChem's `GraphConvModel` or a custom PyTorch model) to predict a property like molecular energy or dipole moment. Evaluate its performance against a hold-out test set.
- **Learnings:** This project reinforces ML skills and lets you articulate how ML can approximate quantum chemistry – a great talking point since it shows synergy between domains (and is a current research direction <sup>2</sup>). If you want to extend it, you could compare different feature representations (e.g., classic descriptors vs. learned graph features) to see which yields better accuracy.
- **Project Idea 2: Structure-Based Drug Screening with AI** – *Integrating molecular docking with ML re-scoring.* Take a protein target of interest (perhaps one with a known set of ligands, like an enzyme from the PDB database). Use a docking tool to dock a small library of compounds (you can use a subset of, say, the ZINC database or some provided set). Then, train a simple ML model to predict the docking score or binding outcome from the ligand structure. Hypothesis example: *"A machine learning model can learn to predict docking scores for a protein-ligand pair, potentially correcting some errors in the docking scoring function."*
- **Hands-on:** This is a bit more involved. You might use a tool like AutoDock Vina to dock, or if installing that is a hassle, use already published docking results (some papers provide datasets of docking scores). Alternatively, use a simpler approach: generate 3D conformations with RDKit and use a simplified scoring (like shape or pharmacophore overlap) as a proxy, then see if ML can predict that.
- **Learnings:** This project connects classical modeling (docking) with ML. It gives you talking points about how AI can assist virtual screening (which is a key task in drug discovery <sup>13</sup>). It also teaches some practical issues: e.g., the need for negative data (decoys) in training binding predictors, etc.
- **Project Idea 3: Quantum-Inspired Machine Learning** – *Incorporating quantum chemical features into an ML model.* Building on the idea that quantum details improve ML models <sup>19</sup>, design a small experiment: take a property prediction task and include a feature computed from quantum chemistry. For instance, predict acidity (pKa) of molecules using an ML model, with input features being a mix of classical descriptors plus one quantum descriptor (like the energy of the highest occupied molecular orbital, HOMO, computed via a semi-quantum method for each molecule). Hypothesis: *"Adding quantum chemical features (like HOMO energy) to the descriptor set will improve the ML model's accuracy in predicting pKa compared to using only classical descriptors."*

- **Hands-on:** Select ~50 molecules with known pKa (or another property – could be solubility, etc.). For each, use a quantum chemistry tool (Psi4 or even an online server or simpler semi-empirical tool) to compute a property like HOMO energy or dipole moment. Construct a feature matrix that includes this quantum feature plus others (like MW, logP from RDKit). Train a regression model (even a simple one like random forest or neural network) to predict the property. Compare performance with/without the quantum feature.
- **Learnings:** This project is a mini-version of actual research efforts to combine quantum mechanics and ML <sup>2</sup> <sup>19</sup>. You'll learn about feature importance and get experience combining disparate tools (running a calc, then feeding into ML).
- **Project Idea 4: Exploring Quantum Computing for a Toy Problem – Applying a quantum algorithm and comparing to classical.** If you're particularly excited about quantum computing, you could do a small project like: use Qiskit to estimate the energy of a simple molecule ( $H_2$  or LiH) using VQE, and compare it to a classical calculation (from Phase 1). Hypothesis: *"A hybrid quantum-classical algorithm (VQE) can compute the ground-state energy of  $H_2$  within chemical accuracy given idealized conditions, showcasing potential for quantum advantage in the future."*
- **Hands-on:** Follow a Qiskit VQE tutorial (which you may have done partly in Phase 1). Record the energy output and compare to the known true value (from classical methods). You can then experiment: how does the result change if you simulate some noise or use fewer circuit iterations, etc., to mimic real conditions? Document what challenges you encounter.
- **Learnings:** This will deepen your practical understanding of quantum algorithms and give you a nuanced perspective on what's needed to make them competitive. It's a good discussion point for interviews in groups working on quantum computing for chemistry.

These ideas are just starting points. **When selecting your project, keep it bounded:** something you can show results for in a few weeks. Importantly, **formulate a clear hypothesis** or question at the start (as in examples above). This habit of articulating a hypothesis will train you to think like a researcher and will impress interviewers. Even if your mini project is simpler (e.g., "I reimplemented a known algorithm and got it to work"), frame it as answering a question or demonstrating a concept.

## Implementing the Mini Project (4-week plan)

Once you've chosen a project, structure the implementation roughly as follows (adjust as needed):

- **Week 1: Project Planning and Setup**
- Refine your project scope and make sure you have all needed tools installed. Write a one-page project plan: background of the problem, your hypothesis, and the steps you will take. This is basically a mini-proposal. It helps clarify your thoughts and can be something you share with mentors for feedback.
- Gather necessary data or resources. For example, if your project involves a dataset, download and inspect it. If it involves running simulations or calculations, prepare a small test case to ensure your environment (libraries, hardware) is working.
- **Literature check:** Do a brief literature search to see if others have attempted something similar. For instance, if you're predicting solubility with ML, find a reference model or paper to know what accuracy is expected. A quick check on arXiv or Google Scholar with relevant keywords can prevent

you from reinventing the wheel and provides context <sup>14</sup> <sup>32</sup> . Keep notes of 1-2 key references – you might cite them in your project report.

#### • **Week 2: Initial Implementation**

- Start coding the core parts of your project. If it's an ML model, set up your training pipeline (data loading, model definition, training loop). If it's a computational experiment, write the scripts to run simulations or quantum calculations as needed.
- Use **GitHub Copilot** actively to accelerate development. For example, if you need to write a function to compute a particular descriptor or to parse a file, let Copilot suggest it and then adjust as necessary. Copilot can also help with API usage (for instance, how to call a specific DeepChem function) based on its training. Always test small components as you write them.
- Aim to get a **minimum viable product** working this week – e.g., a model that trains (even if performance is bad), or one round of computation done. Don't worry about perfect results yet. The goal is to have the pipeline in place.

#### • **Week 3: Testing and Iteration**

- Now, refine and iterate. Run your full experiment and start analyzing results. Did your model converge? If not, adjust hyperparameters or fix bugs. Did your simulation produce output? If it crashed or was unstable, debug the setup (maybe the time step was too large in MD, or the basis set too high-level in quantum calc – adjust accordingly).
- This is the time to also perform any comparisons or ablation studies relevant to your hypothesis. For instance, if you want to see the effect of a quantum feature in an ML model, train two versions (with and without it) and compare metrics. Keep track of these results.
- If results are negative or inconclusive, think scientifically: does it reject your hypothesis, or do you need to change approach? It's okay if your mini project doesn't "succeed" in a traditional sense (e.g., maybe your ML model didn't improve with that quantum feature) – that's still a valuable finding and learning experience you can articulate.

#### • **Week 4: Conclusion and Write-up**

- Wrap up the project by producing a short report or presentation. This could be a Jupyter Notebook that contains your analysis, or a markdown/PDF report with the key sections: **Introduction, Methods, Results, Discussion**. Writing this up is extremely beneficial – it forces you to clarify what you did and learned. In a postdoc interview, you might even present slides on this project, so consider preparing 2-3 slides summarizing the problem and findings.
- In your discussion, include what you would do next if you had more time. This shows you can think beyond the immediate project. For example: "The model showed moderate accuracy in predicting solubility. With more time, I would try a transformer-based approach or gather more data on larger molecules to improve generalization." This forward-thinking aligns with designing publication-level projects (Phase 3).
- Finally, consider sharing your code on GitHub. Even if it's a small project, having it on GitHub shows you are serious about open science and allows others (like your potential postdoc advisor) to see your work. Document it well (a good README, clear instructions to run) – this is good practice for research code. You can mention in interviews that you have a project repository available.

By the end of Phase 2, you will have **hands-on project experience** that ties together multiple domains. This is invaluable for interviews – you can discuss challenges you faced, how you integrated techniques, and what you learned. It demonstrates initiative and the ability to drive a small project independently, which are qualities of a successful postdoc. Moreover, you might even stumble upon a novel idea during this mini project that can be expanded into a full research project (which leads into Phase 3).

## Phase 3 – Publication-Level Project Planning

**Goal:** Synthesize your knowledge and experience to formulate one or more potential research projects that could be pursued during your postdoc and lead to a publication. In this phase, you'll act like you're preparing a research proposal: identifying gaps in current research, proposing innovative approaches (perhaps combining ML, quantum chem, etc.), and planning how to execute and evaluate the project. This will not only prepare you for discussing future work in interviews but will also give you a head start in your postdoc (having a clear plan to execute).

### Literature Review and Gap Analysis

A strong research idea is built on understanding what's been done and what open problems remain. Start with a focused **literature review** in the area that interests you most at the intersection of these fields. For example, if you are fascinated by quantum computing for drug discovery, gather recent papers or reviews on that topic. If your interest leans more to applying ML to improve molecular simulations, focus your reading there.

- **Identify Key Publications:** Find 5-10 recent publications (within last ~3 years) in the interdisciplinary space you're targeting. Some tips:
- Look for **review articles or perspective pieces** – these give overviews of subfields. For example, a *2022 perspective on AI in drug discovery* or a *recent review on quantum computing applications in pharma*. These often highlight challenges and future directions. A viewpoint noted that while ML has made progress in drug discovery, challenges like limited data and model interpretability remain <sup>14</sup>. Similarly, a review might note that quantum computing is mostly in proof-of-concept stage for chemistry, with scalability and noise being key issues <sup>3</sup> <sup>29</sup>. Such stated challenges can be opportunities for you to address.
- Use academic search engines: Google Scholar with keywords (e.g., “machine learning drug discovery 2024 review”, “quantum computing drug discovery proof-of-concept”). Also check arXiv for the latest preprints (often a good source for very up-to-date ideas).
- Don't forget to search the intersection: e.g., “machine learning quantum chemistry drug design” might yield papers where ML is used to predict quantum chem outcomes for druglike molecules <sup>1</sup> <sup>19</sup>, or “quantum machine learning drug discovery” for works that bring quantum computing and ML together <sup>22</sup>.
- **Analyze for Gaps:** As you read, note down: What are the open problems or limitations mentioned? What do authors say “remains an unsolved challenge” or “future work is needed”? These statements often directly point to research opportunities. For instance:

- If a paper on ML for drug discovery mentions that *lack of interpretability* is an issue, maybe a project could focus on making models more explainable for chemists (e.g., highlighting which substructures lead a model to predict toxicity).
- If a review on quantum computing says *scaling to larger molecules is hard due to noise*, perhaps a project idea is to develop a hybrid algorithm that uses ML to reduce the quantum circuit depth needed (some teams are already looking at this <sup>33</sup>, but there is room for innovation).
- If a molecular simulation study notes that *force fields are imperfect for certain interactions*, maybe you can propose an ML correction term trained on quantum chemistry data (some work exists on ML force fields, but you could carve a niche in, say, improving protein-ligand force fields with ML).
- **Refine Your Interest:** Based on this survey, decide on 1-2 areas you're most excited about. It could be something like *"Generative models for covalent drug design"* or *"Quantum-enhanced prediction of protein-ligand binding"* – whatever aligns with your passion and what you think the postdoc lab might be interested in. It's fine if it's ambitious; a postdoc project often is. The key is that you now understand the context and can articulate why it matters.

Throughout this literature review, keep organized notes. You might create a simple table or mind map of identified gaps and potential approaches to address them. Also note any datasets or open-source tools mentioned in papers that could be useful for your project (e.g., a new dataset published in 2023 for protein-ligand complexes, or a new quantum computing library). These can all become parts of your proposal.

## Formulating a Research Proposal

Now, take your refined idea and formulate a **mini research proposal**. This doesn't have to be a formal document unless you want to use it for fellowship applications, but it should outline the project as if you were going to execute it in the postdoc. Include:

- **Title and Hypothesis:** Give your project a working title and state the core hypothesis or research question. For example: *"Hybrid Quantum-Classical Workflow for Drug Binding Affinity Prediction"* – Hypothesis: *Combining quantum chemistry calculations for key interaction motifs with machine learning will yield more accurate binding affinity predictions than purely classical or purely quantum approaches, within feasible computational cost.* Make sure the hypothesis is testable and specific. This will form the narrative of your project.
- **Background and Significance:** Write a brief background (a few paragraphs) explaining the problem and why it's important. Cite a couple of sources from your Phase 3 literature review to show context. For instance: "Drug-target affinity prediction remains a challenge; current ML models achieve X accuracy <sup>28</sup>, but often miss certain quantum-mechanical effects. Recent work suggests including quantum-derived features could improve predictions <sup>2</sup>. Meanwhile, quantum computing approaches have been demonstrated for tiny systems <sup>3</sup>, but are not yet scalable <sup>29</sup>. This project aims to bridge these developments...". This section is essentially a mini intro like you'd see in a paper. It convinces the reader (and yourself) why your project should be done.
- **Approach/Methods:** Outline how you plan to tackle the problem. Break it into tasks or aims. Example: Aim 1 – **Data curation:** compile a dataset of known binding affinities for a set of protein-ligand complexes, along with quantum chemistry calculations for each complex's interaction energy

(small model system calculations). Aim 2 – **Model development**: build a hybrid model that inputs both classical descriptors and quantum features into a ML model (perhaps a neural network) to predict binding affinity. Aim 3 – **Evaluation**: evaluate on a benchmark set (e.g., PDBbind core set) and compare against standard scoring methods <sup>13</sup>. Aim 4 – **Quantum scalability analysis**: use a quantum simulator or small quantum hardware via Qiskit to test the approach on a very limited case, to identify bottlenecks for future quantum advantage.

- For each aim, list what tools or techniques you would use (this shows you know how to implement it). For instance: “Use RDKit and MD simulations (OpenMM) to generate binding poses; use Psi4 for quantum energy calc on ligand fragments; use PyTorch for model training; possibly Qiskit for a proof-of-concept quantum calculation.”
- Also consider the timeline and feasibility: acknowledge what can be done in 1-2 years of a postdoc. You’re not committing to this plan, but it’s good to show you can assess scope. You might say, “In the first 6 months, focus on dataset and classical model; months 6-12, integrate quantum calculations; year 2, refine and attempt a small quantum computing experiment.” This level of planning is impressive if conveyed.
- **Expected Outcomes and Impact**: Describe what success looks like. How will you measure it? For example: “Success would be a hybrid model that improves prediction accuracy of binding affinities by 20% over baseline ML models on the benchmark.” <sup>9</sup> We would also demonstrate the approach on at least one prospective case (a new molecule not in the training data). If successful, this approach could be a new paradigm for integrating quantum mechanics in AI-driven drug design, guiding more accurate lead optimization.” Also, be honest about risks and challenges: “We anticipate challenges in computational cost; if full quantum calculations are too slow, we will explore lower-level quantum methods or train a surrogate ML model for the quantum part, etc.” Discussing mitigation strategies shows maturity.
- **References**: Even in your own notes, list the key references that support your idea (you likely have them from the literature review). This will allow you to quickly cite prior work when writing papers or proposals later. It also reminds you to give credit and context in discussions.

This research proposal exercise will crystallize your understanding. By writing it out, you’ll discover if there are any knowledge gaps you still need to fill (maybe you realize you’re not sure how to obtain a certain kind of data – you can then quickly research that). It’s better to identify and address these now, before you’re in an interview or starting the project under a supervisor.

## Planning and Next Steps

With a proposal drafted, you have essentially prepared a roadmap for **your future research**. Here are final steps and tips to maximize the benefit of Phase 3:

- **Feedback**: If possible, seek feedback on your project idea. This could be from a current advisor, a colleague, or even online forums. You don’t need to give away your whole idea (if you’re worried about novelty), but discussing with someone experienced can provide valuable perspective. They might point out an aspect you overlooked or suggest a reference. If you don’t have a person to ask, consider posting a high-level question on communities like Stack Exchange (e.g., the Computational

Chemistry StackExchange) or relevant subreddits, to test if the idea sounds compelling. Feedback can help refine your plan.

- **Prepare to Communicate:** Now that you have this plan, prepare to **present it**. As a postdoc applicant, you may be asked about what you want to do or how you'd approach a problem. You could prepare a 5-minute talk (with a couple of slides) summarizing your proposal, as if you were presenting a conference poster idea. Practicing explaining it in simple terms (to a non-expert friend, for example) can help. You want to be able to convey the essence: *"I see an opportunity to solve X using Y, building on my skills in ML/chemistry. Here's my idea..."*. The clarity in your plan will translate to clarity in explanation.
- **Stay Updated:** Science moves fast. Continue to stay informed about new developments even as you transition to your postdoc. Set Google Scholar alerts for keywords of your project, or follow key people/organizations on Twitter or LinkedIn. For instance, follow DeepChem's community updates (their forum or newsletter), or IBM Quantum's news for any breakthrough that could influence your plan. This habit will keep your knowledge fresh and allow you to incorporate the latest techniques into your research. It will also impress colleagues that you're on top of new trends (which is expected of a postdoc who will often be guiding grad students in cutting-edge directions).

By the end of Phase 3, you will have a concrete research idea and a plan to execute it. You've essentially trained yourself to go from broad learning to pinpointing a novel research question – a journey very similar to what you'll do in a postdoc (and indeed throughout a research career). This will give you **confidence in interviews**, because you can speak not only about what you've learned and done (Phase 1 and 2), but also about what you **plan to do** and why it matters. You'll come across as a well-prepared, independent, and motivated researcher – exactly what postdoc advisors are looking for.

## Additional Resources and Tips

To conclude, here are some general resources and tips to support you throughout this learning journey and into your postdoc:

- **Communities and Forums:** Engage with the computational chemistry and ML community. The **DeepChem forum** <sup>34</sup> is a great place to ask questions about using ML for molecules. The **RDKit mailing list/Slack** can help with cheminformatics queries. There are also specialized communities for quantum computing (IBM Quantum has a Slack/Discord, Qiskit has forums). Don't hesitate to ask for help – the community can often save you time by pointing to the right resource.
- **GitHub Repositories to Explore:** Browse GitHub for projects related to your interests. For example:
  - **Open Source Models:** Check out **ChemBerta** or **MolBERT** (transformer models for molecules) if interested in NLP-like approaches to chemistry. Or **OpenForceField** initiative if force fields interest you. These repos often have examples and issues you can learn from.
  - **Notebooks and Tutorials:** DeepChem's GitHub has a "examples" directory with notebooks on various topics. Similarly, Qiskit has example repos (like the Qiskit Nature tutorials <sup>24</sup>). Exploring these and running them on your machine is a great way to learn by example.

- **Datasets:** The Therapeutic Data Commons (TDC) is a newer initiative that curates datasets for AI in drug discovery <sup>32</sup>. Their GitHub ( [THUDM/TDC](#) ) provides a Python API to access tasks like ADMET prediction, DTI, etc. This could be a goldmine for future projects or benchmarking your models.
- **Time Management:** With only 4 hours per day, make sure to **pace yourself and avoid burnout**. It's better to be consistent (4 hours daily) than to cram 12 hours on weekends. Each week, plan what you aim to achieve and adjust if you find things taking longer. There will be tough weeks (quantum mechanics might feel daunting, or debugging code might eat time) – that's normal. The key is to keep moving and not get discouraged by occasional roadblocks. Use the variety in the curriculum to your advantage: if you're stuck on a math-heavy concept one day, perhaps switch to writing some code for a different task and come back later with fresh eyes.
- **Leverage Copilot and AI assistants wisely:** You have a powerful tool in Copilot – it can generate code and even pseudo-explanations, but remember it's not always correct. Use it to accelerate routine coding and to see examples of library usage, but always validate the output. As you get into advanced topics, you might also experiment with prompting ChatGPT (or similar) for explanations of papers or debugging help. For example, you can paste an error message or ask for a simpler explanation of a concept from a paper – this can sometimes give quick insights. However, treat AI suggestions as *assistants*, not authorities.
- **Documentation and Notes:** Maintain a personal knowledge base. This could be a Notion page, a GitHub wiki, or just a folder of markdown files. Keep track of commands you often use, tricky installation steps you solved, summaries of papers you read, etc. This will be incredibly useful reference in the future. Plus, writing notes aids retention.
- **Soft Skills:** While the focus is on technical skills, remember aspects like **communication and collaboration**. If you have opportunities to present what you're learning to others (even informally), do it. Explaining “here's how a GNN works” to a friend or writing a blog post about your mini project can reinforce your understanding and reveal any gaps. It also prepares you for the clear communication expected in a postdoc role (where you may mentor students or present at lab meetings). Similarly, if you can contribute to an open-source project (even a minor documentation fix in DeepChem or a tutorial in Qiskit), it shows collaborative spirit and gets you accustomed to the standards of research code.

Finally, stay curious and enjoy the **process of learning**. The field of computational drug discovery is vibrant and impactful – from designing life-saving therapeutics to pioneering algorithms that push the boundaries of technology. By following this roadmap, you are equipping yourself with a rare combination of skills. Take pride in small victories along the way (your first working molecular model, your first quantum circuit run, etc.). Each is a step toward mastery. Good luck on your journey, and may it lead to exciting breakthroughs in your postdoc and beyond!



- 3 29 30 A hybrid quantum computing pipeline for real world drug discovery | Scientific Reports  
[https://www.nature.com/articles/s41598-024-67897-8?error=cookies\\_not\\_supported&code=a6208bd2-1d94-48fb-8e1b-f487f9baadaa](https://www.nature.com/articles/s41598-024-67897-8?error=cookies_not_supported&code=a6208bd2-1d94-48fb-8e1b-f487f9baadaa)
- 4 RDKit for Beginners: A Gentle Introduction to Cheminformatics — ChemCopilot: Copilot for Chemical Formulation  
<https://www.chemcopilot.com/blog/rdkit-for-beginners>
- 5 6 GitHub - rdkit/rdkit: The official sources for the RDKit library  
<https://github.com/rdkit/rdkit>
- 7 Basic docking — Autodock Vina 1.2.0 documentation - Read the Docs  
[https://autodock-vina.readthedocs.io/en/latest/docking\\_basic.html](https://autodock-vina.readthedocs.io/en/latest/docking_basic.html)
- 8 OpenMM - RCC User Guide  
<https://docs.rcc.uchicago.edu/software/apps-and-envs/openmm/>
- 9 arxiv.org  
<http://www.arxiv.org/pdf/1703.00564v1>
- 10 DeepChem  
<https://deepchem.io/tutorials/an-introduction-to-moleculenet/>
- 11 Datasets - MoleculeNet  
<https://moleculenet.org/datasets-1>
- 12 34 GitHub - deepchem/deepchem: Democratizing Deep-Learning for Drug Discovery, Quantum Chemistry, Materials Science and Biology  
<https://github.com/deepchem/deepchem>
- 13 32 Machine Learning for Drug Development - Zitnik Lab  
<https://zitniklab.hms.harvard.edu/drugml/>
- 14 Transforming Computational Drug Discovery with Machine Learning ...  
<https://pubs.acs.org/doi/10.1021/acsmmedchemlett.8b00437>
- 15 16 17 GitHub - psi4/psi4numpy: Combining Psi4 and Numpy for education and development.  
<https://github.com/psi4/psi4numpy>
- 20 21 22 23 25 26 31 33 Synergizing quantum techniques with machine learning for advancing drug discovery challenge | Scientific Reports  
[https://www.nature.com/articles/s41598-024-82576-4?error=cookies\\_not\\_supported&code=852b2601-44bd-45c1-9437-8f6249e2c2ac](https://www.nature.com/articles/s41598-024-82576-4?error=cookies_not_supported&code=852b2601-44bd-45c1-9437-8f6249e2c2ac)
- 24 Qiskit Nature Tutorials - Qiskit Nature 0.7.2  
<https://qiskit-community.github.io/qiskit-nature/tutorials/index.html>
- 27 Advancing Drug Discovery with Quantum Computing Breaking ...  
<https://nhjsj.com/2025/advancing-drug-discovery-with-quantum-computing-breaking-artificial-intelligence-barriers/>
- 28 Quantum computing in drug discovery  
<https://ojs.acad-pub.com/index.php/ISSC/article/view/294>