

## **What Should Be ADDED to Improve the Pipeline (Biology + Bioinformatics)**

Below is a **strictly additive list**, grouped by pipeline stage.

### **A. Target Biology & Context (Currently Missing Entirely)**

#### **What to Add**

##### **1. Target annotation module**

- ⑩ UniProt ID resolution
- ⑩ Protein function, pathway, disease relevance

##### **2. Biological rationale check**

- ⑩ Is the target druggable?
- ⑩ Is it intracellular / extracellular?

##### **3. Isoform awareness**

- ⑩ Which isoform is being docked?
- ⑩ Does the PDB match the biologically relevant isoform?

#### **Why This Matters**

Right now, the pipeline assumes the protein is biologically meaningful.  
**Docking without biological justification is not publishable.**

### **B. Protein Structure Quality Control (Partially Implemented, Incomplete)**

You already have:

- ⑩ protein\_prep.py
- ⑩ ramachandran.py

#### **What to Add**

##### **1. Pre-docking structure quality gate**

- ⑩ % residues in favored Ramachandran regions
- ⑩ Missing residues count
- ⑩ Clash score (even approximate)

##### **2. pH-aware protonation**

- ⑩ Histidine state selection (HIS, HID, HIE)
- ⑩ Salt bridge preservation

##### **3. Multi-structure support**

- ⑩ Dock against:

- ⑩ Crystal structure
- ⑩ AlphaFold model
- ⑩ At least one alternative conformation

### Why This Matters

Docking accuracy depends more on **protein quality** than ligand quality. Your pipeline currently **checks geometry but does not enforce acceptance criteria**.

## C. Binding Site Biology (Weak Validation)

You use **P2Rank**, which is good.

### What to Add

#### 1. Pocket biological plausibility filter

- ⑩ Pocket near catalytic residues?
- ⑩ Pocket conserved across homologs?

#### 2. Known ligand cross-mapping

- ⑩ Does predicted pocket overlap with:
  - ⑩ Co-crystallized ligands?
  - ⑩ BindingDB annotations?

#### 3. Pocket consensus

- ⑩ Combine P2Rank + geometric pocket detection (e.g., fpocket-style logic)

### Why This Matters

Automatically predicted pockets often include false positives.  
Right now, the pipeline **trusts P2Rank blindly**.

## D. Ligand Biology & Chemistry (Major Weakness)

This is the **single biggest biological gap**.

### What to Add

#### 1. Ligand sanity checks

- ⑩ Molecular weight
- ⑩ Formal charge
- ⑩ Rotatable bond count

#### 2. Protonation and tautomer enumeration

- ⑩ Physiological pH ≠ neutral molecule

#### 3. Chemical class awareness

⑩ Peptides vs small molecules

⑩ Macrocycles vs fragments

#### 4. Redundancy removal

⑩ Same scaffold, same binding mode

### Why This Matters

You are docking **FDA drugs**, but:

⑩ Some are biologics

⑩ Some are prodrugs

⑩ Some are not meant to bind intracellular targets

Docking them **without filtering is biologically misleading.**

### E. Docking Strategy (Technically Works, Scientifically Thin)

You use AutoDock Vina correctly.

#### What to Add

##### 1. Multiple docking runs per ligand

⑩ Different random seeds

##### 2. Pose stability check

⑩ RMSD clustering of poses

##### 3. Flexible side chains at binding site

⑩ Even 3–5 residues improves realism

##### 4. Decoy docking

⑩ Dock inactive compounds as negative controls

### Why This Matters

Single-run docking = **high variance, low confidence.**

### F. Post-Docking Biological Validation (Missing)

#### What to Add

##### 1. Interaction fingerprinting

⑩ H-bonds

⑩ Salt bridges

⑩  $\pi-\pi$  stacking

##### 2. Residue importance mapping

⑩ Are key catalytic residues involved?

### 3. Water mediation awareness

⑩ Flag buried polar interactions that require water

## Why This Matters

Binding energy alone does not imply biological inhibition.

## G. ADMET Module (Conceptually Correct, Scientifically Underspecified)

You do ADMET prediction.

### What to Add

#### 1. Model transparency

⑩ Which dataset?

⑩ Which ML model?

#### 2. Confidence intervals

⑩ Not just a single value

#### 3. Contradiction detection

⑩ High affinity + poor bioavailability

#### 4. Human vs animal relevance

⑩ Species-specific metabolism flags

## Why This Matters

ADMET predictions must be interpreted, not reported raw.

## H. Cross-Stage Intelligence (Entirely Missing)

This is where the pipeline can become **research-grade**.

### What to Add

#### 1. Decision engine

⑩ Drop ligands that fail  $\geq 2$  biological criteria

#### 2. Explainability layer

⑩ Why was ligand A ranked higher than B?

#### 3. Iterative refinement loop

⑩ Redock top hits with stricter settings

### **3. Efficiency Improvements (Pipeline-Level)**

#### **Add These to Improve Throughput and Reliability**

- ⑩ Parallel docking with job scheduling
- ⑩ Checkpointing (resume after crash)
- ⑩ Hash-based result caching
- ⑩ Versioning of results with metadata