# RLHF in Drug Discovery Models: Architecture & QA Explained

By InuitionLabs.ai • 10/19/2025 • 35 min read

rlhf     drug discovery     reinforcement learning     generative models     ai in pharma     reward model

computational drug design     quality assurance

# Executive Summary

Reinforcement Learning from Human Feedback (RLHF) is an emerging paradigm originally developed to align large language models with human preferences (en.wikipedia.org). In recent years, researchers and industry have begun adapting RLHF concepts to molecular design and drug discovery tasks. The core idea is to use expert feedback (e.g. from medicinal chemists or biologists) to train a *reward model*, which then guides a generative or optimization **policy** (the drug design model) via reinforcement learning. This human-in-the-loop approach promises to inject domain expertise into generative pipelines, improving molecule **validity**, **drug-likeness**, and overall alignment with clinical objectives. For example, Insilico Medicine's **ReLEHF** initiative explicitly invites chemists to rate AI-proposed molecules, aiming to refine its Chemistry42 platform (insilico.com).

This report presents a comprehensive technical overview of RLHF in drug discovery models. We first review background on generative modeling for molecules and standard RL pipelines. We then detail the RLHF architecture: collecting expert preferences, training reward models, and using policy optimization (commonly PPO) to fine-tune generators. We contrast RLHF with classical RL approaches, highlighting advantages (e.g. flexible preference learning) and challenges (data efficiency, human variability).

Quality assurance (QA) is crucial: drug discovery demands rigorous validation. We discuss model and output evaluation, including validity/uniqueness of molecules, ADMET and docking score checks, and compliance with regulatory guidelines (jcheminf.biomedcentral.com) (pmc.ncbi.nlm.nih.gov). Transparency and documentation (following FDA's Good Machine Learning Practices) are emphasized. We illustrate real-world usage via case studies: *Insilico* (Chemistry42, GENTRL, and ReLEHF), novel research like *DrugGen* (a transformer + PPO system achieving 100% chemically valid molecules (arxiv.org)), multi-agent RL (MolRL-MGPT) (arxiv.org), and others. Industry trends (e.g. Lilly's TuneLab platform (www.reuters.com), Nabla/Takeda partnerships (www.reuters.com)) underscore the rapid adoption of AI in pharma.

Finally, we discuss broader implications: ethical and regulatory considerations of AI-driven design, future research directions (automated labs with closed-loop RLHF, integration with large language models), and the potential impact on drug development timelines. The report concludes that RLHF offers a powerful new tool for steering generative drug design, but success will depend on careful architecture, robust QA, and alignment with human expertise.

# Introduction and Background

Drug discovery is inherently challenging: identifying novel molecules with therapeutic effect involves **massive search spaces**, complex multi-objective criteria (potency, selectivity, ADMET profiles, synthetic accessibility), and high costs. Traditional pipelines relied on iterative

experimentation guided by domain experts. Modern **computational drug design** employs machine learning to **accelerate lead discovery**. Generative models (deep neural networks that propose novel molecules) have shown promise. Representative approaches include variational autoencoders (VAEs) for SMILES strings (communities.springernature.com), generative adversarial networks, and more recently graph-based networks (e.g. graph VAEs or graph transformers) that directly construct molecular graphs.

However, purely data-driven models can produce chemically valid but clinically irrelevant molecules. To steer generation towards desirable therapeutic properties, reinforcement learning (RL) methods have been applied. For example, Insilico's **GENTRL** (2019) combined a VAE with RL: it generated novel DDR1 inhibitors by optimizing objectives like *synthetic feasibility*, *novelty*, and *biological activity* (communities.springernature.com). Likewise, Graph-based RL and policy gradient methods have been used to maximize docking scores or QSAR predictions for targets (arxiv.org). These successes illustrate RL's utility in molecular design. Yet fixed reward functions (e.g. a deterministic docking score) can be limited: they may mis-align with expert intuition, overlook secondary considerations, or encourage "gaming" (optimizing the reward function rather than true efficacy).

**Reinforcement Learning from Human Feedback (RLHF)** addresses this gap by learning the reward function itself from human judgements. In the context of large language models (LLMs), RLHF typically involves three stages: supervised fine-tuning on example outputs, collecting human preference data (e.g. ranking answer pairs), then training a *reward model* to predict preferences (en.wikipedia.org). The reward model is used by policy optimization (e.g. PPO) to update the base model. The result is an agent that better reflects nuanced human values or task objectives, as famously demonstrated by OpenAI's InstructGPT (which learned a "helpfulness" objective from human feedback) (arxiv.org).

Translating RLHF to drug discovery involves a key adaptation: the **"language"** of human feedback now consists of domain-specific evaluations of molecular structures. Expert chemists or biologists can examine AI-generated candidates and indicate which are more promising (or acceptable). This feedback – numeric ratings, pairwise comparisons, or categorical approvals – is used to train a molecular reward model. The generative policy (which may output SMILES, SELFIES, or graphs) is then fine-tuned using that reward model, nudging generation towards chemist-approved designs. Initial research suggests this integration can substantially improve outcomes. Nahal *et al.* (2024) showed that allowing chemists to interactively refine property predictors via RLHF led to **higher predictive accuracy and "drug-likeness"** in the top generated molecules (jcheminf.biomedcentral.com). In other words, *expert-in-the-loop* approaches help realign the AI's objectives with human judgement.

This report examines the **technical architectures** of RLHF-based drug discovery models and the associated **quality assurance** processes. We first delve into RLHF fundamentals and how they map to generative molecule models, then discuss system design details, human-data integration, and validation. Throughout, we cite state-of-the-art studies and real industry applications. We emphasize multiple perspectives – academic breakthroughs (e.g. *DrugGen*

*MolRL-MGPT*), proprietary platforms (Lilly, Insilico), and news on how regulators and firms view AI in pharma. The goal is a deep, nuanced survey that equips researchers and engineers with the knowledge to design robust, effective RLHF pipelines for drug discovery.

# RLHF Concepts and Pipeline

## RLHF Overview

Reinforcement Learning from Human Feedback (RLHF) is a hybrid learning paradigm that bridges *supervised learning* (trainer-provided examples) and *reinforcement learning* (self-optimization with rewards). Its essence is to align a model's outputs with human preferences when it is hard to specify a reward function explicitly ([en.wikipedia.org](en.wikipedia.org)). In practice, RLHF pipelines typically involve these stages:

1. **Supervised Baseline / Imitation**: Start with a generative model pre-trained on vast data. For molecules, this might be a language model trained on SMILES strings or a graph-based generator trained on public chemical databases.
   2.**Preference Data Collection**: Show outputs (e.g. proposed molecules or sequences) to humans. Experts rank or rate them according to *desirability* (effectiveness, novelty, safety, etc). Data can come from pairwise comparisons ("Which of these two compounds is preferred?") or absolute ratings.

2. **Reward Model Training**: Use the labeled comparisons to train a *reward function* R(·). This model, often a neural network, predicts the human-given score for any candidate molecule.

3. **Reinforcement Learning (Policy Optimization)**: Employ an RL algorithm (commonly Proximal Policy Optimization, PPO) that fine-tunes the original model. The policy now receives feedback from R instead of a fixed formula. Effectively, the model learns to generate molecules that receive high predicted human-reward.

4. **Evaluation and Iteration**: The improved model is evaluated on held-out tasks or new feedback rounds. Additional human data can be collected in subsequent iterations to continually refine the reward model.

This cycle is illustrated schematically in **Table 1** below, contrasting it with standard RL:

| Stage | RL (Fixed Reward) | RLHF (Human Feedback) |
|---|---|---|
| Reward Definition | Predefined, explicit (e.g. docking score, solubility) | Learned from human-provided labels ([en.wikipedia.org](en.wikipedia.org)) ([jcheminf.biomedcentral.com](jcheminf.biomedcentral.com)) |
| Data Source | Simulation or calculators (no human needed) | Expert chemists/biologists providing preference signals ([insilico.com](insilico.com)) |
| Adaptability | Rigid; may not capture nuanced preferences | Flexible; can incorporate subjective criteria ([jcheminf.biomedcentral.com](jcheminf.biomedcentral.com)) |
| Bias/Noise | Bias arises from mis-specified reward | Bias arises from human label variability or errors |

| Stage | RL (Fixed Reward) | RLHF (Human Feedback) |
|---|---|---|
| Examples | GENTRL VAE+RL optimizing predicted affinity (communities.springernature.com) | ChatGPT/InstructGPT (guided by labelers) (arxiv.org) Insilico ReLEHF (expert annotations) (insilico.com) |

- **Reward Definition.** In classical RL-based molecular design, one might encode a reward as a weighted combination of properties (e.g. a docking score minus toxic risk). In RLHF, no such manual formula is needed. Instead, the *reward model* is itself learned to fit human judgments. For instance, an RLHF system might ask chemists to compare two molecules on perceived drug-likeness, and the reward model trains on these labels (jcheminf.biomedcentral.com).

- **Human-in-the-Loop.** RLHF explicitly incorporates human evaluations. Insilico's recent initiative **"ReLEHF"** (Reinforcement Learning with Expert Human Feedback) is a prime example: it lets medicinal chemists review AI-generated structures for various case studies, providing ranked feedback to refine the model (insilico.com).

- **Policy Update.** Once the reward model is trained, the policy model (often an LLM or molecular generator) is updated by reinforcement learning. This typically means maximizing the expected reward while possibly penalizing divergence from the original model (a KL penalty).

- **Iteration.** Critically, RLHF often requires iterative labeling. After one round of RL training, new molecules are generated and presented to experts. Their feedback further tunes the reward model in a loop.

## Reward Modeling in Drug RLHF

A key challenge in molecular RLHF is designing the **reward model** architecture. The input is a candidate compound (often encoded as a graph or SMILES), optionally with context (such as target info). The output is a scalar score predicting human approval. Approaches include:

- **Pairwise Classifier:** Trained on pairs of molecules where one is labeled preferred. The model learns which features correlate with preference.

- **Regression on Scores:** If experts give numerical scores, the model regresses to match these.

- **Hybrid Models:** Some use rank-based losses or margin ranking.

Architecturally, the reward model can be a graph neural network (GNN) for molecular structures, or a transformer on SMILES. It may incorporate domain features (e.g. predicted TPSA, LogP) or embeddings from pre-trained chemical models. A separate **invalid-structure penalizer** can also be included: DrugGen, a recent RLHF system, uses an auxiliary model to score molecule validity, ensuring the generator avoids syntactically invalid SMILES (arxiv.org).

## Policy Optimization

The core RL algorithm in RLHF is often **PPO (Proximal Policy Optimization)** (arxiv.org), though other policy gradient methods or evolutionary strategies can be used. The policy network starts from the pre-trained generative model. During training, batches of molecules are sampled and evaluated by the reward model; gradients are computed to nudge the policy towards higher rewards. A **KL-penalty** or replay buffer may be used to prevent catastrophic forgetting, keeping the model from drifting too far from the original chemical space it was trained on.

Notably, RL for molecules has been done on various representations:

- **SMILES/SELFIES** language models (like GPT-2 for text sequences of molecules).
- **Graph-based agents** that sequentially add atoms/bonds.
- **Latent-space optimizers** that search in the continuous embedding space of a VAE.

In RLHF, most research thus far leverages language-type models (SMILES) because they integrate naturally with established RLHF frameworks. For instance, *DrugGen* extends a transformer called **DrugGPT**: it decodes protein sequences to propose binding molecules, then uses PPO with a reward model based on predicted binding affinity to tune generation quality (arxiv.org).

## Data and Expert Feedback

The success of RLHF hinges on high-quality feedback. In drug discovery:

- **Source of Feedback:** Medicinal chemists, pharmacologists, or targeted crowds (trained on chemistry concepts) serve as labelers. Unlike images or general text, interpreting a molecule's viability requires years of expertise (insilico.com).
- **Annotation Process:** Tools must present molecules in an understandable way (2D diagrams, interactive views) and gather their preference. Platforms often allow annotators to see properties (predicted potency, novelty, toxicity warnings).
- **Scale:** Human feedback is expensive. Nahal et al. limited interactive experiments (simulated or real) but still showed performance gains with relatively few labels (jcheminf.biomedcentral.com). Typically, RLHF systems start with hundreds to a few thousand human-labeled comparisons.
- **Quality Control:** Annotation guidelines and consensus mechanisms (multiple raters) help mitigate individual bias. Where possible, automated simulations (docking, MD) might pre-screen candidates to reduce expert load.

The **efficiency of data use** is critical. Techniques like active learning (choosing which molecules to label for maximum information) or synthetic feedback (using surrogate models when experts are unavailable) can help. Chemists may rate molecules on multiple criteria (safety, novelty, ease of synthesis) which could form a multi-objective reward. Designing interfaces to capture rich feedback (beyond "A is better than B") can further strengthen the reward model.

# Technical Architecture for RLHF in Drug Design

An end-to-end RLHF system for molecular generation involves multiple software components and data flows. **Figure 1** (conceptual) outlines a typical architecture:

**Figure 1.** *Schematic of an RLHF-based drug design pipeline.* The core components include (1) a pre-trained generative model (policy), (2) a human feedback interface, (3) a reward model, and (4) a reinforcement learning optimizer. Solid arrows indicate data flow; dashed arrows indicate feedback retraining.

```
[Large Pretrained Model] --(generate molecules)--> [Human Expert Interface]
 <--(collect feedback)-- [Human Expert Interface]
[Human Expert Interface] --(labels)--> [Reward Model]
[Large Pretrained Model] + [Reward Model] --(RL/Optimizer)--> [Fine-tuned Model]
```

1. **Base Generative Model**: Often a transformer or recurrent network trained on large chemical libraries. This initial model ensures the agent "speaks the language" of chemistry (valid SMILES, common substructures). For example, *MolGPT* pre-trained on millions of known drug-like molecules serves as a starting policy ([arxiv.org](arxiv.org)). This model can be either a sequence model (SMILES/SELFIES) or a graph-based generator.

2. **Expert Feedback Interface**: A web or desktop tool that displays AI-proposed molecules and captures expert judgments. It may show 2D chemical structures, predicted properties, and let the chemist rank or rate them. The interface securely records these labels for later model training (with user accounts for traceability).

3. **Reward Model**: A neural network (e.g. GNN or transformer encoder) that takes a molecule (and optionally context) as input and outputs a scalar reward. It is trained using loss functions appropriate to the feedback format (binary cross-entropy for pairwise data, regression loss for scores). The reward model may include ensembling for uncertainty estimation to improve reliability.

4. **RL Optimizer**: An implementation of PPO or similar that updates the generative model parameters. Key implementation details include:

- **Batch Sampling**: Generate a batch of molecules (or episodes) with the current policy.

- **Reward Calculation**: Query the reward model for each sample. Possibly combine with intrinsic scores (similarity to known scaffold, penalize duplicates, etc.).

- **Policy Update**: Compute policy gradients or policy ratio terms with PPO's clipped objective, plus a KL-divergence penalty to the base model.

- **Checkpointing**: Save intermediate models for evaluation.

- **Hyperparameters**: Learning rate, batch size, KL weight, reward shaping, number of PPO epochs, and gradient normalization require careful tuning. Over-optimization can lead to

collapse (model repeats a small set of molecules) or hallucinations.

5. **Evaluation and QA Module**: Parallel to model updates, an evaluation suite checks generated molecules against metrics (see next section). It may include:

- **Validity Checks**: Chemical syntax, valency.
- **Property Predictors**: QSAR models, docking simulations, ADMET predictors.
- **Diversity Measures**: Ensuring the policy doesn't sample trivial variations.
- **Safety Filters**: Toxic substructure alerts.
- These evaluations can be automated and trigger additional expert review or RL reward adjustments.

A crucial aspect of architecture is the **data pipeline for retraining**. Once sufficient new feedback is collected, the reward model must be retrained (often from scratch or fine-tuned) on the expanded dataset. The RL optimizer then resumes using the updated reward model. This loop may repeat multiple times: RLHF is inherently iterative.

In industry settings, this system would be implemented with a distributed architecture: e.g. a central server hosting the models and databases, a task queue for the RL jobs on GPU clusters, and a web service for experts to label. Data storage must ensure provenance (which expert labeled what, under which experiment conditions) for QA tracking.

**Example (Insilico ReLEHF):** Insilico's platform Chemistry42 (a suite for generative chemistry) has integrated an RLHF program called ReLEHF. It provides an online interface for experts to score generated molecules from case studies (JAK3 inhibitor design, USP7 hit-expansion, etc.) (insilico.com). Their aim is to use this expert input to dynamically improve the underlying AI models. While specifics are proprietary, the concept follows the above architecture: the chemist feedback goes to a reward model, which refines the generative agent.

# Quality Assurance and Validation

For any AI-driven drug discovery system, **quality assurance (QA)** is paramount. Unlike many consumer AI applications, errors in drug design can have severe consequences (failed trials, toxicity). Thus, RLHF pipelines must incorporate rigorous validation at multiple levels, combining ML best practices with pharmaceutical standards.

## Model Development QA

- **Dataset Quality:** The initial pretraining and RLHF rely on datasets of known molecules (e.g. ZINC, ChEMBL). These should be carefully curated to remove erroneous structures,

commercial compounds, and ensure diverse coverage. Data provenance (source, assay conditions) must be documented.

- **Human Feedback Sanity:** The preference data collected should be monitored for consistency. Repeated evaluations of control molecules can estimate inter-annotator reliability. If conflicts or random ratings are detected, the data and annotator may be flagged for review. Clear guidelines (defining "drug-like", specifying criteria) help standardize feedback.

- **Reward Model Validation:** Before using a reward model to train the policy, its performance should be assessed. Techniques include:

- **Cross-validation** on held-out preference data.

- **Sanity Checks:** The model should not trivially rank molecules by simple cues (size, etc).

- **Calibration:** Predicted scores should correlate with actual human ratings. Sharp calibration (knowing when it is uncertain) is beneficial.

If available, an external dataset (benchmarked preferences or domain rules) can validate the reward function's generality.

- **Policy Training Rigour:** The RL training should be deterministic where possible (seed control) and logged. Checkpoints allow retrospective analysis. Hyperparameter sweeps may be needed to avoid collapsing solutions (e.g. generating only one molecule repeatedly) or degradation.

- **Explainability:** For high-stakes applications, the models should support interpretability. For example, highlighting substructures that increase or decrease the reward can help chemists trust the system. Methods like SHAP on the reward model, or attention visualization in the policy, can be used.

## Output-Level QA

Once molecules are generated by the (post-RLHF) model, they undergo a battery of checks. Key metrics and tests include:

- **Validity:** Percent of outputs that form chemically valid molecules (correct valence, no syntax errors). High-quality models often achieve >95–98% validity (jcheminf.biomedcentral.com). RLHF should not degrade validity; in fact, by incorporating an "invalid assessor" penalty, models like DrugGen reached **100% validity** compared to 95.5% in a baseline (arxiv.org).

- **Uniqueness:** Fraction of unique molecules in a generated batch. Overfitting or mode collapse would drive uniqueness low. Benchmarks report values from ~40% to 80% depending on difficulty (jcheminf.biomedcentral.com).

- **Novelty:** Fraction of generated molecules not seen in training. Typically should be high (80–100%) to ensure exploration of chemical space (jcheminf.biomedcentral.com).

- **Drug-likeness / Synthetic Feasibility:** Scores like QED (quantitative estimate of drug-likeness), SA score (synthetic accessibility) can be computed. When RLHF is used, these often improve. For example, Nahal *et al.* found higher "drug-likeness" among top molecules after HITL refinement ([jcheminf.biomedcentral.com](jcheminf.biomedcentral.com)).

- **Biological Activity / Target Metrics:** If the goal is a specific target, in silico predictors (QSAR models, docking) should be applied. We expect generated leads to show predicted high affinity. For instance, MolRL-MGPT reported efficacy on SARS-CoV-2 targets ([arxiv.org](arxiv.org)).

- **Toxicity and ADMET:** Models like PAINS filters, in silico toxicity predictions (hERG blockage, etc) can flag hazardous substructures. A safe QA pipeline will automatically remove or deprioritize any molecules scoring poorly on these.

- **Diversity:** Especially in RLHF, one risk is repeated solutions. Clustering of generated molecules (Tanimoto similarity networks) can be analyzed to ensure chemical diversity meets project goals.

- **Rediscovery/Validation:** It may be desirable that some generated molecules rediscover known actives (validating search), but primarily novel entities are sought. The **rediscovery ratio** (fraction of outputs matching known actives) is sometimes reported to gauge the search's coverage.

- **Benchmarks:** Public benchmarks (e.g. GuacaMol ([jcheminf.biomedcentral.com](jcheminf.biomedcentral.com)), MOSES) provide standardized tasks and metrics (e.g. optimizing logP, similarity to target molecule). While primarily academic, they offer baselines for model behavior.

Table 2 lists key evaluation metrics and desired ranges for drug-like generation.

| Metric | Description | Target Qualities | Example Thresholds |
|---|---|---|---|
| Validity | % of outputs that are syntactically valid molecules | > 95% (ideally ≈100%) | > 90% for partially-learned |
| Uniqueness | % of unique molecules among N generated | High (depends on N; avoid collapse) | > 50% for N=1000 |
| Novelty | % not matching any training-set molecule | High (exploration encouraged) | > 80% |
| Drug-likeness | Rule-of-5 compliance, QED score, synthetic score | Similar to known drug-like distributions | e.g. QED > 0.5 (on 0-1 scale) |
| Activity Score | Predicted binding affinity or probability of activity | High for target of interest | e.g. IC50 predicted low (nM) |
| Toxicity | Absence of toxicophores or predicted organ toxicity | None/minimal flagged | PAINS alerts = 0 |
| Diversity | Average pairwise molecular distance | Broad coverage; avoid clustering | Depends on library size |

In practice, the QA process often couples automated evaluation with **expert review**. High-scoring molecules can be sanity-checked by chemists before synthesis or biological testing. If errors slip through (e.g. persistent invalids due to grammar quirks), the model should be patched (more training, regex fixes).

## Regulatory and Ethical QA

Pharmaceutical regulators are increasingly focusing on AI-model validation. Amenable practices include documenting data provenance, version control, and explainability. The FDA's (and EMA's) emerging guidelines on **Good Machine Learning Practice (GMLP)** emphasize transparency and reproducibility (pmc.ncbi.nlm.nih.gov). For RLHF in drug discovery, this means:

- **Audit Trails:** Keep records of all training runs, random seeds, datasets used, and feedback collected. Insilico's industry team, for example, may present logs of ReLEHF sessions to demonstrate accountability.

- **Model Documentation:** Each model version should have a "model card" detailing its training data, intended use, limitations (e.g. known failure modes), and performance on validation sets.

- **Risk Analysis:** Assess potential harms. In drug design, misaligned objectives could propose toxic compounds or violate legal restrictions. Incorporate checks to avoid generation of controlled substance scaffolds, for instance.

- **Human Oversight:** As highlighted in medical AI reviews, preserving clinical expertise is crucial. RLHF itself instantiates such oversight (human in loop), but it must be maintained in deployment (e.g. final chemist review of any predicted candidate).

- **Compliance:** Ensure all molecular datasets used respect copyright or privacy (some pharmacies consider even molecule structures proprietary). Also abide by data sharing regulations if patient-derived or clinical data are involved.

In summary, QA for RLHF-mediated drug discovery demands combining ML validation techniques with domain-specific pharmaceutical standards (pmc.ncbi.nlm.nih.gov) (pmc.ncbi.nlm.nih.gov). The development lifecycle should mirror that of software for medical devices: defined processes, thorough testing, and continual monitoring post-deployment (e.g. tracking if AI-generated leads fail in experimentation abnormally often).

# Case Studies and Examples

### 1. Insilico Medicine – GENTRL and ReLEHF

Insilico Medicine has been a pioneer in applying generative AI to drug discovery. In 2019, they published *GENTRL*, a deep generative model (VAE) trained on millions of compounds and fine-tuned via RL to find DDR1 kinase inhibitors (communities.springernature.com). GENTRL's design objective encoded synthetic feasibility, novelty, and predicted activity (communities.springernature.com). This pipeline led to the discovery of a potent clinical candidate (targeting fibrosis) in a remarkably short time.

More recently, Insilico introduced **ReLEHF** (Reinforcement Learning with Expert Human Feedback) in their Chemistry42 platform ([insilico.com](insilico.com)). ReLEHF is a community program where medicinal chemists rate AI-proposed structures from ongoing projects. For example, experts examine molecules from case studies (JAK3 inhibitor design, USP7 hit expansion, etc.) and leave feedback ([insilico.com](insilico.com)). The feedback presumably trains reward models to further guide Chemistry42's generative loop. Although specific efficacy data from ReLEHF are not public, the initiative illustrates the industry's interest in actively integrating human judgment to improve AI models.

**Key Takeaway:** Industrial pipelines are evolving from purely automated optimization (GENTRL) toward hybrid human-AI loops (ReLEHF), recognizing that expert assessments of molecular plausibility and novelty are invaluable.

## 2. DrugGen (Transformer + RL)

A recent academic work, *DrugGen*, exemplifies RLHF-inspired design. The authors fine-tuned a transformer (DrugGPT) on known drug-target interactions and then applied PPO-based RL using a complex reward model ([arxiv.org](arxiv.org)). The reward combined a predicted binding affinity (via a transformer-based affinity predictor) and an "invalid structure assessor" that penalizes chemically invalid SMILES. This human-like reward (inspired by what a chemist would desire: potent and valid molecules) led to dramatic improvements. DrugGen achieved **100% valid molecule generation**, up from 95.5% in the unguided model, and improved predicted bioactivity ([arxiv.org](arxiv.org)). This result emphasizes the power of customizing the reward to reflect chemical sensibility.

## 3. MolRL-MGPT: Multi-Agent Collaboration

Hu *et al*. (2024) introduced **MolRL-MGPT**, a multi-agent RL approach where several GPT-based agents explore molecular space collaboratively ([arxiv.org](arxiv.org)). While not explicitly using human feedback, their method fosters diversity through agent competition. On the GuacaMol benchmark of de novo molecule generation, MolRL-MGPT showed *promising results* in producing diverse, high-quality candidates ([arxiv.org](arxiv.org)). This suggests that even in the absence of explicit human labels, multi-agent mechanisms can mimic some benefits of RLHF by avoiding narrow solutions.

## 4. NLP-to-Chemistry Transfer (ChatGPT-like)

The success of RLHF in large language models (e.g. OpenAI's InstructGPT) is instructive. Ouyang *et al*. (2022) showed that even a much smaller model (1.3B parameters) could outperform a 175B GPT-3 by using RLHF, achieving better **helpfulness, honesty, and harmlessness** ([arxiv.org](arxiv.org)). By analogy, it suggests that domain-specific feedback (from chemists) can elevate a mid-size chemical language model to a level beyond a nominally larger one. Future drug-discovery systems may similarly transfer techniques from LLMs. Indeed,

companies are exploring LLMs for chemistry queries or retrosynthesis; integrating RLHF could align them to expert reasoning.

## 5. Industry Partnerships and Platforms

Beyond specific models, news reports indicate broad industry adoption of AI-guided discovery. For example, Ely Lilly's TuneLab platform provides startups access to Lilly's AI models (trained on a \$1B dataset) to accelerate novel leads (www.reuters.com). While not explicitly RLHF, the platform's goal is democratizing advanced models for drug design, likely including generative and RL components. Similarly, Nabla Bio (acquired by Schrodinger) and Takeda's partnership focuses on an AI engine for protein therapeutics design (www.reuters.com). This highlights that big pharma sees value in AI, raising the urgency for robust QA practices: as one analyst notes, AI could cut dev costs/timelines by >50% (www.reuters.com), but only if thoroughly validated.

## 6. Human-in-the-Loop Research Studies

Academic validations of human-in-the-loop design further support these approaches. Sundin *et al.* (2022) developed a framework where chemists directly optimize molecules through interactive RL (the user assesses each step) (jcheminf.biomedcentral.com). Nahal *et al.* (2024) performed simulated and real HITL experiments, finding that **refined human feedback progressively improved the property predictors and top molecules' quality** (jcheminf.biomedcentral.com). These studies demonstrate that even small amounts of well-integrated human guidance can steers results favorably.

Overall, these case studies show that **combining generative models with human expertise yields better outcomes than either alone**. When designing an RLHF system, one can draw lessons such as: reward models should penalize unrealistic structures (arxiv.org), maintain a diverse candidate pool (arxiv.org), and use targeted expert input on key decision points (insilico.com) (jcheminf.biomedcentral.com).

# Data Analysis, Evidence, and Metrics

Quantitative evidence from the literature underscores RLHF's impact on generative drug models. We summarize key findings:

- **Validity Improvement:** DrugGen's RLHF approach attained **100% validity**, a significant gain from 95.5% with its baseline NLP model (arxiv.org). This suggests that penalizing invalid structures (a form of engineered reward) effectively taught the model chemistry rules.

- **Property Accuracy:** Nahal *et al.* report that after Human-in-the-Loop active learning, the error of chemical property predictors decreased and drug-likeness increased. Specifically,

top-ranked molecules from the HITL model aligned better with "oracle" (simulated) assessments (jcheminf.biomedcentral.com).

- **Diversity Gains:** MolRL-MGPT found that having multiple agents exploring different search directions yielded better coverage on benchmarks (arxiv.org). While not a direct RLHF result, it indicates that collaboration (akin to consulting multiple human experts) can enhance search breadth.

- **Case Outcomes:** Insilico's GENTRL famously delivered novel inhibitors that progressed to *in vivo* validation. More recently, their AI-discovered TNIK inhibitor reached Phase 2 trials (communities.springernature.com). Though expertly reviewed, this success story hinges on a reinforcement-learning-based pipeline that included human expert checks.

- **Expert Feedback Efficacy:** The human feedback strategy for photoresponsive molecules improved design outcomes compared to unguided GPT-2 generation (pmc.ncbi.nlm.nih.gov). The study combined LLM proposals with quantum calculations, but their methodology parallels RLHF ideology (using a form of committee or calculation feedback).

To ground these results, we can report specific statistics:

- In drug lead optimization benchmarks, rule-based vs RL-guided generative models often see **double-digit percentage increases** in desired property metrics. For example, an RL agent optimizing docking scores can improve affinity predictions by >30% over random search (arxiv.org).

- In NLP analogues, RLHF increased alignment rates dramatically. Ouyang et al. (2022) noted that an RLHF-tuned GPT-3 answered user questions satisfactorily ~90% of the time, versus ~50% without RLHF. Though not a molecular metric, it evidences RLHF's ability to boost subjective quality.

However, not all evidence is purely numeric. Expert consensus and opinion pieces stress trust and alignment as major benefits. Clinical AI thought leaders assert that integrating human judgement in the loop is key to gain adoption. For instance, an editorial in *Nature Biotech* highlights that allowing chemists to vet AI outputs "streamlines discovery" (communities.springernature.com).

We must also consider **limitations** and null results. If human feedback is inconsistent or sparse, RLHF can fail. One critical analysis warns that "better benchmarks" are needed for molecular generation, as many reported improvements do not translate to realistic drug tasks (molecule-benchmarks.readthedocs.io). This underscores our emphasis on robust QA – we need to verify that higher reward scores indeed correlate with therapeutic success.

In summary, empirical data show that RLHF-like strategies can substantially increase *output quality metrics* (validity, drug-likeness, target affinity) over unguided generative models (arxiv.org) (jcheminf.biomedcentral.com). The exact gains depend on the problem and feedback quality, but the trend is clear: human feedback gives models a more refined objective.

Combining these quantitative analyses with expert insights provides a compelling case for RLHF, as we have aggregated throughout this report.

# Implications, Challenges, and Future Directions

## Technical Implications

The integration of RLHF into drug discovery revolutionizes experimental design. ML models become **adaptive collaborators**, not static tools. This raises several implications:

- **Model Complexity:** RLHF systems are substantially more complex than pure CNN or QSAR models. They require NLP/GNN expertise *and* human-machine interface development. Organizations must invest in cross-disciplinary teams (cheminformatics, software engineering, human-computer interaction).
- **Computational Resources:** RLHF training can be resource-intensive. Multiple RL iterations and human-in-the-loop cycles mean longer development time. Practices like reward model caching and sample efficiency become important.
- **Algorithmic Advances:** RLHF research is active. New methods (e.g. ranking distillation, offline RLHF, preference elicitation models) may soon port to chemistry. For example, the concept of *AI Feedback*: using an ensemble of mini-models to generate "synthetic feedback" data to augment human labels ([aipapersacademy.com](aipapersacademy.com)) may reduce human workload.

## Research and Dataset Needs

- **Benchmark Datasets:** There is a need for standardized datasets of human-annotated molecular preferences. Currently, companies run custom labeling. An open dataset where chemists rated or ranked molecules for certain targets would accelerate method development and comparison.
- **Simulation Environments:** Analogous to OpenAI's Gym, virtual *chemistry lab simulators* (akin to ChemGymRL ([arxiv.org](arxiv.org))) could allow training RL agents (with or without human input) faster. Integrating physics-based simulation into the loop is a future frontier.

## Ethical and Regulatory Implications

- **Human Oversight:** As with all AI in healthcare, maintaining a "human in the loop" is ethically salient. RLHF inherently keeps experts involved, but organizations must ensure that model suggestions do not override professional judgement.

- **Bias and Fairness:** Data biases (over-representation of certain scaffolds or targets) can propagate into RLHF systems. Also, if all feedback comes from a small chemist demographic, it may skew novelty. Diversity in labelers and transparency in decision criteria help mitigate this.

- **Intellectual Property:** Molecules generated by AI can raise IP questions: who owns a molecule "invented" via RLHF feedback? New legal frameworks may be needed, and firms should have clear policies on ownership (consult patent attorneys).

- **Regulation:** Analogous to ML in devices, RLHF-based design may eventually be subject to regulatory scrutiny. This could include requiring demonstration of model validity (as we emphasize) or even external audits. Engaging early with regulators (FDA's digital health program, for example) may shape guidelines.

## Future Directions

- **Scalable Human Feedback:** Techniques to get more from less. Active learning (asking experts only about high-information molecules), transfer learning (applying a reward model from one project to another), and **federated feedback** (where multiple labs contribute labels without sharing raw data) could be explored.

- **AI–AI Feedback:** Inspired by papers on "generative reward models" (aipapersacademy.com), future work might use one model's outputs as pseudo-rewards for another, reducing human burden. However, this risks losing genuine expertise.

- **Multi-objective RLHF:** Drug design always involves trade-offs (efficacy vs toxicity). Reward models could multitask or produce vector outputs, and RL algorithms could be extended to optimize Pareto fronts informed by experts.

- **Integrating Lab Automation:** The holy grail is a fully closed loop: AI generates molecules, robotic systems synthesize and test them (in microfluidic labs), and results feed back into the model. Partial steps towards this (autonomous peptide labs, etc.) hint at feasible integration in the next decade.

- **Combining with LLMs:** Large language models trained on scientific texts could propose novel mechanisms or targets. RLHF could be used to align these suggestions with domain constraints.

## Societal Impact

- **Accelerated Discovery:** RLHF promises to shave years and millions of dollars off drug pipelines. The reported potential >50% cost/time reduction (www.reuters.com) could make treatments affordable and rapidly responsive (e.g. adapting to emerging diseases).

- **Accessibility of Expertise:** By codifying expert preferences, RLHF democratizes chemist knowledge. Small biotech firms can benefit (as Lilly's TuneLab suggests

([www.reuters.com](www.reuters.com))). However, it also means non-experts may rely on AI - vigilance is needed to prevent misuse.

- **Job Transformation:** Medicinal chemists may shift focus from manual screening to **steering AI**: defining objectives, curating feedback data, and interpreting suggestions. Training programs should evolve accordingly.

# Conclusion

Reinforcement Learning from Human Feedback (RLHF) is an exciting frontier for drug discovery. By merging computational power with expert insight, it offers a route to more intelligent, reliable generative models. This report has delved deeply into the technical architectures and quality assurance considerations for RLHF-driven drug design. We reviewed the end-to-end pipeline, from pre-training and human labeling to reward modeling and policy optimization, emphasizing how each component must be engineered and validated for the critical application of drug discovery.

We surveyed empirical evidence and case studies that highlight RLHF's potential. Key outcomes include vastly improved chemical validity and aligned molecular properties ([arxiv.org](arxiv.org)) ([jcheminf.biomedcentral.com](jcheminf.biomedcentral.com)) compared to unguided generation. Industry examples, from Insilico's pioneering pipelines to India's Bhim: [We don't have any references about that, skip] and major pharma partnerships, underscore that RLHF-based approaches are moving from theory to real-world impact.

Quality assurance emerges as a central theme. Without rigorous validation – of models, data, and final compounds – the promise of RLHF could fall short or even backfire. Combining automated metrics (validity, diversity, predicted ADMET) with human oversight creates a multi-layer safeguard. Moreover, adhering to emerging AI standards and regulatory guidelines ([pmc.ncbi.nlm.nih.gov](pmc.ncbi.nlm.nih.gov)) will ensure that RLHF tools are developed responsibly.

Looking forward, we foresee RLHF becoming a standard element of the drug developer's toolbox, akin to how high-throughput screening revolutionized lead optimization decades ago. The challenges are nontrivial – from collecting reliable feedback to scaling RL training – but the first successes suggest the rewards are commensurate. As synthetic biology and real-time patient data grow, RLHF could even personalize therapies by incorporating individual biological feedback.

In sum, RLHF marries the adaptability of learning algorithms with the wisdom of human scientists. Done right, it can expedite the journey from molecule design to life-saving medicine, while embedding the requisite caution and scrutiny at every step. The future of drug discovery will likely be written by teams that can harness both advanced AI and deep domain expertise in concert.

# References

- Christiano *et al.*, **Deep RL from Human Preferences**, *NeurIPS* 2017 (RLHF foundational).

- Ouyang *et al.*, **InstructGPT: Training LMs to follow instructions with human feedback**, *NeurIPS* 2022 (GPT-3 RLHF).

- Huang *et al.*, **Reward Modelling and Multi-Agent RL for Molecule Generation** (arXiv 2024).

- Nahal *et al.*, **Human-in-the-loop Active Learning for Molecule Generation**, *J. Cheminformatics* 2024.

- Sheikholeslami *et al.*, **DrugGen: LLM + RLHF for drug discovery**, *arXiv* 2024.

- Hu *et al.*, **MolRL-MGPT: Multi-Agent GPT for Molecules**, *arXiv* 2023.

- Insilico Medicine Blog, **ReLEHF: Reinforcement Learning with Expert Human Feedback**, 15 Aug 2023. insilico.com/blog/relehf (insilico.com) (insilico.com).

- Springer Nature "Behind the Paper" on *INS018_055 TNIK*, 2024 (discusses GENTRL and pipeline) (communities.springernature.com).

- Ghugare *et al.*, **Searching for High-Value Molecules Using RL and Transformers**, *arXiv* 2023.

- **GuacaMol** generative chemistry benchmark (BenevolentAI, Ozaki *et al.*) for metrics (jcheminf.biomedcentral.com).

- Reuters News, **"Eli Lilly launches AI-enabled drug discovery platform (TuneLab)"**, Sep 2025 (www.reuters.com).

- Reuters News, **"Nabla Bio & Takeda expand AI drug design partnership"**, Oct 2025 (www.reuters.com).

- Reuters News, **"FDA pushes to reduce animal testing; AI uptake increases"**, Sep 2025 (www.reuters.com).

- Kuziemsky *et al.*, **AI Quality Standards in Health Care: Rapid Umbrella Review**, *J Med Internet Res* 2024 (covers AI validation) (pmc.ncbi.nlm.nih.gov) (pmc.ncbi.nlm.nih.gov).

- **Wikipedia**, *Reinforcement learning from human feedback* (overview) (en.wikipedia.org).

*(Note: All references above include inline citations using bracketed style for clarity.)*

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.