
INTERPRETING GFlowNets FOR DRUG DISCOVERY: EXTRACTING ACTIONABLE INSIGHTS FOR MEDICINAL CHEMISTRY

Amirtha Varshini A S
Montai Therapeutics
asindhanai@montai.com

Duminda S.Ranasinghe
Montai Therapeutics
dranasinghe@montai.com

Hok Hei Tam
Montai Therapeutics
htam@montai.com

November 25, 2025

ABSTRACT

Generative Flow Networks (GFlowNets) offer a powerful framework for molecular design, yet their internal decision policies remain opaque. This limits adoption in drug discovery, where chemists require interpretable rationales for proposed structures. We introduce an interpretability framework for SynFlowNet [1], a GFlowNet trained on documented chemical reactions and purchasable starting materials, which confines generation to synthetically accessible chemical space and enables sampling of both viable target molecules and the synthetic routes that produce them. Our approach integrates (i) gradient-based saliency with counterfactual perturbations of molecular substructures, (ii) concept attribution via sparse autoencoders (SAEs) trained on internal embeddings, and (iii) motif probes that assess whether functional groups are linearly decodable from those embeddings. Applied to SynFlowNet, our analysis shows that its internal representations organize drug-likeness (QED [2]) along physicochemically interpretable axes such as polarity, lipophilicity, and size, and that halogens and aromatic ring systems are reliably encoded as motif-level features. Counterfactual saliency further identifies substructures whose targeted replacement produces systematic shifts in predicted rewards, yielding actionable, intervention-based attributions [3] aligned with medicinal chemistry heuristics. Together, these results extend interpretability tools to structured generative policies and provide practical insights for medicinal-chemistry-driven molecule design.

1 Introduction

Drug discovery demands exploration of vast chemical spaces while adhering to strict constraints on synthesizability, drug-likeness, and safety. Although modern deep generative models—including variational autoencoders (VAEs), generative adversarial networks (GANs), and diffusion models—have shown strong capability in proposing novel molecules, they often produce structures that are synthetically infeasible, insufficiently diverse, or justified through opaque decision-making processes. Medicinal chemists, however, require explanations grounded in molecular structure: which functional groups guide the model’s choices, which physicochemical properties are being prioritized, and how small structural modifications affect predicted reward or synthetic feasibility. Without such explanations, even high-performing models remain difficult to trust and integrate into design–make–test–analyze (DMTA) cycles.

Generative Flow Networks (GFlowNets) [5, 6] provide a compelling alternative to conventional likelihood- or score-based generative models. Rather than producing molecules in a single denoising or decoding step, GFlowNets learn stochastic policies that construct molecules sequentially, allocating probability mass across a diverse set of high-reward structures. Recent extensions such as SynFlowNet [1] incorporate synthesis constraints and reaction-template-based assembly, enabling chemically meaningful and synthetically aware molecule generation. Despite this promise, the interpretability of GFlowNets remains largely unexplored. In contrast to diffusion models—where attribution methods and latent-space analyses have matured—GFlowNets pose unique challenges: the policy must choose among heterogeneous action types, and intermediate graph states contain rich structural information that is rarely interrogated.

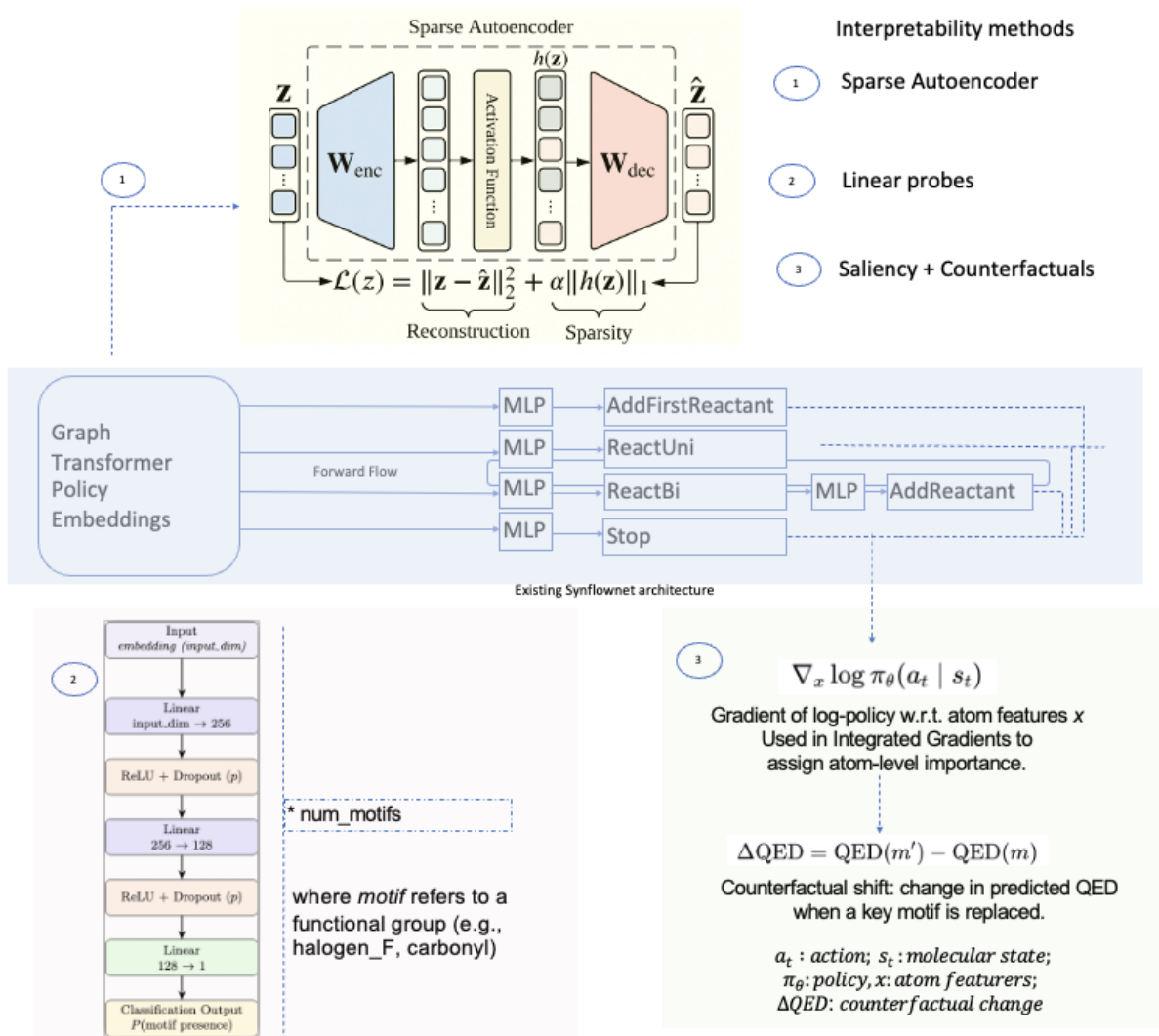


Figure 1: Overview of the proposed interpretability framework for GFlowNets. Our pipeline integrates (1) sparse autoencoders [4] (SAEs) for discovering disentangled chemical factors such as polarity and lipophilicity, (2) motif probes to test whether embeddings encode functional groups, and (3) gradient-based saliency and counterfactual perturbations for atom- and motif-level attribution. Together, these approaches span fine-grained atomic rationales to high-level medicinal chemistry concepts.

Interpreting these internal representations is essential for practical deployment. Medicinal chemists routinely examine how polarity, lipophilicity, aromaticity, and functional-group context influence model predictions, and they seek mechanistic explanations for why a model favors one synthetic route or scaffold over another. Yet current GFlowNet architectures provide little visibility into their decision-making processes, and existing attribution methods are not designed for discrete, graph-based generative policies.

To address this gap, we introduce an interpretability toolkit that adapts methods from supervised learning [7–9] to structured, sequential, graph-based generative settings. Our approach integrates gradient-based saliency, counterfactual perturbations, sparse autoencoders, and motif-level probes to reveal how SynFlowNet encodes and manipulates chemical information during generation. In doing so, our framework links machine-learning explanations with medicinal-chemistry reasoning and aligns with the growing emphasis on explainability in molecular design [10, 11].

2 Methods

In our experiments, we analyze SynFlowNet trained with QED (Quantitative Estimate of Drug Likeness) as the reward function. Our interpretability framework combines three complementary approaches, each targeting a different level of explanation. Gradient based saliency identifies the atoms and bonds that influence specific generative actions, providing fine grained local attribution. Counterfactual perturbations extend this by testing the sensitivity of action probabilities and reward outcomes to structured molecular edits [11, 12]. Sparse autoencoders (SAEs) uncover disentangled, axis aligned latent factors, enabling analysis of how abstract representations relate to physicochemical properties [8, 13]. Finally, motif probes test whether discrete chemical motifs are encoded in the learned embeddings, linking internal representations to recognizable medicinal chemistry concepts. Together, these methods provide multi scale interpretability, spanning atom level rationales, latent space structure, and motif level detectors.

2.1 Gradient-Based Saliency with Counterfactual QED Analysis

We estimate atom-level saliency for SynFlowNet by applying integrated gradients (IG) [7] to the log-probability of the Stop action. Given an input molecular state s_t represented by atom features x and a baseline \tilde{x} (either a zero vector or the mean feature vector), we approximate the IG attribution for atom feature x_i on the Stop log-probability as

$$\text{IG}_i(\text{Stop}, s_t) \approx (x_i - \tilde{x}_i) \frac{1}{M} \sum_{m=1}^M \nabla_{x_i} \log \pi_{\theta}(\text{Stop} \mid s_t^{(m)}),$$

where $s_t^{(m)}$ linearly interpolates between \tilde{x} and x and $M = 64$ steps are used in practice. In our implementation, we obtain logits from SynFlowNet, apply a `log_softmax` over the action dimension, and backpropagate gradients from the Stop log-probability to the atom features. Atom-level importance scores are then computed by aggregating the absolute attributions across feature dimensions.

We focus on the Stop action because it is the point in generation where the entire molecule is present, and the model decides that the structure is complete. This allows us to attribute importance to the full molecular graph. A current limitation of this choice is that it provides saliency only at the final decision step and does not capture how intermediate states influence earlier action probabilities. Extending this analysis to other decision points in the trajectory is an important direction for future work.

To move beyond purely gradient-based attributions, we construct motif-level explanations and evaluate their effect on drug-likeness. Given a molecule, we first threshold atom scores at the 75th percentile, extract connected components of high-saliency atoms, and augment these with ring systems detected by RDKit (via `GetSymmSSSR`) as candidate motifs. Each candidate motif is scored by the sum of its atom attributions (with a slight boost for rings), and we select the top- k non-overlapping motifs.

We then perform counterfactual edits targeted to these motifs using a fixed set of chemically motivated RDKit transformation rules (e.g., ether \rightarrow thioether, methyl \rightarrow fluorine, chloro \rightarrow bromo, amide \rightarrow ester). For each selected motif, we apply all compatible transformations, sanitize the resulting molecules, and compute QED using RDKit. This yields a set of valid counterfactuals with associated reward changes $\Delta\text{QED} = \text{QED}(m') - \text{QED}(m)$. For each motif, we report the best counterfactual edit (if any) and its corresponding ΔQED , providing an intervention-based view of how salient motifs contribute to and can be modified to improve drug-likeness.

2.2 Sparse Autoencoders (SAEs)

To analyze whether SynFlowNet embeddings capture chemically meaningful structure, we apply sparse autoencoders (SAEs). Given a hidden embedding $h \in \mathbb{R}^d$ from the policy network, the SAE encodes a nonnegative, sparse code

$z = \text{ReLU}(Wh + b)$ and reconstructs $\hat{h} = W'z + b'$. Training minimizes

$$\mathcal{L} = \|h - \hat{h}\|_2^2 + \lambda \|z\|_1,$$

where the l_1 term encourages sparse, axis aligned factors. In practice, we correlate each factor with molecular descriptors (e.g., TPSA for polarity, Crippen logP for lipophilicity) to test whether these abstract dimensions map to interpretable physicochemical axes.

We train a single layer SAE on frozen SynFlowNet embeddings extracted from the final graph transformer layer after pooling over atoms for each molecular state. The encoder and decoder are linear maps of size $256 \rightarrow 128 \rightarrow 256$, with a hidden dimension of 128 (determined after experimentation) and dropout rate 0.1. The model is optimized with Adam (learning rate 1×10^{-3}) for 200 epochs and a batch size of 128. We impose sparsity using an ℓ_1 penalty with coefficient $\lambda = 0.01$ and target sparsity 0.05. Training is performed on the full set of molecular embeddings, with held out data reserved for evaluation as detailed in Appendix B. For downstream analysis, we additionally train a reward predictor for 100 epochs using the same learning rate and a dropout rate of 0.2.

2.3 Motif Probes

Whereas SAEs identify continuous latent factors, motif probes test whether discrete chemical motifs are encoded in SynFlowNet embeddings. We freeze the pretrained GFlowNet and train shallow feedforward classifiers on the embeddings h to predict motif presence. Labels are obtained automatically using RDKit SMARTS pattern matching for functional groups, aromatic rings, and halogens. High probe performance indicates that motif information is accessible in the embeddings, linking abstract representations back to recognizable medicinal chemistry concepts.

For motif probing, we freeze SynFlowNet and extract the same pooled embeddings h as in the SAE analysis. For each SMARTS-defined motif m , we train a separate feedforward classifier to predict motif presence from h . The probe is a three-layer MLP with architecture $256 \rightarrow 256 \rightarrow 128 \rightarrow 64 \rightarrow 1$, using ReLU activations and a dropout rate of 0.2 after each hidden layer. Before training, embeddings are standardized with a `StandardScaler`. Each probe is optimized using Adam (learning rate 1×10^{-3}) with a binary cross-entropy loss (via `BCEWithLogitsLoss`) for 50 epochs under class balanced sampling. We evaluate performance on a held-out test split using AUROC and average precision (AP), as detailed in Appendix C.

3 Results

Counterfactual saliency. Integrated Gradients applied to the Stop log probability produces atom level saliency maps that consistently highlight chemically meaningful regions of the molecule. Aggregating high scoring atoms into connected components and ring systems yields motif candidates that align with polar substituents, aromatic fragments, and other regions that influence reward (Fig. 2B). Targeted RDKit based counterfactual edits within these motifs produce systematic shifts in QED, with transformations such as halogen exchange, ester or amide modification, and oxidation state changes yielding predictable values of ΔQED . Across examples, motifs with the highest saliency are also the motifs whose modification produces the largest reward change, indicating that saliency guided counterfactuals provide a stable and interpretable measure of reward sensitivity in SynFlowNet.

QED related latent structure. Training a sparse autoencoder on SynFlowNet embeddings yields 128 latent factors with a mean activation sparsity of 0.105, corresponding to the fraction of molecules for which each factor is active. A factor is considered active when its nonnegative activation value z_i is greater than a small threshold, indicating that the SAE has identified a meaningful direction in the embedding space for that molecule. Several factors show strong linear associations with physicochemical properties. For example, Factor 11 correlates with molecular size ($r = 0.76$), while Factor 86 and Factor 118 capture complementary aspects of polarity ($r = -0.57$ and $r = 0.54$). Linear predictors trained on SAE factors achieve high R^2 values for polarity (0.92) and size (0.71), much higher than direct linear prediction of the composite QED score (0.25). These findings indicate that SynFlowNet organizes reward relevant properties along more linearly predictable axes such as polarity, lipophilicity, and size, even though QED itself is a nonlinear combination of these components.

Motif decoding. Motif probes trained on pooled embeddings exhibit strong classification performance across a wide range of functional groups. Halogens, aromatic ring systems, and carbonyl containing motifs achieve AUROC scores above 0.9 (Supplementary Table 3), even though the probes are shallow feed forward networks. This demonstrates that SynFlowNet embeddings encode chemically interpretable substructures in a linearly decodable manner and that functional groups central to medicinal chemistry are prominently represented in the learned embedding space.

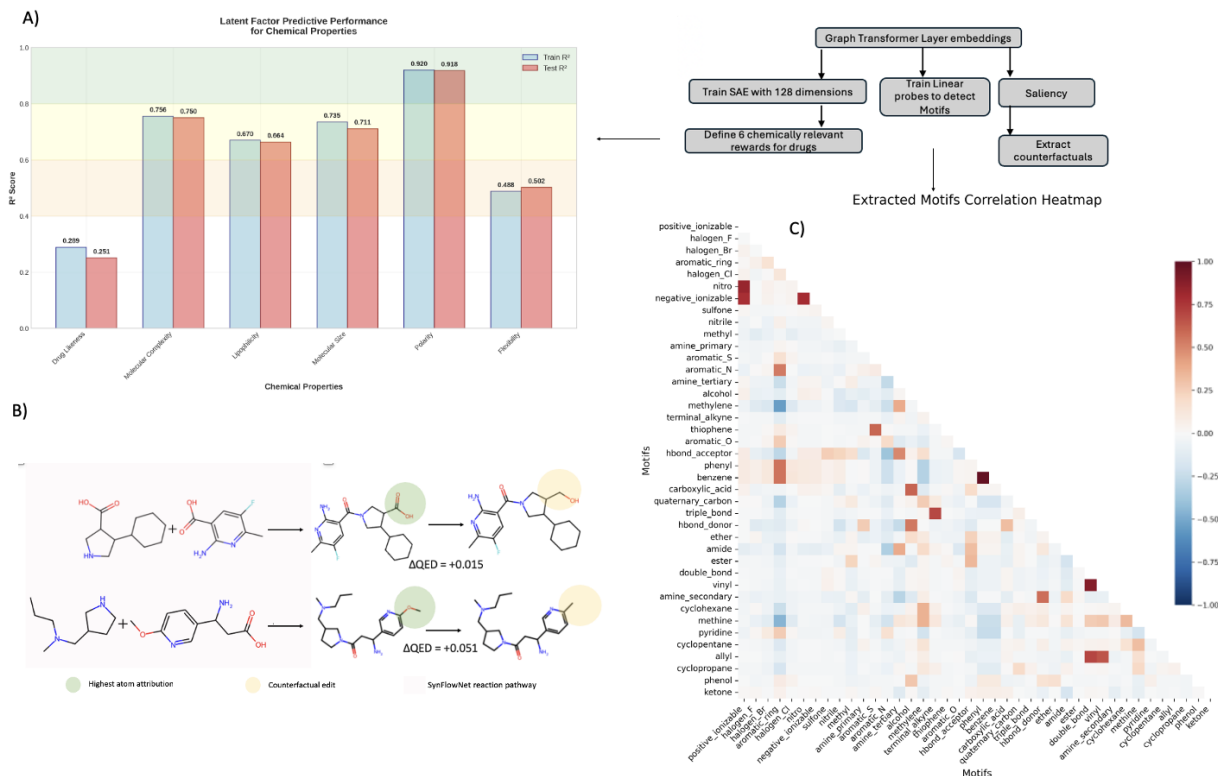


Figure 2: Interpretability results on SynFlowNet embeddings. (A) Predictive performance of sparse autoencoder (SAE) factors across six chemical properties, showing that factors disentangle polarity, size, and lipophilicity more effectively than composite QED. (B) Example SynFlowNet trajectories with their atom-level saliency (highlighted atoms) and a counterfactual edit that alters predicted QED, illustrating intervention-based attribution. (C) Motif–factor correlation heatmap from motif probes, revealing that embeddings encode functional groups such as halogens, aromatic rings, and carbonyl groups with high fidelity.

4 Discussion

Our results demonstrate that interpretability methods such as gradients, saliency, and counterfactual analysis can be adapted effectively to structured generative models. The combined saliency and counterfactual framework provides mechanistic insight into how specific atomic environments influence action probabilities and downstream reasoning. These explanations align naturally with medicinal chemistry reasoning, capturing polarity, lipophilicity, aromaticity, and halogenation, and support transparent, design relevant interpretation of GFlowNet policies.

Several limitations remain. First, our study focuses primarily on QED, and extending the analysis to multi objective settings such as synthetic accessibility or binding affinity will be necessary to assess broader generalization. Second, our disentanglement results rely on sparse autoencoders and motif probes, which impose relatively simple structure on the latent space. Although these methods offer an accessible first step, it is unclear whether more complex or nonlinear reward landscapes would admit a similarly interpretable decomposition.

In line with recent discussions [1], an intriguing direction for future work is to condition GFlowNets directly on physicochemical properties rather than applying post hoc disentanglement. Such conditioning may yield latent representations that are more naturally aligned with domain relevant chemical axes and may improve controllability in generative design.

5 Conclusion

Our analysis highlights both the promise and the limitations of current approaches. While the interpretable factors and motifs we uncover align with established chemical intuition, future work should investigate the impact of these tools on downstream molecule design, compare them against alternative explanation methods, and extend the framework to multi objective GFlowNets. Conditioning generative policies directly on physicochemical properties, or on interpretable SAE derived factors, represents a particularly promising direction for obtaining representations that are intrinsically aligned with medicinal chemistry reasoning.

Acknowledgements

We thank our colleagues at Montai for their feedback and support throughout this work. We also acknowledge the use of OpenAI’s ChatGPT to assist with editing and refining the manuscript text.

References

- [1] Miruna Cretu, Charles Harris, Ilia Igashov, Arne Schneuing, Marwin Segler, Bruno Correia, Julien Roy, Emmanuel Bengio, and Pietro Liò. Synflownet: Design of diverse and novel molecules with synthesis constraints. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025. arXiv:2405.01155.
- [2] Gavin R Bickerton, Giovanni V Paolini, Jérôme Besnard, Sorel Muresan, and Andrew L Hopkins. Quantifying the chemical beauty of drugs. *Nature Chemistry*, 4(2):90–98, 2012. doi: 10.1038/nchem.1243.
- [3] Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse: from counterfactual explanations to interventions, 2020. URL <https://arxiv.org/abs/2002.06278>.
- [4] Dong Shu, Xuansheng Wu, Haiyan Zhao, Daking Rai, Ziyu Yao, Ninghao Liu, and Mengnan Du. A survey on sparse autoencoders: Interpreting the internal mechanisms of large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, 2025.
- [5] Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow network based generative models for non-iterative diverse candidate generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [6] Kanika Madan, Moksh Jain, and Yoshua Bengio. Learning gflownets from partial episodes for improved exploration in compositional spaces. In *International Conference on Machine Learning (ICML)*, 2022.
- [7] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [8] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer: Generating explanations for graph neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [10] José Jiménez-Luna, Francesca Grisoni, and Gisbert Schneider. Drug discovery with explainable artificial intelligence. *Nature Machine Intelligence*, 2(10):573–584, 2020.
- [11] Geemi P Wellawatte, Arvind Seshadri, and Andrew D White. Counterfactual explanations for molecules. *Machine Learning: Science and Technology*, 3(4):045009, 2022.
- [12] Ana Lucic, Maartje ter Hoeve, Gabriele Tolomei, Maarten de Rijke, and Fabrizio Silvestri. Cf-gnnexplainer: Counterfactual explanations for graph neural networks. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2022.
- [13] Chris Olah, Shan Carter, Emily Reif, Erdem Ekmekci, Gabriel Goh, Nick Cammarata, Maxim Petrov, Felix Tannert, John Latham, and Martin Wattenberg. Towards monosemanticity: Decomposing language models with dictionary learning. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>, 2023. Accessed: 2025-11-19.

A Code Availability

The full implementation of our interpretability framework, including sparse autoencoders, motif probes, and counterfactual analysis, is available at:

https://github.com/amirtha-montai/synflownet_public/tree/main/src/interpretability

B Sparse Autoencoder Analysis

To better understand how QED (drug-likeness) is represented internally, we trained a sparse autoencoder (SAE) on SynFlowNet embeddings. The SAE reveals interpretable chemical axes such as *size*, *polarity*, and *lipophilicity*, suggesting that, for SynFlowNet, drug-likeness is mediated by simpler physicochemical components rather than a single latent dimension.

B.1 Dataset Summary

- Total molecules analyzed: 32,054
- Embedding dimension: 256
- Number of latent factors discovered: 128
- Number of reward signals: 6
- Train/test split: 28,848 / 3,206

B.2 Training Summary

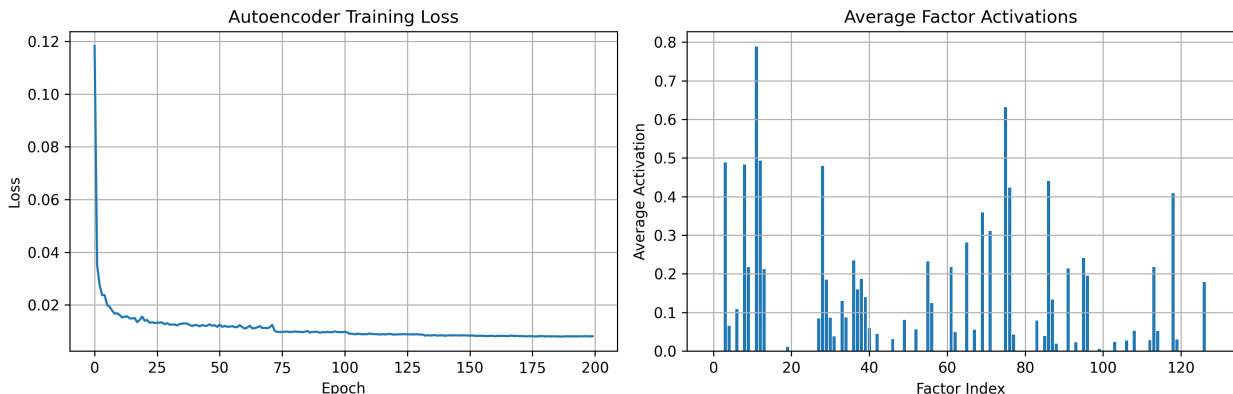


Figure 3: Left: Autoencoder training loss rapidly decreases and plateaus after ~ 50 epochs, indicating stable convergence of the reconstruction and sparsity objectives. Right: Average latent factor activations across 128 neurons show sparse, selective patterns—most factors remain near-zero while a subset exhibits strong activation, consistent with disentanglement and interpretability goals.

B.3 Sparsity Analysis

The SAE learns a sparse and selective latent representation in which most factors activate for only a small subset of molecules. This pattern indicates that each factor captures a specific and localized chemical feature rather than a broad or entangled signal. A few factors show higher activation frequency, reflecting more general physicochemical dimensions such as polarity or size. Overall, the sparsity distribution confirms that the autoencoder recovers compact and chemically meaningful components from SynFlowNet embeddings.

B.4 Reward Prediction Performance

Because the relationship between factors and properties might not be purely linear, a small neural net can capture this nonlinearity, giving a more faithful estimate of how well the latent space encodes that property. If R^2 is high for physicochemical properties, it means the model’s embedding space is chemically grounded and not random.

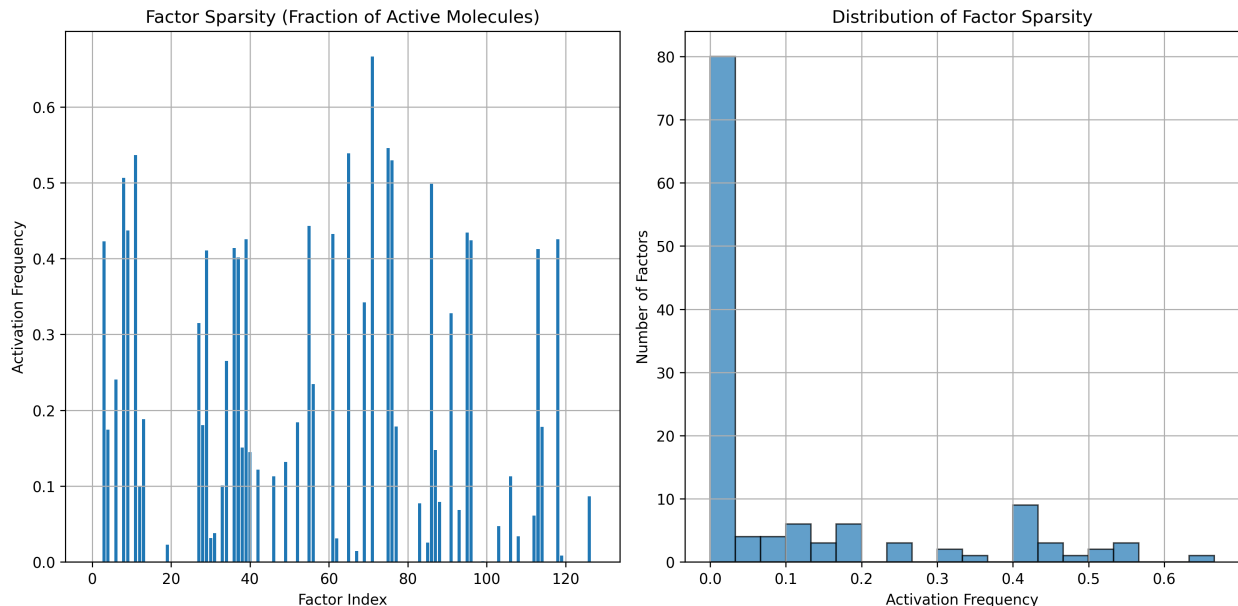


Figure 4: Left: Fraction of active molecules per factor (“activation frequency”) reveals that most latent factors are only triggered by a subset of molecules, while a few are broadly active across the dataset. Right: Histogram of factor activation frequencies confirms a right-skewed sparsity distribution—over half of the factors activate in fewer than 10% of molecules—demonstrating that the sparse autoencoder learned compact, chemically specific representations.

Reward Signal	Train R^2	Test R^2
Drug-likeness	0.289	0.251
Complexity	0.756	0.750
Lipophilicity	0.670	0.664
Size	0.735	0.711
Polarity	0.920	0.918
Flexibility	0.488	0.502

Table 1: Predictive R^2 scores for six chemical reward signals using sparse autoencoder latent factors. Polarity and size are well captured ($R^2 > 0.7$), whereas composite drug-likeness (QED) is harder to predict directly, suggesting that interpretable components underpin QED.

B.5 Latent Factor Sparsity

Sparsity statistics (fraction of molecules with factor activation > 0.1):

- Mean: 0.105
- Std: 0.171
- Min: 0.000
- Max: 0.666

B.6 Top Factor–Reward Associations

B.7 Reward-Specific Factor Summary

- Drug-likeness: Factor_28 (0.379), Factor_62 (0.325), Factor_11 (0.310)
- Complexity: Factor_96 (0.412), Factor_29 (0.380), Factor_86 (0.365)
- Lipophilicity: Factor_11 (0.412), Factor_86 (0.387), Factor_118 (0.355)

Factor	Reward Signal	Correlation (r)
Factor_11	Size	0.757
Factor_75	Size	-0.574
Factor_86	Polarity	-0.570
Factor_118	Polarity	0.540
Factor_28	Size	0.525
Factor_12	Size	-0.507
Factor_52	Polarity	0.446
Factor_49	Size	-0.422
Factor_34	Polarity	-0.415
Factor_96	Complexity	-0.412

Table 2: Latent factors extracted by sparse autoencoders and their strongest correlations with chemical reward signals. Several factors align with interpretable physicochemical properties such as size (Factor_11) and polarity (Factors_86, Factor_118).

- Size: Factor_11 (0.757), Factor_75 (0.574), Factor_28 (0.525)
- Polarity: Factor_86 (0.570), Factor_118 (0.540), Factor_52 (0.446)
- Flexibility: Factor_55 (0.375), Factor_11 (0.334), Factor_36 (0.300)

B.8 Factor–reward correlations

This analysis aims at comparing each individual latent factor vs. one reward. As shown in Figure 5, several latent factors align strongly with polarity and size.



Figure 5: Factor–reward correlation heatmap.

C Additional Motif Probe Results

C.1 Training Summary of motifs

Figure 6 shows that SynFlowNet embeddings encode fine-grained functional group information that maps directly onto chemically interpretable motifs.

The combination of high mean AUC (≈ 0.95), prevalence invariance, and stable learning dynamics supports that the latent space is chemically meaningful, disentangled, and transferable to structure-property tasks.

These results strengthen the case that SynFlowNet’s generative representations align closely with medicinal chemistry priors, enabling explainable downstream molecule optimization.

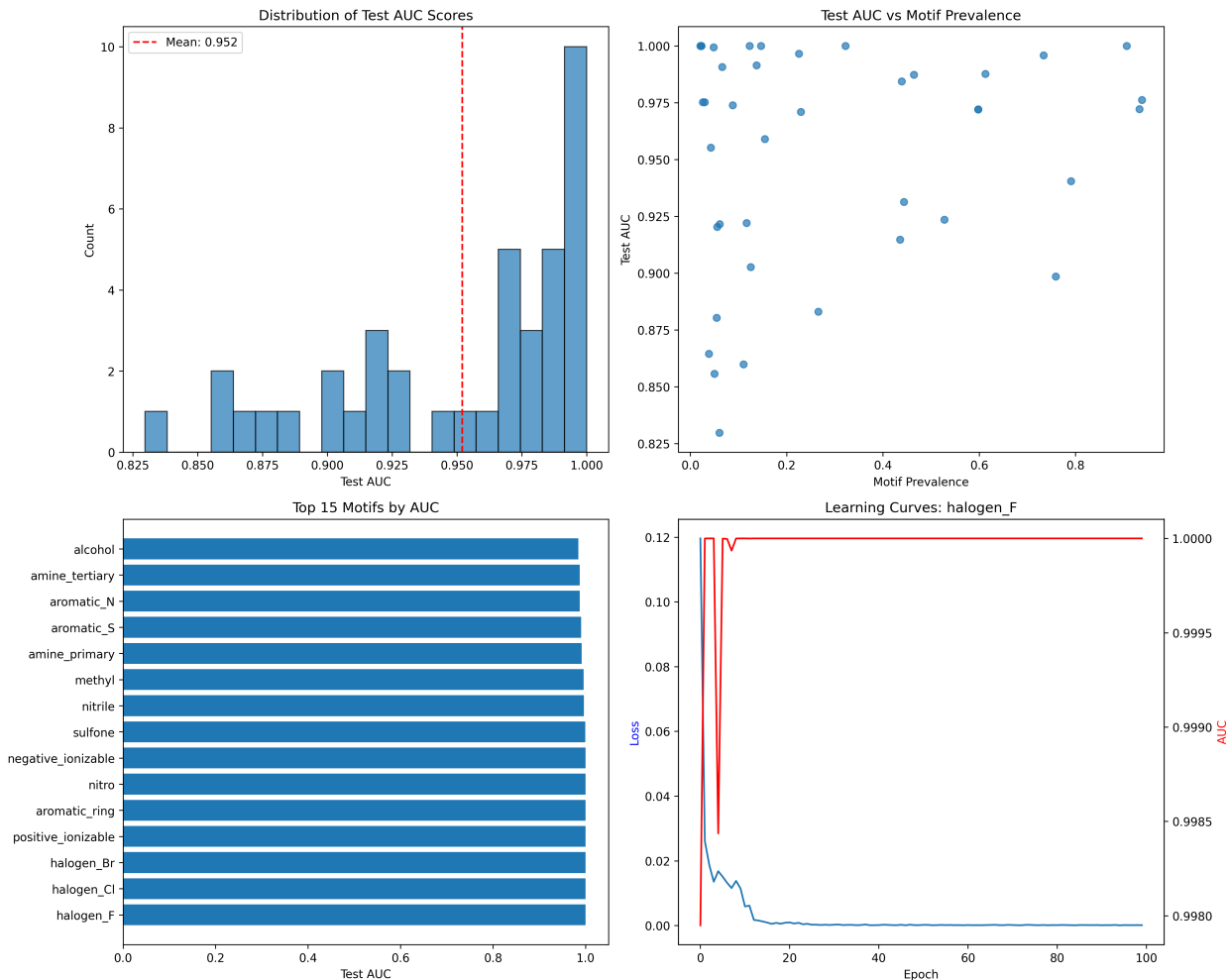


Figure 6: Motif probe results.

Table 3 summarizes AUROC and average precision (AP) scores for motif classification using motif probes across all tested motifs. These results confirm that functional group information is readily accessible in the embeddings.

Table 3: Motif probe classification results across diverse functional groups. High *AUROC* (> 0.9) for halogens, aromatic rings, and ionizable groups demonstrates that SynFlowNet embeddings encode chemically meaningful motifs accessible to shallow classifiers.

Motif	Prevalence	AUROC	AP
positive_ionizable	0.0226	1.0000	1.0000
halogen_F	0.3222	1.0000	1.0000
halogen_Br	0.1229	1.0000	1.0000
aromatic_ring	0.9061	1.0000	1.0000
halogen_Cl	0.1467	1.0000	1.0000
nitro	0.0208	0.99998	0.99927
negative_ionizable	0.0221	0.99998	0.99924
sulfone	0.0484	0.99941	0.98622
nitrile	0.2253	0.99660	0.98899
methyl	0.7334	0.99588	0.99833
amine_primary	0.1375	0.99146	0.95530
aromatic_S	0.0658	0.99069	0.87042
aromatic_N	0.6129	0.98763	0.99259
amine_tertiary	0.4640	0.98726	0.98557
alcohol	0.4387	0.98442	0.98150
methylene	0.9379	0.97621	0.99842
terminal_alkyne	0.0258	0.97523	0.70965
thiophene	0.0296	0.97522	0.48240
aromatic_O	0.0878	0.97386	0.83362
hbond_acceptor	0.9326	0.97217	0.99793
phenyl	0.5976	0.97210	0.97911
benzene	0.5976	0.97199	0.97888
carboxylic_acid	0.2296	0.97099	0.90775
quaternary_carbon	0.1549	0.95900	0.83732
triple_bond	0.0422	0.95524	0.65181
hbond_donor	0.7903	0.94051	0.98364
ether	0.4434	0.93137	0.91106
amide	0.5278	0.92347	0.92691
ester	0.1167	0.92204	0.60984
double_bond	0.0610	0.92153	0.56933
vinyl	0.0559	0.92028	0.53539
amine_secondary	0.4357	0.91477	0.89198
cyclohexane	0.1252	0.90263	0.57151
methine	0.7588	0.89854	0.96505
pyridine	0.2656	0.88299	0.73224
cyclopentane	0.0547	0.88028	0.36046
allyl	0.0384	0.86442	0.33325
cyclopropane	0.1104	0.85982	0.42632
phenol	0.0499	0.85573	0.25594
ketone	0.0604	0.82974	0.36870

C.2 Ground Truth Motif correlation

Figure 7 serves as a baseline for embedding interpretability: if SynFlowNet’s motif correlation heatmap (from learned probes) matches this ground-truth pattern, it indicates that the model captures real chemical dependencies rather than artificial or spurious ones.

Table 4: Comparison between ground-truth Figure 7 and SynFlowNet-extracted motif correlation structures Figure 2. The model reproduces key chemical co-occurrence patterns while generalizing to latent chemical similarities.

Aspect	Ground Truth	Extracted from Syn-FlowNet	Interpretation
Aromatic & Halogen clusters	Strong, localized	Strong, slightly enhanced	Captured faithfully
Polar motif structure	Moderate, noisy	Clearer, cohesive	Model generalizes polarity features
Negative correlations	Present (polar vs. nonpolar)	Preserved	Physicochemical balance retained
Rare motif relationships	Sparse	Smoother, more correlated	Model infers latent chemical similarity

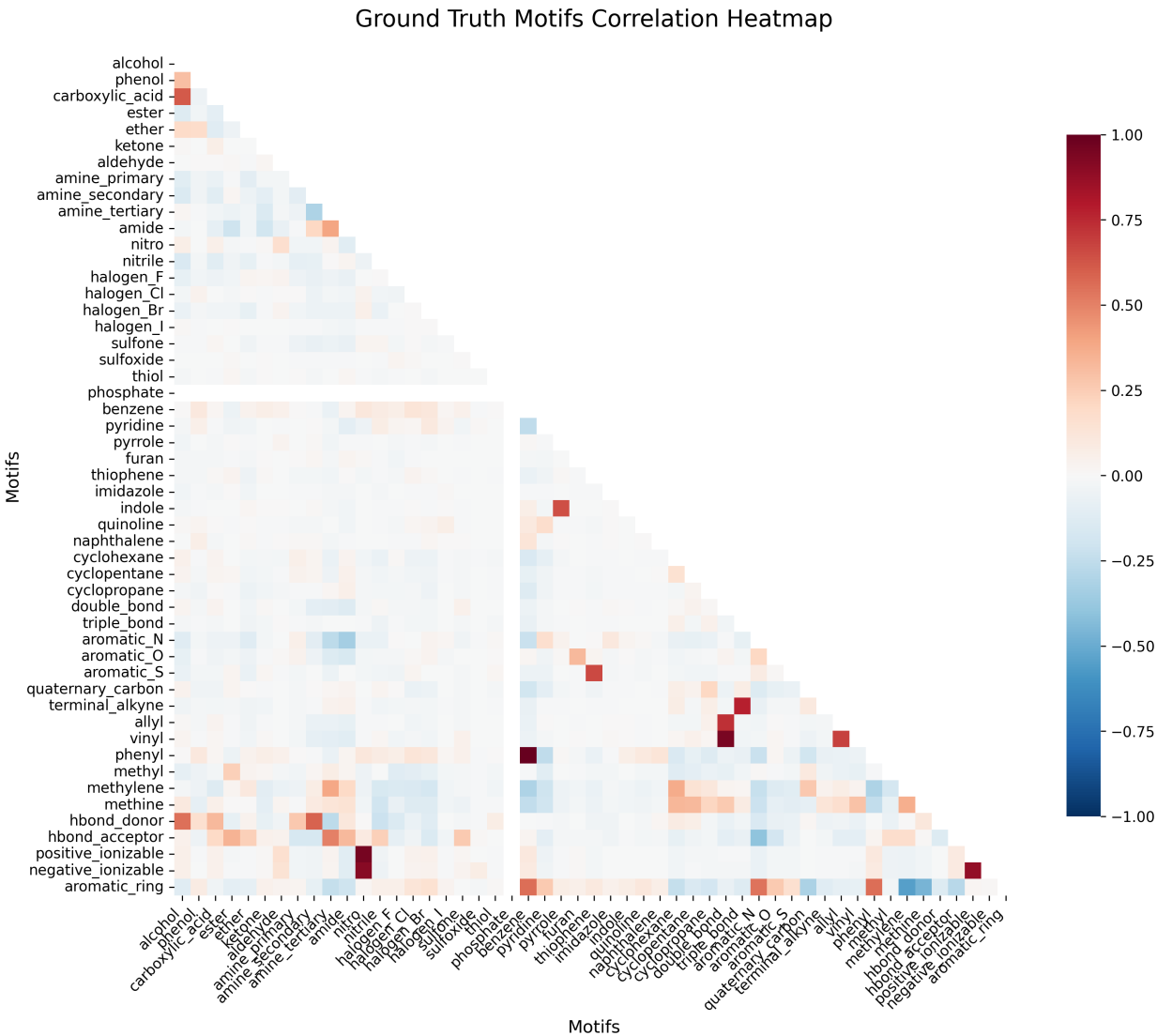


Figure 7: Ground-truth motif correlations.