

# Generative Diffusion Models for De Novo Catalyst Design Modulating Hydrocarbon Fuel Production

## Generative Diffusion Models for De Novo Catalyst Design: Modulating Hydrocarbon Fuel Production

```
tags: [[AI]] [[Catalysis]] [[DiffusionModels]] [[MolecularGeneration]] [[DeNovoDesign]] [[FischerTropsch]]  
[[HydrocarbonSynthesis]] [[ComputationalChemistry]] [[InfrastructureCosts]]
```

### Abstract

This paper outlines a research strategy for leveraging advanced generative diffusion models, inspired by breakthroughs in protein design, to accelerate the *de novo* discovery of novel catalysts. The primary focus is on modulating the output of the Fischer-Tropsch (F-T) synthesis, a critical reaction for converting syngas (CO and H<sub>2</sub>) into hydrocarbon chains for fuel applications. We propose a methodology for generating catalyst isomers and novel structures with predetermined attributes, such as hydrocarbon chain length and density, crucial for specific fuel applications like aviation. The paper details the methodological setup, infrastructure requirements, training strategies, and provides a multi-order-of-magnitude cost analysis for hardware acquisition or rental. It also estimates the feasible scale of research compound sets and the necessary training data, extending beyond simple diffusion models to incorporate cutting-edge methodologies for conditional and controllable generation.

### 1. Introduction

The efficient and selective production of hydrocarbons from abundant feedstocks remains a grand challenge in chemical engineering, with profound implications for energy security and sustainability. The Fischer-Tropsch (F-T) synthesis, which converts syngas (CO and H<sub>2</sub>) into liquid hydrocarbons, stands as a cornerstone technology. However, controlling the product distribution—specifically, the length and density of the hydrocarbon chains—is paramount for maximizing the yield of high-value fuels. Traditional catalyst discovery is often a laborious, trial-and-error process.

Recent advancements in generative artificial intelligence, particularly denoising diffusion probabilistic models (DDPMs), have revolutionized *de novo* design in fields such as protein engineering (e.g., RFDiffusion). This report posits that an analogous paradigm can be applied to heterogeneous catalyst discovery, enabling the rational, goal-oriented design of materials with bespoke properties. This research aims to develop a computational framework capable of generating novel catalyst structures and their isomers, precisely tuned to yield hydrocarbons of desired length and density from the F-T reaction, and to explore its potential for paraffin cracking.

### 2. Reaction Chemistry Clarification

The initial inquiry referenced a reaction between hydrogen (H<sub>2</sub>) and nitrogen monoxide (NO) to form hydrocarbon chains. While H<sub>2</sub> and NO can react, their primary products are typically nitrogen (N<sub>2</sub>) and water (H<sub>2</sub>O), or ammonia (NH<sub>3</sub>), in processes aimed at NO<sub>x</sub> reduction. The formation of hydrocarbon chains from H<sub>2</sub> and NO is not a chemically straightforward or established pathway for fuel synthesis.

Based on the attached paper and the industrial goal of hydrocarbon fuel production, the Fischer-Tropsch (F-T) synthesis is the chemically accurate and industrially relevant reaction for this research:



This reaction directly produces hydrocarbon chains ( $C_nH_{(2n+2)}$ ) and water from syngas (CO and  $H_2$ ). The ability to modulate the output density and chain length is a well-known challenge in F-T catalysis, making it an ideal target for AI-driven design. This paper will proceed with the F-T synthesis as the core reaction. Additionally, the potential for cracking paraffins as a further source of hydrocarbon elements will be considered as a related, but distinct, application requiring similar generative design principles for catalysts.

### 3. Proposed Methodology: CatDiff for Hydrocarbon Modulation

Drawing inspiration from RFDiffusion's success in protein design, our proposed methodology, termed "CatDiff," will employ a conditional generative diffusion framework for heterogeneous catalyst discovery.

#### 3.1. Catalyst Representation: The Multi-Modal $C_{rep}$

A critical component is the development of a comprehensive, multi-modal representation for catalysts, denoted as  $C_{rep}$ . This goes beyond simple atomic coordinates to capture the intricate interplay of factors governing catalytic activity:

- **Geometric Representation (3D Coordinates):** A 3D point cloud defining the precise atomic positions of the active metal nanoparticles (e.g., Fe, Co, Ni, Ru) and their interaction with the support material (e.g.,  $SiO_2$ ,  $Al_2O_3$ ). This captures the physical structure of the active sites.
- **Electronic Representation (Volumetric/Per-Atom Features):** A volumetric grid or per-atom feature vectors encoding key electronic properties. This includes the d-band center of surface atoms, local electrostatic potential, and electron density, which are known to dictate chemical reactivity.
- **Energetic Representation (Adsorption Energy Distribution - AED):** A functional fingerprint representing the calculated binding energies of critical F-T reaction intermediates (e.g., CO, H,  $CH_x$ , OH) across various potential active sites on the catalyst surface. This directly informs how the catalyst interacts with reactants and intermediates.

This holistic  $C_{rep}$  allows the diffusion model to learn deep, non-linear relationships between structure, electronic properties, and energetic landscapes, enabling the generation of catalysts with specific functional outcomes.

#### 3.2. Denoising Network and Pre-training: The "Catalyst Foundation Model"

The core of CatDiff will be a 3D-equivariant Graph Neural Network (GNN) or transformer architecture serving as the denoising network ( $\epsilon_{\theta}$ ). Equivariance is crucial to ensure that the model's predictions are independent of the molecule's orientation in space.

The most significant undertaking is the pre-training of this network on a massive database of catalyst structures and their corresponding properties, primarily generated via high-throughput Density Functional Theory (DFT) simulations. This pre-training will create a "Catalyst Foundation Model"—a network imbued with the fundamental "physics" of surface chemistry, catalysis, and materials stability, analogous to how AlphaFold learned protein folding. This learned prior will guide the generative process towards chemically plausible and catalytically relevant structures.

#### 3.3. Conditional Generation and Guidance

The generative process will be a conditional diffusion model, learning the distribution  $p_{\theta}(C_{rep,t+1}|C_{rep,t},c)$ . The condition,  $c$ , will be the desired functional outcome:

- **Target Hydrocarbon Properties:** Quantifiable metrics such as the target Anderson-Schulz-Flory (ASF) chain growth probability ( $\alpha_{target}$ ), indicating desired chain length distribution (e.g.,  $\alpha=0.9$  for longer chains like diesel).
- **Density Alignment:** Specific target density ranges for the final hydrocarbon product, ensuring alignment with existing fuel specifications for industrial engines (e.g., jet fuel density ranges). This would be achieved by integrating a property predictor for product density into the guidance mechanism.
- **Paraffin Cracking Specificity:** For the paraffin cracking application, conditions would specify desired product distribution (e.g., light olefins, gasoline range hydrocarbons) and resistance to coking.

A separate, smaller "property prediction" network will be trained alongside the main model. This predictor will take a generated  $C_{rep}$  as input and predict its resulting performance (e.g.,  $\alpha_{predicted}$ ,  $density_{predicted}$ ). During the iterative generative sampling, the gradient of the error between the target and predicted properties (e.g.,  $\alpha_{target} - \alpha_{predicted}$ ) will be

used to "steer" each denoising step. This "guidance" mechanism will push the generation towards structures predicted to achieve the desired catalytic function.

### 3.4. Iterative Workflow

- ① **Define Target:** Specify desired hydrocarbon properties (e.g.,  $\alpha=0.9$ , density range) and base materials (e.g., Co on SiO<sub>2</sub>).
- ② **Initialize:** Start with a random, multi-modal `C_rep` (random atom positions, electronic/energetic values).
- ③ **Iterative Denoising:** Apply guided reverse diffusion steps. The Catalyst Foundation Model proposes physically plausible corrections, and the property guidance adjusts these corrections to steer towards the functional target.
- ④ **Output:** A novel, multi-modal `C_rep` representing a computationally predicted catalyst, including its 3D atomic structure, electronic properties, and energetic landscape.
- ⑤ **Post-processing & Validation:** Convert `C_rep` to standard formats (e.g., CIF, XYZ), perform classical energy minimization, and conduct rapid DFT calculations on generated candidates for further validation.
- ⑥ **Experimental Synthesis & Testing:** The most promising candidates are synthesized and experimentally validated, with results feeding back into the training data for continuous model improvement (active learning).

## 4. Infrastructure Requirements and Cost Analysis

The computational demands for training large-scale diffusion models, especially those operating on complex 3D molecular data and requiring extensive DFT pre-computation, are substantial.

### 4.1. Hardware Components

- **GPUs (Graphics Processing Units):** The primary compute resource for deep learning. High-end data center GPUs (e.g., NVIDIA A100, H100) are essential for efficient training.
- **High-Bandwidth Memory (HBM):** Crucial for large models and batch sizes.
- **CPUs:** For data loading, preprocessing, and orchestrating GPU tasks.
- **RAM:** Sufficient system RAM to feed data to GPUs and handle large intermediate tensors.
- **Storage (SSD/NVMe):** Fast storage for datasets and model checkpoints. DFT data can be terabytes in size.
- **Networking:** High-speed interconnects (e.g., InfiniBand) for multi-GPU and multi-node training.

### 4.2. Cost Orders of Magnitude and Compute Usage

Scenario 1: Entry-Level Research (\$500 - \$5,000 expenditure / Couple of Days - Weeks compute)

- **Hardware:** Access to a single high-end consumer GPU (e.g., NVIDIA RTX 4090, ~\$1,800 - \$2,500) or a cloud instance with a single A100 40GB (rental: ~\$1-3/hour).
- **Compute Usage:** A few days to 2-3 weeks of continuous training.
- **Model Scale:** This budget would realistically support training a smaller-scale molecular diffusion model, likely in the order of tens of millions to a few hundred million parameters (e.g., 50M-500M parameters). This is significantly smaller than the LLM scales (7B+) but still substantial for molecular data due to its complexity.
- **Training Data:** Limited to pre-existing, publicly available datasets (e.g., subsets of QM9, GDB-17, or smaller custom DFT datasets). Perhaps 10,000 - 100,000 unique catalyst structures with basic properties.
- **Expected Results:** Proof-of-concept generation of simple catalyst structures or isomers with coarse control over a few attributes. The generated compound set size would be in the hundreds to low thousands of unique, chemically valid candidates. The utility would be primarily for demonstrating feasibility and exploring basic design principles.

Scenario 2: Mid-Range Research (\$5,000 - \$50,000 expenditure / Weeks - 2-3 Months compute)

- **Hardware:** Purchasing 2-4 NVIDIA A100 80GB GPUs (~\$15,000 - \$30,000 per GPU, so ~\$30,000 - \$120,000 for a server) or renting cloud instances with 4-8 A100 80GB GPUs (rental: ~\$10-25/hour).
- **Compute Usage:** Several weeks to 2-3 months of continuous training.
- **Model Scale:** This budget allows for training a more sophisticated CatDiff model, potentially in the range of 500 million to 2-3 billion parameters. This would enable learning more complex relationships and finer-grained control.
- **Training Data:** Access to larger, curated public datasets, and the ability to perform a moderate amount of *de novo* DFT calculations (e.g., 100,000 - 500,000 DFT calculations) to enrich the dataset with specific catalyst types or reaction intermediates.

- **Expected Results:** Generation of tens of thousands to hundreds of thousands of novel catalyst candidates. The model could show promising control over hydrocarbon chain length distribution (alpha) and initial density modulation. This level would likely produce results of "usable type" for *in silico* screening and lead generation, requiring experimental validation for top candidates.

### Scenario 3: High-End Research / Dedicated Cluster (\$50,000+ expenditure / Several Months compute)

- **Hardware:** Purchasing a dedicated server with 8x NVIDIA H100 80GB GPUs (~\$300,000 - \$500,000+) or renting a large cluster with 16-64 H100 GPUs (rental: ~\$50-200+/hour). Significant investment in fast networking and petabytes of storage.
- **Compute Usage:** Several months (3-6+ months) of continuous training.
- **Model Scale:** This scale could support training a CatDiff model with 5-10 billion parameters, pushing towards the lower end of "large" models, specifically tailored for molecular generation. Models of this size for molecular data are still highly experimental and at the cutting edge.
- **Training Data:** The ability to generate massive, high-fidelity DFT datasets (e.g., millions of DFT calculations) specifically tailored to the F-T reaction and desired hydrocarbon properties. This is a multi-million dollar undertaking in itself if done from scratch.
- **Expected Results:** Generation of millions to tens of millions of chemically diverse and highly optimized catalyst candidates. The model could achieve very fine-grained control over hydrocarbon length and density, potentially discovering truly novel catalytic mechanisms. This level aims for direct impact on industrial processes, significantly reducing experimental costs and time.

**Note on Model Parameter Scales:** The parameter counts (7B, 14B, etc.) are typically associated with Large Language Models (LLMs) which operate on discrete tokens and have different architectural requirements. Molecular diffusion models, even "large" ones, are usually in the hundreds of millions to low billions of parameters. A 7B parameter molecular diffusion model would be exceptionally large and would require compute resources far exceeding the \$50,000 budget for more than a few days. The provided ranges reflect realistic scales for molecular generative models given the budget constraints.

## 5. Training Data Necessary for Usable Results

The quality and quantity of training data are paramount. For CatDiff, this involves:

- **Existing Catalysis Databases:** Initial bootstrapping from public databases of known catalysts (e.g., Open Catalyst Project, Materials Project, NOMAD) for general material properties and stability.
- **High-Throughput DFT Simulations:** This is the most crucial and expensive data generation step. To train a robust CatDiff model capable of modulating hydrocarbon output, a massive dataset of DFT-calculated properties for various catalyst compositions and active site geometries is required. This includes:
  - **Adsorption Energies:** For key F-T intermediates (CO, H, CH<sub>x</sub>, O, H<sub>2</sub>O) on different catalyst facets and defect sites.
  - **Reaction Barriers:** Activation energies for elementary steps of the F-T mechanism (e.g., C-O bond dissociation, C-C coupling, H-addition).
  - **Electronic Descriptors:** d-band centers, work functions, charge transfer.
  - **Structural Information:** Atomic coordinates, lattice parameters, surface reconstructions.
  - **Simulated Reaction Outcomes:** For training the property prediction network, simulated or experimental data linking catalyst structure to hydrocarbon product distribution (ASF alpha, density).
- **Data Size:**
  - **Usable Type (Mid-Range):** For "usable type" results (leading to promising candidates for experimental validation), a dataset of hundreds of thousands to a few million DFT-calculated catalyst-reaction snapshots would be necessary. This would likely involve a diverse set of transition metal catalysts (Fe, Co, Ru, Ni) on various supports, with different promoter loadings.
  - **Cutting-Edge (High-End):** To achieve truly novel, highly optimized catalysts with fine-grained control, a dataset spanning tens of millions of DFT calculations would be ideal, covering a much broader chemical space and more complex active site motifs.

## 6. Cutting-Edge Diffusion Model Methodologies

Beyond the basic DDPM framework, several cutting-edge methodologies will be crucial for CatDiff's success:

- **Equivariant Diffusion Models (E(3) EDMs):** As highlighted in the attached paper's discussion of AlphaFold's architecture, molecular data possesses inherent symmetries (rotation, translation). E(3) EDMs are designed to be invariant to these transformations, ensuring that the generated molecules are chemically meaningful regardless of their orientation. This is critical for 3D catalyst generation.
- **Conditional and Controllable Generation:** Advanced conditioning mechanisms (e.g., classifier-free guidance, prompt-based conditioning similar to DpDNet or OmniGen) will be essential to steer the generation towards specific hydrocarbon properties (length, density). This

allows the model to interpret complex instructions like "Generate a cobalt-based catalyst that produces C<sub>12</sub>-C<sub>20</sub> hydrocarbons with a density of 0.78 g/cm<sup>3</sup>."

- **Latent Diffusion Models:** Operating in a compressed latent space can significantly improve computational efficiency and allow for larger effective model sizes, making training and inference faster.
- **Hierarchical Diffusion:** Generating catalysts in a multi-scale fashion (e.g., first a rough scaffold, then refining active site details) can improve stability and chemical validity.
- **Integration with Property Predictors:** Tightly coupling the diffusion model with a high-fidelity property prediction model (e.g., a GNN trained on DFT data) allows for real-time guidance during the generative process, pushing the model towards desired functional outcomes.
- **Active Learning Loops:** Implementing a feedback loop where computationally generated candidates are screened (via rapid DFT or even experimental synthesis), and the most informative results are used to retrain or fine-tune the diffusion model, continuously improving its performance and reducing the need for purely random exploration.

## 7. Potential Results and Impact

The successful implementation of CatDiff would yield transformative results:

- **Novel Catalyst Discovery:** Generation of entirely new catalyst compositions and active site geometries previously unexplored by traditional methods.
- **Targeted Hydrocarbon Production:** Catalysts precisely designed to produce hydrocarbon chains of specific lengths and densities, enabling the direct synthesis of aviation fuels, diesel, or specialized lubricants from syngas.
- **Enhanced Selectivity and Efficiency:** Catalysts optimized for higher conversion rates and reduced undesirable byproducts (e.g., methane), improving the economic viability of F-T synthesis.
- **Paraffin Cracking Catalysts:** Extension of the methodology to design catalysts for cracking heavier paraffins into lighter, more valuable hydrocarbons, providing additional flexibility in fuel production.
- **Accelerated R&D:** Dramatically reducing the time and cost associated with experimental catalyst screening, shifting the paradigm from trial-and-error to rational, *de novo* design.
- **Fundamental Insights:** The trained CatDiff model would implicitly encode a deep understanding of structure-property relationships in catalysis, providing new scientific insights into reaction mechanisms and material design principles.

The ability to precisely modulate fuel density is critical for industrial engines like those in aviation, where weight and energy density are paramount. A generative model that can design catalysts for this purpose would offer a significant strategic advantage in fuel formulation.

## 8. Conclusion

The proposed research into generative diffusion models for *de novo* catalyst design represents a high-impact frontier in chemical engineering and materials science. By adapting the successful paradigm from protein engineering (RFdiffusion) and leveraging cutting-edge AI methodologies, we aim to create a powerful computational engine for discovering catalysts capable of precisely modulating hydrocarbon fuel properties from the Fischer-Tropsch synthesis and potentially enabling efficient paraffin cracking. While the computational demands for data generation and model training are significant, ranging from tens of thousands to millions of dollars in hardware and compute over several months, the potential for accelerating discovery and revolutionizing fuel production justifies this investment. The resultant ability to design catalysts that align hydrocarbon outputs with specific density requirements for industrial applications underscores the profound practical utility of this research.