
Emergent SO(3)-Invariant Molecular Representations from Multimodal Alignment

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning molecular representations robust to 3D rotations typically relies on
2 symmetry-aware architectures or extensive augmentation. Here, we show that
3 contrastive multimodal pretraining alone can induce SO(3) invariance in molecular
4 embeddings. We jointly train a 3D electron density encoder, based on a VQGAN,
5 and a SMILES-based transformer encoder on 855k molecules, using CLIP-style and
6 SigLIP objectives to align volumetric and symbolic modalities. Because SMILES
7 embeddings are rotation-invariant, the contrastive loss implicitly enforces rotation-
8 consistency in the 3D encoder. To assess geometric generalization, we introduce a
9 benchmark of 1,000 molecules with five random SO(3) rotations each. Our model
10 retrieves rotated variants with 77% Recall@10 (vs. 9.8% for a unimodal baseline)
11 and organizes latent space by chemical properties, achieving functional group-wise
12 Recall@10 above 98% and a Davies–Bouldin index of 2.35 (vs. 34.46 baseline).
13 Fine-tuning with rotated data reveals a trade-off between retrieval precision and
14 pose diversity. These results demonstrate that contrastive multimodal pretraining
15 can yield symmetry-aware molecular representations without explicit equivariant
16 design.

17 1 Introduction

18 Learning molecular representations that are both chemically expressive and geometrically invariant
19 remains a central challenge in molecular machine learning [1, 2]. Most 3D molecular models achieve
20 invariance to spatial transformations by explicitly encoding symmetry through architectural design or
21 by leveraging rotation-based data augmentation [3, 4, 5]. These methods assume that symmetry priors
22 must be built into the model to preserve physical consistency, particularly under SO(3) rotations. This
23 raises a fundamental question: *Can pose-invariant representations instead emerge implicitly from the*
24 *training objective, without enforcing geometric priors through model design?* [6, 7].

25 We hypothesize that contrastive alignment between invariant symbolic descriptors (e.g., SMILES) and
26 spatially variant 3D fields (e.g., electron densities) can induce pose-consistent molecular embeddings,
27 even in the absence of symmetry-aware architectures [6, 8, 9]. This builds on the intuition that
28 multimodal contrastive learning can serve as a *functional regularizer*, promoting semantic alignment
29 across heterogeneous modalities despite differences in spatial representation [10, 11, 12].

30 Multimodal contrastive learning has shown promise in molecular domains by aligning symbolic and
31 topological views of a molecule [9, 13, 14]. However, existing approaches predominantly operate on
32 graph-based or discrete representations and do not evaluate whether learned embeddings are robust
33 to arbitrary spatial transformations [15]. In particular, it remains unexplored whether contrastive
34 pretraining over unaligned continuous 3D fields can give rise to emergent SO(3) invariance.

35 In this work, we investigate whether a CLIP-style model trained to align SMILES strings with ab
36 initio-derived 3D electron density grids can learn pose-invariant representations, despite lacking
37 architectural equivariance or rotation augmentation. Our model is pretrained on a dataset of 855,000
38 molecules, each presented in a canonical orientation, and jointly embeds both symbolic and volumetric
39 views.

40 To evaluate generalization under spatial transformations, we construct a benchmark of 1,000
41 molecules, each paired with five randomly rotated SO(3) variants. Our contrastive model retrieves
42 at least one rotated instances in the top-10 for 77.3% of queries, approaching the performance of
43 the SE(3)-equivariant Pos-EGNN baseline (79.1%). Pos-EGNN is a large-scale foundation model
44 trained on 1.4M ab initio simulation snapshots from the Materials Project Trajectory dataset to predict
45 energies, forces, and stress tensors using symmetry-aware message passing [16].

46 Beyond retrieval, we probe the latent space for chemical coherence. Without any supervision on
47 quantum properties, the model organizes molecules based on HOMO energies and functional groups:
48 for example, nitrogen-containing species cluster tightly in HOMO-aligned regions. In contrast, the
49 Pos-EGNN latent space—while geometrically grounded—exhibits weaker clustering around frontier
50 orbital descriptors, suggesting that symbolic anchoring plays a critical role in inducing chemically
51 meaningful structure. This organization is quantified by a Davies–Bouldin index of 2.35, compared
52 to 34.46 for a unimodal 3D baseline and 5.53 for the SE(3)-equivariant Pos-EGNN model, indicating
53 superior alignment between geometry and electronic structure.

54 These findings demonstrate that multimodal contrastive pretraining can induce symmetry-aware
55 molecular representations through emergent behavior, without hard-coded inductive biases. While
56 our approach assumes rotational equivalence across poses—an idealization that may not hold in
57 stereochemically sensitive tasks—it offers a flexible and scalable alternative to equivariant model-
58 ing. All code and pretrained models are available at: [https://anonymous.4open.science/r/
59 anonymous-BOBB/README.md](https://anonymous.4open.science/r/anonymous-BOBB/README.md).

60 2 Related Work

61 Learning molecular representations that incorporate 3D structure has been a longstanding objective
62 in machine learning for chemistry. Early approaches relied on graph-based models augmented with
63 spatial features [17, 18], while more recent methods leverage equivariant neural networks [2, 3, 5, 19].
64 These architectures enforce rotational and translational symmetry by design, often using group
65 convolutions or tensor representations. Although effective, these methods hard-code geometric priors
66 into the model, which may limit flexibility across tasks where symmetries are not strictly preserved.

67 Beyond equivariance, several works explore data-driven approaches to learning molecular 3D struc-
68 ture. Models such as GemNet [20] and DimeNet++ [21, 22] use angle and distance information
69 explicitly, while diffusion-based models [23, 24] attempt to generate 3D conformers in a probabilistic
70 manner. These methods assume access to accurate conformations or focus on generating new 3D
71 geometries, rather than studying robustness to transformations applied to known structures.

72 Multimodal learning in molecular domains has focused largely on combining symbolic and graph-
73 based modalities [25, 26, 27]. Works such as MolCLR [9] and Smiclr [28] demonstrate that contrastive
74 pretraining over graphs or SMILES can improve downstream property prediction. AMOLE [29]
75 applies a CLIP-style objective to graphs and text but does not incorporate continuous 3D field-
76 based inputs. As a result, existing multimodal methods primarily operate over discrete structural
77 abstractions, limiting their capacity to exploit fine-grained geometric information available in physical
78 electron density fields.

79 Invariance learning without explicit symmetry enforcement has been explored in vision [30, 31],
80 where models trained without augmentations nonetheless exhibit partial viewpoint robustness. In
81 molecular machine learning, such emergent invariance remains largely unexplored, with most models
82 enforcing rotational symmetry by design [2, 3]. Recent works on SO(3)-equivariant diffusion [23]
83 primarily address generative modeling rather than retrieval robustness under unseen transformations.

84 Our work contributes to this landscape by demonstrating that contrastive multimodal pretraining
85 over symbolic descriptors and continuous 3D grids can induce pose invariance without requiring
86 symmetry-aware architectures. We provide systematic evaluation over rotated benchmarks and relate
87 retrieval stability to chemical and geometric consistency.

88 3 Methodology

89 We propose a multimodal contrastive pretrain-
 90 ing framework to learn molecular representa-
 91 tions that align symbolic descriptors and contin-
 92 uous 3D fields—without relying on symmetry-
 93 aware architectural priors. The model jointly
 94 embeds SMILES strings and electron density
 95 grids derived from ab initio calculations using
 96 independent encoders optimized under a con-
 97 trastive loss. As illustrated in Figure 1, our
 98 architecture combines a transformer-based en-
 99 coder for SMILES with a 3D VQGAN-style
 100 convolutional encoder for electron densities. All
 101 parameters—including those from SMI-TED and the 3DGrid-VQGAN encoder—are trained jointly
 102 from scratch.

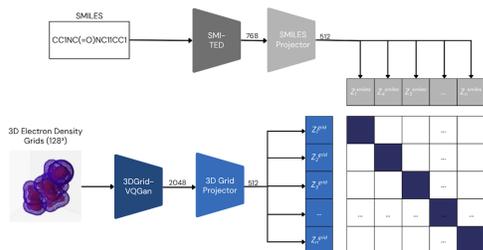


Figure 1: Architecture of the multimodal contrastive model.

103 3.1 Pretraining Dataset

104 We curate a dataset of 855,000 molecules from PubChem, filtered to include: (i) only main-group
 105 elements up to Barium; (ii) a maximum of 30 heavy atoms; (iii) zero net charge; and (iv) no formal
 106 charge separation.

107 Each SMILES string is converted into 50 conformers using RDKit’s distance geometry and MMFF94
 108 optimization [32]. The five lowest-energy conformers are reoptimized using MINDO3 in PySCF [33],
 109 and the conformer with the lowest energy is retained. This structure is further evaluated at the
 110 RHF/STO-3G level, and its electron density is projected onto a $128 \times 128 \times 128$ voxel grid, yielding
 111 a physically grounded 3D representation without relying on classical graph approximations.

112 3.2 Multimodal Contrastive Pretraining

113 We align SMILES and 3D electron density representations via contrastive learning. Let $g : \mathcal{X} \rightarrow \mathbb{R}^d$
 114 and $h : \mathcal{T} \rightarrow \mathbb{R}^d$ denote the 3D and SMILES encoders, respectively. For a batch of N molecule pairs
 115 $\{(x_i, t_i)\}_{i=1}^N$, we compute embeddings as $\mathbf{z}_i^{\text{grid}} = \text{Proj}_g(g(x_i))$ and $\mathbf{z}_i^{\text{smiles}} = \text{Proj}_h(h(t_i))$, where
 116 Proj denotes a learnable projection head.

117 **SMILESDFt-CLIP** uses the symmetric InfoNCE loss:

$$\mathcal{L}_{\text{CLIP}} = \frac{1}{2N} \sum_i [\ell(\mathbf{z}_i^{\text{grid}}, \mathbf{z}_i^{\text{smiles}}) + \ell(\mathbf{z}_i^{\text{smiles}}, \mathbf{z}_i^{\text{grid}})],$$

118 where $\ell(z, z') = -\log \frac{\exp(\text{sim}(z, z')/\tau)}{\sum_j \exp(\text{sim}(z, z'_j)/\tau)}$ and $\text{sim}(z, z')$ is cosine similarity.

119 **SMILESDFt-SigLIP** employs a sigmoid-based contrastive loss. After normalization $\tilde{\mathbf{z}} = \mathbf{z}/\|\mathbf{z}\|_2$,
 120 we define:

$$\text{logits}_{ij} = \exp(\tau) \cdot \langle \tilde{\mathbf{z}}_i^{\text{grid}}, \tilde{\mathbf{z}}_j^{\text{smiles}} \rangle + b, \quad \mathcal{L}_{\text{SigLIP}} = -\frac{1}{N} \sum_{i,j} \log \sigma(y_{ij} \cdot \text{logits}_{ij}),$$

121 where σ is the sigmoid function and $y_{ij} = 1$ for positive pairs, -1 otherwise.

122 3.3 3D Electron Density Encoder

123 We use a 3DGrid-VQGAN adapted for volumetric inputs to encode electron density grids [16]. The
 124 encoder $E(\cdot)$ maps G to a latent tensor:

$$z_e(G) \in \mathbb{R}^{\frac{H}{s} \times \frac{W}{s} \times \frac{D}{s} \times k}, \quad s = 4, \quad k = 512.$$

125 Latents are quantized using a learned codebook $\{e_j\}_{j=1}^{16384}$:

$$z_q(G) = e_{k^*}, \quad \text{where } k^* = \arg \min_j \|z_e(G) - e_j\|_2.$$

126 The 3DGrid-VQGAN is trained with:

$$\mathcal{L}_{\text{VQGAN}} = \mathcal{L}_{\text{rec}} + \beta \mathcal{L}_{\text{commit}} + \gamma \mathcal{L}_{\text{adv}},$$

127 where \mathcal{L}_{rec} is an L_1 reconstruction loss, $\mathcal{L}_{\text{commit}}$ encourages codebook usage, and \mathcal{L}_{adv} is a 3D
128 PatchGAN adversarial loss. During contrastive training, we use the encoder output *before quantization*
129 and fine-tune all encoder parameters jointly.

130 3.4 SMILES Encoder

131 The SMILES modality is encoded using SMI-TED_{289M} [16], a pretrained transformer encoder trained
132 on 91 million canonical SMILES strings. Input tokens $X \in \mathbb{R}^{D \times L}$ are processed via RoFormer-style
133 attention:

$$\text{Attention}_m(Q, K, V) = \frac{\sum_{n=1}^N \langle \varphi(R_m q_m), \varphi(R_n k_n) \rangle v_n}{\sum_{n=1}^N \langle \varphi(R_m q_m), \varphi(R_n k_n) \rangle},$$

134 where R_m is a position-specific rotation matrix and $\varphi(\cdot)$ is a Fourier feature mapping. A pooled
135 embedding is computed as:

$$\mathbf{z} = \text{LayerNorm}(\text{GELU}(X \mathbf{W}_1 + \mathbf{b}_1)) \mathbf{W}_2.$$

136 Unlike prior work, we fine-tune the SMI-TED encoder during contrastive learning, which we find
137 improves performance in both retrieval and structure–property clustering.

138 3.5 Training Details

139 We train using AdamW with batch size 128 and learning rate 3×10^{-4} , employing a linear warmup
140 over 1,000 steps. Models are trained for 50,000 steps using both CLIP and SigLIP objectives, with
141 checkpoints selected by retrieval accuracy on a held-out validation set. All experiments are conducted
142 on 4 NVIDIA A100 GPUs.

143 4 Experimental Setup

144 We conduct a comprehensive evaluation to assess the extent to which our multimodal model exhibits
145 geometric generalization, chemical organization, and transferability. Our evaluation protocol includes
146 retrieval under both canonical and unseen SO(3) rotations, unsupervised structure–property clustering,
147 and molecular property prediction on the QM9 benchmark.

148 **Retrieval under SO(3) Rotations.** We evaluate retrieval performance in two settings:

- 149 1. **Canonical retrieval** – Each query is matched against a corpus of unrotated (canonical)
150 molecules.
- 151 2. **Unseen rotation retrieval** – Each query is matched against five rotated spatial variants of
152 each molecule, not observed during training.

153 To generate unseen rotations, we apply random rigid-body transformations to the
154 atomic coordinates of each molecule. Rotation axes are sampled from the set
155 $\{(0, 0, 1), (0, 1, 0), (1, 0, 0), (1, 1, 0), (1, 0, 1), (0, 1, 1)\}$, and rotation angles are drawn uni-
156 formly from $[0^\circ, 360^\circ]$. For each rotated conformer, we recompute the electron density using the
157 same RHF/STO-3G procedure as used during pretraining, ensuring physically valid volumetric fields.

158 We report the following metrics as in Table 1:

159 To benchmark invariance, we compare our model against a unimodal 3D electron density baseline
160 (3DGrid-VQGAN) trained without symbolic alignment. This evaluation probes both instance-level
161 and class-level generalization under unseen spatial transformations.

162 **Structure–Property Clustering.** We assess whether the latent space reflects chemically meaningful
163 organization by analyzing clustering behavior of molecules with similar frontier orbital properties.
164 In particular, we focus on nitrogen-containing species with high HOMO energies—chemically
165 important due to lone-pair reactivity. We quantify cluster quality using the **Davies–Bouldin (DB)**
166 **index**, where lower values indicate compact, well-separated clusters. This analysis tests whether the
167 model implicitly learns structure–property relationships without supervision.

Table 1: Evaluation metrics used to assess geometric and chemical generalization.

Metric	Description
Accuracy@10	Proportion of queries retrieving the correct molecule within the top-10 results.
Recall@10	Fraction of retrieved molecules belonging to the same functional group.
Group-wise Recall@10	Recall@10 computed for six chemical classes: amines, aromatics, ethers, ketones, halides, and carboxylic acids.
Pose diversity	Mean number of distinct rotational variants retrieved in the top-10.
Multi-pose retrieval rate	Proportion of queries retrieving at least three distinct rotated variants among the top-10.

168 **Property Prediction on QM9.** To evaluate transferability to downstream tasks, we train linear
 169 regression models on frozen multimodal embeddings to predict 12 molecular properties from the
 170 QM9 dataset [34]. The encoders are not fine-tuned, ensuring that performance reflects the intrinsic
 171 quality of the pretrained representation. We report **mean absolute error (MAE)** on the standard
 172 train/validation/test splits and compare against an equivariant baseline embeddings from Pos-EGNN
 173 encoder.

174 5 Results

175 We evaluate the capacity of our multimodal model to achieve pose-invariant molecular retrieval and
 176 chemically consistent embeddings without architectural equivariance. The evaluation is organized
 177 along two main axes: retrieval under unseen $SO(3)$ rotations and retrieval consistency across known
 178 rotations observed during training. Additional analyses include functional group-specific recall and
 179 structural similarity assessments.

180 5.1 Retrieval under $SO(3)$ Rotations

181 We evaluate the ability of our multimodal model to achieve chemically and geometrically consistent
 182 retrieval without architectural symmetry constraints. Our experiments are organized into two main
 183 settings: retrieval among canonical poses (supplementary materials) and retrieval under unseen $SO(3)$
 184 rotations.

185 **Retrieval under $SO(3)$ Rotations.** In this experiment, we evaluate the ability of pretrained models
 186 to retrieve molecular representations under unseen rigid-body transformations. To simulate $SO(3)$
 187 rotation invariance, molecules are randomly rotated around arbitrary axes, with rotation angles
 188 uniformly sampled from $[0^\circ, 360^\circ]$. Retrieval is performed by querying canonical molecules against
 189 rotated versions in embedding space, testing whether models generalize across poses without having
 190 observed such transformations during training.

191 Table 2 summarizes the performance across four models using three metrics: (i) Accuracy@10, which
 192 captures exact retrieval of a rotated instance; (ii) Recall@10, which measures class-level or functional
 193 group recovery; and (iii) the proportion of queries for which three or more rotated variants appear
 194 among the top-10 candidates.

Table 2: Retrieval performance under unseen $SO(3)$ rotations. **Equivariant** indicates $SE(3)$ -equivariant models. Accuracy@10 measures instance-level retrieval; Recall@10 captures functional group recovery; final column reports the percentage of queries retrieving 3 distinct rotated variants in the top-10.

Model	Equiv.	Modality	Acc@10	Rec@10	3 Rot. Retrieved
Ours					
SMILESDFE-CLIP	✗	3D Grids + SMILES	77.3% ± 0.51	98.4% ± 0.13	45.3% ± 0.57
SMILESDFE-SigLIP	✗	3D Grids + SMILES	46.1% ± 0.57	98.9% ± 0.13	43.0% ± 0.63
SMILESDFE-CLIP (finetuned)	✗	3D Grids + SMILES	85.4% ± 0.42	99.4% ± 0.09	57.9% ± 0.54
SMILESDFE-SigLIP (finetuned)	✗	3D Grids + SMILES	88.4% ± 0.37	99.6% ± 0.08	59.1% ± 0.52
Baselines					
Pos-EGNN	✓	Atom Positions ($SE(3)$)	79.1% ± 0.44	99.2% ± 0.12	51.2% ± 0.51
3DGrid-VQGAN	✗	3D Grids Only	9.1% ± 0.22	2.3% ± 0.02	0.0% ± 0.01

195 Table 2 reports retrieval performance under unseen SO(3) rotations, comparing our multimodal
 196 models to both equivariant and non-equivariant baselines. Fine-tuned variants of SMILESDF-CLIP
 197 and SMILESDF-SigLIP—trained on 1,000 randomly selected molecules with five randomly rotated
 198 poses each—achieve the highest Accuracy@10 (85.4% and 88.4%, respectively), outperforming
 199 the SE(3)-equivariant Pos-EGNN baseline (79.1%) despite lacking explicit symmetry priors. All
 200 multimodal models exhibit strong functional group recovery (Recall@10), with fine-tuned versions
 201 reaching 99.6% (SMILESDF-SigLIP). Furthermore, over 57% of fine-tuned model queries retrieve
 202 at least three distinct rotated variants in the top-10, exceeding the equivariant baseline (51.2%) and
 203 substantially outperforming the unimodal 3DGrid-VQGAN model, which fails under rotation. These
 204 results suggest that contrastive multimodal pretraining, when exposed to a modest set of diverse
 205 poses, can induce rotation-consistent representations without requiring architectural equivariance.

206 This result underscores the central contribution of our approach: emergent rotational invariance
 207 arises from multimodal contrastive pretraining, even in the absence of architectural equivariance or
 208 rotation augmentation. The SMILES representation remains invariant under rotation and acts as a
 209 semantic anchor. Minimizing the contrastive loss aligns the spatially-variant 3D electron density
 210 fields with these invariant anchors, inducing consistent embeddings across different orientations.
 211 Pooling operations further reduce sensitivity to local spatial deformations, contributing to pose-robust
 212 representations.

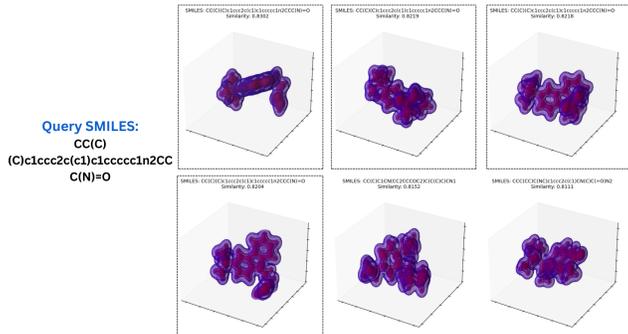


Figure 2: Visualization of retrieval results under unseen rotations using SMILESDF-CLIP. Query SMILES: CC(C)(C)c1ccc2c(c1)c1cccc1n2CCC(N)=O. Retrieved electron density grids are matched with corresponding SMILES and cosine similarity scores. The model retrieves four perfect matches and one close structural analogue, illustrating robustness to SO(3) transformations.

213 Figure 2 illustrates a retrieval example using SMILESDF-CLIP. Among the six closest retrieved
 214 samples, four are exact matches under distinct rotations, and one is a structurally similar analogue.
 215 This highlights the model’s ability to capture both spatial and semantic consistency.

Group-wise Recall@10 scores (Table 3) reveal high retrieval robustness across functional groups. Both multimodal models achieve near-perfect recovery for aromatic and ketone-containing compounds. Slightly lower recall for carboxylic acids may stem from their conformational flexibility and smaller spatial extent in grid representation, which challenges invariant matching.

Table 3: Recall@10 across functional groups under unseen SO(3) rotations.

Functional Group	SMILESDF-CLIP	SMILESDF-SigLIP
Amine	0.987	0.994
Aromatic	0.999	1.000
Ether	0.987	0.981
Ketone	0.987	1.000
Halide	0.961	0.978
Carboxylic Acid	0.893	0.890

216 In summary, our results show that contrastive multimodal pretraining can induce pose-invariant
 217 molecular representations without relying on symmetry-aware inductive biases. By leveraging the
 218 invariant nature of symbolic descriptors during alignment, the model internalizes spatial consistency
 219 across orientations. This emergent behavior bridges the gap between architectural equivariance
 220 and semantic invariance, opening new directions for building chemically robust models from weak
 221 supervision alone.

222 5.2 Structure–Property Relationship

223 To evaluate whether the learned latent representations reflect chemically meaningful struc-
224 ture–property relationships, we analyze clustering behavior based on the HOMO (Highest Occupied
225 Molecular Orbital) energy, a key descriptor of molecular reactivity. Nitrogen-containing species are
226 of particular interest due to the strong influence of nitrogen lone pairs, which elevate HOMO energy
227 and enhance molecular reactivity.

228 In the QM9 dataset, nitrogen-containing molecules comprise only 9.10% of the total population but
229 represent 32.81% of the top decile in HOMO energy. Capturing such functional and electronic trends
230 in the learned embedding space—without direct supervision on quantum properties—is a critical test
231 of the model’s chemical fidelity.

232 To quantify clustering quality, we compute the Davies–Bouldin (DB) index, which penalizes overlap-
233 ping or diffuse clusters (lower is better). Table 4 summarizes the DB scores across models. Notably,
234 the SMILESDFDFT-CLIP-based multimodal model achieves the lowest DB index (2.35), indicating a
235 tightly organized latent space with well-separated clusters that align with HOMO energy variations.
236 In contrast, the position-equivariant Pos-EGNN model, despite its architectural symmetry priors,
237 yields a higher DB index (5.53), suggesting weaker alignment with electronic structure. This is
238 surprising, as equivariant models are expected to encode physically grounded representations, but
239 lack symbolic anchoring to enforce chemical alignment.

Table 4: Davies–Bouldin (DB) index for structure–property clustering by HOMO energy (lower is better). SMILESDFDFT-CLIP achieves the lowest Davies–Bouldin index, indicating tight HOMO-aligned clustering. Symbolic input (SMILES) plays a critical role in structuring latent space, even in the absence of equivariant design.

Model	SMILES	3D Grids / Atom Positions	DB Index
SMILESDFDFT-CLIP	✓	3D Grids	2.35
SMI-TED	✓	✗	2.82
MoLFormer	✓	✗	4.28
Pos-EGNN	✗	Atom Positions (SE(3))	5.53
3DGrid-VQGAN	✗	3D Grids	34.46

240 Figures 3 visualize 2D projections of the learned latent spaces, with colors representing HOMO
241 energy and triangle markers highlighting nitrogen-containing species. The SMILESDFDFT-CLIP
242 latent space reveals compact clusters strongly correlated with HOMO energy and clearly segregated
243 nitrogen-rich regions, supporting the hypothesis that contrastive multimodal pretraining promotes
244 chemically meaningful representation learning.

245 In contrast, Pos-EGNN—despite encoding atom positions in an equivariant manner—produces a more
246 diffuse and intermixed embedding space, with nitrogen-containing species scattered across regions of
247 varying energy. This suggests that architectural symmetry alone does not guarantee property-aligned
248 representations unless supported by complementary semantic signals.

249 These findings emphasize that contrastive multimodal learning acts not merely as a cross-modal
250 alignment strategy, but as a *functional regularizer* that filters and reinforces task-relevant structural
251 patterns. The invariant SMILES anchor encourages the 3D encoder to focus on chemical features
252 consistent across orientations, facilitating the emergence of rotationally robust, semantically grounded
253 embeddings.

254 5.3 Property Prediction on QM9

255 We assess the downstream utility of our pretrained representations on the QM9 benchmark, which
256 comprises 12 regression tasks spanning electronic, thermodynamic, and geometric properties. Mean
257 absolute error (MAE) is reported in QM9-standard units.

258 Using a pre-trained linear probe setup, we evaluate SMILESDFDFT-CLIP and SMILESDFDFT-SigLIP with-
259 out task-specific fine-tuning to isolate representation quality. As shown in Table 5, both SMILESDFDFT-
260 CLIP and SMILESDFDFT-SigLIP consistently outperform the SE(3)-equivariant Pos-EGNN baseline
261 across most tasks. SMILESDFDFT-CLIP achieves the lowest MAE on 8 of 12 properties—including
262 ϵ_{HOMO} , C_v , and $\langle R^2 \rangle$ —while SigLIP is competitive, especially on thermodynamic targets (U , U_0 , H ,
263 G)

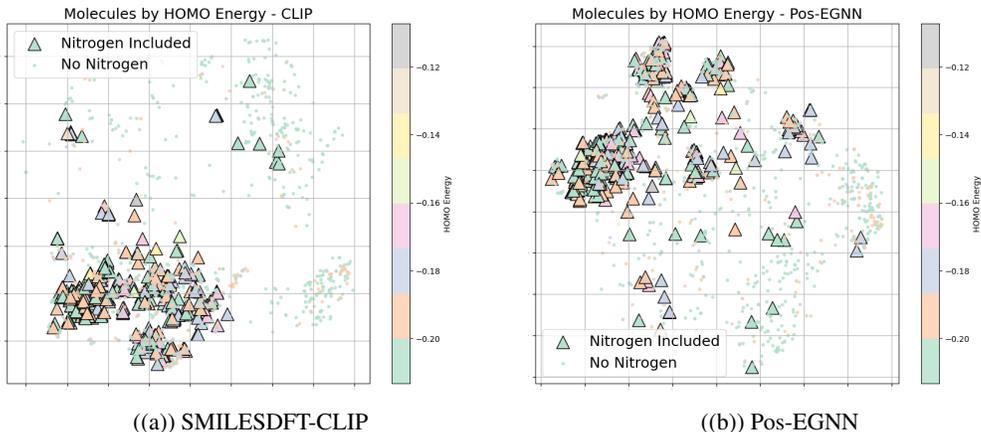


Figure 3: Latent space projections colored by HOMO energy. Triangular markers denote nitrogen-containing molecules.

Table 5: Mean Absolute Error (MAE) on QM9 regression tasks. All models are evaluated in a frozen linear probe setting. **Blue** and **Orange** highlight the best and second-best results, respectively.

Category	Property (Unit)	Pos-EGNN (Equivariant)	SMILESDFE-CLIP (Non-equivariant)	SMILESDFE-SigLIP (Non-equivariant)
Electronic	HOMO energy ϵ_{HOMO} (eV)	0.0093	0.0083	0.0090
	LUMO energy ϵ_{LUMO} (eV)	0.0141	0.0110	0.0118
	Energy gap (eV)	0.0165	0.0135	0.0144
	Dipole moment μ (Debye)	0.6288	0.6836	0.7243
Thermodynamic	Internal energy U (eV)	6.8596	2.8141	2.3437
	Internal energy at 0K U_0 (eV)	6.8308	2.8403	2.3316
	Enthalpy H (eV)	6.8503	2.8137	2.3402
	Free energy G (eV)	6.8335	2.8098	2.3706
Geometric	Heat capacity C_v (cal/mol-K)	0.6622	0.4450	0.4455
	Polarizability α (bohr ³)	1.5346	1.0771	1.1791
	Spatial extent $\langle R^2 \rangle$ (bohr ²)	70.5140	45.2012	46.5561
	ZPVE (eV)	0.0064	0.0038	0.0040

264 Despite lacking symmetry-aware priors, both models outperform Pos-EGNN on geometry-sensitive
 265 metrics such as polarizability (α) and spatial extent ($\langle R^2 \rangle$), suggesting that contrastive symbolic
 266 alignment can induce symmetry-consistent behavior through emergent structure in the latent space.

267 6 Limitations

268 Our approach assumes that rigid SO(3) rotations preserve molecular semantics, which may not
 269 generalize to stereochemically sensitive or highly flexible molecules. The use of RHF/STO-3G
 270 electron densities, while providing quantum-consistent inputs, adds computational cost and may
 271 limit scalability. Although both encoders are fine-tuned jointly, the model is not explicitly trained to
 272 align different conformers of the same molecule. Future extensions could incorporate conformer-
 273 invariant objectives, lightweight electronic representations, or hybrid training with physicochemical
 274 supervision.

275 7 Conclusion

276 We show that contrastive multimodal pretraining between SMILES and 3D electron densities yields
 277 chemically meaningful, pose-invariant representations—without symmetry-aware architectures or
 278 rotation augmentation. The model generalizes across SO(3) rotations and reflects orbital energy
 279 structure, despite no geometric or quantum supervision. Symbolic anchoring emerges as a simple but
 280 effective inductive signal. Future work may explore conformer-aware or property-specific extensions.

281 References

- 282 [1] K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller, “SchNet—a deep learning
283 architecture for molecules and materials,” *The Journal of Chemical Physics*, vol. 148, no. 24, 2018.
- 284 [2] V. G. Satorras, E. Hoogeboom, and M. Welling, “E (n) equivariant graph neural networks,” in *International
285 conference on machine learning*. PMLR, 2021, pp. 9323–9332.
- 286 [3] N. Thomas, T. Smidt, S. Kearnes, L. Yang, L. Li, K. Kohlhoff, and P. Riley, “Tensor field networks:
287 Rotation-and translation-equivariant neural networks for 3d point clouds,” *arXiv preprint arXiv:1802.08219*,
288 2018.
- 289 [4] F. Fuchs, D. Worrall, V. Fischer, and M. Welling, “Se (3)-transformers: 3d roto-translation equivariant
290 attention networks,” *Advances in neural information processing systems*, vol. 33, pp. 1970–1981, 2020.
- 291 [5] B. Anderson, T. S. Hy, and R. Kondor, “Cormorant: Covariant molecular neural networks,” *Advances in
292 neural information processing systems*, vol. 32, 2019.
- 293 [6] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual
294 representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- 295 [7] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” *IEEE
296 transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- 297 [8] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
298 J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International
299 conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- 300 [9] Y. Wang, J. Wang, Z. Cao, and A. Barati Farimani, “Molecular contrastive learning of representations via
301 graph neural networks,” *Nature Machine Intelligence*, vol. 4, no. 3, pp. 279–287, 2022.
- 302 [10] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, Y. Wei, Q. Dai, and H. Hu, “On data scaling in masked image modeling,”
303 in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp.
304 10 365–10 374.
- 305 [11] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, “Multimodal
306 transformer for unaligned multimodal language sequences,” in *Proceedings of the conference. Association
307 for computational linguistics. Meeting*, vol. 2019, 2019, p. 6558.
- 308 [12] J. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka, “Contrastive learning with hard negative samples,” *arXiv
309 preprint arXiv:2010.04592*, 2020.
- 310 [13] X. Zhang, Y. Xu, C. Jiang, L. Shen, and X. Liu, “Moleml: a multi-level contrastive learning framework
311 for molecular pre-training,” *Bioinformatics*, vol. 40, no. 4, p. btac164, 2024.
- 312 [14] B. Kaufman, E. C. Williams, C. Underkoffler, R. Pederson, N. Mardirossian, I. Watson, and J. Parkhill,
313 “Coati: Multimodal contrastive pretraining for representing and traversing chemical space,” *Journal of
314 Chemical Information and Modeling*, vol. 64, no. 4, pp. 1145–1157, 2024.
- 315 [15] S. Takeda, A. Kishimoto, L. Hamada, D. Nakano, and J. R. Smith, “Foundation model for material science,”
316 in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 13, 2023, pp. 15 376–15 383.
- 317 [16] “Anonymized.”
- 318 [17] K. Schütt, P.-J. Kindermans, H. E. Sauceda Felix, S. Chmiela, A. Tkatchenko, and K.-R. Müller, “SchNet:
319 A continuous-filter convolutional neural network for modeling quantum interactions,” *Advances in neural
320 information processing systems*, vol. 30, 2017.
- 321 [18] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, “Neural message passing for quantum
322 chemistry,” in *International conference on machine learning*. PMLR, 2017, pp. 1263–1272.
- 323 [19] Y. Luo, Z. Liu, Y. Zhao, S. Li, K. Kawaguchi, T.-S. Chua, and X. Wang, “Towards unified latent space for
324 3d molecular latent diffusion modeling,” *arXiv preprint arXiv:2503.15567*, 2025.
- 325 [20] J. Klicpera, F. Becker, and S. Günnemann, “Gemnet: Universal directional graph neural networks for
326 molecules,” in *Proceedings of the 35th International Conference on Neural Information Processing Systems*,
327 2021, pp. 6790–6802.
- 328 [21] J. Gasteiger, S. Giri, J. T. Margraf, and S. Günnemann, “Fast and uncertainty-aware directional message
329 passing for non-equilibrium molecules,” *arXiv preprint arXiv:2011.14115*, 2020.

- 330 [22] F. Zhu, M. Futrega, H. Bao, S. B. Eryilmaz, F. Kong, K. Duan, X. Zheng, N. Angel, M. Jouanneaux,
331 M. Stadler *et al.*, "Fastdimenet++: Training dimenet++ in 22 minutes," in *Proceedings of the 52nd*
332 *International Conference on Parallel Processing*, 2023, pp. 274–284.
- 333 [23] E. Hoogeboom, V. G. Satorras, C. Vignac, and M. Welling, "Equivariant diffusion for molecule generation
334 in 3d," in *International conference on machine learning*. PMLR, 2022, pp. 8867–8887.
- 335 [24] M. Xu, L. Yu, Y. Song, C. Shi, S. Ermon, and J. Tang, "Geodiff: A geometric diffusion model for molecular
336 conformation generation," *arXiv preprint arXiv:2203.02923*, 2022.
- 337 [25] Y. Xiao, X. Zhou, Q. Liu, and L. Wang, "Bridging text and molecule: A survey on multimodal frameworks
338 for molecule," *arXiv preprint arXiv:2403.13830*, 2024.
- 339 [26] J. Li, D. Zhang, X. Wang, Z. Hao, J. Lei, Q. Tan, C. Zhou, W. Liu, Y. Yang, X. Xiong *et al.*, "Chemvlm:
340 Exploring the power of multimodal large language models in chemistry area," in *Proceedings of the AAAI*
341 *Conference on Artificial Intelligence*, vol. 39, no. 1, 2025, pp. 415–423.
- 342 [27] P. Liu, Y. Ren, J. Tao, and Z. Ren, "Git-mol: A multi-modal large language model for molecular science
343 with graph, image, and text," *Computers in biology and medicine*, vol. 171, p. 108073, 2024.
- 344 [28] G. A. Pinheiro, J. L. Da Silva, and M. G. Quiles, "Smicl: Contrastive learning on multiple molecular repre-
345 sentations for semisupervised and unsupervised representation learning," *Journal of Chemical Information*
346 *and Modeling*, vol. 62, no. 17, pp. 3948–3960, 2022.
- 347 [29] N. Lee, S. Laghuvarapu, C. Park, and J. Sun, "Vision language model is not all you need: Augmentation
348 strategies for molecule language models," in *Proceedings of the 33rd ACM International Conference on*
349 *Information and Knowledge Management*, 2024, pp. 1153–1162.
- 350 [30] T. Chen, S. Saxena, L. Li, D. J. Fleet, and G. Hinton, "Pix2seq: A language modeling framework for object
351 detection," *arXiv preprint arXiv:2109.10852*, 2021.
- 352 [31] G. Kim, T. Kwon, and J. C. Ye, "Diffusionclip: Text-guided diffusion models for robust image manipula-
353 tion," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp.
354 2426–2435.
- 355 [32] G. Landrum, "Rdkit documentation," *Release*, vol. 1, no. 1-79, p. 4, 2013.
- 356 [33] Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. McClain, E. R. Sayfutyarova,
357 S. Sharma *et al.*, "Pyscf: the python-based simulations of chemistry framework," *Wiley Interdisciplinary*
358 *Reviews: Computational Molecular Science*, vol. 8, no. 1, p. e1340, 2018.
- 359 [34] Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande,
360 "Moleculenet: a benchmark for molecular machine learning," *Chemical science*, vol. 9, no. 2, pp. 513–530,
361 2018.

362 A Supplementary Materials

363 A.1 Retrieval among Canonical Poses

Model	Top-1	Top-10	FG Matches (Top-10)
SMILESDF-SigLIP	68.9% ± 1.96	97.6% ± 0.21	8.33
SMILESDF-CLIP	71.4% ± 0.83	98.8% ± 0.14	8.43

Table 6: Retrieval performance among canonical poses.

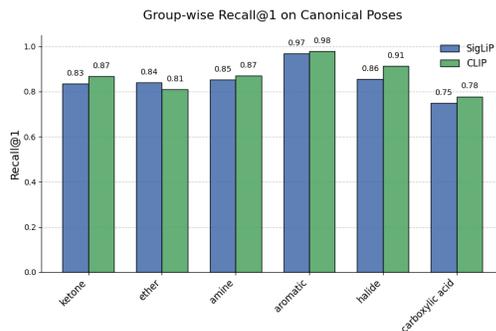


Figure 4: Group-wise Recall@1 for canonical retrieval.

364 In the first setting, retrieval is performed among unrotated (canonical) molecular conformations. Each molecule
365 is embedded in a fixed pose, and retrieval relies solely on feature similarity without any unseen spatial transfor-
366 mations. This setting tests whether the learned embeddings capture molecular identity and functional similarity
367 under ideal alignment.

368 Retrieval metrics include Top- k Match Accuracy (the fraction of queries retrieving the exact molecule) and
369 functional group (FG) recall, measuring how many retrieved molecules share dominant chemical groups with
370 the query.

371 Both SMILESDFT-CLIP and SMILESDFT-SigLIP achieve high retrieval performance. At Top-10, SMILESDFT-
372 CLIP achieves 98.8% accuracy, while SMILESDFT-SigLIP reaches 97.6%. The average number of functional
373 group matches within the Top-10 retrieved molecules exceeds eight for both models, demonstrating chemically
374 aligned latent organization.

375 As shown in Figure 4, aromatic systems exhibit the highest recall (97.9% SMILESDFT-CLIP, 96.8%
376 SMILESDFT-SigLIP), consistent with their distinct electronic signatures. Carboxylic acids, by contrast, show
377 lower recall (77.8% SMILESDFT-CLIP, 75.0% SMILESDFT-SigLIP), likely due to their conformational flexi-
378 bility. Across all groups, SMILESDFT-CLIP consistently outperforms SMILESDFT-SigLIP at Top-1, indicating
379 sharper discrimination from InfoNCE-based alignment.