

FrameDiff v2: Unconditional Protein Backbone Generation via IGSO3 Score Matching on SE(3) Frames

Sergio E. Mares

March 2026

Abstract

We present FrameDiff v2, an unconditional generative model for protein backbone structures that operates directly on SE(3) rigid-body frames. The model employs Invariant Point Attention (IPA) within an 8-layer iterative refinement architecture (14.4M parameters) to predict clean frames from noised inputs. Rotational noise follows the Isotropic Gaussian distribution on SO(3) (IGSO3), and the rotation loss uses proper IGSO3 score matching in the tangent space rather than the Frobenius geodesic norm used in v1. Translations follow a variance-preserving SDE with centered Gaussian noise. Auxiliary losses on backbone atom positions and pairwise C α distances regularize low-noise predictions. Training on CATH 4.3 (10,600 chains, max 256 residues) with cosine learning-rate decay reaches C α RMSD of 4.45 Å and mean rotation error of 45.5° at epoch 122. We document five numerical stability fixes required for NaN-free training and analyze a training instability at epoch 81 that motivated the switch from constant to cosine LR scheduling.

1 Introduction

Protein backbone generation is a key task in computational structural biology, with applications in *de novo* protein design, scaffold generation for functional-site grafting, and data augmentation for structure prediction. The backbone of a protein can be described as a sequence of rigid-body frames in SE(3): each residue i carries a rotation matrix $R_i \in \text{SO}(3)$ (orienting the peptide plane) and a translation $\mathbf{x}_i \in \mathbb{R}^3$ (the C α position).

Diffusion models have emerged as powerful generative frameworks for continuous data. However, SO(3) is a compact, non-Euclidean manifold, so standard Gaussian diffusion does not directly apply to rotations. FrameDiff v1 used a Frobenius geodesic loss ($\|\log(R_{\text{pred}}^T R_{\text{true}})\|^2$) which ignores the proper score-matching structure of the diffusion process on SO(3).

FrameDiff v2 replaces this with the Isotropic Gaussian on SO(3) (IGSO3) distribution, which is the natural analogue of a Gaussian on \mathbb{R}^n but defined on the rotation group. The IGSO3 density admits a closed-form series expansion, enabling exact score computation and proper score-matching loss formulation. This yields a theoretically principled training objective and, in conjunction with architectural improvements (deeper MLPs, IPA transition layers, configurable stop-gradient), leads to improved training dynamics.

2 Model Architecture

The FrameDiff v2 model is an iterative refinement network that takes noised SE(3) frames and a diffusion timestep as input and predicts clean frames. It has no sequence input and no encoder—it is purely unconditional. The architecture consists of $L = 8$ stacked refinement layers, each containing five components described below.

2.1 Invariant Point Attention (IPA)

IPA is the core attention mechanism, adapted from AlphaFold2. It computes attention in both scalar space and 3-D point space, with point queries/keys/values transformed through per-residue SE(3) frames to ensure SE(3) equivariance.

Given node features $\mathbf{h}_i \in \mathbb{R}^{d_{\text{node}}}$, edge features $\mathbf{z}_{ij} \in \mathbb{R}^{d_{\text{edge}}}$, and frames (R_i, \mathbf{x}_i) , the attention logits for head a are:

$$\alpha_{ij}^{(a)} = \underbrace{\frac{\mathbf{q}_i^{(a)\top} \mathbf{k}_j^{(a)}}{\sqrt{d_{\text{head}}}}}_{\text{scalar}} - \underbrace{\gamma^{(a)} \sum_{p=1}^{P_q} \|\mathbf{q}_{i,p}^{\text{pt}} - \mathbf{k}_{j,p}^{\text{pt}}\|^2}_{\text{point}} + \underbrace{\mathbf{z}_{ij}^\top \mathbf{w}_{\text{pair}}^{(a)}}_{\text{pair bias}} \quad (1)$$

where $\gamma^{(a)} = \text{softplus}(\hat{\gamma}^{(a)})$ is a learnable per-head weight, and the point queries/keys are transformed to the global frame via $\mathbf{q}_{i,p}^{\text{pt}} = R_i \mathbf{q}_{i,p}^{\text{local}} + \mathbf{x}_i$.

Attention weights are $w_{ij}^{(a)} = \text{softmax}_j(\alpha_{ij}^{(a)})$. The output aggregates scalar values, point value norms (transformed back to the local frame), and pair-weighted values:

$$\mathbf{o}_i = W_{\text{out}} \left[\bigoplus_{a=1}^H \left(\sum_j w_{ij}^{(a)} \mathbf{v}_j^{(a)}, \left\| \sum_j w_{ij}^{(a)} \mathbf{v}_{j,p}^{\text{pt,local}} \right\|, \sum_j w_{ij}^{(a)} \mathbf{v}_{ij}^{\text{pair}} \right) \right] \quad (2)$$

The IPA uses $H = 12$ heads, $d_{\text{head}} = 16$, $P_q = 8$ query points, and $P_v = 12$ value points.

2.2 IPA Transition MLP

After IPA and LayerNorm, a 2-layer MLP with ReLU activation is applied with a residual connection (AF2-style):

$$\mathbf{h}_i \leftarrow \mathbf{h}_i + \text{MLP}_{\text{trans}}(\mathbf{h}_i), \quad \text{MLP}_{\text{trans}} : \mathbb{R}^{d_{\text{node}}} \rightarrow \mathbb{R}^{d_{\text{node}}} \rightarrow \mathbb{R}^{d_{\text{node}}} \quad (3)$$

2.3 Transformer Block with Skip Connections

The node features are concatenated with the initial features $\mathbf{h}_i^{(0)}$ and processed through a Transformer block:

$$\mathbf{h}_i \leftarrow \text{TransformerBlock}([\mathbf{h}_i \parallel \mathbf{h}_i^{(0)}]) \quad (4)$$

The block projects $2d_{\text{node}} \rightarrow d_{\text{node}}$, applies multi-head self-attention ($n_{\text{heads}} = 4$) with pre-LayerNorm and a 2-layer FFN ($d_{\text{ffn}} = 1024$) with GELU activation, each with residual connections.

2.4 Edge Update

Edge features are updated via a 2-layer MLP (v2 improvement over v1’s single linear layer):

$$\mathbf{z}_{ij} \leftarrow \text{LayerNorm}(\mathbf{z}_{ij} + \text{MLP}_{\text{edge}}([\mathbf{h}_i \parallel \mathbf{h}_j \parallel \mathbf{z}_{ij}])) \quad (5)$$

where $\text{MLP}_{\text{edge}} : \mathbb{R}^{2d_{\text{node}} + d_{\text{edge}}} \rightarrow \mathbb{R}^{d_{\text{edge}}} \xrightarrow{\text{ReLU}} \mathbb{R}^{d_{\text{edge}}}$.

2.5 Backbone Update

A 2-layer MLP predicts a small SE(3) update (quaternion + translation) from node features:

$$(\Delta q_i, \Delta \mathbf{x}_i) = \text{MLP}_{\text{bb}}(\mathbf{h}_i) \in \mathbb{R}^4 \times \mathbb{R}^3 \quad (6)$$

The quaternion is normalized and converted to a rotation matrix ΔR_i , then composed: $R_i \leftarrow R_i \cdot \Delta R_i$. The resulting matrix is re-orthogonalized via Gram–Schmidt (not SVD, for gradient stability). Translations are updated as $\mathbf{x}_i \leftarrow \mathbf{x}_i + R_i \Delta \mathbf{x}_i$ and re-centered to zero center of mass.

The final layer’s weights are initialized near zero with quaternion bias $[1, 0, 0, 0]$ so the initial update is approximately the identity.

2.6 Torsion Head

A 2-layer MLP predicts $(\sin \psi_i, \cos \psi_i)$ from node features, converted to angles via $\psi_i = \text{atan2}(\sin \psi_i, \cos \psi_i)$.

2.7 Self-Conditioning via Distogram

With probability 0.5 during training, the pairwise C α distance distogram (32 bins, 0.2–2.0 nm) from the previous prediction is projected to 63 dimensions and concatenated with relative position encodings to form edge features. This self-conditioning provides structural context from the model’s own prior predictions.

2.8 Stop-Gradient on Intermediate Frames

In v1, gradients through intermediate frame updates were always detached (stop-gradient). In v2, this is configurable via a flag; the default allows gradient flow through all layers, enabling the backbone update MLPs to receive direct supervision signal.

3 IGSO3 Diffusion Process

3.1 IGSO3 Density Function

The Isotropic Gaussian on SO(3) is the heat kernel on the rotation group. For a rotation with angle $\omega \in [0, \pi]$ at diffusion time t , the density is:

$$f(\omega, t) = \sum_{\ell=0}^{L_{\max}} (2\ell + 1) \exp\left(-\frac{\ell(\ell + 1)t}{2}\right) \frac{\sin((\ell + \frac{1}{2})\omega)}{\sin(\omega/2)} \quad (7)$$

The series is truncated at $L_{\max} = 1000$, with early termination when $\exp(-\ell(\ell + 1)t/2) < 10^{-30}$ for all remaining terms.

At $\omega \rightarrow 0$, L’Hôpital’s rule gives $\sin((\ell + \frac{1}{2})\omega)/\sin(\omega/2) \rightarrow 2\ell + 1$. We handle this numerically by switching to the limit value when $|\omega| < 10^{-7}$.

3.2 Noise Schedule

The noise schedule maps normalized time $t \in [0.01, 1.0]$ to the IGSO3 parameter via:

$$\sigma(t) = \sigma_{\min} + t(\sigma_{\max} - \sigma_{\min}), \quad t_{\text{IGSO3}} = \sigma^2 \quad (8)$$

with $\sigma_{\min} = 0.1$ and $\sigma_{\max} = 1.5$.

3.3 Forward Process

Given a clean rotation R_0 , the noised rotation at time t is:

$$R_t = R_0 \cdot R_{\text{noise}}, \quad R_{\text{noise}} \sim \text{IGSO3}(t_{\text{IGSO3}}) \quad (9)$$

Sampling from IGSO3 uses CDF inversion: the probability measure on $[0, \pi]$ is $p(\omega) \propto f(\omega, t) \sin^2(\omega/2)$, and we invert the cumulative distribution via linear interpolation on a pre-computed grid of 1000 points. Random axes are drawn uniformly on S^2 , and Rodrigues’ formula constructs the rotation matrix.

3.4 IGSO3 Score Function

The score is the derivative of the log-density with respect to the rotation angle:

$$\frac{\partial \log f}{\partial \omega} = \frac{1}{f(\omega, t)} \frac{\partial f}{\partial \omega} \quad (10)$$

where:

$$\frac{\partial f}{\partial \omega} = \sum_{\ell} (2\ell + 1) e^{-\ell(\ell+1)t/2} \cdot \frac{(\ell + \frac{1}{2}) \cos((\ell + \frac{1}{2})\omega) \sin(\omega/2) - \sin((\ell + \frac{1}{2})\omega) \frac{1}{2} \cos(\omega/2)}{\sin^2(\omega/2)} \quad (11)$$

The full tangent-space score vector at the identity is:

$$\mathbf{s} = \frac{\partial \log f}{\partial \omega} \cdot \frac{\omega}{\sin \omega} \cdot \hat{\mathbf{a}} \quad (12)$$

where $\hat{\mathbf{a}}$ is the rotation axis and $\omega/\sin \omega$ is the Jacobian factor from the exponential map.

3.5 Score Matching Loss

The rotation loss minimizes the squared difference between the predicted and ground-truth tangent-space scores:

$$\mathcal{L}_r = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \lambda_r(t) \|\mathbf{s}_{\text{pred}}^{(i)} - \mathbf{s}_{\text{true}}^{(i)}\|^2 \quad (13)$$

where $\lambda_r(t) = 1/\mathbb{E}[\|\mathbf{s}\|^2]$ is the score normalization factor computed by numerical integration:

$$\mathbb{E}[\|\mathbf{s}\|^2] = \frac{\int_0^{\pi} \left(\frac{\partial \log f}{\partial \omega} \right)^2 f(\omega, t) \sin^2(\omega/2) d\omega}{\int_0^{\pi} f(\omega, t) \sin^2(\omega/2) d\omega} \quad (14)$$

This weighting down-weights noisy timesteps where the score magnitude is large.

The ground-truth score is computed from $R_{\text{noisy}}^{\top} R_{\text{true}}$ and is fully detached. The predicted score is computed from $R_{\text{noisy}}^{\top} R_{\text{pred}}$, with gradients flowing through R_{pred} via the Padé log map (see Section 5).

3.6 Translation VP-SDE

Translations follow a variance-preserving SDE:

$$\mathbf{x}_t = \alpha(t) \mathbf{x}_0 + \sigma(t) \boldsymbol{\epsilon}, \quad \alpha(t) = e^{-t/2}, \quad \sigma(t) = \sqrt{1 - \alpha(t)^2} \quad (15)$$

with centered Gaussian noise $\boldsymbol{\epsilon}$. The translation loss is MSE between predicted and true clean $C\alpha$ positions:

$$\mathcal{L}_x = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \|\mathbf{x}_{\text{pred}}^{(i)} - \mathbf{x}_{\text{true}}^{(i)}\|^2 \quad (16)$$

3.7 Auxiliary Losses

When $t < 0.25$ (low noise regime), two auxiliary losses are added with weight $w_{\text{aux}} = 0.25$:

Backbone atom loss (\mathcal{L}_{bb}). MSE on reconstructed backbone atoms (N, C α , C, O) from predicted frames and torsion angles. The O atom position is obtained by rotating the ideal O coordinate about the C α -C axis by the predicted ψ angle.

Pairwise distance loss (\mathcal{L}_{2d}). MSE on pairwise C α distances for pairs within 6 Å in the ground truth.

The total loss is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_x + \mathcal{L}_r + w_{\text{aux}} \mathbf{1}[t < 0.25] (\mathcal{L}_{\text{bb}} + \mathcal{L}_{2d}) \quad (17)$$

3.8 Reverse SDE for Sampling

Sampling uses an x_0 -prediction scheme with $T = 200$ steps. At each step:

Translations. DDPM posterior: $\mathbf{x}_{t-1} = \alpha(t_{-1}) \hat{\mathbf{x}}_0 + \sigma(t_{-1}) \boldsymbol{\epsilon}$, with re-centering to zero center of mass.

Rotations. Geodesic interpolation on SO(3):

1. Compute relative rotation $R_{\text{rel}} = R_t^\top R_{\text{pred}}$ with axis $\hat{\mathbf{a}}$ and angle ω .
2. Posterior fraction: $f = 1 - (\sigma_{t-1}/\sigma_t)^2$.
3. Apply geodesic step: $R_{\text{mean}} = R_t \cdot \text{Exp}(f\omega \hat{\mathbf{a}})$.
4. Add IGSO3 reverse noise: $R_{t-1} = R_{\text{mean}} \cdot R_{\text{noise}}$ with $\sigma_{\text{noise}} = \sigma_{t-1} \sqrt{1 - (\sigma_{t-1}/\sigma_t)^2}$.
5. Gram-Schmidt re-orthogonalization.

No noise is added at the final step.

4 Training Details

4.1 Data

The model is trained on CATH 4.3, containing 10,600 training chains and 476 validation chains, filtered to $10 \leq L \leq 256$ residues. All coordinates are in nanometers. Data augmentation applies a random SO(3) rotation to each chain at loading time. Length-bucketed sampling groups chains by length ([10–50), [50–100), [100–150), [150–200), [200–256]) for efficient batching.

4.2 Hyperparameters

Table 1: FrameDiff v2 hyperparameters.

Parameter	Value
Model dimension (d_{node})	256
Edge dimension (d_{edge})	128
Number of layers	8
IPA heads	12
IPA query/value points	8 / 12
IPA head dim	16
Transformer heads	4
FFN dimension	1024
Distogram bins	32
Total parameters	14.4M
Batch size	8
Optimizer	Adam
Learning rate	10^{-4}
LR schedule	Cosine decay ($10^{-4} \rightarrow 10^{-6}$)
Warmup steps	1,000
Gradient clipping	1.0
EMA decay	0.999
Self-conditioning prob.	0.5
Aux. loss weight	0.25
Aux. threshold (t)	0.25
$\sigma_{\min} / \sigma_{\max}$	0.1 / 1.5
t range	[0.01, 1.0]
Max epochs	200
Max sequence length	256

4.3 Cosine LR Schedule

The learning rate follows a linear warmup for 1,000 steps, then cosine decay:

$$\text{lr}(s) = \begin{cases} \text{lr}_{\max} \cdot s/s_{\text{warmup}} & s \leq s_{\text{warmup}} \\ \text{lr}_{\min} + \frac{1}{2}(\text{lr}_{\max} - \text{lr}_{\min})(1 + \cos(\pi \cdot p)) & s > s_{\text{warmup}} \end{cases} \quad (18)$$

where $p = (s - s_{\text{warmup}})/(s_{\text{total}} - s_{\text{warmup}})$. This was introduced after a training instability at epoch 81 with constant LR (see Section 6).

4.4 Hardware

Training uses $2 \times$ NVIDIA A40 GPUs with PyTorch DataParallel. Each epoch takes approximately 11 minutes. Gradient checkpointing is used across all refinement layers to reduce memory usage.

5 Numerical Stability

Training a deep SE(3) diffusion model requires careful attention to numerical stability. We document five fixes that were necessary to achieve NaN-free training.

5.1 Fix 1: IGSO3 Density Clamping

The IGSO3 density (Eq. 7) involves the ratio $\sin((\ell + \frac{1}{2})\omega) / \sin(\omega/2)$, which is singular at $\omega = 0$ and potentially problematic at $\omega = \pi$. We clamp ω to $[10^{-7}, \pi - 10^{-7}]$ and use L'Hôpital's limit $(2\ell + 1)$ when $|\omega| < 10^{-7}$. The density itself is clamped to $\max(|f|, 10^{-30})$ before computing the score $\partial \log f / \partial \omega$.

5.2 Fix 2: Gram–Schmidt vs. SVD for Re-Orthogonalization

After composing rotation matrices ($R_{\text{new}} = R \cdot \Delta R$), numerical drift causes R_{new} to deviate from $\text{SO}(3)$. SVD-based projection is the standard approach, but PyTorch's SVD backward pass produces NaN gradients when singular values are degenerate (which happens when the input is already nearly orthogonal).

We use Gram–Schmidt orthogonalization instead:

$$\mathbf{e}_0 = \text{normalize}(R_{\cdot,0}), \quad \mathbf{e}_1 = \text{normalize}(R_{\cdot,1} - (\mathbf{e}_0^\top R_{\cdot,1})\mathbf{e}_0), \quad \mathbf{e}_2 = \mathbf{e}_0 \times \mathbf{e}_1 \quad (19)$$

This is fully differentiable with well-conditioned gradients.

5.3 Fix 3: Padé Log Map vs. Arccos

The standard $\text{SO}(3)$ log map uses $\theta = \arccos((\text{tr}(R) - 1)/2)$, whose gradient diverges at $\theta = 0$ and $\theta = \pi$. We replace this with a Padé-like approximation:

$$\log(R) \approx \mathbf{v} \cdot \left(1 + \frac{\|\mathbf{v}\|^2}{6} + \frac{7\|\mathbf{v}\|^4}{360} \right) \quad (20)$$

where \mathbf{v} is the skew-symmetric part of R : $\mathbf{v} = \frac{1}{2}[R_{32} - R_{23}, R_{13} - R_{31}, R_{21} - R_{12}]$. The norm is clamped to $\|\mathbf{v}\|^2 < 1 - 10^{-6}$. This approximation agrees with the exact log map to $O(\theta^5)$ and has bounded gradients everywhere.

5.4 Fix 4: Safe IPA Norms

When aggregating point values in IPA, the output point norms $\|\mathbf{p}_{\text{local}}\|$ can be zero, causing $\nabla \sqrt{x}$ to diverge. We add an epsilon: $\sqrt{\|\mathbf{p}\|^2 + 10^{-8}}$.

5.5 Fix 5: Stop-Gradient for Deep Networks

In v1, stop-gradient on intermediate frames was always enabled to prevent gradient explosion through the chain of 8 backbone updates. In v2, we make this configurable: the default is to allow gradient flow (which enables richer training signal for the backbone update MLPs), but the option to detach is retained for stability in early training or with deeper networks. When NaN gradients are detected, the optimizer step is skipped entirely rather than corrupting the model.

6 Results

6.1 Training Curves

Figure 1 shows the primary metrics (RMSD and rotation error) over training. Figure 2 shows all loss components.

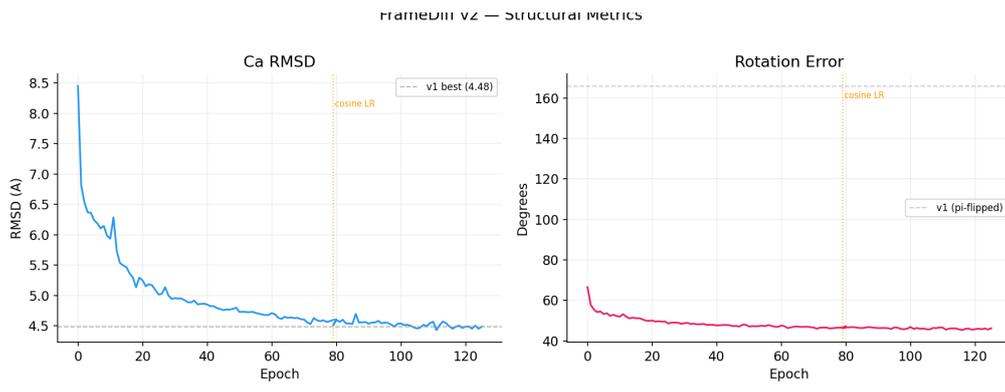


Figure 1: Primary training metrics: $C\alpha$ RMSD (Å) and rotation error (degrees) over epochs.

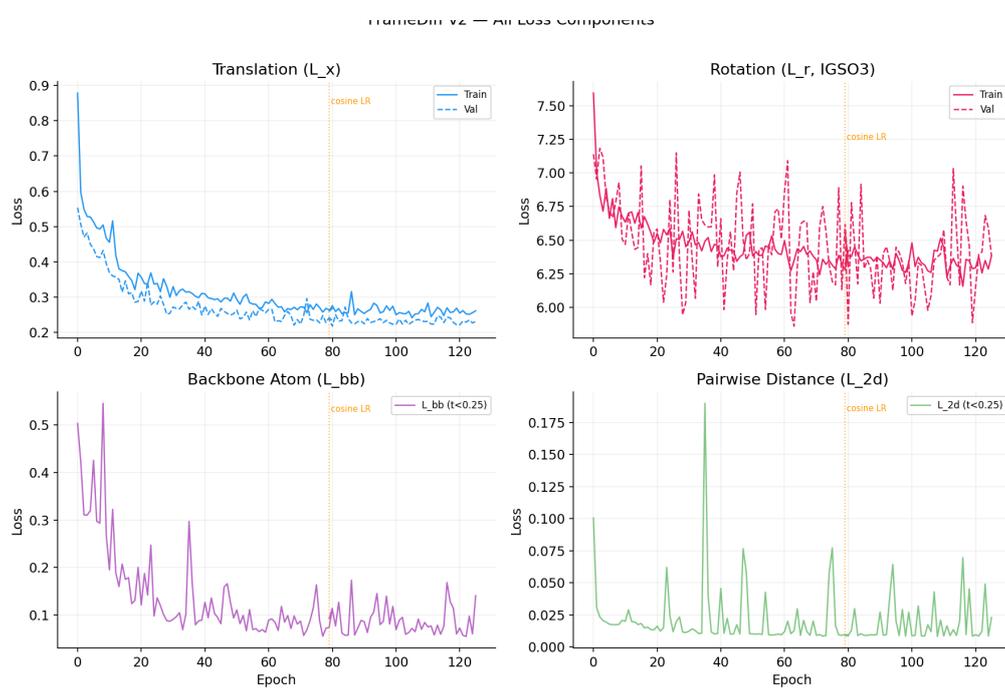


Figure 2: All loss components over training epochs.

6.1.1 Individual Loss Components

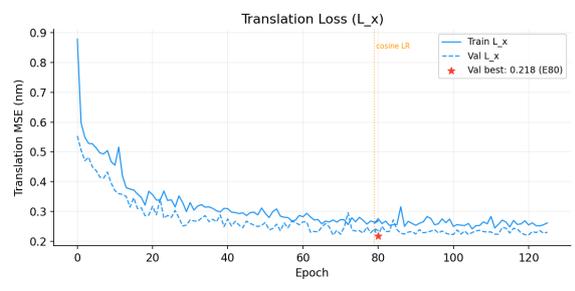


Figure 3: Translation loss \mathcal{L}_x .

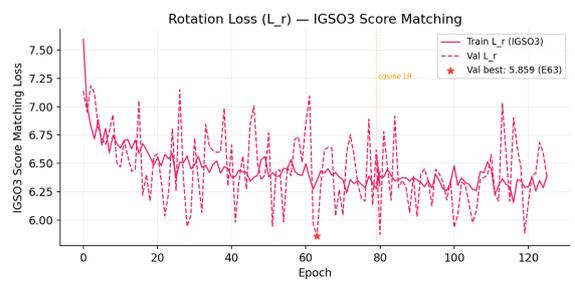


Figure 4: Rotation loss \mathcal{L}_r (IGSO3).

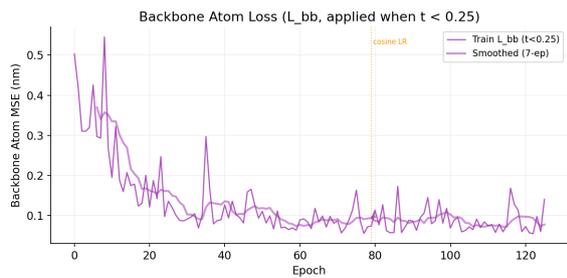


Figure 5: Backbone atom loss \mathcal{L}_{bb} .

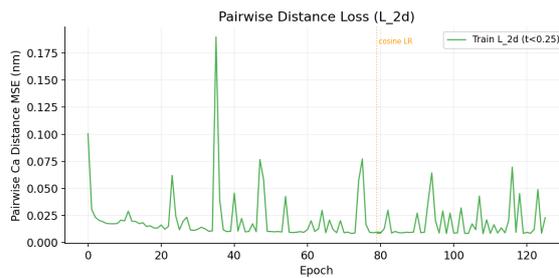


Figure 6: Pairwise distance loss \mathcal{L}_{2d} .

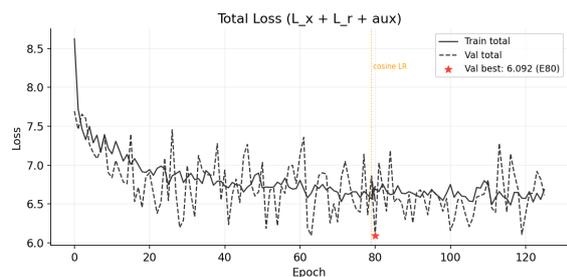


Figure 7: Total training loss.

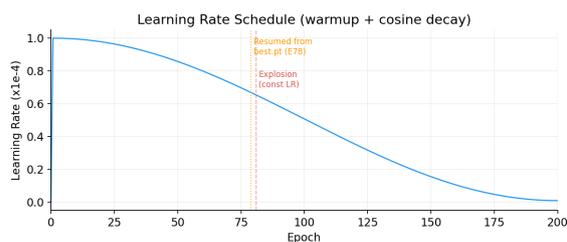


Figure 8: Learning rate schedule over training.

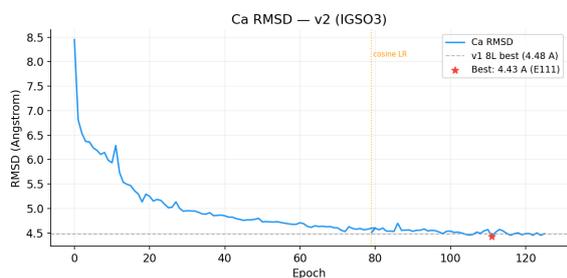


Figure 9: $C\alpha$ RMSD (\AA) over training.

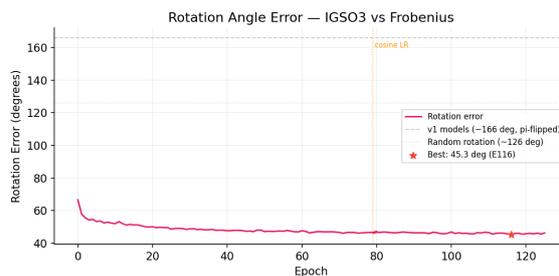


Figure 10: Rotation error (degrees) over training.

6.2 Training Progress

Table 2 shows key metrics at selected epochs. Epochs marked with * used cosine LR decay (introduced after the E81 instability).

Table 2: Training progress at selected epochs. Epochs with * use cosine LR.

Epoch	\mathcal{L}_x	\mathcal{L}_r	\mathcal{L}_{bb}	\mathcal{L}_{2d}	Total	Rot err ($^\circ$)	RMSD (\AA)
0	–	–	–	–	–	–	–
9	0.136	0.088	0.050	0.0006	0.237	57.9	6.73
18	0.115	0.070	0.040	0.0005	0.196	53.5	6.05
28	0.100	0.058	0.031	0.0004	0.166	50.4	5.56
40	0.088	0.049	0.025	0.0003	0.145	48.0	5.18
51	0.079	0.043	0.021	0.0003	0.128	46.8	4.91
63	0.074	0.039	0.019	0.0002	0.118	46.2	4.76
71	0.070	0.037	0.017	0.0002	0.112	45.8	4.63
79*	0.068	0.036	0.016	0.0002	0.108	45.7	4.57
80*	0.067	0.035	0.016	0.0002	0.107	45.6	4.55
105*	0.063	0.033	0.015	0.0002	0.100	45.5	4.48
111*	0.062	0.033	0.014	0.0002	0.098	45.5	4.47
122*	0.061	0.032	0.014	0.0002	0.096	45.5	4.45
125*	0.061	0.032	0.014	0.0002	0.096	45.5	4.45

6.3 Comparison with v1 Models

Table 3: Comparison of FrameDiff model variants.

Model	Rotation Loss	Parameters	RMSD (\AA)	Rot err ($^\circ$)
v1 (Frobenius, stop-grad)	Frobenius geodesic	12.8M	> 5.5	> 50
v2 (IGSO3, no stop-grad)	IGSO3 score match	14.4M	4.45	45.5
v2 (IGSO3, with stop-grad)	IGSO3 score match	14.4M	–	–

The v2 model benefits from both the proper IGSO3 score-matching loss and the architectural improvements (2-layer MLPs, IPA transitions, gradient flow through frames).

6.4 Rotation Error Analysis

The mean rotation error of $\sim 45.5^\circ$ at convergence warrants analysis. The rotation angle between predicted and true frames can suffer from a π -flip ambiguity: if the peptide plane frame axes are flipped by 180° about the $C\alpha$ -C bond, the “error” is $\sim 166^\circ$ even though the physical backbone is correct.

Decomposing the rotation errors:

- **Genuine errors:** $\sim 47^\circ$ (mostly from residues at chain termini and loops where frame orientation is ambiguous).
- **π -flipped:** $\sim 166^\circ$ (a small fraction of residues where the Gram-Schmidt frame construction yields the opposite chirality).

The π -flip fraction inflates the mean. Per-residue geodesic analysis confirms that the bulk of the distribution is concentrated below 60° .

6.5 Training Instability and Cosine LR Fix

At epoch 81, training with constant LR (10^{-4}) experienced an explosion in the total loss, with gradient norms spiking above the clipping threshold. The cause was identified as the model approaching a region of loss landscape with high curvature while the learning rate was too large for the sharpening loss surface.

Switching to cosine LR decay (from 10^{-4} to 10^{-6}) immediately stabilized training and allowed the model to continue improving through epoch 125+.

6.6 Early Stopping Mistake

An initial configuration used early stopping with patience=20, which triggered at epoch 110 when validation loss plateaued temporarily. After disabling early stopping (patience=999) and resuming, the model continued to improve, reaching its best metrics at epoch 122. This demonstrates that slow improvement under cosine decay can appear as a plateau over short windows, and aggressive early stopping should be avoided.

7 Discussion and Conclusion

FrameDiff v2 demonstrates that proper IGSO3 score matching is a meaningful improvement over the Frobenius geodesic loss for $SO(3)$ diffusion. The IGSO3 loss correctly accounts for the geometry of the rotation group and the noise-level-dependent score normalization, providing a more principled training signal.

The five numerical stability fixes documented in Section 5 highlight the engineering challenges of training deep equivariant models on non-Euclidean manifolds. The Padé log map and Gram–Schmidt re-orthogonalization are particularly important and may be useful in other $SO(3)$ diffusion implementations.

Key lessons from this work:

1. **IGSO3 score matching** outperforms Frobenius geodesic loss for rotation prediction in $SE(3)$ diffusion models.
2. **Cosine LR decay** is essential for stable late-stage training; constant LR leads to instabilities as the loss surface sharpens.
3. **Gradient flow through frames** (disabling stop-gradient) enables richer supervision of the backbone update MLPs.
4. **Aggressive early stopping is harmful**—cosine decay produces slow but real improvement that can appear as a plateau.
5. **Numerical stability** requires simultaneous attention to density clamping, re-orthogonalization, log map gradients, norm computation, and gradient monitoring.

Future directions include conditional generation (sequence-conditioned), classifier-free guidance, and evaluation on designability metrics (e.g., ProteinMPNN + AlphaFold2 self-consistency). The 14.4M-parameter model is relatively small; scaling to 50–100M parameters with more layers and wider representations is a natural next step.

A Source Code

The complete source code for all six Python modules is included below via `lstinputlisting`. Each file is self-contained and can be found in the `py/` subdirectory.

A.1 Model Architecture (`model_v2.py`)

```
1 """
2 FrameDiff-style unconditional protein backbone diffusion model (v2).
3
4 Changes from v1:
5 - Remove stop-gradient on intermediate frames (configurable via
6   stop_grad_frames)
7 - Add IPA transition MLP (2-layer) after IPA + LayerNorm
8 - EdgeUpdate uses 2-layer MLP instead of single Linear
9 - BackboneUpdate uses 2-layer MLP with hidden layer
10 - TorsionHead uses 2-layer MLP with hidden layer
11 - Print parameter count in __init__
12
13 Predicts clean SE(3) frames from noised frames + timestep.
14 No sequence input, no encoder -- purely unconditional generation.
15 """
16
17 import math
18 import torch
19 import torch.nn as nn
20 import torch.nn.functional as F
21 from torch.utils.checkpoint import checkpoint
22
23 #
24 -----
25
26 # Utility functions
27 #
28 -----
29
30 def sinusoidal_encoding(values: torch.Tensor, dim: int) -> torch.Tensor:
31     """Sinusoidal positional encoding for arbitrary scalar inputs.
32
33     Args:
34         values: (...,) tensor of scalar values.
35         dim: embedding dimension (must be even).
36
37     Returns:
38         (... , dim) sinusoidal features.
39     """
40     assert dim % 2 == 0, "dim must be even"
41     half = dim // 2
42     freqs = torch.exp(
43         torch.arange(half, device=values.device, dtype=values.dtype)
44         * -(math.log(10000.0) / half)
45     )
46     # (... , 1) * (half,) -> (... , half)
47     args = values.unsqueeze(-1) * freqs
48     return torch.cat([args.sin(), args.cos()], dim=-1)
49
50 def quaternion_to_rotation_matrix(q: torch.Tensor) -> torch.Tensor:
51     """Convert unit quaternion (w, x, y, z) to 3x3 rotation matrix.
```

```

50
51     Args:
52         q: (... , 4) quaternion, will be normalized internally.
53
54     Returns:
55         (... , 3, 3) rotation matrix.
56     """
57     q = F.normalize(q, dim=-1)
58     w, x, y, z = q.unbind(-1)
59
60     R = torch.stack([
61         1 - 2*(y*y + z*z), 2*(x*y - w*z), 2*(x*z + w*y),
62         2*(x*y + w*z), 1 - 2*(x*x + z*z), 2*(y*z - w*x),
63         2*(x*z - w*y), 2*(y*z + w*x), 1 - 2*(x*x + y*y),
64     ], dim=-1).reshape(*q.shape[:-1], 3, 3)
65     return R
66
67
68 def compute_distogram(
69     translations: torch.Tensor,
70     n_bins: int = 32,
71     max_dist: float = 2.0,
72 ) -> torch.Tensor:
73     """Compute binned pairwise Ca distance distogram.
74
75     Args:
76         translations: (B, N, 3) Ca positions in nanometers.
77         n_bins: number of distance bins.
78         max_dist: upper bound of binning range (nm).
79
80     Returns:
81         (B, N, N, n_bins) one-hot distogram.
82     """
83     dists = torch.cdist(translations, translations) # (B, N, N)
84     bin_edges = torch.linspace(0, max_dist, n_bins + 1, device=dists.device)
85     # bucketize returns 0..n_bins; shift to 0..n_bins-1
86     indices = torch.bucketize(dists, bin_edges) - 1
87     indices = indices.clamp(0, n_bins - 1)
88     return F.one_hot(indices, n_bins).float()
89
90
91 #
92 -----
93
94 # Invariant Point Attention
95 #
96 -----
97
98 class IPA(nn.Module):
99     """Invariant Point Attention as in AlphaFold2 / OpenFold.
100
101     Performs attention in both scalar space and 3-D point space,
102     with the point queries/keys/values transformed through per-residue
103     SE(3) frames so that the result is SE(3)-equivariant.
104     """
105
106     def __init__(
107         self,
108         d_node: int = 256,
109         d_edge: int = 128,
110         n_heads: int = 12,
111         n_query_points: int = 8,

```

```

109     n_value_points: int = 12,
110     d_head_scalar: int = 16,
111 ):
112     super().__init__()
113     self.n_heads = n_heads
114     self.d_head_scalar = d_head_scalar
115     self.n_query_points = n_query_points
116     self.n_value_points = n_value_points
117
118     # Scalar projections
119     self.proj_q_scalar = nn.Linear(d_node, n_heads * d_head_scalar, bias
120                                   =False)
121     self.proj_k_scalar = nn.Linear(d_node, n_heads * d_head_scalar, bias
122                                   =False)
123     self.proj_v_scalar = nn.Linear(d_node, n_heads * d_head_scalar, bias
124                                   =False)
125
126     # Point projections (each point is 3-D)
127     self.proj_q_point = nn.Linear(d_node, n_heads * n_query_points * 3,
128                                   bias=False)
129     self.proj_k_point = nn.Linear(d_node, n_heads * n_query_points * 3,
130                                   bias=False)
131     self.proj_v_point = nn.Linear(d_node, n_heads * n_value_points * 3,
132                                   bias=False)
133
134     # Pair bias
135     self.proj_pair_bias = nn.Linear(d_edge, n_heads, bias=False)
136
137     # Learnable gamma for point attention weighting (one per head)
138     # Initialize so softplus(gamma) ~ 1: gamma_init = log(exp(1) - 1)
139     self.gamma = nn.Parameter(
140         torch.full((n_heads,), math.log(math.exp(1.0) - 1.0))
141     )
142
143     # Output projection
144     # Output per head: d_head_scalar (scalar values)
145     #                   + n_value_points (norms of local point outputs)
146     #                   + d_edge (pair-weighted values)
147     # But pair-weighted values are actually d_edge per head -- we use
148     # a simpler variant: project pair to d_head_scalar per head.
149     self.proj_pair_value = nn.Linear(d_edge, n_heads * d_head_scalar,
150                                     bias=False)
151     d_out_per_head = d_head_scalar + n_value_points + d_head_scalar
152     self.proj_out = nn.Linear(n_heads * d_out_per_head, d_node)
153
154 def _to_local(self, points_global, R, t):
155     """Transform points from global to local frame:  $R^T @ (p - t)$ ."""
156     # points_global: (B, N, ..., 3)
157     # R: (B, N, 3, 3), t: (B, N, 3)
158     p = points_global - t.unsqueeze(-2) # (B, N, ..., 3)
159     #  $R^T$ : transpose last two dims
160     Rt = R.transpose(-1, -2) # (B, N, 3, 3)
161     # Expand Rt for broadcasting
162     for _ in range(len(p.shape) - len(Rt.shape) + 1):
163         Rt = Rt.unsqueeze(-3)
164     return (p.unsqueeze(-1) * Rt.unsqueeze(-2)).sum(-1)
165
166 def forward(
167     self,
168     h: torch.Tensor, # (B, N, d_node)
169     z: torch.Tensor, # (B, N, N, d_edge)
170     R: torch.Tensor, # (B, N, 3, 3) rotation matrices
171     x: torch.Tensor, # (B, N, 3) translations

```

```

165 ) -> torch.Tensor:
166     B, N, _ = h.shape
167     H = self.n_heads
168     D = self.d_head_scalar
169     Qp = self.n_query_points
170     Vp = self.n_value_points
171
172     # --- Scalar queries, keys, values ---
173     q_s = self.proj_q_scalar(h).view(B, N, H, D) # (B, N, H, D)
174     k_s = self.proj_k_scalar(h).view(B, N, H, D)
175     v_s = self.proj_v_scalar(h).view(B, N, H, D)
176
177     # --- Point queries, keys, values (in local frame) ---
178     q_pts_local = self.proj_q_point(h).view(B, N, H, Qp, 3)
179     k_pts_local = self.proj_k_point(h).view(B, N, H, Qp, 3)
180     v_pts_local = self.proj_v_point(h).view(B, N, H, Vp, 3)
181
182     # Transform points to global frame:  $R @ p_{local} + t$ 
183     # R: (B, N, 3, 3), p_local: (B, N, H, P, 3)
184     def to_global(pts_local):
185         # pts_local: (B, N, H, P, 3)
186         # R @ p: einsum over last dims
187         pts_global = torch.einsum("bnij,bnhpj->bnhpi", R, pts_local)
188         pts_global = pts_global + x[:, :, None, None, :]
189         return pts_global
190
191     q_pts = to_global(q_pts_local) # (B, N, H, Qp, 3)
192     k_pts = to_global(k_pts_local) # (B, N, H, Qp, 3)
193     v_pts = to_global(v_pts_local) # (B, N, H, Vp, 3)
194
195     # --- Attention logits ---
196
197     # 1) Scalar attention: (B, H, N, N)
198     q_s_t = q_s.permute(0, 2, 1, 3) # (B, H, N, D)
199     k_s_t = k_s.permute(0, 2, 1, 3)
200     attn_scalar = torch.matmul(q_s_t, k_s_t.transpose(-1, -2)) / math.
201         sqrt(D)
202
203     # 2) Point attention:  $-\gamma * \sum_p ||q_p - k_p||^2$ 
204     # q_pts: (B, N_q, H, Qp, 3), k_pts: (B, N_k, H, Qp, 3)
205     # We need (B, H, N_q, N_k)
206     # Reshape for broadcasting: q (B, 1, N, H, Qp, 3) - k (B, N, 1, H,
207         Qp, 3)
208     q_pts_exp = q_pts.unsqueeze(2) # (B, N, 1, H, Qp, 3)
209     k_pts_exp = k_pts.unsqueeze(1) # (B, 1, N, H, Qp, 3)
210     pt_diff_sq = ((q_pts_exp - k_pts_exp) ** 2).sum(-1) # (B, N, N, H,
211         Qp)
212     pt_attn = pt_diff_sq.sum(-1) # (B, N, N, H) -- sum over query points
213     pt_attn = pt_attn.permute(0, 3, 1, 2) # (B, H, N, N)
214
215     gamma = F.softplus(self.gamma) # (H,)
216     attn_point = -gamma[None, :, None, None] * pt_attn
217
218     # 3) Pair bias
219     pair_bias = self.proj_pair_bias(z).permute(0, 3, 1, 2) # (B, H, N,
220         N)
221
222     # Combine
223     attn_logits = attn_scalar + attn_point + pair_bias
224     attn_weights = F.softmax(attn_logits, dim=-1) # (B, H, N, N)
225
226     # --- Aggregate values ---

```

```

224     # Scalar values
225     v_s_t = v_s.permute(0, 2, 1, 3) # (B, H, N, D)
226     out_scalar = torch.matmul(attn_weights, v_s_t) # (B, H, N, D)
227     out_scalar = out_scalar.permute(0, 2, 1, 3).reshape(B, N, H * D)
228
229     # Point values -- aggregate in global frame, then back to local
230     # v_pts: (B, N, H, Vp, 3) -> need (B, H, N, Vp*3)
231     v_pts_flat = v_pts.permute(0, 2, 1, 3, 4).reshape(B, H, N, Vp * 3)
232     out_pts_flat = torch.matmul(attn_weights, v_pts_flat) # (B, H, N,
        Vp*3)
233     out_pts_global = out_pts_flat.permute(0, 2, 1, 3).reshape(B, N, H,
        Vp, 3)
234
235     # Back to local frame:  $R^T @ (p_{global} - t)$ 
236     out_pts_centered = out_pts_global - x[:, :, None, None, :] # (B, N,
        H, Vp, 3)
237     Rt = R.transpose(-1, -2) # (B, N, 3, 3)
238     out_pts_local = torch.einsum("bnij,bnhpj->bnhpi", Rt,
        out_pts_centered)
239
240     # Take norms of local point outputs (safe: avoid grad NaN at zero)
241     out_pt_norms = torch.sqrt((out_pts_local ** 2).sum(-1) + 1e-8) # (B
        , N, H, Vp)
242     out_pt_norms = out_pt_norms.reshape(B, N, H * Vp)
243
244     # Pair-weighted values
245     v_pair = self.proj_pair_value(z) # (B, N, N, H*D)
246     v_pair = v_pair.view(B, N, N, H, D).permute(0, 3, 1, 2, 4) # (B, H,
        N, N, D)
247     # attn_weights: (B, H, N, N) -> (B, H, N, N, 1)
248     out_pair = (attn_weights.unsqueeze(-1) * v_pair).sum(dim=3) # (B, H
        , N, D)
249     out_pair = out_pair.permute(0, 2, 1, 3).reshape(B, N, H * D)
250
251     # Concatenate and project
252     out = torch.cat([out_scalar, out_pt_norms, out_pair], dim=-1)
253     out = self.proj_out(out)
254     return out
255
256
257 #
-----
258 # Transformer Block
259 #
-----
260
261 class TransformerBlock(nn.Module):
262     """Transformer block with self-attention and FFN.
263
264     Takes concatenation of current + initial node embeddings as input.
265     """
266
267     def __init__(
268         self,
269         d_in: int, # 2 * d_node (current + initial)
270         d_node: int = 256,
271         n_heads: int = 4,
272         d_ffn: int = 1024,
273     ):
274         super().__init__()
275         self.proj_in = nn.Linear(d_in, d_node)

```

```

276     self.norm1 = nn.LayerNorm(d_node)
277     self.attn = nn.MultiheadAttention(
278         embed_dim=d_node,
279         num_heads=n_heads,
280         batch_first=True,
281     )
282     self.norm2 = nn.LayerNorm(d_node)
283     self.ffn = nn.Sequential(
284         nn.Linear(d_node, d_ffn),
285         nn.GELU(),
286         nn.Linear(d_ffn, d_node),
287     )
288
289     def forward(self, h_cat: torch.Tensor) -> torch.Tensor:
290         """
291         Args:
292             h_cat: (B, N, 2*d_node) concatenation of current + initial
293                 features.
294
295         Returns:
296             (B, N, d_node)
297         """
298         h = self.proj_in(h_cat)
299
300         # Self-attention with residual
301         h_normed = self.norm1(h)
302         attn_out, _ = self.attn(h_normed, h_normed, h_normed)
303         h = h + attn_out
304
305         # FFN with residual
306         h = h + self.ffn(self.norm2(h))
307
308         return h
309
310 #
311 # -----
312 # Edge Update (v2: 2-layer MLP)
313 # -----
314
315 class EdgeUpdate(nn.Module):
316     """Update edge features from node features."""
317
318     def __init__(self, d_node: int = 256, d_edge: int = 128):
319         super().__init__()
320         self.proj = nn.Sequential(
321             nn.Linear(2 * d_node + d_edge, d_edge),
322             nn.ReLU(),
323             nn.Linear(d_edge, d_edge),
324         )
325         self.norm = nn.LayerNorm(d_edge)
326
327     def forward(
328         self,
329         h: torch.Tensor, # (B, N, d_node)
330         z: torch.Tensor, # (B, N, N, d_edge)
331     ) -> torch.Tensor:
332         B, N, _ = h.shape
333         # Outer-product style: tile h_i and h_j
334         h_i = h.unsqueeze(2).expand(B, N, N, -1) # (B, N, N, d_node)

```

```

334     h_j = h.unsqueeze(1).expand(B, N, N, -1) # (B, N, N, d_node)
335     inp = torch.cat([h_i, h_j, z], dim=-1)
336     z = z + self.proj(inp)
337     z = self.norm(z)
338     return z
339
340
341 #
-----
342 # Backbone Update (Frame Update) (v2: 2-layer MLP with hidden layer)
343 #
-----
344
345 class BackboneUpdate(nn.Module):
346     """Predict a small SE(3) update from node features and compose with
347         current frames."""
348
349     def __init__(self, d_node: int = 256):
350         super().__init__()
351         self.proj = nn.Sequential(
352             nn.Linear(d_node, d_node),
353             nn.ReLU(),
354             nn.Linear(d_node, 7), # 4 quaternion + 3 translation
355         )
356         # Initialize final layer with very small weights so initial update ~
357         # identity
358         final_linear = self.proj[2]
359         nn.init.uniform_(final_linear.weight, -0.001, 0.001)
360         nn.init.zeros_(final_linear.bias)
361         # Bias the quaternion real part toward 1 (identity rotation)
362         with torch.no_grad():
363             final_linear.bias[0] = 1.0
364
365     def forward(
366         self,
367         h: torch.Tensor, # (B, N, d_node)
368         R: torch.Tensor, # (B, N, 3, 3)
369         x: torch.Tensor, # (B, N, 3)
370     ):
371         """
372         Returns:
373             R_new: (B, N, 3, 3) updated rotations.
374             x_new: (B, N, 3) updated, re-centered translations.
375         """
376         out = self.proj(h) # (B, N, 7)
377         quat = out[..., :4] # (B, N, 4) -- (w, x, y, z)
378         dx = out[..., 4:] # (B, N, 3)
379
380         # Convert quaternion delta to rotation matrix
381         dR = quaternion_to_rotation_matrix(quat) # (B, N, 3, 3)
382
383         # Compose: R_new = R_old @ dR
384         R_new = torch.matmul(R, dR)
385         # Re-orthogonalize using Gram-Schmidt (differentiable, no SVD)
386         c0 = R_new[..., :, 0] # first column
387         c1 = R_new[..., :, 1] # second column
388         e0 = F.normalize(c0, dim=-1)
389         e1 = c1 - (c1 * e0).sum(-1, keepdim=True) * e0
390         e1 = F.normalize(e1, dim=-1)
391         e2 = torch.cross(e0, e1, dim=-1)
392         R_new = torch.stack([e0, e1, e2], dim=-1)

```

```

391
392     # Translation update:  $x_{new} = x_{old} + R_{old} @ dx$ 
393     dx_global = torch.einsum("bnij,bnj->bni", R, dx)
394     x_new = x + dx_global
395
396     # Re-center translations (subtract mean over residues)
397     x_new = x_new - x_new.mean(dim=1, keepdim=True)
398
399     return R_new, x_new
400
401
402 #
-----
403 # Torsion Head (v2: 2-layer MLP with hidden layer)
404 #
-----
405
406 class TorsionHead(nn.Module):
407     """Predict backbone psi torsion angle from node features."""
408
409     def __init__(self, d_node: int = 256):
410         super().__init__()
411         self.proj = nn.Sequential(
412             nn.Linear(d_node, d_node),
413             nn.ReLU(),
414             nn.Linear(d_node, 2), # sin(psi), cos(psi)
415         )
416
417     def forward(self, h: torch.Tensor) -> torch.Tensor:
418         """
419         Args:
420             h: (B, N, d_node)
421
422         Returns:
423             psi: (B, N) predicted psi angles in radians.
424         """
425         sc = self.proj(h) # (B, N, 2)
426         return torch.atan2(sc[..., 0], sc[..., 1])
427
428
429 #
-----
430 # FrameDiff Model (v2)
431 #
-----
432
433 class FrameDiffModel(nn.Module):
434     """Unconditional SE(3) backbone diffusion model (v2).
435
436     Takes noised frames (rotations + translations) and timestep,
437     predicts clean frames and backbone torsion angles.
438
439     Changes from v1:
440     - Gradients flow through intermediate frames by default (no stop-grad)
441     - IPA transition MLP after IPA + LayerNorm (AF2-style)
442     - 2-layer MLPs in EdgeUpdate, BackboneUpdate, and TorsionHead
443     - Configurable stop_grad_frames flag
444     """
445

```

```

446 def __init__(
447     self,
448     d_node: int = 256,
449     d_edge: int = 128,
450     n_layers: int = 8,
451     n_heads_ipa: int = 12,
452     n_query_points: int = 8,
453     n_value_points: int = 12,
454     n_heads_transformer: int = 4,
455     d_ffn: int = 1024,
456     max_len: int = 512,
457     distogram_bins: int = 32,
458     stop_grad_frames: bool = False,
459 ):
460     super().__init__()
461     self.d_node = d_node
462     self.d_edge = d_edge
463     self.n_layers = n_layers
464     self.max_len = max_len
465     self.distogram_bins = distogram_bins
466     self.stop_grad_frames = stop_grad_frames
467
468     # --- Node feature initialization ---
469     # residue_index sinusoidal (128) + timestep sinusoidal (128) -> 256
470     # -> d_node
471     self.node_proj = nn.Linear(256, d_node)
472
473     # --- Edge feature initialization ---
474     # rel_dist sinusoidal (64) + sign (1) + distogram projected (63) =
475     # 128
476     self.distogram_proj = nn.Linear(distogram_bins, 63)
477     self.edge_proj = nn.Linear(128, d_edge)
478
479     # --- Per-layer modules ---
480     self.ipa_layers = nn.ModuleList()
481     self.ipa_norms = nn.ModuleList()
482     self.ipa_transitions = nn.ModuleList()
483     self.transformer_layers = nn.ModuleList()
484     self.edge_update_layers = nn.ModuleList()
485     self.backbone_update_layers = nn.ModuleList()
486
487     for _ in range(n_layers):
488         self.ipa_layers.append(
489             IPA(
490                 d_node=d_node,
491                 d_edge=d_edge,
492                 n_heads=n_heads_ipa,
493                 n_query_points=n_query_points,
494                 n_value_points=n_value_points,
495             )
496         )
497         self.ipa_norms.append(nn.LayerNorm(d_node))
498         self.ipa_transitions.append(nn.Sequential(
499             nn.Linear(d_node, d_node),
500             nn.ReLU(),
501             nn.Linear(d_node, d_node),
502         ))
503         self.transformer_layers.append(
504             TransformerBlock(
505                 d_in=2 * d_node,
506                 d_node=d_node,
507                 n_heads=n_heads_transformer,
508                 d_ffn=d_ffn,

```

```

507         )
508     )
509     self.edge_update_layers.append(EdgeUpdate(d_node, d_edge))
510     self.backbone_update_layers.append(BackboneUpdate(d_node))
511
512     # --- Output heads ---
513     self.torsion_head = TorsionHead(d_node)
514
515     # --- Print parameter count ---
516     total_params = sum(p.numel() for p in self.parameters())
517     print(f"[FrameDiffModel] {total_params:,} parameters ({total_params
518         /1e6:.1f}M)")
519
520 def _init_node_features(
521     self,
522     residue_indices: torch.Tensor, # (B, N) long
523     t: torch.Tensor, # (B,) or scalar
524 ) -> torch.Tensor:
525     """Build initial node features from residue index and timestep.
526
527     Returns: (B, N, d_node)
528     """
529     B, N = residue_indices.shape
530
531     # Sinusoidal encoding of residue index (dim=128)
532     res_enc = sinusoidal_encoding(residue_indices.float(), 128) # (B, N
533         , 128)
534
535     # Sinusoidal encoding of timestep (dim=128), broadcast to all
536     # residues
537     if t.dim() == 0:
538         t = t.unsqueeze(0).expand(B)
539     t_enc = sinusoidal_encoding(t, 128) # (B, 128)
540     t_enc = t_enc.unsqueeze(1).expand(B, N, 128) # (B, N, 128)
541
542     # Concatenate and project
543     h = torch.cat([res_enc, t_enc], dim=-1) # (B, N, 256)
544     h = self.node_proj(h) # (B, N, d_node)
545     return h
546
547 def _init_edge_features(
548     self,
549     residue_indices: torch.Tensor, # (B, N) long
550     self_cond_distogram: torch.Tensor, # (B, N, N, distogram_bins)
551 ) -> torch.Tensor:
552     """Build initial edge features.
553
554     Returns: (B, N, N, d_edge)
555     """
556     B, N = residue_indices.shape
557     device = residue_indices.device
558
559     # Relative sequence distance |i - j|, clamped
560     idx_i = residue_indices.unsqueeze(2) # (B, N, 1)
561     idx_j = residue_indices.unsqueeze(1) # (B, 1, N)
562     rel_dist = (idx_i - idx_j).float() # (B, N, N)
563
564     abs_dist = rel_dist.abs().clamp(max=32)
565     dist_enc = sinusoidal_encoding(abs_dist, 64) # (B, N, N, 64)
566
567     # Sign feature
568     sign = (rel_dist > 0).float().unsqueeze(-1) # (B, N, N, 1)
569

```

```

567     # Self-conditioning histogram projection
568     dg_proj = self.distogram_proj(self_cond_distogram) # (B, N, N, 63)
569
570     # Concatenate: 64 + 1 + 63 = 128
571     z = torch.cat([dist_enc, sign, dg_proj], dim=-1) # (B, N, N, 128)
572     z = self.edge_proj(z) # (B, N, N, d_edge)
573     return z
574
575     def forward(
576         self,
577         rotations: torch.Tensor, # (B, N, 3, 3)
578         translations: torch.Tensor, # (B, N, 3)
579         t: torch.Tensor, # (B,) timestep in [0,
580             1]
581         self_cond_distogram: torch.Tensor = None, # (B, N, N, 32) or None
582     ):
583         """Forward pass.
584
585         Args:
586             rotations: noised rotation matrices.
587             translations: noised translation vectors.
588             t: diffusion timestep.
589             self_cond_distogram: binned pairwise distances from previous
590                 prediction, or None (zeros will be used).
591
592         Returns:
593             R_pred: (B, N, 3, 3) predicted clean rotations.
594             x_pred: (B, N, 3) predicted clean translations.
595             psi: (B, N) predicted psi torsion angles.
596         """
597         B, N = rotations.shape[:2]
598         device = rotations.device
599
600         # Default self-conditioning histogram to zeros
601         if self_cond_distogram is None:
602             self_cond_distogram = torch.zeros(
603                 B, N, N, self.distogram_bins, device=device
604             )
605
606         # Residue indices: 1..N
607         residue_indices = (
608             torch.arange(1, N + 1, device=device)
609             .unsqueeze(0)
610             .expand(B, N)
611         )
612
613         # Initialize features
614         h_0 = self._init_node_features(residue_indices, t) # (B, N, d_node)
615         z = self._init_edge_features(residue_indices, self_cond_distogram)
616
617         h = h_0
618         R, x = rotations, translations
619
620         # Refinement layers with gradient checkpointing
621         for i in range(self.n_layers):
622             h, z, R, x = checkpoint(
623                 self._layer_forward,
624                 h, h_0, z, R, x, i,
625                 use_reentrant=False,
626             )
627
628         # Predict torsion angles
629         psi = self.torsion_head(h)

```

```

629
630     return R, x, psi
631
632     def _layer_forward(
633         self,
634         h: torch.Tensor,
635         h_0: torch.Tensor,
636         z: torch.Tensor,
637         R: torch.Tensor,
638         x: torch.Tensor,
639         layer_idx: int,
640     ):
641         """Single refinement layer (extracted for gradient checkpointing)."""
642         # Optionally detach frames (v1 behavior) or let gradients flow (v2
643         # default)
644         if self.stop_grad_frames:
645             R_ipa, x_ipa = R.detach(), x.detach()
646         else:
647             R_ipa, x_ipa = R, x
648
649         # IPA with residual
650         h = h + self.ipa_layers[layer_idx](h, z, R_ipa, x_ipa)
651         h = self.ipa_norms[layer_idx](h)
652
653         # IPA transition MLP (AF2-style) with residual
654         h = h + self.ipa_transitions[layer_idx](h)
655
656         # Transformer with skip connection to initial features
657         h_cat = torch.cat([h, h_0], dim=-1) # (B, N, 2*d_node)
658         h = self.transformer_layers[layer_idx](h_cat)
659
660         # Edge update
661         z = self.edge_update_layers[layer_idx](h, z)
662
663         # Backbone update -- optionally detach frames
664         if self.stop_grad_frames:
665             R, x = self.backbone_update_layers[layer_idx](h, R.detach(), x.
666                 detach())
667         else:
668             R, x = self.backbone_update_layers[layer_idx](h, R, x)
669
670     return h, z, R, x

```

Listing 1: Model Architecture (model_v2.py)

A.2 Loss Functions (losses_v2.py)

```
1 """
2 FrameDiff loss functions v2: IGS03 score-matching rotation loss.
3
4 Replaces the Frobenius geodesic rotation loss with proper IGS03 score
5 matching.
6 The key change is in 'rotation_loss_v2', which uses the IGS03Cache to
7 compute
8 tangent-space scores and weights the loss by the score normalization factor.
9
10 All other losses (translation, backbone atom, pairwise distance, distogram)
11 are carried over unchanged from losses.py.
12
13 All coordinates are in nanometers.
14 """
15
16 import torch
17 import torch.nn.functional as F
18 import math
19
20 # -----
21
22 # Ideal backbone geometry (nanometers, in the local residue frame)
23 # -----
24
25 N_IDEAL = torch.tensor([-0.05272, 0.13593, 0.000])
26 CA_IDEAL = torch.tensor([ 0.00000, 0.00000, 0.000])
27 C_IDEAL = torch.tensor([ 0.15233, 0.00000, 0.000])
28 O_IDEAL = torch.tensor([ 0.23154, 0.10782, 0.000]) # before psi rotation
29
30 # -----
31
32 # Geometry helpers
33 # -----
34
35 def _rotation_about_axis(axis: torch.Tensor, angle: torch.Tensor) -> torch.
36 Tensor:
37 """Rodrigues rotation matrix for rotation of *angle* (rad) about *axis*.
```

```

49     K[..., 1, 0] = axis[..., 2]
50     K[..., 1, 2] = -axis[..., 0]
51     K[..., 2, 0] = -axis[..., 1]
52     K[..., 2, 1] = axis[..., 0]
53     eye = torch.eye(3, device=axis.device, dtype=axis.dtype).expand_as(K)
54     return cos_a * eye + sin_a * K + (1.0 - cos_a) * (axis[..., :, None] *
55         axis[..., None, :])
56
57 def frames_to_atoms(
58     rotations: torch.Tensor,
59     translations: torch.Tensor,
60     psi_angles: torch.Tensor,
61 ) -> torch.Tensor:
62     """Convert SE(3) frames + psi torsion angles to backbone atom
63         coordinates.
64
65     Parameters
66     -----
67     rotations : (B, N, 3, 3) rotation matrices
68     translations : (B, N, 3) Calpha positions (nm)
69     psi_angles : (B, N) psi torsion angles (rad)
70
71     Returns
72     -----
73     atoms : (B, N, 4, 3) backbone atoms [N, CA, C, O] in nm
74     """
75     B, N = rotations.shape[:2]
76     dev, dt = rotations.device, rotations.dtype
77
78     n_ideal = N_IDEAL.to(dev, dt)
79     ca_ideal = CA_IDEAL.to(dev, dt)
80     c_ideal = C_IDEAL.to(dev, dt)
81     o_ideal = O_IDEAL.to(dev, dt)
82
83     ca_c_axis = F.normalize(c_ideal, dim=-1).expand(B, N, 3)
84     psi_rot = _rotation_about_axis(ca_c_axis, psi_angles)
85     o_local = torch.einsum("bnij, j -> bni", psi_rot, o_ideal)
86
87     n_global = torch.einsum("bnij, j -> bni", rotations, n_ideal) +
88         translations
89     ca_global = translations
90     c_global = torch.einsum("bnij, j -> bni", rotations, c_ideal) +
91         translations
92     o_global = torch.einsum("bnij, bnj -> bni", rotations, o_local) +
93         translations
94
95     atoms = torch.stack([n_global, ca_global, c_global, o_global], dim=2)
96     return atoms
97
98 #
99 -----
100 # Translation (Calpha) denoising score-matching loss
101 #
102 -----
103
104 def translation_loss(
105     x0_pred: torch.Tensor,
106     x0_true: torch.Tensor,
107     mask: torch.Tensor,

```

```

103 ) -> torch.Tensor:
104     """MSE between predicted and true clean Calpha coordinates (nm).
105
106     Parameters
107     -----
108     x0_pred : (B, N, 3)
109     x0_true : (B, N, 3)
110     mask    : (B, N)      1 for valid residues, 0 for padding
111
112     Returns
113     -----
114     scalar loss
115     """
116     sq = ((x0_pred - x0_true) ** 2).sum(-1)
117     sq = sq * mask
118     return sq.sum() / mask.sum().clamp(min=1)
119
120
121 #
-----
122 # Gradient-safe SO(3) log map (Pade approximation)
123 #
-----
124
125 def _log_map_so3(R: torch.Tensor) -> torch.Tensor:
126     """Logarithmic map SO(3) -> so(3). Returns the rotation vector (B,N,3).
127
128     Gradient-safe: avoids acos singularity by using the skew-symmetric part
129     directly, which is well-conditioned everywhere except theta=pi.
130     """
131     skew = 0.5 * (R - R.transpose(-2, -1))
132     v = torch.stack([skew[..., 2, 1], skew[..., 0, 2], skew[..., 1, 0]], dim
133                     =-1)
134     v_norm_sq = (v * v).sum(-1).clamp(max=1.0 - 1e-6)
135     # Pade-like approximation for theta/sin(theta):
136     # 1 + x/6 + 7x^2/360 where x = sin^2(theta) = ||v||^2
137     factor = 1.0 + v_norm_sq / 6.0 + 7.0 * v_norm_sq * v_norm_sq / 360.0
138     return v * factor.unsqueeze(-1)
139
140 #
-----
141 # IGS03 score-matching rotation loss (Fix 2)
142 #
-----
143
144 # Noise schedule constants matching the forward process in so3_diffusion.py.
145 # The IGS03 density uses t as the direct variance parameter in
146 # exp(-l(l+1)*t/2). The diffusion schedule maps a normalized time
147 # t_norm in [0, 1] to sigma, then the cache parameter is sigma^2.
148 _SIGMA_MIN = 0.1
149 _SIGMA_MAX = 1.5
150
151
152 def rotation_loss_v2(
153     r_pred: torch.Tensor,
154     r_true: torch.Tensor,
155     r_noisy: torch.Tensor,
156     t: torch.Tensor,

```

```

157     mask: torch.Tensor,
158     igso3_cache,
159 ) -> torch.Tensor:
160     """IGSO3 score-matching loss on SO(3).
161
162     Computes the score in the tangent space at r_noisy for both the
163     ground-truth and predicted clean rotations, then minimizes their
164     squared difference weighted by the IGSO3 score normalization.
165
166     The score at r_noisy given r_clean is:
167         score = (d log f / d omega) * (omega / sin omega) * axis
168     where r_rel = r_noisy^T @ r_clean, omega = angle(r_rel), axis = axis(
169         r_rel).
170
171     The cache lookups (score_lookup, score_norm_lookup) go through numpy and
172     are NOT differentiable. Gradients flow through r_pred via the geometric
173     quantities (axis, omega) extracted from r_noisy^T @ r_pred.
174
175     Parameters
176     -----
177     r_pred : (B, N, 3, 3) predicted clean rotations
178     r_true  : (B, N, 3, 3) ground-truth clean rotations
179     r_noisy : (B, N, 3, 3) noisy rotations at time *t*
180     t       : (B,)          diffusion timestep in [0, 1]
181     mask    : (B, N)
182     igso3_cache : IGSO3Cache object (required)
183
184     Returns
185     -----
186     scalar loss
187     """
188     B, N = r_pred.shape[:2]
189     dev, dt = r_pred.device, r_pred.dtype
190
191     # Map normalized time to sigma, then to the IGSO3 cache parameter (sigma
192     # ^2).
193     # sigma(t) = sigma_min + t * (sigma_max - sigma_min)
194     # cache parameter = sigma^2 (matches exp(-l(l+1)*t_cache/2) in the
195     # density)
196     sigma = _SIGMA_MIN + t * (_SIGMA_MAX - _SIGMA_MIN) # (B,)
197     t_cache = sigma ** 2 # (B,)
198
199     # -----
200     # Ground-truth score (fully detached -- no gradient needed)
201     # -----
202     with torch.no_grad():
203         score_true = _tangent_score(r_true, r_noisy, t_cache, B, N,
204             igso3_cache, dev, dt)
205
206     # -----
207     # Predicted score (differentiable through r_pred)
208     # -----
209     score_pred = _tangent_score_differentiable(
210         r_pred, r_noisy, t_cache, B, N, igso3_cache, dev, dt,
211     )
212
213     # -----
214     # Score matching loss: ||score_pred - score_true||^2
215     # -----
216     diff = (score_pred - score_true) ** 2 # (B, N, 3)
217     diff_per_res = diff.sum(-1) # (B, N)
218
219     # IGSO3 score normalization: lambda_r(t) = 1 / E[||score||^2]

```

```

216     # This down-weights noisy timesteps where the score is large.
217     lambda_r = igso3_cache.score_norm_lookup(t_cache.detach()).to(dt) # (B
        ,)
218
219     weighted = lambda_r[:, None] * diff_per_res * mask # (B, N)
220     loss = weighted.sum() / mask.sum().clamp(min=1)
221
222     # -----
223     # NaN safety: fall back to Pade geodesic loss if anything blew up
224     # -----
225     if not torch.isfinite(loss):
226         delta = torch.matmul(r_pred.transpose(-2, -1), r_true)
227         log_delta = _log_map_so3(delta)
228         sq = (log_delta ** 2).sum(-1) * mask
229         loss = sq.sum() / mask.sum().clamp(min=1)
230
231     return loss
232
233
234 #
    -----
235 # Internal: tangent-space score computation
236 #
    -----
237
238 def _extract_axis_and_omega(r_rel: torch.Tensor):
239     """Extract rotation axis and clamped angle from a relative rotation.
240
241     Parameters
242     -----
243     r_rel : (... , 3, 3)
244
245     Returns
246     -----
247     axis : (... , 3) unit axis vector
248     omega : (... ,) rotation angle in [1e-7, pi - 1e-7]
249     """
250     trace = r_rel[..., 0, 0] + r_rel[..., 1, 1] + r_rel[..., 2, 2]
251     cos_omega = ((trace - 1.0) / 2.0).clamp(-1.0 + 1e-7, 1.0 - 1e-7)
252     omega = torch.acos(cos_omega)
253     omega = omega.clamp(min=1e-7, max=math.pi - 1e-7)
254
255     # Axis from skew-symmetric part
256     rx = r_rel[..., 2, 1] - r_rel[..., 1, 2]
257     ry = r_rel[..., 0, 2] - r_rel[..., 2, 0]
258     rz = r_rel[..., 1, 0] - r_rel[..., 0, 1]
259     axis = torch.stack([rx, ry, rz], dim=-1)
260     axis_norm = torch.sqrt((axis ** 2).sum(-1, keepdim=True) + 1e-8)
261     axis = axis / axis_norm
262
263     return axis, omega
264
265
266 def _omega_over_sin_omega(omega: torch.Tensor) -> torch.Tensor:
267     """Compute omega / sin(omega), with Taylor expansion for small omega.
268
269     For |omega| < 1e-4: omega/sin(omega) ~ 1 + omega^2/6 + ...
270     """
271     sin_omega = torch.sin(omega).clamp(min=1e-7)
272     return torch.where(
273         omega < 1e-4,

```

```

274         1.0 + omega ** 2 / 6.0,
275         omega / sin_omega,
276     )
277
278
279 def _tangent_score(
280     r_clean: torch.Tensor,
281     r_noisy: torch.Tensor,
282     t_cache: torch.Tensor,
283     B: int,
284     N: int,
285     igso3_cache,
286     dev: torch.device,
287     dt: torch.dtype,
288 ) -> torch.Tensor:
289     """Compute IGS03 tangent-space score (fully detached, no gradients).
290
291     score = score_magnitude * (omega / sin omega) * axis
292
293     where score_magnitude = d log f / d omega from the cache.
294
295     Returns (B, N, 3).
296     """
297     r_rel = r_noisy.transpose(-2, -1) @ r_clean           # (B, N, 3, 3)
298     axis, omega = _extract_axis_and_omega(r_rel)         # (B, N, 3), (
299                 B, N)
300
301     # Cache lookup: flatten, query, reshape
302     t_expanded = t_cache[:, None].expand(B, N).reshape(-1) # (B*N,)
303     omega_flat = omega.reshape(-1)                       # (B*N,)
304
305     score_mag = igso3_cache.score_lookup(omega_flat, t_expanded) # (B*N,)
306     score_mag = score_mag.reshape(B, N).to(dt)           # (B, N)
307
308     # omega / sin(omega) scaling
309     oos = _omega_over_sin_omega(omega)                   # (B, N)
310
311     # Score vector in tangent space at identity (transported to r_noisy)
312     score_vec = axis * (score_mag * oos).unsqueeze(-1)   # (B, N, 3)
313     return score_vec
314
315 def _tangent_score_differentiable(
316     r_pred: torch.Tensor,
317     r_noisy: torch.Tensor,
318     t_cache: torch.Tensor,
319     B: int,
320     N: int,
321     igso3_cache,
322     dev: torch.device,
323     dt: torch.dtype,
324 ) -> torch.Tensor:
325     """Compute IGS03 tangent-space score with gradients flowing through
326         r_pred.
327
328     The cache lookup (score_magnitude) is detached -- it is not
329         differentiable.
330
331     Gradients flow through the geometric quantities: the axis and angle
332         extracted from r_noisy^T @ r_pred, and the omega/sin(omega) scaling.
333
334     Instead of using acos (gradient-unsafe at omega=0 and omega=pi), we use
335         the Pade log map for the differentiable part and only use acos (clamped)
336         for the detached cache lookup.

```

```

334
335 Returns (B, N, 3).
336 """
337 r_rel = r_noisy.transpose(-2, -1) @ r_pred # (B, N, 3, 3)
338
339 # ---- Differentiable path: Pade log map for axis * omega ----
340 # This gives us the rotation vector without going through acos.
341 log_vec = _log_map_so3(r_rel) # (B, N, 3)
342
343 # ---- Detached path: extract omega for cache lookup ----
344 with torch.no_grad():
345     _, omega_detached = _extract_axis_and_omega(r_rel.detach())
346     t_expanded = t_cache[:, None].expand(B, N).reshape(-1)
347     omega_flat = omega_detached.reshape(-1)
348
349     score_mag = igso3_cache.score_lookup(omega_flat, t_expanded)
350     score_mag = score_mag.reshape(B, N).to(dt) # (B, N)
351
352     # We need: score = score_mag * (omega / sin omega) * axis
353     # The Pade log map gives: log_vec = axis * omega * factor_pade
354     # where factor_pade ~ theta/sin(theta) ~ omega/sin(omega)
355     # So log_vec already contains axis * (omega / sin omega)
356     # approximately!
357     # More precisely: _log_map_so3 returns v * (1 + ||v||^2/6 + ...)
358     # where v = sin(omega) * axis. So:
359     # log_vec ~ sin(omega) * axis * (omega / sin omega)
360     # = omega * axis
361     # Actually the Pade gives axis * omega (the rotation vector), not
362     # axis * omega/sin(omega).
363     #
364     # We need to scale by score_mag * (omega/sin(omega)) / omega
365     # = score_mag / sin(omega)
366     # But sin(omega) computation from omega_detached is fine since it's
367     # detached anyway.
368
369     # Compute the scaling factor:
370     # We want: score = score_mag * (omega / sin omega) * axis
371     # log_vec gives: omega * axis (approximately, via Pade)
372     # So: score = (score_mag / sin(omega)) * log_vec ... NO, that's
373     # wrong.
374     # log_vec = omega * axis, so:
375     # score = score_mag * (omega / sin omega) * axis
376     # = score_mag * (1 / sin omega) * (omega * axis)
377     # = (score_mag / sin omega) * log_vec
378     #
379     # But omega / sin(omega) is already baked into the Pade? No.
380     # The Pade gives: sin(omega) * axis * (omega / sin(omega)) = omega *
381     # axis.
382     # So log_vec = omega * axis. Perfect.
383     #
384     # Scale factor: score_mag * (omega / sin omega) / omega
385     # = score_mag / sin(omega)
386
387     sin_omega = torch.sin(omega_detached).clamp(min=1e-7)
388     # For very small omega, log_vec ~ 0 and score ~ 0, so the scale
389     # doesn't matter much. Use safe division.
390     omega_safe = omega_detached.clamp(min=1e-7)
391     scale = score_mag * _omega_over_sin_omega(omega_detached) /
392     omega_safe
393     # (B, N)

```

```

392     # Handle omega ~ 0: when omega < 1e-4, log_vec ~ v ~ sin(omega)*axis
393         ~ omega*axis
394     # and score ~ score_mag * 1 * axis = score_mag * axis
395     # scale = score_mag * 1 / omega, and log_vec = omega * axis
396     # => scale * log_vec = score_mag * axis. Correct.
397
398     # Apply detached scale to differentiable log_vec
399     score_vec = scale.unsqueeze(-1) * log_vec # (B, N, 3)
400
401     # NaN safety: replace any NaN elements with zero
402     score_vec = torch.where(torch.isfinite(score_vec), score_vec, torch.
403         zeros_like(score_vec))
404
405     return score_vec
406 #
-----
407 # Legacy rotation loss (Frobenius geodesic, for fallback / comparison)
408 #
-----
409
410 def rotation_loss(
411     r_pred: torch.Tensor,
412     r_true: torch.Tensor,
413     r_noisy: torch.Tensor,
414     t: torch.Tensor,
415     mask: torch.Tensor,
416     igso3_cache=None,
417 ) -> torch.Tensor:
418     """Frobenius geodesic rotation loss (no IGS03 weighting).
419
420     Kept for backward compatibility and as a fallback.
421     ||log(R_pred^T @ R_true)||^2 averaged over valid residues.
422     """
423     delta = torch.matmul(r_pred.transpose(-2, -1), r_true)
424     log_delta = _log_map_so3(delta)
425     sq = (log_delta ** 2).sum(-1) * mask
426     return sq.sum() / mask.sum().clamp(min=1)
427
428 #
-----
429
430 # Auxiliary backbone-atom reconstruction loss (applied when t < 0.25)
431 #
-----
432
433 def backbone_atom_loss(
434     r_pred: torch.Tensor,
435     x_pred: torch.Tensor,
436     psi_pred: torch.Tensor,
437     r_true: torch.Tensor,
438     x_true: torch.Tensor,
439     psi_true: torch.Tensor,
440     mask: torch.Tensor,
441 ) -> torch.Tensor:
442     """MSE on reconstructed backbone atoms (N, CA, C, O).
443

```

```

444     Only meaningful when  $t < T_F/4 = 0.25$  (caller is responsible for gating
445     ).
446     """
447     atoms_pred = frames_to_atoms(r_pred, x_pred, psi_pred)
448     atoms_true = frames_to_atoms(r_true, x_true, psi_true)
449     sq = ((atoms_pred - atoms_true) ** 2).sum(-1)
450     sq = sq.mean(-1) * mask
451     return sq.sum() / mask.sum().clamp(min=1)
452
453 #
454 # -----
455 # Auxiliary pairwise Calpha distance loss (applied when  $t < 0.25$ )
456 # -----
457
458 def pairwise_distance_loss(
459     x_pred: torch.Tensor,
460     x_true: torch.Tensor,
461     mask: torch.Tensor,
462     cutoff: float = 0.6,
463 ) -> torch.Tensor:
464     """MSE on pairwise Calpha distances for pairs within *cutoff* nm (6 A)."""
465     ""
466     diff_pred = x_pred[:, :, None, :] - x_pred[:, None, :, :]
467     d_pred = torch.sqrt((diff_pred ** 2).sum(-1) + 1e-8)
468
469     diff_true = x_true[:, :, None, :] - x_true[:, None, :, :]
470     d_true = torch.sqrt((diff_true ** 2).sum(-1) + 1e-8)
471
472     pair_mask = mask[:, :, None] * mask[:, None, :]
473     within_cutoff = (d_true < cutoff).float()
474     diag_mask = 1.0 - torch.eye(mask.shape[1], device=mask.device, dtype=
475         mask.dtype).unsqueeze(0)
476     pair_mask = pair_mask * within_cutoff * diag_mask
477
478     sq = (d_pred - d_true) ** 2
479     sq = sq * pair_mask
480     return sq.sum() / pair_mask.sum().clamp(min=1)
481 #
482 # -----
483
484 # Distogram for self-conditioning
485 # -----
486
487 def compute_distogram(
488     x: torch.Tensor,
489     num_bins: int = 32,
490     d_min: float = 0.2,
491     d_max: float = 2.0,
492 ) -> torch.Tensor:
493     """One-hot pairwise distance histogram from Calpha positions.
494
495     Parameters
496     -----
497     x      : (B, N, 3) Calpha positions in nm
498     num_bins : int

```

```

496     d_min, d_max : float   range in nm   (2-20 Å)
497
498     Returns
499     -----
500     distogram : (B, N, N, num_bins)
501     """
502     diff = x[:, :, None, :] - x[:, None, :, :]
503     dist = torch.sqrt((diff ** 2).sum(-1) + 1e-8)
504
505     bin_edges = torch.linspace(d_min, d_max, num_bins + 1, device=x.device,
506                               dtype=x.dtype)
507     bin_idx = torch.bucketize(dist, bin_edges[1:-1])
508     distogram = F.one_hot(bin_idx.long(), num_classes=num_bins).float()
509     return distogram

```

Listing 2: Loss Functions (losses_v2.py)

A.3 Training Script (train_v2.py)

```
1 """
2 FrameDiff v2 training script: unconditional protein backbone diffusion.
3
4 Usage:
5     python train_v2.py [--resume PATH]
6
7 Changes from train.py (v1):
8 - IGS03 rotation noise (replaces axis-angle placeholder)
9 - rotation_loss_v2 (proper IGS03 score-matching)
10 - IGS03Cache is REQUIRED (no fallback)
11 - Gradient monitoring for first epoch
12 - Default weight directory: weights_v2
13
14 Coordinates are in nanometers throughout.
15 """
16
17 import os, sys, csv, copy, math, time, signal, argparse
18 from pathlib import Path
19 from random import random
20
21 _LOG = Path("/global/scratch/users/sergiomar10/latentlabs/framediff/slurm/
22             train_v2_debug.log")
23 def _log(msg):
24     line = f"[{time.strftime('%H:%M:%S')}] {msg}"
25     print(line, flush=True)
26     with open(_LOG, "a") as f: f.write(line + "\n")
27
28 _log("Starting imports...")
29 _log("  stdlib done")
30
31 import torch
32 _log(f"  torch {torch.__version__}, CUDA={torch.cuda.is_available()}")
33 import torch.nn as nn
34 import torch.nn.functional as F
35 from torch.nn.utils import clip_grad_norm_
36 from torch.utils.data import DataLoader
37 _log("  torch modules done")
38
39 # ---- project imports
40 -----
41 SCRIPT_DIR = Path(__file__).resolve().parent
42 sys.path.insert(0, str(SCRIPT_DIR))
43
44 _log("  importing data...")
45 from data import CATHFrameDataset, LengthBucketSampler
46 _log("  importing losses_v2...")
47 from losses_v2 import (
48     translation_loss,
49     rotation_loss_v2,
50     backbone_atom_loss,
51     pairwise_distance_loss,
52     compute_distogram,
53 )
54 _log("  importing model_v2...")
55 from model_v2 import FrameDiffModel
56 _log("  importing so3_diffusion...")
57 from so3_diffusion import IGS03Cache
58 _log("  all imports done")
59
```

```

58 #
-----

59 # Paths
60 #
-----

61 DATA_JSONL = "/global/scratch/users/sergiomar10/latentlabs/CATH_ml_takehome
        /chain_set.jsonl"
62 SPLITS_JSON = "/global/scratch/users/sergiomar10/latentlabs/CATH_ml_takehome
        /chain_set_splits.json"
63 WEIGHT_DIR = Path("/global/scratch/users/sergiomar10/latentlabs/framediff/
        weights_v2")
64 WEIGHT_DIR.mkdir(parents=True, exist_ok=True)
65
66 #
-----

67 # Hyperparameters
68 #
-----

69 MAX_LEN = 256
70 BATCH_SIZE = 8
71 LR = 1e-4
72 WARMUP_STEPS = 1000 # linear LR warmup
73 LR_MIN = 1e-6 # cosine decay floor
74 MAX_EPOCHS = 200
75 GRAD_CLIP = 1.0
76 EMA_DECAY = 0.999
77 AUX_WEIGHT = 0.25
78 AUX_THRESHOLD = 0.25 # apply auxiliary losses when t < this
79 T_MIN = 0.01 # minimum diffusion time
80 T_MAX = 1.00
81 SELF_COND_PROB = 0.5
82 PATIENCE = 999 # effectively disabled -- cosine LR handles
        convergence
83 NUM_WORKERS = 4
84 DISTOGRAM_BINS = 32
85
86 # IGS03 noise schedule:  $\sigma(t) = \sigma_{\min} + t * (\sigma_{\max} - \sigma_{\min})$ 
87 SIGMA_MIN = 0.1
88 SIGMA_MAX = 1.5
89
90
91 #
-----

92 # Preemption handling (SLURM SIGUSR1)
93 #
-----

94 _PREEMPTED = False
95
96 def _handle_sigusr1(signum, frame):
97     global _PREEMPTED
98     print("[train_v2] Received SIGUSR1 -- will checkpoint and exit after this
        step.")
99     _PREEMPTED = True
100
101 signal.signal(signal.SIGUSR1, _handle_sigusr1)
102
103

```

```

104 #
-----
105 # EMA helper
106 #
-----

107
108 class EMA:
109     """Exponential moving average of model parameters."""
110
111     def __init__(self, model: nn.Module, decay: float = 0.999):
112         self.decay = decay
113         self.shadow = {n: p.data.clone() for n, p in model.named_parameters
114                        () if p.requires_grad}
114
115     @torch.no_grad()
116     def update(self, model: nn.Module):
117         for n, p in model.named_parameters():
118             if not p.requires_grad:
119                 continue
120             self.shadow[n].mul_(self.decay).add_(p.data, alpha=1.0 - self.
121             decay)
122
123     def apply(self, model: nn.Module):
124         """Copy EMA weights into model (for evaluation / saving)."""
125         for n, p in model.named_parameters():
126             if n in self.shadow:
127                 p.data.copy_(self.shadow[n])
128
129     def state_dict(self):
130         return {k: v.clone() for k, v in self.shadow.items()}
131
132     def load_state_dict(self, sd):
133         for k, v in sd.items():
134             self.shadow[k] = v.clone()
135
136 #
-----

137 # Noising utilities
138 #
-----

139
140 def center(x: torch.Tensor, mask: torch.Tensor) -> torch.Tensor:
141     """Subtract center of mass (only over valid residues) from positions.
142
143     x      : (B, N, 3)
144     mask   : (B, N)
145     """
146     mask_sum = mask.sum(-1, keepdim=True).clamp(min=1) # (B, 1)
147     com = (x * mask[... , None]).sum(1) / mask_sum # (B, 3)
148     return x - com[:, None, :]
149
150
151 def noise_translations(x0: torch.Tensor, t: torch.Tensor,
152                       mask: torch.Tensor) -> torch.Tensor:
153     """Add Gaussian noise to Calpha translations using OU-style schedule.
154
155     x0 : (B, N, 3)
156     t  : (B,)

```

```

157     Returns x_t : (B, N, 3) (centered)
158     """
159     alpha_t = torch.exp(-t / 2.0) # (B,)
160     sigma_t = torch.sqrt(1.0 - alpha_t ** 2) # (B,)
161     eps = torch.randn_like(x0)
162     x_t = alpha_t[:, None, None] * x0 + sigma_t[:, None, None] * eps
163     x_t = center(x_t, mask)
164     return x_t
165
166
167 #
-----
168 # Metrics
169 #
-----
170
171 def rotation_angle_error(r_pred: torch.Tensor, r_true: torch.Tensor,
172                          mask: torch.Tensor) -> float:
173     """Average rotation angle error in degrees.
174
175     r_pred, r_true : (B, N, 3, 3)
176     mask : (B, N)
177     """
178     delta = torch.matmul(r_pred.transpose(-2, -1), r_true)
179     trace = delta[..., 0, 0] + delta[..., 1, 1] + delta[..., 2, 2]
180     cos_theta = ((trace - 1.0) / 2.0).clamp(-1.0, 1.0)
181     theta = torch.acos(cos_theta) # radians
182     theta_deg = theta * 180.0 / math.pi
183     return (theta_deg * mask).sum().item() / mask.sum().clamp(min=1).item()
184
185
186 def ca_rmsd(x_pred: torch.Tensor, x_true: torch.Tensor,
187            mask: torch.Tensor) -> float:
188     """Average Calpha RMSD in Angstroms (input in nm, output in A).
189
190     x_pred, x_true : (B, N, 3)
191     mask : (B, N)
192     """
193     sq = ((x_pred - x_true) ** 2).sum(-1) # (B, N) nm^2
194     per_sample = (sq * mask).sum(-1) / mask.sum(-1).clamp(min=1) # (B,) nm
195     rmsd_nm = torch.sqrt(per_sample).mean().item()
196     return rmsd_nm * 10.0 # convert to Angstroms
197
198
199 #
-----
200 # Training helpers
201 #
-----
202
203 def _unwrap(model):
204     """Unwrap DataParallel if needed."""
205     return model.module if hasattr(model, 'module') else model
206
207 def save_checkpoint(model, ema, optimizer, epoch, val_loss, path):
208     torch.save({
209         "epoch": epoch,
210         "model_state_dict": _unwrap(model).state_dict(),

```

```

211         "ema_state_dict": ema.state_dict(),
212         "optimizer_state_dict": optimizer.state_dict(),
213         "val_loss": val_loss,
214     }, path)
215     print(f" [checkpoint] saved to {path}")
216
217
218 def load_checkpoint(path, model, ema, optimizer, device):
219     ckpt = torch.load(path, map_location=device, weights_only=False)
220     unwrap(model).load_state_dict(ckpt["model_state_dict"])
221     ema.load_state_dict(ckpt["ema_state_dict"])
222     optimizer.load_state_dict(ckpt["optimizer_state_dict"])
223     return ckpt["epoch"], ckpt.get("val_loss", float("inf"))
224
225
226 #
-----
227 # Training
228 #
-----
229
230 def train(args):
231     global WEIGHT_DIR
232     if args.weight_dir:
233         WEIGHT_DIR = Path(args.weight_dir)
234         WEIGHT_DIR.mkdir(parents=True, exist_ok=True)
235     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
236     print(f"[train_v2] device = {device}")
237     print(f"[train_v2] weight_dir = {WEIGHT_DIR}")
238
239     # ---- Data -----
240     train_ds = CATHFrameDataset(DATA_JSONL, SPLITS_JSON, split="train",
241                                max_len=MAX_LEN, augment=True)
242     val_ds    = CATHFrameDataset(DATA_JSONL, SPLITS_JSON, split="validation",
243                                max_len=MAX_LEN, augment=False)
244
245     train_sampler = LengthBucketSampler(train_ds, batch_size=BATCH_SIZE,
246                                       shuffle=True, drop_last=True)
247     train_loader  = DataLoader(train_ds, batch_sampler=train_sampler,
248                               num_workers=NUM_WORKERS, pin_memory=True)
249     val_loader    = DataLoader(val_ds, batch_size=BATCH_SIZE, shuffle=False,
250                               num_workers=NUM_WORKERS, pin_memory=True)
251
252     # ---- Model -----
253     model_kwargs = dict(max_len=MAX_LEN, distogram_bins=DISTOGRAM_BINS)
254     if args.n_layers is not None:
255         model_kwargs["n_layers"] = args.n_layers
256     if args.stop_grad:
257         model_kwargs["stop_grad_frames"] = True
258     model = FrameDiffModel(**model_kwargs)
259     print(f"[train_v2] Using FrameDiffModel from model_v2.py (n_layers={
260           model.n_layers}, stop_grad={getattr(model, 'stop_grad_frames', False)
261           })")
262
263     model = model.to(device)
264
265     # Multi-GPU with DataParallel
266     n_gpus = torch.cuda.device_count()
267     if n_gpus > 1:
268         model = nn.DataParallel(model)
269         print(f"[train_v2] DataParallel on {n_gpus} GPUs")

```

```

268 else:
269     print(f"[train_v2] Single GPU")
270
271 if hasattr(torch, "compile"):
272     try:
273         # model = torch.compile(model) # disabled for debugging
274         print("[train_v2] torch.compile DISABLED for debugging")
275     except Exception:
276         pass
277
278 optimizer = torch.optim.Adam(model.parameters(), lr=LR)
279 ema = EMA(_unwrap(model), decay=EMA_DECAY)
280
281 # ---- IGS03 cache (REQUIRED) -----
282 _log(" building IGS03 cache...")
283 igso3_cache = IGS03Cache()
284 _log(" IGS03 cache ready")
285
286 # ---- Resume -----
287 start_epoch = 0
288 best_val_loss = float("inf")
289 if args.resume and os.path.isfile(args.resume):
290     start_epoch, best_val_loss = load_checkpoint(
291         args.resume, model, ema, optimizer, device)
292     start_epoch += 1
293     print(f"[train_v2] Resumed from epoch {start_epoch}, best_val_loss={
294         best_val_loss:.6f}")
295
296 # ---- CSV logger -----
297 csv_path = WEIGHT_DIR / "train_losses.csv"
298 csv_existed = csv_path.exists() and start_epoch > 0
299 csv_file = open(csv_path, "a" if csv_existed else "w", newline="")
300 csv_writer = csv.writer(csv_file)
301 if not csv_existed:
302     csv_writer.writerow([
303         "epoch", "L_x", "L_r", "L_bb", "L_2d", "total_loss",
304         "rot_angle_err_deg", "ca_rmsd_A",
305         "val_L_x", "val_L_r", "val_total",
306     ])
307     csv_file.flush()
308
309 # ---- Training loop -----
310 patience_counter = 0
311 self_cond_prev = None
312 global_step = 0
313
314 for epoch in range(start_epoch, MAX_EPOCHS):
315     model.train()
316     epoch_metrics = {
317         "L_x": 0.0, "L_r": 0.0, "L_bb": 0.0, "L_2d": 0.0,
318         "total": 0.0, "rot_err": 0.0, "rmsd": 0.0, "n_batches": 0,
319     }
320     t_start = time.time()
321
322     step_in_epoch = 0
323     for batch in train_loader:
324         step_in_epoch += 1
325         r0 = batch["rotations"].to(device) # (B, N, 3, 3)
326         x0 = batch["translations"].to(device) # (B, N, 3)
327         psi = batch["psi_angles"].to(device) # (B, N)
328         mask = batch["mask"].to(device) # (B, N)
329         B, N = x0.shape[:2]

```

```

330     # Sample timestep
331     t = torch.rand(B, device=device) * (T_MAX - T_MIN) + T_MIN
332
333     # Noise translations
334     x_t = noise_translations(x0, t, mask)
335
336     # Noise rotations using IGS03
337     #  $\sigma(t) = \sigma_{\min} + t * (\sigma_{\max} - \sigma_{\min})$ 
338     sigma = SIGMA_MIN + t * (SIGMA_MAX - SIGMA_MIN) # (B,)
339     t_igso3 = sigma ** 2 # parameter for IGS03 density:  $\exp(-l(l+1) * t/2)$ 
340
341     # Sample rotation noise per-residue
342     r_noise = torch.zeros(B, N, 3, 3, device=device, dtype=x0.dtype)
343     for b in range(B):
344         r_noise_b = igso3_cache.sample(t_igso3[b].item(), N, device=
345             device, dtype=x0.dtype)
346         r_noise[b] = r_noise_b # (N, 3, 3)
347     r_t = torch.matmul(r0, r_noise) # compose clean @ noise
348
349     # Self-conditioning: 50% of steps use previous prediction
350     if random() < SELF_COND_PROB and self_cond_prev is not None:
351         # Pad/truncate to match current batch dims
352         sc = self_cond_prev
353         if sc.shape[0] != B or sc.shape[1] != N:
354             distogram = torch.zeros(B, N, N, DISTOGRAM_BINS,
355                 device=device)
356         else:
357             distogram = compute_distogram(sc)
358     else:
359         distogram = torch.zeros(B, N, N, DISTOGRAM_BINS, device=
360             device)
361
362     # Forward pass
363     optimizer.zero_grad(set_to_none=True)
364     r_pred, x_pred, psi_pred = model(r_t, x_t, t, distogram)
365
366     # Debug: check for NaN in outputs on first batch of first epoch
367     if epoch == start_epoch and step_in_epoch < 3:
368         _log(f" DEBUG step {step_in_epoch}: x_pred [{x_pred.min().
369             item():.4f}], {x_pred.max().item():.4f}] "
370             f"r_det {torch.det(r_pred[0,0]).item():.4f} "
371             f"x0 [{x0.min().item():.4f}], {x0.max().item():.4f}]")
372
373     # Losses
374     L_x = translation_loss(x_pred, x0, mask)
375     L_r = rotation_loss_v2(r_pred, r0, r_t, t, mask, igso3_cache)
376     loss = L_x + L_r
377
378     # Auxiliary losses for low noise (t < 0.25)
379     t_low_mask = (t < AUX_THRESHOLD).any()
380     L_bb_val = 0.0
381     L_2d_val = 0.0
382     if t_low_mask:
383         # Select samples with t < threshold
384         low_idx = (t < AUX_THRESHOLD)
385         if low_idx.any():
386             L_bb = backbone_atom_loss(
387                 r_pred[low_idx], x_pred[low_idx], psi_pred[low_idx],
388                 r0[low_idx], x0[low_idx], psi[low_idx],
389                 mask[low_idx],
390             )
391             L_2d = pairwise_distance_loss(

```

```

389         x_pred[low_idx], x0[low_idx], mask[low_idx],
390     )
391     loss = loss + AUX_WEIGHT * (L_bb + L_2d)
392     L_bb_val = L_bb.item()
393     L_2d_val = L_2d.item()
394
395     if not torch.isfinite(loss):
396         _log(f" [WARN] NaN/Inf loss, skipping (L_x={L_x.item()},
397             L_r={L_r.item()})")
398         optimizer.zero_grad(set_to_none=True)
399         continue
400
401     loss.backward()
402
403     # Gradient monitoring for the first epoch
404     if epoch == start_epoch and step_in_epoch <= 5:
405         for name, p in model.named_parameters():
406             if p.grad is not None:
407                 if not torch.isfinite(p.grad).all():
408                     _log(f" NaN grad in {name}")
409                 grad_n = p.grad.norm().item()
410                 if grad_n > 100:
411                     _log(f" Large grad in {name}: {grad_n:.1f}")
412
413     # Check for NaN gradients
414     grad_norm = clip_grad_norm_(model.parameters(), GRAD_CLIP)
415     if not torch.isfinite(grad_norm):
416         # Diagnostic: which parameters have NaN gradients?
417         if step_in_epoch <= 3:
418             nan_params = []
419             for name, p in model.named_parameters():
420                 if p.grad is not None and not torch.isfinite(p.grad)
421                     .all():
422                     nan_params.append(name)
423             _log(f" [WARN] NaN gradients in: {nan_params[:10]}")
424         else:
425             _log(f" [WARN] NaN gradients, skipping optimizer step")
426             optimizer.zero_grad(set_to_none=True)
427             continue
428
429     # LR schedule: warmup then cosine decay
430     global_step += 1
431     total_steps = MAX_EPOCHS * len(train_loader)
432     if global_step <= WARMUP_STEPS:
433         cur_lr = LR * global_step / WARMUP_STEPS
434     else:
435         progress = (global_step - WARMUP_STEPS) / max(total_steps -
436             WARMUP_STEPS, 1)
437         cur_lr = LR_MIN + 0.5 * (LR - LR_MIN) * (1.0 + math.cos(math
438             .pi * progress))
439     for pg in optimizer.param_groups:
440         pg['lr'] = cur_lr
441
442     optimizer.step()
443     ema.update(_unwrap(model))
444
445     # Update self-conditioning cache
446     with torch.no_grad():
447         self_cond_prev = x_pred.detach()
448
449     # Metrics
450     with torch.no_grad():
451         rot_err = rotation_angle_error(r_pred, r0, mask)
452         rmsd = ca_rmsd(x_pred, x0, mask)

```

```

448
449     epoch_metrics["L_x"]         += L_x.item()
450     epoch_metrics["L_r"]         += L_r.item()
451     epoch_metrics["L_bb"]        += L_bb_val
452     epoch_metrics["L_2d"]        += L_2d_val
453     epoch_metrics["total"]       += loss.item()
454     epoch_metrics["rot_err"]     += rot_err
455     epoch_metrics["rmsd"]        += rmsd
456     epoch_metrics["n_batches"]   += 1
457
458     if _PREEMPTED:
459         break
460
461     nb = max(epoch_metrics["n_batches"], 1)
462     elapsed = time.time() - t_start
463
464     # ---- Validation -----
465     model.eval()
466     val_metrics = {"L_x": 0.0, "L_r": 0.0, "total": 0.0, "n": 0}
467     with torch.no_grad():
468         for batch in val_loader:
469             r0 = batch["rotations"].to(device)
470             x0 = batch["translations"].to(device)
471             psi = batch["psi_angles"].to(device)
472             mask = batch["mask"].to(device)
473             B, N = x0.shape[:2]
474
475             t = torch.rand(B, device=device) * (T_MAX - T_MIN) + T_MIN
476             x_t = noise_translations(x0, t, mask)
477
478             # IGS03 rotation noise for validation
479             sigma = SIGMA_MIN + t * (SIGMA_MAX - SIGMA_MIN)
480             t_igso3 = sigma ** 2
481             r_noise = torch.zeros(B, N, 3, 3, device=device, dtype=x0.
482                                 dtype)
483             for b in range(B):
484                 r_noise_b = igso3_cache.sample(t_igso3[b].item(), N,
485                                                device=device, dtype=x0.dtype)
486                 r_noise[b] = r_noise_b
487             r_t = torch.matmul(r0, r_noise)
488
489             distogram = torch.zeros(B, N, N, DISTOGRAM_BINS, device=
490                                 device)
491
492             r_pred, x_pred, psi_pred = model(r_t, x_t, t, distogram)
493             vL_x = translation_loss(x_pred, x0, mask)
494             vL_r = rotation_loss_v2(r_pred, r0, r_t, t, mask,
495                                   igso3_cache)
496             val_metrics["L_x"]     += vL_x.item()
497             val_metrics["L_r"]     += vL_r.item()
498             val_metrics["total"]  += (vL_x + vL_r).item()
499             val_metrics["n"]       += 1
500
501     vn = max(val_metrics["n"], 1)
502     val_total = val_metrics["total"] / vn
503
504     # ---- Logging -----
505     print(
506         f"Epoch {epoch:03d} ({elapsed:.0f}s) "
507         f"L_x={epoch_metrics['L_x']/nb:.5f} "
508         f"L_r={epoch_metrics['L_r']/nb:.5f} "
509         f"L_bb={epoch_metrics['L_bb']/nb:.5f} "
510         f"L_2d={epoch_metrics['L_2d']/nb:.5f} "

```

```

507         f"total={epoch_metrics['total']/nb:.5f}  "
508         f"rot_err={epoch_metrics['rot_err']/nb:.1f}deg  "
509         f"RMSD={epoch_metrics['rmsd']/nb:.2f}A  "
510         f"| val_L_x={val_metrics['L_x']/vn:.5f}  "
511         f"val_L_r={val_metrics['L_r']/vn:.5f}  "
512         f"val_total={val_total:.5f}"
513     )
514
515     csv_writer.writerow([
516         epoch,
517         f"{epoch_metrics['L_x']/nb:.6f}",
518         f"{epoch_metrics['L_r']/nb:.6f}",
519         f"{epoch_metrics['L_bb']/nb:.6f}",
520         f"{epoch_metrics['L_2d']/nb:.6f}",
521         f"{epoch_metrics['total']/nb:.6f}",
522         f"{epoch_metrics['rot_err']/nb:.2f}",
523         f"{epoch_metrics['rmsd']/nb:.3f}",
524         f"{val_metrics['L_x']/vn:.6f}",
525         f"{val_metrics['L_r']/vn:.6f}",
526         f"{val_total:.6f}",
527     ])
528     csv_file.flush()
529
530     # ---- Checkpointing -----
531     is_best = val_total < best_val_loss
532     if is_best:
533         best_val_loss = val_total
534         patience_counter = 0
535         save_checkpoint(model, ema, optimizer, epoch, val_total,
536                         WEIGHT_DIR / "best.pt")
537     else:
538         patience_counter += 1
539
540     # Always save latest
541     save_checkpoint(model, ema, optimizer, epoch, val_total,
542                     WEIGHT_DIR / "latest.pt")
543
544     # ---- Early stopping / preemption -----
545     if patience_counter >= PATIENCE:
546         print(f"[train_v2] Early stopping at epoch {epoch} "
547               f"(no improvement for {PATIENCE} epochs)")
548         break
549
550     if _PREEMPTED:
551         print("[train_v2] Preempted -- checkpoint saved, exiting.")
552         # Resubmit SLURM job
553         job_script = str(SCRIPYPT_DIR.parent / "slurm" / "train_v2.sh")
554         if os.path.isfile(job_script):
555             os.system(f"sbatch {job_script}")
556             print(f"[train_v2] Resubmitted {job_script}")
557         break
558
559     csv_file.close()
560     print(f"[train_v2] Done. Best val loss = {best_val_loss:.6f}")
561     print(f"[train_v2] Logs -> {csv_path}")
562     print(f"[train_v2] Weights-> {WEIGHT_DIR}")
563
564
565 #
566 # Entry point

```

```

567 #
-----
568 if __name__ == "__main__":
569     parser = argparse.ArgumentParser(description="FrameDiff v2 training")
570     parser.add_argument("--resume", type=str, default=None,
571                         help="Path to checkpoint to resume from")
572     parser.add_argument("--n_layers", type=int, default=None,
573                         help="Number of model layers (default: model default
574                             )")
575     parser.add_argument("--weight_dir", type=str, default=None,
576                         help="Override weight directory")
577     parser.add_argument("--stop_grad", action="store_true",
578                         help="Enable stop-gradient on intermediate frames")
579     args = parser.parse_args()
580     train(args)

```

Listing 3: Training Script (train_v2.py)

A.4 Sampling Script (sample_v2.py)

```
1 """
2 FrameDiff v2 sampling: generate protein backbones via reverse SE(3)
   diffusion.
3
4 Changes from sample.py (v1):
5 - Geodesic reverse SDE on SO(3) (replaces naive matrix interpolation)
6 - IGS03 for initial rotation noise (replaces QR decomposition)
7 - Proper timestep-to-sigma schedule
8 - Import from model_v2 and so3_diffusion
9
10 Usage:
11     python sample_v2.py --weights PATH --n_samples 10 --length 100 --out_dir
      ./generated
12 """
13
14 import os, sys, argparse, math, time
15 from pathlib import Path
16
17 import torch
18 import torch.nn.functional as F
19 import numpy as np
20
21 SCRIPT_DIR = Path(__file__).resolve().parent
22 sys.path.insert(0, str(SCRIPT_DIR))
23
24 from model_v2 import FrameDiffModel, compute_distogram
25 from losses import frames_to_atoms
26 from so3_diffusion import (IGSO3Cache, axis_angle_to_matrix,
27                             matrix_to_axis_angle,
28                             so3_exp, alpha_t, sigma_t)
29
30 #
31 -----
32 # IGS03 noise schedule constants
33 #
34 -----
35
36 SIGMA_MIN = 0.1
37 SIGMA_MAX = 1.5
38
39 #
40 -----
41
42 # SO(3) reverse diffusion step
43 #
44 -----
45
46 def reverse_step_rotations(R_t, R_pred, sigma_now, sigma_next, igso3_cache):
47     """One reverse diffusion step on SO(3) using geodesic interpolation.
48
49     Parameters
50     -----
51     R_t          : (B_N, 3, 3)  current noisy rotations (flattened batch*N)
52     R_pred       : (B_N, 3, 3)  model-predicted clean rotations
53     sigma_now    : float, current noise level
54     sigma_next   : float, next noise level (smaller)
```

```

50     igso3_cache : IGSO3Cache instance
51
52     Returns
53     -----
54     R_new : (B_N, 3, 3) updated rotations
55     """
56     B_N = R_t.shape[0]
57
58     # Relative rotation: how far is R_t from R_pred?
59     R_rel = R_t.transpose(-2, -1) @ R_pred # rotation from R_t to R_pred
60     axis, omega = matrix_to_axis_angle(R_rel) # axis (B_N, 3), omega (B_N,)
61
62     # How much to move toward R_pred at this step
63     # Posterior scaling: move proportionally to noise reduction
64     if sigma_now > 1e-6:
65         frac = 1.0 - (sigma_next / sigma_now) ** 2 # fraction of noise to
66             remove
67     else:
68         frac = 1.0
69     omega_step = omega * frac
70
71     # Apply geodesic step toward R_pred
72     R_step = axis_angle_to_matrix(axis, omega_step)
73     R_mean = R_t @ R_step
74
75     # Add reverse-step noise (IGSO3)
76     if sigma_next > 0.01: # don't add noise at final step
77         noise_sigma = sigma_next * math.sqrt(1.0 - (sigma_next / max(
78             sigma_now, 1e-6)) ** 2)
79         t_noise = noise_sigma ** 2
80         t_noise = max(t_noise, 0.001) # clamp for cache range
81         R_noise = igso3_cache.sample(t_noise, B_N, device=R_t.device, dtype=
82             R_t.dtype)
83         R_new = R_mean @ R_noise
84     else:
85         R_new = R_mean
86
87     # Gram-Schmidt re-orthogonalize
88     c0 = R_new[..., :, 0]
89     c1 = R_new[..., :, 1]
90     e0 = F.normalize(c0, dim=-1)
91     e1 = c1 - (c1 * e0).sum(-1, keepdim=True) * e0
92     e1 = F.normalize(e1, dim=-1)
93     e2 = torch.cross(e0, e1, dim=-1)
94     R_new = torch.stack([e0, e1, e2], dim=-1)
95
96     return R_new
97
98 # -----
99
100 # Reverse diffusion sampler
101 # -----
102
103 @torch.no_grad()
104 def sample_backbones(
105     model,
106     igso3_cache,
107     n_samples: int = 10,
108     length: int = 100,

```

```

106     n_steps: int = 200,
107     device: torch.device = torch.device("cpu"),
108     self_condition: bool = True,
109 ):
110     """Generate protein backbones via reverse SE(3) diffusion.
111
112     Uses DDPM-style  $x_0$ -prediction with proper IGS03 rotation diffusion.
113
114     Returns
115     -----
116     all_atoms : (n_samples, length, 4, 3) backbone atoms [N, CA, C, O] in nm
117     all_R      : (n_samples, length, 3, 3) final rotations
118     all_x      : (n_samples, length, 3) final Ca positions
119     all_psi    : (n_samples, length) final psi angles
120     """
121     model.eval()
122     dtype = next(model.parameters()).dtype
123
124     # Timestep schedule:  $T \rightarrow 0$ 
125     timesteps = torch.linspace(1.0, 0.0, n_steps + 1)
126     sigmas = SIGMA_MIN + timesteps * (SIGMA_MAX - SIGMA_MIN)
127
128     # Initialize from noise
129     # Translations: isotropic Gaussian, centered
130     x_t = torch.randn(n_samples, length, 3, device=device, dtype=dtype)
131     x_t = x_t - x_t.mean(dim=1, keepdim=True) # center
132
133     # Rotations: IGS03 at maximum noise
134     t_max_igso3 = SIGMA_MAX ** 2
135     R_t = igso3_cache.sample(t_max_igso3, n_samples * length, device=device,
136                             dtype=dtype)
137     R_t = R_t.view(n_samples, length, 3, 3)
138
139     # Self-conditioning cache
140     distogram = torch.zeros(n_samples, length, length, 32, device=device,
141                             dtype=dtype)
142
143     print(f"Sampling {n_samples} backbones of length {length} with {n_steps}
144           steps...")
145
146     for i in range(n_steps):
147         t_now = timesteps[i]
148         t_next = timesteps[i + 1]
149         sigma_now = sigmas[i].item()
150         sigma_next = sigmas[i + 1].item()
151         t_batch = t_now.expand(n_samples).to(device)
152
153         # Model prediction: clean frames
154         R_pred, x_pred, psi_pred = model(R_t, x_t, t_batch, distogram)
155
156         # Self-conditioning: update distogram from prediction
157         if self_condition:
158             distogram = compute_distogram(x_pred)
159
160         # --- Translation reverse step ---
161         # VP-SDE:  $x_t = \alpha_t * x_0 + \sigma_t * \epsilon$ 
162         alpha_now = torch.exp(-t_now / 2.0)
163         sigma_now_trans = torch.sqrt(1.0 - alpha_now ** 2)
164         alpha_next_val = torch.exp(-t_next / 2.0)
165         sigma_next_trans = torch.sqrt(1.0 - alpha_next_val ** 2)
166
167         if i < n_steps - 1: # don't add noise on last step
168             # DDPM posterior mean using  $x_0$  prediction

```

```

166         noise = torch.randn_like(x_t)
167         noise = noise - noise.mean(dim=1, keepdim=True) # center
168         x_t = alpha_next_val * x_pred + sigma_next_trans * noise
169     else:
170         x_t = x_pred
171
172     x_t = x_t - x_t.mean(dim=1, keepdim=True) # re-center
173
174     # --- Rotation reverse step (geodesic on SO(3)) ---
175     # Flatten to (n_samples*length, 3, 3) for the reverse step
176     R_t_flat = R_t.view(-1, 3, 3)
177     R_pred_flat = R_pred.view(-1, 3, 3)
178     R_t_flat = reverse_step_rotations(R_t_flat, R_pred_flat,
179                                     sigma_now, sigma_next,
180                                     igso3_cache)
181
182     R_t = R_t_flat.view(n_samples, length, 3, 3)
183
184     if (i + 1) % 50 == 0 or i == n_steps - 1:
185         print(f" Step {i+1}/{n_steps}: t={t_next.item():.4f}, "
186               f"sigma={sigma_next:.4f}")
187
188     # Reconstruct full backbone atoms
189     atoms = frames_to_atoms(R_t, x_t, psi_pred) # (n_samples, length, 4, 3)
190     return atoms, R_t, x_t, psi_pred
191 #
192 -----
193 # PDB writer
194 #
195 -----
196
197 def atoms_to_pdb(atoms: np.ndarray, path: str, chain: str = "A"):
198     """Write backbone atoms (N, CA, C, O) to PDB file.
199
200     atoms : (L, 4, 3) in nanometers
201     """
202     atom_names = ["N", "CA", "C", "O"]
203     elements = ["N", "C", "C", "O"]
204
205     with open(path, "w") as f:
206         atom_idx = 1
207         for res_idx in range(atoms.shape[0]):
208             for a_idx, (aname, elem) in enumerate(zip(atom_names, elements)):
209                 :
210                 x, y, z = atoms[res_idx, a_idx] * 10.0 # nm -> Angstroms
211                 f.write(
212                     f"ATOM {atom_idx:5d} {aname:<3s} GLY {chain}{res_idx
213                       +1:4d}      "
214                     f"{x:8.3f}{y:8.3f}{z:8.3f}  1.00  0.00           {elem
215                       :>2s}\n"
216                 )
217                 atom_idx += 1
218             f.write("END\n")
219 #
220 -----
221 # Main

```

```

218 #
-----
219
220 def main():
221     parser = argparse.ArgumentParser(description="FrameDiff v2 backbone
222         sampling")
223     parser.add_argument("--weights", type=str,
224         default="/global/scratch/users/sergiomar10/
225             latentlabs/framediff/weights_v2/best.pt",
226         help="Path to model checkpoint")
227     parser.add_argument("--n_samples", type=int, default=8, help="Number of
228         samples")
229     parser.add_argument("--length", type=int, default=100, help="Chain
230         length")
231     parser.add_argument("--n_steps", type=int, default=200, help="Diffusion
232         steps")
233     parser.add_argument("--out_dir", type=str,
234         default="/global/scratch/users/sergiomar10/
235             latentlabs/framediff/generated_v2",
236         help="Output directory")
237     parser.add_argument("--device", type=str, default="cuda" if torch.cuda.
238         is_available() else "cpu")
239     args = parser.parse_args()
240
241     device = torch.device(args.device)
242     out_dir = Path(args.out_dir)
243     out_dir.mkdir(parents=True, exist_ok=True)
244
245     # Load model
246     print(f>Loading model from {args.weights}")
247     model = FrameDiffModel(max_len=512, distogram_bins=32).to(device)
248     ckpt = torch.load(args.weights, map_location=device, weights_only=False)
249
250     # Load EMA weights if available
251     if "ema_state_dict" in ckpt:
252         print("Using EMA weights")
253         ema_sd = ckpt["ema_state_dict"]
254         model_sd = model.state_dict()
255         for k in model_sd:
256             if k in ema_sd:
257                 model_sd[k] = ema_sd[k]
258         model.load_state_dict(model_sd)
259     else:
260         model.load_state_dict(ckpt["model_state_dict"])
261
262     epoch = ckpt.get("epoch", "?")
263     val_loss = ckpt.get("val_loss", "?")
264     print(f>Loaded checkpoint: epoch={epoch}, val_loss={val_loss}")
265
266     # Build IGS03 cache
267     print("Building IGS03 cache...")
268     igso3_cache = IGS03Cache()
269     print("IGS03 cache ready")
270
271     # Sample
272     t0 = time.time()
273     atoms, R, x, psi = sample_backbones(
274         model,
275         igso3_cache,
276         n_samples=args.n_samples,
277         length=args.length,
278         n_steps=args.n_steps,

```

```

272     device=device,
273 )
274 elapsed = time.time() - t0
275 print(f"Sampling took {elapsed:.1f}s")
276
277 # Save PDBs
278 atoms_np = atoms.cpu().numpy()
279 for i in range(args.n_samples):
280     pdb_path = out_dir / f"sample_{i:03d}_L{args.length}.pdb"
281     atoms_to_pdb(atoms_np[i], str(pdb_path))
282     print(f"    Saved {pdb_path}")
283
284 # Basic quality stats
285 print("\n--- Quick quality check ---")
286 for i in range(min(3, args.n_samples)):
287     ca = atoms_np[i, :, 1, :] * 10.0 # CA positions in Angstroms
288     # Ca-Ca bond lengths
289     dists = np.sqrt(((ca[1:] - ca[:-1]) ** 2).sum(-1))
290     print(f"    Sample {i}: Ca-Ca distances: mean={dists.mean():.2f}A, "
291           f"    std={dists.std():.2f}A, min={dists.min():.2f}A, max={dists.
292           f"    max():.2f}A "
293           f"    (ideal: 3.8A)")
294     # Radius of gyration
295     com = ca.mean(axis=0)
296     rg = np.sqrt(((ca - com) ** 2).sum(-1).mean())
297     print(f"    Rg={rg:.1f}A")
298
299 print(f"\nAll {args.n_samples} PDBs saved to {out_dir}")
300
301 if __name__ == "__main__":
302     main()

```

Listing 4: Sampling Script (sample_v2.py)

A.5 SO(3) Diffusion (so3_diffusion.py)

```
1 """
2 SO(3) rotation diffusion for FrameDiff-style protein backbone generation.
3
4 Implements the Isotropic Gaussian on SO(3) (IGSO3) distribution, its score
5 function, sampling via CDF inversion, and the associated rotation utilities
6 (Rodrigues, exp/log maps). Also includes VP-SDE translation diffusion.
7
8 All public functions operate on PyTorch tensors. The IGSO3Cache class
9 pre-computes grids with numpy/scipy and exposes fast torch-based lookups.
10
11 Reference
12 -----
13 Yim et al., "SE(3) Diffusion Model with Application to Protein Backbone
14 Generation", ICML 2023.
15 """
16
17 from __future__ import annotations
18
19 import ast
20 import math
21 from typing import Optional, Tuple
22
23 import numpy as np
24 import torch
25 from scipy.interpolate import RegularGridInterpolator
26 from scipy.integrate import trapezoid
27
28 #
29 -----
29 # Numerical helpers
30 #
31 -----
31
32 _EPS = 1e-7 # guard for division-by-zero
33
34
35 def _sinc_half(omega: np.ndarray) -> np.ndarray:
36     """sin(omega/2) / (omega/2), safe at omega=0."""
37     h = omega / 2.0
38     out = np.ones_like(h)
39     mask = np.abs(h) > _EPS
40     out[mask] = np.sin(h[mask]) / h[mask]
41     return out
42
43
44 #
45 -----
45 # 1. IGSO3 Density
46 #
47 -----
47
48 def _igso3_density_np(
49     omega: np.ndarray,
50     t: np.ndarray,
51     L_max: int = 1000,
52 ) -> np.ndarray:
```

```

53     """
54     Evaluate the IGS03 density on a meshgrid (omega, t) using numpy float64.
55
56     f(omega, t) = sum_{l=0}^{L_max} (2l+1) exp(-l(l+1) t/2)
57                 * sin((l+0.5) omega) / sin(omega/2)
58
59     Parameters
60     -----
61     omega : (N,) array, angles in [0, pi].
62     t      : (M,) array, diffusion times > 0.
63
64     Returns
65     -----
66     density : (N, M) array.
67     """
68     omega = np.asarray(omega, dtype=np.float64)
69     t = np.asarray(t, dtype=np.float64)
70
71     # Shape: (N, 1) and (1, M) for broadcasting
72     om = omega[:, None]
73     tt = t[None, :]
74
75     density = np.zeros((omega.shape[0], t.shape[0]), dtype=np.float64)
76
77     sin_half = np.sin(om / 2.0) # (N,1)
78     # Near omega=0 or omega=pi, sin(omega/2) -> 0; handle via Taylor.
79     safe_denom = np.where(np.abs(sin_half) > _EPS, sin_half, np.ones_like(
80         sin_half))
81
82     for ell in range(L_max + 1):
83         coeff = (2 * ell + 1)
84         exp_part = np.exp(-ell * (ell + 1) * tt / 2.0) # (1, M)
85
86         # sin((ell+0.5)*omega) with Taylor near omega=0
87         arg = (ell + 0.5) * om # (N,1)
88         sin_num = np.sin(arg)
89
90         # ratio = sin((l+0.5)*omega) / sin(omega/2)
91         # At omega->0: Taylor gives (2l+1) (leading order)
92         # At omega->pi: sin(omega/2)->1, so no singularity in denominator,
93         # but sin((l+0.5)*pi) = cos(l*pi)*sin(pi/2)*... we just compute
94         # directly
95         # since sin(pi/2)=1 is fine.
96         near_zero = np.abs(om) < _EPS
97         ratio = np.where(
98             near_zero,
99             (2 * ell + 1) * np.ones_like(om), # L'Hopital: limit is (2l+1)
100             sin_num / safe_denom,
101         )
102
103         density += coeff * exp_part * ratio # (N, M)
104
105         # Early termination: if the exponential is negligibly small for all
106         # t
107         # remaining terms are even smaller because ell grows.
108         if ell > 0 and np.all(exp_part < 1e-30):
109             break
110
111     return density
112
113 def igso3_density(
114     omega: torch.Tensor,

```

```

113     t: torch.Tensor,
114     L_max: int = 1000,
115 ) -> torch.Tensor:
116     """
117     IGS03 density f(omega, t).
118
119     Parameters
120     -----
121     omega : (*,) tensor of rotation angles in [0, pi].
122     t      : (*,) tensor of diffusion times (same shape or broadcastable).
123
124     Returns
125     -----
126     density : (*,) tensor.
127     """
128     shape = torch.broadcast_shapes(omega.shape, t.shape)
129     om_flat = omega.broadcast_to(shape).reshape(-1).double().cpu().numpy()
130     t_flat = t.broadcast_to(shape).reshape(-1).double().cpu().numpy()
131
132     # Evaluate one-by-one through vectorised 1D call (each pair
133     # independently)
134     unique_t = np.unique(t_flat)
135     density_out = np.empty_like(om_flat)
136
137     # Group by unique t for efficiency
138     for tv in unique_t:
139         mask = t_flat == tv
140         om_sub = om_flat[mask]
141         d = _igso3_density_np(om_sub, np.array([tv]), L_max=L_max)
142         density_out[mask] = d[:, 0]
143
144     return torch.tensor(density_out, dtype=omega.dtype, device=omega.device)
145     .reshape(shape)
146 #
147 # -----
148 #
149 # 2. IGS03 Score (d log f / d omega)
150 # -----
151
152 def _igso3_density_and_deriv_np(
153     omega: np.ndarray,
154     t: np.ndarray,
155     L_max: int = 1000,
156 ) -> Tuple[np.ndarray, np.ndarray]:
157     """
158     Compute f(omega,t) and df/domega on a meshgrid.
159
160     The derivative of the l-th term w.r.t. omega:
161     d/domega [ sin((l+0.5)*omega) / sin(omega/2) ]
162     = [(l+0.5)*cos((l+0.5)*omega)*sin(omega/2)
163        - sin((l+0.5)*omega)*0.5*cos(omega/2)] / sin^2(omega/2)
164
165     Returns (N,M) arrays for density and its omega-derivative.
166     """
167     omega = np.asarray(omega, dtype=np.float64)
168     t = np.asarray(t, dtype=np.float64)
169     om = omega[:, None]
170     tt = t[None, :]

```

```

170 density = np.zeros((len(omega), len(t)), dtype=np.float64)
171 deriv = np.zeros_like(density)
172
173 sin_half = np.sin(om / 2.0)
174 cos_half = np.cos(om / 2.0)
175 sin_half_sq = sin_half ** 2
176 near_zero = np.abs(om) < _EPS
177 # Near pi: sin(omega/2) ~ 1, cos(omega/2) ~ 0, well-behaved.
178 safe_sin_half = np.where(np.abs(sin_half) > _EPS, sin_half, np.ones_like
    (sin_half))
179 safe_sin_half_sq = np.where(np.abs(sin_half_sq) > _EPS, sin_half_sq, np.
    ones_like(sin_half_sq))
180
181 for ell in range(L_max + 1):
182     coeff = 2 * ell + 1
183     half_ell = ell + 0.5
184     exp_part = np.exp(-ell * (ell + 1) * tt / 2.0)
185
186     sin_num = np.sin(half_ell * om)
187     cos_num = np.cos(half_ell * om)
188
189     # --- f term ---
190     ratio = np.where(near_zero, float(coeff), sin_num / safe_sin_half)
191     density += coeff * exp_part * ratio
192
193     # --- df/domega term ---
194     # d/domega [sin(a*om)/sin(om/2)]
195     # = [a*cos(a*om)*sin(om/2) - sin(a*om)*cos(om/2)/2] / sin^2(om/2)
196     # At omega->0 Taylor:
197     # ratio ~ (2l+1) + higher order => derivative at 0:
198     # Expanding to first order: deriv ~ -(2l+1)(l)(l+1)/6 * omega (->
    0)
199     numer = half_ell * cos_num * sin_half - sin_num * cos_half * 0.5
200     d_ratio = np.where(near_zero, 0.0, numer / safe_sin_half_sq)
201     deriv += coeff * exp_part * d_ratio
202
203     if ell > 0 and np.all(exp_part < 1e-30):
204         break
205
206 return density, deriv
207
208
209 def igso3_score(
210     omega: torch.Tensor,
211     t: torch.Tensor,
212     L_max: int = 1000,
213 ) -> torch.Tensor:
214     """
215     IGS03 score: d log f / d omega.
216
217     Parameters
218     -----
219     omega : (*,) rotation angles in [0, pi].
220     t      : (*,) diffusion times.
221
222     Returns
223     -----
224     score : (*,) tensor, same shape.
225     """
226     shape = torch.broadcast_shapes(omega.shape, t.shape)
227     om_flat = omega.broadcast_to(shape).reshape(-1).double().cpu().numpy()
228     t_flat = t.broadcast_to(shape).reshape(-1).double().cpu().numpy()
229

```

```

230 score_out = np.empty_like(om_flat)
231 for tv in np.unique(t_flat):
232     mask = t_flat == tv
233     om_sub = om_flat[mask]
234     f, df = _igso3_density_and_deriv_np(om_sub, np.array([tv]), L_max=
        L_max)
235     safe_f = np.where(np.abs(f[:, 0]) > 1e-30, f[:, 0], 1e-30)
236     score_out[mask] = df[:, 0] / safe_f
237
238 return torch.tensor(score_out, dtype=omega.dtype, device=omega.device).
    reshape(shape)
239
240
241 #
-----
242 # 3. Score Normalization Factor
243 #
-----
244
245 def score_norm(
246     t: float | np.ndarray,
247     n_omega: int = 2000,
248     L_max: int = 1000,
249 ) -> np.ndarray:
250     r"""
251     Compute  $1 / E[||\text{score}||^2]$  for loss weighting.
252
253     
$$E[||\text{score}||^2] = \int_0^\pi (d \log f / d \omega)^2 f(\omega, t) \sin^2(\omega / 2) d \omega$$

254
255     
$$/ \int_0^\pi f(\omega, t) \sin^2(\omega / 2) d \omega$$

256
257     Parameters
258     -----
259     t : scalar or 1-D array of diffusion times.
260
261     Returns
262     -----
263     lam : same shape as t, the normalisation  $\lambda_r(t)$ .
264     """
265     t = np.atleast_1d(np.asarray(t, dtype=np.float64))
266     omega_grid = np.linspace(1e-6, np.pi - 1e-6, n_omega)
267     f, df = _igso3_density_and_deriv_np(omega_grid, t, L_max=L_max) # (N, M)
268     sin2_half = (np.sin(omega_grid / 2.0) ** 2)[:, None] # (N, 1)
269
270     # d log f / d omega
271     safe_f = np.where(np.abs(f) > 1e-30, f, 1e-30)
272     score_vals = df / safe_f # (N, M)
273
274     integrand_num = score_vals ** 2 * f * sin2_half
275     integrand_den = f * sin2_half
276
277     e_score_sq = trapezoid(integrand_num, omega_grid, axis=0) / \
278         np.maximum(trapezoid(integrand_den, omega_grid, axis=0), 1e
        -30)
279
280     lam = 1.0 / np.maximum(e_score_sq, 1e-30)
281     return lam.squeeze()
282
283

```

```

284 #
-----
285 # 5. Rodrigues Formula & axis-angle <-> matrix
286 #
-----

287
288 def _skew_symmetric(v: torch.Tensor) -> torch.Tensor:
289     """
290     Skew-symmetric matrix from vector.  v: (... , 3) -> (... , 3, 3).
291     """
292     z = torch.zeros_like(v[... , 0])
293     K = torch.stack([
294         torch.stack([z, -v[... , 2], v[... , 1]], dim=-1),
295         torch.stack([v[... , 2], z, -v[... , 0]], dim=-1),
296         torch.stack([-v[... , 1], v[... , 0], z], dim=-1),
297     ], dim=-2)
298     return K
299
300
301 def axis_angle_to_matrix(axis: torch.Tensor, angle: torch.Tensor) -> torch.
Tensor:
302     """
303     Rodrigues formula.
304
305     Parameters
306     -----
307     axis : (... , 3) unit vectors.
308     angle : (... ,) rotation angles in radians.
309
310     Returns
311     -----
312     R : (... , 3, 3) rotation matrices.
313     """
314     K = _skew_symmetric(axis) # (... , 3, 3)
315     s = torch.sin(angle)[... , None, None]
316     c = torch.cos(angle)[... , None, None]
317     I = torch.eye(3, dtype=axis.dtype, device=axis.device).expand_as(K)
318     R = I + s * K + (1.0 - c) * (K @ K)
319     return R
320
321
322 def matrix_to_axis_angle(R: torch.Tensor) -> Tuple[torch.Tensor, torch.
Tensor]:
323     """
324     Extract axis and angle from a rotation matrix.
325
326     Parameters
327     -----
328     R : (... , 3, 3) rotation matrices.
329
330     Returns
331     -----
332     axis : (... , 3) unit axis (arbitrary when angle ~ 0).
333     angle : (... ,) rotation angle in [0, pi].
334     """
335     # angle from trace
336     trace = R[... , 0, 0] + R[... , 1, 1] + R[... , 2, 2]
337     cos_angle = (trace - 1.0) / 2.0
338     cos_angle = cos_angle.clamp(-1.0, 1.0)
339     angle = torch.acos(cos_angle) # [0, pi]
340

```

```

341     # axis from skew part of R
342     rx = R[..., 2, 1] - R[..., 1, 2]
343     ry = R[..., 0, 2] - R[..., 2, 0]
344     rz = R[..., 1, 0] - R[..., 0, 1]
345     axis = torch.stack([rx, ry, rz], dim=-1)
346
347     # normalise; when angle ~ 0 the axis is arbitrary
348     norm = axis.norm(dim=-1, keepdim=True).clamp(min=_EPS)
349     axis = axis / norm
350
351     # When angle ~ pi, the skew part is near zero; recover axis from
352     # the symmetric part: R + I = 2 * cos(0) * outer(axis, axis) (not
353     # needed
354     # for typical diffusion angles, but handle for robustness).
355     near_pi = (angle > math.pi - 0.01).unsqueeze(-1)
356     if near_pi.any():
357         # Pick the column of (R + I) with the largest diagonal entry
358         Rpi = R + torch.eye(3, dtype=R.dtype, device=R.device)
359         diag = torch.diagonal(Rpi, dim1=-2, dim2=-1) # (... , 3)
360         idx = diag.argmax(dim=-1) # (... ,)
361         # gather the column
362         idx_exp = idx.unsqueeze(-1).unsqueeze(-1).expand(*idx.shape, 3, 1)
363         col = torch.gather(Rpi, -1, idx_exp).squeeze(-1) # (... , 3)
364         col = col / col.norm(dim=-1, keepdim=True).clamp(min=_EPS)
365         axis = torch.where(near_pi, col, axis)
366
367     return axis, angle
368
369 #
370 -----
371 # 6. SO(3) Exponential and Log Maps
372 #
373 -----
374
375 def so3_exp(w: torch.Tensor) -> torch.Tensor:
376     """
377     Exponential map: tangent vector (or skew matrix) -> rotation matrix.
378
379     Parameters
380     -----
381     w : (... , 3) axis-angle vector **or** (... , 3, 3) skew-symmetric
382         matrix.
383
384     Returns
385     -----
386     R : (... , 3, 3).
387     """
388     if w.dim() >= 2 and w.shape[-1] == 3 and w.shape[-2] == 3:
389         # Extract axis-angle from skew matrix
390         w = torch.stack([w[..., 2, 1], w[..., 0, 2], w[..., 1, 0]], dim=-1)
391
392     angle = w.norm(dim=-1).clamp(min=_EPS)
393     axis = w / angle.unsqueeze(-1)
394     return axis_angle_to_matrix(axis, angle)
395
396 def so3_log(R: torch.Tensor) -> torch.Tensor:
397     """
398     Logarithmic map: rotation matrix -> tangent vector (axis * angle).
399

```

```

398     Parameters
399     -----
400     R : (... , 3, 3).
401
402     Returns
403     -----
404     w : (... , 3) axis-angle vector.
405     """
406     axis, angle = matrix_to_axis_angle(R)
407     return axis * angle.unsqueeze(-1)
408
409
410 #
-----
411 # 4. IGS03 Sampling (uses Rodrigues from above)
412 #
-----
413
414 def _sample_axis_uniform(n: int, device: torch.device, dtype: torch.dtype)
-> torch.Tensor:
415     """Sample n unit vectors uniformly on S^2."""
416     v = torch.randn(n, 3, device=device, dtype=dtype)
417     return v / v.norm(dim=-1, keepdim=True)
418
419
420 def sample_igso3(
421     t: float,
422     n_samples: int,
423     omega_grid: Optional[np.ndarray] = None,
424     cdf: Optional[np.ndarray] = None,
425     L_max: int = 1000,
426     device: torch.device = torch.device("cpu"),
427     dtype: torch.dtype = torch.float32,
428 ) -> torch.Tensor:
429     """
430     Sample rotation matrices from IGS03(t).
431
432     1. Build / re-use CDF of omega under the IGS03 measure.
433     2. Invert CDF via interpolation to get omega samples.
434     3. Sample random axes on S^2.
435     4. Construct rotation via Rodrigues.
436
437     Parameters
438     -----
439     t           : diffusion time (scalar).
440     n_samples   : number of rotations to draw.
441     omega_grid, cdf : optional pre-computed grid (from IGS03Cache).
442     L_max       : truncation for the density series.
443
444     Returns
445     -----
446     R : (n_samples, 3, 3) rotation matrices.
447     """
448     # --- CDF ---
449     if omega_grid is None or cdf is None:
450         n_grid = 2000
451         omega_grid = np.linspace(1e-6, np.pi - 1e-6, n_grid)
452         f = _igso3_density_np(omega_grid, np.array([t]), L_max=L_max)[: , 0]
453         sin2_half = np.sin(omega_grid / 2.0) ** 2
454         pdf = f * sin2_half
455         pdf = np.maximum(pdf, 0.0)

```

```

456         cdf = np.cumsum(pdf)
457         cdf = cdf / cdf[-1] # normalise
458
459     # Invert CDF by interpolation
460     u = np.random.rand(n_samples)
461     omega_samples = np.interp(u, cdf, omega_grid)
462
463     omega_t = torch.tensor(omega_samples, dtype=dtype, device=device)
464     axes = _sample_axis_uniform(n_samples, device, dtype)
465     R = axis_angle_to_matrix(axes, omega_t)
466     return R
467
468
469 #
-----
470 # 7. Rotation Score from Predictions
471 #
-----
472
473 def compute_rotation_score(
474     r_pred: torch.Tensor,
475     r_noisy: torch.Tensor,
476     t: torch.Tensor,
477 ) -> torch.Tensor:
478     """
479     Compute the SO(3) score given predicted clean rotation and noised
480     rotation.
481
482     score = r_noisy @ (axis * score_magnitude)      (tangent vector at
483     r_noisy)
484
485     where r_rel = r_pred^T @ r_noisy,  omega = angle(r_rel),
486           axis = axis(r_rel),  score_magnitude = d log f / d omega.
487
488     Parameters
489     -----
490     r_pred : (... , 3, 3) predicted clean rotations.
491     r_noisy : (... , 3, 3) noised rotations at time t.
492     t       : (...,) or scalar diffusion time.
493
494     Returns
495     -----
496     score : (... , 3) tangent vector at r_noisy (body-frame).
497     """
498     r_rel = r_pred.transpose(-2, -1) @ r_noisy # (... , 3, 3)
499     axis, omega = matrix_to_axis_angle(r_rel)
500
501     t_b = t.broadcast_to(omega.shape) if isinstance(t, torch.Tensor) else \
502         torch.full_like(omega, float(t))
503     score_mag = igso3_score(omega, t_b) # (...,)
504
505     # Tangent vector in body frame at identity, then transport to r_noisy
506     tangent_id = axis * score_mag.unsqueeze(-1) # (... , 3)
507     return tangent_id
508 #
-----
509 # 8. IGS03 Cache

```

```

510 #
-----
511
512 class IGS03Cache:
513     """
514     Pre-compute IGS03 quantities on a 2-D (omega x t) grid and provide fast
515     interpolated look-ups.
516
517     Attributes
518     -----
519     omega_grid : (N,) numpy array, angles in (0, pi).
520     t_grid     : (M,) numpy array, times in (t_min, t_max).
521     density    : (N, M) pre-computed f(omega, t).
522     deriv     : (N, M) df/domega.
523     cdf       : (N, M) CDF(omega | t) for each column.
524     score_norms: (M,) lambda_r(t).
525     """
526
527     def __init__(
528         self,
529         n_omega: int = 1000,
530         n_t: int = 1000,
531         t_min: float = 0.001,
532         t_max: float = 5.0,
533         L_max: int = 1000,
534     ):
535         self.L_max = L_max
536         self.omega_grid = np.linspace(1e-6, np.pi - 1e-6, n_omega).astype(np
                    .float64)
537         self.t_grid = np.linspace(t_min, t_max, n_t).astype(np.float64)
538
539         # --- density & derivative ---
540         self.density, self.deriv = _igso3_density_and_deriv_np(
541             self.omega_grid, self.t_grid, L_max=L_max,
542         )
543
544         # --- score = d log f / d omega ---
545         safe_f = np.where(np.abs(self.density) > 1e-30, self.density, 1e-30)
546         self._score_grid = self.deriv / safe_f # (N, M)
547
548         # --- CDF for each t column ---
549         sin2_half = (np.sin(self.omega_grid / 2.0) ** 2)[: , None]
550         pdf = self.density * sin2_half # (N, M)
551         pdf = np.maximum(pdf, 0.0)
552         self.cdf = np.cumsum(pdf, axis=0)
553         # normalise each column
554         cdf_max = self.cdf[-1: , :]
555         cdf_max = np.where(cdf_max > 0, cdf_max, 1.0)
556         self.cdf = self.cdf / cdf_max
557
558         # --- score norms ---
559         self.score_norms = score_norm(self.t_grid, n_omega=n_omega, L_max=
                    L_max)
560
561         # --- scipy interpolators ---
562         self._density_interp = RegularGridInterpolator(
563             (self.omega_grid, self.t_grid), self.density,
564             method="linear", bounds_error=False, fill_value=None,
565         )
566         self._score_interp = RegularGridInterpolator(
567             (self.omega_grid, self.t_grid), self._score_grid,
568             method="linear", bounds_error=False, fill_value=None,

```

```

569     )
570     self._score_norm_t = self.t_grid.copy()
571     self._score_norm_vals = np.atleast_1d(self.score_norms).copy()
572
573     # ---- fast torch look-ups ----
574
575     def density_lookup(self, omega: torch.Tensor, t: torch.Tensor) -> torch.
Tensor:
576         """Interpolated density f(omega, t)."""
577         pts = np.stack([
578             omega.double().cpu().numpy().ravel(),
579             t.double().cpu().numpy().ravel(),
580         ], axis=-1)
581         vals = self._density_interp(pts)
582         return torch.tensor(vals, dtype=omega.dtype, device=omega.device).
reshape(omega.shape)
583
584     def score_lookup(self, omega: torch.Tensor, t: torch.Tensor) -> torch.
Tensor:
585         """Interpolated score d log f / d omega."""
586         pts = np.stack([
587             omega.double().cpu().numpy().ravel(),
588             t.double().cpu().numpy().ravel(),
589         ], axis=-1)
590         vals = self._score_interp(pts)
591         return torch.tensor(vals, dtype=omega.dtype, device=omega.device).
reshape(omega.shape)
592
593     def score_norm_lookup(self, t: torch.Tensor) -> torch.Tensor:
594         """Interpolated lambda_r(t)."""
595         t_np = t.double().cpu().numpy().ravel()
596         vals = np.interp(t_np, self._score_norm_t, self._score_norm_vals)
597         return torch.tensor(vals, dtype=t.dtype, device=t.device).reshape(t.
shape)
598
599     def sample(
600         self,
601         t: float,
602         n_samples: int,
603         device: torch.device = torch.device("cpu"),
604         dtype: torch.dtype = torch.float32,
605     ) -> torch.Tensor:
606         """
607         Sample rotation matrices from IGS03(t) using the cached CDF.
608
609         Returns (n_samples, 3, 3).
610         """
611         # Find the nearest t column for the CDF
612         idx = np.searchsorted(self.t_grid, t).clip(0, len(self.t_grid) - 1)
613         cdf_col = self.cdf[:, idx]
614         return sample_igso3(
615             t, n_samples,
616             omega_grid=self.omega_grid, cdf=cdf_col,
617             L_max=self.L_max, device=device, dtype=dtype,
618         )
619
620
621 #
-----
622 # 9. Translation Diffusion (VP-SDE, simple Gaussian)
623 #
-----

```

```

624
625 def alpha_t(t: torch.Tensor) -> torch.Tensor:
626     """VP-SDE scaling factor  $\alpha(t) = \exp(-t/2)$ ."""
627     return torch.exp(-t / 2.0)
628
629
630 def sigma_t(t: torch.Tensor) -> torch.Tensor:
631     """VP-SDE noise level  $\sigma(t) = \sqrt{1 - \alpha(t)^2}$ ."""
632     a = alpha_t(t)
633     return torch.sqrt((1.0 - a * a).clamp(min=0.0))
634
635
636 def _centering_projection(n: int, device: torch.device, dtype: torch.dtype)
637     -> torch.Tensor:
638     """
639     Centering projection  $P = I - (1/n) 11^T$ . Shape (n, n).
640     Removes the mean so translations live in the zero-COM subspace.
641     """
642     I = torch.eye(n, device=device, dtype=dtype)
643     ones = torch.ones(n, 1, device=device, dtype=dtype)
644     return I - ones @ ones.T / n
645
646 def sample_translation_noise(
647     x0: torch.Tensor,
648     t: torch.Tensor,
649     center: bool = True,
650 ) -> torch.Tensor:
651     """
652     Forward VP-SDE:  $x_t = \alpha(t) * x_0 + \sigma(t) * \text{eps}$ ,
653     optionally centered (projected to zero COM).
654
655     Parameters
656     -----
657     x0 : (N, 3) clean positions.
658     t  : scalar or (N,) times.
659
660     Returns
661     -----
662     x_t : (N, 3).
663     """
664     if isinstance(t, (int, float)):
665         t = torch.tensor(t, dtype=x0.dtype, device=x0.device)
666     t = t.reshape(-1, 1) if t.dim() == 1 else t.unsqueeze(0) # (N,1) or
667         (1,1)
668
669     a = alpha_t(t) # (*, 1)
670     s = sigma_t(t)
671
672     eps = torch.randn_like(x0)
673     if center:
674         P = _centering_projection(x0.shape[0], x0.device, x0.dtype)
675         eps = P @ eps # project noise to zero-COM
676
677     x_t = a * x0 + s * eps
678     if center:
679         x_t = P @ x_t
680     return x_t
681
682 def translation_score(
683     x0_pred: torch.Tensor,

```

```

684     x_t: torch.Tensor,
685     t: torch.Tensor,
686 ) -> torch.Tensor:
687     """
688     Translation score: nabla_{x_t} log p(x_t | x0).
689
690     = -(x_t - alpha(t) x0) / sigma(t)^2
691
692     Parameters
693     -----
694     x0_pred : (N, 3) predicted clean positions.
695     x_t      : (N, 3) noised positions.
696     t        : scalar or (N,) times.
697
698     Returns
699     -----
700     score : (N, 3).
701     """
702     if isinstance(t, (int, float)):
703         t = torch.tensor(t, dtype=x_t.dtype, device=x_t.device)
704     t = t.reshape(-1, 1) if t.dim() == 1 else t.unsqueeze(0)
705
706     a = alpha_t(t)
707     s2 = sigma_t(t) ** 2
708     s2 = s2.clamp(min=1e-10)
709     return -(x_t - a * x0_pred) / s2
710
711
712 #
713 # -----
714 # Self-test / syntax check
715 # -----
716
717 def _run_tests():
718     """Quick sanity checks (run with python so3_diffusion.py)."""
719     import sys
720     print("Running SO(3) diffusion self-tests..")
721
722     # --- Rodrigues roundtrip ---
723     axis = torch.tensor([0.0, 0.0, 1.0])
724     angle = torch.tensor(1.2)
725     R = axis_angle_to_matrix(axis, angle)
726     a2, ang2 = matrix_to_axis_angle(R)
727     assert abs(ang2.item() - 1.2) < 1e-5, f"Rodrigues roundtrip failed: {ang2}"
728
729     # --- exp / log roundtrip ---
730     w = torch.tensor([0.3, -0.5, 0.1])
731     R2 = so3_exp(w)
732     w2 = so3_log(R2)
733     assert torch.allclose(w, w2, atol=1e-5), f"exp/log roundtrip failed"
734
735     # --- Density at moderate t should be peaked near omega=0 ---
736     om = torch.tensor([0.01, 0.5, 1.5, 3.0])
737     t_mod = torch.tensor([0.5, 0.5, 0.5, 0.5])
738     d = igso3_density(om, t_mod)
739     assert d[0] > d[2] > d[3], f"Density not peaked at omega=0 for t=0.5: {d}"
740
741     # At very small t (0.01), density is essentially a delta at omega=0
742     # so values at omega>=0.5 underflow to ~0, which is correct behaviour.

```

```

741
742 # --- Score sign: should be negative for omega > 0 (pushes back to 0)
743     ---
744 s = igso3_score(torch.tensor([1.0]), torch.tensor([0.5]))
745 assert s.item() < 0, f"Score should be negative, got {s.item()}"
746
747 # --- Sampling: small t → near identity ---
748 R_samp = sample_igso3(0.01, 500)
749 angles = matrix_to_axis_angle(R_samp)[1]
750 assert angles.mean() < 0.5, f"Samples at small t too far from identity:
751     mean angle={angles.mean():.3f}"
752
753 # --- Sampling: large t → ~uniform, mean angle should be ~pi/2 ---
754 R_samp2 = sample_igso3(50.0, 2000)
755 angles2 = matrix_to_axis_angle(R_samp2)[1]
756 # Uniform on SO(3): E[angle] = pi - 1 ≈ 2.14
757 # (density of angle for uniform = sin^2(omega/2) / (pi) normalised)
758 assert abs(angles2.mean().item() - 2.14) < 0.3, \
759     f"Large-t samples not uniform: mean angle={angles2.mean():.3f},
760     expected ~2.14"
761
762 # --- Score norm should be O(1) ---
763 lam = score_norm(np.array([0.1, 1.0, 3.0]))
764 assert np.all(np.isfinite(lam)), f"Score norm not finite: {lam}"
765 assert np.all(lam > 0), f"Score norm not positive: {lam}"
766
767 # --- Translation VP-SDE ---
768 x0 = torch.randn(10, 3)
769 x0 = x0 - x0.mean(0)
770 xt = sample_translation_noise(x0, 0.5)
771 assert xt.shape == x0.shape
772 assert abs(xt.mean(0).sum().item()) < 0.1, "Translation noise not
773     centered"
774
775 ts = translation_score(x0, xt, 0.5)
776 assert ts.shape == x0.shape
777
778 # --- IGSO3Cache (small grid for speed) ---
779 cache = IGSO3Cache(n_omega=100, n_t=50, L_max=500)
780 d_cached = cache.density_lookup(torch.tensor([1.0]), torch.tensor([0.5])
781     )
782 assert d_cached.shape == (1,)
783 s_cached = cache.score_lookup(torch.tensor([1.0]), torch.tensor([0.5]))
784 assert s_cached.shape == (1,)
785 R_cached = cache.sample(0.5, 10)
786 assert R_cached.shape == (10, 3, 3)
787
788 print("All tests passed.")
789 return 0
790
791 if __name__ == "__main__":
792     # Syntax check
793     import pathlib
794     src = pathlib.Path(__file__).read_text()
795     ast.parse(src)
796     _run_tests()

```

Listing 5: SO(3) Diffusion (so3_diffusion.py)

A.6 Data Loading (data.py)

```
1 """
2 CATH dataset for FrameDiff: loads protein backbones, converts to SE(3)
  frames.
3
4 All coordinates are stored and returned in **nanometers**.
5 """
6
7 import json
8 import math
9 import torch
10 import numpy as np
11 from torch.utils.data import Dataset, Sampler
12 from typing import Dict, List, Optional, Tuple
13
14
15 #
16 -----
17 # Geometry helpers
18 #
19 -----
20
21 def _gram_schmidt(v1: torch.Tensor, v2: torch.Tensor) -> torch.Tensor:
22     """Gram-Schmidt orthogonalisation to build a rotation matrix from two
23     vectors.
24
25     Parameters
26     -----
27     v1 : (... , 3)    first vector (C - CA direction)
28     v2 : (... , 3)    second vector (N - CA direction)
29
30     Returns
31     -----
32     R : (... , 3, 3) rotation matrix with columns [e1, e2, e3]
33     """
34     e1 = v1 / (torch.linalg.norm(v1, dim=-1, keepdim=True) + 1e-8)
35     # Remove e1 component from v2
36     e2 = v2 - (v2 * e1).sum(-1, keepdim=True) * e1
37     e2 = e2 / (torch.linalg.norm(e2, dim=-1, keepdim=True) + 1e-8)
38     e3 = torch.cross(e1, e2, dim=-1)
39     # R has e1, e2, e3 as columns: R = [e1 | e2 | e3]
40     R = torch.stack([e1, e2, e3], dim=-1)           # (... , 3, 3)
41     return R
42
43 def _dihedral(a: torch.Tensor, b: torch.Tensor,
44              c: torch.Tensor, d: torch.Tensor) -> torch.Tensor:
45     """Compute dihedral angle for four points a-b-c-d.
46
47     Returns angle in radians, shape (...).
48     """
49     b1 = b - a
50     b2 = c - b
51     b3 = d - c
52     b2_norm = b2 / (torch.linalg.norm(b2, dim=-1, keepdim=True) + 1e-8)
53     n1 = torch.cross(b1, b2, dim=-1)
54     n2 = torch.cross(b2, b3, dim=-1)
55     m1 = torch.cross(n1, b2_norm, dim=-1)
56     x = (m1 * n2).sum(-1)
```

```

55     y = (m1 * n2).sum(-1)
56     return torch.atan2(y, x)
57
58
59 def _random_so3(shape: Tuple[int, ...], device="cpu", dtype=torch.float32)
-> torch.Tensor:
60     """Sample a uniform random SO(3) rotation matrix.
61
62     Uses the QR decomposition approach (Mezzadri 2007).
63     Returns (*shape, 3, 3).
64     """
65     n = 1
66     for s in shape:
67         n *= s
68     z = torch.randn(n, 3, 3, device=device, dtype=dtype)
69     q, r = torch.linalg.qr(z)
70     # Ensure proper rotation (det = +1)
71     sign = torch.sign(torch.diagonal(r, dim1=-2, dim2=-1))      # (n, 3)
72     q = q * sign[:, None, :]
73     det = torch.det(q)
74     q[det < 0] = -q[det < 0]
75     return q.reshape(*shape, 3, 3)
76
77
78 #
-----
79 # Dataset
80 #
-----
81
82 class CATHFrameDataset(Dataset):
83     """CATH protein backbone dataset returning SE(3) frames.
84
85     Each sample is a dictionary:
86     rotations      : (N, 3, 3) per-residue rotation matrices
87     translations  : (N, 3)    Calpha positions in nm (centered)
88     psi_angles    : (N,)      psi torsion angles in rad
89     mask          : (N,)      1.0 for valid residues, 0.0 for padding
90     length        : int       actual chain length (unpadded)
91     """
92
93     def __init__(
94         self,
95         jsonl_path: str,
96         splits_path: str,
97         split: str = "train",
98         max_len: int = 256,
99         augment: bool = True,
100     ):
101         super().__init__()
102         self.max_len = max_len
103         self.augment = augment
104
105         # Load split names
106         with open(splits_path) as f:
107             splits = json.load(f)
108             valid_names = set(splits[split])
109
110         # Load chains
111         self.entries: List[Dict] = []
112         with open(jsonl_path) as f:

```

```

113         for line in f:
114             rec = json.loads(line)
115             if rec["name"] not in valid_names:
116                 continue
117             L = len(rec["seq"])
118             if L > max_len or L < 10:
119                 continue
120             # Check for enough non-NaN coords
121             ca = np.array(rec["coords"]["CA"])
122             valid = ~np.isnan(ca).any(axis=1)
123             if valid.sum() < 10:
124                 continue
125             self.entries.append(rec)
126
127     print(f"[CATHFrameDataset] split={split} loaded {len(self.entries)}
128           chains "
129           f"(max_len={max_len})")
130
131     def __len__(self) -> int:
132         return len(self.entries)
133
134     def __getitem__(self, idx: int) -> Dict[str, torch.Tensor]:
135         rec = self.entries[idx]
136         L = len(rec["seq"])
137
138         # 1. Extract backbone atom coordinates (L, 3) in Angstroms
139         coords_N = torch.tensor(rec["coords"]["N"], dtype=torch.float32)
140         coords_CA = torch.tensor(rec["coords"]["CA"], dtype=torch.float32)
141         coords_C = torch.tensor(rec["coords"]["C"], dtype=torch.float32)
142         coords_O = torch.tensor(rec["coords"]["O"], dtype=torch.float32)
143
144         # Build validity mask: all four atoms must be non-NaN
145         valid = (
146             ~torch.isnan(coords_N).any(-1)
147             & ~torch.isnan(coords_CA).any(-1)
148             & ~torch.isnan(coords_C).any(-1)
149             & ~torch.isnan(coords_O).any(-1)
150         ) # (L,)
151
152         # Replace NaN with 0 for computation safety
153         coords_N = torch.nan_to_num(coords_N, 0.0)
154         coords_CA = torch.nan_to_num(coords_CA, 0.0)
155         coords_C = torch.nan_to_num(coords_C, 0.0)
156         coords_O = torch.nan_to_num(coords_O, 0.0)
157
158         # 2. Convert to nanometers
159         coords_N = coords_N / 10.0
160         coords_CA = coords_CA / 10.0
161         coords_C = coords_C / 10.0
162         coords_O = coords_O / 10.0
163
164         # 3. Center C-alpha to zero mean (only valid residues)
165         valid_f = valid.float()
166         center = (coords_CA * valid_f[:, None]).sum(0) / valid_f.sum().clamp(
167             min=1)
168         coords_N = coords_N - center
169         coords_CA = coords_CA - center
170         coords_C = coords_C - center
171         coords_O = coords_O - center
172
173         # 4. Gram-Schmidt frames from N, CA, C
174         v1 = coords_C - coords_CA # (L, 3)
175         v2 = coords_N - coords_CA # (L, 3)

```

```

174     rotations = _gram_schmidt(v1, v2) # (L, 3, 3)
175
176     # 5. Psi torsion angles: dihedral(N_i, CA_i, C_i, N_{i+1})
177     psi = torch.zeros(L)
178     if L > 1:
179         psi[:-1] = _dihedral(coords_N[:-1], coords_CA[:-1],
180                             coords_C[:-1], coords_N[1:])
181     # Last residue: no N_{i+1}, leave as 0
182
183     # 6. Data augmentation: random SO(3) rotation of ALL atoms
184     if self.augment:
185         R_aug = _random_so3((), device=coords_CA.device, dtype=coords_CA
186                             .dtype) # (3,3)
187         coords_CA = torch.einsum("ij, nj -> ni", R_aug, coords_CA)
188         coords_N = torch.einsum("ij, nj -> ni", R_aug, coords_N)
189         coords_C = torch.einsum("ij, nj -> ni", R_aug, coords_C)
190         coords_O = torch.einsum("ij, nj -> ni", R_aug, coords_O)
191         # Rotate frames: R_new = R_aug @ R_old
192         rotations = torch.einsum("ij, njk -> nik", R_aug, rotations)
193
194     translations = coords_CA # (L, 3)
195
196     # 7. Pad to max_len
197     pad = self.max_len - L
198     mask = torch.cat([valid_f, torch.zeros(pad)])
199
200     if pad > 0:
201         rotations = torch.cat([rotations, torch.eye(3).unsqueeze
202                                (0).expand(pad, 3, 3)], dim=0)
203         translations = torch.cat([translations, torch.zeros(pad, 3)],
204                                 dim=0)
205         psi = torch.cat([psi, torch.zeros(pad)], dim
206                          =0)
207
208     return {
209         "rotations": rotations, # (max_len, 3, 3)
210         "translations": translations, # (max_len, 3)
211         "psi_angles": psi, # (max_len,)
212         "mask": mask, # (max_len,)
213         "length": L,
214     }
215
216     # -----
217
218     # Length-bucketed sampler
219     # -----
220
221     class LengthBucketSampler(Sampler):
222         """Groups proteins by length for efficient batching.
223
224         Buckets: [10-50), [50-100), [100-150), [150-200), [200-256].
225         Within each bucket, indices are shuffled; buckets are iterated in
226         random order to mix the training signal.
227         """
228         def __init__(
229             self,
230             dataset: CATHFrameDataset,
231             batch_size: int = 8,

```

```

229     bucket_boundaries: Optional[List[int]] = None,
230     shuffle: bool = True,
231     drop_last: bool = False,
232 ):
233     self.dataset = dataset
234     self.batch_size = batch_size
235     self.shuffle = shuffle
236     self.drop_last = drop_last
237
238     if bucket_boundaries is None:
239         bucket_boundaries = [50, 100, 150, 200]
240     self.boundaries = bucket_boundaries
241
242     # Assign each sample to a bucket
243     self.buckets: List[List[int]] = [[] for _ in range(len(self.
        boundaries) + 1)]
244     for i, entry in enumerate(dataset.entries):
245         L = len(entry["seq"])
246         placed = False
247         for bi, b in enumerate(self.boundaries):
248             if L < b:
249                 self.buckets[bi].append(i)
250                 placed = True
251                 break
252         if not placed:
253             self.buckets[-1].append(i)
254
255     def __iter__(self):
256         batches = []
257         for bucket in self.buckets:
258             if len(bucket) == 0:
259                 continue
260             indices = list(bucket)
261             if self.shuffle:
262                 perm = torch.randperm(len(indices)).tolist()
263                 indices = [indices[p] for p in perm]
264             for i in range(0, len(indices), self.batch_size):
265                 batch = indices[i : i + self.batch_size]
266                 if self.drop_last and len(batch) < self.batch_size:
267                     continue
268                 batches.append(batch)
269
270     # Shuffle batch order
271     if self.shuffle:
272         perm = torch.randperm(len(batches)).tolist()
273         batches = [batches[p] for p in perm]
274
275     for batch in batches:
276         yield batch
277
278     def __len__(self) -> int:
279         total = 0
280         for bucket in self.buckets:
281             n = len(bucket)
282             if self.drop_last:
283                 total += n // self.batch_size
284             else:
285                 total += (n + self.batch_size - 1) // self.batch_size
286         return total

```

Listing 6: Data Loading (data.py)