
Flow Matching for Respiratory Motion Interpolation from Dual-Phase CT

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Accurate respiratory motion characterization is essential for thoracic radiotherapy, yet 4D-CT remains inaccessible to many centers due to infrastructure and
2 resource requirements, in addition to concerns of increased radiation dose exposure
3 compared to 3D-CT. This accessibility gap creates treatment disparities across
4 healthcare systems globally. To address this disparity, we present the first applica-
5 tion of flow matching to reconstruct respiratory dynamics from standard dual-phase
6 CT: routine images captured at maximum inhalation and exhalation. Our method
7 combines MedVAE compression with a teacher-forced Diffusion Transformer that
8 learns breathing patterns from 4D-CT data. Evaluating 587 sequences from 20
9 patients, reconstructions maintain anatomical consistency with the strongest gains
10 at mid-respiratory positions where uncertainty peaks. Performance gains of +2.9%
11 MS-SSIM and -6.3% LPIPS validate that generative motion synthesis from rou-
12 tine imaging is feasible. This foundation enables future extensions to volumetric
13 reconstruction and tumor tracking.
14

15 1 Introduction

16 1.1 Clinical Motivation and Challenges

17 Precision radiation therapy for thoracic and abdominal tumors requires accurate targeting while
18 minimizing exposure to healthy tissue. This requirement becomes particularly challenging due to
19 respiratory motion. Multi-institution audits report four-dimensional computed tomography (4D-CT)
20 reconstruction amplitude deviations of about 2 mm during regular breathing and up to 6 mm during
21 irregular patterns [1]. These motion uncertainties correlate with worse local control rates of lung and
22 liver metastases treated with SBRT [2].

23 4D-CT characterizes motion across respiratory phases and represents the current standard for tumor
24 motion visualization in radiation treatment planning. However, adoptability and practice patterns
25 vary widely across institutions. 4D-CT radiation dose exposure is often two to four times higher than
26 conventional CT [3], and typically requires longer acquisition times, motion monitoring hardware,
27 and trained staff [4]. Protocol choices also differ by scanner, reconstruction method, and phase
28 binning, which leads to heterogeneous motion representations [1]. These technical and resource
29 requirements create accessibility gaps—both across institutions within health systems and globally,
30 where many low- and middle-income countries lack access to advanced imagings [5–7]. Our goal is
31 to close this gap with methods that work under conventional CT acquisitions without 4D availability
32 while the field progresses toward more standardized protocols.

33 **1.2 Our Approach and Contributions**

34 To address these disparities, we ask whether it is possible to recover clinically reliable mid-cycle
35 motion from routine dual-phase CT without extra dose or workflow change. Our approach infers
36 intermediate phases from end-inspiration and end-expiration. This enables motion-aware planning
37 without 4D-CT and can lower imaging exposure and shorten simulation time. The benefit is greatest
38 for pediatric patients and for people with cancer predisposition syndromes such as Li–Fraumeni [8, 9].
39 To our knowledge, this is the first adaptation of flow matching [10, 11] for medical image interpo-
40 lation between sparse temporal observations. Predictions are anchored at the measured phases to
41 prevent drift and respect observed anatomy. This boundary-constrained approach improves temporal
42 consistency and strengthens training signals at mid-cycle where uncertainty is typically highest.
43 Our implementation combines MedVAE [12] for diagnostically-aware compression with a lightweight
44 Diffusion Transformer [13] optimized for temporal reconstruction under first-last frame conditioning.
45 The training process leverages ground truth phases (0%, 10%, 20%, 30%, 40%, 50%) with optical
46 flow interpolation [14] for dense temporal supervision. At inference, the model requires only the
47 endpoint phases (0% and 50%) to generate complete respiratory motion.
48 We validate our approach on 587 breathing sequences from the TCIA 4D-Lung dataset [15, 16]
49 and demonstrate improvements over standard interpolation baselines, particularly at challenging
50 mid-cycle phases where clinical uncertainty is highest. This work represents the foundational step in
51 a comprehensive motion modeling framework. We begin with slice-wise temporal reconstruction
52 to establish the feasibility of flow matching for respiratory motion. Future work will extend this
53 approach to volumetric reconstruction that incorporates information from adjacent slices and to target-
54 aware motion modeling that captures crano-caudal dynamics essential for treatment planning. By
55 establishing that accurate motion completion from universally available dual-phase CT is achievable,
56 we provide a practical pathway toward democratizing advanced motion management across diverse
57 clinical settings without requiring protocol changes or additional radiation exposure.

58 **2 Related Work**

59 **Respiratory motion modeling and DIR.** In radiotherapy, respiratory motion is often estimated
60 with deformable image registration (DIR). Classical variational optical flow (Horn–Schunck, TV-L1,
61 Farnebäck) and free-form B-spline registration remain strong baselines but need careful regularization
62 and can degrade under large deformations or in homogeneous regions [14, 17–19]. Learning-based
63 DIR, including LungRegNet and unsupervised variants, improves pairwise alignment on respiratory-
64 correlated datasets [20–22]. In routine planning only end-inspiration and end-expiration may be
65 available. Pairwise DIR does not directly produce temporally consistent mid-cycle frames from
66 dual-phase inputs. Chaining deformations to fill frames can introduce drift. Simple frame filling with
67 linear or cubic interpolation or optical flow is common but does not guarantee anatomically plausible
68 trajectories [14, 17, 23, 24].

69 **Flow matching and video generation.** Flow Matching (and rectified-flow variants) learns a
70 continuous-time velocity field or probability-flow ODE that transports one distribution to another, pro-
71 viding an efficient alternative to stochastic diffusion sampling [10, 11, 25]. For video, diffusion-style
72 models learn spatiotemporal dynamics directly (often in latent space) and scale to high resolution, e.g.,
73 Video Diffusion Models, Pyramidal Flow Matching, and Diffusion Transformers (DiT) [13, 26, 27].

74 **First-last-frame conditioning (FLF2V).** Conditioning on the first and last frames has emerged as a
75 practical interface for image-to-video generation. Large systems such as Meta MovieGen and Wan2.1-
76 FLF2V-14B demonstrate high-fidelity trajectories between two keyframes, while diffusion-based
77 academic in-betweening (e.g., VIDIM) explores sparse-anchor interpolation [28–30]. We adapt this
78 endpoint-aware design to medical CT, where “keyframes” are clinically acquired end-inspiratory and
79 end-expiratory phases.

80 **Domain-specific latent compression.** Operating in a compact latent space is crucial for high-
81 resolution medical imagery. Latent diffusion established the efficacy of autoencoder latents for
82 scalable generation, and medical foundation autoencoders such as MedVAE preserve diagnostically
83 relevant detail while enabling efficient training [12, 31].

84 **Supervising mid-trajectory dynamics.** To strengthen intermediate states, we draw on the spirit of
 85 teacher forcing and related sequence-training strategies (scheduled sampling, Professor Forcing),
 86 which provide ground-truth context to stabilize one-step predictions [32–34]. Conceptually related
 87 “consistency” objectives target one/few-step state agreement, complementary to our velocity-based
 88 teacher-forced supervision [35, 36]. For evaluation we use the TCIA 4D-Lung collection, a standard
 89 resource in radiotherapy motion studies [15, 16].

90 3 Method

91 3.1 Problem Formulation

92 Let $\mathcal{X} \subset \mathbb{R}^{H \times W}$ denote the space of 2D CT slices. We consider dual-phase inputs $\mathbf{x}_{\text{input}} =$
 93 $\{I^{(0)}, I^{(50)}\}$, where $I^{(0)} \in \mathcal{X}$ corresponds to end-inspiration (0% breathing phase) and $I^{(50)} \in \mathcal{X}$
 94 corresponds to the mid-cycle 50% phase (commonly corresponding to end-expiration) acquired in 4D-
 95 CT. Our objective is to generate the missing intermediate phases $\mathbf{I} = \{I^{(10)}, I^{(20)}, I^{(30)}, I^{(40)}\} \in \mathcal{X}^4$,
 96 providing a temporally complete motion sequence between the two observed anchors.

97 **Slice Selection Strategy.** Clinical CT acquisitions consist of hundreds of axial slices that form a
 98 complete 3D thoracic volume. Rather than processing the entire volume, which is currently com-
 99 putationally expensive for our configuration, we select three representative axial slices capturing
 100 respiratory motion across different lung regions. As shown in Figure 1, we first identify the range
 101 where lung tissue is visible, then select slices at relative positions: inferior (Slice 0), middle (Target
 102 Slice), and superior (Slice 60). In our experiments, these correspond to slices 12, 25, and 37 respec-
 103 tively within the lung-containing range. This slice-based approach allows us to focus computational
 104 resources on learning high-quality temporal dynamics while maintaining anatomical coverage from
 105 the lung base to the apex.

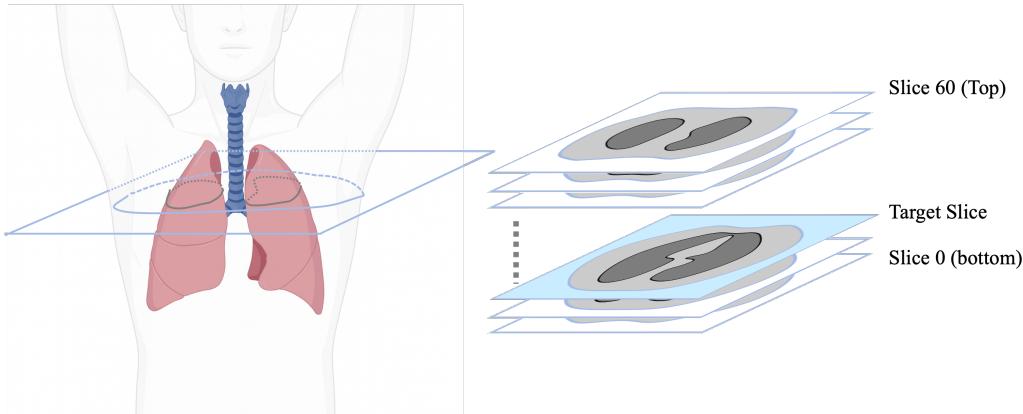


Figure 1: Anatomical slice selection strategy. From a complete thoracic CT volume, we identify slices containing visible lung tissue and select three representative levels: inferior (Slice 0), middle (Target Slice), and superior (Slice 60). This provides comprehensive coverage of respiratory motion patterns from the lung base to the apex.

106 3.2 Overall Architecture

107 Our approach combines three key components: (1) a medical VAE (MedVAE) for efficient latent
 108 space compression while preserving clinically relevant features, (2) a Diffusion Transformer (DiT)
 109 for learning respiratory motion dynamics, and (3) flow matching with teacher-forced supervision for
 110 generating smooth temporal interpolations.

111 3.2.1 Latent Space Compression with MedVAE

112 To address computational constraints while preserving medical image fidelity, we employ MedVAE
 113 with compression factor $f = 64$:

$$\mathbf{z}^{(i)} = E_\phi(I^{(i)}) \in \mathbb{R}^{H' \times W' \times C} \quad (1)$$

114 where E_ϕ is the MedVAE encoder with $H' = H/8$, $W' = W/8$. We extend the latent channels from
 115 1 to 8 through a learned projection to provide sufficient representational capacity for motion modeling
 116 while computing KL divergence on the original single-channel distribution.

117 3.2.2 Respiratory Motion Modeling with DiT

118 We model respiratory motion using a Diffusion Transformer adapted for temporal interpolation.
 119 The architecture processes 3D latent sequences $\mathbf{z} \in \mathbb{R}^{C \times D \times H' \times W'}$ where D represents the tem-
 120 poral dimension. We use a configuration with hidden dimension 512 and 24 transformer layers
 121 (approximately 140M parameters).

122 The DiT receives conditioning from both endpoint frames [28–30]:

$$\mathbf{v}_\theta(\mathbf{z}(t), t, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) = \text{DiT}(\mathbf{z}(t), \phi_{\text{time}}(t), \phi_{\text{flf}}(\mathbf{z}^{(0)}, \mathbf{z}^{(50)})) \quad (2)$$

123 where $\phi_{\text{time}}(t)$ provides sinusoidal time embeddings and ϕ_{flf} encodes first-last-frame conditioning.

124 3.2.3 Flow Matching with Adapted Supervision

125 We model respiratory motion as a continuous normalizing flow between end-inspiratory and end-
 126 expiratory latent states. Let $\mathbf{z}^{(0)}$ and $\mathbf{z}^{(50)}$ denote the VAE latents of the 0% (end-inspiration) and
 127 50% (end-expiration) phases. For curriculum and initialization we use the linear path

$$\mathbf{z}(t) = (1 - t) \mathbf{z}^{(0)} + t \mathbf{z}^{(50)}, \quad t \in [0, 1], \quad (3)$$

128 and inject a small Gaussian noise during training/inference for numerical stability. The velocity field
 129 $\mathbf{v}_\theta(\cdot, t, \mathbf{z}^{(0)}, \mathbf{z}^{(50)})$ parameterizes the ODE $\dot{\mathbf{z}}(t) = \mathbf{v}_\theta(\mathbf{z}(t), t, \mathbf{z}^{(0)}, \mathbf{z}^{(50)})$ integrated by an explicit
 130 solver.

131 3.2.4 Standard Flow Matching Loss:

132 The baseline objective minimizes:

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t \sim \mathcal{U}[0, 1]} \left[\left\| \mathbf{v}_\theta(\mathbf{z}(t), t, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) - \mathbf{v}_{\text{target}}(t) \right\|_2^2 \right] \quad (4)$$

133 While \mathcal{L}_{FM} fits the global field, it can be weak at mid-cycle times far from the two anchors (0%, 50%),
 134 where small local errors may accumulate during integration.

135 3.2.5 Teacher-forced supervision (local one-step consistency).

136 Given ground-truth latents $\{\mathbf{z}_{\text{gt}}^{(j)}\}_{j=0}^{D-1}$ at times $\{t_j\}$ with $\Delta t = t_{j+1} - t_j$, we randomly sample
 137 $j \in \{0, \dots, D-2\}$ and enforce a one-step prediction consistency:

$$\mathcal{L}_{\text{TF}} = \frac{1}{N} \sum_j \left\| \mathbf{z}_{\text{gt}}^{(j)} + \Delta t \cdot \mathbf{v}_\theta(\mathbf{z}_{\text{gt}}^{(j)}, t_j, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) - \mathbf{z}_{\text{gt}}^{(j+1)} \right\|_2^2 \quad (5)$$

138 Equivalently (up to a constant factor Δt^2):

$$\mathcal{L}_{\text{TF}} \propto \frac{1}{N} \sum_j \left\| \mathbf{v}_\theta(\mathbf{z}_{\text{gt}}^{(j)}, t_j, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) - \frac{\mathbf{z}_{\text{gt}}^{(j+1)} - \mathbf{z}_{\text{gt}}^{(j)}}{\Delta t} \right\|_2^2 \quad (6)$$

139 thus providing *pointwise* constraints on the field at sparse mid-cycle times where supervision is
 140 typically scarce. The finite difference $(\mathbf{z}_{\text{gt}}^{(j+1)} - \mathbf{z}_{\text{gt}}^{(j)})/\Delta t$ represents the true anatomical velocity
 141 between phases, capturing nonlinear tissue deformations. This direct supervision is strongest precisely
 142 where \mathcal{L}_{FM} 's gradient weakens (around $t \approx 0.5$), which is validated by our phase-specific results
 143 showing maximum improvement at 20–30% phases (Table 3).

144 **3.2.6 First–Last Frame (FLF) endpoint consistency.**

145 To keep the endpoints faithful, we add a small penalty that matches the predicted velocity to the target
 146 velocity *only at the first and last frames*:

$$\mathcal{L}_{\text{FLF}} = \frac{1}{2} \left\| \mathbf{v}_\theta(\mathbf{z}(0), 0, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) - \mathbf{v}_{\text{target}}(0) \right\|_2^2 + \frac{1}{2} \left\| \mathbf{v}_\theta(\mathbf{z}(1), 1, \mathbf{z}^{(0)}, \mathbf{z}^{(50)}) - \mathbf{v}_{\text{target}}(1) \right\|_2^2. \quad (7)$$

147 **Densified intermediate supervision.** As detailed in Section 4.1, we augment training sequences by
 148 applying optical flow interpolation between consecutive GT phases (0%→10%, 10%→20%, etc.)
 149 to create 41-frame sequences. These densified frames provide additional supervision for \mathcal{L}_{FM} and
 150 \mathcal{L}_{TF} , but critically differ from end-to-end optical flow (0%→50%) which would need to hallucinate
 151 all intermediate dynamics. Our approach leverages optical flow only for local temporal smoothness
 152 within 10% intervals while maintaining anatomical ground truth anchors at every respiratory phase,
 153 combining dense sampling benefits with guaranteed clinical fidelity. We emphasize that we never
 154 regress to optical flow *vectors*, and all quantitative evaluation uses exclusively ground-truth phases.

155 **3.2.7 VAE regularization.**

156 When the encoder/decoder are trained or fine-tuned jointly, we include a standard VAE objective to
 157 preserve anatomical fidelity in the latent space:

$$\mathcal{L}_{\text{VAE}} = \lambda_{\text{rec}} \mathbb{E} \left[\| \mathbf{x} - g_\psi(f_\phi(\mathbf{x})) \|_2^2 \right] + \beta \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| \mathcal{N}(\mathbf{0}, \mathbf{I})) . \quad (8)$$

158 with f_ϕ the encoder and q_ϕ the approximate posterior. If the VAE is frozen during flow training, we
 159 set $\mathcal{L}_{\text{VAE}} = 0$.

160 **3.2.8 Total objective.**

161 The overall training objective combines the flow-matching loss, the teacher-forced local consistency
 162 loss, the first–last frame endpoint penalty, and the VAE regularization term. Each component is
 163 weighted by a scalar coefficient that controls its relative contribution:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{FM}} \mathcal{L}_{\text{FM}} + \lambda_{\text{TF}} \mathcal{L}_{\text{TF}} + \lambda_{\text{FLF}} \mathcal{L}_{\text{FLF}} + \lambda_{\text{VAE}} \mathcal{L}_{\text{VAE}}. \quad (9)$$

164 **3.3 Training Strategy**

165 **3.3.1 Implementation Details**

166 We employ a two-stage training strategy: (1) joint VAE-DiT training for initial learning, then (2)
 167 frozen VAE with focused DiT training. Key hyperparameters include:

- 168 • Loss weights: $\lambda_{\text{FM}} = 1.0$, $\lambda_{\text{TF}} = 1.0$, $\lambda_{\text{FLF}} = 0.1$
- 169 • AdamW optimizer with learning rate 10^{-4} and cosine annealing
- 170 • Mixed precision training (FP16) with gradient clipping at norm 1.0
- 171 • Batch size 2-4 with gradient accumulation for effective batch size 8-16
- 172 • Training on phases 0-50% (full inspiration to full expiration) for consistency
- 173 • Intensity windowing (e.g., $[-600, 900]$ HU) for lung parenchyma; rescale to $[-1, 1]$

174 **3.4 Inference Procedure**

175 During inference, we generate intermediate frames through ODE solving with linear initialization
 176 (rather than pure noise) for stability. The initial latent state for integration is constructed as:

$$\mathbf{z}_{\text{init}} = (1 - \alpha) \mathbf{z}^{(0)} + \alpha \mathbf{z}^{(50)} + \epsilon \quad (10)$$

177 where $\alpha \in [0, 1]$ specifies the target temporal position, $\mathbf{z}^{(0)}$ and $\mathbf{z}^{(50)}$ are the encoded 0% and
 178 50% phase latents, and $\epsilon \sim \mathcal{N}(0, \sigma_{\min}^2 I)$ adds small Gaussian noise. With this parameterization,
 179 $\alpha = 0$ corresponds to initialization at the inhalation phase while $\alpha = 1$ corresponds to the 50%
 180 mid-cycle phase. We integrate the velocity field using Euler steps (50 iterations) and apply optional
 181 classifier-free guidance [37] (scale 1.5) to produce the full sequence of intermediate latents, which
 182 are then decoded by the VAE.

183 **4 Experimental Setup**

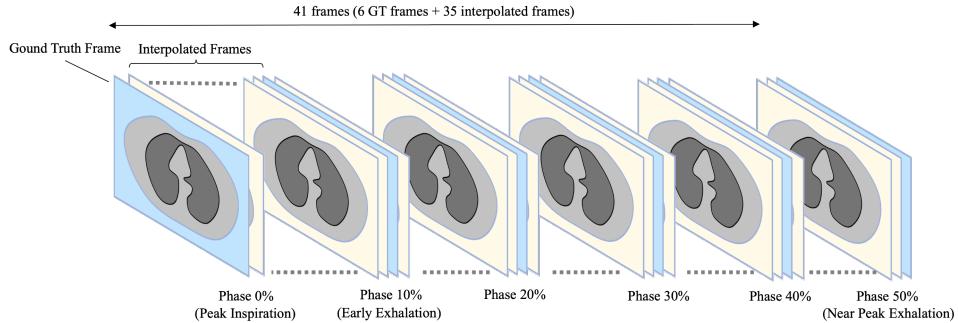
184 **4.1 Dataset and Preprocessing**

185 We evaluate on the 4D-Lung dataset containing 587 complete breathing cycles from 20 patients. Each
 186 cycle provides 10 respiratory phases (0%, 10%, ..., 90%) [15, 16]. We use patient-level splitting: 13
 187 patients for training (392 cycles), 3 for validation (96 cycles), and 4 for testing (99 cycles).

188 Preprocessing includes: (1) Hounsfield unit windowing $[-600, 900]$ for lung tissue, (2) spatial
 189 resampling to 128×128 for computational efficiency during training (512×512 for final generation),
 190 (3) extraction of three representative axial slices per volume, and (4) phase selection focusing on
 191 0-50% (end-inspiration to near end-expiration) for training stability.

192 **Training Data Construction.** As illustrated in Figure 2(a), we construct dense training sequences
 193 from the sparse ground truth phases. Starting with 6 GT phases from the 4D-Lung dataset (0%,
 194 10%, 20%, 30%, 40%, 50%), we apply optical flow interpolation between consecutive phases to
 195 generate intermediate frames. This creates 41-frame sequences combining original GT frames with
 196 35 interpolated frames. The optical flow provides smooth transitions between phases while the GT
 197 frames serve as anchor points for our teacher-forced supervision. During training, the model learns to
 198 predict velocity fields across this dense temporal sequence while being specifically constrained at GT
 199 frame positions.

(a) Training: Dense 41-Frame Sequence with 6 GT Phases



(b) Inference: First-Last Frame Conditioning



Figure 2: Training and inference data structure. (a) During training, we use dense 41-frame sequences with 6 ground truth phases from the 4D-Lung dataset, creating interpolated frames between available phases. (b) During inference, our method takes sparse dual-phase inputs (first-last frame conditioning) to generate intermediate respiratory phases.

200 At inference (Figure 2(b)), our method requires only the first and last frames (0% and 50%) as input.
 201 The model then generates all 39 intermediate frames through learned velocity field integration without
 202 requiring optical flow computation or access to other GT phases. For all reported metrics, supervision
 203 and evaluation use only ground-truth phases; interpolated frames are used to densify inputs but are
 204 not treated as targets.

205 **4.2 Evaluation Metrics**

206 We assess performance using standard image quality metrics and report mean \pm std across all
207 evaluated slices and phases:

- 208 • **MS-SSIM** (Multi-Scale Structural Similarity [38]): Measures structural preservation across
209 multiple scales.
- 210 • **PSNR** (Peak Signal-to-Noise Ratio): Quantifies pixel-level reconstruction accuracy.
- 211 • **LPIPS** (Learned Perceptual Image Patch Similarity): Evaluates perceptual quality in deep
212 feature space (AlexNet trunk).

213 We evaluate across three anatomical slice levels (inferior, middle, superior; slices 12, 25, 37) and four
214 intermediate breathing phases (10%, 20%, 30%, 40%) to assess both spatial and temporal robustness.
215 Error analysis includes signed error maps to identify systematic biases in intensity prediction.

216 **4.3 Baseline Comparisons**

217 We compare four established interpolation/registration baselines:

- 218 • **Linear interpolation**: Direct pixel-wise linear blending between frames [23, 24].
- 219 • **B-spline interpolation (Spline)**: Smooth cubic spline-based frame interpolation [24].
- 220 • **Optical flow**: Dense motion field estimation using Farnebäck’s algorithm [14].
- 221 • **Deformable image registration (DIR)**: 2D B-spline registration from 0% \rightarrow 50% (SimpleITK;
222 Mattes MI; coarse-to-fine), followed by backward warping. Intermediate phases
223 are generated by linearly scaling the displacement field (e.g., 20% phase uses 0.4 \times the full
224 deformation), which assumes linear motion between endpoints.

225 **5 Results**

226 **5.1 Quantitative Evaluation**

227 We evaluate four baselines (Linear, Spline, Optical Flow, DIR) across three anatomical slice levels (12,
228 25, 37) and four intermediate phases (10-40%), and summarize results in Table 1. Spline and Linear
229 are comparable overall, Optical Flow trails slightly, and DIR performs substantially lower, consistent
230 with a single endpoint deformation field failing to capture phase-dependent nonlinearity/hysteresis
231 and with 2D registration ignoring through-plane sliding interfaces.

Table 1: Quantitative comparison on the 4D-Lung test set (mean \pm std). Higher is better for MS-
SSIM/PSNR; lower is better for LPIPS.

Method	MS-SSIM \uparrow	PSNR (dB) \uparrow	LPIPS \downarrow
Linear	0.762 ± 0.135	25.31 ± 5.10	0.103 ± 0.035
Optical Flow	0.718 ± 0.127	23.58 ± 4.57	0.111 ± 0.038
Spline	0.767 ± 0.135	25.22 ± 5.10	0.095 ± 0.035
DIR	0.660 ± 0.161	19.85 ± 2.28	0.124 ± 0.044
Flow Matching (Ours)	0.789 ± 0.138	25.80 ± 5.12	0.089 ± 0.035

232 Compared to the strongest baseline in each metric, our method improves MS-SSIM by +0.022
233 ($\sim 2.9\%$ vs. Spline at 0.767), increases PSNR by +0.49 dB (25.80 vs. Linear at 25.31 dB), and
234 reduces LPIPS by $\sim 6.3\%$ (0.095 \rightarrow 0.089 vs. Spline).

235 **5.2 Per-Slice Analysis**

236 Performance is stable across anatomical levels, with Spline/Linear leading the classical baselines and
237 Optical Flow slightly trailing. DIR is consistently lower across slices. Our model retains its per-slice
238 advantage.

Table 2: Per-slice Performance (MS-SSIM / PSNR).

Method	Slice 12 (Base)	Slice 25 (Middle)	Slice 37 (Apex)
Linear	0.753 / 24.96	0.768 / 25.45	0.764 / 25.51
Optical Flow	0.709 / 23.37	0.721 / 23.67	0.724 / 23.69
Spline	0.758 / 24.87	0.773 / 25.37	0.770 / 25.43
DIR	0.647 / 19.42	0.667 / 20.04	0.666 / 20.09
Flow Matching (Ours)	0.764 / 24.99	0.791 / 26.02	0.811 / 26.38

239 5.3 Phase-Specific Analysis

240 All methods peak near the endpoints (10%, 40%) and degrade in the mid-cycle (20%, 30%), reflecting
 241 increased difficulty away from keyframes. DIR departs from this trend at 40%, consistent with the
 242 limitations discussed for endpoint-only deformation fields.

Table 3: Phase-specific PSNR (dB).

Method	Phase 10%	Phase 20%	Phase 30%	Phase 40%
Linear	26.32	24.28	24.29	26.34
Optical Flow	25.13	23.21	22.83	23.15
Spline	26.17	24.25	24.26	26.21
DIR	22.55	19.82	18.79	18.25
Flow Matching (Ours)	26.11	25.43	25.26	26.38

243 5.4 Qualitative Evaluation

244 Figure 3 presents visual comparisons of our flow matching interpolation across different breathing
 245 phases. We show ground truth, interpolated results, and signed error maps for phases at 10%, 20%,
 246 30%, and 40% of the respiratory cycle.

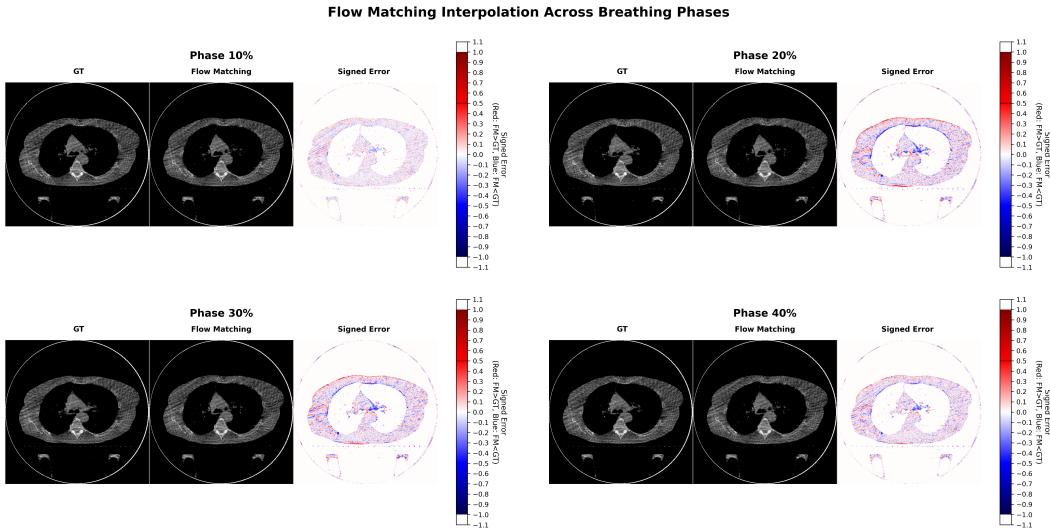


Figure 3: Flow matching interpolation across breathing phases. Ground truth (GT), flow matching results, and signed error maps (blue: under-estimation, red: over-estimation, scale: $[-1.1, 1.1]$) are shown for respiratory phases at 10%, 20%, 30%, and 40%. Anatomical structures remain well-preserved across all phases with balanced error distribution.

247 The signed error maps reveal balanced positive and negative deviations without systematic bias,
 248 indicating accurate intensity mapping between phases. Error magnitude follows the expected
 249 interpolation pattern for first-last frame conditioning, with minimal errors at 10% and 40% (close to

250 input frames at 0% and 50%) and maximum errors at 20%-30% phases where the model must infer
251 intermediate states furthest from both conditioning frames. Notably, errors concentrate primarily at
252 tissue boundaries and regions of respiratory motion, while homogeneous regions maintain minimal
253 deviation. The preservation of anatomical structures across all phases, particularly at the challenging
254 20%-30% midpoints where uncertainty is highest, demonstrates the method’s robustness and
255 clinical viability. These visual results corroborate our quantitative findings, confirming that flow
256 matching provides superior temporal interpolation while maintaining diagnostic image quality even
257 at maximum interpolation distance.

258 **5.5 Discussion**

259 Our results demonstrate successful adaptation of flow matching to medical respiratory motion
260 modeling. While improvements over traditional interpolation are modest ($\sim 2.9\%$ MS-SSIM, $\sim 6\%$
261 LPIPS, and $+0.49$ dB PSNR), they are consistent across all phases and most pronounced at challenging
262 mid-cycle phases (20%-30%) where accurate motion characterization is clinically critical. Balanced
263 errors and preserved structures support clinical plausibility; full clinical validation is future work.

264 The effectiveness stems from key domain adaptations: MedVAE preserves diagnostic features during
265 compression, teacher-forced supervision anchors predictions to ground truth anatomy, and training
266 on respiratory-specific phase ranges ensures physiologically realistic motion. These adaptations
267 demonstrate that video generation techniques developed for natural images can effectively handle the
268 unique requirements of medical imaging when properly modified.

269 **5.6 Limitations and Future Work**

270 This study is slice-wise and therefore lacks explicit volumetric coupling, which limits fidelity
271 for crano-caudal (through-plane) motion and sliding interfaces. The evaluation also focuses on
272 regularized breathing patterns; irregular cycles and pronounced hysteresis are underexplored, and
273 we did not assess downstream contour propagation or dose impact. Future work will extend the
274 model to volumetric latent flows that incorporate information from neighboring slices to capture 3D
275 consistency, and will further integrate target-aware cues (e.g., tumor or organ contours) to enable
276 trajectory estimation and contour propagation for radiotherapy planning.

277 **6 Conclusion**

278 We present a novel application of flow matching to respiratory motion synthesis in medical imaging.
279 By combining flow matching with domain-specific adaptations including medical-aware compression
280 and teacher-forced supervision, we generate anatomically plausible intermediate respiratory phases
281 from limited images at select phases. Our approach demonstrates that modern video generation
282 techniques can address clinical imaging challenges when thoughtfully adapted to preserve diagnostic
283 quality and respect physiological constraints. This work provides a foundation for applying generative
284 modeling to medical motion, with potential impact on radiation therapy planning and delivery.

285 **Ethics and Trust Statement**

286 This work uses a public, de-identified dataset (TCIA 4D-Lung). The method is intended as an
287 assistive tool; it does not replace clinical judgment. To mitigate risks from implausible motion or
288 artifacts, we (i) constrain supervision to ground-truth phases, (ii) report signed error maps, and (iii)
289 recommend clinician review and uncertainty visualization before downstream use. No protected
290 health information was used.

291 **Reproducibility Statement**

292 We provide implementation details (architectures, losses, hyperparameters), patient-level splits, and
293 evaluation protocols in the main text.

294 **References**

- 295 [1] J. Dhont M. Kyndt A. Gulyban J. Szkitsak E. Bogaert D. van Gestel M. Burghelea, J.
296 Bakkali Tahiri and N. Reynaert. Results of a multicenter 4d computed tomography qual-
297 ity assurance audit: Evaluating image accuracy and consistency. *Physics and Imaging in
298 Radiation Oncology*, 2023.
- 299 [2] A.K. Ozga C. Petersen L. Pinnschmidt R. Werner T. Sentker, V. Schmidt and T. Gauer. 4d
300 ct image artifacts affect local control in sbrt of lung and liver metastases. *Radiotherapy and
301 Oncology*, 148:229–234, 2020.
- 302 [3] T. Ishii S. Mori, S. Ko and K. Nishizawa. Effective doses in four-dimensional computed
303 tomography for lung radiotherapy planning. *Medical Dosimetry*, 2009.
- 304 [4] J. M. Balter R. S. Emery K. M. Forster S. Jiang et al M. J. Keall, G. S. Mageras. The management
305 of respiratory motion in radiation oncology (aapm task group 76 report). *Medical Physics*, 2006.
- 306 [5] World Health Organization. *Global atlas of medical devices*. World Health Organization,
307 Geneva, 2022. ISBN 978-92-4-006220-7. URL <https://iris.who.int/bitstream/handle/10665/364709/9789240062207-eng.pdf>.
- 309 [6] Melissa Silverberg. How radiologists overcome barriers to provide imaging in low to
310 middle income countries, July 2024. URL <https://www.rsna.org/news/2024/july/imaging-in-lmics>. States: “There is less than one CT scanner per million inhabitants
311 in LMICs compared to 40 per million in high-income countries.”.
- 313 [7] Emma Beede, Elizabeth Elliott Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan
314 Raumviboonsuk, and Laura Vardoulakis. A human-centered evaluation of a deep learning
315 system deployed in clinics for the detection of diabetic retinopathy. In *Proceedings of the 2020
316 CHI Conference on Human Factors in Computing Systems*, CHI ’20, pages 1–12, New York,
317 NY, USA, 2020. Association for Computing Machinery. doi: 10.1145/3313831.3376718. URL
318 <https://dl.acm.org/doi/10.1145/3313831.3376718>.
- 319 [8] R. K. Otto R. S. Iyer G. S. Philips J. O. Swanson C. Zacharias, A. M. Alessio and M. M. Thapa.
320 Pediatric ct: strategies to lower radiation dose. *AJR American Journal of Roentgenology*, 2013.
- 321 [9] K. E. Nichols A. S. Levine K. Schneider, K. Zelley and J. Garber. Li-fraumeni syndrome.
322 genereviews. *NCBI Bookshelf*, 2025.
- 323 [10] T. Karras T. Aila S. Laine Y. Lipman, R. T. Q. Chen and J. Lehtinen. Flow matching for
324 generative modeling. In *NeurIPS*, 2022.
- 325 [11] Z. Liu T. Liu, T. Zhao and Y. Cao. Flow straight and fast: Learning to generate and refine
326 samples in one network. In *ICLR*, 2023.
- 327 [12] et al R. Varma. Medvae: Foundation autoencoders for medical imaging, 2024.
- 328 [13] W. Peebles and S. Xie. Scalable diffusion models with transformers. In *ICCV*, 2023.
- 329 [14] G. Farnebäck. Two-frame motion estimation based on polynomial expansion. *SCIA*, 2003.
- 330 [15] et al K. Clark. The cancer imaging archive (tcia) – 4d lung imaging of nsclc patients, 2016.
331 URL <https://doi.org/10.7937/K9/TCIA.2016.ELN8YGL>. Dataset.
- 332 [16] et al G. Hugo. 4d lung ct in radiotherapy: Dataset and motion analysis. *Medical Physics*, 2017.
- 333 [17] B. K. P. Horn and B. G. Schunck. Determining optical flow. *Artificial Intelligence*, 1981.
- 334 [18] T. Pock C. Zach and H. Bischof. A duality based approach for realtime tv-11 optical flow.
335 *Pattern Recognition (DAGM)*, 2007.
- 336 [19] et al D. Rueckert, L. I. Sonoda. Nonrigid registration using free-form deformations: Application
337 to breast mr. *IEEE TMI*, 18(8):712–721, 1999.
- 338 [20] et al M. Rogers. Lungregnet: Fast learning-based deformable registration for 4d lung ct. In
339 *MICCAI*, 2020.

- 340 [21] et al A. Hering, S. Heldmann. Learn2reg: Benchmark for learning-based medical image
341 registration. *TMI*, 2021.
- 342 [22] et al H. Xiao. A survey of deep learning-based medical image registration. *Medical Image
343 Analysis*, 2023.
- 344 [23] R. G. Keys. Cubic convolution interpolation for digital image processing. *IEEE TAP*, 29(6):
345 1153–1160, 1981.
- 346 [24] T. Blu P. Thévenaz and M. Unser. Interpolation revisited. *IEEE TMI*, 19(7):739–758, 2000.
- 347 [25] et al Y. Luo, S. Zheng. Training improvements for rectified flows. In *NeurIPS*, 2024.
- 348 [26] W. Chan et al J. Ho, A. Saharia. Video diffusion models. In *NeurIPS*, 2022.
- 349 [27] et al Z. Jin, C. Yu. Pyramidal flow matching for efficient video generation. In *CVPR*, 2024.
- 350 [28] Meta AI. Moviegen: A large-scale text-to-video system, 2024.
- 351 [29] Alibaba DAMO Academy. Wan2.1-flf2v-14b-720p: First-last-frame to video generation. *Tech
352 Report / Model Card*, 2025.
- 353 [30] et al P. Jain. Video interpolation with diffusion models. In *CVPR*, 2024.
- 354 [31] D. Lorenz P. Esser R. Rombach, A. Blattmann and B. Ommer. High-resolution image synthesis
355 with latent diffusion models. In *CVPR*, 2022.
- 356 [32] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural
357 networks. *Neural Computation*, 1(2):270–280, 1989.
- 358 [33] N. Jaitly S. Bengio, O. Vinyals and N. Shazeer. Scheduled sampling for sequence prediction. In
359 *NeurIPS*, 2015.
- 360 [34] et al A. Lamb, A. Goyal. Professor forcing: A new algorithm for training rnns. In *NeurIPS*,
361 2016.
- 362 [35] et al T. Song, J. Vahdat. Consistency models. In *ICML*, 2023.
- 363 [36] et al T. Song. Improved techniques for training consistency models. In *NeurIPS*, 2023.
- 364 [37] J. Ho and T. Salimans. Classifier-free diffusion guidance, 2022.
- 365 [38] E. P. Simoncelli Z. Wang and A. C. Bovik. Multiscale structural similarity for image quality
366 assessment. In *Proc. IEEE Asilomar Conf. Signals, Systems, and Computers*, 2003.