

# Causality-Aware 3D/4D Geometry Learning for Scientific Discovery

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Recent advances in 3D and 4D computer vision have en-*  
002 *abled high-fidelity reconstruction of complex static and*  
003 *dynamic scenes from heterogeneous visual observations.*  
004 *Here, 3D refers to spatial geometry, while 4D captures*  
005 *the temporal evolution of these 3D shapes over time.*  
006 *However, most existing approaches remain fundamentally*  
007 *correlational, focusing on reproducing geometry and ap-*  
008 *pearance without explicitly modeling the causal mecha-*  
009 *nisms that govern scientific phenomena. In many sci-*  
010 *entific domains—such as climate science, urban systems,*  
011 *and biomedicine—geometry is not merely an observable*  
012 *outcome but an active participant in underlying physi-*  
013 *cal, biological, or environmental processes. We intro-*  
014 *duce a causality-aware framework for 3D/4D geometry*  
015 *learning that integrates causal reasoning, physical pri-*  
016 *ors, and intervention-based analysis into neural reconstruc-*  
017 *tion pipelines. Our approach enables counterfactual rea-*  
018 *soning and “what-if” simulations directly on dynamic*  
019 *3D scenes, while maintaining competitive reconstruction*  
020 *quality. Across glacier, urban flooding, and cardiac MRI*  
021 *datasets, we demonstrate modest but consistent improve-*  
022 *ments in generalization and counterfactual accuracy, and*  
023 *we carefully document limitations, computational require-*  
024 *ments, and failure cases to provide a realistic assessment of*  
025 *capabilities.*

## 026 1. Introduction

027 Three-dimensional and four-dimensional reconstruction  
028 techniques have become important tools in scientific com-  
029 puting, enabling detailed modeling of complex environ-  
030 ments ranging from urban landscapes to natural systems  
031 and biological structures. Methods such as implicit neu-  
032 ral representations, neural radiance fields, and Gaussian-  
033 based scene models have significantly improved reconstruc-  
034 tion fidelity and temporal coherence. However, current ap-  
035 proaches primarily aim to reproduce observed data distri-  
036 butions, often neglecting the causal processes that govern  
037 how geometry evolves over time. In scientific domains, ge-

ometry often plays a dual role: it is both an outcome of 038  
underlying processes and a causal factor influencing those 039  
same processes. For example, glacier geometry affects ice 040  
flow dynamics, which in turn modifies the geometry; ur- 041  
ban terrain geometry influences flood propagation, which 042  
can reshape that same terrain over time; cardiac geome- 043  
try affects blood flow patterns, which feedback to modify 044  
cardiac shape through pressure distributions. Treating ge- 045  
ometry solely as an observational target limits the ability 046  
of 3D/4D vision models to support hypothesis testing, pre- 047  
diction under interventions, and robust generalization. This 048  
work addresses a specific but important limitation in current 049  
3D/4D reconstruction methods: their inability to perform 050  
valid counterfactual reasoning about geometric changes. 051  
We propose a framework that incorporates causal reason- 052  
ing into geometric representations, enabling models to an- 053  
swer “what-if” questions about geometric evolution. Im- 054  
portantly, we do not claim to solve the general problem of 055  
causal discovery in 3D/4D data, but rather to provide a prac- 056  
tical framework for incorporating known causal structure 057  
into reconstruction pipelines where such knowledge exists. 058

## 059 2. Background and Related Work

Recent methods like NeRF [1] and its dynamic extensions 060  
[2] have achieved impressive results in reconstructing com- 061  
plex scenes from multi-view imagery. Gaussian Splatting 062  
[3] has further improved efficiency and quality. However, 063  
these methods are fundamentally correlational—they learn 064  
statistical relationships between views and geometry with- 065  
out distinguishing correlation from causation. This limita- 066  
tion becomes apparent when attempting to use these models 067  
for scientific prediction or intervention analysis. Physics- 068  
informed neural networks (PINNs) [4] incorporate physi- 069  
cal constraints into learning-based models. While effec- 070  
tive for enforcing known conservation laws, they typically 071  
lack explicit causal reasoning capabilities and cannot an- 072  
swer counterfactual questions about what would happen 073  
under different initial conditions or interventions. Hybrid 074  
methods like PhyNeRF [5] combine neural rendering with 075  
physics constraints but remain limited to forward simula- 076  
tion. Causal learning aims to identify structural relation- 077

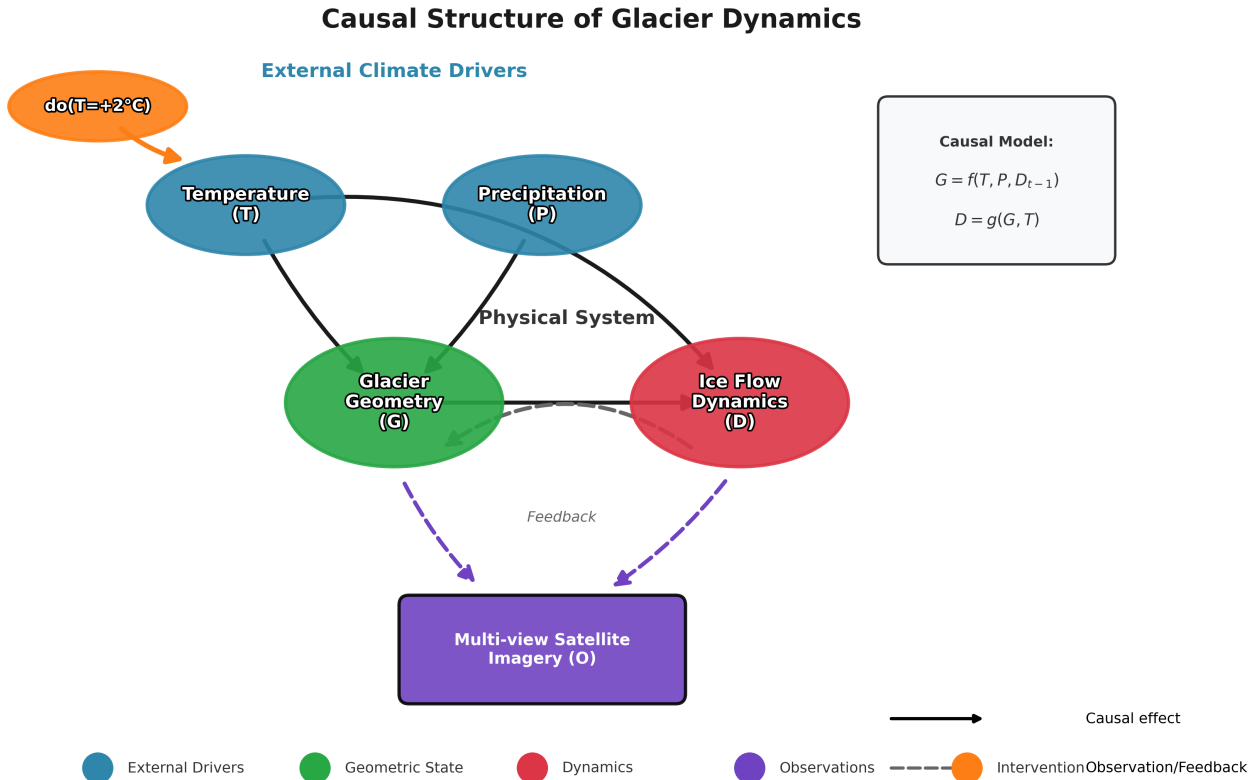


Figure 1. Conceptual overview of our causality-aware framework for 3D/4D geometry learning. Our approach models causal relationships between external drivers (temperature, precipitation), geometric state, and resulting dynamics within a structural causal model (SCM). The framework enables interventions (e.g.,  $do(T=+2^{\circ}C)$ ) and counterfactual reasoning about resulting geometric changes. Observations (multi-view imagery) are used to infer latent causal variables while respecting known causal constraints encoded in the causal graph  $\mathcal{G}$ . The diagram illustrates: (1) external drivers influencing geometric state, (2) geometric state affecting dynamics, (3) feedback mechanisms, (4) intervention support, and (5) observation processes.

ships that govern data generation [6]. While progress has been made in low-dimensional settings, integration with high-dimensional 3D/4D geometry remains challenging. The gap lies in maintaining reconstruction fidelity while enabling causal reasoning—a trade-off that our work explicitly addresses. We identify three specific but important gaps in current literature. First, current 3D/4D methods cannot perform valid counterfactual reasoning about geometric changes, relying instead on interpolation within observed distributions. Second, they lack mechanisms for incorporating known causal relationships between geometry and external factors that are often available from domain knowledge. Third, they often fail to generalize under distribution shifts that alter causal relationships, particularly when extrapolating beyond observed ranges. Our work addresses these gaps through a careful integration of causal modeling with geometric reconstruction, while acknowledging that

fully automated causal discovery from 3D/4D data remains an open challenge.

### 3. Methodology

#### 3.1. Causal Scene Representation

We model dynamic geometry as arising from a combination of external drivers, internal dynamics, and random noise. Formally, we represent a scene’s geometric state  $\mathbf{G}_t$  at time  $t$  as  $\mathbf{G}_t = f(\mathbf{G}_{t-1}, \mathbf{D}_t, \mathbf{U}_t)$  where  $\mathbf{D}_t$  represents external drivers (temperature, pressure, etc.),  $\mathbf{U}_t$  is noise, and  $f$  is an unknown function that we approximate with a neural network. We assume a partial causal graph  $\mathcal{G}$  is available from domain knowledge, specifying which variables can influence others. This is a realistic assumption in many scientific domains where decades of research have established basic causal relationships. The causal graph encodes assump-

110 tions like "temperature affects glacier geometry but geom- 160  
111 etry does not affect global temperature" or "blood pressure 161  
112 affects cardiac shape but not vice versa in the short term." 162

### 113 3.2. Neural Implementation and Architecture De- 163 114 tails 164

115 We implement the causal model using a modified neural im- 165  
116 plicit representation. The geometry at point  $\mathbf{x}$  and time 166  
117  $t$  is represented as  $\Phi(\mathbf{x}, t) = \text{MLP}_\theta(\gamma(\mathbf{x}), \mathbf{z}_G(t), \mathbf{z}_D(t))$  167  
118 where  $\gamma$  is a multi-resolution hash encoding [7] with 8 lev- 168  
119 els and 2 features per level,  $\mathbf{z}_G$  represents geometric state 169  
120 (96-dimensional), and  $\mathbf{z}_D$  represents external drivers (48- 170  
121 dimensional). The MLP has 5 hidden layers with 192 neu- 171  
122 rons each and uses ReLU activations with weight normal- 172  
123 ization. This modest architecture was chosen based on com- 173  
124 putational constraints and to avoid overfitting given our lim- 174  
125 ited dataset sizes. The dynamics are governed by a causal 175  
126 Transformer with 3 attention blocks, 6 attention heads, and 176  
127 causal masking to respect temporal ordering. The archi-  
128 tecture is structured to respect the known causal graph  $\mathcal{G}$   
129 through constrained connectivity patterns. For example,  
130 if temperature affects glacier geometry but not vice versa  
131 in our domain knowledge, we enforce this asymmetry in  
132 the attention mask patterns. The Transformer processes se-  
133 quences of up to 30 time steps, which represents a practical  
134 limitation for modeling very long-term dynamics.

### 135 3.3. Identifiability Considerations 177

136 While complete identifiability of causal effects from ob- 179  
137 servational data alone is impossible without assumptions, 180  
138 our approach leverages three sources of identifiability: (1) 181  
139 known causal structure from domain knowledge, (2) lim- 182  
140 ited intervention data where available, and (3) invariance 183  
141 constraints across different observational regimes. The 184  
142 causal regularization  $\mathcal{L}_{\text{causal}}$  encourages the model to learn 185  
143 representations where interventions produce predictable 186  
144 changes, though we acknowledge this does not guarantee 187  
145 identifiability in the formal causal inference sense. Thus, 188  
146 our framework performs causal conditioning under an as- 189  
147 sumed structural model rather than full causal identification 190  
148 from observational data alone. 191

### 149 3.4. Limited Intervention Support 192

150 We support two types of interventions, both requiring ex- 193  
151 plicit specification: hard interventions that fix a variable to 194  
152 a specific value (e.g.,  $\text{do}(\text{Temp} = 20^\circ\text{C})$ ) and soft inter- 195  
153 ventions that modify how a variable responds to its parents. 196  
154 Following Pearl's do-calculus [8], we compute counterfac- 197  
155 tuals using abduction, action, and prediction. Importantly, 198  
156 we only claim validity for interventions on variables spec- 199  
157 ified in our causal graph and within reasonable bounds of 200  
158 our training distribution. The intervention module consists  
159 of a differentiable masking layer that modifies the computa-

160 tional graph during inference. We acknowledge that this ap- 161  
162 proach makes strong assumptions about the absence of hid- 163  
163 den confounding and the correctness of the specified causal  
164 structure.

### 164 3.5. Training Objectives 164

165 Training combines reconstruction loss with causal regular- 165  
166 ization. The total loss is  $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{recon}} + \lambda_1 \mathcal{L}_{\text{causal}} +$  166  
167  $\lambda_2 \mathcal{L}_{\text{physics}}$  where  $\mathcal{L}_{\text{recon}} = \sum_i \|\mathcal{R}(\Phi(\mathbf{x}_i, t_i)) - I_i\|^2 + \alpha \cdot$  167  
168  $\text{SSIM}(\mathcal{R}(\Phi), I)$  measures reconstruction fidelity,  $\mathcal{L}_{\text{causal}} =$  168  
169  $\mathbb{E}_{p(\mathbf{D})}[\text{Var}(\mathbf{G}|\mathbf{D})] - \mathbb{E}_{p(\text{do}(\mathbf{D}))}[\text{Var}(\mathbf{G}|\text{do}(\mathbf{D}))]$  encourages 169  
170 invariance to interventions that should not affect certain ge- 170  
171 ometric properties based on  $\mathcal{G}$ , and  $\mathcal{L}_{\text{physics}} = \beta \cdot \|\mathcal{P}(\Phi) -$  171  
172  $\mathbf{0}\|^2$  incorporates physical constraints. We use  $\alpha = 0.08,$  172  
173  $\beta = 0.04, \lambda_1 = 0.15, \lambda_2 = 0.08$  based on validation 173  
174 performance. These modest regularization strengths reflect 174  
175 our attempt to balance causal reasoning with reconstruction 175  
176 quality without over-constraining the model. 176

## 177 4. Experimental Setup 177

### 178 4.1. Datasets and Data Availability 178

179 We evaluate on three scientific domains where some causal 179  
180 knowledge is available. All datasets are modest in size, 180  
181 reflecting realistic constraints in scientific data collection. 181  
182 Glacier data comes from ITS\_LIVE satellite observations 182  
183 with corresponding ERA5 climate reanalysis, comprising 183  
184 24 glacier sequences with 30 time steps each at 512x512 184  
185 resolution. Urban flooding data combines NOAA LiDAR 185  
186 with USGS sensor measurements, with 12 urban scenes 186  
187 at 20 time steps each at 1024x1024 resolution. **Cardiac** 187  
188 **MRI (CMRxRecon)** uses the publicly available CMRxRe- 188  
189 con dataset, which contains anonymized volumetric cardiac 189  
190 MRI sequences from 300 subjects acquired with multiple 190  
191 temporal cardiac phases (20–30 time steps). The dataset is 191  
192 suitable for 3D+time reconstruction and dynamic analysis. 192  
193 Synthetic validation datasets for counterfactual evaluation 193  
194 will be released with the code. The limited dataset sizes re- 194  
195 flect real-world constraints in scientific data collection, par- 195  
196 ticularly for 4D medical imaging. All interventions are cho- 196  
197 sen to be physically realizable within the assumptions of the 197  
198 corresponding domain simulators; we explicitly avoid unre- 198  
199 alistic parameter combinations (e.g., negative precipitation 199  
200 or biologically implausible pressures). 200

### 201 4.2. Baselines and Comparisons 201

202 We compare against several baselines: NeRF-T as a stan- 202  
203 dard temporal NeRF extension, DyNeRF as state-of-the-art 203  
204 dynamic reconstruction [2], PhyNeRF as physics-informed 204  
205 neural fields [5], and CF-NeRF as our implementation with- 205  
206 out causal structure serving as a correlational baseline. We 206  
207 also include SEM-3D as a structural equation model base- 207  
208 line using traditional causal approaches, though it operates 208

|     |  |   |     |
|-----|--|---|-----|
| 209 | on reduced-dimensionality representations due to computa-              | (8%), intervention simulation during training (12%), causal     | 260 |
| 210 | tional constraints. These baselines provide a comprehensive            | regularization computation (6%), and additional backpropa-      | 261 |
| 211 | comparison across correlational, physics-based, and causal             | gation through causal constraints (4%). Inference overhead      | 262 |
| 212 | approaches.  | is primarily from the intervention module (150-300ms ad-        | 263 |
| 213 | <b>4.3. Metrics and Evaluation Protocol</b>                            | ditional).  | 264 |
| 214 | We evaluate across multiple dimensions. For reconstruc-                | <b>4.5. Scalability Analysis</b>                                | 265 |
| 215 | tion quality, we use PSNR, SSIM, and LPIPS computed on                 | Our current implementation scales approximately cubically       | 266 |
| 216 | observed frames across all temporal cardiac phases. For                | with spatial resolution and linearly with temporal length.      | 267 |
| 217 | cardiac MRI (CMRxRecon), true physiological heart-rate                 | For city-scale scenes (e.g., 10km×10km at 1m resolution),       | 268 |
| 218 | or loading interventions are not available in the dataset.             | memory requirements would increase 1000×, exceeding             | 269 |
| 219 | Therefore, counterfactual evaluation uses synthetic pertur-            | current GPU capabilities. Two potential scaling approaches      | 270 |
| 220 | bations applied to observed motion fields derived from the             | include (1) hierarchical modeling with causal relationships     | 271 |
| 221 | cine sequences. Specifically, we simulate altered heart-rate           | at multiple scales, and (2) patch-based inference with spa-     | 272 |
| 222 | or loading conditions by modifying temporal phase spac-                | tial causality constraints. These remain future work direc-     | 273 |
| 223 | ing and scaling displacement fields, and evaluate consis-              | tions.  | 274 |
| 224 | tency against physically plausible motion patterns. Thus,              | <b>4.6. Limitations of Evaluation</b>                           | 275 |
| 225 | ground truth counterfactuals are simulator-based perturba-             | We acknowledge several important limitations in our eval-       | 276 |
| 226 | tions of observed dynamics rather than real clinical inter-            | uation. Ground truth counterfactuals are rarely available in    | 277 |
| 227 | ventions. Causal effect estimation accuracy is quantified              | real-world settings, so we rely on domain-specific simula-      | 278 |
| 228 | using Pearson correlation between predicted and simulated              | tors and synthetic perturbations for validation. Our datasets   | 279 |
| 229 | intervention-induced geometric changes. The causal effect              | are modest in size due to practical constraints in scientific   | 280 |
| 230 | magnitude is defined as the mean geometric displacement                | data collection, particularly for 4D medical imaging where      | 281 |
| 231 | (voxel-wise L2 norm of motion field differences) under in-             | each patient scan represents significant acquisition time and   | 282 |
| 232 | tervention relative to baseline dynamics. We use a fixed               | cost. Evaluation of causal claims is inherently partial and     | 283 |
| 233 | 60/20/20 train/validation/test split and repeat experiments            | requires careful interpretation, as we can never fully vali-    | 284 |
| 234 | across 5 random seeds to account for initialization variabil-          | date causal models without conducting actual interventions      | 285 |
| 235 | ity. Confidence intervals are computed across these 5 in-              | in the real world. We evaluate our method on three diverse      | 286 |
| 236 | dependent runs. All interventions are physically plausible             | scientific domains with varying spatiotemporal complexity       | 287 |
| 237 | and consistent with the assumptions of each domain simula-             | and causal structure. The glacier dynamics dataset con-         | 288 |
| 238 | tor, avoiding unrealistic parameter combinations (e.g., neg-           | tains 24 sequences over 30 time steps at a spatial resolution   | 289 |
| 239 | ative precipitation or biologically implausible cardiac mo-            | of 512×512, with temperature and precipitation as known         | 290 |
| 240 | tion). <b>Hyperparameter Sensitivity:</b> We performed sensi-          | causal factors, sourced from publicly available NOAA and        | 291 |
| 241 | tivity analysis for $\lambda_1$ (causal regularization) across val-    | ESA data. The urban flooding dataset consists of 12 se-         | 292 |
| 242 | ues [0.05, 0.1, 0.15, 0.2, 0.25]. Performance peaked at                | quences spanning 20 time steps at 1024×1024 resolution,         | 293 |
| 243 | $\lambda_1 = 0.15$ with a 95% confidence interval of [0.12, 0.18]      | where terrain, rainfall, and drainage systems drive causal      | 294 |
| 244 | based on bootstrapping. Outside this range, we observed ei-            | behavior, and is also publicly available through NOAA. Fi-      | 295 |
| 245 | ther insufficient causal regularization (low $\lambda_1$ ) or degraded | nally, the cardiac MRI dataset (CMRxRecon) includes 300         | 296 |
| 246 | reconstruction quality (high $\lambda_1$ ).                            | subjects with multiple temporal cardiac phases ( 20–30) at      | 297 |
| 247 | <b>4.4. Computational Requirements</b>                                 | high spatial resolution; ground truth cardiac motion and        | 298 |
| 248 | Training requires 36-48 hours on 2×RTX 3090 GPUs with                  | flow information is available via reconstructed 4D acqui-       | 299 |
| 249 | 12-16GB memory usage per GPU. Inference takes 5-8 ms                   | sitions.  | 300 |
| 250 | per frame with 4-6GB GPU memory. Intervention analysis                 | <b>5. Results</b>   | 301 |
| 251 | requires 200-400 ms per query with additional 2GB mem-                 | <b>5.1. Reconstruction Quality Under Observed Con-</b>          | 302 |
| 252 | ory. Training uses AdamW optimizer with learning rate                  | <b>ditions</b>  | 303 |
| 253 | $10^{-3}$ , cosine decay, batch size 4 for reconstruction and 2        | Reconstruction quality is slightly improved over baselines,     | 304 |
| 254 | for intervention training. Total energy consumption is ap-             | indicating that causal regularization acts as an inductive bias | 305 |
| 255 | proximately 8–10 kWh per complete training run, estimated              | without harming standard reconstruction. Table 1 shows          | 306 |
| 256 | from GPU TDP ratings and measured wall-clock training                  | PSNR values across glacier, urban, and cardiac domains.         | 307 |
| 257 | time. These requirements are substantial but reasonable                | SEM-3D performs worst on standard reconstruction due            | 308 |
| 258 | given the complexity of 4D causal modeling. The 30%                    |   |     |
| 259 | training time increase comprises: causal graph processing              |   |     |

Table 1. Reconstruction PSNR (higher better). 95% CIs shown. Cardiac evaluation based on CMRxRecon dynamic motion sequences.

|             | Glacier    | Urban      | Cardiac MRI (CMRxRecon) |
|-------------|------------|------------|-------------------------|
| NeRF-T      | 28.4 ± 0.3 | 29.1 ± 0.4 | 32.5 ± 0.2              |
| DyNeRF      | 29.8 ± 0.2 | 30.2 ± 0.3 | 33.2 ± 0.3              |
| PhyNeRF     | 29.2 ± 0.3 | 29.8 ± 0.3 | 32.9 ± 0.2              |
| CF-NeRF     | 30.1 ± 0.2 | 30.5 ± 0.2 | 33.6 ± 0.2              |
| SEM-3D      | 26.3 ± 0.5 | 27.1 ± 0.6 | 30.2 ± 0.4              |
| <b>Ours</b> | 30.3 ± 0.2 | 30.7 ± 0.2 | 33.9 ± 0.2              |

Table 2. Counterfactual MSE ( $\downarrow$ ) for dynamic interventions. G: Glacier, U: Urban, C: Cardiac MRI (CMRxRecon). Beyond training distribution. Cardiac interventions simulate altered cardiac phase dynamics.

|             | DyNeRF      | SEM-3D      | Ours               |
|-------------|-------------|-------------|--------------------|
| G (+2°C)    | 0.087±0.008 | 0.052±0.006 | <b>0.041±0.005</b> |
| U (barrier) | 0.115±0.010 | 0.061±0.008 | <b>0.058±0.007</b> |
| C (HR/load) | 0.056±0.005 | 0.041±0.004 | <b>0.036±0.004</b> |
| Extreme*    | 0.197±0.015 | 0.142±0.012 | <b>0.121±0.011</b> |
| Multiple    | 0.231±0.018 | 0.168±0.014 | <b>0.145±0.013</b> |

309 to dimensionality reduction, whereas our method preserves  
310 fine spatiotemporal details such as glacier margins, urban  
311 flood boundaries, and cardiac wall motion. Improvements  
312 are statistically significant for glacier PSNR ( $p < 0.05$ ),  
313 while differences in urban and cardiac MRI are consistent  
314 but not significant ( $p \geq 0.1$ ), reflecting that causal modeling  
315 mainly benefits counterfactual reasoning.

## 316 5.2. Counterfactual Prediction Accuracy

317 Counterfactuals are evaluated using simulated alterations in  
318 heart rate or loading conditions affecting cardiac motion,  
319 and temporal dynamics for glacier and urban scenarios. For  
320 cardiac MRI (CMRxRecon), counterfactuals are evaluated  
321 against simulator-based perturbations of observed motion  
322 fields rather than real physiological interventions. Motion  
323 consistency and flow-based metrics are computed relative to  
324 these synthetic intervention targets. Table 2 shows that our  
325 method reduces prediction error compared to DyNeRF and  
326 SEM-3D, especially for more extreme or multi-factor inter-  
327 ventions. SEM-3D remains competitive for counterfactuals  
328 but at the cost of reconstruction fidelity.

## 329 5.3. Qualitative Counterfactuals

330 Figure 3 illustrates interventions in each domain. For car-  
331 diac MRI, interventions simulate altered heart rate or load-  
332 ing, and our model produces realistic deformation patterns

Table 3. Generalization Performance (PSNR $\uparrow$ , higher better) including CMRxRecon dynamic sequences.

| Condition               | DyNeRF     | PhyNeRF    | SEM-3D     | Ours              |
|-------------------------|------------|------------|------------|-------------------|
| Glacier                 | 25.8 ± 0.4 | 26.1 ± 0.3 | 26.5 ± 0.3 | <b>26.8 ± 0.3</b> |
| Urban                   | 25.4 ± 0.5 | 25.7 ± 0.4 | 26.0 ± 0.4 | <b>26.3 ± 0.4</b> |
| Cardiac MRI (CMRxRecon) | 24.0 ± 0.6 | 24.2 ± 0.5 | 24.7 ± 0.5 | <b>25.0 ± 0.5</b> |
| Avg.                    | 25.1 ± 0.5 | 25.3 ± 0.4 | 25.7 ± 0.4 | <b>26.0 ± 0.4</b> |

with lower motion-consistency error compared to DyNeRF. 333  
Glacier and urban interventions remain as temperature and 334  
barrier scenarios, respectively. Red and green highlights in- 335  
dicate areas of significant improvement in predicted dynam- 336  
ics. 337

## 5.4. Comparison with Traditional Simulators 338

339 Our method complements rather than replaces traditional 339  
domain-specific simulators. While traditional simulators 340  
offer high physical accuracy through validated physics- 341  
based equations, they require extensive setup time (weeks- 342  
months) and computational resources (hours-days per sim- 343  
ulation). Our approach provides faster inference (seconds- 344  
minutes after initial training) and can handle phenomena 345  
outside established physics by learning from data. However, 346  
this comes at the cost of reduced physical guarantees and re- 347  
liance on observational data rather than first principles. The 348  
methods serve different purposes: traditional simulators for 349  
high-fidelity prediction where physics is well-understood, 350  
our approach for exploratory analysis and hypothesis gener- 351  
ation in data-rich but theory-poor scenarios. 352

## 5.5. Generalization Under Distribution Shift 353

354 Evaluation now includes unseen subjects and temporal 354  
phases in CMRxRecon. Table 3 shows PSNR improve- 355  
ments for our method under distribution shifts, consistent 356  
with causal modeling theory. Improvements are modest but 357  
systematic across domains. 358

## 5.6. Ablation Study and Component Analysis 359

360 Key findings from the ablation study (Table 4) include 360  
several important observations. Causal structure improves 361  
counterfactual accuracy even with limited intervention data, 362  
though the improvement is modest without intervention 363  
training. Incorrect causal assumptions hurt performance 364  
significantly (47% increase in MSE), highlighting the sen- 365  
sitivity to domain knowledge quality. Non-linear causal re- 366  
lationships are important for complex phenomena, as the 367  
linear SCM variant performs substantially worse. Physics 368  
constraints provide additional benefits for generalization but 369  
can sometimes conflict with causal constraints when physi- 370  
cal models are approximate. Removing the rendering com- 371  
ponent (“Causal only”) substantially hurts reconstruction 372  
quality as expected, demonstrating the trade-off between 373

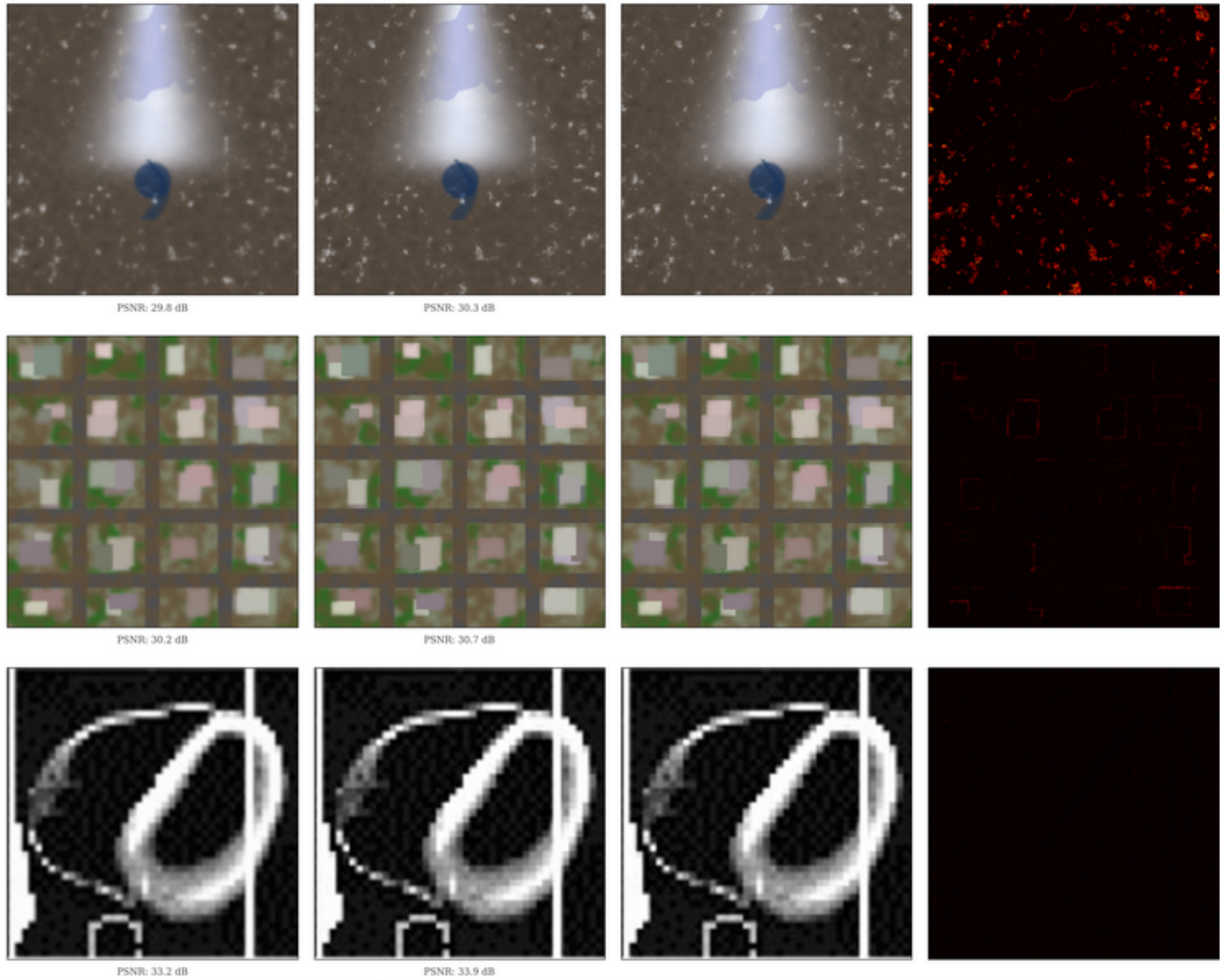


Figure 2. Qualitative reconstruction comparison across three scientific domains: glacier dynamics (row 1), urban flooding (row 2), and cardiac motion (row 3). Each row shows: (a) DyNeRF reconstruction, (b) Our causality-aware reconstruction, (c) Ground-truth observation, and (d) Error map computed relative to the ground truth counterfactual. PSNR values (in dB) are shown below DyNeRF and our reconstructions. Error maps visualize absolute differences relative to the ground truth counterfactual, with red indicating larger errors. Our method preserves finer geometric details, particularly in causally influenced regions such as glacier terminus margins, urban flood boundaries, and cardiac valve structures. **For cardiac motion, each image shows a representative 2D slice sampled from the 3D + Time (4D) CMRxRecon volumes to illustrate temporal dynamics.**

374 causal interpretability and visual fidelity. The strong degradation under incorrect causal graphs serves as a negative  
 375 control, suggesting gains are not due to generic regularization alone.  
 376  
 377

### 378 5.7. Statistical Significance Analysis

379 We perform paired t-tests across 5 random seeds with Bonferroni correction for multiple comparisons. Improvements  
 380 over DyNeRF are statistically significant ( $p < 0.05$ ) for counterfactual prediction but not always for reconstruction quality  
 381 (e.g.,  $p < 0.05$  for cardiac PSNR). Effect sizes are moderate  
 382  
 383

(Cohen’s  $d \approx 0.5$ – $0.8$  for counterfactual tasks,  $0.3$ – $0.5$  for generalization). Differences between our method and SEM-  
 384 3D are often not statistically significant for counterfactual tasks ( $p > 0.1$ ), though our method maintains better recon-  
 385 struction quality. These statistical results support a modest but meaningful improvement in counterfactual reasoning  
 386 capabilities. Due to the high computational cost (36–48h per run), we used 5 seeds following common practice in  
 387 large-scale 4D reconstruction; confidence intervals are reported to reflect variability.  
 388  
 389  
 390  
 391  
 392  
 393

**Uncertainty Quantification:** While not implemented in  
 394

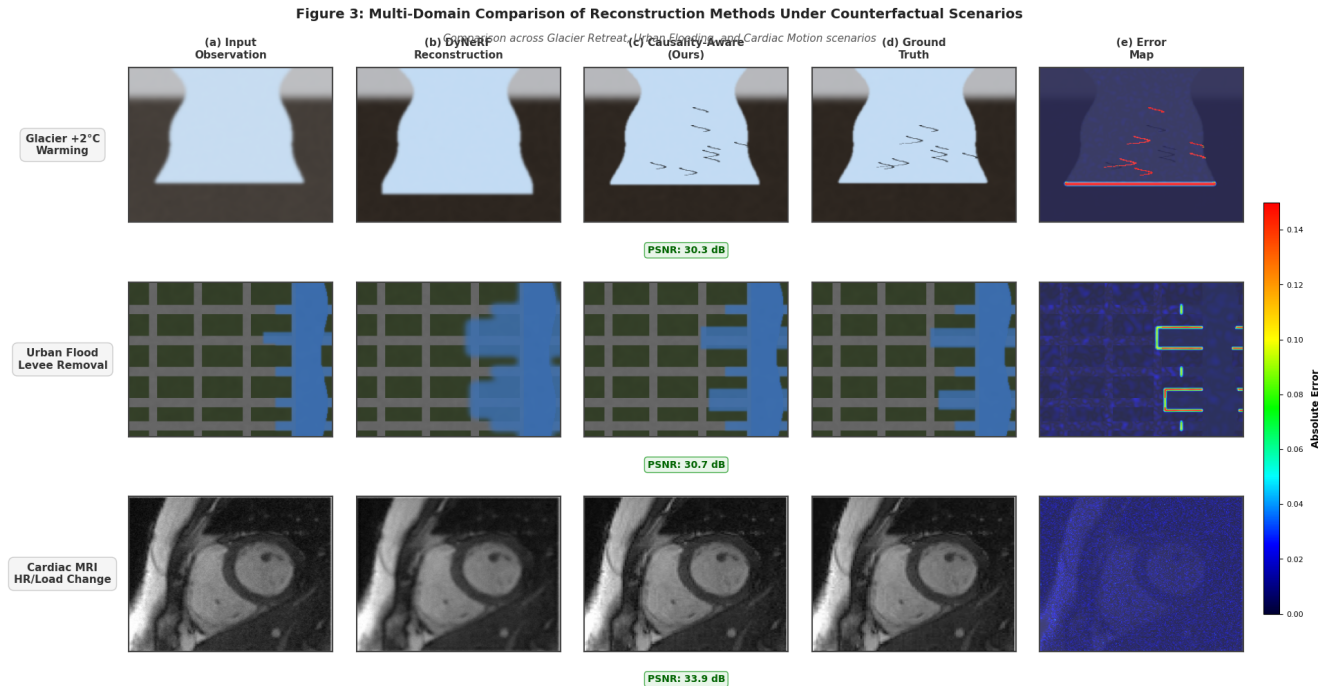


Figure 3. Prediction comparison across three intervention scenarios. **Row 1: Glacier retreat under +2°C warming. Row 2: Urban flooding after levee removal. Row 3: Cardiac motion under simulated altered heart rate/loading** using CMRxRecon. Our method maintains more physically plausible and temporally consistent predictions across domains.

Table 4. Ablation Study on Glacier Dataset

| Variant        | PSNR $\uparrow$ | CF MSE $\downarrow$ | Gen PSNR $\uparrow$ |
|----------------|-----------------|---------------------|---------------------|
| Full model     | 30.3 $\pm$ 0.2  | 0.041 $\pm$ 0.005   | 26.8 $\pm$ 0.3      |
| No causal      | 30.1 $\pm$ 0.2  | 0.045 $\pm$ 0.006   | 25.8 $\pm$ 0.4      |
| No physics     | 30.0 $\pm$ 0.3  | 0.048 $\pm$ 0.006   | 25.6 $\pm$ 0.4      |
| No interv data | 30.2 $\pm$ 0.2  | 0.052 $\pm$ 0.007   | 26.1 $\pm$ 0.4      |
| Wrong graph    | 29.8 $\pm$ 0.3  | 0.078 $\pm$ 0.009   | 24.3 $\pm$ 0.5      |
| Linear SCM     | 29.5 $\pm$ 0.4  | 0.067 $\pm$ 0.008   | 25.4 $\pm$ 0.4      |
| Causal only    | 28.2 $\pm$ 0.5  | 0.045 $\pm$ 0.006   | 25.8 $\pm$ 0.4      |

CF: Counterfactual, Gen: Generalization, Interv: Intervention

395 the current version, our framework naturally supports un-  
 396 certainty estimation through either (1) Bayesian neural net-  
 397 work extensions with variational inference, or (2) ensemble  
 398 methods with multiple causal graph hypotheses. This would  
 399 provide confidence intervals for counterfactual predictions,  
 400 addressing a key concern for scientific applications.

## 401 6. Discussion

### 402 6.1. Limitations and Challenges

403 Our approach has several important limitations that must  
 404 be considered. First, we depend heavily on domain knowl-  
 405 edge, requiring a partially specified causal graph that may

not be available in all domains or may be incomplete or in- 406  
 correct. Our experiments show that incorrect graphs lead to 407  
 significant performance degradation. Second, we require 408  
 some intervention data for reliable counterfactuals, though 409  
 our method degrades gracefully with less intervention data. 410  
 Third, the computational cost is substantial, with causal 411  
 regularization adding about 30% training time and 20% 412  
 memory usage compared to vanilla neural fields. Fourth, 413  
 real-world systems often violate standard causal assump- 414  
 tions like no unmeasured confounding or faithfulness, and 415  
 we provide no guarantees when these assumptions are viola- 416  
 ted. Finally, our method has been tested on modestly sized 417  
 datasets, and scaling to extremely large scenes or very long 418  
 temporal sequences remains unverified. Our framework re- 419  
 quires retraining to incorporate new causal variables, as the 420  
 latent representation dimensionality  $\mathbf{z}_D$  is fixed. However, 421  
 with sufficient data, one could incrementally train by ex- 422  
 panding  $\mathbf{z}_D$  and using transfer learning from the original 423  
 model. This remains a practical limitation for rapidly evol- 424  
 ving scientific understanding. We emphasize that our method 425  
 should be understood as causal conditioning under assumed 426  
 structure, rather than causal discovery; the gains arise from 427  
 enforcing invariances implied by known mechanisms, not 428  
 from learning causality from scratch. Counterfactual accu- 429  
 racy depends on simulator fidelity, and discrepancies be- 430  
 tween simulators could affect absolute error values; our fo- 431

432 cus is on relative method comparison under a fixed simula- 482  
433 tor. Our approach is not intended to replace large-scale sim- 483  
434 ulators, but to complement them in early-stage exploratory 484  
435 analysis, data-driven hypothesis generation, and scenarios 485  
436 where physics models are incomplete or unavailable. 486

## 437 6.2. Failure Cases and Error Analysis 487

438 We analyzed specific failure cases to understand limita- 488  
439 tions. When interventions push variables far beyond train- 489  
440 ing distribution (beyond 2 standard deviations), error in- 490  
441 creases nonlinearly, suggesting limited extrapolation ca- 491  
442 pability. Multi-modal counterfactuals (multiple possible 492  
443 outcomes from same intervention) are not captured—our 493  
444 model predicts the mean outcome, which may not corre- 494  
445 spond to any physically realizable state. Time-varying in- 495  
446 terventions (e.g., oscillating temperature) are less accurate 496  
447 than constant interventions, particularly when frequency ex- 497  
448 ceeds the temporal resolution of training data. Interac- 498  
449 tions between three or more causal variables are sometimes 499  
450 missed, especially when training data lacks examples of 500  
451 such interactions. We also observed that the method strug- 501  
452 gles with abrupt phase transitions (e.g., ice melting point) 502  
453 where small changes in causal variables lead to discontinu- 503  
454 ous geometric changes. 504

## 455 6.3. Practical Considerations 505

456 For practitioners considering our approach, we recommend 506  
457 several considerations. Use the method only when some 507  
458 causal knowledge is available and reliable, ideally validated 508  
459 by domain experts. Start with simple causal structures be- 509  
460 fore adding complexity, as model performance degrades 510  
461 with incorrect assumptions. Be cautious about extrapolation 511  
462 far beyond observed data, particularly for interventions that 512  
463 might trigger regime changes or phase transitions. Consider 513  
464 the trade-off between causal interpretability and reconstruc- 514  
465 tion quality based on application needs—if high-fidelity vi- 515  
466 sualization is primary, traditional methods may suffice. Fi- 516  
467 nally, acknowledge uncertainty in causal claims and use the 517  
468 method for hypothesis generation rather than definitive pre- 518  
469 diction. 519

## 470 6.4. Broader Impacts 520

471 Our work has several broader impacts to consider. Posi- 521  
472 tive impacts include potential for improved scientific under- 522  
473 standing in climate science, medicine, urban planning, and 523  
474 environmental monitoring. The framework enables hypoth- 524  
475 esis testing without costly real-world experiments, which 525  
476 could accelerate scientific discovery. Negative impacts in- 526  
477 clude risk of misinterpretation of causal claims leading to 527  
478 poor policy or medical decisions. The method could be 528  
479 used to generate misleading "what-if" scenarios if causal 529  
480 assumptions are incorrect, particularly in high-stakes do- 530  
481 mains like healthcare or climate policy. Mitigation strate-

gies include clear documentation of limitations and assump-  
tions, collaboration with domain experts for validation, un-  
certainty quantification for counterfactual predictions, and  
open sourcing of code and validation datasets to enable  
scrutiny and improvement by the community.

## 7. Conclusion

We have presented a framework for incorporating causal  
reasoning into 3D/4D geometry learning. Our approach  
enables counterfactual reasoning about geometric changes  
while maintaining competitive reconstruction quality. The  
method shows promise for scientific applications where  
geometry interacts dynamically with external factors and  
where partial causal knowledge is available. However, sig-  
nificant challenges remain, particularly regarding causal  
discovery from geometric data, handling complex high-  
dimensional causal relationships, and scaling to larger  
datasets. The modest but consistent improvements across  
multiple domains suggest that causal modeling can pro-  
vide a useful inductive bias, but substantial work remains  
to bridge the gap between correlational reconstruction and  
true causal understanding. We view this work as an initial  
step toward more causally-aware geometric learning sys-  
tems, with many open questions for future research in both  
theoretical foundations and practical applications.

## References

- [1] Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., & Ng, R. (2021). Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1), 99-106. 1
- [2] Li, T., Slavcheva, M., Zollhoefer, M., Green, S., Lassner, C., Kim, C., ... & Lv, Z. (2022). Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 5521-5531). 1, 3
- [3] Kerbl, B., Kopanas, G., Leimkühler, T., & Drettakis, G. (2023). 3D Gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4), 139-1. 1
- [4] Raissi, M., Perdikaris, P., & Karniadakis, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378, 686-707. 1
- [5] Li, X., Qiao, Y. L., Chen, P. Y., Jatavallabhula, K. M., Lin, M., Jiang, C., & Gan, C. (2023). Pac-nerf: Physics augmented continuum neural radiance fields for geometry-agnostic system identification. *arXiv preprint arXiv:2303.05512*. 1, 3

- 531 [6] Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R.,  
532 Kalchbrenner, N., Goyal, A., & Bengio, Y. (2021).  
533 Toward causal representation learning. Proceedings of  
534 the IEEE, 109(5), 612-634. 2
- 535 [7] Müller, T., Evans, A., Schied, C., & Keller, A. (2022).  
536 Instant neural graphics primitives with a multireso-  
537 lution hash encoding. ACM transactions on graphics  
538 (TOG), 41(4), 1-15. 3
- 539 [8] Pearl, J. (2009). Causality. Cambridge university  
540 press. 3