

# FROM CIRCUITS TO DYNAMICS: UNDERSTANDING AND STABILIZING FAILURE IN 3D DIFFUSION TRANSFORMERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reliable surface completion from sparse point clouds underpins many applications spanning content creation and robotics. While 3D diffusion transformers attain state-of-the-art results on this task, we uncover that they exhibit a catastrophic mode of failure: arbitrarily small on-surface perturbations to the input point cloud can fracture the output into multiple disconnected pieces – a phenomenon we call *meltdown*. Using activation-patching from mechanistic interpretability, we localize meltdown to a single early denoising cross-attention activation. We find that the singular-value spectrum of this activation provides a scalar proxy: its spectral entropy rises when fragmentation occurs and returns to baseline when patched. Interpreted through diffusion dynamics, we show that this proxy tracks a symmetry-breaking bifurcation of the reverse process. Guided by this insight, we introduce `PowerRemap`, a drop-in, test-time control that stabilizes sparse point-cloud conditioning. On Google Scanned Objects, `PowerRemap` has a stabilization rate of 98.3% for the state-of-the-art diffusion transformer `WALA`. Overall, this work is a case study on how diffusion model behavior can be understood and guided based on mechanistic analysis, linking a circuit-level cross-attention mechanism to diffusion-dynamics accounts of trajectory bifurcations.

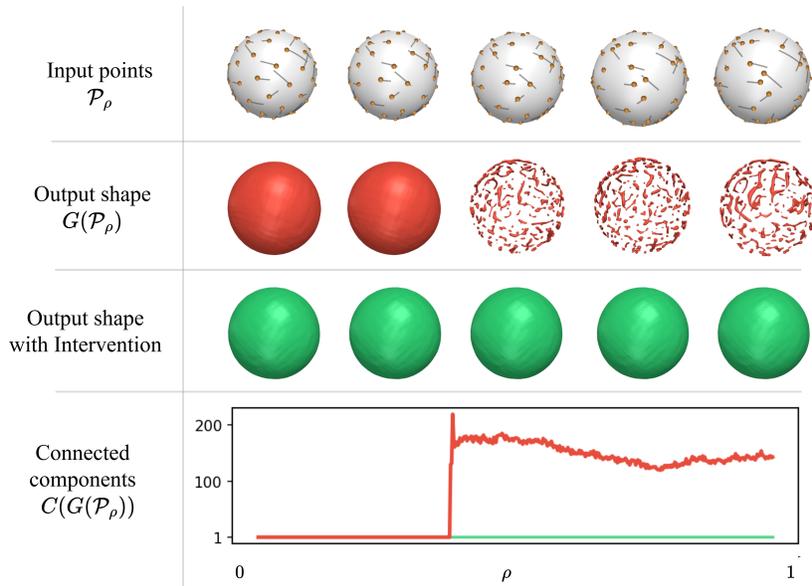


Figure 1: We investigate diffusion transformers on the task of surface reconstruction from sparse point clouds. We find that arbitrarily small on-surface perturbations to a point cloud can turn a shape into a speckle. We call this failure *meltdown* and study it through mechanistic interpretability and diffusion dynamics. Based on this analysis, we propose a test-time intervention, `PowerRemap`, which unlocks diffusion-based surface reconstruction under sparse conditions at test-time.

---

# 1 INTRODUCTION

From virtual content creation to dexterous manipulation, many core applications in vision and robotics hinge on reliable recovery of 3D surfaces from incomplete observations. This problem is called *surface reconstruction* from point clouds (Huang et al., 2022). In real-world applications, the available point clouds are often *sparse*, exacerbating the challenge of surface reconstruction. This sparsity motivates generative priors: *diffusion transformers* attain state-of-the-art results in generative tasks for many modalities (Chen et al., 2024; Sahoo et al., 2024; Jia et al., 2025; Lu et al., 2024b). Recently, they have been introduced to the 3D domain, overcoming challenges in surface reconstruction from sparse point clouds due to learned priors from large-scale datasets (Sanghi et al., 2024; Hui et al., 2024; Wu et al., 2024; Cao et al., 2024).

In this work, we investigate two state-of-the-art diffusion transformers for 3D surface reconstruction, WALA (Sanghi et al., 2024) and MAKE-A-SHAPE (Hui et al., 2024), on the task of surface reconstruction from sparse point clouds. We observe an intriguing *failure mode*: an imperceptible on-surface perturbation to the input point cloud can fracture the output into many disconnected pieces. We call this failure mode *meltdown* and analyze it through two lenses – *mechanistic interpretability* and *diffusion dynamics* – to understand its cause.

First, we employ *activation patching* to test the causal role of activations with respect to meltdown. We find that a single cross-attention activation, early in the denoising process, controls the failure. By investigating the effect of this activation, we find that its spectral entropy constitutes a *proxy* of the observed failure phenomenon. Second, we find that the identified proxy tracks a symmetry-breaking bifurcation of the reverse diffusion process. Based on these insights, we reverse-engineer a test-time intervention, `PowerRemap`, which unlocks diffusion-based surface reconstruction under sparse conditions at test-time.

Our contributions are summarized as follows:

1. **Failure Phenomenon: Meltdown.** We show that the state-of-the-art 3D diffusion transformers WALA (Sanghi et al., 2024) and MAKE-A-SHAPE (Hui et al., 2024) perform surface reconstruction from sparse point clouds in a brittle manner: small on-surface perturbations to the input point cloud can fracture the output into multiple disconnected pieces. We call this failure phenomenon *meltdown*.
2. **Interpretability.** This work provides a case study on how diffusion model behavior can be understood and guided based on mechanistic interpretability. We link a circuit-level cross-attention mechanism to diffusion-dynamics accounts of trajectory bifurcations.
3. **Test-Time Intervention: PowerRemap.** We propose a drop-in, test-time control to stabilize sparse point-cloud conditioning in diffusion transformers. `PowerRemap` averts meltdown in 98.3% of cases on the *Google Scanned Object (GSO)* dataset (Downs et al., 2022) for WALA.

We introduce the failure phenomenon, meltdown, in Section 2. In Section 3, we analyze meltdown from the perspectives of mechanistic interpretability. Section 3.4 introduces our method, `PowerRemap` and presents results for it on the GSO dataset. Finally, we link meltdown to diffusion dynamics in Section 4 and discuss on current limitations in Section 6.

## 2 FAILURE PHENOMENON: MELTDOWN

We study a state-of-the-art diffusion transformer for 3D shape generation, namely WALA (Sanghi et al., 2024). In Appendix B, we show that many observations and insights transfer to another 3D diffusion transformer, namely MAKE-A-SHAPE (Hui et al., 2024). Such models can generate surfaces from several input modalities, including point clouds, thus solving the *surface reconstruction* task: given a set  $\mathcal{P} = \{p_i\}_{i=1}^N \subset \mathcal{S} \subset \mathbb{R}^3$  of  $N$  points sampled from an underlying surface  $\mathcal{S}$ , the model  $G$  should reconstruct a surface consistent with the input and approximating the underlying surface  $G(\mathcal{P}) \approx \mathcal{S}$ . In many real-world scenarios (e.g., fast scene capture),  $N$  can be small, i.e. the point cloud is *sparse*.

As illustrated in Figure 1, we observe that there exist two sparse point clouds  $\mathcal{P}, \mathcal{Q}$  that are close in the input space, but the corresponding outputs differ severely:  $G(\mathcal{P})$  is a connected surface while

$G(\mathcal{Q})$  is a fragmented “speckle” of disconnected pieces. We will refer to this sudden catastrophic fracture as *meltdown*.

To study this failure phenomenon systematically, let us first introduce the topological quantity  $C$  that counts the connected components of the output surface and serves as a quantifiable identifier of the healthy ( $C = 1$ ) versus unhealthy output ( $C > 1$ ). Furthermore, let us consider a running example where the points are sampled from a simple sphere:  $\mathcal{S} = \{x : \|x\|_2 = 1\}$ . This allows us to perform experiments that precisely control for the distribution of the points. Specifically, we fix the random seed and first identify two point clouds of the same size  $N = 400$ :  $\mathcal{P}_0$  which produces a sphere output  $C(G(\mathcal{P}_0)) = 1$  and  $\mathcal{P}_1$  which produces a speckle output  $C(G(\mathcal{P}_1)) \gg 1$  (typically around 100). Using spherical interpolation (geodesics on general surfaces), we can construct a continuous family of point clouds  $\mathcal{P}_\rho \subset \mathcal{S}$ . We sweep  $\rho \in [0, 1]$  and record  $C(\rho) := C(G(\mathcal{P}_\rho))$ .

Figure 1 illustrates the outcome of this experiment. As we sweep  $\rho$  from 0 to 1, we first observe a long plateau of  $C(\rho) = 1$ , followed by a sudden jump to  $C(\rho) \gg 1$  over a very narrow range of  $\rho$ . Refining the steps around this transition, we observe an effectively discontinuous jump in the macroscopic descriptor  $C(\rho)$ .

In Appendix B, we report observing meltdown across different diffusion transformers, i.e., WALA (Sanghi et al., 2024) and MAKE-A-SHAPE (Hui et al., 2024), denoising strategies, i.e., DDIM (Song et al., 2021) and DDPM (Ho et al., 2020b), and Google Scanned Objects (Downs et al., 2022).

### 3 MECHANISTIC INTERPRETATION AND INTERVENTION

After observing and quantifying *meltdown* in WALA, we ask: What is the root cause of this phenomenon? To address this question, we turn to the growing field of *mechanistic interpretability*.

Mechanistic interpretability (Geiger et al., 2021; Wang et al., 2023b; Sharkey et al., 2025) studies the internal mechanisms by which networks generalize, aiming to reverse-engineer representations so that model behavior can be predicted and guided. Core workflows of mechanistic interpretability involve decomposing models into analyzable components, and describing their roles and the flow of information between them. In this regard, *activation patching* (Heimersheim & Nanda, 2024; Zhang & Nanda, 2024) is one of the most prominent techniques. It tests the causal role of activations by swapping them between a *healthy* and a *unhealthy* run and measuring the respective outcome. In our context, we transition between a healthy (sphere) and unhealthy (speckle) run by continuously moving along a meltdown path quantified by our control parameter  $\rho$ . This enables us to transition between a healthy and unhealthy run in a controlled manner, which are ideal conditions for systematic activation patching to identify the root cause of meltdown.

#### 3.1 WALA: DIFFUSION TRANSFORMER

Before we investigate the mechanistic behavior, we briefly summarize the relevant parts of the WALA diffusion transformer. A more detailed description is available in Appendix A, the original work (Sanghi et al., 2024), and the references therein.

**Transformer.** WALA is a latent diffusion model with a point-net encoder  $E$ , U-ViT-style (Hoogeboom et al., 2023) denoising backbone  $B$ , and VQ-VAE decoder (van den Oord et al., 2017)  $D$ . The U-ViT  $B = B^{K-1} \circ \dots \circ B^0$  has  $K = 32$  transformer *blocks*  $B^k$ . The condition  $\mathbf{C} \in \mathbb{R}^{1024 \times 1024}$  enters via both AdaLN modulation Esser et al. (2024) and cross-attention. Denoting by  $\mathbf{Z}^k \in \mathbb{R}^{1728 \times 1152}$  the tokens entering the  $k$ -th block, it computes  $B^k : \mathbf{Z}^k \mapsto \mathbf{Z}^{k+1}$  as a combination of multi-head self-attention SA and cross-attention CA layers (col. 2) with residual connections (col. 3):

$$\overset{\circ}{\mathbf{Z}} = \text{AdaLN}(\overset{\circ}{\mathbf{Z}}^k, \mathbf{C}), \quad \overset{\circ}{\mathbf{Y}} = \text{SA}(\overset{\circ}{\mathbf{Z}}), \quad \overset{\circ}{\mathbf{R}} = \overset{\circ}{\mathbf{Y}} + \overset{\circ}{\mathbf{Z}}^k, \quad (1a)$$

$$\overset{\times}{\mathbf{Z}} = \text{AdaLN}(\overset{\times}{\mathbf{R}}, \mathbf{C}), \quad \overset{\times}{\mathbf{Y}} = \text{CA}(\overset{\times}{\mathbf{Z}}, \mathbf{C}), \quad \overset{\times}{\mathbf{R}} = \overset{\times}{\mathbf{Y}} + \overset{\circ}{\mathbf{R}}, \quad (1b)$$

$$\bar{\mathbf{Z}} = \text{AdaLN}(\bar{\mathbf{R}}, \mathbf{C}), \quad \bar{\mathbf{Y}} = \text{MLP}(\bar{\mathbf{Z}}), \quad \mathbf{Z}^{k+1} = \bar{\mathbf{Y}} + \overset{\times}{\mathbf{R}}. \quad (1c)$$

**Diffusion.** WALA is trained in the standard DDPM (Ho et al., 2020b) framework. At inference, the reverse diffusion maps an initial Gaussian latent  $\mathbf{Z}_T \sim \mathcal{N}(0, I)$  to  $\mathbf{Z}_0$  by iterating over a fixed

---

162 **Algorithm 1** Localizing Meltdown via Activation Patching

---

163 **Require:** Encoder  $E$ ; latent diffusion transformer  $B$ ; decoder  $D$ ; healthy point-cloud  $\mathcal{P}$ ; unhealthy

164 point-cloud  $\mathcal{Q}$

165 1:  $\mathbf{Z}_T^0 \sim \mathcal{N}(0, I)$  ▷ sample initial noise

166 **Record healthy activations:**

167 2:  $\mathbf{C}_{\mathcal{P}} \leftarrow E(\mathcal{P})$  ▷ embed the healthy point-cloud

168 3: **for**  $t = T : 1$  **do** ▷ denoising loop

169 4: **for**  $k = 0 : K - 1$  **do** ▷ block loop

170 5:  $\mathbf{Z}_t^{k+1} \leftarrow B^k(\mathbf{Z}_t^k, \mathbf{C}_{\mathcal{P}})$  and record  $\mathbf{Y}_{k,t}^{\text{healthy}} \leftarrow \overset{\times}{\mathbf{Y}}$

171 6: **end for**

172 7:  $\mathbf{Z}_{t-1}^0 \leftarrow \text{DDIM}(\mathbf{Z}_t^{K-1})$  ▷ discrete denoising update

173 8: **end for**

174 **Patch unhealthy activations:**

175 9:  $\mathbf{C}_{\mathcal{Q}} \leftarrow E(\mathcal{Q})$  ▷ embed the unhealthy point-cloud

176 10: **for**  $t' = T : 1$  **do** ▷ denoising substitution loop

177 11: **for**  $k' = 0 : K - 1$  **do** ▷ block substitution loop

178 12: **for**  $t = T : 1$  **do** ▷ denoising loop

179 13: **for**  $k = 0 : K - 1$  **do** ▷ block loop

180 14:  $\mathbf{Z}_t^{k+1} \leftarrow B^k(\mathbf{Z}_t^k, \mathbf{C}_{\mathcal{Q}})$  but patch  $\overset{\times}{\mathbf{Y}} \leftarrow \mathbf{Y}_{k,t}^{\text{healthy}}$  **if**  $t' = t$  **and**  $k' = k$

181 15: **end for**

182 16:  $\mathbf{Z}_{t-1}^0 \leftarrow \text{DDIM}(\mathbf{Z}_t^{K-1})$  ▷ discrete denoising update

183 17: **end for**

184 18:  $C_{k,t} \leftarrow C(D(\mathbf{Z}_0^{K-1}))$  ▷ decode shape and count connected components after patch

185 19: **end for**

186 20: **end for**

187 21: **return** repair map  $\{C_{k,t}\}_{k=0:K-1, t=1:T}$

---

188

189 schedule of denoising steps  $t \in \mathcal{T} = \{T, \dots, 0\}$ , where at each step the denoiser conditioned on  $\mathbf{C}$

190 updates  $\mathbf{Z}_t \rightarrow \mathbf{Z}_{t-1}$ . At inference-time, we can sample using DDIM (Song et al., 2021) or DDPM

191 (Ho et al., 2020b).

### 192 193 194 3.2 LOCALIZING MELTDOWN VIA ACTIVATION PATCHING

195 In a diffusion transformer, activations span two axes: network depth (*blocks*)  $k \in \mathcal{K} = \{0, \dots, 31\}$

196 and diffusion time (*denoising steps*)  $t \in \mathcal{T} = \{7, \dots, 0\}$ . This depth–time grid  $\mathcal{K} \times \mathcal{T}$  constitutes

197 the search space for the activation patching. In our search, we specifically target the token-wise

198 cross-attention write  $\overset{\times}{\mathbf{Y}} \in \mathbb{R}^{1728 \times 1152}$  which serves as the additive interface to the residual stream

199 through which the context  $\mathbf{C}$  influences the latent tokens. We scan the depth-time grid as given in

200 Algorithm 1 and record healthy activations to patch them for forward passes based on an unhealthy

201 point cloud.

202 We depict the result of the search procedure in Figure 2. We identify a *single* cross-attention write,

203  $\mathbf{Y} \equiv \overset{\times}{\mathbf{Y}}_{4,7}$ , to be responsible for meltdown. The location of this activation is consistent with prior

204 work that finds cross-attention to have the strongest impact on the generated output in the *early*

205 denoising time-steps (Liu et al., 2025). An intuition for this is that at early denoising steps the

206 diffusion model overly relies on the conditioning as there is no other meaningful signal present.

207

208 The goal of mechanistic interpretability is to reverse-engineer internal mechanisms that are human-

209 understandable functions. Since we observe the meltdown as we increase  $\rho$ , we ask whether there is

210 an interpretable function of  $\mathbf{Y}(\rho)$  that allows us to understand the mechanistic cause of meltdown.

### 211 212 213 3.3 INVESTIGATING THE EFFECT OF PATCHING

214 Having localized the failure to a single cross-attention write  $\mathbf{Y}$ , we ask which properties of this write

215 predict meltdown and its rescue under patching. Empirically, we find a single scalar that captures

the effect: the spectral entropy (Powell & Percival, 1979).

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

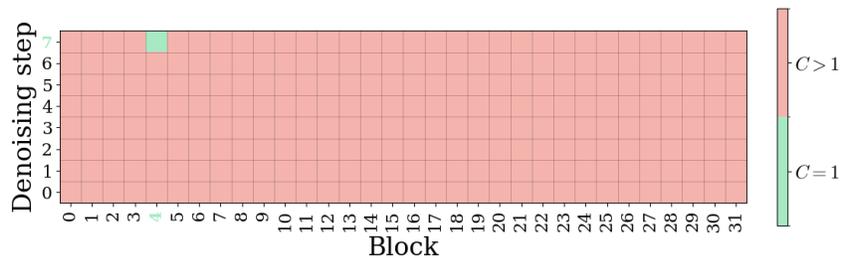


Figure 2: Our search in activation space finds that a single cross-attention write  $\mathbf{Y}_{4,7}$  controls melt-down.

**Definition 1** (Spectral entropy). Let  $\sigma_i$  be the singular values of the matrix  $\mathbf{Y}$ . The spectral entropy of  $\mathbf{Y}$  is

$$H = - \sum_i p_i \log p_i \quad \text{with} \quad p_i = \frac{\sigma_i^2}{\sum_j \sigma_j^2} \quad (2)$$

being the normalized directional energies.

Low  $H$  indicates that the update concentrates its energy in a few directions; high  $H$  indicates a more isotropic, spread-out update. Figure 3a plots connectivity  $C(\rho)$  and Figure 3b plots spectral entropy  $H(\rho)$  along the same path  $\rho$  for three diffusion seeds. In the baseline run,  $H(\rho)$  increases roughly monotonically and smoothly with  $\rho$ ; at some  $\rho = \rho_{\text{melt}}$  the surface fragments ( $C > 1$ ). When we patch  $\mathbf{Y}$  with its healthy value,  $H(\rho)$  remains flat at the healthy level and the surface remains connected for all  $\rho$ . Thus  $H(\rho)$  serves as a simple proxy that tracks failure (baseline) and rescue (patched) at the causal site.

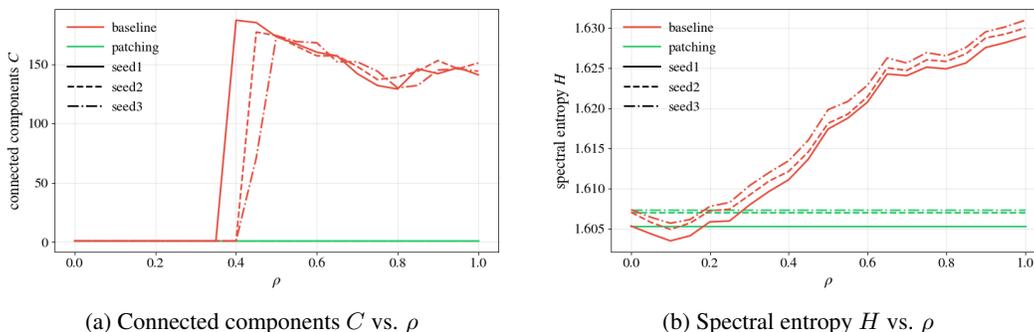


Figure 3: As we move from a healthy to an unhealthy run, we observe that the **baseline** case shows a smooth rise in spectral entropy and a sudden jump in connectivity. **Patching** our  $\mathbf{Y}$  keeps the spectral entropy at healthy levels and preserves connectivity. This behavior is consistent across diffusion seeds.

### 3.4 POWERREMAP: A TEST-TIME SPECTRAL INTERVENTION

Having identified the cross-attention write  $\mathbf{Y}$  and the spectral entropy  $H$  of its singular spectrum as a proxy that tracks failure and rescue, we now ask whether we can *directly* steer this quantity at test-time without access to healthy activations, which in-turn presumes access to a healthy point-cloud which is not the practice case. Our approach is simple: modify  $\mathbf{Y}$  so that  $H$  decreases while leaving the feature-directions (singular vectors of  $\mathbf{Y}$ ) intact.

#### 3.4.1 METHOD

We introduce a minimal intervention that changes only the *magnitudes* of the singular values of  $\mathbf{Y}$  and keeps its singular vectors fixed.

**Definition 2** (`PowerRemap`). Let  $\mathbf{Y} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$  be the SVD of  $\mathbf{Y}$  with singular values  $\sigma_i \geq 0$ . Let  $\sigma_{\max} = \max_i \sigma_i > 0$  and  $\gamma > 1$ . Then define

$$\mathbf{\Sigma}' = \sigma_{\max} \left( \frac{\mathbf{\Sigma}}{\sigma_{\max}} \right)^\gamma, \quad \text{PowerRemap}(\mathbf{Y}) = \mathbf{U} \mathbf{\Sigma}' \mathbf{V}^\top. \quad (3)$$

**Proposition 1** (`PowerRemap` lowers spectral entropy). Let  $H$  and `PowerRemap` be defined as above. For any  $\gamma > 1$ ,

$$H(\text{PowerRemap}(\mathbf{Y})) \leq H(\mathbf{Y}),$$

with equality iff all  $\sigma_i > 0$  are equal.

This compresses the spectrum (smaller singular values shrink faster), which by construction targets the proxy parameter without changing the singular vectors of the features. We refer to the Appendix C for a corresponding proof.

### 3.4.2 EVALUATION AT SCALE

In this section, we assess whether the meltdown phenomenon and the effectiveness of `PowerRemap` generalize across diverse input geometries. To that end, we consider the *Google Scanned Objects (GSO)* (Downs et al., 2022) and *SimJEB* (Whalen et al., 2021) datasets and describe our evaluation protocol in Appendix B.3.

**GSO.** GSO (Downs et al., 2022) is a diverse corpus of 1,030 scanned household objects and was *not* used to train either WALA or MAKE-A-SHAPE. As depicted in Table 1 (top), we identify meltdown in 89.9% out of the 1,030 shapes and find that `PowerRemap` stabilizes failure in 98.3% of cases for the WALA model. We depict qualitative examples of the baseline meltdown and the corresponding `PowerRemap` in Figure 4. For MAKE-A-SHAPE, we identify meltdown in 88.3% out of the 1,030 shapes (cf. Table 4). We evaluate `PowerRemap` using a grid-search over  $\gamma$  for a subset of 130 GSO shapes. We find that our method achieves a stabilization rate of 84.6% (cf. Table 2).

**SimJEB.** SimJEB (Whalen et al., 2021) is a curated benchmark of 381 3D jet-engine bracket CAD models that was *not* included in the training data of either WALA or MAKE-A-SHAPE. As depicted in Table 1, we identify meltdown in 92.4% out of the 381 shapes and find that `PowerRemap` stabilizes failure in 97.7% for the WALA model. For MAKE-A-SHAPE, we identify meltdown in 95.3% out of the 381 shapes (cf. Table 7). We evaluate `PowerRemap` using a grid-search over  $\gamma$  over subset of 30 category-representative SimJEB shapes. We find that our method achieves a stabilization rate of 83.3% (cf. Table 2).

We note that we also experimented with alternative ways of reducing the spectral entropy, for example, reducing the temperature of the cross-attention block responsible for meltdown. However, those attempts did not alleviate the failure mode.

Table 1: Category-wise evaluation of PowerRemap on GSO and SimJEB (WALA) for a global  $\gamma = 100$ . Our method stabilizes failures in 98.3% of cases on GSO and 97.7% on SimJEB.

Category	Shapes	Meltdown occurs [%]	PowerRemap rescues [%]	Avg. areal density
<b>GSO</b>				
Shoe	254	97.2	99.6	–
Consumer goods	248	97.6	99.2	–
Unknown	216	88.4	95.8	–
Other	112	92.9	99.0	–
<b>Total</b>	<b>1030</b>	<b>89.9</b>	<b>98.3</b>	<b>–</b>
<b>SimJEB</b>				
Arch	37	89.2	100.0	1.225e-02
Beam	46	100.0	100.0	1.065e-02
Block	99	87.9	97.7	8.413e-03
Butterfly	43	93.0	95.0	1.282e-02
Flat	147	93.2	97.8	9.785e-03
Other	9	100.0	88.9	1.400e-02
<b>Total</b>	<b>381</b>	<b>92.4</b>	<b>97.7</b>	<b>1.024e-02</b>

Table 2: Category-wise evaluation of PowerRemap on a 130-shape subset of GSO and on a 30-shape subset of SimJEB (MAKE-A-SHAPE), using an adaptive- $\gamma$  strategy. Our method stabilizes failures in 84.6% of cases on the MAS subset and 83.3% on SimJEB (MAKE-A-SHAPE).

Category	Shapes	Meltdown occurs [%]	PowerRemap rescues [%]	Avg. areal density
<b>GSO (subset)</b>				
Consumer goods	60	100.0	90.0	–
Bottles, cans & cups	23	100.0	95.7	–
Unknown	18	100.0	83.3	–
Other	29	100.0	65.5	–
<b>Total</b>	<b>381</b>	<b>100.0</b>	<b>84.6</b>	<b>–</b>
<b>SimJEB (subset)</b>				
Arch	2	100.0	50.0	1.920e-02
Beam	4	100.0	75.0	1.856e-02
Block	9	100.0	100.0	1.371e-02
Butterfly	3	100.0	66.7	1.491e-02
Flat	11	100.0	90.9	1.370e-02
Other	1	100.0	0.0	8.706e-03
<b>Total</b>	<b>30</b>	<b>100.0</b>	<b>83.3</b>	<b>1.488e-02</b>

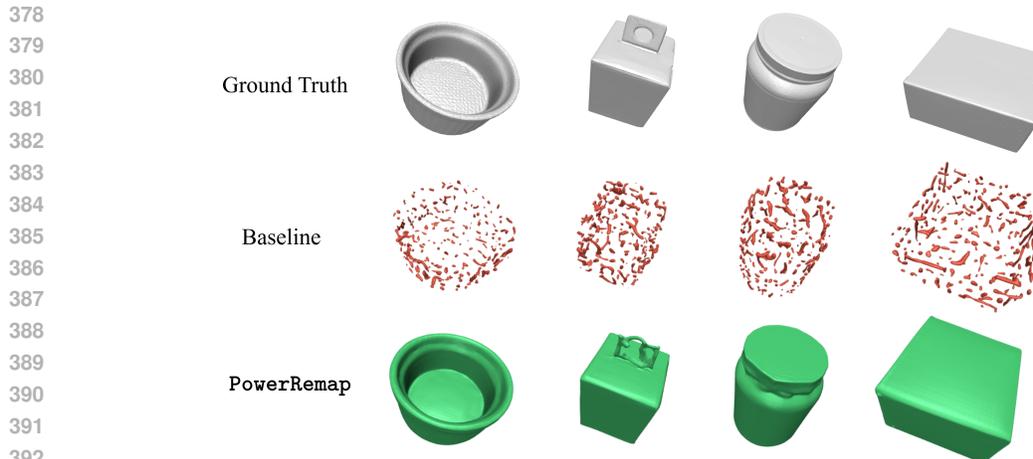
## 4 DIFFUSION DYNAMICS

*Diffusion dynamics* refers to a collection of ideas describing the generative diffusion process using established theory from statistical physics (Raya & Ambrogioni, 2023; Biroli et al., 2024; Yu & Huang, 2025; Ambrogioni, 2025), information theory (Ambrogioni, 2025), information geometry (Chen et al., 2023; Ventura et al., 2025), random-matrix theory (Ventura et al., 2025), and dynamical systems (Ambrogioni, 2025). Key concepts from diffusion dynamics allow us to frame both the observed failure phenomenon and the intervention, ultimately connecting the mechanistic analysis to a theoretically established interpretation of the generative diffusion process.

### 4.1 PRELIMINARIES

We introduce key ideas of diffusion dynamics adapted from Raya & Ambrogioni (2023); Biroli et al. (2024); Ambrogioni (2025). The reverse-time diffusion can be viewed as a noisy gradient flow in a time-dependent potential  $u(\cdot, s)$ :

$$d\mathbf{X}_t = -\nabla_{\mathbf{x}}u(\mathbf{X}_t, s)dt + g(s)d\mathbf{W}_t, \quad u(\mathbf{x}, s) = -g^2(s) \log p(\mathbf{x}, s) + \Phi(\mathbf{x}, s), \quad (4)$$



393  
394  
395  
396  
397

Figure 4: Example results on the *Google Scanned Objects* dataset. We identify meltdown behavior in the WALA diffusion transformer for 89.9% of shapes. Out of these, the `PowerRemap` intervention rescues 98.3%, producing semantically valid outputs.

398  
399  
400  
401  
402  
403  
404  
405  
406

where  $p(\cdot, s)$  is the forward marginal,  $g$  is the noise scale, and  $\Phi(\mathbf{x}, s) = \int_0^{\mathbf{x}} f(\mathbf{z})d\mathbf{z}$  integrates the forward drift  $f$ . The potential  $u$  is essentially a scaled and shifted marginal. The critical points  $x^*$  of this potential  $\nabla u(x^*, s) = 0$  are the *attractors* of the dynamics. Early in the generation ( $t \approx 0$ ), there is a global symmetric basin with a stable central fixed point, and the trajectories exhibit mean-reverting fluctuations around it. As noise decreases, the energy landscape deforms and, at a critical time  $\tau^*$ , the fixed point loses stability and the landscape *bifurcates* into two basins. Such bifurcations repeat until at  $t \approx T$  the potential has many fixed points aligning with the data modes (i.e., the data points under an exact score assumption). These bifurcation times  $\tau^*$  can be interpreted as *decision* times where the sample trajectory is committed to a future attractor basin.

407  
408  
409  
410

Around the degenerate critical point  $x^*(\tau^*)$ , two paths that are nearby for  $t < \tau^*$  may diverge exponentially for  $t > \tau^*$  due to the Lyapunov exponent becoming positive (the smallest eigenvalue of  $\nabla^2 u$  obtained from linearizing the reverse dynamics around the critical point). This can amplify tiny input differences and is the mechanism behind sending trajectories to different attractors.

411  
412  
413  
414  
415

This selection of one among many symmetry-equivalent states is called *spontaneous symmetry breaking*. A canonical example is a ferromagnet: at high temperature ( $t \approx T$ ) spins are disordered, while as  $t \rightarrow 0$  they align. Any magnetization direction is a priori equivalent, yet each realization picks one. The underlying symmetry is visible only in the ensemble over many realizations.

#### 416 4.2 APPLICATION TO MELTDOWN AND INTERVENTION

417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427

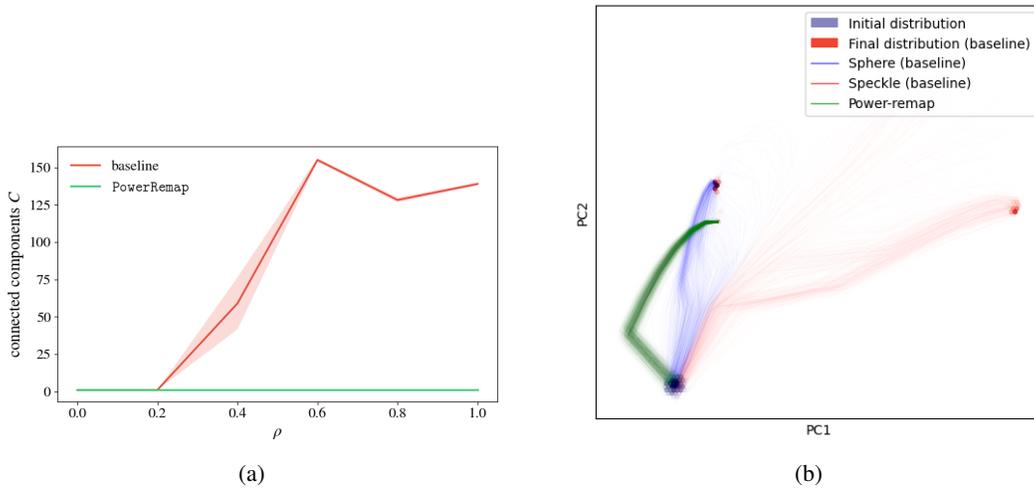
To test the diffusion dynamic perspective of the meltdown phenomenon and the intervention, we perform several experiments predicted by this view. However, we must first introduce conditioning in the above diffusion dynamics view. For a fixed condition  $\mathbf{C}$ , this extension is trivial: simply modify the marginal  $p(\cdot, s) = p(\cdot, s|\mathbf{C})$ . However, a family of conditions, like the univariate interpolation  $\{\mathbf{C}(\rho)|\rho \in [0, 1]\}$ , introduces an additional dependence in the above formalism, and it is not obvious how to analyze the evident bifurcation around  $\mathbf{C}^*$  instead of  $\tau^*$ . Fortunately, since the symmetry breaking originates *locally* around the bifurcation time  $\tau^*$  and point  $x^*$ , a small change in the condition  $d\mathbf{C}$  can be related to a small change in the initial condition  $dx_T$  through the total differential of the reverse path  $\gamma : (t, x_T, \mathbf{C}) \mapsto x_t$ . Qualitatively, this allows us to consider different  $x_T$  for a fixed  $\mathbf{C}$  in place of different  $\mathbf{C}$  for a fixed  $x_T$ .

428  
429  
430  
431

**Ensemble.** Spontaneous symmetry breaking suggests that even though a single trajectory commits to a single attractor (sphere versus speckle), both “symmetric” configurations are visited over an ensemble of random trajectories. We record the trajectories for 100 initial conditions  $x_T \sim \mathcal{N}(0, I)$  over the  $\rho \in [0, 1]$  range and plot the resulting shape connected component distribution in Figure 5a. The extremes  $\rho = 0, 1$  are far from a critical condition, and all trajectories converge to the

432 respective attractors. However, at the intermediate conditions, the ensemble of trajectories visits  
 433 both attractors, with the ratio of fractured shapes increasing steadily with  $\rho$ . In expectation, the  
 434 component curve  $\mathbb{E}_{x_T}[C(\rho)]$  exhibits a smooth behavior, relaxing the discrete jump in  $C(\rho)$  for a  
 435 single  $x_T$ .  
 436

437 **Trajectories.** We can illustrate and qualitatively confirm the explanation for meltdown suggested  
 438 by the diffusion dynamics perspective by considering an intermediate  $\rho = 0.4$  for which the ensem-  
 439 ble visits both attractors. This is visualized in Figure 5b for 1000 trajectories, projected onto a 2D  
 440 linear subspace spanned by the first two principal components of the final distribution  $p(\cdot, 0)$ , which  
 441 shows two major modes and some minor modes. The baseline trajectories are colored blue/red for  
 442 sphere/speckle results. The initial distribution  $p(\cdot, T)$  is a Gaussian from which all the trajectories  
 443 originate. The first denoising step resembles the initial mean-reverting stage of denoising (Biroli  
 444 et al., 2024; Ventura et al., 2025) as the trajectories remain close. The second step marks the sym-  
 445 metry breaking. We can also imagine a “decision boundary” between the blue and red trajectories:  
 446 this is the projected *separatrix* that demarcates the two attractor basins. The intervention alters the  
 447 first step, after which the close bundle of trajectories flows smoothly to a tight minor mode of the  
 448 baseline.  
 449



465 Figure 5: A collection of diffusion trajectories reveals additional insights about the meltdown phe-  
 466 nomenon. (a) In expectation over the initial noise, both the sphere and speckle shapes are produced  
 467 at intermediate conditions, relaxing the sharp meltdown behavior for a fixed initial noise. (b) Latent  
 468 diffusion trajectories projected onto a 2D linear subspace spanned by the first two principal compo-  
 469 nents of the final distribution of the baseline. The `PowerRemap` trajectories in green form a tight  
 470 bundle following a different path that converges to a minor mode of the baseline distribution.  
 471

472  
 473 **Potential.** We calculate the potential similar to the procedure introduced by Raya & Ambro-  
 474 gioni (2023). We select a pair of representative trajectories from each attractor and interpo-  
 475 late between them along a variance-preserving curve  $x_t(\alpha) = \cos(\alpha)x_t^{\text{sphere}} + \sin(\alpha)x_t^{\text{speckle}}$  for  
 476  $\alpha \in [-0.2\pi, 1.2\pi]$ . Figure 6 reveals the two diffusion stages separated by the bifurcation time  
 477  $\tau^* \approx 5$ , where the single potential well flattens and splits into the two attractor basins.  
 478

## 479 5 RELATED WORK

480

481 **Activation patching** In the literature, activation patching has been used to locate mechanistic  
 482 modules of causal interest in many modalities, ranging from language (Wang et al., 2023a; Meng  
 483 et al., 2022; Conmy et al., 2023) to vision-language (Golovanevsky et al., 2025) and audio (Fac-  
 484 chiano et al., 2025). To the best of our knowledge, this is the first work that leverages activation  
 485 patching in the geometry domain. The differentiating factor of geometry over the modalities men-  
 tioned before is that it is objectively measurable at the output (connectedness) as well as at the input

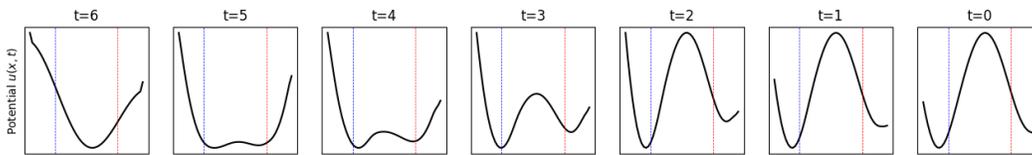


Figure 6: The potential  $u$  (related to the marginal probability via Eq. (4)) reveals the two diffusion stages separated by the bifurcation time  $\tau^* \approx 5$ , where the single potential well flattens and splits into the two attractor basins. The particle’s location just before this early bifurcation commits it to the final attractor and ultimately determines the generated shape. Small perturbations around this time become amplified, giving the appearance of discrete jumps that characterize the observed meltdown.

(point cloud). Using point-cloud inputs as opposed to text prompts allows us to continuously interpolate between a healthy and an unhealthy run. As a result, we can quantify how the causal influence of an identified site changes smoothly along this interpolation.

**Spectral entropy.** Spectral measures are increasingly used in interpretability. Skean et al. (2025) study the matrix entropy of hidden layers, while Yunis et al. (2024) and Lu et al. (2024a) relate weight spectra to generalization. Our analysis provides further evidence that the spectral properties of model internals are indicative of behavior.

## 6 DISCUSSION AND LIMITATIONS

While we identify a simple interpretable proxy for the meltdown (spectral entropy) and also reduce the effect of activation patching to a simple intervention (`PowerRemap`), we cannot argue for the uniqueness of these, for a concrete example, how to select the strength  $\gamma$  of the remap. There exists a phenomenological hierarchy: the generated shape, the diffusion process, and the transformer circuits. While we link the sharp meltdown behavior to bifurcations in the diffusion dynamics, and a specific circuit to the generated shape, a question remains: *why* does a decreased spectral entropy reduce invalid outputs? [From the SVD of the  \$\mathbf{Y}\$  matrix, which stacks outputs of all cross-attention heads, the first singular vector is by construction the direction that maximizes the shared variance across heads and can thus be interpreted as a feature identified by multiple heads. Correspondingly, the relative boost of this dominant shared direction may facilitate a “consensus” which is required for a valid output.](#) A fully satisfactory answer would presume this link to be universal across models, not a mere particularity of the circuit. Preliminary results on the `MAKE-A-SHAPE` model suggest the link between the spectral entropy and the validity of the output transfers to a large degree, but more work is needed to draw strong conclusions. [Overall, we find that the optimal `PowerRemap` strength  \$\gamma\$  is model-dependent.](#)

## 7 CONCLUSION

We identified an intriguing failure mode in the state-of-the-art diffusion transformer `WALA` (Sanghi et al., 2024) for 3D surface reconstruction, namely meltdown: small on-surface perturbations to the input point cloud result in fragmented output shapes. This failure mode can be traced back to activations after a cross-attention branch at early stages of the reverse diffusion process. We found that the spectral entropy of the activations is an indicator of impending meltdown. Based on these insights, we proposed a simple but efficient test-time remedy, `PowerRemap`, that reduces spectral entropy and successfully rescues meltdown in the majority of cases on the widely used *Google Scanned Objects* (Downs et al., 2022) dataset. While our analysis was derived from Sanghi et al. (2024), we observed that meltdown and its mechanistic cause transfer to another 3D diffusion transformer, `MAKE-A-SHAPE` (Hui et al., 2024), suggesting broader relevance beyond a single architecture. In addition to this practical remedy, we established a connection between meltdown and bifurcations in diffusion dynamics, offering a mechanistic lens on how instabilities arise during the reverse process. We believe this perspective not only advances the robustness of 3D surface reconstruction but also opens new avenues for interpretability research in diffusion models.

---

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## REPRODUCIBILITY STATEMENT

We describe all experimental setups in the main text. Appendix B provides exact reproduction protocols and lists the random seeds for every result.

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

---

## REFERENCES

- Luca Ambrogioni. The information dynamics of generative diffusion. *arXiv preprint arXiv:2508.19897*, 2025.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 22669–22679. IEEE, 2023. doi: 10.1109/CVPR52729.2023.02171. URL <https://doi.org/10.1109/CVPR52729.2023.02171>.
- Giulio Biroli, Tony Bonnaire, Valentin de Bortoli, and Marc Mézard. Dynamical regimes of diffusion models. *Nature Communications*, 15(1), November 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-54281-3. URL <http://dx.doi.org/10.1038/s41467-024-54281-3>.
- Wei Cao, Chang Luo, Biao Zhang, Matthias Nießner, and Jiapeng Tang. Motion2vecsets: 4d latent vector set diffusion for non-rigid shape reconstruction and tracking, 2024.
- Defang Chen, Zhenyu Zhou, Jian-Ping Mei, Chunhua Shen, Chun Chen, and Can Wang. A geometric perspective on diffusion models. *arXiv preprint arXiv:2305.19947*, 2023.
- Xinwang Chen, Ning Liu, Yichen Zhu, Feifei Feng, and Jian Tang. Edt: An efficient diffusion transformer framework inspired by human-like sketching. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 134075–134106. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/f1f9962f76581ce8bf38d04c6d6c96b1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/f1f9962f76581ce8bf38d04c6d6c96b1-Paper-Conference.pdf).
- Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 16318–16352. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/34e1dbe95d34d7ebaf99b9bcaeb5b2be-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/34e1dbe95d34d7ebaf99b9bcaeb5b2be-Paper-Conference.pdf).
- Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B. McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items, 2022. URL <https://arxiv.org/abs/2204.11918>.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. URL <https://arxiv.org/abs/2403.03206>.
- Simone Facchiano, Giorgio Strano, Donato Crisostomi, Irene Tallini, Tommaso Mencattini, Fabio Galasso, and Emanuele Rodolà. Activation patching for interpretable steering in music generation, 2025. URL <https://arxiv.org/abs/2504.04479>.
- Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 9574–9586. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/4f5c422f4d49a5a807eda27434231040-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/4f5c422f4d49a5a807eda27434231040-Paper.pdf).
- Michal Golovanevsky, William Rudman, Vedant Palit, Ritambhara Singh, and Carsten Eickhoff. What do vlms notice? a mechanistic interpretability pipeline for gaussian-noise-free text-image corruption and evaluation, 2025. URL <https://arxiv.org/abs/2406.16320>.
- Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching, 2024. URL <https://arxiv.org/abs/2404.15255>.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. URL <https://arxiv.org/abs/2207.12598>.

- 
- 648 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
649 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-*  
650 *ral Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc.,  
651 2020a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf)  
652 [file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf).
- 653 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020b. URL  
654 <https://arxiv.org/abs/2006.11239>.
- 655  
656 Emiel Hooeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: End-to-end diffusion for  
657 high resolution images, 2023. URL <https://arxiv.org/abs/2301.11093>.
- 658  
659 Zhangjin Huang, Yuxin Wen, Zihao Wang, Jinjuan Ren, and Kui Jia. Surface reconstruction from  
660 point clouds: A survey and a benchmark, 2022. URL [https://arxiv.org/abs/2205.](https://arxiv.org/abs/2205.02413)  
661 [02413](https://arxiv.org/abs/2205.02413).
- 662  
663 Ka-Hei Hui, Aditya Sanghi, Arianna Rampini, Kamal Rahimi Malekshan, Zhengzhe Liu, Hooman  
664 Shayani, and Chi-Wing Fu. Make-a-shape: a ten-million-scale 3d shape model, 2024. URL  
665 <https://arxiv.org/abs/2401.11067>.
- 666  
667 Dongya Jia, Zhuo Chen, Jiawei Chen, Chenpeng Du, Jian Wu, Jian Cong, Xiaobin Zhuang, Chumin  
668 Li, Zhen Wei, Yuping Wang, and Yuxuan Wang. Ditar: Diffusion transformer autoregressive  
669 modeling for speech generation. *CoRR*, abs/2502.03930, 2025. doi: 10.48550/ARXIV.2502.  
670 03930. URL <https://doi.org/10.48550/arXiv.2502.03930>.
- 671 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the de-  
672 sign space of diffusion-based generative models. In S. Koyejo, S. Mohamed,  
673 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Infor-*  
674 *mation Processing Systems*, volume 35, pp. 26565–26577. Curran Associates, Inc.,  
675 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf)  
676 [file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf).
- 677 Haozhe Liu, Wentian Zhang, Jinheng Xie, Francesco Faccio, Mengmeng Xu, Tao Xiang,  
678 Mike Zheng Shou, Juan-Manuel Perez-Rua, and Jürgen Schmidhuber. Faster diffusion via  
679 temporal attention decomposition. *Transactions on Machine Learning Research*, 2025. URL  
680 <https://openreview.net/forum?id=xXs2GKXPnH>.
- 681  
682 Haiquan Lu, Yefan Zhou, Shiwei Liu, Zhangyang Wang, Michael W. Mahoney, and Yao-  
683 qing Yang. Alphapruning: Using heavy-tailed self regularization theory for improved  
684 layer-wise pruning of large language models. In A. Globerson, L. Mackey, D. Bel-  
685 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural*  
686 *Information Processing Systems*, volume 37, pp. 9117–9152. Curran Associates, Inc.,  
687 2024a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/10fc83943b4540a9524af6fc67a23fef-Paper-Conference.pdf)  
688 [file/10fc83943b4540a9524af6fc67a23fef-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/10fc83943b4540a9524af6fc67a23fef-Paper-Conference.pdf).
- 689 Haoyu Lu, Guoxing Yang, Nanyi Fei, Yuqi Huo, Zhiwu Lu, Ping Luo, and Mingyu Ding. VDT:  
690 general-purpose video diffusion transformers via mask modeling. In *The Twelfth International*  
691 *Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenRe-  
692 view.net, 2024b. URL <https://openreview.net/forum?id=Un0rgm9f04>.
- 693  
694 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and  
695 editing factual associations in gpt. In S. Koyejo, S. Mohamed, A. Agarwal,  
696 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Process-*  
697 *ing Systems*, volume 35, pp. 17359–17372. Curran Associates, Inc., 2022. URL  
698 [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf)  
699 [6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf).
- 700 G. E. Powell and I. C. Percival. A spectral entropy method for distinguishing regular and irregular  
701 motion of hamiltonian systems. *Journal of Physics A: Mathematical and General*, 12(11):2053–  
2071, 1979. doi: 10.1088/0305-4470/12/11/017.

- 
- 702 Gabriel Raya and Luca Ambrogioni. Spontaneous symmetry breaking in generative diffusion  
703 models. In *Thirty-seventh Conference on Neural Information Processing Systems, 2023*. URL  
704 <https://openreview.net/forum?id=lxGFGMMSV1>.  
705
- 706 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
707 ical image segmentation, 2015. URL <https://arxiv.org/abs/1505.04597>.
- 708 Subham Sekhar Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marro-  
709 quin, Justin T Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and ef-  
710 fective masked diffusion language models. In A. Globerson, L. Mackey, D. Bel-  
711 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-  
712 formation Processing Systems*, volume 37, pp. 130136–130184. Curran Associates, Inc.,  
713 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/  
714 file/eb0b13cc515724ab8015bc978fdde0ad-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/eb0b13cc515724ab8015bc978fdde0ad-Paper-Conference.pdf).
- 715 Aditya Sanghi, Aliasghar Khani, Pradyumna Reddy, Arianna Rampini, Derek Cheung, Ka-  
716 mal Rahimi Malekshan, Kanika Madan, and Hooman Shayani. Wavelet latent diffusion (wala):  
717 Billion-parameter 3d generative model with compact wavelet encodings, 2024. URL <https://arxiv.org/abs/2411.08017>.
- 719 Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas  
720 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, Stella Biderman, Adria  
721 Garriga-Alonso, Arthur Conmy, Neel Nanda, Jessica Rumbelow, Martin Wattenberg, Nandi  
722 Schoots, Joseph Miller, Eric J. Michaud, Stephen Casper, Max Tegmark, William Saunders,  
723 David Bau, Eric Todd, Atticus Geiger, Mor Geva, Jesse Hoogland, Daniel Murfet, and Tom Mc-  
724 Grath. Open problems in mechanistic interpretability, 2025. URL [https://arxiv.org/  
725 abs/2501.16496](https://arxiv.org/abs/2501.16496).
- 726 Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Patel, Jalal Naghiyev, Yann LeCun, and Ravid  
727 Shwartz-Ziv. Layer by layer: Uncovering hidden representations in language models. In *Pro-  
728 ceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceed-  
729 ings of Machine Learning Research*. PMLR, 2025. doi: 10.48550/arXiv.2502.02013. URL  
730 <https://arxiv.org/abs/2502.02013>. Camera-ready; see arXiv:2502.02013.
- 731 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
732 learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Pro-  
733 ceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings  
734 of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL  
735 <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.  
736
- 737 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In  
738 *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria,  
739 May 3-7, 2021*. OpenReview.net, 2021. URL [https://openreview.net/forum?id=  
740 StlgIarCHLP](https://openreview.net/forum?id=StlgIarCHLP).
- 741 Viacheslav Surkov, Chris Wendlar, Antonio Mari, Mikhail Terekhov, Justin Deschenaux, Robert  
742 West, Caglar Gulcehre, and David Bau. One-step is enough: Sparse autoencoders for text-to-  
743 image diffusion models, 2025. URL <https://arxiv.org/abs/2410.22366>.
- 744 Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stene-  
745 torp, Jimmy Lin, and Ferhan Ture. What the daam: Interpreting stable diffusion using cross  
746 attention, 2022. URL <https://arxiv.org/abs/2210.04885>.  
747
- 748 Berk Tinaz, Zalan Fabian, and Mahdi Soltanolkotabi. Emergence and evolution of interpretable  
749 concepts in diffusion models, 2025. URL <https://arxiv.org/abs/2504.15473>.
- 750 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-  
751 ing. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.  
752
- 753 Enrico Ventura, Beatrice Achilli, Gianluigi Silvestri, Carlo Lucibello, and Luca Ambrogioni. Man-  
754 ifolds, random matrices and spectral gaps: The geometric phases of generative diffusion. In *The  
755 Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April  
24-28, 2025*. OpenReview.net, 2025.

---

756 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. In-  
757 terpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The*  
758 *Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda,*  
759 *May 1-5, 2023*. OpenReview.net, 2023a. URL <https://openreview.net/forum?id=NpsVSN6o4u1>.  
760  
761 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. In-  
762 terpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The*  
763 *Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda,*  
764 *May 1-5, 2023*. OpenReview.net, 2023b. URL <https://openreview.net/forum?id=NpsVSN6o4u1>.  
765  
766 E. Whalen, A. Beyene, and C. Mueller. Simjeb: Simulated jet engine bracket dataset. *Computer*  
767 *Graphics Forum*, 40(5):9–17, August 2021. ISSN 1467-8659. doi: 10.1111/cgf.14353. URL  
768 <http://dx.doi.org/10.1111/cgf.14353>.  
769  
770 Shuang Wu, Youtian Lin, Feihu Zhang, Yifei Zeng, Jingxi Xu, Philip Torr, Xun Cao, and Yao Yao.  
771 Direct3d: Scalable image-to-3d generation via 3d latent diffusion transformer. In A. Globerson,  
772 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neu-*  
773 *ral Information Processing Systems*, volume 37, pp. 121859–121881. Curran Associates, Inc.,  
774 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/dc970c91c0a82c6e4cb3c4af7bff5388-Paper-Conference.pdf)  
775 [file/dc970c91c0a82c6e4cb3c4af7bff5388-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/dc970c91c0a82c6e4cb3c4af7bff5388-Paper-Conference.pdf).  
776  
777 Zhendong Yu and Haiping Huang. Nonequilibrium physics of generative diffusion models. *Phys.*  
778 *Rev. E*, 111:014111, Jan 2025. doi: 10.1103/PhysRevE.111.014111. URL [https://link.](https://link.aps.org/doi/10.1103/PhysRevE.111.014111)  
779 [aps.org/doi/10.1103/PhysRevE.111.014111](https://link.aps.org/doi/10.1103/PhysRevE.111.014111).  
780  
781 David Yunis, Kumar Kshitij Patel, Samuel Wheeler, Pedro Savarese, Gal Vardi, Karen Livescu,  
782 Michael Maire, and Matthew R. Walter. Approaching deep learning through the spectral dynamics  
783 of weights, 2024. URL <https://arxiv.org/abs/2408.11804>.  
784  
785 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:  
786 Metrics and methods, 2024. URL <https://arxiv.org/abs/2309.16042>.  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

---

## 810 A BACKGROUND

811  
812 In this section, we provide background information on the diffusion transformers WALA (Sanghi  
813 et al., 2024) and MAKE-A-SHAPE (Hui et al., 2024) along the dimensions *diffusion* (Appendix A.1)  
814 and *transformer* (Appendix A.2).

### 816 A.1 DIFFUSION

817  
818 **Forward transition.** Diffusion generative models synthesize data by inverting a *Markov* chain  
819 that gradually corrupts an observation  $\mathbf{x}_0 \sim p_{\text{data}}$  with Gaussian noise over  $T$  discrete timesteps  
820 (Sohl-Dickstein et al., 2015; Ho et al., 2020a). The forward (noising) transition is

$$821 \quad q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad t = 1, \dots, T, \quad (5)$$

$$822 \quad q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad \text{with} \quad \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s), \quad (6)$$

823  
824 where the variance schedule  $\{\beta_t\}_{t=1}^T \subset (0, 1)$  is chosen so that  $\mathbf{x}_T$  is nearly *i.i.d.*  $\mathcal{N}(0, I)$ . Both  
825 architectures are associated with a cosine variance schedule.

826  
827 **Noise prediction objective.** Instead of directly regressing  $\mathbf{x}_0$ , the denoising neural networks  $\epsilon_\theta$  of  
828 WALA and MAKE-A-SHAPE have been trained to predict the added noise:

$$829 \quad \mathcal{L}_{\text{simple}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[ \left\| \epsilon - \epsilon_\theta \left( \underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{\mathbf{x}_t}, t \right) \right\|_2^2 \right], \quad \epsilon \sim \mathcal{N}(0, I). \quad (7)$$

830  
831 This “ $\epsilon$ -parameterization” empirically stabilizes training and is adopted by nearly all modern models  
832 (Ho et al., 2020a).

833  
834 **Denoising (DDPM).** Given a trained  $\epsilon_\theta$ , the original *Denoising Diffusion Probabilistic Model*  
835 (DDPM) (Ho et al., 2020a) samples via the stochastic reverse transition

$$836 \quad p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N} \left( \mathbf{x}_{t-1}; \underbrace{\frac{1}{\sqrt{1 - \beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)}_{\mu_\theta(\mathbf{x}_t, t)}, \sigma_t^2 \mathbf{I} \right), \quad \sigma_t^2 = \tilde{\beta}_t, \quad (8)$$

837  
838 where  $\tilde{\beta}_t = \beta_t \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$ . Iterating Eq. equation 8 from  $t=T$  to 1 produces  $\mathbf{x}_0$  in  $T$  noisy steps.

839  
840 **Denoising (DDIM).** Song et al. (2021) showed that the same model admits a *deterministic* implicit  
841 sampler (DDIM) obtained by setting the variance term to zero:

$$842 \quad \mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left( \underbrace{\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}}_{\mathbf{x}_0} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(\mathbf{x}_t, t). \quad (9)$$

843  
844 Eq. equation 9 preserves the marginal  $q(\mathbf{x}_{t-1} | \mathbf{x}_0)$ , enabling user-specified *inference* schedules  
845 (e.g.,  $t = T, \dots, 1$  with  $T \gg 1$  for high fidelity or sparse subsets for speed) without retraining.  
846 Crucially, Eqs. equation 8–equation 9 share the same  $\epsilon_\theta$  trained via Eq. equation 7. Hence one can  
847 *train* with the log-likelihood-consistent DDPM objective but *sample* using DDPM or DDIM.

848  
849 **Classifier-free guidance.** Both architectures employ *classifier-free guidance* (CFG) (Ho & Sal-  
850 imans, 2022). CFG biases the denoising direction toward a user condition without requiring an  
851 external classifier. For a current latent  $\mathbf{x}_t$  and the shared noise predictor  $\epsilon_\theta$ , we obtain two esti-  
852 mates at the same step  $t$ : the unconditional prediction  $\epsilon_{\text{uncond}}$  (with the condition omitted) and the  
853 conditional prediction  $\epsilon_{\text{cond}}$  (under the desired condition). We then form a guided estimate

$$854 \quad \tilde{\epsilon} = \epsilon_{\text{uncond}} + s(\epsilon_{\text{cond}} - \epsilon_{\text{uncond}}), \quad s \geq 0,$$

855  
856 and substitute  $\tilde{\epsilon}$  in place of  $\epsilon_\theta(\mathbf{x}_t, t)$  in the DDPM/DDIM updates (Eqs. equation 8–equation 9).  
857 Setting  $s = 1$  at inference-time ignores updates from the unconditional stream, i.e.,  $\tilde{\epsilon} = \epsilon_{\text{cond}}$ .

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

## A.2 TRANSFORMER

The noise-prediction objective equation 7 only specifies what to learn but leaves open how the denoiser  $\epsilon_\theta$  is parameterised. Classical DDPMs adopt a convolutional UNET encoder-decoder (Ronneberger et al., 2015; Ho et al., 2020a), whereas modern large-scale models (e.g. Stable Diffusion, Imagen) replace the convolutional blocks with *Transformer* layers, yielding the **UViT** (“U-shaped Vision Transformer”) backbone that now dominates high-fidelity diffusion systems (Bao et al., 2023; Karras et al., 2022). Both architectures, WALA and MAKE-A-SHAPE, implement a *diffusion generative model* via *Transformer* layers.

### A.2.1 OVERVIEW

Both methods adopt a wavelet-latent diffusion pipeline in which 3D shapes are represented as multi-scale wavelet coefficients and a U-ViT-style denoising backbone (Hoogeboom et al., 2023) is trained in the DDPM (Ho et al., 2020b) framework. The key difference lies in how the wavelet data are fed to the diffusion core.

1. WALA first compresses the full wavelet tree with a convolutional VQ-VAE (stage 1), mapping the diffusible wavelet tree to a latent grid. The latent grid is then modeled by a 32-layer U-ViT (stage 2), where each Transformer layer runs self-attention and cross-attention, totaling 32 cross-attention calls.
2. MAKE-A-SHAPE skips the auto-encoder and instead packs selected wavelet coefficients into a compact grid. The U-ViT backbone then downsamples this tensor to a bottleneck volume. The bottleneck is traversed by a 16-layer U-ViT core—8 self-attention layers immediately followed by 8 cross-attention layers—before up-sampling restores the packed grid.

### A.2.2 CONDITIONING PATHWAY (POINT-CLOUD)

In general, both MAKE-A-SHAPE and WALA share a common pipeline: PointNet encoding followed by aggregation and injecting the resulting latent vectors into the U-ViT generator via (i) affine modulation of normalization layers and (ii) cross-attention.

In particular, MAKE-A-SHAPE injects the conditioning latent vectors into the U-ViT generator at three stages: (1) *concatenation*: the latent vectors are aggregated and concatenated as additional channels of the input noise coefficients, (2) *affine modulation*: the latent vectors are aggregated and subsequently utilized to condition the convolution (down-sampling) and deconvolution (up-sampling) layers via modulating the affine parameters of the group normalization layers, (3) *cross-attention*: each condition latent vector is augmented with an element-wise positional encoding and then fed into a cross-attention module alongside the bottleneck volume.

WALA injects the conditioning latent vectors into the U-ViT generator at two stages: (1) *affine modulation*: the latent vectors are linearly projected via a global projection network and used to modulate the scale and bias parameters of GroupNorm layers in both the ResNet and attention blocks (AdaGN) (Esser et al., 2024), (2) *cross-attention*: each latent vector, augmented with an element-wise positional encoding, is employed as the key and value in cross-attention modules interleaved within each transformer block.

---

## 918 B EXPERIMENTS

919  
920 This section shows that many observations and insights gained through studying the diffusion trans-  
921 former WALA (Sanghi et al., 2024) under DDIM sampling on spheres transfer (i) to the diffusion  
922 transformer MAKE-A-SHAPE (Hui et al., 2024) (ii) sampling under DDPM and (iii) other shapes  
923 (GSO and SimJEB). Additionally, this section details the experimental setup to reproduce our re-  
924 sults. In particular:

- 925 1. GENERAL (B.1): This section provides an overview on our experimental setup.
- 926 2. SPHERE EXPERIMENTS (B.2):
  - 927 (a) GENERAL (B.2.1): This section provides an overview on the setup for the sphere ex-  
928 periments.
  - 929 (b) WALA, DDIM (B.2.2): This section reports the experimental setup for the sphere  
930 experiments in the main text.
  - 931 (c) WALA, DDPM (B.2.3): This sections reports additional results for WALA under  
932 DDPM sampling.
  - 933 (d) MAKE-A-SHAPE, DDIM (B.2.4): This section provides results for MAKE-A-SHAPE  
934 under DDIM sampling.
  - 935 (e) MAKE-A-SHAPE, DDPM B.2.5: This section provides results for MAKE-A-SHAPE  
936 under DDPM sampling.
- 937 3. Datasets: GOOGLE SCANNED OBJECTS (GSO) (B.3) and SIMJEB: This section details  
938 our experiments on GSO (Downs et al., 2022) and SimJEB (Whalen et al., 2021).

### 941 B.1 GENERAL

942  
943 This section provides a general overview on the experimental setup for all results reported in this  
944 work.

945  
946 **Restrict Analysis to Conditional Stream.** As the failure behavior, *meltdown*, is independent of  
947 the unconditional stream, we exclusively investigate the conditional prediction stream. That is, we  
948 set the CFG scale  $s = 1.0$  and restrict our mechanistic analysis (e.g. activations) and diffusion  
949 dynamics analysis (e.g., latents  $x_T$ ) to the conditional stream.

950  
951 **Seeding.** Randomness regarding a diffusion trajectory is controlled globally by seed-  
952 ing Python, NumPy, and PyTorch (`torch.backends.cudnn.deterministic=True`,  
953 `benchmark=False`) so that every evaluation at a given  $\rho$  starts from the same terminal noise  
954  $x_T$ .

### 955 B.2 SPHERE EXPERIMENTS

956  
957 This section (i) provides a detailed account on our setup for the sphere experiments and (ii) reports  
958 additional results for *meltdown* on MAKE-A-SHAPE and DDPM sampling.

#### 960 B.2.1 SETUP

961  
962 We detail the minimal, fully reproducible setup used to produce the sphere experiments for WALA  
963 and MAKE-A-SHAPE. Throughout, the control parameter is  $\rho \in [0, 1]$ , and the phenomenological  
964 order parameter is the number of connected components  $C(\rho)$  in the generated mesh.

965  
966 **Conditioning clouds on the sphere.** We work on  $\mathcal{S} = \{x : \|x\|_2 = 1\}$  and fix  $N = 400$  points for  
967 WALA and  $N = 1200$  for MAKE-A-SHAPE. The base cloud  $\mathcal{P}(0)$  uses a golden-angle (Fibonacci)  
968 sphere distribution:

$$969 g = \pi(3 - \sqrt{5}), \quad i \in \{0, \dots, N - 1\}, \quad y_i = 1 - \frac{2(i + 0.5)}{N}, \quad r_i = \sqrt{1 - y_i^2}, \quad \theta_i = g i,$$
$$970 p_i(0) = (r_i \cos \theta_i, y_i, r_i \sin \theta_i).$$

971

972 A second target cloud  $\mathcal{P}(1)$  is produced by jittering each  $p_i(0)$  with i.i.d. Gaussian noise  $n_i \sim$   
 973  $\mathcal{N}(0, 0.1^2 \mathbf{I}_3)$  and renormalizing to the unit sphere:

974  
 975 
$$\tilde{p}_i = \frac{p_i(0) + n_i}{\|p_i(0) + n_i\|_2}, \quad \mathcal{P}(1) = \{\tilde{p}_i\}_{i=1}^N.$$

976  
 977 We then move each point *along* the surface via per-point spherical linear interpolation (SLERP)  
 978 between corresponding pairs:

979  
 980 
$$p_i(\rho) = \text{slerp}(p_i(0), \tilde{p}_i; \rho) = \frac{\sin((1-\rho)\omega_i)}{\sin \omega_i} p_i(0) + \frac{\sin(\rho\omega_i)}{\sin \omega_i} \tilde{p}_i, \quad \omega_i = \arccos(\langle p_i(0), \tilde{p}_i \rangle).$$

981  
 982 This yields the cloud path  $\mathcal{P}(\rho) = \{p_i(\rho)\}_{i=1}^N$  used throughout.

983  
 984 **Decoding and component counting.** Given  $\mathcal{P}(\rho)$ , we compute a conditioning code via the  
 985 model’s encoder and sample a latent with the diffusion sampler to yield  $G(\mathcal{P}(\rho))$ , i.e. a mesh.  
 986 We report

987 
$$C(\rho) = \text{len}(\text{trimesh.split}(\text{mesh})),$$

988 i.e., the number of connected components in `trimesh`.

989  
 990 **Grid over the control parameter.** We sweep a uniform grid of  $\rho$  values, i.e.,  $\rho \in$   
 991  $\{0, 0.05, \dots, 1.0\}$ .

992  
 993 **Connectivity curve  $C(\rho)$ .** We evaluate  $C(\rho)$  on the uniform  $\rho$  grid. For each  $\rho$ , we reseed the  
 994 RNGs to reproduce the identical terminal noise  $x_T$ . The curve reported is the set

995 
$$\{(\rho, C(\rho))\}_{\rho \in \{0, 0.05, \dots, 1.0\}},$$

996 from which the observed plateau at  $C(\rho) = 1$  and the subsequent jump to  $C(\rho) > 1$  over a narrow  
 997  $\rho$ -interval (a connectivity bifurcation) are directly obtained.

998  
 999 **Spectral entropy curve  $H(\rho)$ .** We evaluate the spectral entropy of the localized cross-attention  
 1000 write on the same uniform control grid  $\rho \in \{0, 0.05, \dots, 1.0\}$  and with identical terminal noise  
 1001 across  $\rho$ .

1002 For each  $\rho$ , we encode the cloud  $\mathcal{P}(\rho)$ , run a single sampling trace, and read out the token-wise  
 1003 cross-attention write at the chosen site,  $\mathbf{Y}(\rho)$ . Let  $\{\sigma_i(\rho)\}_i$  be the singular values of  $\mathbf{Y}(\rho)$  (SVD of  
 1004 the matrix with shape tokens  $\times$  features). We form normalized directional energies

1005  
 1006 
$$p_i(\rho) = \frac{\sigma_i(\rho)^2}{\sum_j \sigma_j(\rho)^2}, \quad H(\rho) = - \sum_i p_i(\rho) \log p_i(\rho),$$

1007 using the natural logarithm. The reported curve is the set

1008  
 1009 
$$\{(\rho, H(\rho))\}_{\rho \in \{0, 0.05, \dots, 1.0\}}.$$

1010  
 1011  
 1012 **B.2.2 WALA, DDIM**

1013  
 1014 **Key hyperparameters.**

1015

1016	Model	ADSKAILab/WaLa-PC-1B
1017	Sampler	DDIM ( $\eta = 0$ )
1018	Diffusion rescale	8 steps ( <code>diffusion_rescale_timestep=8</code> ), i.e., default
1019	CFG weight	1.0 ( <code>scale=1.0</code> ), i.e., we consider only conditional stream
1020	Points per cloud	$N = 400$
1021	Cloud source	Unit sphere, golden-angle placement
1022	Target cloud	Gaussian jitter $\sigma = 0.1$ on $\mathbb{R}^3$ , renormalize to $\mathbb{S}^2$
1023	Interpolation	Per-point SLERP, control $\rho \in [0, 1]$
1024	$\rho$ grid	21 values: 0, 0.05, $\dots$ , 1.0
1024	Seeds	0 – 1000 for all RNG calls
1025	Order parameter	$C(\rho) = \#$ connected components ( <code>trimesh.split</code> )
	Device	cuda (CPU is functionally equivalent but slower)

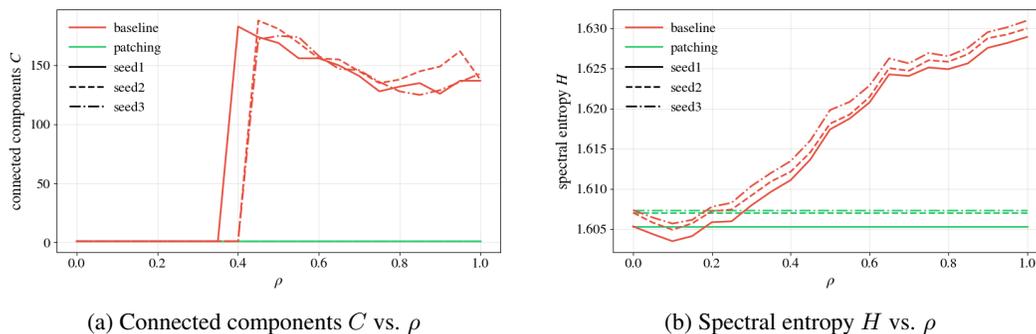
1026 B.2.3 WALA, DDPM

1027  
1028 The activation-patching grid for WALA under DDPM is equivalent to Figure 2, i.e., WALA under  
1029 DDIM. The corresponding for curves can be found in Figure 7.

1030  
1031 **Key hyperparameters.**

1032

1033	Model	ADSKAILab/WaLa-PC-1B
1034	Sampler	DDPM
1035	Diffusion rescale	8 steps (diffusion_rescale_timestep=8)
1036	CFG weight	1.0 (scale=1.0) (we consider only conditional stream)
1037	Points per cloud	$N = 400$
1038	Cloud source	Unit sphere, golden-angle placement
1039	Target cloud	Gaussian jitter $\sigma = 0.1$ on $\mathbb{R}^3$ , renormalize to $\mathbb{S}^2$
1040	Interpolation	Per-point SLERP, control $\rho \in [0, 1]$
1041	$\rho$ grid	21 values: 0, 0.05, ..., 1.0
1042	Seeds	0 for all RNG calls
1043	Order parameter	$C(\rho) = \#$ connected components (trimesh.split)
1044	Device	cuda (CPU is functionally equivalent but slower)



1055  
1056  
1057 Figure 7: [WALA, DDPM]. Our results from WALA under DDIM sampling transfer to WALA under  
1058 DDPM sampling.

1059  
1060 B.2.4 MAKE-A-SHAPE, DDIM

1061  
1062 We report the result for the activation search procedure for MAKE-A-SHAPE under DDIM in Figure  
1063 8. The corresponding curves are depicted in Figure 9.

1064  
1065 **Key hyperparameters.**

1066

1067	Model	ADSKAILab/Make-A-Shape-point-cloud-20m
1068	Sampler	DDIM
1069	Diffusion rescale	100 steps (diffusion_rescale_timestep=8), i.e., default
1070	CFG weight	1.0 (scale=1.0), i.e., we consider only conditional stream
1071	Points per cloud	$N = 1200$
1072	Cloud source	Unit sphere, golden-angle placement
1073	Target cloud	Gaussian jitter $\sigma = 0.1$ on $\mathbb{R}^3$ , renormalize to $\mathbb{S}^2$
1074	Interpolation	Per-point SLERP, control $\rho \in [0, 1]$
1075	$\rho$ grid	21 values: 0, 0.05, ..., 1.0
1076	Seeds	0 for all RNG calls
1077	Order parameter	$C(\rho) = \#$ connected components (trimesh.split)
1078	Device	cuda (CPU is functionally equivalent but slower)

1079

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

### B.2.5 MAKE-A-SHAPE, DDPM

The activation-patching grid for MAKE-A-SHAPE under DDPM is equivalent to Figure 8, i.e., MAKE-A-SHAPE under DDIM. The corresponding curves can be found in Figure 10.

#### Key hyperparameters.

Model	ADSKAILab/Make-A-Shape-point-cloud-20m
Sampler	DDPM
Diffusion rescale	100 steps ( <code>diffusion_rescale_timestep=8</code> ), i.e., default
CFG weight	1.0 ( <code>scale=1.0</code> ) (we consider only conditional stream)
Points per cloud	$N = 1200$
Cloud source	Unit sphere, golden-angle placement
Target cloud	Gaussian jitter $\sigma = 0.1$ on $\mathbb{R}^3$ , renormalize to $\mathbb{S}^2$
Interpolation	Per-point SLERP, control $\rho \in [0, 1]$
$\rho$ grid	21 values: 0, 0.05, ..., 1.0
Seeds	0 for all RNG calls
Order parameter	$C(\rho) = \#$ connected components ( <code>trimesh.split</code> )
Device	cuda (CPU is functionally equivalent but slower)

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

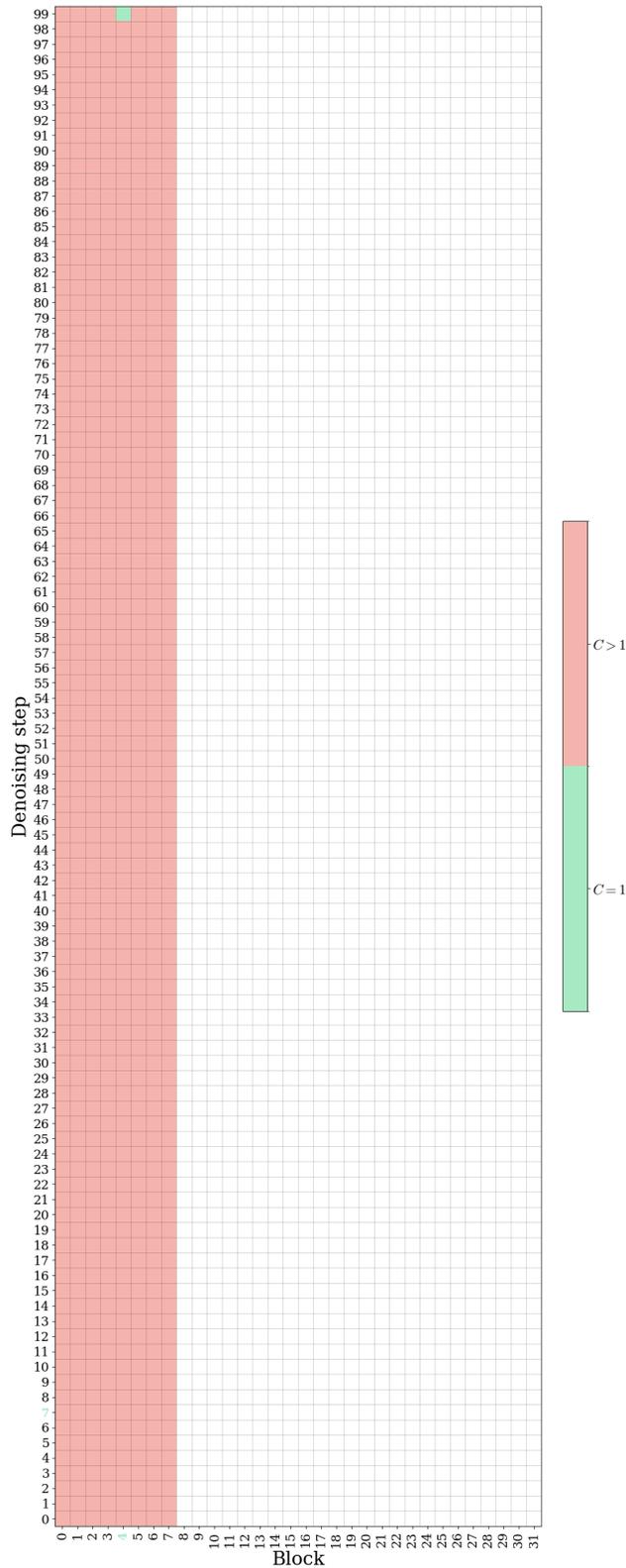


Figure 8: Activation-patching result for MAKE-A-SHAPE. Analogous to our result for WALA, we find an early denoising cross-attention activation that controls meltdown behavior.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

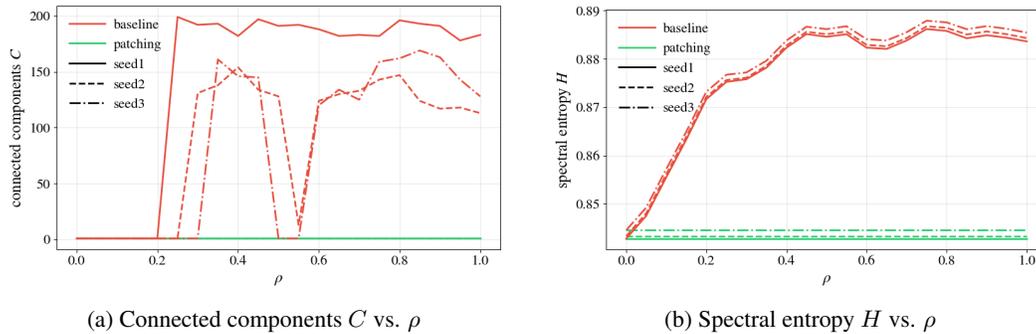


Figure 9: [MAKE-A-SHAPE, DDIM]. Our results from WALA transfer to MAKE-A-SHAPE: As we move from a healthy to an unhealthy run, we observe that the **baseline** case shows a smooth rise in spectral entropy and a sudden jump in connectivity. **Patching** our  $\mathbf{Y}$  keeps the spectral entropy at healthy levels and preserves connectivity. This behavior is consistent across diffusion seeds.

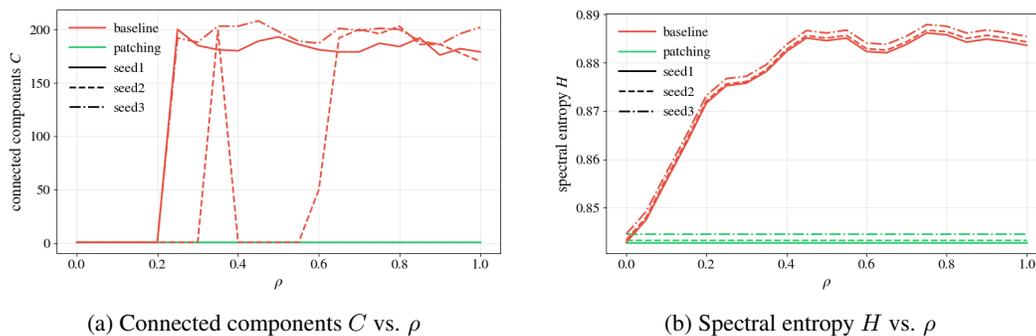


Figure 10: [MAKE-A-SHAPE, DDPM]. Our results from [MAKE-A-SHAPE under DDIM sampling transfer to WALA under DDPM sampling.

---

### 1242 B.3 DATASET EVALUATION

1243  
1244 This section provides (i) evidence that *meltdown* exists across a variety of shapes, i.e., across the  
1245 GSO and SimJEB corpora, and diffusion transformers, i.e., WALA and MAKE-A-SHAPE. Further-  
1246 more it details our setup to evaluate our method `PowerRemap` on GSO (Downs et al., 2022) and  
1247 SimJEB as well as the results of these evaluations.

1248  
1249 **General.** We evaluate `PowerRemap` on the WALA and MAKE-A-SHAPE architectures, using  
1250 DDIM sampling Song et al. (2021). For each object, we load the corresponding mesh as the  
1251 ground-truth surface  $\mathcal{S} \subset \mathbb{R}^3$ , center it and scale it to the unit cube.

1252  
1253 **Find meltdown.** We reuse the notation of §2. Given a mesh  $\mathcal{S}$  and generator  $G$ , we first determine  
1254 a sparse point budget by searching the smallest  $N$  over a grid for which a Poisson-disk sample  
1255  $A = \{a_i\}_{i=1}^N \subset \mathcal{S}$  yields a healthy output  $C(G(A)) = 1$ . We then define a surface-constrained  
1256 *meltdown path* by jittering  $A$  and projecting back to  $\mathcal{S}$  to obtain  $B = \{b_i\}_{i=1}^N$ , and interpolate  
1257 on-manifold

$$1258 \mathcal{P}_\rho = \Pi_{\mathcal{S}}((1 - \rho)A + \rho B), \quad \rho \in [0, 1],$$

1259 where  $\Pi_{\mathcal{S}}$  is nearest-point projection. With a fixed random seed (reseeded before every inference),  
1260 we sweep  $\rho$  on a geometric grid to bracket a jump in connectivity, then refine by bisection to the  
1261 smallest  $\varepsilon$  such that  $C(G(\mathcal{P}_\varepsilon)) \gg 1$ .

---

#### 1262 **Algorithm 2** Adversarial meltdown search

---

1263 **Require:** surface  $\mathcal{S}$ , generator  $G$ , component counter  $C(\cdot)$ , seed  $s = 0$   
1264 1: normalize  $\mathcal{S}$ ; find smallest  $N$  s.t.  $A \sim \text{Poisson}(\mathcal{S}, N)$  gives  $C(G(A))=1$  (reseed  $s$ )  
1265 2:  $B \leftarrow \Pi_{\mathcal{S}}(A + \xi)$  ▷ jitter & project  
1266 3: sweep  $\rho$  on a geometric grid; find  $\rho_{\text{lo}} < \rho_{\text{hi}}$  with  $C(\rho_{\text{lo}})=1, C(\rho_{\text{hi}})>1$   
1267 4:  $\varepsilon \leftarrow \text{bisection}(\rho_{\text{lo}}, \rho_{\text{hi}})$  with reseeding to  $s$   
1268 5: **return**  $(N, \varepsilon, C_0=C(G(A)), C_\varepsilon=C(G(\mathcal{P}_\varepsilon)))$   
1269

---

1270  
1271 **Evaluate `PowerRemap`.** The task of reconstructing a global surface from a sparse point cloud has  
1272 only two possible outcomes: success or failure. Thus, we assess the effectiveness `PowerRemap` on  
1273 GSO by counting the number of times it succeeded in reducing  $C_\varepsilon$  to 1, i.e., turning a speckle into  
1274 a shape. Hence, we treat each shape as a Bernoulli trial under our adversarial search (Algorithm 2).  
1275 Each trail has an outcome  $p \in \{0, 1\}$ , where  $p = 1$  iff the reconstruction meets the criterion  $C_\varepsilon = 1$ ;  
1276 otherwise  $p = 0$ . We first identify baseline failures as those with  $C_{0,\text{baseline}} = 1$  and  $C_{\varepsilon,\text{baseline}} > 1$ .  
1277 We then apply `PowerRemap` only to these failures and count a remedy when  $C_{\varepsilon,\text{PowerRemap}} = 1$ .  
1278

#### 1279 B.3.1 GSO

1280  
1281 This section provides a quantitative evaluation of the *meltdown* phenomenon on the GSO  
1282 dataset (Downs et al., 2022) and assesses the effectiveness of `PowerRemap` as a mitigation strategy.  
1283 SimJEB is a curated benchmark of 381 3D jet-engine bracket CAD models that was *not* included  
1284 in the training data of either WALA or MAKE-A-SHAPE. All results in this section are obtained by  
1285 applying the protocol described in Appendix B.3 to the SimJEB dataset.

1286  
1287 **Evaluate `PowerRemap` for global  $\gamma = 100$ .** For WALA, we found meltdown in 926/1,030  
1288 (89.9%) shapes. Our method `PowerRemap` remedies failure in 910/926 (98.3%) for  $\gamma = 100$ .  
1289 Table 3 depicts the performance of our method across all shape categories. For MAKE-A-SHAPE,  
1290 we found meltdown in 910/1,030 (88.9%) shapes. For a *global*  $\gamma = 100$ , our method `PowerRemap`  
1291 remedies failure in 92/910 (10.1%) cases (cf. Table 4).

1292  
1293  
1294 **Evaluate `PowerRemap` for adaptive- $\gamma$  strategy.** We further investigate whether *per-instance*  
1295 *tuning* of  $\gamma$  can recover meltdown cases that are not resolved by a single global setting. Specifically,  
we consider 130 GSO shapes for which meltdown persists under the global choice  $\gamma = 100$  (cf.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

Table 4) and re-evaluate `PowerRemap` on these shapes over the grid

$$\gamma \in \{1.05, 1.1, 1.15, 1.2, 1.25, 1.3, 1.35, 1.4, 1.5, 2\}.$$

For 110 out of 130 shapes (84.6%), we find at least one value of  $\gamma$  in this range that successfully remedies meltdown. Among the rescued cases, the median effective value of  $\gamma$  is 1.10 with a standard deviation of 0.13.

Category	#Shapes	Meltdown	Meltdown (%)	PR Success	PR@Melt
<b>All</b>	1030	<b>926</b>	<b>89.9</b>	910	<b>98.3</b>
Shoe	254	247	97.2	246	99.6
Consumer Goods	248	242	97.6	240	99.2
Unknown	216	191	88.4	183	95.8
Toys	147	89	60.5	85	95.5
Bottles and Cans and Cups	53	53	100.0	53	100.0
Bag	28	26	92.9	26	100.0
Media Cases	21	21	100.0	21	100.0
Action Figures	17	16	94.1	15	93.8
Board Games	17	16	94.1	16	100.0
Legos	10	6	60.0	6	100.0
Headphones	4	4	100.0	4	100.0
Keyboard	4	4	100.0	4	100.0
Mouse	4	4	100.0	4	100.0
Stuffed Toys	3	3	100.0	3	100.0
Hat	2	2	100.0	2	100.0
Camera	1	1	100.0	1	100.0
Car Seat	1	1	100.0	1	100.0

Table 3: [WALA (Sanghi et al., 2024)] Evaluation of `PowerRemap` on GSO for  $\gamma = 100$ . *Meltdown* counts shapes where the baseline generator yields a single component on the unperturbed input but multiple components after a small, on-surface perturbation ( $C_0=1$ ,  $C_\epsilon>1$ ). *Meltdown (%)* is Meltdown divided by #Shapes. *PR Success* counts—among meltdown shapes only—cases where `PowerRemap` restores a single component ( $C_{\epsilon, \text{PowerRemap}} = 1$ ). *PR@Melt* is the success rate on meltdown shapes (PR Success / Meltdown, shown in %).

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

Category	#Shapes	Meltdown	Meltdown (%)	PR Success	PR@Melt (%)
Global	1030	910	88.3	92	10.1
Consumer Goods	248	235	94.8	26	11.1
Shoe	254	234	92.1	17	7.3
<b>Unknown</b>	216	198	91.7	24	12.1
Toys	147	83	56.5	9	10.8
Bottles and Cans and Cups	53	50	94.3	1	2.0
<b>Bag</b>	28	28	100.0	8	28.6
<b>Media Cases</b>	21	21	100.0	3	14.3
Board Games	17	17	100.0	0	0.0
Action Figures	17	15	88.2	1	6.7
Legos	10	6	60.0	1	16.7
Mouse	4	4	100.0	1	25.0
Headphones	4	3	75.0	0	0.0
Stuffed Toys	3	3	100.0	1	33.3
Keyboard	4	2	50.0	0	0.0
Hat	2	2	100.0	0	0.0
Camera	1	1	100.0	0	0.0
Car Seat	1	1	100.0	0	0.0
Macro avg	–	–	88.4	–	9.9

Table 4: [MAKE-A-SHAPE (Hui et al., 2024)] Evaluation of PowerRemap on GSO for  $\gamma = 100$ . *Meltdown* counts shapes where the baseline generator yields a single component on the unperturbed input but multiple components after a small, on-surface perturbation ( $C_0=1, C_\epsilon>1$ ). *Meltdown (%)* is Meltdown divided by #Shapes. *PR Success* counts—among meltdown shapes only—cases where PowerRemap restores a single component ( $C_{\epsilon, \text{PowerRemap}} = 1$ ). *PR@Melt* is the success rate on meltdown shapes (PR Success / Meltdown, shown in %).

Table 5: Category-wise evaluation of PowerRemap on the 130-shape MAS subset (Top-3 Categories). Our method stabilizes failures in 84.6% of cases.

Category	Shapes	Meltdown occurs [%]	PowerRemap rescues [%]
Consumer goods	60	100.0	90.0
Bottles, cans & cups	23	100.0	95.7
Unknown	18	100.0	83.3
Other	29	100.0	65.5
<b>Total</b>	<b>130</b>	<b>100.0</b>	<b>84.6</b>

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

## B.4 SIMJEB

This section provides a quantitative evaluation of the *meltdown* phenomenon on the SimJEB dataset (Whalen et al., 2021) and assesses the effectiveness of `PowerRemap` as a mitigation strategy. SimJEB is a curated benchmark of 381 3D jet-engine bracket CAD models that was *not* included in the training data of either WALA or MAKE-A-SHAPE. All results in this section are obtained by applying the protocol described in Appendix B.3 to the SimJEB dataset.

**Meltdown** We detect meltdown in 352 out of 381 shapes (92.4%) generated by WALA and in 363 out of 381 shapes (95.3%) generated by MAKE-A-SHAPE. The distribution of meltdown across SimJEB shape categories is reported in Table 6 for WALA and in Table 7 for MAKE-A-SHAPE. These results indicate that meltdown is pervasive across categories and affects both generative pipelines at similar rates.

Table 6: Category-wise evaluation of meltdown on SimJEB (WALA).

Category	Shapes	Meltdowns	Meltdown [%]	Avg. areal density
Arch	37	33	89.2	1.225e-02
Beam	46	46	100.0	1.065e-02
Block	99	87	87.9	8.413e-03
Butterfly	43	40	93.0	1.282e-02
Flat	147	137	93.2	9.785e-03
Other	9	9	100.0	1.400e-02
<b>Total</b>	<b>381</b>	<b>352</b>	<b>92.4</b>	<b>1.024e-02</b>

Table 7: Category-wise evaluation of meltdown on SimJEB (MAKE-A-SHAPE).

Category	Shapes	Meltdowns	Meltdown [%]	Avg. areal density
Arch	37	36	97.3	1.920e-02
Beam	46	46	100.0	1.856e-02
Block	99	87	87.9	1.371e-02
Butterfly	43	42	97.7	1.491e-02
Flat	147	143	97.3	1.370e-02
Other	9	9	100.0	8.706e-03
<b>Total</b>	<b>381</b>	<b>363</b>	<b>95.3</b>	<b>1.488e-02</b>

**Evaluate `PowerRemap`** We evaluate `PowerRemap` on the meltdown cases identified. For WALA, we apply `PowerRemap` with a fixed hyperparameter  $\gamma = 100$  to each meltdown instance and measure the fraction of cases in which the failure is resolved. Under this setting, `PowerRemap` remedies 97.7% of meltdown cases; the category-wise breakdown is reported in Table 8. For MAKE-A-SHAPE, we consider a category-representative subset of 30 SimJEB shapes and evaluate `PowerRemap` over the hyperparameter grid  $\gamma \in \{1.05, 1.1, 1.15, 1.2, 1.25, 1.3, 1.35, 1.4, 1.5, 2\}$ . Across this range, `PowerRemap` successfully remedies meltdown in 83.3% of the evaluated cases with a median effective  $\gamma$  of 1.05 (standard deviation 0.063).

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

Table 8: Category-wise evaluation of PowerRemap on SimJEB (WALA).

Category	Shapes	Meltdowns	PowerRemap rescues	PowerRemap rescues [%]
Arch	37	3	33	100.0
Beam	46	46	46	100.0
Block	99	87	85	97.7
Butterfly	43	40	38	95.0
Flat	147	137	134	97.8
Other	9	9	8	88.9
<b>Total</b>	<b>381</b>	<b>352</b>	<b>344</b>	<b>97.7</b>

Table 9: Category-wise evaluation of PowerRemap on SimJEB (MAKE-A-SHAPE).

Category	Shapes	Meltdowns	PowerRemap rescues	PowerRemap rescues [%]
Arch	2	2	1	50.0
Beam	4	4	3	75.0
Block	9	9	9	100.0
Butterfly	3	3	2	66.7
Flat	11	11	10	90.9
Other	1	1	0	0.0
<b>Total</b>	<b>30</b>	<b>30</b>	<b>25</b>	<b>83.3</b>

---

C POWERREMAP

*Proof of Proposition 1.* Let  $\sigma_i \geq 0$  be the singular values of  $\mathbf{Y}$  and set  $z_i := \sigma_i^2$  (unnormalized directional energies). The baseline normalized spectrum is  $p_i := z_i / \sum_j z_j$ , with entropy  $H(\mathbf{Y}) = -\sum_i p_i \log p_i$ .

By definition,  $\sigma'_i = \sigma_{\max}(\sigma_i / \sigma_{\max})^\gamma$ , hence  $(\sigma'_i)^2 = \sigma_{\max}^{2-2\gamma} \sigma_i^{2\gamma} = \kappa z_i^\gamma$  with a common  $\kappa > 0$  that cancels upon normalization. Therefore the post-intervention normalized spectrum is

$$p_i^{(\gamma)} = \frac{z_i^\gamma}{\sum_j z_j^\gamma}.$$

We note that indices with  $z_i = 0$  keep  $p_i^{(\gamma)} = 0$  for all  $\gamma > 0$  and can be excluded without loss. Write  $s_i := \log z_i$  (for the retained indices). Then

$$p_i^{(\gamma)} = \frac{e^{\gamma s_i}}{\sum_j e^{\gamma s_j}} \quad \text{and} \quad \phi(\gamma) := \log \sum_j e^{\gamma s_j}.$$

Thus  $\log p_i^{(\gamma)} = \gamma s_i - \phi(\gamma)$  and the spectral entropy after the intervention is

$$H(\gamma) := -\sum_i p_i^{(\gamma)} \log p_i^{(\gamma)} = \phi(\gamma) - \gamma \sum_i p_i^{(\gamma)} s_i. \quad (10)$$

Using the standard identity for the softmax measure,

$$\phi'(\gamma) = \frac{\sum_i s_i e^{\gamma s_i}}{\sum_j e^{\gamma s_j}} = \sum_i p_i^{(\gamma)} s_i =: \mathbb{E}_{p^{(\gamma)}}[s],$$

equation 10 simplifies to

$$H(\gamma) = \phi(\gamma) - \gamma \phi'(\gamma).$$

Differentiating once gives

$$H'(\gamma) = \phi'(\gamma) - (\phi'(\gamma) + \gamma \phi''(\gamma)) = -\gamma \phi''(\gamma).$$

It remains to identify  $\phi''(\gamma)$ . A direct computation shows

$$\frac{d}{d\gamma} p_i^{(\gamma)} = p_i^{(\gamma)} (s_i - \mathbb{E}_{p^{(\gamma)}}[s]), \quad \text{hence} \quad \phi''(\gamma) = \frac{d}{d\gamma} \mathbb{E}_{p^{(\gamma)}}[s] = \sum_i s_i \frac{d}{d\gamma} p_i^{(\gamma)} = \text{Var}_{p^{(\gamma)}}(s) \geq 0.$$

Therefore

$$H'(\gamma) = -\gamma \text{Var}_{p^{(\gamma)}}(s) \leq 0,$$

with equality iff all  $\sigma_i > 0$  are equal.

We conclude that  $H(\gamma)$  is nonincreasing in  $\gamma$ ; in particular, for any  $\gamma > 1$ ,

$$H(\text{PowerRemap}(\mathbf{Y})) = H(\gamma) \leq H(1) = H(\mathbf{Y}).$$

□

---

1566 D USE OF LARGE LANGUAGE MODELS (LLMs)  
1567

1568 To a limited extent, we used LLMs to **aid or polish writing**, i.e., in some instances, we used LLMs  
1569 to reformulate sentences.  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

## E DENSITY

We quantify the sparsity of an input point cloud of size  $N$  by its *areal density*

$$\eta = \frac{N}{A_S}, \tag{11}$$

where  $A_S$  denotes the surface area of the underlying surface  $\mathcal{S}$ . Higher values of  $\eta$  correspond to denser samplings of  $\mathcal{S}$ .

In Figure 11, we examine how the prevalence of meltdown depends on  $\eta$  for SimJEB shape 492. For each target areal density, we run Algorithm 2 (parameterized by  $\eta$ ) for 10 independent trials and record how often a meltdown configuration is identified. We observe that meltdown is particularly frequent in the low- $\eta$  regime.

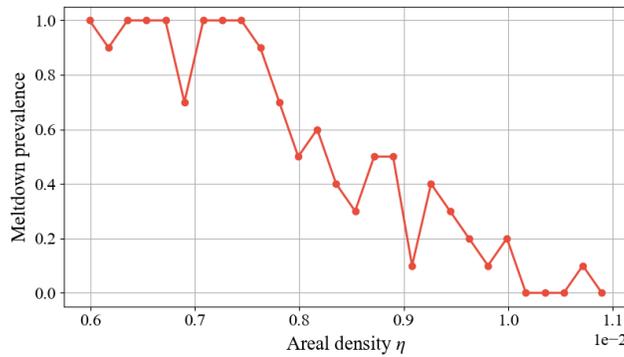


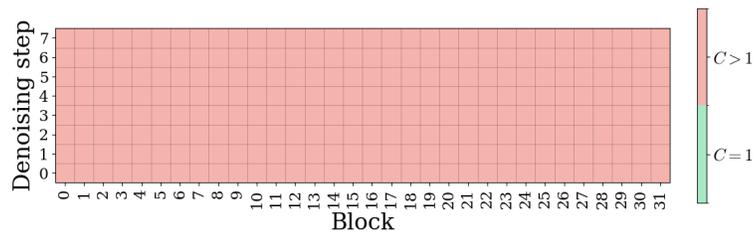
Figure 11: *Incidence of meltdown* on SimJEB shape 492 as a *function of areal density*  $\eta$ . Meltdown events are especially common for low areal densities, underscoring the difficulty of robust surface reconstruction from sparse point clouds.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

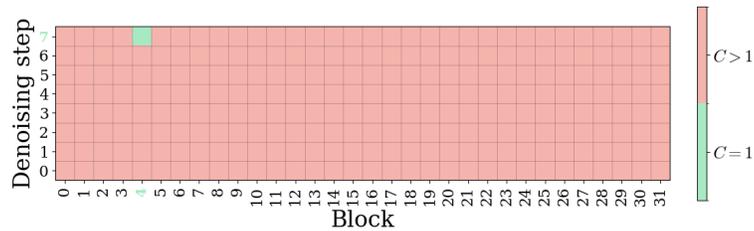
## F ACTIVATION PATCHING

The decision to restrict the search-space for activation-patching in Section 3.2 to cross-attention activations is motivated by prior work (Surkov et al., 2025; Tang et al., 2022; Tinaz et al., 2025). As depicted in Figure 12 extending the search for WALA in <sup>1</sup> to

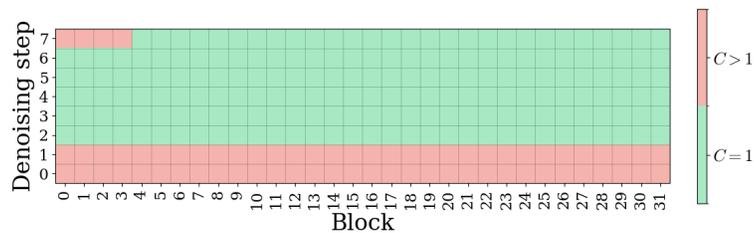
self-attention  $\bar{\mathbf{Y}}^{\circ}$ , residual  $\bar{\mathbf{R}}^{\times}$ , and MLP activations  $\bar{\mathbf{Y}}^{\times}$  reveals that the early cross-attention block identified in Figure 2 is the source of meltdown and other components merely carry the signal downstream. Additionally, the authors found in their experiments cross-attention to be the only interpretable component (cf. Figure 14 and Figure 15b).



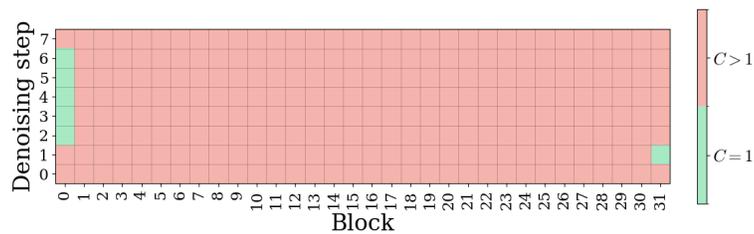
(a) Self-attention activations  $\bar{\mathbf{Y}}^{\circ}$



(b) Cross-attention activations  $\bar{\mathbf{Y}}^{\times}$



(c) Residual-stream activations  $\bar{\mathbf{R}}^{\times}$



(d) MLP activations  $\bar{\mathbf{Y}}^{\bar{\mathbf{Y}}}$

Figure 12: Results for depth-time grid search for self-attention, cross-attention, residual-stream and MLP activations.

<sup>1</sup>Results are averaged over three trials using the setup described in Section 2.

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

## G MORE DATAPOINTS

In this section, we provide additional evidence that the patterns observed in Section 3.2-3.3 generalize when evaluated on more data points and random seeds. Figure 13 and Figure 14 show that the behavior transfers to diverse shapes from the GSO (Downs et al., 2022) and SimJEB (Whalen et al., 2021) corpora as well as diffusion seeds for the WALA model. Figure 15 shows that the average behavior over a population of 150 diffusion seeds for SimJEB shape 492 is consistent with the observations reported in Section 3.3 for the WALA model.

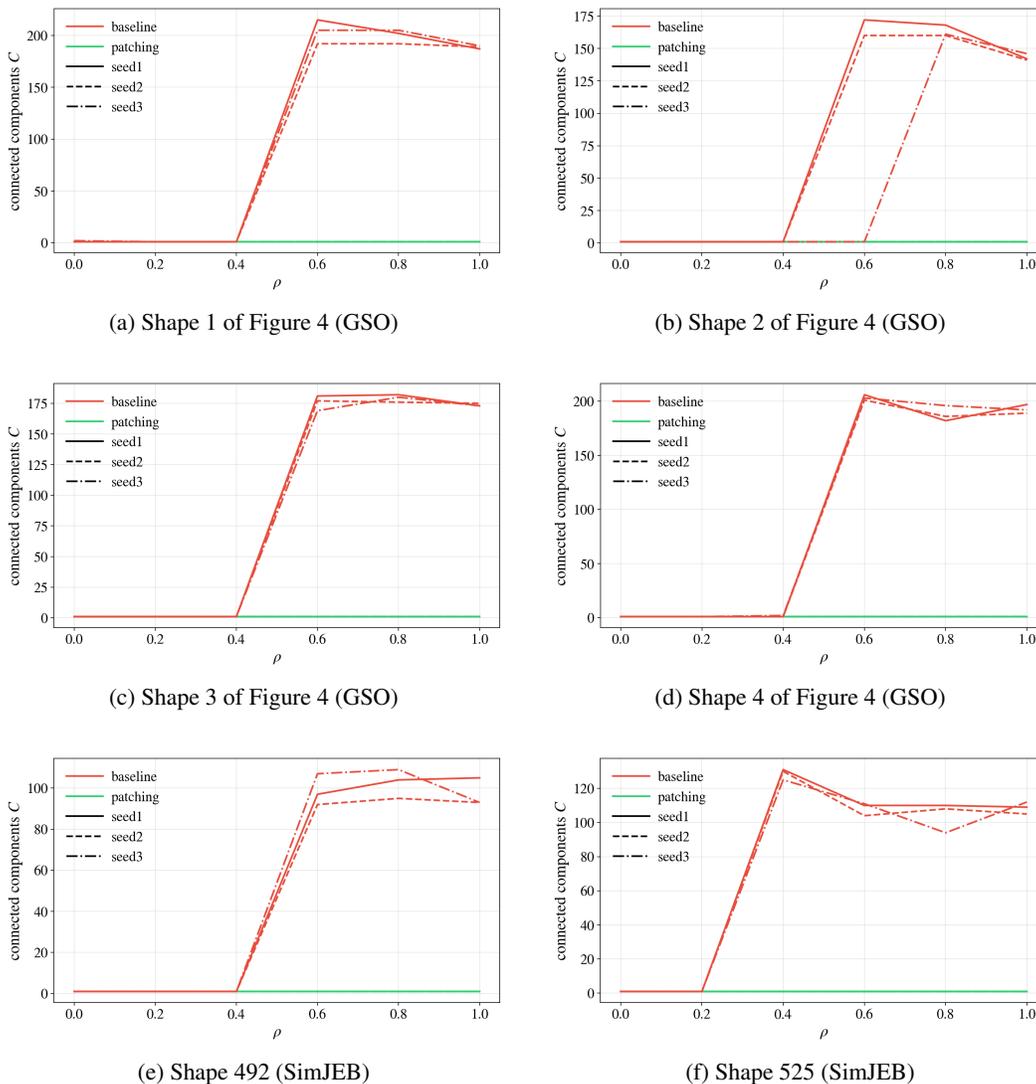


Figure 13: Connected components  $C$  vs.  $\rho$ .

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

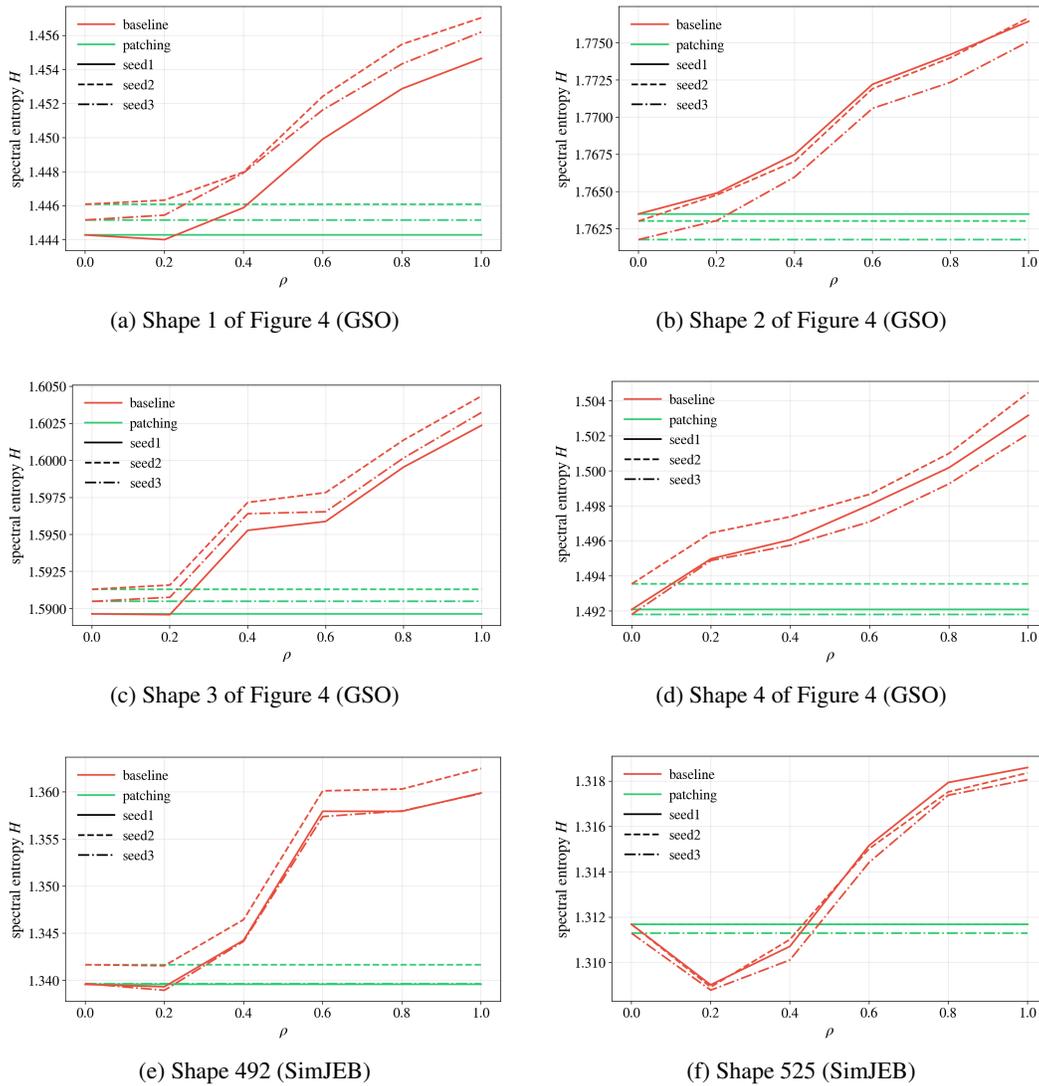


Figure 14: Spectral entropy  $H$  vs.  $\rho$ .

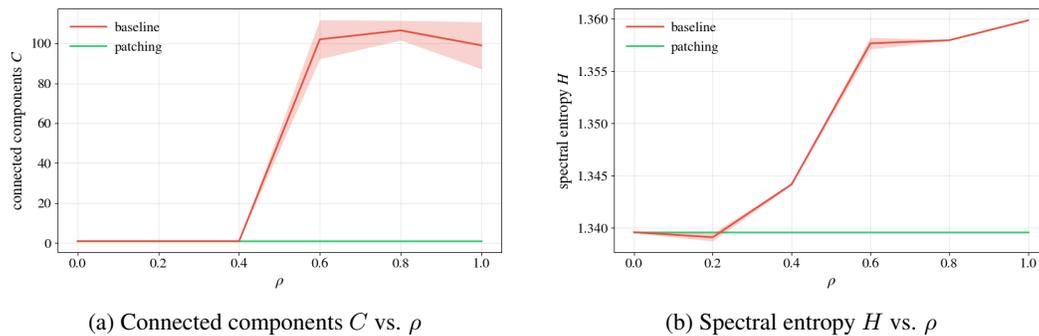


Figure 15: Patterns at population level for SimJEB shape 492, using 150 diffusion seeds.

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

## H ADDITIONAL SPECTRAL METRICS

In this section, we analyze additional spectral metrics to assess their suitability as indicators of meltdown for the WALA model. In particular, Figure 13 reports the effective rank  $r_{\text{eff}} = \exp(H)$  (cf. Definition 1) and Figure 17 the condition number  $\kappa = \sigma_{\text{max}}/\sigma_{\text{min}}$  (cf. Definition 1) as alternatives to spectral entropy for a diverse set of shapes. We observe that the effective rank—which is a monotonic transformation of spectral entropy—provides an equally informative indicator of meltdown. By contrast, the condition number exhibits no apparent correlation with the failure phenomenon, suggesting that it is not a suitable diagnostic metric in this setting.

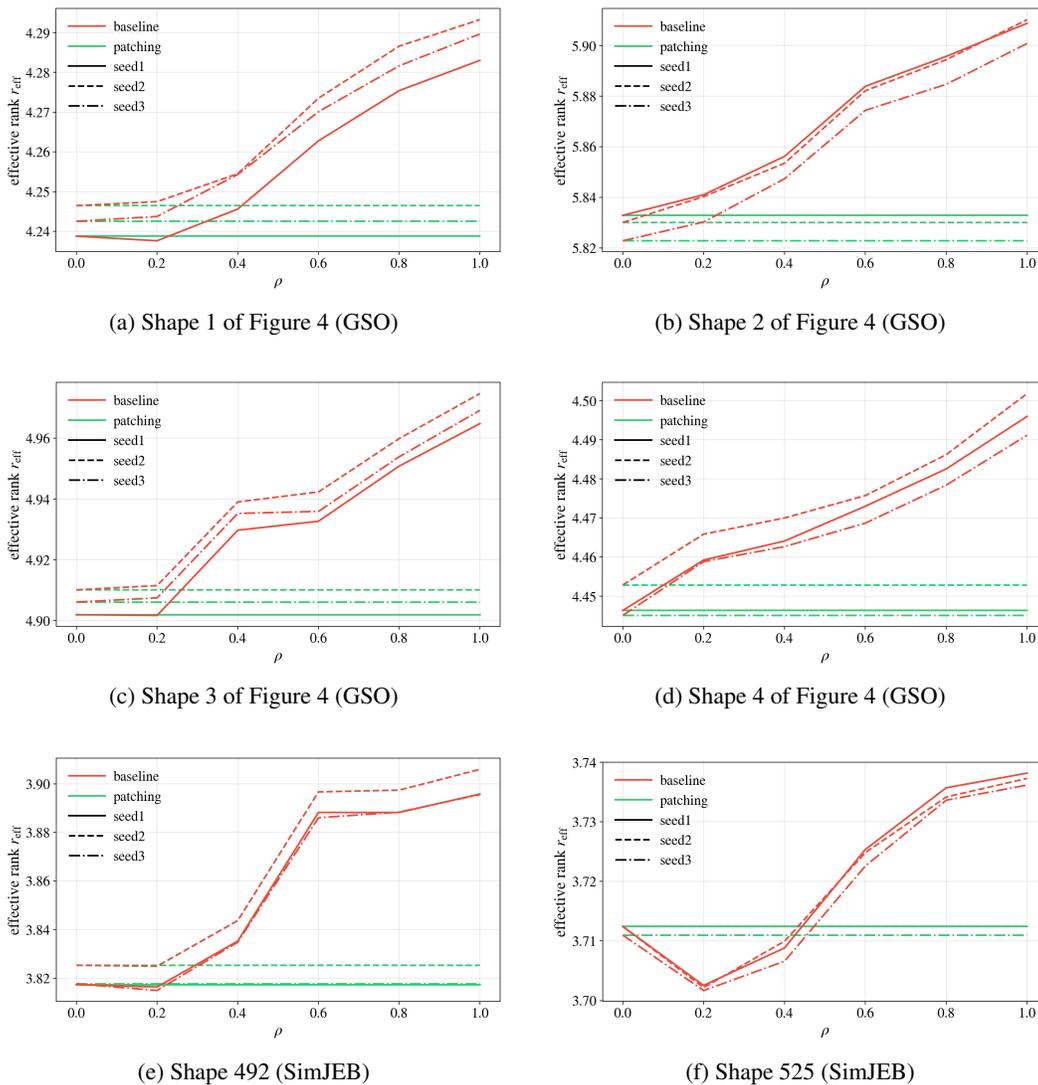
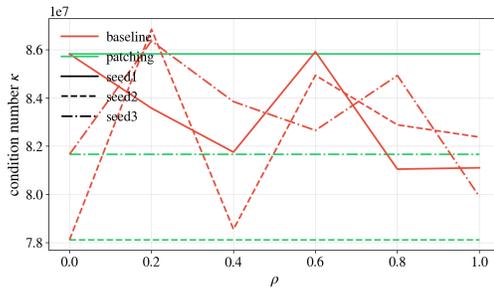
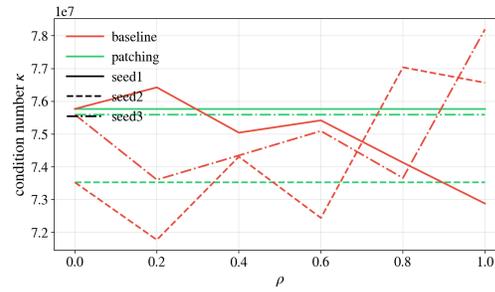


Figure 16: Effective rank  $r_{\text{eff}}$  vs.  $\rho$ .

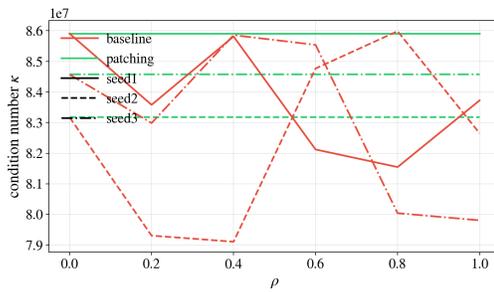
1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943



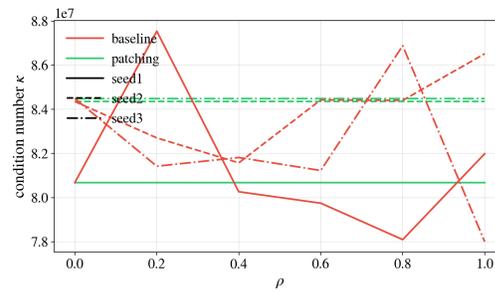
(a) Shape 1 of Figure 4 (GSO)



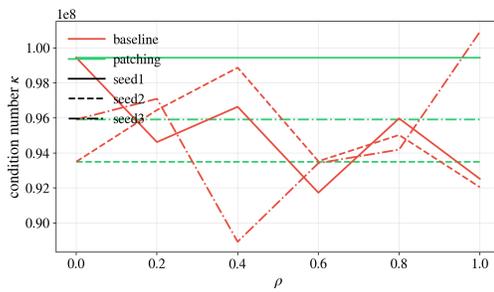
(b) Shape 2 of Figure 4 (GSO)



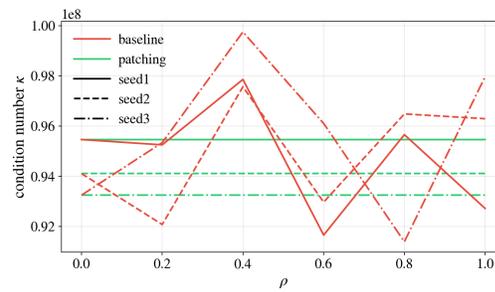
(c) Shape 3 of Figure 4 (GSO)



(d) Shape 4 of Figure 4 (GSO)



(e) Shape 492 (SimJEB)



(f) Shape 525 (SimJEB)

Figure 17: Condition number  $\kappa$  vs.  $\rho$ .

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

## I POWERREMAP EVALUATION

### I.1 MULTIPLE OBJECTS

We further assess whether the meltdown phenomenon and the effectiveness of `PowerRemap` extend beyond single-object inputs. Figure 18 provides a qualitative evaluation on a scene containing multiple objects for the `WALA` model. We observe that meltdown still occurs in this multi-object setting, while using `PowerRemap` reliably suppresses the failure and preserves a plausible reconstruction of all objects in the scene.

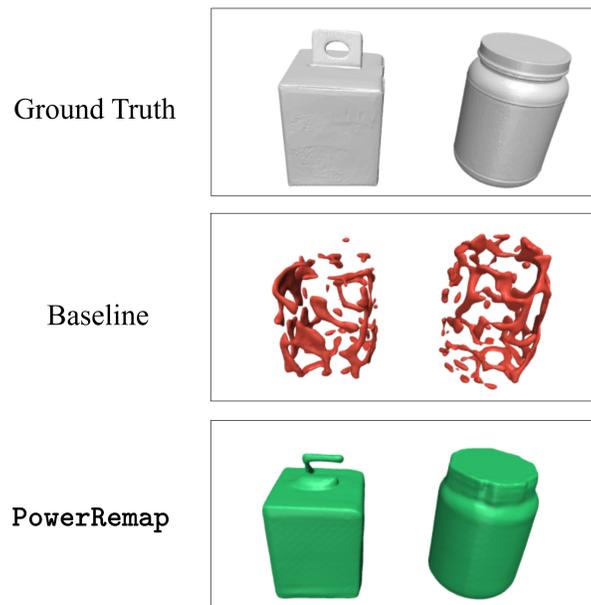


Figure 18: Qualitative evaluation of multi-object inputs . Meltdown persists in scenes with multiple objects, leading to severe degradation of the reconstruction, whereas `PowerRemap` effectively prevents this failure mode and yields a stable reconstruction of all objects.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

## I.2 NON-MELTDOWN CASES

We empirically verify that `PowerRemap` does not interfere with non-meltdown runs. Specifically, we run Algorithm 2 on SimJEB (WALA) shape 492 for 10 independent random seeds and select a configuration with  $C_0 = 1$  (a single connected component). We then apply `PowerRemap` to this configuration over the hyperparameter grid  $\gamma \in \{2, 5, 10, 100\}$ . As can be seen in Table 10, the reconstructed surface retains  $C_0 = 1$  in all cases and we do not observe any change in the final topology. The results in Table 10 indicate that `PowerRemap` is effectively topologically neutral on this non-meltdown instance.

Table 10: Evaluation of `PowerRemap` on a non-meltdown SimJEB shape (shape 492). Across 10 random seeds and  $\gamma \in \{2, 5, 10, 100\}$ , `PowerRemap` preserves the original topology ( $C_0 = 1$ ) in all cases.

Shapes	Seeds	Meltdown occurs [%]	Topology preserved [%]
SimJEB 492	10	0.0	100.0
<b>Total</b>	<b>10</b>	<b>0.0</b>	<b>100.0</b>

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

### I.3 CHOICE OF $\gamma$ TO REMEDY MELTDOWN

We empirically investigate the influence of the `PowerRemap` strength  $\gamma$  on reconstruction connectivity.

**WALA** We run Algorithm 2 on SimJEB (WALA) shape 492 for 10 independent random seeds and select the meltdown configuration with  $C_\epsilon > 1$ . We then apply `PowerRemap` to this configuration over the hyperparameter grid  $\gamma \in \{1, 2, 5, 10, 100\}$ , where  $\gamma = 1$  denotes the identity mapping. As can be seen in 19, our `PowerRemap` method achieves a high success rate for  $\gamma > 2$ .

**MAKE-A-SHAPE** For MAKE-A-SHAPE, we investigate the distribution of `PowerRemap` strengths  $\gamma$  for the 130-shape subset of GSO as discussed in Table 5. We find that  $\gamma$  values around 1.05 are effective to remedy meltdown.

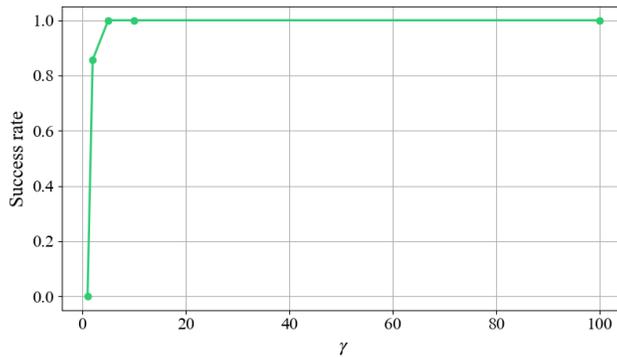


Figure 19: For the `WaLa` model, we find that a `PowerRemap` strength of  $\gamma > 2$  remedies meltdown.

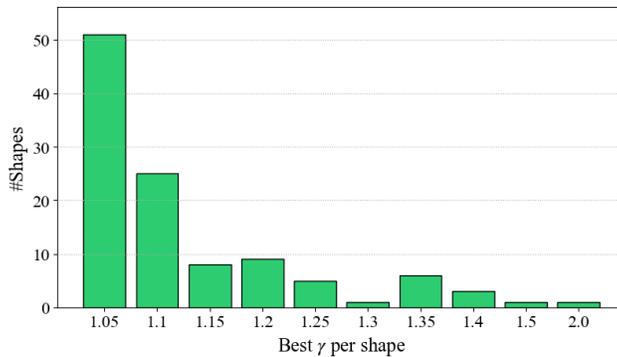


Figure 20: For the 130-shape subset in Table 5 (MAKE-A-SHAPE), we find that  $\gamma$  values around 1.05 are effective to remedy meltdown.