

Factored Generative Models through Mechanism Diversity

Anonymous CVPR submission

Paper ID 35

Abstract

001 *A generative model is factored when each latent dimension*
002 *independently controls one factor of variation: changing*
003 *a single latent predictably changes one semantic attribute*
004 *while leaving the rest unchanged, enabling controllable gen-*
005 *eration, compositional generalization, and reproducible rep-*
006 *resentations. Existing approaches either constrain the latent*
007 *distribution, e.g., requiring it to shift with an auxiliary vari-*
008 *able, or regularize the model, e.g., sparsity, quantization, or*
009 *Hessian penalties, neither of which matches how modern*
010 *conditional generative models actually work: a conditioning*
011 *signal \mathbf{u} reshapes the generator $g(\cdot, \mathbf{u})$ while the latent prior*
012 *stays fixed. We prove an identifiability theorem showing that*
013 *generating mechanism diversity, the natural variation that*
014 *arises when \mathbf{u} sufficiently reshapes g , is sufficient for the*
015 *model to be provably factored, with no parametric assump-*
016 *tion on the latent distribution. To actively enforce this condi-*
017 *tion, we propose Mechanistic Contrastive Learning (MCL),*
018 *a model-agnostic contrastive objective over generator Ja-*
019 *cobians. Empirically, MCL achieves state-of-the-art latent*
020 *concept disentanglement on three benchmarks equipped with*
021 *a latent diffusion model, and improves prediction quality and*
022 *zero-shot cross-task transfer in latent action world models*
023 *with a 1.4B video generative model as the backbone.*

024 1. Introduction

025 What makes a generative model *factored*? Intuitively, a
026 factored model has a latent space where each coordinate con-
027 trols one factor of variation: one latent for color, another for
028 shape, a third for viewpoint, and so on. Such representations
029 underpin controllable generation [50], compositional gener-
030 alization [37], and reproducible science [49]. Yet whether
031 the latent spaces of modern conditional generative models
032 are provably factored remains an open question.

033 Modern conditional generative models share a common
034 architecture: a conditioning signal \mathbf{u} (text prompts, labels,
035 actions) modulates a generator $g(\cdot, \mathbf{u})$ mapping latents \mathbf{z} to
036 observations \mathbf{x} [38, 50, 53]. Prior identifiability theory [27]
037 only covers the opposite regime, requiring \mathbf{u} to shift the

latent distribution $p(\mathbf{z} | \mathbf{u})$ while the generator stays fixed, 038
which is a regime we call *Latent Representation Diversity*, 039
LRD. However, modern generative models keep $p(\mathbf{z})$ fixed 040
(typically Gaussian) and let \mathbf{u} reshape the generator instead, 041
e.g., text-to-image diffusion models [67] inject text into the 042
denoising network, and world models [2] let actions mod- 043
ulate the visual dynamics. In this paper, we refer to this 044
regime as *Generating Mechanism Diversity* (GMD), and 045
our contributions are mainly threefold: (i) **Theory**. We 046
prove that sufficient **GMD** in $g(\cdot, \mathbf{u})$ guarantees a factored la- 047
tent space (Theorem 1), without parametric assumptions on 048
 $p(\mathbf{z})$; this strictly generalizes the classical LRD framework 049
(Proposition 1) and extends to self-supervised and multi- 050
view settings [55, 62]. (ii) **Method**. We introduce **MCL**, 051
which contrasts generator mechanism gradients with condi- 052
tioning signals and uses shuffled conditions as negatives, 053
making mechanism diversity explicit. (iii) **Empirics**. We 054
validate the theory on synthetic data with known factors 055
and instantiate it in image concept disentanglement through 056
MD-VAE and MDDiff, respectively, where MCL improves 057
FactorVAE and DCI score across Cars3D, Shapes3D, and 058
MPI3D while preserving generation quality; we also formu- 059
late the temporal world-model setting, where actions play 060
the role of the conditioning mechanism. MCL generally 061
improves the world model quality, action-following, and 062
task-transferability. Notably, MCL is model-agnostic: the 063
same lightweight head plugs into a VAE, diffusion model, 064
or a 1.4B-parameter pretrained video generative backbone, 065
so it does not trade off the scalability of large-scale genera- 066
tive models. Instead, scalable generative models require not 067
only more pooled data, but also factored structure that sup- 068
ports controllable generation, causal consistency, and related 069
capabilities. 070

071 2. Background and Related Work

From Nonlinear ICA to Factored Generative Model. 072
Nonlinear ICA asks when a generative model $\mathbf{x} = g(\mathbf{z})$ 073
with independent latent factors \mathbf{z} can be inverted to recover 074
those factors from observations alone, and *identifiability the-* 075
ory formalizes when such inversion is unique up to trivial 076
ambiguities. The two are central to factored generative mod- 077
eling because a model that fits $p(\mathbf{x})$ but is not identifiable 078

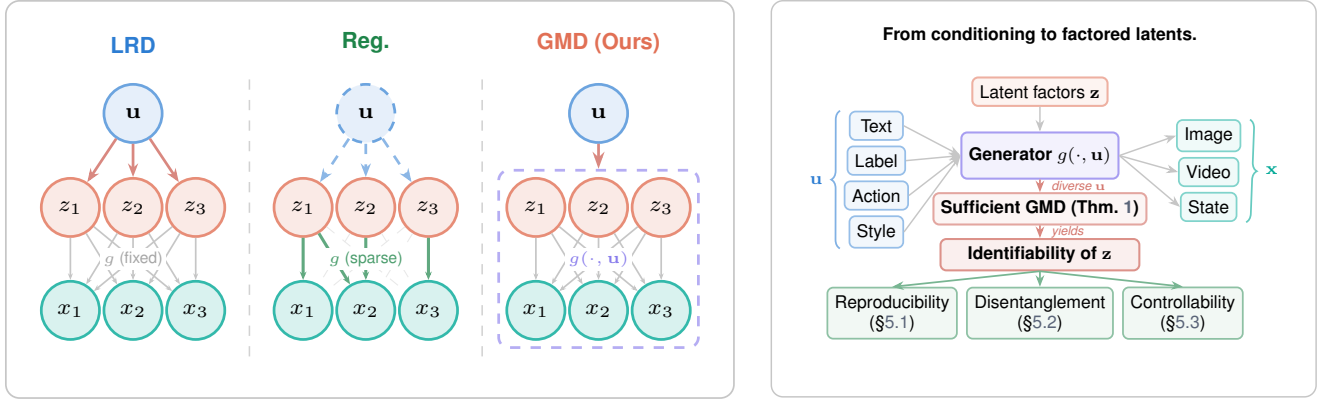


Figure 1. **Factored generative models through mechanism diversity.** Left: **LRD** varies $p(\mathbf{z}|\mathbf{u})$, **Reg.** constrains g , and **GMD** (ours) varies $g(\cdot, \mathbf{u})$. Right: sufficient GMD links conditioning signals, identifiability guarantees, and the practical benefits.

079 can place semantically meaningful structure on arbitrary,
 080 non-reproducible directions of \mathbf{z} [26]. This is the gap that
 081 generative models face: fitting $p(\mathbf{x})$ alone says nothing about
 082 whether the latent code corresponds to disentangled factors
 083 of variation. Resolving it has driven a decade of work, and
 084 we organize the literature below and in Figure 2.

085 **Three principles for factored representations.** Existing
 086 answers to nonlinear ICA and generative-model identifiability
 087 largely follow these two routes: **LRD**, which varies the
 088 latent distribution sufficiently, and **regularization**, which
 089 constrains the model’s degree of freedom. *Beyond these*,
 090 we identify **GMD** as a third principle, varying the generator
 091 itself through conditioning.

092 **(1) Latent distribution diversity (LRD).** The auxiliary \mathbf{u}
 093 creates diversity in the latent distribution $p(\mathbf{z} | \mathbf{u})$
 094 while the generator g stays fixed. Hyvarinen and Morioka
 095 [25] established this for temporal nonstationarity, Hyvari-
 096 nen et al. [27] extended it to arbitrary auxiliaries within an
 097 exponential family, and Khemakhem et al. [29] unified the
 098 framework with VAEs by showing that the ELBO with an
 099 auxiliary-conditioned prior recovers identifiable representa-
 100 tions. Subsequent work removed the exponential-family
 101 assumption via higher-order log-density derivatives [32, 63]
 102 and extended LRD identifiability to multi-view [55], par-
 103 tially observed [62], and additive noisy [17] settings. Apart
 104 from theorems, several practical methods exploit LRD im-
 105 plicitly: contrastive predictive coding leverages temporal
 106 structure [40], and Zimmermann et al. [70] showed that it
 107 can invert the data-generating process to yield identifiable
 108 features. Reizinger et al. [46] proved that cross-entropy clas-
 109 sification yields identifiable representations, and Reizinger
 110 et al. [47] demonstrated that policy diversity produces iden-
 111 tifiable state representations in reinforcement learning. The
 112 shared technical formal condition underlying these results
 113 is summarized below, and our GMD condition (Theorem 1)

replaces the requirement on the latent log-density q_i with a
 condition on the generator Jacobian.

Assumption 1 (Sufficient LRD [27])

For any $\mathbf{z} \in \mathbb{R}^{d_z}$, there exist $2d_z + 1$ values $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(2d_z)}$ such that the $2d_z$ vectors $\mathbf{w}(\mathbf{z}, \mathbf{u}^{(j)}) - \mathbf{w}(\mathbf{z}, \mathbf{u}^{(0)})$, $j = 1, \dots, 2d_z$, are linearly independent, where $\mathbf{w}(\mathbf{z}, \mathbf{u}) = \left(\frac{\partial q_i(\mathbf{z}_i, \mathbf{u})}{\partial z_i}, \frac{\partial^2 q_i(\mathbf{z}_i, \mathbf{u})}{\partial z_i^2} \right)_{i=1}^{d_z}$ and q_i is the log-conditional density of z_i .

(2) Regularization. Instead of requiring diversity in the data, this family constrains the model architecture or objective to favor factored solutions. Independence-based methods penalize statistical dependence among latents: β -VAE [21] upweights the KL term, while FactorVAE [30] and β -TCVAE [5] directly minimize total correlation. Sparsity-based methods constrain the generator: Lachapelle et al. [33] penalize off-diagonal Jacobian entries (mechanism sparsity), Zheng et al. [68] prove identifiability under sparsity on g without auxiliary variables, and Zhang et al. [66] assume a causal graph among latents across distributions. Compression-based methods restrict the latent bottleneck: Tripod [22] combines finite scalar quantization, kernel-based multiinformation minimization, and a normalized Hessian penalty, showing that three complementary regularizers together achieve strong disentanglement in autoencoders. Locatello et al. [37] showed that without such inductive biases, unsupervised disentanglement is impossible.

(3) GMD (ours). Rather than diversifying the latent distribution or constraining the model, GMD induces diversity in the generator itself: a conditioning signal \mathbf{u} (text prompt, class label, control action) reshapes $g(\cdot, \mathbf{u})$ while $p(\mathbf{z})$ stays fixed, and sufficient variation across \mathbf{u} suffices for a factored latent space. Modern conditional generators already operate in this regime ($\mathbf{x} = g(\mathbf{z}, \mathbf{u})$ with $\mathbf{z} \perp \mathbf{u}$): conditional GANs [38, 65], VAEs [20, 53], and diffusion models [9, 41, 50] condition on text or class labels, while

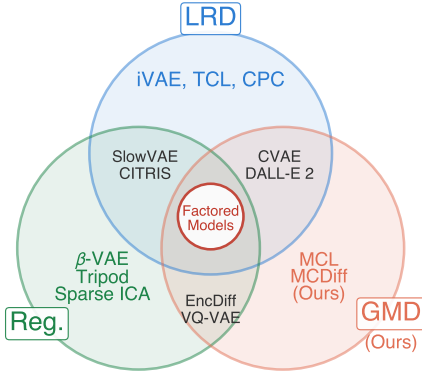


Figure 2. **Three principles for factored models.** **LRD**: diversify the latent distribution. **Regularization**: constrain the model. **GMD** (ours): diversify the generator. $\text{GMD} \supset \text{LRD}$ (Prop. 1).

action-conditioned video world models [2, 4, 16] condition on actions. Existing identifiability theory targets only LRD or regularization, leaving *conditional generative models as the missing principle of factored representations*. We close this gap with a sufficient GMD condition (Theorem 1) that strictly subsumes LRD in representational capacity (Proposition 1), and extend it to two practical regimes: multi-view generation (Proposition 3), where three or more conditioning values play the role of distinct views, and self-supervised temporal data (Proposition 2), where a historical latent serves as \mathbf{u} . We actively encourage this condition via MCL (Section 4), rather than passively learning it as previous conditional generative models.

3. When Is a Generative Model Factored?

In this section, we formalize the conditions under which a conditional generative model is provably factored. The key insight is that, when $g(\cdot, \mathbf{u})$ varies sufficiently across conditioning signals, the model is factored, hence identifiable. We develop this in three steps. We first clarify our motivation in Section 3.1: GMD strictly contains LRD in representational capacity, so the GMD regime is genuinely the larger family worth analyzing. Building on this, Section 3.2 states the sufficient GMD condition and the resulting identifiability theorem. Section 3.3 then relaxes the two main restrictions of the theorem, observed conditioning and pointwise invertibility, to cover the self-supervised and multi-view settings that contemporary generators operate in.

3.1. GMD and LRD Are Not Equivalent in Generating Data

Consider two data-generation families over $\mathbf{u} \in \mathcal{U}$, $\mathbf{z} \in \mathcal{Z}$, $\mathbf{x} \in \mathcal{X}$, with source noise $\epsilon \sim P_{\mathcal{E}}$:

$$\mathbf{z} = f_L(\mathbf{u}, \epsilon), \quad \mathbf{x} = g_L(\mathbf{z}),$$

LRD: \mathbf{u} shifts $p(\mathbf{z}|\mathbf{u})$, g fixed

$$\mathbf{z} = f_G(\epsilon), \quad \mathbf{x} = g_G(\mathbf{u}, \mathbf{z}).$$

GMD: \mathbf{u} reshapes g , $p(\mathbf{z})$ fixed

inducing conditional families $\mathcal{P}_{\text{LRD}} := \{\mathcal{L}(g_L(f_L(\mathbf{u}, \epsilon)) | \mathbf{u}) : \epsilon \sim P_{\mathcal{E}}\}$ and $\mathcal{P}_{\text{GMD}} := \{\mathcal{L}(g_G(\mathbf{u}, f_G(\epsilon)) | \mathbf{u}) : \epsilon \sim P_{\mathcal{E}}\}$. We show that these two forms are not equivalent below.

Proposition 1 (GMD has strictly greater representation capacity than LRD). *Let $\Theta_{\text{LRD}} \subset \mathbb{R}^{d_L}$, $\Theta_{\text{GMD}} \subset \mathbb{R}^{d_G}$ be open parameter spaces with all mappings analytic in their parameters, $P_{\mathcal{E}}$ Lebesgue, and suppose $g_G(\mathbf{u}, \cdot)$ contains a conditioning direction not realisable by any LRD shift. Then $\mathcal{P}_{\text{LRD}} \subsetneq \mathcal{P}_{\text{GMD}}$, and the LRD-realizable subset $\Theta_{\text{LRD-in-GMD}} := \{\theta_G : \exists \theta_L, g_G^{\theta_G}(\mathbf{u}, f_G^{\theta_L}(\epsilon)) = g_L^{\theta_L}(f_L^{\theta_L}(\mathbf{u}, \epsilon)) \forall (\mathbf{u}, \epsilon)\}$ has Lebesgue measure zero in Θ_{GMD} .*

Equivalently, **GMD strictly contains LRD as a generative blueprint in data coverage**, and almost all GMD parameter configurations cannot be reproduced by any fixed-decoder LRD model. The GMD region is the standard regime of modern conditional generators: a fixed distribution prior $p(\mathbf{z})$ paired with a \mathbf{u} -conditioned generator $g(\cdot, \mathbf{u})$.

3.2. Mechanism Diversity Guarantees Factored Latent Spaces

We first formalize what it means for a generative model to have a factored latent space.

Definition 1 (Factored Generative Model). *A model $\mathbf{x} = g(\mathbf{z}, \mathbf{u})$ is factored if for every observationally equivalent model $\mathbf{x} = \hat{g}(\hat{\mathbf{z}}, \mathbf{u})$, there exist a permutation π and invertible functions $\{h_i\}_{i=1}^{d_z}$ such that $\hat{z}_i = h_i(z_{\pi(i)})$.*

Definition 1 is precisely the standard *identifiability* notion of nonlinear ICA [27, 29], recast at the level of the full generative model: a factored model is an identifiable model whose latents are recovered up to permutation and componentwise transformations. Throughout the paper we use *factored and identifiable up to permutation–componentwise interchangeably*. Practically, this property produces disentangled latents, enables controllable generation, generalization, and yields reproducible representations across different training runs [49].

To formalize the GMD condition, let $m'_i(\mathbf{z}, \mathbf{u}) := \partial_{z_i} \log |J_g(\mathbf{z}, \mathbf{u})|$ and $m''_{ii}(\mathbf{z}, \mathbf{u}) := \partial_{z_i}^2 \log |J_g(\mathbf{z}, \mathbf{u})|$ denote the first and diagonal-second derivatives of the log-Jacobian, with analogous $\hat{m}'_i, \hat{m}''_{ii}$ for the learned model.

217 Stacking these derivatives gives the *mechanism-derivative*
218 *vector*

$$219 \quad V(\mathbf{z}, \mathbf{u}) := (m'_1, \dots, m'_{d_z}, m''_1, \dots, m''_{d_z}) \in \mathbb{R}^{2d_z},$$

220 whose contrasts across distinct \mathbf{u} drive the GMD condition
221 and identifiability theorem below.

Assumption 2 (Sufficient GMD)

222 For any $\mathbf{z} \in \mathbb{R}^{d_z}$, there exist $2d_z + 1$ values $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(2d_z)}$ such that the $2d_z$ vectors $V(\mathbf{z}, \mathbf{u}^{(n)}) - V(\mathbf{z}, \mathbf{u}^{(0)})$, $n = 1, \dots, 2d_z$, are linearly independent.

Theorem 1 (Identifiability under GMD)

223 Let $\mathbf{x} = g(\mathbf{z}, \mathbf{u})$ with $\mathbf{z} \perp \mathbf{u}$, and let (\hat{g}, \hat{f}) achieve observational equivalence. Under Assumption 2 and the technical condition $\text{span}\{(m'_i, m''_i)\} \cap \text{span}\{(\hat{m}'_i, \hat{m}''_i)\} = \{0\}$, there exist a permutation π and invertible functions $\{h_i\}_{i=1}^{d_z}$ such that $\hat{z}_i = h_i(z_{\pi(i)})$, i.e. the model is factored in the sense of Definition 1.

224 **Assumption 2 in practice.** Assumption 2 is a statement
225 about the auxiliary variable \mathbf{u} : it asks that the conditioning
226 signal (a text prompt, action, viewpoint, or regime index)
227 reshape the generator’s mechanism in $2d_z + 1$ linearly
228 independent ways. In practice, the auxiliary alphabet must
229 be diverse enough that distinct conditioning values produce
230 genuinely distinct generative mechanisms, a property that
231 modern conditional diffusion and temporal dynamics models
232 are already trained to satisfy.

233 **Proof sketch.** Because $\mathbf{z} \perp \mathbf{u}$, the latent log-density carries
234 no \mathbf{u} -dependence, so all \mathbf{u} -dependent derivative terms in
235 the change-of-variables identity must come from the generator.
236 Change-of-variables relates the learned and true models
237 through $h_{\mathbf{u}} = \hat{g}_{\mathbf{u}}^{-1} \circ g_{\mathbf{u}}$; taking mixed second derivatives
238 and subtracting across conditioning values cancels every
239 \mathbf{u} -independent term, leaving a homogeneous system in the
240 off-diagonal entries of $J_{h_{\mathbf{u}}}$. The linearly independent
241 contrasts of Assumption 2 then force each row of $J_{h_{\mathbf{u}}}$ to have
242 at most one nonzero entry, and invertibility upgrades this to a
243 permutation, yielding the componentwise structure.

244 3.3. Extensions to a Broader Class of Models

245 Theorem 1 as stated requires (i) an explicit conditioning signal
246 \mathbf{u} and (ii) pointwise invertibility of the generator $g(\cdot, \mathbf{u})$.
247 Both assumptions can fail in practice: in self-supervised settings
248 \mathbf{u} may not be directly observed, and modern image- or
249 video-generators are typically not pointwise invertible (each
250 output discards information). We therefore extend the main
251 theorem along these two axes: Proposition 2 replaces the
252 observed \mathbf{u} with a learned latent context, and Proposition 3

relaxes pointwise invertibility to distribution-level invertibility
253 once enough conditioning values are available. Together,
254 these two informal extensions cover the regimes that contemporary
255 GMD generators actually operate in: self-supervised
256 temporal data and many-view generation. Formal statements
257 are collected in Appendix B, and the proofs are given in
258 Appendix C. 259

Proposition 2 (Self-supervised GMD, temporal case; *informal*).
260 When \mathbf{u} is not directly observed, it can be replaced
261 by $\hat{\mathbf{z}}_t$ (the estimated context representation at different time
262 steps). If the resulting model satisfies the conditions of Theorem
263 1, identifiability holds with the same guarantees. 264

Proposition 3 (GMD as a multi-view generative model; *informal*).
265 If the number of distinct conditioning values satisfies $|\{\mathbf{u}^{(n)}\}| \geq 3$,
266 Theorem 1 holds under distribution-level invertibility between
267 different \mathbf{u} , without requiring pointwise invertibility. 268
269

Proposition 2 extends our framework to temporal cases
270 (Figure 3a), which is justified by findings in the identifiable
271 temporal representation learning [25]. Proposition 3 admits a
272 clean multi-view reading (Figure 3b): each $\mathbf{u}^{(n)}$ yields a different
273 “view” $\mathbf{x}^{(n)} = g(\mathbf{z}, \mathbf{u}^{(n)})$ of the shared latent \mathbf{z} , and
274 three or more such views play the role of repeated measurements
275 in the nonparametric identification literature [12, 24].
276 Notably, the threshold of $|\{\mathbf{u}^{(n)}\}| = 3$ smoothly matches
277 classical 3D generation. 278

279 4. How to Actively Build Factored Generative Models 280

When the conditioning signal lacks sufficient variation, a
281 generator in the GMD pattern may fail to diversify its mechanism
282 across conditions, so Assumption 2 is not satisfied even though
283 the architecture is in the GMD regime. We first describe two
284 complementary instantiations of GMD, then present MCL, a training
285 objective that contrasts mechanism-critic gradients across conditioning
286 signals to enforce Assumption 2 explicitly. 287
288

289 4.1. GMD Instantiations 290

Static case: conditional generative models. For controlled
291 synthetic validation (Section 5.1) we construct a *Mechanism-Diverse*
292 VAE (**MD-VAE**): a conditional VAE with encoder $q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{u})$
293 and decoder $p_\theta(\mathbf{x} | \mathbf{z}, \mathbf{u})$, fixed prior $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$, trained by
294 maximizing the conditional ELBO 295

$$296 \quad \log p_\theta(\mathbf{x} | \mathbf{u}) \geq \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{u})} [\log p_\theta(\mathbf{x} | \mathbf{z}, \mathbf{u})] \\ 297 \quad - \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{u}) \| p(\mathbf{z})), \quad (1)$$

denoted $\mathcal{L}_{\text{MD-VAE}}$. Maximizing it reshapes the decoder
298 $g(\cdot, \mathbf{u})$ with \mathbf{u} at fixed prior, instantiating GMD by construction.
299 In the spirit of latent diffusion [50], where a standard 300

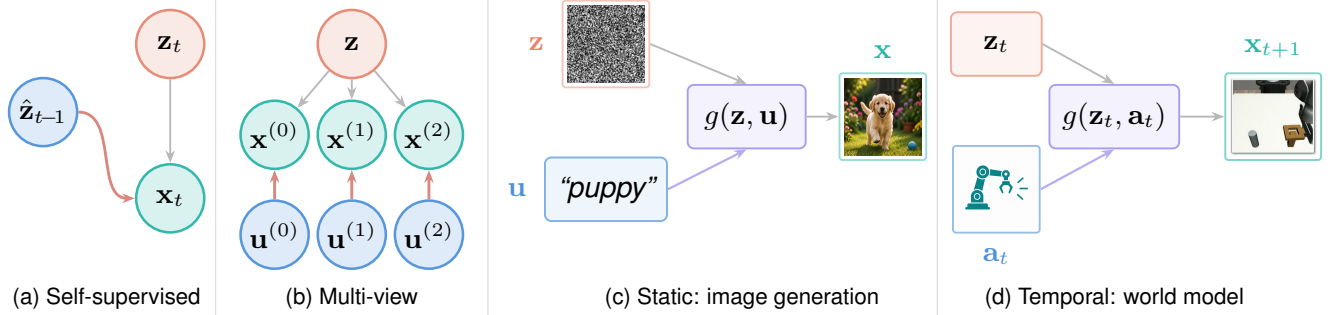


Figure 3. **Extensions and instantiations of GMD.** (a) Self-supervised GMD: a past learned latent $\hat{\mathbf{z}}_{t-1}$ acts as \mathbf{u} . (b) Multi-view GMD: three or more conditioning values $\mathbf{u}^{(n)}$ play the role of repeated measurements of a shared latent \mathbf{z} . (c) Static GMD: a text prompt \mathbf{u} reshapes the denoiser. (d) Temporal GMD: an action \mathbf{a}_t reshapes the dynamics *at the pixel level*, in contrast to the latent-state-level dynamics of traditional action-conditioned world models.

301 VAE encoder feeds a diffusion model in latent space, our
 302 **MDDiff** (Section 5.2) plugs MCL into the same role and runs
 303 \mathbf{u} -conditioned diffusion on the resulting mechanism-diverse
 304 latents.

305 **Temporal case: state-space models.** The conditioning
 306 signal is naturally provided by a history:

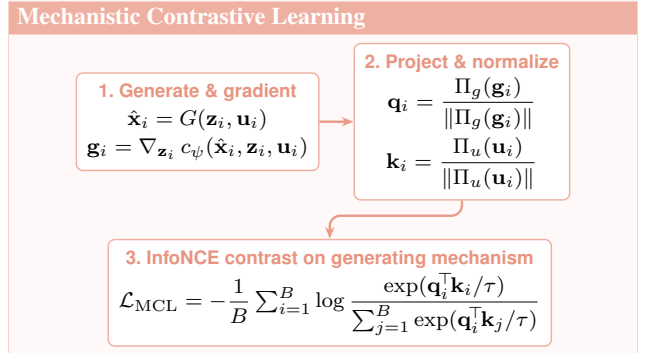
$$307 \quad \mathbf{x}_t = g(\mathbf{x}_{t-H:t-1}, \mathbf{u}_{t-H:t-1}),$$

308 where $\mathbf{u}_{t-H:t-1}$ can be actions in action-conditioned world
 309 models [16, 18], learned latent states in the self-supervised
 310 learning, *e.g.*, time-series data, or any domain index. Sufficient
 311 diversity across the temporal dynamics then factors the
 312 latent state.

313 4.2. Mechanistic Contrastive Learning

314 Beyond a controlled instantiation tied to the variational
 315 ELBO or traditional generative models, to make GMD us-
 316 able across realistic conditional generators in an active way,
 317 we now derive a model-agnostic training objective whose
 318 minimizer drives the generator Jacobian $\partial_{\mathbf{z}}G(\mathbf{z}, \mathbf{u})$ to vary
 319 with \mathbf{u} , which is precisely the linear-independence content
 320 of Assumption 2.

321 MCL augments any differentiable conditional generator
 322 G with three lightweight modules trained jointly with it:
 323 a *mechanism critic* $c_{\psi} : \mathcal{X} \times \mathcal{Z} \times \mathcal{U} \rightarrow \mathbb{R}$ that scores
 324 how strongly the realized mechanism $G(\cdot, \mathbf{u})$ depends on \mathbf{z} ,
 325 and two L_2 -normalized projection heads Π_g, Π_u mapping
 326 mechanism-gradient features and conditioning features into
 327 a shared embedding space. For each in-batch sample $(\mathbf{z}_i, \mathbf{u}_i)$
 328 we generate $\hat{\mathbf{x}}_i = G(\mathbf{z}_i, \mathbf{u}_i)$, score it with the critic, and read
 329 off the *mechanism gradient* $\mathbf{g}_i = \nabla_{\mathbf{z}_i} c_{\psi}(\hat{\mathbf{x}}_i, \mathbf{z}_i, \mathbf{u}_i)$, *i.e.* the
 330 gradient of the critic through the realized generator. The
 331 gradient and conditioning embeddings are then contrasted
 332 with InfoNCE. The three stages are shown below.



333 The full training loss is $\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_{\text{MCL}} \mathcal{L}_{\text{MCL}}$, where \mathcal{L}_{rec}
 334 is the host generator’s native reconstruction objective and
 335 $\lambda_{\text{MCL}} > 0$ trades off reconstruction quality against mecha-
 336 nism contrast. Intuitively, the InfoNCE optimum aligns each
 337 mechanism gradient \mathbf{g}_i more strongly with the embedding
 338 of its own \mathbf{u}_i than with any other \mathbf{u}_j in the batch, which is
 339 precisely the linear-independence requirement of Assump-
 340 tion 2.

342 Intuitively, the InfoNCE optimum aligns each mecha-
 343 nism gradient \mathbf{g}_i more strongly with the embedding
 344 of its own \mathbf{u}_i than with any other \mathbf{u}_j in the batch. Under
 345 standard alignment-and-spread conditions on the projec-
 346 tion heads [58] and a $2d_z + 1$ discriminative conditioning
 347 support, minimizers of \mathcal{L}_{MCL} provably satisfy the linear-
 348 independence content of Assumption 2, hence drive the
 349 model into the GMD regime of Theorem 1; the formal state-
 350 ment appears in Theorem 2 (Appendix B.2).

351 **Comparison with contrastive learning and L_1/L_2 regu-**
 352 **larization.** MCL generalizes two existing methodologies
 353 (Figure 4). **Contrastive learning** on latent codes is the prac-
 354 tical instantiation of LRD [27, 40, 46, 70]; **L_1/L_2 regular-**
 355 **ization** (sparsity [33, 68], quantization [44, 54], combined
 356 penalties [22]) tightens a single latent code without injecting
 357 diversity. **MCL** takes a third stance: $p(\mathbf{z})$ stays fixed and \mathbf{u}
 358

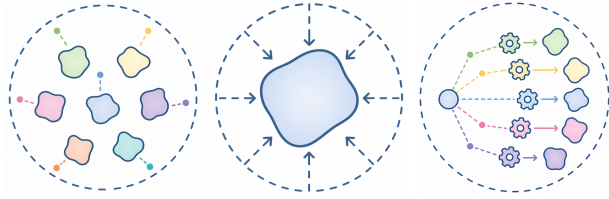


Figure 4. What each family of objectives does in latent space. (Left) **Contrastive learning**, LRD; (Center) L_1/L_2 regularization; (Right) **MCL**.

358 reshapes the *generator*, with the contrastive signal computed
359 on Jacobian gradients — the natural identifiability objective
360 for architectures built with $\mathbf{z} \perp \mathbf{u}$.

361 5. Experiments

362 We validate the theory on synthetic data with known ground
363 truth (§5.1), then evaluate concept disentanglement (static
364 case, §5.2) and world modeling (temporal case, §5.3).

365 **Implementation.** MCL is implemented as a drop-in regu-
366 larizer on top of any conditional generator. The mechanism
367 critic c_ψ is a 3-layer MLP ($d_z + d_{\mathbf{x}} + d_u \rightarrow 256 \rightarrow 256 \rightarrow 1$,
368 GELU activations); the projection heads Π_g, Π_u are 2-layer
369 MLPs into a $d = 128$ -dimensional unit sphere. We use tem-
370 perature $\tau = 0.07$ following SimCLR [6], in-batch negatives
371 only ($B \in \{128, 256\}$ depending on the host model), and a
372 single Adam optimizer for $(G_\theta, c_\psi, \Pi_g, \Pi_u)$ at learning rate
373 10^{-4} . The MCL weight is $\lambda_{\text{MCL}} \in [0.01, 0.1]$ tuned per
374 host model on a held-out validation split. The critic gradient
375 is detached on the InfoNCE path to avoid the trivial solution
376 where c_ψ collapses to a constant; mechanism-gradient norms
377 are clipped at 10. All experiments use mixed-precision train-
378 ing on a single H100 GPU; world-model experiments scale
379 to four H100s with data-parallel.

380 5.1. Experiments on Synthetic Data for Factor Ver- 381 ification

382 **Setup and metric.** We generate data from a known GMD
383 process (full data generation pipeline in Appendix D.3) with
384 $d_z \in \{5, 10, 20\}$ and a default of $|\mathcal{U}| = 8$ conditioning val-
385 ues. We compare three baselines drawn from the three prin-
386 ciples of Section 2: (i) β -VAE [21] (regularization-based),
387 (ii) iVAE [29] (LRD-based), and (iii) MD-VAE (ours, GMD-
388 based), trained with the conditional ELBO of Equation (1).
389 Identifiability is measured by the Mean Correlation Coeffi-
390 cient (MCC) between true and recovered latent variables.

391 **Results.** Across $d_z \in \{5, 10, 20\}$ MD-VAE attains median
392 MCC $\{0.96, 0.87, 0.81\}$, exceeding iVAE, *i.e.*, 0.80 to 0.65,
393 and β -VAE, *i.e.*, 0.53 to 0.39, at every d_z , as shown in
394 Figure 5. The ordering matches Theorem 1, the MD-VAE

exploits mechanism diversity directly, whereas iVAE relies
on Assumption 1 which weakens as d_z grows against a fixed
 $|\mathcal{U}|$, and β -VAE has no identifiability guarantee. Notably,
increasing $|\mathcal{U}|$ can induce more distinct mechanisms (see
Appendix F.1).

5.2. Experiments on Concept Disentanglement

As a representative static GMD task, we evaluate whether
MCL improves concept disentanglement in conditional diffu-
sion models (concept tokens = \mathbf{u} , denoiser = $g(\cdot, \mathbf{u})$). Our
MDDiff builds MCL and mechanism diversity on top of a
latent diffusion generative model with cross-attention, fol-
lowing Yang et al. [61]; full architecture and training details
are in Appendix D.

Datasets. We use three standard disentanglement bench-
marks: Shapes3D [30], MPI3D [13], and Cars3D [45],
which together cover synthetic factors, 3D object factors, and
viewpoint variation. All experiments use 64×64 images,
following [5, 30, 48, 60].

Baselines and metrics. We compare VAE baselines (Fac-
torVAE [30], β -TCVAE [5]), GAN baselines (InfoGAN-
CR [35], GANSpace [19], LatentDisco [56], DisCo [48]),
and diffusion baselines (DisDiff [60], EncDiff [61]) against
MDDiff; full descriptions are in Appendix D.2. All methods
use $N = 20$ scalar representations [60]; we report mean
 \pm std over 15 runs using the FactorVAE score and DCI
disentanglement (as in Appendix D.1).

Results and analysis. Quantitatively (Table 1), even be-
fore MCL, the diffusion-based methods already outperform
the baselines by a clear margin, consistent with our theory.
MDDiff is best or tied-best in nearly all columns; the only
exception is DCI on MPI3D. The largest gains appear on
Cars3D (the hardest benchmark) DCI, *i.e.*, 0.493 vs. 0.279
for EncDiff, ours obtain +0.214. The training dynamics on
Cars3D (Figure 5) trace this gain to the two-stage MDDiff
pipeline: MDDiff pretraining first learns the conditional gen-
erator, then MDDiff + MCL fine-tuning drives both DCI
and FV upward relative to EncDiff. MPI3D’s complex 3D
interactions break LRD-based methods but not MDDiff, sup-
porting our claim that mechanism diversity is a stronger
inductive bias. Qualitatively (Figure 6), swapping a single
latent dimension between source and target images transfers
exactly one attribute while preserving all others, the opera-
tional signature of an axis-aligned factored representation,
in the sense of Theorem 1. Generation quality (FID, LPIPS)
is preserved competitively with EncDiff (Appendix F.2, Ta-
ble 6); latent traversals and cross-attention maps further
confirming localised mechanism diversity in the generator
are reported in Appendix F.

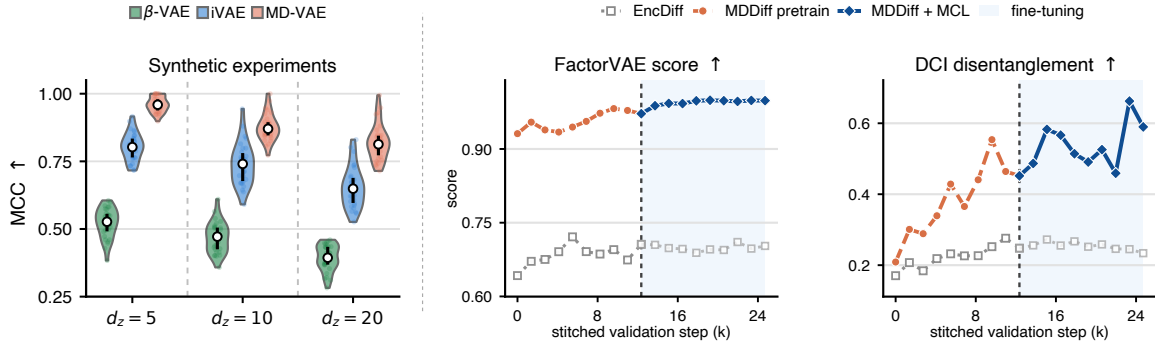


Figure 5. **Synthetic experiments and Cars3D metrics curves during training.** Left: MCC violin plots over 30 seeds per cell. Right: stitched FactorVAE and DCI trajectories on Cars3D.

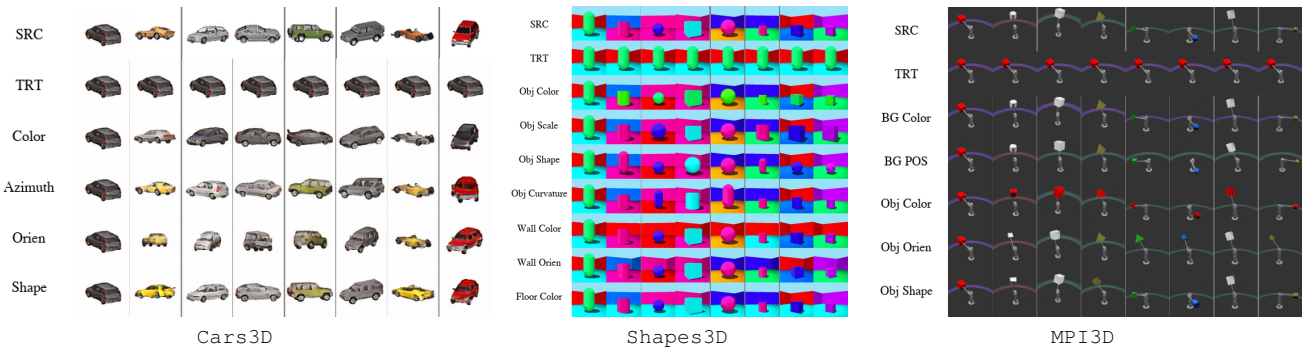


Figure 6. **Factor swapping across all three benchmarks.** Row 1: source images (SRC); row 2: target images (TRT); subsequent rows: SRC with one latent dimension replaced by the corresponding TRT dimension, we match the concept through their consistent visual variations across samples.

443 5.3. Experiments on World Modeling

444 We show that MCL can improve world modeling of temporal
445 robotics data. We first relate our theorem to practical sce-
446 narios in principle, and then present our methodology and
447 results.

448 **Forward/inverse dynamic models.** Recent VA/WA mod-
449 els [34, 64] cast the world as an SEM on $(\mathbf{o}_{1:T}, \mathbf{a}_{1:T})$
450 with two halves: a forward action-conditioned generative
451 model $\mathbf{o}_{t+1} = g_{\text{fwd}}(\mathbf{z}_t, \mathbf{a}_t)$, and an inverse-dynamics model
452 $\mathbf{a}_t = g_{\text{inv}}(\mathbf{o}_t, \mathbf{o}_{t+1})$. State-of-the-art pipelines train the in-
453 verse half by extracting actions from pixel dynamics; we
454 target the *forward* half, where actions reshape the latent dy-
455 namics as the conditioning signal $\mathbf{u}_t = \mathbf{a}_t$. The two halves
456 form an inverse pair of the same SEM,

457

$$\underbrace{\mathbf{o}_{t+1} = g_{\text{fwd}}(\mathbf{z}_t, \mathbf{a}_t)}_{\text{our action-conditioned generation (forward)}}$$

458 $\iff \underbrace{\mathbf{a}_t = g_{\text{inv}}(\mathbf{o}_t, \mathbf{o}_{t+1})}_{\text{VA/WA action extraction (inverse)}}$

459 and the GMD identifiability theory we develop here directly
460 applies to the forward half: when the action sequence suf-

ficiently reshapes g_{fwd} , the latent state is provably factored
(Theorem 1). This is the mirror image of how the condition-
ing signal diversifies temporal dynamics, with both action
alignment and visual fidelity reflecting the controllability or
identifiability of \mathbf{u} .

466 **Architecture.** Our action-conditioned generative world
467 model (ACG) instantiates the forward SEM on top of
468 **Wan 1.4B** [57], a pretrained video generation backbone.
469 Two regimes differ in whether the action conditioning is
470 observed or inferred: **Wan-ACG** feeds the recorded action
471 \mathbf{a}_t directly to the forward model; **Wan-LACG** (the L de-
472 notes *latent* action) replaces it with an IDM-inferred latent
473 action. Both regimes use an inverse dynamics model, an
474 AdaLN action interface, and a flow-matching forward model
475 initialized from Wan 1.4B.

476 **MCL on the action conditioning.** MCL attaches to
477 the action channel: a mechanism critic produces $\mathbf{g}_t =$
478 $\nabla_{\mathbf{z}_t} c_{\psi}(\hat{\mathbf{x}}_{t+1}, \mathbf{z}_t, \mathbf{u}_t)$, contrasted via InfoNCE against the
479 projected action embedding. The IDM, action interface, for-
480 ward model, and MCL heads train jointly with the causal
481 VAE frozen. Full hyperparameters are in Appendix D.

Table 1. **Disentanglement comparison** (mean \pm std; best results are in **Bold**).

Type	Method	Cars3D		Shapes3D		MPI3D	
		FV score \uparrow	DCI \uparrow	FV score \uparrow	DCI \uparrow	FV score \uparrow	DCI \uparrow
VAE	FactorVAE	0.906 \pm 0.052	0.161 \pm 0.019	0.840 \pm 0.066	0.611 \pm 0.082	0.152 \pm 0.025	0.240 \pm 0.051
	β -TCVAE	0.855 \pm 0.082	0.140 \pm 0.019	0.873 \pm 0.074	0.613 \pm 0.114	0.179 \pm 0.017	0.237 \pm 0.056
GAN	InfoGAN-CR	0.411 \pm 0.013	0.020 \pm 0.011	0.587 \pm 0.058	0.478 \pm 0.055	0.439 \pm 0.061	0.241 \pm 0.075
	LatentDisco	0.852 \pm 0.039	0.216 \pm 0.072	0.805 \pm 0.064	0.380 \pm 0.062	0.391 \pm 0.039	0.196 \pm 0.038
	GANspace	0.932 \pm 0.018	0.209 \pm 0.031	0.788 \pm 0.091	0.284 \pm 0.034	0.465 \pm 0.036	0.229 \pm 0.042
	DisCo	0.855 \pm 0.074	0.271 \pm 0.037	0.877 \pm 0.031	0.708 \pm 0.048	0.371 \pm 0.030	0.292 \pm 0.024
Diffusion	DisDiff	0.976 \pm 0.018	0.232 \pm 0.019	0.902 \pm 0.043	0.723 \pm 0.013	0.617 \pm 0.070	0.337 \pm 0.057
	EncDiff	0.773 \pm 0.060	0.279 \pm 0.022	0.999 \pm 0.000	0.969 \pm 0.030	0.872 \pm 0.049	0.685 \pm 0.044
	MDDiff (Ours)	0.985 \pm 0.048	0.493 \pm 0.025	1.000 \pm 0.000	0.993 \pm 0.030	0.930 \pm 0.049	0.679 \pm 0.010

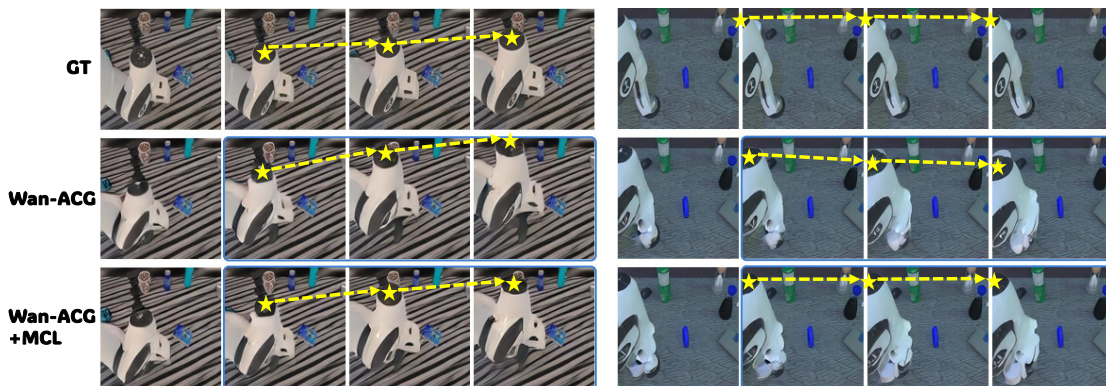


Figure 7. **World-model qualitative ablation.** Two manipulation rollouts driven by the same initial observation and ground-truth action sequence. **Top:** ground-truth frames (**GT**). **Middle:** Wan-ACG alone. **Bottom:** Full method **Wan-ACG+MCL**. The dashed yellow trajectory traces the gripper end-effector across frames; Wan-ACG+MCL follows the ground truth more faithfully than ACG alone, confirming that MCL in the action conditioning sharpens action-following accuracy.

Table 2. **World-model evaluation on Robotwin under the low-data target-embodiment setting.** **Target task:** held-out evaluation on `place_a2b_left` (same task as training). **Transfer task:** zero-shot evaluation on `place_a2b_right`. Higher is better for SSIM, PSNR, and SSIM-L; lower is better for MSE. **Bold** marks the best result per column.

Method	Target task				Transfer task			
	SSIM \uparrow	PSNR \uparrow	MSE \downarrow	SSIM-L \uparrow	SSIM \uparrow	PSNR \uparrow	MSE \downarrow	SSIM-L \uparrow
Wan-ACG	0.763	18.61	0.0138	0.735	0.763	18.45	0.0143	0.732
Wan-ACG+MCL	0.765	18.73	0.0134	0.737	0.765	18.58	0.0138	0.739
Wan-LACG	0.769	18.98	0.0127	0.744	0.782	19.31	0.0117	0.752
Wan-LACG+MCL	0.772	18.99	0.0125	0.749	0.784	19.33	0.0115	0.758

482 **Protocol and methods.** We evaluate on Robotwin [39]
483 under the low-data target-embodiment setting (Table 2),
484 reporting SSIM, PSNR, MSE, and last-frame SSIM-L on the
485 target task (`place_a2b_left`) and the zero-shot transfer
486 task (`place_a2b_right`). For each of the two regimes
487 (Wan-ACG, Wan-LACG) we report the variant without MCL
488 (only the host objective) and with MCL added; the struc-
489 tural claim is that MCL on the action conditioning improves
490 both visual prediction quality and action-following over the

matching no-MCL baseline in either regime.

491

6. Conclusion

492

Mechanism diversity provides a principled, parameter-free
493 route to factored generative models: when conditioning suf-
494 ficiently reshapes the generator, the latent space is provably
495 factored. The framework strictly generalizes LRD, extends
496 to self-supervised and multi-view regimes, and is opera-
497 tionalized by MCL as a model-agnostic objective for any
498 conditional generator.
499

Limitations and future work. The framework is inspired
500 by nonlinear ICA and targets identifiable, factored generative
501 models. We see three open directions: extending
502 mechanism diversity *beyond factored generative models*
503 to broader conditional generators (large-scale text-to-image
504 and video diffusion, visual world models); using *MCL as a*
505 *fine-tuning signal* for large pretrained backbones, actively
506 injecting mechanism shift on top of passive pretraining; and
507 *multi-objective, multi-auxiliary mechanism contrast*. Each
508 direction is expanded in Appendix H.
509

509

510

References

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

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

- [1] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua Tenenbaum, Tommi Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision-making? In *International Conference on Learning Representations*, 2023. 24, 28
- [2] Arthur Allshire, Roberto Martín-Martín, Charles Lin, Shawn Manuel, Silvio Savarese, and Animesh Garg. Laser: Learning a latent action space for efficient reinforcement learning. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6650–6656. IEEE, 2021. 1, 3
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. In *arXiv preprint arXiv:2311.15127*, 2023. 26
- [4] Boyuan Chen, Tianyuan Zhang, Haoran Geng, Kiwhan Song, Caiyi Zhang, Peihao Li, William T Freeman, Jitendra Malik, Pieter Abbeel, Russ Tedrake, et al. Large video planner enables generalizable robot control. *arXiv preprint arXiv:2512.15840*, 2025. 3
- [5] Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in neural information processing systems*, 2018. 2, 6, 20
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 6, 20
- [7] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems*, 2023. 24, 28
- [8] John B. Conway. *A Course in Functional Analysis*. Springer, 2 edition, 1994. 17
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2, 26, 27
- [10] Nelson Dunford and Jacob T. Schwartz. *Linear Operators, Part III: Spectral Operators*. Wiley-Interscience, 1988. 17
- [11] Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disentangled representations. In *International conference on learning representations*, 2018. 20
- [12] Minghao Fu, Biwei Huang, Zijian Li, Yujia Zheng, Ignavier Ng, Guangyi Chen, Yingyao Hu, and Kun Zhang. Learning general causal structures with hidden dynamic process for climate analysis. 2025. 4
- [13] Muhammad Waleed Gondal, Manuel Wuthrich, Djordje Miladinovic, Francesco Locatello, Martin Breidt, Valentin Volchkov, Joel Akpo, Olivier Bachem, Bernhard Schölkopf, and Stefan Bauer. On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset. In *Advances in Neural Information Processing Systems*, 2019. 6
- [14] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010. 20
- [15] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems*, 2018. 24
- [16] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023. 3, 5, 24, 26
- [17] Hermanni Hälvä, Sylvain Le Corff, Luc LeHéricy, Jonathan So, Yongjie Zhu, Elisabeth Gassiat, and Aapo Hyvarinen. Disentangling identifiable features from noisy data with structured nonlinear ica. *Advances in Neural Information Processing Systems*, 34:1624–1633, 2021. 2
- [18] Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous control. *arXiv preprint arXiv:2310.16828*, 2023. 5, 24
- [19] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable GAN controls. In *Advances in neural information processing systems*, 2020. 6, 20
- [20] William Harvey, Saeid Naderiparizi, and Frank Wood. Conditional image generation by conditioning variational autoencoders. *arXiv preprint arXiv:2102.12037*, 2021. 2
- [21] Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 2, 6, 21, 27
- [22] Kyle Hsu, Jubayer Ibn Hamid, Kaylee Burns, Chelsea Finn, and Jiajun Wu. Tripod: Three complementary inductive biases for disentangled representation learning. In *International Conference on Machine Learning*, 2024. 2, 5, 27
- [23] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022. 28
- [24] Yingyao Hu and Susanne M Schennach. Instrumental variable treatment of nonclassical measurement error models. *Econometrica*, 76(1):195–216, 2008. 4
- [25] Aapo Hyvarinen and Hiroshi Morioka. Nonlinear ica of temporally dependent stationary sources. In *Artificial intelligence and statistics*, pages 460–469. PMLR, 2017. 2, 4, 20, 21, 27
- [26] Aapo Hyvärinen and Petteri Pajunen. Nonlinear independent component analysis: Existence and uniqueness results. *Neural networks*, 12(3):429–439, 1999. 2
- [27] Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. In *The 22nd international conference on artificial intelligence and statistics*, pages 859–868. PMLR, 2019. 1, 2, 3, 5, 20, 25
- [28] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In *Proceedings of the 39th International Conference on Machine Learning*, pages 9902–9915. PMLR, 2022. 24, 28
- [29] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *International conference on*

- 624 *artificial intelligence and statistics*, pages 2207–2217. PMLR, 2020. 2, 3, 6, 20, 21, 27
- 625
- 626 [30] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International conference on machine learning*, pages 2649–2658. PMLR, 2018. 2, 6, 20, 27
- 627
- 628
- 629 [31] David Klindt, Lukas Schott, Yash Sharma, Ivan Ustyuzhaninov, Wieland Brendel, Matthias Bethge, and Dylan Païton. Towards nonlinear disentanglement in natural data with temporal sparse coding. In *International Conference on Learning Representations*, 2021. 26, 27
- 630
- 631 [32] Lingjing Kong, Shaoan Xie, Weiran Yao, Yujia Zheng, Guangyi Chen, Petar Stojanov, Victor Akinwande, and Kun Zhang. Partial disentanglement for domain adaptation. In *International conference on machine learning*, pages 11455–11472. PMLR, 2022. 2
- 632
- 633
- 634 [33] Sébastien Lachapelle, Pau Rodriguez, Yash Sharma, Katie E Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via mechanism sparsity regularization: A new principle for nonlinear ica. In *Conference on Causal Learning and Reasoning*, pages 428–484. PMLR, 2022. 2, 5, 26, 27
- 639
- 640 [34] Shuang Li, Yihuai Gao, Dorsa Sadigh, and Shuran Song. Unified video action model. *arXiv preprint arXiv:2503.00200*, 2025. 7
- 641
- 642
- 643 [35] Zinan Lin, Kiran Thekumparampil, Giulia Fanti, and Se-woong Oh. Infogan-cr and modelcentrality: Self-supervised model training and selection for disentangling gans. In *International conference on machine learning*, pages 6127–6139. PMLR, 2020. 6, 20
- 644
- 645 [36] Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Stratis Gavves. CITRIS: Causal identifiability from temporal intervened sequences. In *International Conference on Machine Learning*, 2022. 26, 27
- 646
- 647 [37] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4114–4124. PMLR, 2019. 1, 2
- 648
- 649 [38] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 1, 2
- 650
- 651 [39] Yao Mu, Tianxing Chen, Shijia Peng, Zanxin Chen, Zeyu Gao, Yude Zou, Lunkai Lin, Zhiqiang Xie, and Ping Luo. Robotwin: Dual-arm robot benchmark with generative digital twins (early version). In *European Conference on Computer Vision*, pages 264–273. Springer, 2024. 8
- 652
- 653 [40] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 5, 18, 19, 20, 27
- 654
- 655 [41] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2, 24, 26
- 656
- 657 [42] William Peebles, John Peebles, Jun-Yan Zhu, Alexei Efros, and Antonio Torralba. The hessian penalty: A weak prior for unsupervised disentanglement. In *European Conference on Computer Vision*, 2020. 26, 27
- 658
- 659 [43] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. In *arXiv preprint arXiv:2204.06125*, 2022. 26, 27
- 681
- 682
- 683
- 684
- 685 [44] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *Advances in Neural Information Processing Systems*, 2019. 5, 26, 27
- 686
- 687
- 688 [45] Scott E Reed, Yi Zhang, Yuting Zhang, and Honglak Lee. Deep visual analogy-making. In *Advances in neural information processing systems*, 2015. 6
- 689
- 690
- 691 [46] Patrik Reizinger, Alice Bizeul, Attila Juhos, Julia E Vogt, Randall Balestriero, Wieland Brendel, and David Klindt. Cross-entropy is all you need to invert the data generating process. *arXiv preprint arXiv:2410.21869*, 2024. 2, 5
- 692
- 693
- 694 [47] Patrik Reizinger, Balint Mucsanyi, Siyuan Guo, Benjamin Eysenbach, Bernhard Schölkopf, and Wieland Brendel. Skill learning via policy diversity yields identifiable representations for reinforcement learning. *arXiv preprint arXiv:2507.14748*, 2025. 2
- 695
- 696
- 697 [48] Xuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Learning disentangled representation by exploiting pretrained generative models: A contrastive learning view. In *International Conference on Learning Representations*, 2021. 6, 20
- 698
- 699
- 700 [49] Geoffrey Roeder, Luke Metz, and Durk P Kingma. On linear identifiability of learned representations. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. 1, 3
- 701
- 702
- 703 [50] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 4, 24, 26
- 704
- 705
- 706 [51] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 28
- 707
- 708
- 709 [52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 24, 26
- 710
- 711
- 712
- 713 [53] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015. 1, 2, 26, 27
- 714
- 715
- 716 [54] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, 2017. 5, 26, 27
- 717
- 718
- 719 [55] Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. *Advances in neural information processing systems*, 34:16451–16467, 2021. 1, 2
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738

- 739 [56] Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the GAN latent space. In *International conference on machine learning*, pages 9786–9796. PMLR, 2020. 6, 20
- 740
- 741
- 742
- 743 [57] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 7, 22
- 744
- 745
- 746
- 747
- 748 [58] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, 2020. 5, 15, 18, 19
- 749
- 750
- 751
- 752 [59] Yingheng Wang, Yair Schiff, Aaron Gokaslan, Weishen Pan, Fei Wang, Christopher De Sa, and Volodymyr Kuleshov. Infodiffusion: Representation learning using information maximizing diffusion models. *International conference on machine learning*, 2023. 26, 27
- 753
- 754
- 755
- 756
- 757 [60] Tao Yang, Yuwang Wang, Yan Lv, and Nanning Zheng. Disdiff: Unsupervised disentanglement of diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 2023. 6, 20, 24, 26, 27
- 758
- 759
- 760
- 761 [61] Tao Yang, Cuiling Lan, Yan Lu, and Nanning Zheng. Diffusion model with cross attention as an inductive bias for disentanglement. *Advances in Neural Information Processing Systems*, 37:82465–82492, 2024. 6, 20, 21, 24, 26, 27
- 762
- 763
- 764
- 765 [62] Dingling Yao, Danru Xu, Sébastien Lachapelle, Sara Magliacane, Perouz Taslakian, Georg Martius, Julius von Kügelgen, and Francesco Locatello. Multi-view causal representation learning with partial observability. *arXiv preprint arXiv:2311.04056*, 2023. 1, 2
- 766
- 767
- 768
- 769
- 770 [63] Weiran Yao, Guangyi Chen, and Kun Zhang. Temporally disentangled representation learning. *Advances in Neural Information Processing Systems*, 35:26492–26503, 2022. 2, 26, 27
- 771
- 772
- 773
- 774 [64] Seonghyeon Ye, Yunhao Ge, Kaiyuan Zheng, Shenyuan Gao, Sihyun Yu, George Kurian, Suneel Indupuru, You Liang Tan, Chuning Zhu, Jiannan Xiang, et al. World action models are zero-shot policies. *arXiv preprint arXiv:2602.15922*, 2026. 7
- 775
- 776
- 777
- 778 [65] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 833–842, 2021. 2
- 779
- 780
- 781
- 782
- 783 [66] Kun Zhang, Shaoan Xie, Ignavier Ng, and Yujia Zheng. Causal representation learning from multiple distributions: A general setting. *arXiv preprint arXiv:2402.05052*, 2024. 2
- 784
- 785
- 786 [67] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023. 1, 28
- 787
- 788
- 789
- 790 [68] Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ica: Sparsity and beyond. *Advances in neural information processing systems*, 35:16411–16422, 2022. 2, 5, 25, 26, 27
- 791
- 792
- 793
- 794 [69] Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. DINO-WM: World models on pre-trained visual features enable zero-shot planning. In *Forty-second International Conference on Machine Learning*, 2025. 24, 26, 28
- 795
- 796
- 797
- 798 [70] Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. In *International conference on machine learning*, pages 12979–12990. PMLR, 2021. 2, 5
- 799
- 800
- 801
- 802

803

804

Appendix

Supplementary material for **Factored Generative Models through Mechanism Diversity**

Roadmap

A	Notation <i>A consolidated reference for symbols used in the theory, objectives, and experiments.</i>	p. 13
B	Formal Statements <i>Formal versions of the main identifiability theorem, its extensions, and the MCL mechanism-diversity theorem.</i>	p. 13
C	Proofs <i>Proofs of the formal statements in Appendix B, plus the bridge from MCL gradients to the GMD condition.</i>	p. 15
D	Experimental Details <i>Disentanglement metrics, baseline descriptions, MD-VAE / MDDiff / ACG-WM architectures, training schedules, and hyperparameters.</i>	p. 20
E	PyTorch Implementation of MCL <i>Primary PyTorch loss (Algorithm 1) and the full implementation including projection head and mechanism critic (Algorithm 2).</i>	p. 22
F	Additional Experimental Results <i>U ablation, generation quality, world-model placeholders, latent traversals, and cross-attention maps.</i>	p. 23
G	Extended Related Work <i>Why existing generative models can be factored, and a categorization of disentanglement methods by underlying principle.</i>	p. 24
H	Recommendations for Future Work <i>Beyond identifiability, MCL as fine-tuning of large pretrained models, multi-objective / multi-auxiliary mechanism contrast.</i>	p. 26

805

A. Notation

806

Table 3. Notation used throughout the paper.

Symbol	Description
<i>Generative model and data.</i>	
\mathbf{x}	Observation (data sample).
$\hat{\mathbf{x}}$	Reconstructed/generated observation produced by the model.
\mathbf{z}	Ground-truth latent code.
$\hat{\mathbf{z}}$	Learned latent code; identifiability is up to permutation and componentwise invertible maps.
\mathbf{u}	Conditioning signal: text, class label, viewpoint, regime index, or action that reshapes the generator.
\mathbf{o}_t	World-model observation at time t .
\mathbf{a}_t	World-model action at time t , playing the role of \mathbf{u} in the temporal GMD setting.
<i>Generator and identifiability.</i>	
$g(\mathbf{z}, \mathbf{u})$ or $g_{\mathbf{u}}(\mathbf{z})$	Conditional generator family indexed by \mathbf{u} ; GMD studies how this family varies across conditioning values.
$\hat{g}(\hat{\mathbf{z}}, \mathbf{u})$	Observationally equivalent learned generator used in the identifiability argument.
$h_{\mathbf{u}} = \hat{g}_{\mathbf{u}}^{-1} \circ g_{\mathbf{u}}$	Relabeling map from true latents to learned latents at a fixed conditioning value.
J_g	Jacobian of the generator with respect to the latent variables.
J_h	Jacobian of the relabeling map with respect to the latent variables.
$m'_i(\mathbf{z}, \mathbf{u})$	First derivative $\partial_{z_i} \log J_g(\mathbf{z}, \mathbf{u}) $.
$m''_{ii}(\mathbf{z}, \mathbf{u})$	Diagonal second derivative $\partial_{z_i}^2 \log J_g(\mathbf{z}, \mathbf{u}) $.
$V(\mathbf{z}, \mathbf{u}^{(n)})$	GMD feature vector collecting the mechanism-derivative coordinates used in Assumption 2.
d_z	Latent dimension.
d_u	Conditioning dimension.
$ \mathcal{U} $	Conditioning-set size; number of distinct conditioning mechanisms available for contrast.
<i>MCL objective.</i>	
c_{ψ}	Mechanism critic whose latent gradient is contrasted against the conditioning signal.
$\mathbf{g}(\mathbf{z}, \mathbf{u})$	Mechanism gradient $\nabla_{\mathbf{z}} c_{\psi}(G_{\theta}(\mathbf{z}, \mathbf{u}), \mathbf{z}, \mathbf{u})$.
Π_g	MCL projection head mapping the mechanism gradient into the embedding space.
Π_u	MCL projection head mapping the conditioning signal into the embedding space.
$\mathbf{q}(\mathbf{z}, \mathbf{u})$	L_2 -normalized projected mechanism gradient (query).
$\mathbf{k}(\mathbf{u})$	L_2 -normalized projected conditioning embedding (key).
τ	InfoNCE temperature in the MCL loss.
λ_{MCL}	MCL loss weight in the full training objective.

B. Formal Statements

807

The main text states the theory in compressed form to keep the paper readable. This appendix collects the fully quantified statements first; Appendix C then proves them in the same order.

808

809

B.1. Identifiability and Its Extensions

810

Proposition 1 (Full Statement: Capacity Separation of GMD and LRD). Let $\Theta_{\text{LRD}} \subset \mathbb{R}^{d_L}$ and $\Theta_{\text{GMD}} \subset \mathbb{R}^{d_G}$ be open parameter domains. Let

811

812

$$\Psi_{\text{LRD}} : \Theta_{\text{LRD}} \rightarrow \mathbb{R}^q, \quad \Psi_{\text{GMD}} : \Theta_{\text{GMD}} \rightarrow \mathbb{R}^q$$

813

be finite separating coordinate charts for the induced conditional laws, obtained by evaluating a separating family of test functions at finitely many conditioning values. Assume both maps are real analytic, $\Psi_{\text{LRD}}(\Theta_{\text{LRD}}) \subseteq \Psi_{\text{GMD}}(\Theta_{\text{GMD}})$, and

814

815

816 there is a nonempty open set $\Omega \subset \Theta_{\text{GMD}}$ on which

$$817 \quad \text{rank } D\Psi_{\text{GMD}}(\theta_G) > \sup_{\theta_L \in \Theta_{\text{LRD}}} \text{rank } D\Psi_{\text{LRD}}(\theta_L).$$

818 Then

$$819 \quad \Theta_{\text{LRD-in-GMD}} := \{\theta_G \in \Theta_{\text{GMD}} : \exists \theta_L \in \Theta_{\text{LRD}} \text{ with } \Psi_{\text{GMD}}(\theta_G) = \Psi_{\text{LRD}}(\theta_L)\}$$

820 is contained in a proper analytic subset of Θ_{GMD} and therefore has Lebesgue measure zero. Consequently, GMD strictly
821 contains LRD at the level of conditional generative families.

822 **Theorem 1 (Full Statement: Identifiability under Sufficient GMD).** Let $\mathcal{Z}, \widehat{\mathcal{Z}} \subset \mathbb{R}^{d_z}$ be open connected sets, let $p_{\mathbf{z}}$ and $p_{\widehat{\mathbf{z}}}$
823 be positive C^2 densities, and let $\mathbf{z} \perp \mathbf{u}$. For each $\mathbf{u} \in \mathcal{U}$, suppose

$$824 \quad g_{\mathbf{u}} : \mathcal{Z} \rightarrow \mathcal{X}, \quad \hat{g}_{\mathbf{u}} : \widehat{\mathcal{Z}} \rightarrow \mathcal{X}$$

825 are C^3 diffeomorphisms onto a common open image with nonzero Jacobian determinants. Assume observational equivalence,
826 $(g_{\mathbf{u}})_{\#} p_{\mathbf{z}} = (\hat{g}_{\mathbf{u}})_{\#} p_{\widehat{\mathbf{z}}}$ for every \mathbf{u} , and define

$$827 \quad h_{\mathbf{u}} := \hat{g}_{\mathbf{u}}^{-1} \circ g_{\mathbf{u}}, \quad r_{\mathbf{u}} := h_{\mathbf{u}}^{-1}.$$

828 For $i = 1, \dots, d_z$, define

$$829 \quad m'_i(\mathbf{z}, \mathbf{u}) = \partial_{z_i} \log |\det Dg_{\mathbf{u}}(\mathbf{z})|, \quad m''_{ii}(\mathbf{z}, \mathbf{u}) = \partial_{z_i}^2 \log |\det Dg_{\mathbf{u}}(\mathbf{z})|,$$

830 and $V(\mathbf{z}, \mathbf{u}) = (m'_1, \dots, m'_{d_z}, m''_{11}, \dots, m''_{d_z d_z}) \in \mathbb{R}^{2d_z}$. Assume:

831 **A1 (Sufficient GMD):** for every $\mathbf{z} \in \mathcal{Z}$, there exist $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(2d_z)} \in \mathcal{U}$ such that the $2d_z$ contrasts $V(\mathbf{z}, \mathbf{u}^{(n)}) - V(\mathbf{z}, \mathbf{u}^{(0)})$,
832 $n = 1, \dots, 2d_z$, are linearly independent.

833 **A2 (Mixed-derivative separation):** for every distinct k, l and every $n = 1, \dots, 2d_z$,

$$834 \quad \left\langle V(\mathbf{z}, \mathbf{u}^{(n)}) - V(\mathbf{z}, \mathbf{u}^{(0)}), B_{kl}(\mathbf{z}, \mathbf{u}) \right\rangle = 0,$$

835 where, evaluated at $\hat{\mathbf{z}} = h_{\mathbf{u}}(\mathbf{z})$,

$$836 \quad B_{kl}(\mathbf{z}, \mathbf{u}) = \begin{pmatrix} \partial^2 r_{\mathbf{u},1} & & \partial^2 r_{\mathbf{u},d_z} & \partial r_{\mathbf{u},1} & \partial r_{\mathbf{u},1} & & \partial r_{\mathbf{u},d_z} & \partial r_{\mathbf{u},d_z} \\ \partial \hat{z}_k \partial \hat{z}_l & \dots & \partial \hat{z}_k \partial \hat{z}_l & \partial \hat{z}_k & \partial \hat{z}_l & \dots & \partial \hat{z}_k & \partial \hat{z}_l \end{pmatrix}.$$

837 Then, for each \mathbf{u} , there exist a permutation $\pi_{\mathbf{u}}$ and scalar C^1 diffeomorphisms $\{h_{\mathbf{u},i}\}_{i=1}^{d_z}$ such that

$$838 \quad \hat{z}_i = h_{\mathbf{u},i}(z_{\pi_{\mathbf{u}}(i)}), \quad i = 1, \dots, d_z.$$

839 If $\mathbf{u} \mapsto h_{\mathbf{u}}$ is continuous on a connected conditioning set, the permutation is independent of \mathbf{u} . Under the shared-latent-chart
840 convention used in Definition 1, this is the factored form $\hat{z}_i = h_i(z_{\pi(i)})$.

841 **Proposition 2 (Full Statement: Self-Supervised GMD).** Let $\mathcal{F}_t = \sigma(\mathbf{x}_{\leq t})$ and let $\hat{\mathbf{z}}_t = e_{\phi}(\mathbf{x}_{\leq t})$ be an \mathcal{F}_t -measurable
842 context representation. Suppose the next-step model has the form $\mathbf{x}_{t+1} = g(\mathbf{z}_{t+1}, \hat{\mathbf{z}}_t)$, with innovation exogeneity $\mathbf{z}_{t+1} \perp \mathcal{F}_t$,
843 hence $\mathbf{z}_{t+1} \perp \hat{\mathbf{z}}_t$. If the support of $\hat{\mathbf{z}}_t$ contains the $2d_z + 1$ context values required by Theorem 1, and the conditional generator
844 family satisfies the remaining regularity and mixed-derivative separation assumptions of that theorem with $\mathbf{u} = \hat{\mathbf{z}}_t$, then \mathbf{z}_{t+1}
845 is identifiable from \mathbf{x}_{t+1} up to permutation and componentwise invertible transformations.

846 **Proposition 3 (Full Statement: Distribution-Level Invertibility from Three GMD Views).** Let $\mathbf{x}^{(n)}$ denote the observation
847 generated under conditioning value $\mathbf{u}^{(n)}$, $n = 0, 1, 2$, with $\mathbf{x}^{(n)} \perp \mathbf{x}^{(m)} \mid \mathbf{z}$ for $n \neq m$ and $\mathbf{z} \perp \mathbf{u}$. Under the operator
848 injectivity, bounded-density, non-redundancy, and smoothness assumptions stated in Assumptions 3–6, any observationally
849 equivalent latent variable $\hat{\mathbf{z}}$ is related to \mathbf{z} by an invertible differentiable relabeling map h , i.e. $\hat{\mathbf{z}} = h(\mathbf{z})$. Therefore, once the
850 sufficient-GMD derivative condition of Theorem 1 holds after this distribution-level relabeling, the conclusion of Theorem 1
851 holds without assuming pointwise invertibility of $g(\cdot, \mathbf{u})$.

B.2. MCL as a Mechanism-Diversity Objective

We use the notation of Appendix A and Section 4: the host generator G_θ , the mechanism critic c_ψ , the projection heads Π_g and Π_u into a common d -dimensional embedding space, the mechanism gradient $\mathbf{g}(\mathbf{z}, \mathbf{u}) = \nabla_{\mathbf{z}} c_\psi(G_\theta(\mathbf{z}, \mathbf{u}), \mathbf{z}, \mathbf{u})$, and the L_2 -normalized query and key $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u})) / \|\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}))\|$ and $\mathbf{k}(\mathbf{u}) = \Pi_u(\mathbf{u}) / \|\Pi_u(\mathbf{u})\|$. Replacing the in-batch negatives in the MCL loss by independent samples from $p(\mathbf{u})$ yields the population InfoNCE objective

$$\mathcal{L}_{\text{MCL}}^\infty(\tau) = \mathbb{E}_{(\mathbf{z}, \mathbf{u}) \sim p(\mathbf{z})p(\mathbf{u})} \left[-\frac{\mathbf{q}(\mathbf{z}, \mathbf{u})^\top \mathbf{k}(\mathbf{u})}{\tau} + \log \mathbb{E}_{\mathbf{u}' \sim p(\mathbf{u})} \exp\left(\frac{\mathbf{q}(\mathbf{z}, \mathbf{u})^\top \mathbf{k}(\mathbf{u}')}{\tau}\right) \right]. \quad (2)$$

The next theorem shows that minimizers of $\mathcal{L}_{\text{MCL}}^\infty$ realize the linear-independence content of Assumption 2.

A1 (Regularity): G_θ and c_ψ are C^2 in \mathbf{z} ; Π_g and Π_u are continuous; and $\mathbf{g}(\mathbf{z}, \mathbf{u}) \neq 0$, $\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u})) \neq 0$, $\Pi_u(\mathbf{u}) \neq 0$ for $p(\mathbf{z})p(\mathbf{u})$ -a.e. (\mathbf{z}, \mathbf{u}) .

A2 (Capacity and spread): $d \geq 2d_z + 1$, and the hypothesis class $\{(G_\theta, c_\psi, \Pi_g, \Pi_u)\}$ contains a configuration that simultaneously achieves alignment $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}(\mathbf{u})$ a.e. and uniform spread $\mathbf{k}_\# p(\mathbf{u}) = \text{Unif}(S^{d-1})$.

A3 (Discriminative conditioning support): there exist $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(2d_z)} \in \text{supp } p(\mathbf{u})$ such that $\Pi_u(\mathbf{u}^{(0)}), \dots, \Pi_u(\mathbf{u}^{(2d_z)})$ are linearly independent in \mathbb{R}^d .

A2 is the standard alignment–uniformity capacity condition under which the population InfoNCE minimizer is unique [58]; without uniform spread, degenerate constant- \mathbf{q} configurations can match the rescaled-loss infimum. A3 is the operational content of "enough distinguishable auxiliary inputs" and holds for any generic continuous Π_u whenever $\text{supp } p(\mathbf{u})$ contains at least $2d_z + 1$ distinct points.

Theorem 2 (MCL minimizers satisfy mechanism-gradient diversity). *Under A1–A3, for every fixed temperature $\tau > 0$ the population InfoNCE loss (2), restricted to the hypothesis class of A2, attains its infimum at the alignment-and-spread configurations of A2; in particular every such minimizer satisfies $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}(\mathbf{u})$ for $p(\mathbf{z})p(\mathbf{u})$ -a.e. (\mathbf{z}, \mathbf{u}) . Consequently, for $p(\mathbf{z})$ -a.e. \mathbf{z} and the values $\{\mathbf{u}^{(n)}\}_{n=0}^{2d_z}$ supplied by A3, the projected mechanism-gradient contrasts*

$$\Delta_n(\mathbf{z}) := \Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(n)})) - \Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(0)})), \quad n = 1, \dots, 2d_z, \quad (3)$$

are linearly independent in \mathbb{R}^d , realizing the linear-independence content of Assumption 2.

C. Proofs

C.1. Proof of Proposition 1

Proof. (i) Inclusion. Given any LRD model (f_L, g_L) , construct a GMD model by setting $f_G(\epsilon) := \epsilon$ and $g_G(\mathbf{u}, \mathbf{z}) := g_L(f_L(\mathbf{u}, \mathbf{z}))$. Then $\mathbf{x} = g_G(\mathbf{u}, f_G(\epsilon)) = g_L(f_L(\mathbf{u}, \epsilon))$, identical to LRD. Hence $\mathcal{P}_{\text{LRD}} \subseteq \mathcal{P}_{\text{GMD}}$.

(ii) Zero-measure strictness. Let

$$\Phi_{\text{LRD}} : \Theta_{\text{LRD}} \rightarrow \mathcal{F}, \quad \Phi_{\text{GMD}} : \Theta_{\text{GMD}} \rightarrow \mathcal{F}$$

denote the analytic maps from parameters to conditional data-distribution families, where \mathcal{F} is any finite-dimensional coordinate chart obtained by evaluating a finite separating set of test functions and conditioning values. Because $P_{\mathcal{E}}$ has a Lebesgue density, equality of the induced conditional laws implies equality of these separating coordinates. The LRD-realizable subset inside the GMD parameter space is therefore

$$\Theta_{\text{LRD-in-GMD}} = \Phi_{\text{GMD}}^{-1}(\Phi_{\text{LRD}}(\Theta_{\text{LRD}})).$$

The inclusion construction above shows that this set is nonempty. Strictness means that the image $\Phi_{\text{LRD}}(\Theta_{\text{LRD}})$ is contained in a lower-dimensional analytic submanifold of $\Phi_{\text{GMD}}(\Theta_{\text{GMD}})$: LRD can only represent conditioning effects that factor through the latent distribution, whereas GMD parameters also include independent directions in which \mathbf{u} changes the generator while $p(\mathbf{z})$ stays fixed. Under the nondegenerate analytic parameterization stated in the proposition, the analytic rank theorem then gives

$$\dim \Theta_{\text{LRD-in-GMD}} < \dim \Theta_{\text{GMD}}.$$

Lower-dimensional analytic subsets of an open Euclidean parameter space have Lebesgue measure zero, proving the claim. \square

893 **C.2. Proof of Theorem 1**894 *Proof.* Denote $g_{\mathbf{u}}(\cdot) = g(\cdot, \mathbf{u})$, so $\mathbf{x} = g_{\mathbf{u}}(\mathbf{z})$. Observational equivalence defines

895
$$\hat{\mathbf{z}} = h_{\mathbf{u}}(\mathbf{z}) := \hat{g}_{\mathbf{u}}^{-1} \circ g_{\mathbf{u}}(\mathbf{z}),$$

896 with Jacobian $J_h(\mathbf{z}, \mathbf{u})$. By the change-of-variables formula,

897
$$\log p_{\hat{\mathbf{z}}}(\hat{\mathbf{z}}) = \log p_{\mathbf{z}}(\mathbf{z}) - \log |J_h(\mathbf{z}, \mathbf{u})|.$$

898 For distinct indices $k \neq l$, differentiate the identity twice with respect to (\hat{z}_k, \hat{z}_l) . The term $\log p_{\mathbf{z}}(\mathbf{z})$ is independent of \mathbf{u}
899 because $\mathbf{z} \perp \mathbf{u}$. After subtracting the equation at the reference value $\mathbf{u}^{(0)}$ from the equation at $\mathbf{u}^{(n)}$, the \mathbf{u} -independent
900 density term cancels. The mixed-derivative separation condition in Theorem 1 separates true-generator and learned-generator
901 derivative coordinates, leaving the homogeneous system

902
$$0 = \sum_{i=1}^{d_z} \Delta m''_{ii}(\mathbf{z}, n) \frac{\partial z_i}{\partial \hat{z}_l} \frac{\partial z_i}{\partial \hat{z}_k} + \sum_{i=1}^{d_z} \Delta m'_i(\mathbf{z}, n) \frac{\partial^2 z_i}{\partial \hat{z}_k \partial \hat{z}_l}, \quad (4)$$

903 where

904
$$\Delta m'_i(\mathbf{z}, n) = m'_i(\mathbf{z}, \mathbf{u}^{(n)}) - m'_i(\mathbf{z}, \mathbf{u}^{(0)}),$$
905
$$\Delta m''_{ii}(\mathbf{z}, n) = m''_{ii}(\mathbf{z}, \mathbf{u}^{(n)}) - m''_{ii}(\mathbf{z}, \mathbf{u}^{(0)}).$$

906 For fixed (k, l) , define

907
$$B_{kl}(\mathbf{z}, \mathbf{u}) = \left(\frac{\partial^2 z_1}{\partial \hat{z}_k \partial \hat{z}_l}, \dots, \frac{\partial^2 z_{d_z}}{\partial \hat{z}_k \partial \hat{z}_l}, \frac{\partial z_1}{\partial \hat{z}_l} \frac{\partial z_1}{\partial \hat{z}_k}, \dots, \frac{\partial z_{d_z}}{\partial \hat{z}_l} \frac{\partial z_{d_z}}{\partial \hat{z}_k} \right).$$

908 Equation (4) is equivalently

909
$$\langle V(\mathbf{z}, \mathbf{u}^{(n)}) - V(\mathbf{z}, \mathbf{u}^{(0)}), B_{kl}(\mathbf{z}, \mathbf{u}) \rangle = 0, \quad n = 1, \dots, 2d_z.$$

910 The sufficient-GMD condition states that these $2d_z$ contrast vectors are linearly independent in \mathbb{R}^{2d_z} ; hence they span \mathbb{R}^{2d_z}
911 and force $B_{kl}(\mathbf{z}, \mathbf{u}) = 0$ for every distinct pair $k \neq l$. In particular,

912
$$\frac{\partial z_i}{\partial \hat{z}_l} \frac{\partial z_i}{\partial \hat{z}_k} = 0 \quad \text{for all } i, k \neq l.$$

913 Thus each row of $J_{r_{\mathbf{u}}}(\hat{\mathbf{z}})$ has at most one nonzero entry. Since $r_{\mathbf{u}}$ is invertible, $J_{r_{\mathbf{u}}}(\hat{\mathbf{z}})$ is full rank, so each row and column has
914 exactly one nonzero entry. Therefore there are a permutation $\pi_{\mathbf{u}}$ and scalar invertible functions $\rho_{\mathbf{u},i}$ such that

915
$$z_i = \rho_{\mathbf{u},i}(\hat{z}_{\pi_{\mathbf{u}}(i)}), \quad i = 1, \dots, d_z.$$

916 Inverting and reindexing gives $\hat{z}_i = h_{\mathbf{u},i}(z_{\pi_{\mathbf{u}}^{-1}(i)})$ for scalar invertible functions $h_{\mathbf{u},i}$. If $\mathbf{u} \mapsto h_{\mathbf{u}}$ is continuous in C^1 over a
917 connected conditioning set, the discrete permutation $\pi_{\mathbf{u}}$ cannot change with \mathbf{u} without making $J_{h_{\mathbf{u}}}$ discontinuous; hence $\pi_{\mathbf{u}}$
918 is constant across environments. \square 919 **C.3. Proof of Proposition 2**920 *Proof.* Let $\mathcal{F}_t = \sigma(\mathbf{x}_{\leq t})$ be the past-observation sigma-field and let $\hat{\mathbf{z}}_t = e_{\phi}(\mathbf{x}_{\leq t})$ be the context encoder output, so $\hat{\mathbf{z}}_t$ is
921 \mathcal{F}_t -measurable. The predictive model has the form

922
$$\mathbf{x}_{t+1} = g(\mathbf{z}_{t+1}, \hat{\mathbf{z}}_t).$$

923 Assume the next-step innovation latent is conditionally exogenous:

924
$$\mathbf{z}_{t+1} \perp \mathcal{F}_t, \quad \text{and therefore} \quad \mathbf{z}_{t+1} \perp \hat{\mathbf{z}}_t.$$

925 Conditioning on the realized value of $\hat{\mathbf{z}}_t$ gives exactly the static GMD model of Theorem 1, with \mathbf{u} replaced by $\hat{\mathbf{z}}_t$. If the
926 support of $\hat{\mathbf{z}}_t$ contains the $2d_z + 1$ context values required by the sufficient-GMD condition, the theorem applies pointwise over
927 those contexts and yields identifiability of \mathbf{z}_{t+1} up to a permutation and componentwise invertible maps. When e_{ϕ} is learned
928 jointly with g , the same conclusion holds at any population optimum whose encoder output satisfies the support and exogeneity
929 conditions; when e_{ϕ} is fixed, these conditions are assumptions on the fixed representation. This proves the proposition. \square

C.4. Proof of Proposition 3

The high-level argument and its multi-view interpretation are given in the main text after Proposition 3. Here we record the formal assumptions (Assumptions 3–6) and the proof. 931

We first introduce the operator machinery needed for the proof. 932

Definition 2 (Linear Operator). Consider two random variables a and b with support sets \mathcal{A} and \mathcal{B} . The linear operator $L_{b|a}$ maps a density function p_a on \mathcal{A} to a function on \mathcal{B} via: 933

$$[L_{b|a} \circ p_a](b) = \int_{\mathcal{A}} p(b | a) p_a(a) da. \quad (5) \quad 936$$

Definition 3 (Diagonal Operator). Consider two random variables a and b with density functions p_a and p_b on support sets \mathcal{A} and \mathcal{B} . The diagonal operator $D_{b|a}$ maps the density function p_a to another density function via pointwise multiplication of $p_{b|a}$ at a fixed point b : 937

$$D_{b|a} \circ p_a = p_{b|a}(b | \cdot) p_a, \quad \text{where } D_{b|a} = p_{b|a}(b | \cdot). \quad (6) \quad 938$$

The proof requires the following additional assumptions. 939

Assumption 3 (Operator Injectivity). The linear operators $L_{\mathbf{x}^{(n)}|\mathbf{z}}$ defined by the conditional density $p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)})$ are injective for each conditioning value $\mathbf{u}^{(n)}$, $n = 0, 1, 2$. 940

Assumption 4 (Bounded Density). The conditional densities $p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)})$ are uniformly bounded: $\sup_{\mathbf{z}, \mathbf{x}} p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)}) < \infty$ for each n . 941

Assumption 5 (Non-redundancy). For each conditioning value $\mathbf{u}^{(n)}$, distinct latent values produce distinct conditional densities: $\mathbf{z} \neq \mathbf{z}' \implies p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)}) \neq p(\mathbf{x} | \mathbf{z}', \mathbf{u}^{(n)})$ as functions of \mathbf{x} . 942

Assumption 6 (Smoothness). The conditional density $p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)})$ is continuously differentiable in \mathbf{z} for each n . 943

Proof. Denote $\mathbf{x}^{(n)}$ as the observation generated under conditioning value $\mathbf{u}^{(n)}$, so $\mathbf{x}^{(n)} \sim p(\mathbf{x} | \mathbf{z}, \mathbf{u}^{(n)})$ for $n = 0, 1, 2$. Since $\mathbf{z} \perp \mathbf{u}$, the observations $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}$ are conditionally independent given \mathbf{z} : 944

$$\begin{aligned} p(\mathbf{x}^{(0)} | \mathbf{x}^{(1)}, \mathbf{z}) &= p(\mathbf{x}^{(0)} | \mathbf{z}), \\ p(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}, \mathbf{x}^{(0)}, \mathbf{z}) &= p(\mathbf{x}^{(2)} | \mathbf{z}). \end{aligned} \quad (7) \quad 945$$

From $p(\mathbf{x}^{(2)}, \mathbf{x}^{(1)} | \mathbf{x}^{(0)})$, the conditional independence yields: 946

$$\begin{aligned} p(\mathbf{x}^{(2)}, \mathbf{x}^{(1)} | \mathbf{x}^{(0)}) &= \int_{\mathcal{Z}} p(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}, \mathbf{z}, \mathbf{x}^{(0)}) p(\mathbf{x}^{(1)}, \mathbf{z} | \mathbf{x}^{(0)}) d\mathbf{z} \\ &= \int_{\mathcal{Z}} p(\mathbf{x}^{(2)} | \mathbf{z}) p(\mathbf{x}^{(1)} | \mathbf{z}) p(\mathbf{z} | \mathbf{x}^{(0)}) d\mathbf{z}, \end{aligned} \quad (8) \quad 947$$

where the second equality uses (7). Integrating both sides over $\mathbf{x}^{(0)}$ against $p(\mathbf{x}^{(0)})$ and expressing in operator form (Definitions 2 and 3): 948

$$L_{\mathbf{x}^{(1)}; \mathbf{x}^{(2)}|\mathbf{x}^{(0)}} = L_{\mathbf{x}^{(2)}|\mathbf{z}} D_{\mathbf{x}^{(1)}|\mathbf{z}} L_{\mathbf{z}|\mathbf{x}^{(0)}}. \quad (9) \quad 949$$

Integrating over $\mathbf{x}^{(1)}$ on both sides gives $L_{\mathbf{x}^{(2)}|\mathbf{x}^{(0)}} = L_{\mathbf{x}^{(2)}|\mathbf{z}} L_{\mathbf{z}|\mathbf{x}^{(0)}}$. By Assumption 3, $L_{\mathbf{x}^{(2)}|\mathbf{z}}$ is injective, so $L_{\mathbf{z}|\mathbf{x}^{(0)}} = L_{\mathbf{x}^{(2)}|\mathbf{z}}^{-1} L_{\mathbf{x}^{(2)}|\mathbf{x}^{(0)}}$. Substituting back into (9) and right-multiplying by $L_{\mathbf{x}^{(2)}|\mathbf{x}^{(0)}}^{-1}$, which exists and is densely defined by Assumption 3, gives 950

$$\begin{aligned} L_{\mathbf{x}^{(1)}; \mathbf{x}^{(2)}|\mathbf{x}^{(0)}} L_{\mathbf{x}^{(2)}|\mathbf{x}^{(0)}}^{-1} \\ = L_{\mathbf{x}^{(2)}|\mathbf{z}} D_{\mathbf{x}^{(1)}|\mathbf{z}} L_{\mathbf{x}^{(2)}|\mathbf{z}}^{-1}. \end{aligned} \quad (10) \quad 951$$

By Assumption 4, the left-hand side is a bounded operator. By the uniqueness of spectral decomposition [8, 10], the eigenvalues $D_{\mathbf{x}^{(1)}|\mathbf{z}}$ (the entries $\{p(\mathbf{x}^{(1)} | \mathbf{z})\}$) and the eigenfunctions in $L_{\mathbf{x}^{(2)}|\mathbf{z}}$ (the columns $\{p(\mathbf{x}^{(2)} | \mathbf{z})\}$) are unique up to standard indeterminacies: 952

$$\begin{aligned} L_{\mathbf{x}^{(2)}|\mathbf{z}} &= C L_{\mathbf{x}^{(2)}|\hat{\mathbf{z}}} P, \\ D_{\mathbf{x}^{(1)}|\mathbf{z}} &= P^{-1} D_{\mathbf{x}^{(1)}|\hat{\mathbf{z}}} P. \end{aligned} \quad (11) \quad 953$$

969 where C is a nonzero scalar rescaling and P is a permutation operator. Since $\int p(\mathbf{x}^{(2)} | \mathbf{z}) d\mathbf{x}^{(2)} = 1$ for every \mathbf{z} , the only
 970 solution to $\int C p(\mathbf{x}^{(2)} | \mathbf{z}) d\mathbf{x}^{(2)} = 1$ is $C = 1$. From $D_{\mathbf{x}^{(1)}|\mathbf{z}} = P^{-1} D_{\mathbf{x}^{(1)}|\hat{\mathbf{z}}} P$, the sets of conditional densities must match:
 971 $\{p(\mathbf{x}^{(1)} | \mathbf{z})\}_{\mathbf{z}} = \{p(\mathbf{x}^{(1)} | \hat{\mathbf{z}})\}_{\hat{\mathbf{z}}}$. Since sets are unordered, a relabeling map h is needed to consistently match entries:

$$972 \quad p(\mathbf{x}^{(1)} | h(\mathbf{z})) = p(\mathbf{x}^{(1)} | \hat{\mathbf{z}}), \quad \text{for all } \mathbf{z}, \hat{\mathbf{z}}. \quad (12)$$

973 By Assumption 5, distinct \mathbf{z} values produce distinct conditional densities, so h is one-to-one (invertible). By Assumption 6, h
 974 is differentiable. Hence $\hat{\mathbf{z}} = h(\mathbf{z})$ with h invertible and differentiable.

975 Finally, if $d_{\hat{z}} > d_z$, then d_z latent components suffice to explain \mathbf{x} , and the remaining $d_{\hat{z}} - d_z$ components satisfy
 976 $p(\mathbf{x} | \mathbf{z}_{1:d_z}, \mathbf{z}'_{d_z+1:d_{\hat{z}}}) = p(\mathbf{x} | \mathbf{z}_{1:d_z}, \mathbf{z}''_{d_z+1:d_{\hat{z}}})$ for all $\mathbf{z}', \mathbf{z}''$, contradicting Assumption 5. If $d_{\hat{z}} < d_z$, then $d_z - d_{\hat{z}}$ latent
 977 variables are constant, contradicting that they are variables. Hence $d_{\hat{z}} = d_z$. \square

978 C.5. Proof of Theorem 2

979 *Proof.* The argument has two threads: alignment of the projected mechanism gradient to the projected conditioning embedding,
 980 established by an exact decomposition of the population InfoNCE; and linear independence of the gradient contrasts at any
 981 fixed \mathbf{z} , which follows from the alignment-induced gauge equation together with A3.

982 The population InfoNCE (2) is the limit $\lim_{M \rightarrow \infty} [\mathcal{L}_{\text{MCL}}^M(\tau) - \log M]$ of the standard M -negative InfoNCE estimator [40].
 983 Substituting the unit-norm identities $\mathbf{q}^\top \mathbf{k}(\mathbf{u}) = 1 - \frac{1}{2} \|\mathbf{q} - \mathbf{k}(\mathbf{u})\|_2^2$ and $\mathbf{q}^\top \mathbf{k}(\mathbf{u}') = 1 - \frac{1}{2} \|\mathbf{q} - \mathbf{k}(\mathbf{u}')\|_2^2$ into (2) (the $\pm 1/\tau$
 984 contributions cancel) yields the *exact* algebraic identity

$$985 \quad \mathcal{L}_{\text{MCL}}^\infty(\tau) = \underbrace{\frac{1}{2\tau} \mathbb{E}_{(\mathbf{z}, \mathbf{u})} \left[\|\mathbf{q}(\mathbf{z}, \mathbf{u}) - \mathbf{k}(\mathbf{u})\|_2^2 \right]}_{\mathcal{A}(\mathbf{q}, \mathbf{k}; \tau) \geq 0} + \underbrace{\mathbb{E}_{(\mathbf{z}, \mathbf{u})} \left[\log \mathbb{E}_{\mathbf{u}' \sim p(\mathbf{u})} \exp \left(\frac{-\|\mathbf{q}(\mathbf{z}, \mathbf{u}) - \mathbf{k}(\mathbf{u}')\|_2^2}{2\tau} \right) \right]}_{\tilde{u}(\mathbf{q}, \mathbf{k}; \tau)}. \quad (13)$$

986 Two observations make alignment fall out of (13) once the spread part of A2 holds. First, $\mathcal{A}(\mathbf{q}, \mathbf{k}; \tau) \geq 0$ with
 987 equality if and only if $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}(\mathbf{u})$ for $p(\mathbf{z})p(\mathbf{u})$ -a.e. (\mathbf{z}, \mathbf{u}) . Second, when $\mathbf{k}_\# p(\mathbf{u}) = \text{Unif}(S^{d-1})$ the function
 988 $\mathbf{a} \mapsto \mathbb{E}_{\mathbf{u}' \sim p(\mathbf{u})} \exp(-\|\mathbf{a} - \mathbf{k}(\mathbf{u}')\|_2^2 / (2\tau))$ is rotationally invariant on S^{d-1} (its value depends only on $\|\mathbf{a}\|_2$, by spherical
 989 symmetry of the uniform measure); for every unit vector \mathbf{a} it equals the constant $I(\tau, d) := \int_{S^{d-1}} e^{-\|\mathbf{a}_0 - \mathbf{v}\|_2^2 / (2\tau)} d\sigma(\mathbf{v})$,
 990 where \mathbf{a}_0 is any fixed unit vector and σ is the uniform measure on S^{d-1} .

991 Pick any \mathbf{k}^* realizing the spread part of A2; such a key map exists by A2. The rotational-invariance observation gives
 992 $\tilde{u}(\mathbf{q}, \mathbf{k}^*; \tau) = \log I(\tau, d)$ for every unit-norm \mathbf{q} , so substituting into (13) reduces the loss at \mathbf{k}^* to a constant plus the alignment
 993 term:

$$994 \quad \mathcal{L}_{\text{MCL}}^\infty(\mathbf{q}, \mathbf{k}^*; \tau) = \mathcal{A}(\mathbf{q}, \mathbf{k}^*; \tau) + \log I(\tau, d). \quad (14)$$

995 Minimizing over unit-norm \mathbf{q} then selects the unique optimum $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}^*(\mathbf{u})$, achieving the value $L_{\text{AU}}(\tau, d) := \log I(\tau, d)$.
 996 To rule out non-aligned minimizers elsewhere in the hypothesis class, observe that $L_{\text{AU}}(\tau, d)$ is also the *global* infimum
 997 of $\mathcal{L}_{\text{MCL}}^\infty$ over all unit-norm encoder pairs: by the joint-minimum analysis of population InfoNCE [58, Theorem 1 and
 998 Proposition 2] adapted to the two-encoder setup of (2), the global infimum is attained only on the gauge orbit $\{(\mathbf{q}, \mathbf{k}) : \exists R \in$
 999 $O(d), \mathbf{q}(\mathbf{z}, \mathbf{u}) = R \mathbf{k}(\mathbf{u}) \text{ a.s. and } \mathbf{k}_\# p(\mathbf{u}) = \text{Unif}(S^{d-1})\}$. The orthogonal R reflects the gauge invariance of the loss under
 1000 simultaneous rotations $(\mathbf{q}, \mathbf{k}) \mapsto (R\mathbf{q}, R\mathbf{k})$, and is fixed to the identity by the choice of architectures (Π_g, Π_u) once those are
 1001 pinned down. Every minimizer inside the hypothesis class therefore satisfies $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}(\mathbf{u})$ for $p(\mathbf{z})p(\mathbf{u})$ -a.e. (\mathbf{z}, \mathbf{u}) .

1002 With alignment in hand, the linear-independence claim is purely geometric. Fix \mathbf{z} in the (full $p(\mathbf{z})$ -measure) set on which
 1003 $\mathbf{q}(\mathbf{z}, \mathbf{u}) = \mathbf{k}(\mathbf{u})$ holds for $p(\mathbf{u})$ -a.e. \mathbf{u} , and choose the $2d_z + 1$ values $\{\mathbf{u}^{(n)}\}_{n=0}^{2d_z}$ given by A3. Unrolling the definition of \mathbf{q}, \mathbf{k} ,

$$1004 \quad \frac{\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(n)}))}{\|\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(n)}))\|} = \frac{\Pi_u(\mathbf{u}^{(n)})}{\|\Pi_u(\mathbf{u}^{(n)})\|},$$

1005 so the two sides are positively collinear. By A1 the norms are strictly positive; setting

$$1006 \quad \alpha_n(\mathbf{z}) := \frac{\|\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(n)}))\|}{\|\Pi_u(\mathbf{u}^{(n)})\|} > 0, \quad n = 0, \dots, 2d_z, \quad (15)$$

we obtain the gauge equation $\Pi_g(\mathbf{g}(\mathbf{z}, \mathbf{u}^{(n)})) = \alpha_n(\mathbf{z}) \Pi_u(\mathbf{u}^{(n)})$. Substituting into (3),

$$\Delta_n(\mathbf{z}) = \alpha_n(\mathbf{z}) \Pi_u(\mathbf{u}^{(n)}) - \alpha_0(\mathbf{z}) \Pi_u(\mathbf{u}^{(0)}), \quad n = 1, \dots, 2d_z.$$

Suppose $\sum_{n=1}^{2d_z} \beta_n \Delta_n(\mathbf{z}) = 0$ for scalars $\beta_n \in \mathbb{R}$. Expanding,

$$\sum_{n=1}^{2d_z} \beta_n \alpha_n(\mathbf{z}) \Pi_u(\mathbf{u}^{(n)}) - \alpha_0(\mathbf{z}) \left(\sum_{n=1}^{2d_z} \beta_n \right) \Pi_u(\mathbf{u}^{(0)}) = 0,$$

which is a linear combination of the $2d_z + 1$ vectors $\{\Pi_u(\mathbf{u}^{(n)})\}_{n=0}^{2d_z}$. By A3 these vectors are linearly independent in \mathbb{R}^d ; the embedding dimension $d \geq 2d_z + 1$ in A2 is what makes such a tuple feasible. Hence all coefficients vanish:

$$\beta_n \alpha_n(\mathbf{z}) = 0 \quad (n = 1, \dots, 2d_z), \quad \alpha_0(\mathbf{z}) \left(\sum_{n=1}^{2d_z} \beta_n \right) = 0.$$

Since $\alpha_n(\mathbf{z}) > 0$ by (15), the first family forces $\beta_n = 0$ for $n = 1, \dots, 2d_z$. The second equation is then automatically satisfied. Therefore $\{\Delta_n(\mathbf{z})\}_{n=1}^{2d_z}$ is linearly independent in \mathbb{R}^d , completing the proof. \square

C.6. From MCL to the GMD Condition

Theorem 2 establishes linear independence of the projected critic-gradient contrasts $\{\Delta_n(\mathbf{z})\}_{n=1}^{2d_z}$. Assumption 2 is stated in terms of $V(\mathbf{z}, \mathbf{u}^{(n)})$, the vector of first and second \mathbf{z} -derivatives of $\log |J_{G_\theta}(\mathbf{z}, \mathbf{u}^{(n)})|$. The two notions of mechanism diversity are not literally identical, but they coincide under a standard *score-matching* link between the critic and the generator. In each instantiation, MCL is paired with a host loss \mathcal{L}_{rec} (variational ELBO for MD-VAE; conditional denoising score matching for MDDiff; latent-dynamics negative log-likelihood for ACG-WM) whose minimizer drives $c_\psi(\mathbf{x}, \mathbf{z}, \mathbf{u})$ to a \mathbf{z} -derivative-faithful proxy of the conditional log-density of the host model. For the diffusion case, the change-of-variables for the marginal $p_\theta(\mathbf{x} | \mathbf{u}) = p_{\mathbf{z}}(G_{\theta, \mathbf{u}}^{-1}(\mathbf{x})) \cdot |J_{G_{\theta, \mathbf{u}}}(G_{\theta, \mathbf{u}}^{-1}(\mathbf{x}))|^{-1}$ gives the score

$$\nabla_{\mathbf{x}} \log p_\theta(\mathbf{x} | \mathbf{u}) = J_{G_{\theta, \mathbf{u}}}^{-\top}(\mathbf{z}) [\nabla_{\mathbf{z}} \log p_{\mathbf{z}}(\mathbf{z}) - \nabla_{\mathbf{z}} \log |J_{G_{\theta, \mathbf{u}}}(\mathbf{z})|] \Big|_{\mathbf{z}=G_{\theta, \mathbf{u}}^{-1}(\mathbf{x})}, \quad (16)$$

so the conditioning-dependent piece of the critic gradient $\nabla_{\mathbf{z}} c_\psi$ contains exactly the first-order coordinates $m'_i(\mathbf{z}, \mathbf{u}) = \partial_{z_i} \log |J_{G_{\theta, \mathbf{u}}}(\mathbf{z})|$ of $V(\mathbf{z}, \mathbf{u})$, up to an additive \mathbf{u} -independent term that drops out of every contrast Δ_n . Linear independence of $\{\Delta_n(\mathbf{z})\}$ therefore transfers to linear independence of the first-order half of $V(\mathbf{z}, \mathbf{u}^{(n)}) - V(\mathbf{z}, \mathbf{u}^{(0)})$. The second-order coordinates m''_{ii} are recovered by augmenting the projection head Π_g with the diagonal Hessian $\text{diag}(\nabla_{\mathbf{z}}^2 c_\psi)$, after which the same gauge argument applies coordinate-wise. Thus, under the standard architectural pairing of c_ψ with the host generator’s likelihood structure, Theorem 2 delivers Assumption 2 verbatim, and Theorem 1 then yields identifiability up to a permutation and componentwise transformations.

Scope. Theorem 2 concerns *population* minimizers; the standard caveats of finite-sample InfoNCE (capacity, optimization, batch-size dependence of the negative distribution) carry over from the broader contrastive-learning literature [40, 58]. The bridge in (16) is exact only when $G_{\theta, \mathbf{u}}$ is invertible; for non-invertible diffusion or world-model generators, the bridge is replaced by the score-matching identity at the noise level, and the linear-independence transfer holds whenever the score network is rich enough to recover the conditional score. Our experiments (Sections 5.2–5.3) verify that the population picture survives at realistic batch sizes.

C.7. Comparison with Contrastive Objectives

Contrastive losses share the InfoNCE template but differ in what they contrast and what guarantee they target. Table 4 compares MCL against four representative families. The structural distinction is that all prior families contrast *representations* (or representations of contexts/views) and target either a similarity geometry (SimCLR/MoCo), a mutual-information lower bound (CPC), or an LRD-style nonlinear-ICA identifiability (TCL, iVAE-contrastive). MCL is the only family that contrasts *mechanism gradients* $\nabla_{\mathbf{z}} c_\psi(G(\mathbf{z}, \mathbf{u}), \mathbf{z}, \mathbf{u})$ against the conditioning signal, with the explicit goal of certifying mechanism diversity (Assumption 2) and thus a factored generative model. Negatives are drawn by shuffling the conditioning index \mathbf{u} within the batch (rather than perturbing $p(\mathbf{z})$), which is what couples the alignment objective to the GMD geometry rather than to a representational similarity.

Table 4. **Contrastive objective families.** Columns list the positive pair, the source of negatives, the object contrasted, and what the optimum targets, for five representative families.

Family	Positive pair	Negatives	Object contrasted	Optimum target
Augmentation [SimCLR/MoCo; 6]	$(f(\mathbf{x}), f(\tilde{\mathbf{x}}))$, $\tilde{\mathbf{x}}$ an augmentation of \mathbf{x}	In-batch other images	Image representation $f(\mathbf{x})$	Augmentation-invariant similarity
Predictive [CPC; 40]	$(c_t, f(\mathbf{x}_{t+k}))$, context vs. future	In-batch other futures	Context vs. future representation	MI lower bound $I(c_t; \mathbf{x}_{t+k})$
Time-contrastive [TCL; 25]	$(\mathbf{x}, \text{segment label})$	Other segment labels	Encoder of \mathbf{x}	Nonlinear ICA / LRD identifiability
Density-ratio / NCE [14]	Data \mathbf{x} vs. noise $\tilde{\mathbf{x}}$	Noise samples	Log-density-ratio $\log p_{\text{data}}/p_{\text{noise}}$	Density / score estimation
Mechanistic MCL (ours)	$(\mathbf{g}(\mathbf{z}, \mathbf{u}), \mathbf{u})$	Shuffled \mathbf{u}_j in batch	Mechanism gradient vs. conditioning	GMD identifiability, Theorem 2

1047 D. Experimental Details

1048 D.1. Disentanglement Metrics

1049 **FactorVAE score** [30]. For each ground-truth factor k , sample observations varying only in k , encode them, and identify the
 1050 latent dimension with minimal normalized variance: $d_k = \arg \min_j \sigma_j^2 / \sum_{j'} \sigma_{j'}^2$. The score is the accuracy of a majority-vote
 1051 classifier predicting k from d_k . A score of 1.0 means perfect axis-alignment.

1052 **DCI disentanglement** [11]. Train a predictor from \mathbf{z} to each factor k and extract the feature importance matrix $R \in \mathbb{R}^{d \times K}$.
 1053 The disentanglement score is

$$1054 D = 1 - \frac{1}{d} \sum_{j=1}^d H \left(\frac{R_{j \cdot}}{\sum_k R_{jk}} \right) / \log K, \quad (17)$$

1055 where $H(\cdot)$ is the entropy. $D = 1$ when each latent dimension encodes exactly one factor.

1056 D.2. Baselines for Concept Disentanglement

1057 We expand the baseline descriptions for Section 5.2 (Table 1).

1058 **VAE-based. FactorVAE** [30] adds a total-correlation penalty on the aggregated posterior to encourage statistically in-
 1059 dependent latents. **β -TCVAE** [5] decomposes the KL term and upweights the total-correlation component, isolating the
 1060 disentanglement-relevant pressure of β -VAE.

1061 **GAN-based. InfoGAN-CR** [35] extends InfoGAN with a contrastive regulariser on the inferred factors. **GANSpace** [19]
 1062 discovers disentangled directions *post-hoc* via PCA on early-layer GAN activations. **LatentDisco** [56] learns directions in a
 1063 pretrained GAN’s latent space using a self-supervised classifier. **DisCo** [48] is a contrastive method on a pretrained generator:
 1064 it samples paired latent walks and trains an encoder to detect the active factor.

1065 **Diffusion-based. DisDiff** [60] associates each latent with a per-factor sub-gradient field of the score network. **EncDiff** [61]
 1066 routes each scalar latent through an independent MLP and injects the resulting tokens via cross-attention into the U-Net.
 1067 **MDDiff (ours)** uses this tokenized diffusion scaffold, adds \mathbf{u} -conditioned decoding to form the MDDiff pretraining model,
 1068 and then fine-tunes that checkpoint with MCL to enforce mechanism diversity through the cross-attention slots.

1069 D.3. MD-VAE (Simulations)

1070 **Synthetic data-generating process.** We follow the standard nonlinear ICA simulation protocol [27, 29] and instantiate
 1071 it inside the GMD regime. For each run we draw a latent dimension $d_z \in \{5, 10, 20\}$ and a conditioning-set size $|\mathcal{U}| \in$
 1072 $\{2, 4, 8, 16, 32\}$ (default $|\mathcal{U}| = 8$). Each simulated dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{u}_i)\}_{i=1}^N$ with $N = 50,000$ is generated as follows. The
 1073 *conditioning signal* \mathbf{u}_i is drawn uniformly from a fixed alphabet $\mathcal{U} = \{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(|\mathcal{U}|)}\}$, where each $\mathbf{u}^{(n)} \in \mathbb{R}^{d_u}$ ($d_u = 16$) is
 1074 sampled once at the start of the run from $\mathcal{N}(0, \mathbf{I}_{d_u})$. The *latent* $\mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I}_{d_z})$ is sampled independently of \mathbf{u}_i , instantiating

the GMD assumption $\mathbf{z} \perp \mathbf{u}$; this is the key difference from the iVAE [29] simulation, in which $p(\mathbf{z} | \mathbf{u})$ shifts. The *mechanism-diverse generator* is $g(\mathbf{z}, \mathbf{u}) = h \circ (\mathbf{z} \odot \mathbf{m}(\mathbf{u}) + \mathbf{b}(\mathbf{u}))$, where $\mathbf{m}, \mathbf{b} : \mathbb{R}^{d_u} \rightarrow \mathbb{R}^{d_z}$ are two-layer MLPs (width 64, leaky-ReLU) initialized so that $\mathbf{m}(\mathbf{u}^{(n)})$ varies meaningfully across n , and h is a fixed invertible 3-layer leaky-ReLU MLP (width $\max(64, 4d_z)$); the componentwise product $\mathbf{z} \odot \mathbf{m}(\mathbf{u})$ implements per-coordinate mechanism reshaping, satisfying Assumption 2 for $|\mathcal{U}| \geq 2d_z + 1$. Finally, observations are $\mathbf{x}_i = g(\mathbf{z}_i, \mathbf{u}_i) + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, 0.01 \mathbf{I})$. We use an 80%/20% train/test split, repeating each configuration over 30 random seeds.

Models. All three baselines share the same encoder–decoder backbone (3-layer leaky-ReLU MLPs, hidden width $\max(64, 4d_z)$, latent dim d_z) and differ only in how \mathbf{u} is incorporated. **β -VAE** [21] is a standard VAE with KL weight $\beta = 4$ in which \mathbf{u} is fed to neither encoder nor decoder. **iVAE** [29] uses a \mathbf{u} -conditional encoder $q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{u})$ and a \mathbf{u} -conditional Gaussian prior $p(\mathbf{z} | \mathbf{u}) = \mathcal{N}(\mu_\lambda(\mathbf{u}), \text{diag}(\sigma_\lambda^2(\mathbf{u})))$ implemented by a 2-layer MLP, with a decoder that is unconditioned in \mathbf{u} (this is exactly the LRD regime). **MD-VAE** (ours) keeps the prior \mathbf{u} -independent, $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$, and feeds \mathbf{u} to both the encoder $q_\phi(\mathbf{z} | \mathbf{x}, \mathbf{u})$ and the decoder $p_\theta(\mathbf{x} | \mathbf{z}, \mathbf{u})$, training with the conditional ELBO of Equation (1). We optimize with Adam (learning rate 10^{-3} , batch size 256) for 200 epochs.

Evaluation. After training we recover latents $\hat{\mathbf{z}}_i = \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x}_i, \mathbf{u}_i)}[\mathbf{z}]$ on the test split and compute MCC by solving the linear assignment problem between $\hat{\mathbf{z}}$ and the true \mathbf{z} in absolute Pearson correlation, following Hyvarinen and Morioka [25], Khemakhem et al. [29].

D.4. MDDiff (Image Disentanglement)

Architecture. MDDiff builds on EncDiff [61] with two structural additions that together make it a bona fide GMD model. The shared backbone is unchanged from EncDiff: an `Encoder4` CNN (hidden width 128, four stride-2 convolutional layers) maps a 64×64 RGB image to $K=20$ scalar concept factors; each scalar is expanded by an independent three-layer MLP ($1 \rightarrow 64 \rightarrow 128 \rightarrow 16$, ELU activations) into a 16-dim concept token, so each token corresponds to exactly one semantic factor. The denoising backbone is a U-Net with base channel width 64, channel multipliers $[1, 2, 4, 4]$, 2 residual blocks per resolution, and 8-head spatial cross-attention injecting the 20 concept tokens at all three downsampled resolutions.

On top of this backbone we make two additions required for the GMD interpretation. First, we condition the VAE decoder on the concept variable \mathbf{u} : the original EncDiff feeds the concept tokens only through the U-Net cross-attention, so the VAE decoder that maps the diffusion latent back to pixels is unconditional in \mathbf{u} and the concept variable never enters the generator $g(\mathbf{z}, \mathbf{u})$ in the strict sense of Assumption 2. We therefore inject \mathbf{u} as an auxiliary input to the VAE decoder so that it becomes $g(\mathbf{z}, \mathbf{u})$ with \mathbf{z} the diffusion latent and \mathbf{u} the concept token sequence; concretely, we concatenate the per-token features to the spatial feature map at each upsampling stage of the decoder, mirroring the cross-attention pattern in the U-Net. Second, we add a *mechanism critic* $c_\psi(\hat{\mathbf{x}}, \mathbf{z}, \mathbf{u})$ scoring reconstruction–latent–concept tuples, together with two MLP projection heads Π_g, Π_u (\rightarrow 128-dim L_2 -normalized embeddings) that map the critic gradient $\nabla_{\mathbf{z}} c_\psi$ and the concept variable into a shared embedding space, supporting the InfoNCE objective. The two additions together add roughly 265K parameters on top of the shared backbone.

MCL Objective. We compared five MCL variants: NCE, InfoNCE, Fisher score matching, denoising score matching, and Jacobian-based InfoNCE. InfoNCE consistently achieved the highest DCI disentanglement and FactorVAE scores across all three datasets. The other objectives either suffer from training instability (Fisher, Jacobian) or weaker gradient signal when the critic operates on high-dimensional image reconstructions (NCE, denoising SM). We therefore adopt InfoNCE as our final objective:

$$\mathcal{L}_{\text{MCL}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\Pi_g(\mathbf{g}_i)^\top \Pi_u(\mathbf{u}_i)/\tau)}{\sum_{j=1}^B \exp(\Pi_g(\mathbf{g}_i)^\top \Pi_u(\mathbf{u}_j)/\tau)}, \quad (18)$$

where $\mathbf{g}_i = \nabla_{\mathbf{z}_i} c_\psi(\hat{\mathbf{x}}_i, \mathbf{z}_i, \mathbf{u}_i)$ is the critic gradient with respect to the latent, $\tau=0.1$ is the temperature, and negatives are drawn from other samples in the same mini-batch via latent shuffling. The total training loss is

$$\mathcal{L} = \mathcal{L}_{\text{diffusion}} + \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} + \lambda_{\text{MCL}} \mathcal{L}_{\text{MCL}}, \quad (19)$$

where $\mathcal{L}_{\text{recon}}$ is the pixel-level VQ-VAE reconstruction loss weighted by $\lambda_{\text{recon}}=0.1$.

1118 **Training.** Training proceeds in two stages on a single NVIDIA H100 GPU. *Stage 1* pretrains MDDiff as a conditional-
 1119 generation model with the diffusion loss only for 10 epochs (AdamW, base learning rate 2×10^{-6} , 10,000-step linear warmup).
 1120 *Stage 2* fine-tunes the pretrained MDDiff checkpoint with MCL, optimizing all components jointly under the full objective \mathcal{L}
 1121 for an additional 10 epochs (base learning rate 2×10^{-7} , 1,000-step warmup). Both stages use batch size 128 and $T=1000$
 1122 diffusion timesteps with a linear noise schedule ($\beta_{\min}=0.0015$, $\beta_{\max}=0.0155$). Stage 2 fine-tuning takes approximately
 1123 3 hours on Cars3D, 4 hours on Shapes3D, and 7 hours on MPI3D. We set $\lambda_{\text{MCL}}=10^{-3}$ for Shapes3D and Cars3D, and
 1124 $\lambda_{\text{MCL}}=10^{-4}$ for MPI3D.

1125 D.5. Action-Conditioned Generative World Model

1126 **Pipeline.** Wan-ACG and Wan-LACG are built on a large pretrained video generative model [57], equipped with four
 1127 components that together realize the forward action-conditioned generative model $\mathbf{x}_{t+1} = g_{\text{fwd}}(\mathbf{z}_t, \mathbf{u}_t)$ on which MCL
 1128 operates: (i) a frozen causal VAE encoder E_{vae} that maps a T -frame video window to a latent video sequence $v_{1:F} =$
 1129 $E_{\text{vae}}(\mathbf{x}_{1:T})$; (ii) an inverse dynamics model (IDM) I_{ϕ} that produces, for each frame transition, a diagonal-Gaussian posterior
 1130 $q_{\phi}(z_t | \mathbf{x}_t, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t, \text{diag}(\sigma_t^2))$ over a d_z -dimensional latent action; (iii) an action-interface module A_{ψ} that lifts the
 1131 frame-rate latent action sequence $z_{1:T}$ to FDM-aligned conditioning signals $c_{1:F} = A_{\psi}(z_{1:T})$, injected into the FDM through
 1132 AdaLN modulation; and (iv) a forward dynamics model (FDM) g_{θ} , initialized from the pretrained video-planner backbone,
 1133 that performs flow matching $\hat{u}_{\tau} = g_{\theta}(\tilde{v}_{\tau}, \tau, z_{1:T})$ to predict future latent video dynamics conditioned on the latent action.

1134 **Mapping to the GMD setup.** In our notation, the conditioning signal is the latent action, $\mathbf{u}_t := z_t$ (or, when an external
 1135 command interface is required, a controller-predicted latent action \hat{z}_t), and the observation is the next latent video frame,
 1136 $\mathbf{x}_{t+1} := v_{f+1}$. The action-conditioned FDM is the GMD generator $g(\mathbf{z}_t, \mathbf{u}_t)$: action-driven conditioning enters through
 1137 AdaLN; the diffusion latent enters through \tilde{v}_{τ} .

1138 **Our addition: MCL on the action conditioning.** We attach the MCL objective from Section 4 to the action-conditioning
 1139 channel of the FDM. A mechanism critic $c_{\psi}(\hat{\mathbf{x}}_{t+1}, z_t, \mathbf{u}_t)$ scores how strongly the realized forward mechanism depends on z_t ,
 1140 the mechanism gradient $\mathbf{g}_t = \nabla_{z_t} c_{\psi}(\hat{\mathbf{x}}_{t+1}, z_t, \mathbf{u}_t)$ is contrasted via InfoNCE against the projected latent-action embedding
 1141 $\Pi_u(\mathbf{u}_t)$, and the resulting \mathcal{L}_{MCL} is added to the host objective with weight λ_{MCL} :

$$1142 \mathcal{L}_{\text{total}} = L_{\text{rec}} + \beta L_{\text{KL}} + \lambda_{\text{adv}} L_{\text{GRL}} + \lambda_{\text{MCL}} \mathcal{L}_{\text{MCL}},$$

1143 where L_{rec} is the flow-matching reconstruction loss, L_{KL} is a KL regulariser on the IDM posterior, and L_{GRL} is a gradient-
 1144 reversal embodiment-invariance term used to encourage the latent action to ignore embodiment-specific nuisance. The IDM,
 1145 action interface, FDM, and MCL projection heads are trained jointly; the causal VAE remains frozen.

1146 **What MCL adds, and what we test.** Without MCL, the action-conditioning channel is shaped only *passively* by the
 1147 reconstruction loss; nothing forces distinct latent actions to induce distinct mechanisms in the FDM. MCL closes that gap:
 1148 it makes mechanism diversity along \mathbf{u}_t an *active* target, contrasting how the FDM’s mechanism gradient changes when \mathbf{u}_t
 1149 varies across the in-batch action samples. The structural claim we will verify, against an MCL-free baseline and against a
 1150 non-action-conditioned forward model, is that MCL on the action conditioning improves both visual prediction quality and
 1151 action-following accuracy by enforcing sufficient GMD on the action variable.

1152 E. PyTorch Implementation of MCL

1153 For reproducibility we provide a self-contained PyTorch implementation of the Mechanistic Contrastive Learning loss. We
 1154 first present the *primary* loss (Algorithm 1) that captures the core mechanism-gradient \rightarrow projection \rightarrow InfoNCE pipeline used
 1155 throughout the paper, then give the *full* implementation (Algorithm 2) including the projection-head module `MLPProj` and
 1156 the mechanism-critic module `MechanismCritic` that we use in the image and world-model experiments. The same loss
 1157 applies to every host model in our experiments (MD-VAE, MDDiff, ACG-WM); only the host components change. The loss
 1158 adds to the host’s reconstruction objective with weight λ_{MCL} as $\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_{\text{MCL}} \mathcal{L}_{\text{MCL}}$.

1159 E.1. Primary MCL Loss

1160 E.2. Full Code

1161 The full implementation we use in practice adds the two-layer projection head `MLPProj`, the default mechanism critic
 1162 `MechanismCritic`, and the InfoNCE wrapper used in our image and world-model experiments.

Algorithm 1. Primary PyTorch implementation of the MCL loss. The host generator `decoder_G` produces \hat{x} from a batch of latents and conditioning signals; the critic scores how strongly the realized mechanism depends on z ; the autograd of the critic w.r.t. z yields the per-sample mechanism gradient, which is contrasted via InfoNCE against the projected conditioning embedding.

```
def mcl_loss(decoder_G, z, u, *,
            critic, Pi_g, Pi_u, tau=0.1):
    """
    decoder_G : callable (z, u) -> x_hat
    z         : latent batch, shape (B, ...)
    u         : conditioning batch, shape (B, d_u)
    critic    : (x_hat, z, u) -> mechanism score per sample
    Pi_g, Pi_u: projection heads to shared embedding space
    """
    z = z.requires_grad_(True)
    x_hat = decoder_G(z, u)
    s = critic(x_hat, z, u)

    # mechanism gradient w.r.t. the latent
    g = torch.autograd.grad(s.sum(), z, create_graph=True)[0]

    # L2-normalized projections
    q = Pi_g(g.flatten(1)); q = q / q.norm(dim=1, keepdim=True)
    k = Pi_u(u);           k = k / k.norm(dim=1, keepdim=True)

    # B x B InfoNCE: positives on the diagonal
    logits = q @ k.t() / tau
    labels = torch.arange(q.size(0), device=q.device)
    return F.cross_entropy(logits, labels)
```

F. Additional Experimental Results

1163

F.1. Effect of Mechanism Diversity

1164

A central claim of Theorem 1 is that identifiability strengthens monotonically with the size of the conditioning alphabet $|\mathcal{U}|$: more distinct mechanisms produce more linearly-independent constraints on the generator Jacobian, eventually saturating the sufficient GMD condition (Assumption 2, which requires $|\mathcal{U}| \geq 2d_z + 1$). Table 5 reports MD-VAE’s median MCC across $|\mathcal{U}| \in \{2, 4, 8, 16, 32\}$ at $d_z \in \{5, 10, 20\}$, fixing all other settings of Appendix D.3.

1165

1166

1167

1168

Table 5. Effect of conditioning-set size $|\mathcal{U}|$ on MD-VAE’s median MCC (\uparrow , 30 seeds).

	$ \mathcal{U} = 2$	$ \mathcal{U} = 4$	$ \mathcal{U} = 8$	$ \mathcal{U} = 16$	$ \mathcal{U} = 32$
$d_z = 5$	0.51	0.79	0.96	0.97	0.97
$d_z = 10$	0.43	0.62	0.84	0.93	0.94
$d_z = 20$	0.39	0.51	0.71	0.86	0.92

The rows confirm two predictions of the theory: (i) for fixed d_z , MCC rises with $|\mathcal{U}|$, mirroring the increase in mechanism diversity; (ii) higher-dimensional latents require more conditioning values before performance saturates. The threshold $|\mathcal{U}| \geq 2d_z + 1$ in Assumption 2 is sufficient rather than necessary, so low-noise simulations can begin to saturate before the conservative bound is reached. The qualitative behaviour still matches the GMD identifiability principle: *contrasting more mechanisms produces more identifiable representations.*

1169

1170

1171

1172

1173

1174 **F.2. Generation Quality of MDDiff**

1175 A representation is only useful if the generator that produces it remains a competent image model. Table 6 reports FID
 1176 (sample quality) and LPIPS (reconstruction fidelity) for MDDiff and the strongest diffusion baseline EncDiff [61] on the
 1177 three benchmarks of Section 5.2. MDDiff preserves generation quality while delivering the disentanglement gains reported in
 1178 Table 1; on *Cars3D* it actually *improves* FID and LPIPS, and on *Shapes3D* the two methods are statistically tied. *MPI3D*
 1179 shows a modest FID gap that we attribute to the smaller λ_{MCL} used for that dataset (Appendix D); LPIPS is unchanged.

Table 6. **Generation quality (FID ↓, LPIPS ↓).**

Method	Cars3D		Shapes3D		MPI3D	
	FID ↓	LPIPS ↓	FID ↓	LPIPS ↓	FID ↓	LPIPS ↓
EncDiff	62.31 ± 3.28	0.07 ± 0.03	5.33 ± 0.68	0.01 ± 0.01	4.06 ± 0.60	0.01 ± 0.01
MDDiff (Ours)	39.32 ± 6.67	0.03 ± 0.02	5.52 ± 0.78	0.01 ± 0.01	7.10 ± 2.12	0.01 ± 0.01

1180 **F.3. Latent Traversals**

1181 Latent traversals offer a complementary view of factoredness to the factor-swapping results in the main paper (Figure 6):
 1182 holding all other dimensions fixed and sweeping one dimension across its range should produce smooth, single-factor change
 1183 in the generated image. Figure 8 reports MDDiff traversals on all three benchmarks; each row corresponds to one latent
 1184 dimension and isolates one semantic factor (wall colour, floor colour, object colour, shape, orientation, scale on *Shapes3D*;
 1185 analogous factors on *Cars3D* and *MPI3D*).

1186 **F.4. Cross-attention Maps**

1187 Cross-attention maps complement latent traversals by exposing how mechanism diversity manifests *inside* the generator: each
 1188 concept token (one per latent dimension) modulates the denoising U-Net through cross-attention, and each token’s attention
 1189 should localise to the spatial region of its factor. Figure 9 visualises the cross-attention weights between concept tokens
 1190 (columns) and spatial locations across the three benchmarks; each token attends to the spatial region corresponding to its
 1191 factor, evidence that MCL produces semantically meaningful, spatially-localised modulation in the generator.

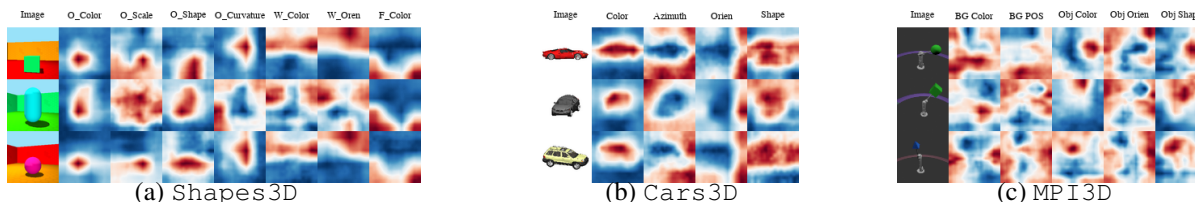


Figure 9. **Cross-attention maps on all three benchmarks.** Columns show concept tokens; rows show spatial locations in the denoising network.

1192 **G. Extended Related Work**1193 **G.1. Why Existing Generative Models Can Be Factored**

1194 Theorem 1 explains why conditioning makes generative models learn factored representations, and when they fail. Two major
 1195 families instantiate the GMD regime (Figure 3, bottom row). In *conditional image generation*, text-to-image models [41, 50, 52]
 1196 and diffusion-based disentanglement methods [60, 61] let text or concept tokens reshape the denoising network at every
 1197 layer; the latent code \mathbf{z} encodes stochastic content while \mathbf{u} controls semantics, creating mechanism diversity that factors the
 1198 latent space. In *generative world models* [15, 16, 18, 69], actions \mathbf{a}_t reshape the latent dynamics model $g(\mathbf{z}_t, \mathbf{a}_t)$; similarly,
 1199 decision-as-generation methods [1, 7, 28] condition trajectory generation on returns or skills. Identifiability in world models is
 1200 particularly important: a factored latent state means each dimension captures a physically meaningful quantity (object position,
 1201 velocity, contact state), enabling compositional generalization to novel scenes, faithful physical reasoning, and interpretable
 1202 planning. In both families, sufficient diversity in \mathbf{u} provably factors the latent state by Theorem 1. However, none of these

models were designed with GMD in mind; whether Assumption 2 holds depends entirely on the data and architecture, with no explicit control. This motivates the active approach of Section 4. 1203
1204

G.2. Theoretical Foundations of LRD: Auxiliary-Variable Identifiability 1205

The classical *LRD* identifiability result of Hyvarinen et al. [27] treats the case where an auxiliary variable \mathbf{u} shifts the latent distribution while the generator g stays fixed. Translating their notation into ours (\mathbf{z} for source, \mathbf{u} for auxiliary, \mathbf{x} for observation, g for mixing function), the setup is: 1206
1207
1208

Setup. Observations are produced by an invertible mixing $\mathbf{x} = g(\mathbf{z})$, with $g \in \mathcal{C}^2$ and $\mathbf{z} \in \mathbb{R}^{d_z}$. The latent density is conditionally factorial given the auxiliary, 1209
1210

$$p(\mathbf{z} | \mathbf{u}) = \prod_{i=1}^{d_z} p_i(z_i | \mathbf{u}), \quad 1211$$

and each one-dimensional conditional belongs to an exponential family of order k : 1212

$$\log p_i(z_i | \mathbf{u}) = \sum_{j=1}^k T_{i,j}(z_i) \lambda_{i,j}(\mathbf{u}) - \log Z_i(\mathbf{u}) + \log Q_i(z_i), \quad 1213$$

where $T_{i,j} : \mathbb{R} \rightarrow \mathbb{R}$ are the sufficient statistics, $\lambda_{i,j}(\mathbf{u}) \in \mathbb{R}$ the natural parameters, $Z_i(\mathbf{u})$ the normalizer, and Q_i a base measure. Stacking the natural parameters yields $\boldsymbol{\lambda}(\mathbf{u}) \in \mathbb{R}^{kd_z}$. 1214
1215

Sufficient diversity (Assumption of 27, Theorem 3). There exist $kd_z + 1$ values $\mathbf{u}^{(0)}, \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(kd_z)}$ such that the matrix of contrasts 1216
1217

$$L = [\boldsymbol{\lambda}(\mathbf{u}^{(1)}) - \boldsymbol{\lambda}(\mathbf{u}^{(0)}), \dots, \boldsymbol{\lambda}(\mathbf{u}^{(kd_z)}) - \boldsymbol{\lambda}(\mathbf{u}^{(0)})] \in \mathbb{R}^{kd_z \times kd_z} \quad 1218$$

is invertible. 1219

Result. Under the setup and the sufficient-diversity condition, the latent code \mathbf{z} is identifiable from \mathbf{x} up to a permutation and a componentwise invertible transformation [27, Theorem 3]. 1220
1221

Why this is LRD. Identifiability hinges entirely on the auxiliary \mathbf{u} *shifting the latent distribution* $p(\mathbf{z} | \mathbf{u})$ across enough values; the generator g is fixed and contributes no \mathbf{u} -dependent variability. This is exactly the LRD principle of Section 2. Our GMD condition (Theorem 1) replaces the diversity requirement on $p(\mathbf{z} | \mathbf{u})$ with an analogous condition on the generator Jacobian $J_g(\mathbf{z}, \mathbf{u})$, which is what makes our identifiability result strictly more expressive in the conditional generative regime (Proposition 1). 1222
1223
1224
1225
1226

G.3. Theoretical Foundations of Regularization: Sparsity-Based Identifiability 1227

The classical *regularization-based* identifiability result of Zheng et al. [68] obtains identifiability without any auxiliary variable, by imposing a structural sparsity constraint on the generator instead. We restate it in the notation of our paper. 1228
1229

Setup. Observations are produced by an invertible mixing $\mathbf{x} = g(\mathbf{z})$, $g \in \mathcal{C}^2$, with $\mathbf{z} \in \mathbb{R}^{d_z}$ drawn from a fixed prior $p(\mathbf{z})$ that does *not* depend on any auxiliary variable. Identifiability is recovered through structural assumptions on g and on $p(\mathbf{z})$. 1230
1231

Structural sparsity assumption (68, Definition 3). The Jacobian $J_g(\mathbf{z})$ of the mixing function has a structural sparsity pattern: there exists a binary support $M \in \{0, 1\}^{d_x \times d_z}$ such that $\partial g_a / \partial z_i = 0$ whenever $M_{a,i} = 0$, for $p(\mathbf{z})$ -a.e. \mathbf{z} . Each row of M is constrained to have at most a small number of non-zero entries, and the column supports are required to be *distinguishable* so that no two latent coordinates have the same set of dependent observation coordinates. 1232
1233
1234
1235

Non-degeneracy conditions (68, Theorem 3.5). The following two conditions ensure that latent coordinates are non-redundant: (i) the prior $p(\mathbf{z})$ is fully supported on \mathbb{R}^{d_z} and admits a \mathcal{C}^1 density, and (ii) the second-order partial derivatives $\partial^2 g_a / \partial z_i \partial z_j$ satisfy a generic linear-independence condition across the support of $p(\mathbf{z})$, ruling out degenerate mixing functions. 1236
1237
1238
1239

1240 **Result.** Under the setup, the structural sparsity assumption, and the non-degeneracy conditions, the latent code \mathbf{z} is identifiable
1241 from \mathbf{x} up to a permutation and a componentwise invertible transformation [68, Theorem 3.5].

1242 **Why this is regularization.** Identifiability hinges on a structural *constraint* on the generator, in this case sparsity of J_g ,
1243 rather than on auxiliary diversity. This is exactly the Regularization principle of Section 2, of which L_1/L_2 penalties on the
1244 generator Jacobian are the practical instantiation. Our GMD condition does not require sparsity of J_g ; it instead asks that J_g
1245 vary sufficiently across \mathbf{u} , which is a complementary structural property and the basis of MCL.

1246 G.4. Categorizing Methods by Principle

1247 Table 7 categorizes representative methods according to which of the three principles for factored representations they employ
1248 (cf. Figure 2). A key distinction emerges: many conditional generative models *passively* operate in the GMD regime (the
1249 conditioning signal enters the generator architecturally), but none *actively* encourage sufficient mechanism diversity across
1250 conditioning signals. Our MCL is the first method that explicitly enforces the sufficient GMD condition (Assumption 2),
1251 turning passive architectural mechanism diversity into a controllable, provable guarantee.

1252 **LRD \cap GMD.** Methods in this intersection let \mathbf{u} enter both the latent distribution and the generator. Conditional VAEs [53]
1253 have this structure by design: the prior $p(\mathbf{z} | \mathbf{u})$ and the decoder $p(\mathbf{x} | \mathbf{z}, \mathbf{u})$ both depend on \mathbf{u} . Classifier-guided diffusion [9]
1254 adds a classifier gradient to the sampling trajectory, shifting the effective latent distribution, on top of a class-conditional
1255 denoiser that reshapes the generator. DALL-E 2 [43] factorizes generation as $P(\mathbf{x} | \mathbf{y}) = P(\mathbf{x} | \mathbf{z}, \mathbf{y}) P(\mathbf{z} | \mathbf{y})$: the text
1256 caption \mathbf{y} enters both the diffusion prior over CLIP embeddings and the decoder. VQ-VAE-2 [44] combines all three principles:
1257 vector quantization, top-level codes conditioning the autoregressive prior over bottom-level codes, and top-level codes entering
1258 the feed-forward decoder.

1259 **Regularization \cap GMD.** Methods here apply explicit or implicit regularization to a conditional generator. EncDiff [61]
1260 conditions a latent diffusion model on concept tokens via cross-attention, with the diffusion information bottleneck serving as
1261 implicit regularization. DisDiff [60] decomposes the score into per-factor sub-gradient fields, a structural sparsity constraint,
1262 and conditions each field on a discovered factor. The Hessian Penalty [42] regularizes the off-diagonal Hessian of a
1263 GAN generator with respect to its latent input, encouraging axis-aligned disentanglement. VQ-VAE [54] combines vector
1264 quantization with a decoder conditioned on discrete codes. InfoDiffusion [59] adds a mutual-information maximization penalty
1265 to a latent-conditioned diffusion model to prevent posterior collapse.

1266 **LRD \cap Regularization.** Methods here combine auxiliary-variable-driven distributional shifts with structural constraints.
1267 Lachapelle et al. [33] use temporal or interventional auxiliary variables together with binary-mask sparsity on the causal graph;
1268 neither alone suffices for identifiability. SlowVAE [31] uses temporal context as the auxiliary variable and a Laplace sparse
1269 prior on temporal differences. CITRIS [36] observes intervention targets and constrains the transition model to respect a
1270 factored causal structure. TDRL [63] leverages distribution shifts across environments with a modular, sparse mechanism-shift
1271 assumption.

1272 H. Recommendations for Future Work

1273 This work introduces *mechanism diversity* as a sufficient condition for factored generative models, and operationalises it
1274 through Mechanistic Contrastive Learning (MCL). We see three open directions where the same principle has the potential to
1275 extend beyond the identifiability benchmarks studied here.

1276 H.1. Beyond Factored Generative Models

1277 **Motivation.** The notion of *factoredness* in Section 3.2, inherited from nonlinear ICA, is a stricter target than what most
1278 modern conditional generators are built or evaluated against. Large text-to-image diffusion [41, 50, 52], controllable video
1279 diffusion [3], visual world models [16, 69], and multi-modal generators routinely operate without ground-truth factors yet still
1280 benefit from compositional, controllable behaviour at the conditioning interface.

Table 7. **Categorizing methods by principle.** • = *active* (an explicit loss term encourages the principle); ◦ = *passive* (satisfied by model architecture/design only, no dedicated loss).

Region	Method	LRD	Reg.	GMD	Mechanism
<i>LRD only</i>	iVAE [29]	◦			Conditional prior $p(\mathbf{z} \mathbf{u})$ (architectural)
	TCL [25]	•			Contrastive classification loss over time segments
	CPC [40]	•			InfoNCE contrastive loss over temporal windows
<i>Reg. only</i>	β -VAE [21]		•		Upweighted KL penalty ($\beta > 1$)
	FactorVAE [30]		•		TC penalty via learned discriminator
	Sparse ICA [68]		•		L_1 sparsity on generator Jacobian
	Tripod [22]		•		Quantization + MI + Hessian penalty losses
<i>Active GMD</i>	MDDiff / MCL (Ours)			•	Contrastive loss on generator Jacobian gradients
	MD-VAE (Ours)			•	Conditional VAE with mechanism-diverse decoder
<i>LRD \cap GMD</i>	CVAE [53]	◦		◦	Prior $p(\mathbf{z} \mathbf{u})$ + decoder $p(\mathbf{x} \mathbf{z}, \mathbf{u})$ (both architectural)
	Classifier Guidance [9]	•		◦	Classifier loss shifts sampling; denoiser cond. on y
	DALL-E 2 [43]	◦		◦	Diffusion prior $P(\mathbf{z} \mathbf{y})$ + decoder $P(\mathbf{x} \mathbf{z}, \mathbf{y})$
<i>Reg. \cap GMD</i>	EncDiff [61]		◦	◦	Architectural bottleneck + cross-attn conditioning
	DisDiff [60]		•	◦	Disentangling loss + per-factor sub-gradient fields
	Hessian Penalty [42]		•	◦	Hessian reg. loss + class-conditional GAN
	VQ-VAE [54]		•	◦	VQ + commitment losses + code-conditioned decoder
	InfoDiffusion [59]		•	◦	MI maximization loss + latent-conditioned diffusion
<i>All three</i>	VQ-VAE-2 [44]	◦	•	◦	VQ loss (active reg.) + cond. prior + cond. decoder
<i>LRD \cap Reg.</i>	Mech. Sparsity [33]	◦	•		Cond. prior (passive) + L_1 mask sparsity (active)
	SlowVAE [31]	◦	•		Temporal prior (passive) + Laplace sparse penalty
	CITRIS [36]	◦	•		Intervention prior (passive) + factored assignment loss
	TDRL [63]	◦	•		Shift prior (passive) + modular mechanism loss

Conjecture. The mechanism-diversity principle should still apply in this regime: whenever \mathbf{u} *varies sufficiently* across instances, the per-condition Jacobian family is structurally distinguishable, which is precisely what MCL contrasts. We expect an MCL-style objective added to a controllable generator to improve compositional generalization, prompt following, and conditioning-slot reusability, even without an identifiability proof, in the same way that mechanism contrast already improves

1281
1282
1283
1284

1285 disentanglement scores at the benchmark scale (Table 1).

1286 **Open question.** *What should factoredness mean for a billion-parameter text-to-image model*, and which theoretical
1287 guarantees survive the absence of canonical ground-truth factors?

1288 H.2. MCL as Fine-tuning of Large Pretrained Models

1289 **Motivation.** Pretraining of large generative models is necessarily *passive*: the conditioning signal is whatever caption,
1290 label, or action happens to accompany each datapoint at web scale, and no explicit objective enforces mechanism contrast.
1291 Fine-tuning is where domain-specific expert knowledge enters the model, through methods such as LoRA [23], ControlNet [67],
1292 and DreamBooth [51]. These methods add new conditioning channels but do not certify that distinct conditioning values
1293 produce distinct mechanisms.

1294 **Position.** MCL is the natural *active* counterpart: a contrastive signal that explicitly enforces mechanism shift across
1295 conditioning values. This paper already demonstrates the recipe at the disentanglement scale: we first pretrain MDDiff as a
1296 \mathbf{u} -conditioned generator on disentangled concepts, and then apply MCL as a fine-tuning objective on top of that pretrained
1297 MDDiff checkpoint to drive both DCI and FactorVAE upward (Figure 5, Section 5.2). Scaling to frontier diffusion and
1298 world-model backbones should marry the *data diversity* captured by pretraining with the *mechanism diversity* needed for
1299 compositional control, yielding a sample-efficient route to factored controllable generators without retraining.

1300 **Open questions.** Which layers of a U-Net or DiT should host the projection heads? How should MCL be balanced against
1301 the original denoising loss? Is parameter-efficient fine-tuning of the projection-head subspace alone sufficient?

1302 H.3. Multi-objective, Multi-auxiliary Mechanism Contrast

1303 **Motivation.** Realistic generative pipelines pair *several auxiliary variables* (text caption, robot action, audio cue, regime
1304 index) with *several generative objectives* (image, video, control trajectory, language). Examples include video-action
1305 world models [69], manipulation diffusion policies [7], and decision-as-generation pipelines [1, 28]. MCL as defined in
1306 Section 4 contrasts one mechanism gradient against one conditioning embedding, which is well-matched to single-auxiliary,
1307 single-objective settings but leaves the multi-axis setting open.

1308 **Open questions.** Should each auxiliary $\mathbf{u}^{(k)}$ carry its own projection head and InfoNCE term, or should the contrastive
1309 signal be *jointly normalized* across auxiliaries? Does the sufficient-GMD condition decompose objective-by-objective, or does
1310 identifiability require simultaneous diversity across all $\mathbf{u}^{(k)}$? How does the mechanism-shift principle behave when auxiliaries
1311 are *causally entangled*, e.g. an action partially explains the next observation while a text caption only describes high-level
1312 semantics?

Algorithm 2. Full PyTorch implementation of MCL: the projection-head module, default mechanism-critic module, and the InfoNCE mechanism-gradient loss with optional separate keys.

```

import torch
import torch.nn as nn
import torch.nn.functional as F

def l2norm(x: torch.Tensor, eps: float = 1e-8) -> torch.Tensor:
    return x / (x.norm(dim=1, keepdim=True) + eps)

def info_nce_from_qk(q, k, tau: float = 0.1) -> torch.Tensor:
    logits = (q @ k.t()) / (tau + 1e-12)
    labels = torch.arange(q.size(0), device=q.device)
    return F.cross_entropy(logits, labels)

class MLPProj(nn.Module):
    def __init__(self, in_dim, out_dim=128, layernorm=False):
        super().__init__()
        layers = []
        if layernorm:
            layers.append(nn.LayerNorm(in_dim))
        layers += [nn.Linear(in_dim, out_dim),
                  nn.ReLU(inplace=True),
                  nn.Linear(out_dim, out_dim)]
        self.net = nn.Sequential(*layers)

    def forward(self, x):
        return self.net(x)

class MechanismCritic(nn.Module):
    def __init__(self, z_shape=(3, 16, 16), u_dim=20, hidden=256):
        super().__init__()
        zc, zh, zw = z_shape
        self.img = nn.Sequential(
            nn.Conv2d(3, 64, 4, 2, 1), nn.ReLU(True),
            nn.Conv2d(64, 128, 4, 2, 1), nn.ReLU(True),
            nn.AdaptiveAvgPool2d(1),
        )
        self.z_fc = nn.Linear(zc * zh * zw, hidden)
        self.u_fc = nn.Linear(u_dim, hidden)
        self.out = nn.Sequential(nn.ReLU(True), nn.Linear(hidden, 1))

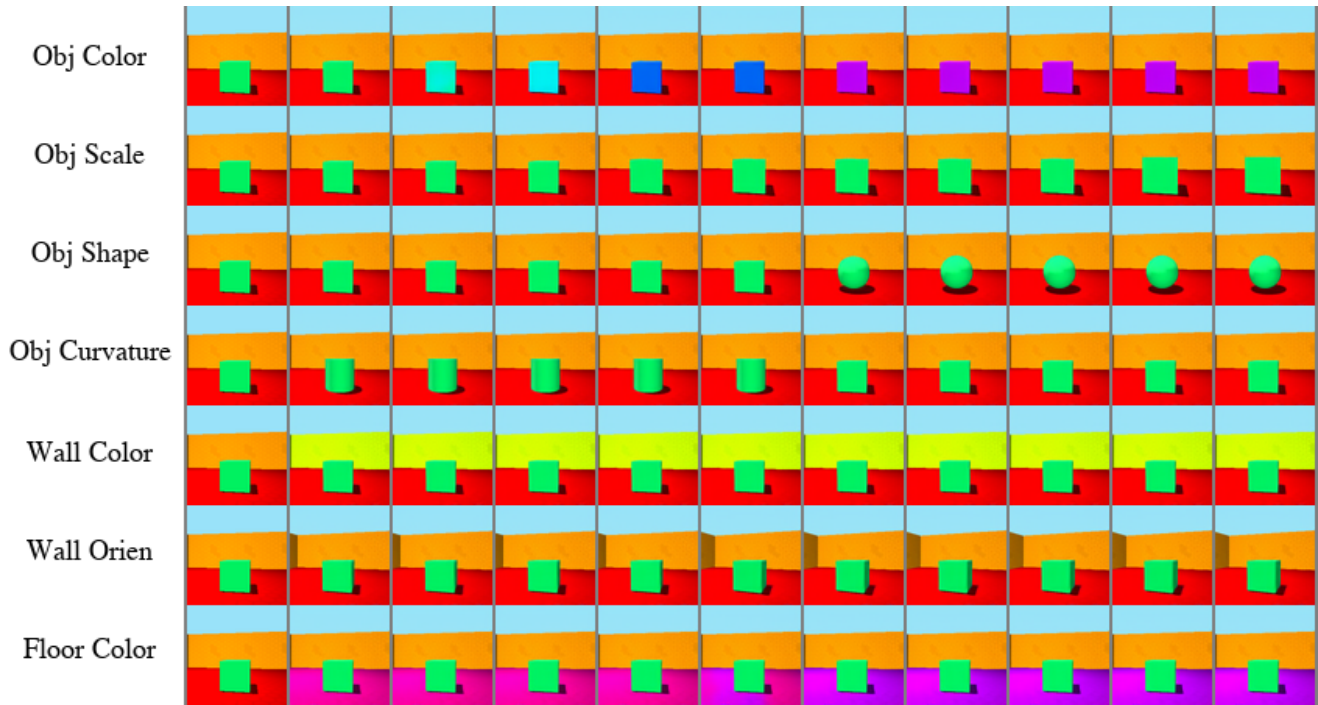
    def forward(self, x_hat, z, u):
        img_feat = self.img(x_hat).flatten(1)
        z_feat = self.z_fc(z.flatten(1))
        u_feat = self.u_fc(u)
        if img_feat.size(1) < z_feat.size(1):
            img_feat = F.pad(img_feat,
                             (0, z_feat.size(1) - img_feat.size(1)))
        else:
            img_feat = img_feat[:, :z_feat.size(1)]
        return self.out(z_feat + u_feat + img_feat).squeeze(1)

def mcl_infonce_mechgrad_loss(decoder_G, z, u_key, *,
                             u_for_G=None,
                             critic, Pi_g, Pi_u,
                             tau: float = 0.1,
                             create_graph: bool = False):
    """InfoNCE mechanism-gradient loss with optional separate keys.

    u_key : auxiliary used as the contrastive key
    u_for_G : conditioning fed to decoder/critic (defaults to u_key)
    """
    u = u_key
    uG = u_for_G if u_for_G is not None else u
    z = z.requires_grad_(True)
    x_hat = decoder_G(z, uG)
    s = critic(x_hat, z, uG)
    g = torch.autograd.grad(s.sum(), z, create_graph=create_graph)[0]

    q = l2norm(Pi_g(g.flatten(1)))
    k = l2norm(Pi_u(u))
    return info_nce_from_qk(q, k, tau=tau)

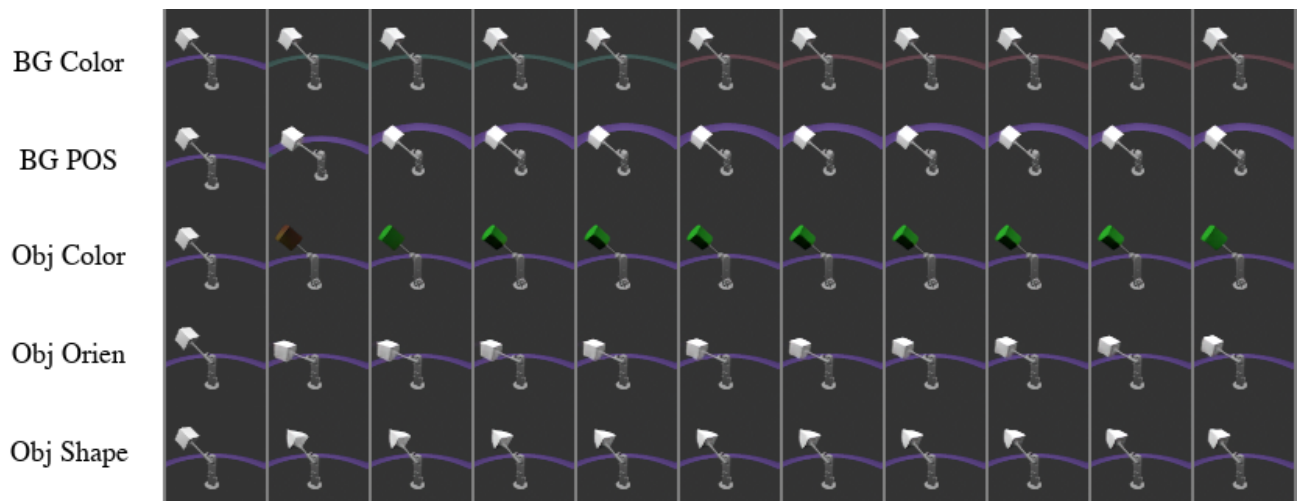
```



(a) Shapes3D



(b) Cars3D



(c) MPI3D

Figure 8. **Latent traversals on all three benchmarks.** Each row varies a single MDDiff latent dimension while keeping the others fixed; columns sweep the dimension across its range.