

MATLAT: Material Latent Space for PBR Texture Generation

Anonymous CVPR workshop submission

Paper ID



Figure 1. **PBR Textures Generated by MATLAT.** Our method produces PBR textures that accurately represent rough materials (left), metallic surfaces (middle), and complex mixed materials (right).

Abstract

001 Text-guided PBR texture generation is essential for re-
002 lightable 3D asset creation, but high-quality PBR supervi-
003 sion is scarce and pretrained RGB latent priors do not read-
004 ily extend to roughness and metallic channels. We present
005 MATLAT, a framework that learns a material latent space
006 by fine-tuning a pretrained latent diffusion model for multi-
007 view PBR generation. Unlike prior methods that freeze
008 the pretrained encoder and suffer from distribution shifts
009 when encoding additional PBR channels, MATVAE fine-
010 tunes the encoder to incorporate material channels through
011 residual latent prediction while regularizing the resulting
012 posterior toward the original RGB latent distribution. We
013 further show that correspondence-aware attention alone is
014 insufficient for cross-view consistency unless the latent-to-
015 image mapping preserves spatial locality; to enforce this,
016 we introduce a regularization that crops latent patches, de-
017 codes them, and aligns the results with corresponding im-
018 age regions.

On Objaverse-XL, MATLAT improves the strongest
multi-view baseline from 6.309 to 3.083 shaded FID and
from 9.630 to 4.599 albedo FID while retaining competitive
roughness and metallic accuracy. Ablations confirm that
both distributional alignment and locality preservation are
individually necessary for high-fidelity PBR texture synthe-
sis.

019
020
021
022
023
024
025

1. Introduction

Generative modeling has driven remarkable progress in 3D
asset creation, spanning geometry [2, 12, 20, 31], tex-
tures [1, 13, 21, 30], motions [8, 19, 25], and complete
scenes [11, 15, 34]. Among these, Physically Based Ren-
dering (PBR) textures—comprising albedo, roughness, and
metallic channels—are essential for production-ready as-
sets, as they enable physically accurate relighting under ar-
bitrary illumination. Recent PBR texture generation meth-
ods [1, 3, 21, 26] predict material maps for relightable as-
sets, yet progress is constrained by the scarcity of large-

026
027
028
029
030
031
032
033
034
035
036

037 scale PBR datasets [4, 5]. A natural remedy is to lever- 089
 038 age pretrained image generative models with strong priors 090
 039 from large-scale RGB data [22]. Among existing strategies, 091
 040 multi-view generation followed by mesh projection [9, 10] 092
 041 has proven most effective, as SDS-based approaches [16,
 042 20, 27] often produce saturation artifacts. However, this
 043 multi-view approach still faces two challenges for PBR syn-
 044 thesis: (1) directly applying diffusion models trained on
 045 RGB latents to PBR generation is non-trivial, since sim-
 046 ply mapping additional material channels to the pretrained
 047 encoder introduces a substantial domain gap; and (2) pre-
 048 serving multi-view consistency remains difficult, with fail-
 049 ures causing blurring and artifacts in overlapping regions
 050 upon unprojection. To address both challenges, we intro-
 051 duce MATLAT, a framework that learns a **Material Latent**
 052 space for high-quality PBR texture generation. Our two-
 053 stage pipeline first fine-tunes the pretrained VAE on PBR
 054 images to obtain Material VAE (MATVAE), then fine-tunes
 055 a diffusion model for multi-view material generation in the
 056 adapted latent space. Our first contribution is a latent-space
 057 adaptation module in MATVAE that extends the pretrained
 058 space to incorporate roughness and metallic channels. Un-
 059 like prior works [3, 7, 9, 10, 28, 36, 37] that freeze the
 060 encoder and use zero-channel padding—causing distribu-
 061 tional mismatch and suboptimal diffusion fine-tuning—we
 062 fine-tune the encoder with a distributional regularization
 063 that constrains deviations from the original latent distribu-
 064 tion. While similar extensions exist for transparent [33] or
 065 RGB-D synthesis [14], this is the *first* adaptation to PBR
 066 texture generation, and we identify crucial components for
 067 its practical effectiveness. Our second contribution is lo-
 068 cality regularization in MATVAE, which enforces spatial
 069 alignment between latent tokens and image pixels. This
 070 alignment is subsequently leveraged by correspondence-
 071 aware attention (CAA) [24] in the diffusion model to en-
 072 force multi-view consistency. While spatial locality largely
 073 holds for pretrained RGB encoders [9], it often breaks for
 074 additional material channels. We address this by cropping
 075 latent patches, decoding them, and applying an ℓ_2 recon-
 076 struction loss against corresponding image regions. Exper-
 077 iments show that MATLAT outperforms baselines trained
 078 from scratch, SDS-based methods, and multi-view diffusion
 079 models on most quantitative metrics, achieving state-of-the-
 080 art performance. Ablation studies validate each proposed
 081 component.

082 2. Method

083 Given a mesh \mathcal{M} and text prompt y , our goal is to gener- 128
 084 ate multi-view PBR material images $\mathbf{x}_i = [\mathbf{a}_i, \mathbf{r}_i, \mathbf{m}_i]$ 129
 085 (albedo/roughness/metallic) for viewpoints $\{c_i\}_{i=1}^N$, and 130
 086 then project them to a consistent texture map. Our pipeline 131
 087 has two stages: (1) prior-preserving material latent learning 132
 088 with MATVAE, and (2) multi-view diffusion fine-tuning in 133

the learned latent space, so that the first stage restores com-
 patibility with pretrained RGB priors and the second stage
 exploits that compatible latent space for multi-view genera-
 tion.

Prior-preserving material latent. Prior multi-view PBR
 pipelines often rely on a frozen RGB VAE [9], which en-
 codes albedo reasonably well but places roughness and
 metallic maps off the pretrained latent manifold. We in-
 stead adapt the latent space while explicitly anchoring it to
 pretrained RGB statistics.

Let $q(\mathbf{z}_{\text{base}} | \mathbf{a}) = \mathcal{N}(\boldsymbol{\mu}_{\text{base}}, \boldsymbol{\sigma}_{\text{base}}^2)$ be the pretrained la-
 tent distribution from albedo. We add a residual encoder \mathcal{E}_{res}
 that predicts $(\boldsymbol{\mu}_{\text{res}}, \boldsymbol{\sigma}_{\text{res}})$ from full PBR input \mathbf{x} and com-
 pose:

$$q(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\text{base}} + \boldsymbol{\mu}_{\text{res}}, \boldsymbol{\sigma}_{\text{base}}^2 \odot \boldsymbol{\sigma}_{\text{res}}^2). \quad (1)$$

To preserve compatibility with pretrained diffusion pri-
 ors, we apply distribution alignment:

$$\mathcal{L}_{\text{reg}} = \lambda_{\text{reg}} \text{KL}(q(\mathbf{z} | \mathbf{x}) \| q(\mathbf{z}_{\text{base}} | \mathbf{a})). \quad (2)$$

Combined with standard VAE reconstruction, KL, and
 discriminator losses, this objective gradually injects rough-
 ness and metallic information into an RGB-aligned latent
 manifold. Compared with residual-only prediction [33]
 or direct KL-aligned prediction [14], MATVAE combines
 residual initialization and distribution alignment in a single
 stochastic posterior.

Locality-preserved multi-view diffusion. We fine-tune
 a multi-view diffusion backbone in the learned mate-
 rial latent space and adopt correspondence-aware attention
 (CAA) [24] so that each latent token aggregates features
 only from geometrically corresponding regions in other
 views (Fig. 2(a)). However, CAA is only meaningful if the
 adapted latents remain locally aligned with decoded pixels.
 We therefore enforce locality during MATVAE training by
 matching cropped image patches before and after encode-
 decode (Fig. 2(b)):

$$\mathcal{L}_{\text{local}} = \lambda_{\text{local}} d(\mathcal{T}(\mathbf{x}), \mathcal{D}(\mathcal{T}(\mathcal{E}(\mathbf{x}))). \quad (3)$$

This patch-level constraint keeps decoded pixels tied to spa-
 tially aligned latent tokens, reducing projection seams and
 making CAA translate into real pixel-space consistency.

Training and inference. We first train MATVAE on ren-
 dered PBR material images, where the pretrained encoder
 is kept frozen and only the residual encoder \mathcal{E}_{res} and de-
 coder \mathcal{D} are optimized. The final convolution of \mathcal{E}_{res} is ze-
 ro-initialized so that training begins at zero residual, exactly
 reproducing the pretrained latent and allowing the material

134 channels to be incorporated gradually. We then fine-tune
135 the multi-view diffusion model with CAA modules inserted
136 alongside the original attention layers, using a conditional
137 flow-matching objective [17] in the learned material latent
138 space. At inference, the model generates six canonical-
139 view PBR images that are projected and blended into tex-
140 ture space for relightable mesh rendering.

141 3. Experiments

142 **Setup.** We fine-tune the STABLE DIFFUSION 3.5-
143 MEDIUM VAE and diffusion checkpoints. Training uses
144 40,723 Objaverse-XL assets [4] with captions from Boot-
145 strap3D and Cap3D [18, 23]. For each training mesh we
146 render 26 fixed views, and at inference we generate six
147 canonical views (front, back, left, right, top, and bottom)
148 following prior multi-view texturing methods. For evalua-
149 tion, we follow PacTure [7] on 128 held-out meshes with
150 20 views per mesh: eight at 0° elevation, six at 30° , four at
151 45° , and two polar views at $\pm 90^\circ$. For each view, we ren-
152 der shaded RGB images under matched HDR environment
153 maps as well as the corresponding albedo, roughness, and
154 metallic maps.

155 **Evaluation metrics.** We report shaded/albedo FID and
156 KID in CLIP feature space, CLIP text-image similarity, and
157 roughness/metallic RMSE. Because these metrics do not
158 directly measure multi-view agreement, we also report c-
159 PSNR, which computes PSNR between each pixel and its
160 geometric correspondences across views. Runtime in sec-
161 onds captures whether the method is practical for iterative
162 text-guided asset creation.

163 **Main comparison.** We compare against three baseline
164 families using official released checkpoints: trained-from-
165 scratch methods (SC) [3, 29], SDS-based optimization
166 pipelines (SDS) [6, 32, 35], and multi-view diffusion meth-
167 ods (MV) [9, 10]. As prior work differs substantially in
168 pretraining, datasets, and optimization procedures, a strictly
169 unified retraining protocol is impractical; we therefore ren-
170 der and evaluate all methods under the same pipeline and
171 treat the comparison as a realistic positioning against exist-
172 ing systems.

173 Table 1 shows that MATLAT achieves the best shaded
174 and albedo fidelity among all compared methods. Com-
175 pared with MaterialMVP [9], shaded FID improves from
176 6.309 to 3.083 and albedo FID from 9.630 to 4.599, while
177 retaining competitive roughness and metallic RMSE. SC
178 methods generally exhibit suboptimal performance, indi-
179 cating that limited PBR supervision alone is insufficient
180 to learn both material realism and text alignment. SDS
181 methods benefit from pretrained 2D priors but, aside from
182 DreamMat [35] which incurs a prohibitive runtime of

approximately 2400s, show only marginal improvements
over SC baselines. Among multi-view diffusion meth-
ods, which are overall the strongest baseline family, MAT-
LAT outperforms on most metrics, with the exception of
the albedo CLIP score which is marginally behind Mate-
rialAnything [10].

183 **Qualitative results.** Figure 3 presents representative out-
184 puts of MATLAT across diverse prompts and geometries,
185 spanning glossy food, mixed-material footwear, patterned
186 upholstery, and wood-and-metal props. Across these cases,
187 MATLAT preserves sharp local boundaries and plausible
188 albedo/roughness/metallic interactions, yielding coherent
189 relighting and materially distinct parts within a single as-
190 set.

191 **Ablation.** We isolate the two core claims of the paper:
192 (i) latent adaptation must preserve the pretrained RGB dis-
193 tribution, and (ii) CAA requires locality-preserved latents to
194 be effective.

195 We first compare encoder-side adaptation strategies in
196 Table 2. All variants share the same diffusion archi-
197 tecture with both CAA and locality regularization en-
198 abled; as Frozen VAE does not require encoder fine-tuning,
199 we allocate additional diffusion training steps to match
200 overall compute. Frozen VAE [9] degrades shaded FID
201 ($3.083 \rightarrow 3.419$) and albedo FID ($4.599 \rightarrow 4.926$), con-
202 firming that encoding roughness and metallic channels with
203 zero-padding introduces a distribution shift that hinders
204 downstream generation. LayerDiffuse-style residual
205 prediction with identity loss [33] and Orchid-style direct
206 prediction with KL alignment [14] each recover part of the im-
207 provement but neither matches MATVAE, which combines
208 residual initialization with distributional regularization in a
209 single stochastic posterior.

210 We next isolate the multi-view consistency components.
211 Removing locality regularization ($\mathcal{L}_{\text{local}}$) while keeping
212 CAA hurts both appearance and consistency (albedo FID
213 $4.599 \rightarrow 5.873$, c-PSNR $21.934 \rightarrow 19.437$): without spa-
214 tially local latent-pixel mappings, CAA propagates features
215 across unrelated regions. Conversely, removing CAA while
216 keeping $\mathcal{L}_{\text{local}}$ preserves appearance metrics but still lowers
217 c-PSNR to 18.687, as the model must infer correspondences
218 implicitly via dense attention. MATLAT, which combines
219 both, achieves the highest c-PSNR without degrading other
220 metrics, confirming that geometry-aware attention and lo-
221 cality regularization are complementary.

222 4. Conclusion

223 We presented MATLAT, a two-stage pipeline that learns a
224 material latent space for text-guided PBR texture genera-
225 tion by leveraging pretrained RGB diffusion priors. Our ap-
226 plication is available at <https://github.com/valentinlombardi/matl原因>.
227

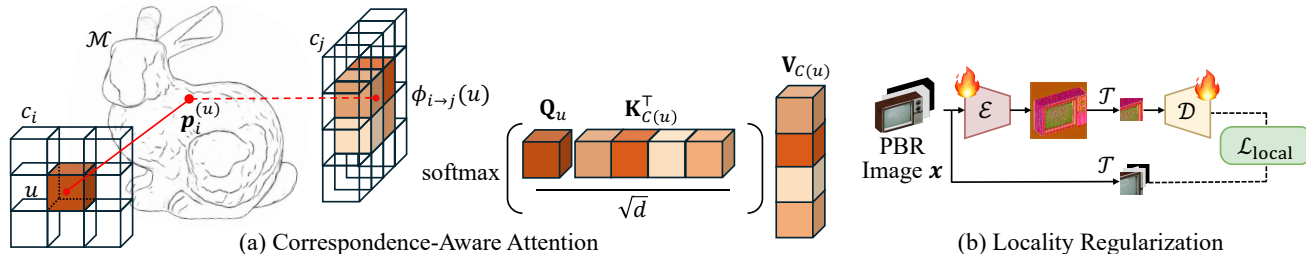


Figure 2. **Locality-preserved multi-view generation.** (a) Correspondence-aware attention (CAA) restricts feature exchange to geometrically corresponding tokens across views. (b) Our locality regularizer enforces patch-wise reconstruction so that decoded pixels depend on spatially aligned latent tokens. Together, they turn the adapted material latent space into a reliable representation for multi-view PBR texture synthesis.

Table 1. **Main quantitative comparison on held-out Objaverse-XL meshes.** Best scores are **bold**, and the runner-up scores are underlined. We highlight shaded/albedo FID (overall fidelity), roughness/metallic RMSE (material faithfulness), c-PSNR in Tab. 2 (cross-view consistency), and runtime as the primary metrics. MATLAT improves the strongest multi-view baseline MaterialMVP by over 50% in both shaded and albedo FID while keeping similar runtime and competitive material-map accuracy.

Method	Type	Shaded			Albedo			Rough.	Metal.	Time ↓
		FID _{CLIP} ↓	KID ↓	CLIP ↑	FID _{CLIP} ↓	KID ↓	CLIP ↑	RMSE ↓	RMSE ↓	
MeshGen [3]	SC	8.637	11.105	0.282	11.322	16.193	0.282	0.143	0.201	195s
TexGaussian [29]	SC	6.025	3.571	0.301	12.119	9.381	0.299	<u>0.145</u>	0.243	73s
Paint-it [32]	SDS	8.547	5.382	0.309	13.063	10.665	0.299	0.168	0.200	1260s
DreamMat [35]	SDS	<u>5.422</u>	<u>2.668</u>	0.311	<u>9.621</u>	<u>6.002</u>	0.311	0.167	0.165	2400s
FlashTex [6]	SDS	7.119	5.354	0.305	12.320	10.441	0.298	0.143	0.186	285s
MaterialAnything [10]	MV	6.582	5.287	<u>0.312</u>	12.691	9.325	0.317	0.233	0.200	500s
MaterialMVP [9]	MV	6.309	5.744	0.294	9.630	8.811	0.290	0.175	0.133	<u>35s</u>
MATLAT (Ours)	MV	3.083	1.327	0.318	4.599	1.574	<u>0.314</u>	0.158	<u>0.134</u>	34s



Figure 3. **Representative qualitative results of MATLAT.** We show generated shaded renderings for diverse prompts and geometries. MATLAT produces coherent PBR textures across glossy, metallic, wooden, and mixed-material assets while preserving fine local details and consistent relighting.

232
233
234
235
236
237

proach addresses the domain gap introduced by roughness and metallic channels through MATVAE, which incorporates PBR material information via residual posterior prediction while regularizing the learned distribution toward the pretrained latent space. Combined with locality regularization that maintains spatial alignment between latent

Table 2. **Detailed ablation summary.** The first block compares latent adaptation strategies, and the second isolates locality regularization and CAA. We report shaded FID, albedo FID, and c-PSNR; lower is better for FID, and higher is better for c-PSNR.

Variant	Shaded FID	Albedo FID	c-PSNR
Frozen VAE [9]	3.419	4.926	19.869
Res. Pred. + \mathcal{L}_{id}	3.210	4.871	20.977
Direct Pred. + \mathcal{L}_{reg}	3.192	4.768	21.468
w/o \mathcal{L}_{local}	3.419	5.873	19.437
w/o CAA	3.110	4.732	18.687
MATLAT	3.083	4.599	21.934

tokens and decoded pixels, correspondence-aware attention reliably enforces multi-view consistency in pixel space. Experiments on Objaverse-XL show that MATLAT achieves superior fidelity and multi-view consistency over strong baselines at practical runtime, with ablations confirming that both distributional alignment and locality preservation are essential.

238
239
240
241
242
243
244

245

References

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

- [1] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. In *ICCV*, 2023. 1
- [2] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *ICCV*, 2023. 1
- [3] Zilong Chen, Yikai Wang, Wenqiang Sun, Feng Wang, Yiwen Chen, and Huaping Liu. Meshgen: Generating pbr textured mesh with render-enhanced auto-encoder and generative data augmentation. In *CVPR*, 2025. 1, 2, 3, 4
- [4] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. In *NeurIPS*, 2023. 2, 3
- [5] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *CVPR*, pages 13142–13153, 2023. 2
- [6] Kangle Deng, Timothy Omernick, Alexander Weiss, Deva Ramanan, Jun-Yan Zhu, Tinghui Zhou, and Maneesh Agrawala. Flashtex: Fast relightable mesh texturing with lightcontrolnet. In *ECCV*, 2024. 3, 4
- [7] Fan Fei, Jiajun Tang, Fei-Peng Tian, Boxin Shi, and Ping Tan. Pactice: Efficient pbr texture generation on packed views with visual autoregressive models. *arXiv preprint arXiv:2505.22394*, 2025. 2, 3
- [8] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. In *SIGGRAPH*, 2023. 1
- [9] Zebin He, Mingxin Yang, Shuhui Yang, Yixuan Tang, Tao Wang, Kaihao Zhang, Guanying Chen, Yuhong Liu, Jie Jiang, Chunchao Guo, and Wenhan Luo. Materialmvp: Illumination-invariant material generation via multi-view pbr diffusion. In *ICCV*, 2025. 2, 3, 4
- [10] Xin Huang, Tengfei Wang, Ziwei Liu, and Qing Wang. Material anything: Generating materials for any 3d object via diffusion. In *CVPR*, 2024. 2, 3, 4
- [11] Lukas Höllein, Ang Cao, Andrew Owens, Justin Johnson, and Matthias Nießner. Text2room: Extracting textured 3d meshes from 2d text-to-image models. In *ICCV*, 2023. 1
- [12] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023. 1
- [13] Jaihoon Kim, Juil Koo, Kyeongmin Yeo, and Minhyuk Sung. Synctweedies: A general generative framework based on synchronized diffusions. In *NeurIPS*, 2024. 1
- [14] Akshay Krishnan, Xinchun Yan, Vincent Casser, and Abhijit Kundu. Orchid: Image latent diffusion for joint appearance and geometry generation. In *ICCV*, 2025. 2, 3
- [15] Haoran Li, Haolin Shi, Wenli Zhang, Wenjun Wu, Yong Liao, Lin Wang, Lik-Hang Lee, and Peng Yuan Zhou. Dreamscene: 3d gaussian-based text-to-3d scene generation via formation pattern sampling. In *ECCV*, 2024. 1
- [16] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *CVPR*, 2023. 2
- [17] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In *ICLR*, 2023. 3
- [18] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. In *NeurIPS*, 2023. 3
- [19] Mathis Petrovich, Michael J. Black, and Gül Varol. Tmr: Text-to-motion retrieval using contrastive 3d human motion synthesis. In *ICCV*, 2023. 1
- [20] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 1, 2
- [21] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *SIGGRAPH*, 2023. 1
- [22] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022. 2
- [23] Zeyi Sun, Tong Wu, Pan Zhang, Yuhang Zang, Xiaoyi Dong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Bootstrap3d: Improving multi-view diffusion model with synthetic data. In *ICCV*, 2025. 3
- [24] Shitao Tang, Fuyang Zhang, Jiacheng Chen, Peng Wang, and Yasutaka Furukawa. Mvdifffusion: Enabling holistic multi-view image generation with correspondence-aware diffusion. In *NeurIPS*, 2023. 2
- [25] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H. Bermano. Human motion diffusion model. In *ICLR*, 2023. 1
- [26] Yitong Wang, Xudong Xu, Li Ma, Haoran Wang, and Bo Dai. Boosting 3d object generation through pbr materials. In *SIGGRAPH*, 2024. 1
- [27] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *NeurIPS*, 2023. 2
- [28] Xiaokang Wei, Bowen Zhang, Xianghui Yang, Yuxuan Wang, Chunchao Guo, Xi Zhao, and Yan Luximon. Pbr3dgen: A vlm-guided mesh generation with high-quality pbr texture. *arXiv preprint arXiv:2503.11368*, 2025. 2
- [29] Bojun Xiong, Jialun Liu, Jiakui Hu, Chenming Wu, Jinbo Wu, Xing Liu, Chen Zhao, Errui Ding, and Zhouhui Lian. Texgaussian: Generating high-quality pbr material via octree-based 3d gaussian splatting. In *CVPR*, 2025. 3, 4
- [30] Kyeongmin Yeo, Jaihoon Kim, and Minhyuk Sung. Stochsync: Stochastic diffusion synchronization for image generation in arbitrary spaces. In *ICLR*, 2025. 1
- [31] Taoran Yi, Jiemin Fang, Junjie Wang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang

- 359 Wang. Gaussiandreamer: Fast generation from text to 3d
360 gaussians by bridging 2d and 3d diffusion models. In *CVPR*,
361 2024. 1
- 362 [32] Kim Youwang, Tae-Hyun Oh, and Gerard Pons-Moll. Paint-
363 it: Text-to-texture synthesis via deep convolutional texture
364 map optimization and physically-based rendering. In *CVPR*,
365 2024. 3, 4
- 366 [33] Lvmin Zhang and Maneesh Agrawala. Transparent image
367 layer diffusion using latent transparency. In *SIGGRAPH*,
368 2024. 2, 3
- 369 [34] Qihang Zhang, Chaoyang Wang, Aliaksandr Siarohin, Peiye
370 Zhuang, Yinghao Xu, Ceyuan Yang, Dahua Lin, Bolei Zhou,
371 Sergey Tulyakov, and Hsin-Ying Lee. Scenewiz3d: Towards
372 text-guided 3d scene composition. In *CVPR*, 2024. 1
- 373 [35] Yuqing Zhang, Yuan Liu, Zhiyu Xie, Lei Yang, Zhongyuan
374 Liu, Mengzhou Yang, Runze Zhang, Qilong Kou, Cheng Lin,
375 Wenping Wang, and Xiaogang Jin. Dreammat: High-quality
376 pbr material generation with geometry- and light-aware dif-
377 fusion models. In *SIGGRAPH*, 2024. 3, 4
- 378 [36] Lingting Zhu, Jingrui Ye, Runze Zhang, Zeyu Hu, Yingda
379 Yin, Lanjiang Li, Jinnan Chen, Shengju Qian, Xin Wang,
380 Qingmin Liao, and Lequan Yu. Muma: 3d pbr texturing
381 via multi-channel multi-view generation and agentic post-
382 processing. *arXiv preprint arXiv:2503.18461*, 2025. 2
- 383 [37] Shenhao Zhu, Lingteng Qiu, Xiaodong Gu, Zhengyi Zhao,
384 Chao Xu, Yuxiao He, Zhe Li, Xiaoguang Han, Yao Yao,
385 Xun Cao, Siyu Zhu, Weihao Yuan, Zilong Dong, and Hao
386 Zhu. Mcmat: Multiview-consistent and physically accurate
387 pbr material generation. *arXiv preprint arXiv:2412.14148*,
388 2024. 2