

# 3M-TI: High-Quality Mobile Thermal Imaging via Calibration-free Multi-Camera Cross-Modal Diffusion

Anonymous CVPR submission

Paper ID 16226

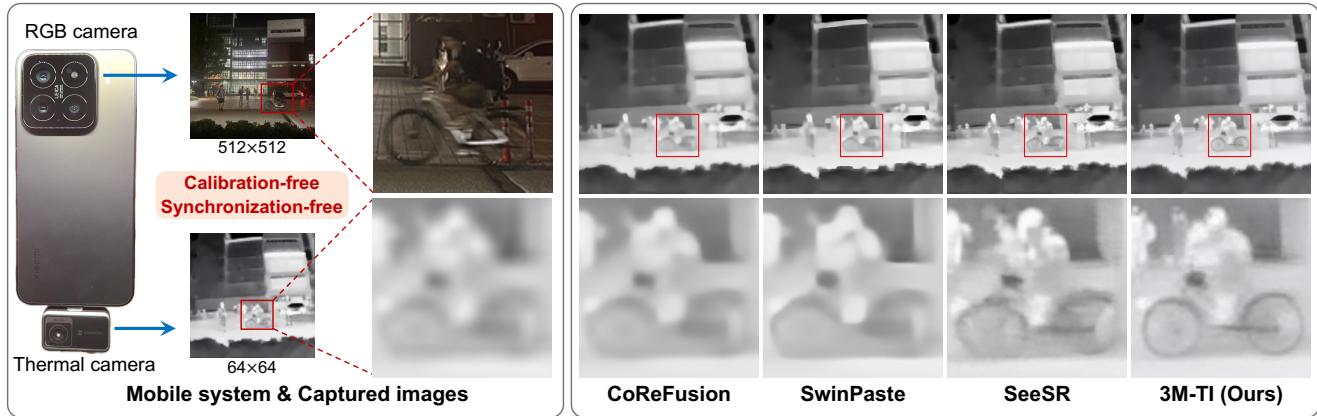


Figure 1. A smartphone-based mobile imaging system integrating calibration-free and synchronization-free RGB and thermal cameras. The proposed 3M-TI method delivers superior thermal image quality compared with state-of-the-art restoration approaches.

## Abstract

001 *The miniaturization of thermal sensors for mobile platforms*  
 002 *inherently limits their spatial resolution and textural fi-*  
 003 *delicacy, leading to blurry and less informative images. Existing*  
 004 *thermal super-resolution (SR) methods can be grouped*  
 005 *into single-image and RGB-guided approaches: the former*  
 006 *struggles to recover fine structures from limited informa-*  
 007 *tion, while the latter relies on accurate and laborious cross-*  
 008 *camera calibration, which hinders practical deployment*  
 009 *and robustness. Here, we propose 3M-TI, a calibration-*  
 010 *free Multi-camera cross-Modality diffusion framework for*  
 011 *Mobile Thermal Imaging. At its core, 3M-TI integrates a*  
 012 *cross-modal self-attention module (CSM) into the diffusion*  
 013 *UNet, replacing the original self-attention layers to adap-*  
 014 *tively align thermal and RGB features throughout the de-*  
 015 *noising process, without requiring explicit camera calibra-*  
 016 *tion. This design enables the diffusion network to lever-*  
 017 *age its generative prior to enhance spatial resolution, struc-*  
 018 *tural fidelity, and texture detail in the super-resolved ther-*  
 019 *mal images. Extensive evaluations on real-world mobile*  
 020 *thermal cameras and public benchmarks validate our supe-*  
 021 *rior performance, achieving state-of-the-art results in both*  
 022 *visual quality and quantitative metrics. More importantly,*

*the thermal images enhanced by 3M-TI lead to substantial gains in critical downstream tasks like object detection and segmentation, underscoring its practical value for robust mobile thermal perception systems.*

023  
024  
025  
026

## 1. Introduction

027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037

Thermal imaging detects infrared radiation emitted by objects, extending perception beyond the visible spectrum for enhanced scene understanding. It is particularly effective under challenging conditions such as darkness, fog, or smoke, providing reliable scene information that complements RGB imaging for multimodal perception and fusion [30, 31, 43]. Owing to its robustness, thermal imaging has become increasingly valuable in safety-critical applications, such as autonomous driving, robotic navigation, and situational awareness [40, 41, 44].

Despite rapid advances in infrared sensing, mobile thermal imaging remains fundamentally constrained by hardware. Miniaturization for compact devices reduces aperture size, lowering spatial resolution and textural detail. In addition, the long wavelength of thermal radiation limits the minimum pixel size, while the high cost of infrared sensors further restricts pixel counts compared with RGB cameras.

038  
039  
040  
041  
042  
043  
044

045 These factors jointly result in blurry and less informative  
046 thermal images on mobile platforms, motivating extensive  
047 efforts to enhance their resolution and perceptual quality  
048 through computational reconstruction.

049 A representative example of such computational recon-  
050 struction is thermal image super-resolution (SR), which  
051 aims to reconstruct high-resolution spatial details from low-  
052 resolution thermal observations [49]. However, the inher-  
053 ently limited information in low-resolution observations,  
054 aggravated by the lack of large-scale thermal datasets,  
055 makes thermal super-resolution highly underconstrained  
056 and hinders network training and generalization compared  
057 with RGB-based tasks. To compensate for this deficiency,  
058 RGB-guided SR methods exploit cross-modal cues to en-  
059 hance reconstruction [1, 12, 17]. Nevertheless, due to the  
060 fundamental differences in imaging principles, thermal and  
061 RGB images often exhibit substantial discrepancies, and di-  
062 rectly merging their features may introduce unrealistic de-  
063 tails [29, 59, 62]. Moreover, most existing methods rely on  
064 pixel-level aligned RGB-thermal pairs, requiring laborious  
065 calibration, which severely limits robustness and scalability  
066 in practical deployment.

067 To address these challenges, we propose 3M-TI,  
068 a calibration-free Multi-camera cross-Modal diffusion  
069 framework for Mobile Thermal Imaging. Firstly, rather  
070 than relying on pixel-level calibration between RGB and  
071 thermal cameras, 3M-TI performs alignment in the latent  
072 space of a variational autoencoder (VAE) [19], where con-  
073 tinuous and disentangled representations facilitate robust  
074 cross-modal correspondence. An attention-based fusion  
075 module further integrates RGB and thermal representations  
076 in this latent domain, achieving reliable alignment and fu-  
077 sion without explicit registration. Secondly, because ex-  
078 isting RGB-thermal datasets are typically small, homoge-  
079 neous, and strictly aligned [16, 28], we design a data aug-  
080 mentation strategy that intentionally offsets RGB-thermal  
081 pairs to simulate real-world conditions. This encourages  
082 the attention module to learn robust cross-modal corre-  
083 spondence and fusion without relying on explicit pixel-level  
084 alignment. Furthermore, we incorporate a one-step latent-  
085 space diffusion [39] to suppress fusion artifacts, enhance  
086 visual fidelity, and improve robustness under limited-data  
087 conditions.

088 Finally, we validate the proposed 3M-TI framework on a  
089 real mobile thermal imaging system. Experiments show that  
090 it significantly outperforms existing state-of-the-art meth-  
091 ods, highlighting both its effectiveness and practical poten-  
092 tial. The main contributions of 3M-TI can be summarized  
093 as follows:

094 • **Calibration-free Fusion:** Fully automatic alignment and  
095 fusion of uncalibrated RGB-thermal pairs via a cross-  
096 modal self-attention module, without requiring pixel-  
097 level image registration.

- **Misalignment Augmentation:** Camera-pose augmentation that generates intentionally misaligned RGB counterparts, enabling the cross-modal self-attention module to handle multi-camera parallax and unsynchronized captures. 098 099 100 101 102 103 104 105 106 107 108 109
- **Cross-modal Diffusion:** Integration of latent-space diffusion to suppress fusion artifacts, enhance detail realism, and leverage pretrained priors to compensate for limited thermal data. 103 104 105 106
- **Practical Validation:** Evaluation on both public datasets and a real mobile thermal imaging system, demonstrating effectiveness and practical potential. 107 108 109

## 2. Related Work

110 Computational thermal imaging intersects several areas of  
111 computer vision and image processing. In this section,  
112 we review the most relevant work in three directions:  
113 classical and learning-based image restoration, reference-  
114 guided image super-resolution, and recent diffusion-based  
115 approaches. 116

### 2.1. Image Restoration

117 Image restoration seeks to recover a high-quality image  
118 from a degraded observation by inverting an imaging  
119 model. Classical image restoration methods are typically  
120 model-based, formulated as inverse problems with explicit  
121 priors. Variational methods regularize the solution with  
122 smoothness or edge-preserving constraints [38]. Sparse rep-  
123 resentation approaches exploit dictionary learning and spar-  
124 sity priors to recover fine details [10, 32]. Non-local self-  
125 similarity models, such as non-local means [3] and BM3D  
126 [8], leverage the redundancy of similar patches across the  
127 image for collaborative filtering. Deep learning has funda-  
128 mentally reshaped image restoration by learning powerful  
129 data-driven priors directly from large-scale datasets. CNN-  
130 based models (e.g., SRCNN[9], EDSR[26]) pioneered end-  
131 to-end restoration pipelines, while recent Transformer ar-  
132 chitectures such as SwinIR[25] further improved long-  
133 range dependency modeling. These advances have led to  
134 state-of-the-art performance in denoising [24, 25, 45], de-  
135 blurring [20, 27], and super-resolution [5, 6, 9, 21, 25]. Nev-  
136 ertheless, single-image restoration remains challenging un-  
137 der severe degradation or large upscaling factors, due to the  
138 limited high-frequency information in a single observation,  
139 motivating approaches that leverage additional cues. 140

### 2.2. Reference-Guided Image Restoration

141 Reference-guided image restoration leverages an auxiliary  
142 high-resolution modality (typically RGB) to recover fine  
143 details lost in degraded observations. This paradigm has  
144 been explored across multiple tasks, including RGB im-  
145 age super-resolution [11, 42, 58, 61], hyperspectral im-  
146 age super-resolution [12, 57], and thermal image enhance-  
147

148 ment [1, 63]. Early approaches, such as CoReFusion [17],  
149 employ a dual-encoder UNet to jointly extract and fuse  
150 RGB–thermal features. A contrastive loss is further in-  
151 troduced to improve cross-modal consistency and over-  
152 all restoration quality. With the advent of Transfor-  
153 mers [46], attention-based fusion has become the dominant  
154 paradigm. MGNet [60] mines multi-level cues (appear-  
155 ance, edge, and semantics) from RGB guidance to enhance  
156 UAV thermal super-resolution. SwinFuSR [1] introduces a  
157 lightweight Swin-Transformer backbone with robust train-  
158 ing under missing guidance, and SwinPaste [63] further  
159 improves it via data mixing and multi-scale supervision.  
160 MSFFCT [35] combines channel-based Transformers with  
161 multi-scale feature fusion to capture long-range dependen-  
162 cies and rich thermal-RGB correlations. However, the large  
163 domain shift between RGB and thermal images poses a  
164 significant challenge to establishing reliable cross-modal  
165 alignment, especially under uncalibrated conditions.

### 166 2.3. Diffusion-Powered Image Restoration

167 Diffusion models have recently demonstrated remarkable  
168 capabilities in image restoration and super-resolution by  
169 learning powerful generative priors that capture rich tex-  
170 tures and structural details [13]. This generative nature  
171 allows diffusion models to recover fine structures and  
172 high-frequency details that are challenging for CNN- and  
173 Transformer-based approaches. However, the original dif-  
174 fusion model is computationally expensive due to its long  
175 iterative denoising process, limiting practical deployment.  
176 To address this, several efficient variants have been pro-  
177 posed: latent-space diffusion methods [47, 51, 52, 54] per-  
178 form denoising within the compact latent representation of  
179 a VAE, reducing computation overhead while preserving  
180 high-fidelity reconstructions; distillation-based approaches  
181 [15, 33, 39] transfer the generative ability of a large teacher  
182 model to a lightweight student, enabling faster inference  
183 with minimal quality loss.

184 Diffusion models have also been applied to thermal  
185 image super-resolution [7, 23]. DifIISR [23] incorpo-  
186 rates gradient-based priors and a thermal spectrum dis-  
187 tribution regulation into the diffusion process, guiding the  
188 model to capture the unique frequency characteristics of  
189 infrared images and achieving strong visual quality and  
190 downstream task performance. [7] adapts the ResShift dif-  
191 fusion model[53] to thermal imaging with an uncertainty-  
192 aware approach. This method produces a confidence map  
193 that evaluates the reliability of each reconstructed region,  
194 thereby enhancing pixel-level fidelity and facilitating tex-  
195 ture reconstruction.

196 However, existing methods still rely on strictly pixel-  
197 aligned datasets and explicit calibration between cameras,  
198 making them sensitive to spatial and temporal misalign-  
199 ment and difficult to generalize to real-world mobile sce-

200 narios. These limitations motivate our design of 3M-TI, a  
201 calibration-free latent-space diffusion framework that lever-  
202 ages multi-modal information for robust mobile thermal  
203 imaging.

## 204 3. Method

We propose 3M-TI, a cross-modal diffusion frame, for high-  
205 quality thermal image reconstruction from uncalibrated  
206 RGB-thermal image pairs (Fig. 2). 3M-TI specifically ad-  
207 dresses three key challenges: (1) uncalibrated and unsyn-  
208 chronized RGB-thermal captures, (2) fusion of heteroge-  
209 neous RGB-thermal representations, and (3) limited data  
210 scale and diversity of RGB-thermal datasets.

211 Section 3.1 provides an overview of the 3M-TI frame-  
212 work. Section 3.2 introduces the cross-modal self-attention  
213 (CSA) module, which performs automatic alignment and  
214 fusion within the VAE latent space. Section 3.3 presents  
215 a misalignment augmentation strategy that simulates realis-  
216 tic inter-camera parallax and temporal offset to improve the  
217 robustness of CSA. Finally, Section 3.4 details the practical  
218 implementation and validation of 3M-TI on a smartphone-  
219 based multi-camera thermal imaging system.

### 221 3.1. 3M-TI Framework Overview

The 3M-TI framework reconstructs a high-resolution ther-  
222 mal image from a low-resolution thermal input and an un-  
223 calibrated high-resolution RGB reference, as illustrated in  
224 Fig. 2(a). Our approach is built upon the one-step diffu-  
225 sion model, SD-Turbo [39], to ensure efficient inference.  
226 The process begins by encoding the thermal and RGB im-  
227 ages into latent representations using the frozen VAE en-  
228 coder. To establish cross-modal correspondence and fu-  
229 sion in the latent space, we introduce the cross-modal self-  
230 attention module (CSM). This module replaces the original  
231 self-attention layers in the diffusion UNet with our cross-  
232 modal self-attention layers, which are designed to learn  
233 multiscale correspondences between the RGB and thermal  
234 latents (Fig. 2(b)). During training, we apply misalign-  
235 ment augmentation to the RGB images to enhance robust-  
236 ness against uncalibration and temporal unsynchronization  
237 (Fig. 2(c)). Furthermore, a skip connection [34] is incor-  
238 porated to enhance structural consistency and mitigate geo-  
239 metric distortions. Specifically, the feature maps from 4 en-  
240 coder downsampling blocks are passed through a  $1 \times 1$  zero-  
241 initialized convolutional layer before being added to the  
242 corresponding decoder upsampling blocks. Since extract-  
243 ing reliable semantics directly from low-resolution thermal  
244 images is suboptimal, we generate text prompts from the  
245 corresponding RGB images using the Recognize Anything  
246 Model (RAM) [56]. Finally, we employ low-rank adapta-  
247 tion (LoRA) [14] to fine-tune both the UNet and the VAE  
248 decoder.

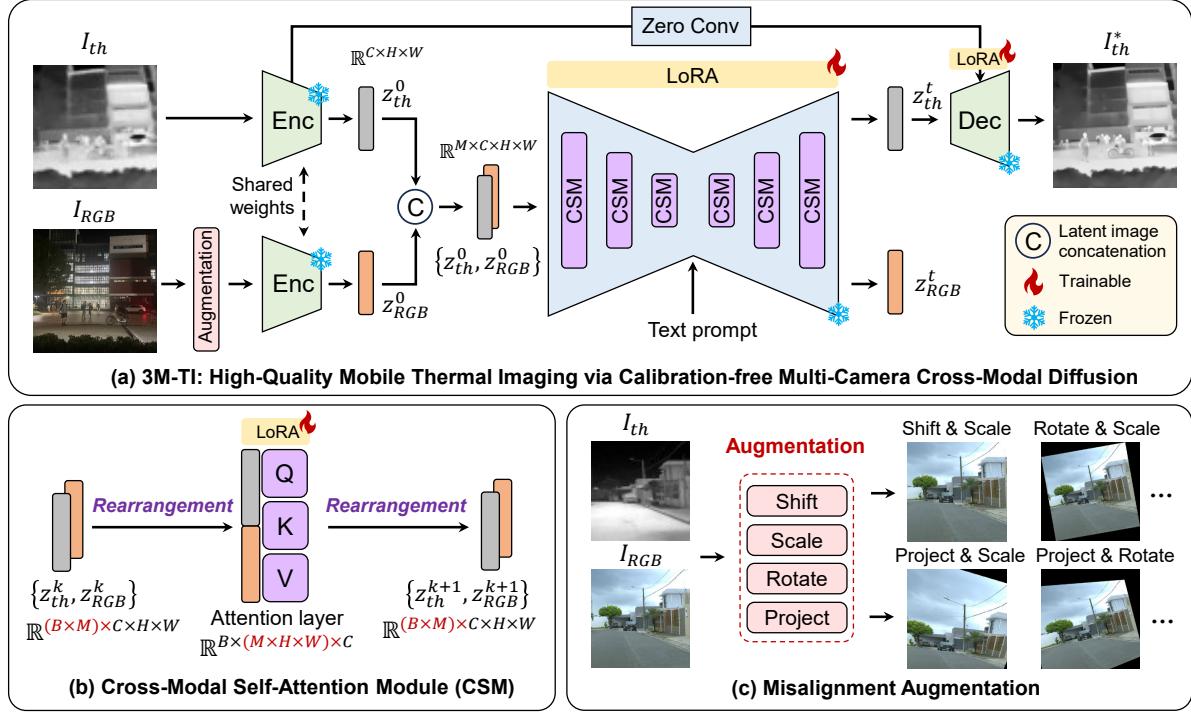


Figure 2. Overview of the 3M-TI architecture. (a) 3M-TI framework. The core of 3M-TI is a one-step diffusion-based model equipped with a cross-modal self-attention module (CSM) and a misalignment augmentation strategy. LoRA fine-tuning is applied to both the UNet and the VAE decoder. (b) Cross-modal self-attention module (CSM). Two rearrangement layers are inserted before and after the original self-attention layers to capture cross-modal correspondences. (c) Misalignment augmentation. A data augmentation strategy designed to enhance model robustness against camera parallax and temporal misalignment between RGB and thermal inputs.

### 3.2. Cross-Modal Self-Attention Module

To establish cross-modal correspondence between uncalibrated RGB and thermal images, we introduce the cross-modal self-attention module (CSM). Inspired by video- and multi-view diffusion models [2, 50], our key idea is to repurpose the transformer blocks within the pretrained diffusion UNet for cross-modal RGB-thermal correspondence and fusion. The RGB and thermal inputs are first encoded into a concatenated latent tensor  $\{z_{RGB}^0, z_{th}^0\} \in \mathbb{R}^{B \times M \times C \times H \times W}$ , where  $B$  is the batch size,  $M = 2$  denotes the number of images (RGB and thermal, can be extended to multiple RGB images), and  $C, H, W$  are the channel and spatial dimensions.

For non-transformer blocks in the diffusion UNet, the image-number dimension  $M$  is folded into the batch dimension so that RGB and thermal images are processed independently. Within the transformer blocks (Fig. 2b), two tensor rearrangements are applied. Before entering a block, the latent tensor is reshaped to  $\mathbb{R}^{B \times (M \times H \times W) \times C}$ , treating each pixel as a token and merging all pixels from both modalities into a single sequence with a length of  $M \times H \times W$ . Self-attention then computes pairwise dependencies among all RGB and thermal pixels, naturally enabling the transformer

to learn cross-modal correspondences and integrate complementary cues. After the block, the latent tensor is reshaped back to  $\mathbb{R}^{(B \times M) \times C \times H \times W}$  for subsequent layers.

At the output stage, the refined thermal latent is merged with the initial latent in the VAE encoder via a zero-initialized skip connection, producing the final high-fidelity thermal image. The cross-modal self-attention module introduces no additional parameters, allowing the UNet and VAE decoder to be efficiently fine-tuned using LoRA [14].

### 3.3. Misalignment Augmentation

Existing RGB-thermal datasets are typically small and strictly pixel-aligned, limiting generalization to practical multi-camera configurations. In practice, spatial misalignment arise from camera parallax and temporal unsynchronization: parallax varies with scene depth and baseline, while temporal offsets are determined by object motion and capture delay.

Since our goal is to learn a model that is robust to these geometric variations, rather than overfitting to a specific calibrated setup, we propose a misalignment augmentation strategy without physical simulation. 3M-TI applies a set of controlled spatial transformations to the RGB images, including translation, scaling, rotation, and perspec-

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296 tive warping. The transformation parameters are chosen  
297 to reflect typical deviations encountered in handheld and  
298 multi-camera scenarios (see Figure 2(c)).

299 This misalignment augmentation encourages the CSM  
300 to learn robust cross-modal correspondences under uncali-  
301 brated and unsynchronized conditions, effectively bridging  
302 the gap between constrained training data and diverse de-  
303 ployment environments.

### 304 3.4. Practical Multi-Camera System

305 To validate our method under realistic conditions, we con-  
306 structed a multi-camera system comprising a HIKVISION  
307 P09 thermal camera module (costing less than 100 US Dol-  
308 lars) and a Xiaomi 15 smartphone, connected via a Type-C  
309 interface (Fig. 1). The thermal module provides a native res-  
310 olution of  $96 \times 96$ , which is resized to  $64 \times 64$  in this work,  
311 with a  $50^\circ \times 50^\circ$  field of view (FOV) and a pixel pitch of  
312  $12 \mu\text{m}$ . For the RGB reference, we utilize the primary cam-  
313 era of the Xiaomi 15, which is equipped with an OV50H  
314 sensor and captures images at a resolution of  $4096 \times 3072$ .  
315 This RGB camera has an equivalent focal length of 23 mm,  
316 yielding a FOV of  $74^\circ \times 59^\circ$ , slightly wider than that of the  
317 thermal module. This system offers a realistic platform to  
318 evaluate the performance of 3M-TI under practical, uncali-  
319 brated mobile imaging conditions.

## 320 4. Experiments

### 321 4.1. Experimental Setup

322 **Datasets Preparation.** Our training and test sets are com-  
323 piled from four public RGB-thermal datasets: IRVI [22],  
324 LLVIP [16], M<sup>3</sup>FD [28], and the PBVS 2025 TISR Chal-  
325 lenge Track 2 [37]. Among them, IRVI, LLVIP, and PBVS  
326 2025 provide official training-validation/test splits, while  
327 M<sup>3</sup>FD does not. To ensure a balanced and diverse dataset,  
328 we sample images from each source as follows: 3,200 pairs  
329 from IRVI, 3,200 from LLVIP, 3,822 from M<sup>3</sup>FD, and 700  
330 from the PBVS 2025 training split, forming a combined  
331 training set of 10,922 image pairs. For evaluation, we con-  
332 struct a test set of 1,176 pairs by sampling 300 pairs from  
333 the IRVI test set, 300 from LLVIP, 376 from M<sup>3</sup>FD, and  
334 200 from the PBVS 2025 validation set. Since the original  
335 datasets contain strictly aligned image pairs, we further  
336 augment the test set by applying controlled spatial trans-  
337 formations to simulate camera parallax and temporal offset.  
338 All thermal images are center-cropped to a square aspect ra-  
339 tio and resized to  $64 \times 64$  pixels. To better mimic real sensor  
340 noise, we add Gaussian noise. Corresponding RGB images  
341 are cropped identically and resized to  $512 \times 512$  pixels to  
342 serve as high-resolution references.

343 **Smartphone Dataset.** For real-world evaluation, we col-  
344 lected 40 image pairs across 30 diverse scenes using  
345 our smartphone-based multi-camera system, including 37

346 nighttime and the remaining daytime captures. Similarly,  
347 all thermal images are resized to  $64 \times 64$ . The correspond-  
348 ing RGB images are center-cropped to a square aspect ratio  
349 and resized to  $512 \times 512$  as high-resolution references.

350 **Implementation Details.** The model is trained with the  
351 loss in Eq. (1), combining L2 and learned perceptual im-  
352 age patch similarity (LPIPS) [55] terms with  $\lambda = 1$ :

$$\mathcal{L} = \mathcal{L}_2 + \lambda \cdot \mathcal{L}_{\text{LPIPS}}. \quad (1)$$

353 Training is performed using the Adam optimizer [18] with  
354 a learning rate of  $2 \times 10^{-5}$  and a batch size of 4 on a single  
355 NVIDIA A800 (80 GB) GPU for about 4 hours (8000 iter-  
356 ations). LoRA is applied with ranks of 16 for the UNet and  
357 4 for the VAE decoder.

358 **Comparative Methods.** We compare our model against 5  
359 representative baselines: CoReFusion [17], SwinFuSR [1],  
360 SwinPaste [63], SeeSR [52], and OSEDiff [51]. CoReFu-  
361 sion is an RGB-guided super-resolution model built upon  
362 the UNet architecture. SwinFuSR employs a Swin Trans-  
363 former backbone with RGB guidance, while SwinPaste is  
364 an enhanced variant. SeeSR and OSEDiff are diffusion-  
365 based image super-resolution methods without reference  
366 guidance. For fair comparison, all baseline models are re-  
367 trained on our dataset using their publicly released codes.

368 **Evaluation Metrics.** We evaluate model performance us-  
369 ing reference metrics, including PSNR and SSIM [48] for  
370 reconstruction fidelity and LPIPS [55] for perceptual qual-  
371 ity. We also report no-reference metrics, MUSIQ [4] and  
372 MANIQA [4], to quantify overall image quality.

### 373 4.2. Image Quality Assessment

375 Table 1 summarizes the quantitative results. Our 3M-TI  
376 achieves the best performance on perceptual metrics, in-  
377 cluding LPIPS, MANIQA, and MUSIQ. For the reference-  
378 based LPIPS metric, 3M-TI, OSEDiff, and SeeSR out-  
379 perform non-diffusion methods (CoReFusion, SwinFuSR,  
380 and SwinPaste), highlighting the power of diffusion pri-  
381 ors for perceptual enhancement. In terms of fidelity met-  
382 rics (PSNR, SSIM), 3M-TI ranks a close second, demon-  
383 strating better structural preservation than OSEDiff, SeeSR,  
384 SwinFuSR, and SwinPaste. The proposed misalignment  
385 augmentation significantly improves 3M-TI, particularly on  
386 perceptual metrics, while other methods see only marginal  
387 gains. This indicates that 3M-TI effectively learns cross-  
388 modal correspondences and generalizes well to diverse un-  
389 calibrated conditions.

390 As shown in Fig. 3, qualitative results reveal a key trend:  
391 non-diffusion methods tend to achieve high PSNR/SSIM  
392 yet produce overly smooth reconstructions that miss fine  
393 high-frequency structures. OSEDiff and SeeSR can synthe-  
394 size high-frequency content but often introduces details that  
395 are inconsistent with the ground truth. In contrast, 3M-TI  
396 reconstructs fine details that are both visually plausible and

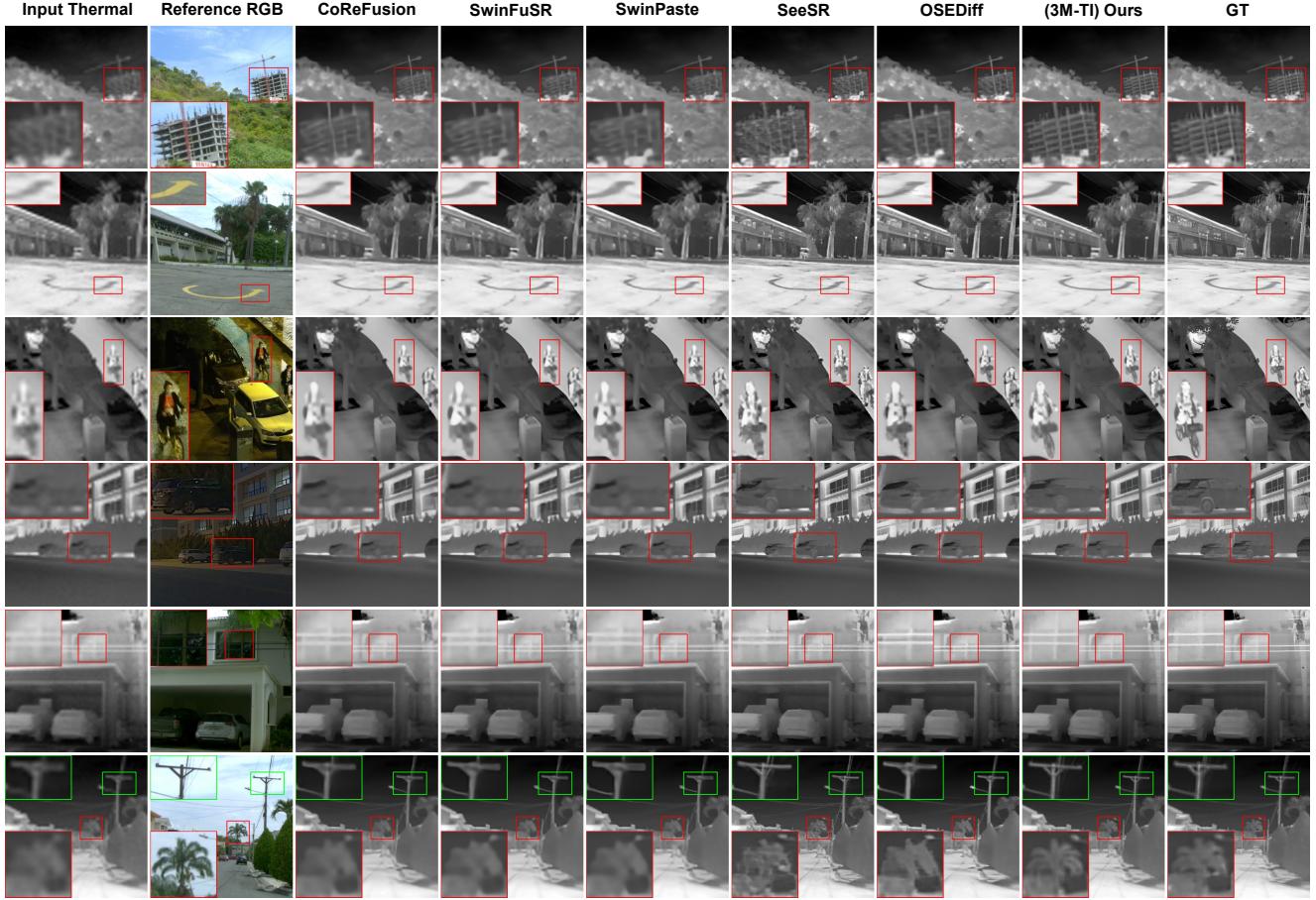


Figure 3. Qualitative comparison on our test set (zoom in for details). 3M-TI achieves the most faithful and visually consistent results, exhibiting sharp structures and accurate thermal patterns that best align with the GT.

397 closely aligned with the ground truth. For example, 3M-TI  
 398 faithfully transfers geometric details such as partially con-  
 399 structed buildings (Row 1), road markings (Row 2), win-  
 400 dows and wires (Row 5), and tree branches and utility poles  
 401 (Row 6) from the RGB references to the thermal outputs,  
 402 whereas other methods either miss these structures or re-  
 403 produce them inaccurately. 3M-TI also yields more natural  
 404 object contours (e.g., cyclists and cars in the Row 3 and  
 405 Row 4), resulting in visually realistic thermal images.

406 Figure 4 presents qualitative results on our smartphone  
 407 dataset, confirming 3M-TI’s robustness to real-world, un-  
 408 calibrated pairs. While CoReFusion, SwinFuSR, and Swin-  
 409 Paste yield blurry outputs; OSEDiff and SeeSR introduce  
 410 unrealistic artifacts, our method produces geometrically  
 411 consistent and visually plausible details. A key example  
 412 is the bicycle wheels (Row 3), where 3M-TI accurately re-  
 413 constructs their circular shape despite significant positional  
 414 misalignment between the RGB and thermal captures due  
 415 to unsynchronization. Moreover, 3M-TI can leverage weak  
 416 cues from degraded references, as shown by its ability to

417 recover the railing structure (Row 5) from an RGB im-  
 418 age heavily corrupted by light flare. Table 2 reports no-  
 419 reference evaluation on the captured smartphone dataset,  
 420 where 3M-TI achieves the highest scores. More results and  
 421 technical details are given in the supplementary material.

### 4.3. Downstream Applications

422 The high-fidelity outputs of 3M-TI provide structurally and  
 423 semantically rich representations that directly benefit down-  
 424 stream vision tasks. To validate this claim, we evaluate two  
 425 representative tasks: open-vocabulary object detection and  
 426 semantic segmentation. In these experiments, we apply the  
 427 pretrained Grounded-SAM model [36] in a zero-shot man-  
 428 nner (no fine-tuning) and use identical text prompts for all  
 429 methods to ensure fair comparisons.

**Object detection.** For the examples in Fig. 5, the text  
 430 prompts are “automobile” and “person”. 3M-TI demon-  
 431 strates superior detection robustness and object discovery  
 432 capability compared to SwinPaste and SeeSR. Table 3 sum-  
 433 marizes detection performance (Precision, Recall, F1-score,  
 434 and IoU) on the test set (LLVIP and M<sup>3</sup>FD with GT an-



Figure 4. Qualitative comparison on our real-world smartphone dataset (zoom in for details). 3M-TI exhibits remarkable generalization capability, producing sharp and faithful thermal details that closely align with RGB images.

| Metrics               | PSNR↑        | SSIM↑         | LPIPS↓        | MANIQA↑       | MUSIQ↑       |
|-----------------------|--------------|---------------|---------------|---------------|--------------|
| CoReFusion [17]       | 30.11        | 0.8588        | 0.3214        | 0.2771        | 28.35        |
| CoReFusion w/ Augment | <b>30.30</b> | <b>0.8634</b> | 0.3174        | 0.2748        | 28.95        |
| SwinFuSR [1]          | 29.85        | 0.8549        | 0.3085        | 0.2740        | 29.86        |
| SwinFuSR w/ Augment   | 29.98        | 0.8581        | 0.3134        | 0.2753        | 30.33        |
| SwinPaste [63]        | 29.83        | 0.8545        | 0.3075        | 0.2719        | 29.63        |
| SwinPaste w/ Augment  | 29.91        | 0.8575        | 0.3084        | 0.2745        | 30.66        |
| SeeSR [52]            | 29.41        | 0.8495        | 0.1828        | 0.4278        | 35.22        |
| OSEDiff [51]          | 28.05        | 0.8422        | 0.2113        | 0.4014        | <b>36.30</b> |
| Ours w/o Augment      | 30.04        | 0.8597        | 0.1917        | 0.4265        | 34.94        |
| Ours w/ Augment       | 30.09        | 0.8610        | <b>0.1787</b> | <b>0.4443</b> | <b>36.66</b> |

Table 1. Performance comparison of different methods on public datasets. Gray cells indicate the best result, and light gray cells indicate the second-best result for each metric.

notations). 3M-TI achieves the best overall detection performance, marginally surpassing the reference RGB results and closely approaching the scores obtained on GT thermal images.

| Metrics | CoReFusion | SwinFuSR | SwinPaste | SeeSR  | OSEDiff | 3M-TI         |
|---------|------------|----------|-----------|--------|---------|---------------|
| MUSIQ↑  | 25.74      | 25.99    | 26.17     | 30.15  | 29.85   | <b>30.62</b>  |
| MANIQA↑ | 0.2701     | 0.2754   | 0.2749    | 0.3454 | 0.3285  | <b>0.3589</b> |

Table 2. Performance comparison on real-world smartphone dataset. Gray cells indicate the best.

**Semantic segmentation.** For the examples in Fig. 6, the prompts are “sky, tree, automobile, window” and “wheel, automobile, rider”. 3M-TI produces more accurate and coherent segmentation maps than other methods, while complementing the RGB reference. In the first example, compared with the RGB results, 3M-TI segments trees more precisely though it misses a small automobile in the bottom-left. In the second low-light example, 3M-TI even outperforms the RGB reference (e.g., wheels and the upper automobile).

Overall, these results confirm that 3M-TI generates se-

441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451

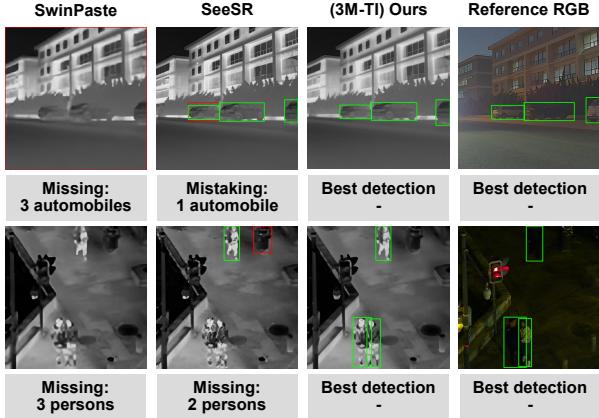


Figure 5. Visualization of detection results, where green bounding boxes indicate the correct detection, red bounding boxes indicate the wrong detection.

452  
453  
454 mantically meaningful details that enhance performance in critical vision tasks, demonstrating its practical utility beyond perceptual quality.

| Metrics \ Methods   | Precision↑    | Recall↑       | F1-score↑     | IoU↑          |
|---------------------|---------------|---------------|---------------|---------------|
| Thermal (SwinPaste) | 0.1800        | 0.2109        | 0.1765        | 0.1941        |
| Thermal (SeeSR)     | 0.3832        | 0.4637        | 0.3849        | 0.3022        |
| Thermal (3M-TI)     | <b>0.4565</b> | <b>0.5455</b> | <b>0.4724</b> | <b>0.3427</b> |
| Reference RGB       | 0.4322        | <b>0.5708</b> | 0.4643        | 0.3359        |
| Thermal (GT)        | 0.4582        | 0.5793        | 0.4887        | 0.3494        |

Table 3. Detection performance comparison across different methods, reference RGB, and GT, evaluated by Precision, Recall, F1-score, and IoU.

| Metrics \ Methods | PSNR↑        | SSIM↑         | LPIPS↓        | MUSIQ↑       |
|-------------------|--------------|---------------|---------------|--------------|
| w/o Reference     | 29.97        | 0.8592        | 0.2106        | 32.85        |
| w/o Augmentation  | 30.04        | 0.8597        | 0.1917        | 34.94        |
| w/o Skip          | 29.86        | 0.8572        | 0.1795        | 36.58        |
| Ours w/ All       | <b>30.09</b> | <b>0.8610</b> | <b>0.1787</b> | <b>36.66</b> |

Table 4. Ablation study of 3M-TI components. Gray cells indicate the best result for each metric.

#### 4.4. Ablation Study

We conduct an ablation study to investigate the contributions of each component. Table 4 and Fig. 7 report quantitative drops and representative visual examples when individual modules are removed. Removing the RGB reference produces markedly blurrier reconstructions (e.g., bicycle spokes and shrub textures), indicating the importance of cross-modal cues. Removing the misalignment augmentation reduces robustness to geometric and temporal offsets, leading to visibly degraded high-frequency details. Removing the skip connection degrades structural fidelity: circular wheels become distorted and geometric consistency is lost. With all components enabled, 3M-TI attains the best balance of fidelity and perceptual quality, producing sharp,

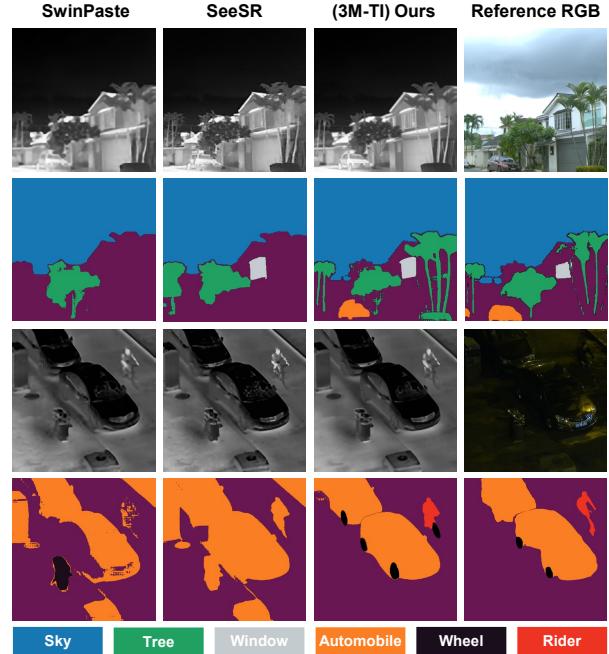


Figure 6. Visualization of segmentation results, where different colors represent different object categories.

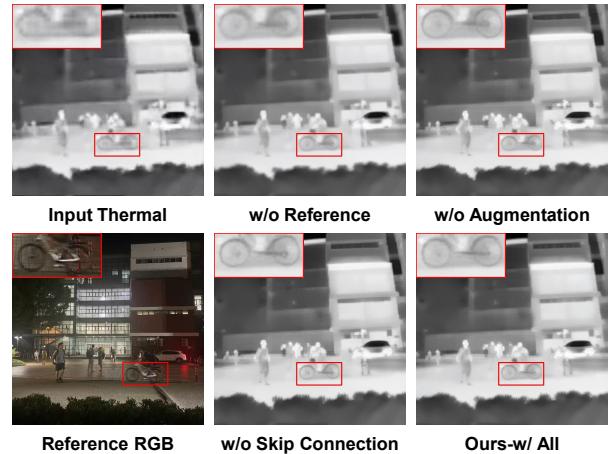


Figure 7. Qualitative results of the ablation study.

semantically consistent thermal reconstructions.

469

## 5. Conclusion

470

We proposed 3M-TI, a cross-modal diffusion framework that integrates calibration-free RGB guidance to enhance thermal image reconstruction. With the cross-modal self-attention module and misalignment augmentation, 3M-TI effectively aligns and fuses information across modalities, producing thermal images with sharper, more realistic, and semantically consistent details. Comprehensive experiments and real-system evaluations demonstrate that 3M-TI not only achieves superior perceptual quality but also significantly benefits downstream machine vision tasks.

471

472

473

474

475

476

477

478

479

480

481 **References**

- [1] Cyprien Arnold, Philippe Jouvet, and Lama Seoud. Swin-fusr: an image fusion-inspired model for rgb-guided thermal image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3027–3036, 2024. 2, 3, 5, 7
- [2] 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. *arXiv preprint arXiv:2311.15127*, 2023. 4
- [3] Antoni Buades, Bartomeu Coll, and Jean-Michel Morel. Non-local means denoising. *Image processing on line*, 1: 208–212, 2011. 2
- [4] Chaofeng Chen, Jiadi Mo, Jingwen Hou, Haoning Wu, Liang Liao, Wenxiu Sun, Qiong Yan, and Weisi Lin. Topiq: A top-down approach from semantics to distortions for image quality assessment. *IEEE Transactions on Image Processing*, 33:2404–2418, 2024. 5
- [5] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22367–22377, 2023. 2
- [6] Zheng Chen, Yulun Zhang, Jinjin Gu, Linghe Kong, Xiaokang Yang, and Fisher Yu. Dual aggregation transformer for image super-resolution. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12312–12321, 2023. 2
- [7] Carlos Cortés-Mendez and Jean-Bernard Hayet. Exploring the usage of diffusion models for thermal image super-resolution: A generic uncertainty-aware approach for guided and non-guided schemes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3123–3130, 2024. 3
- [8] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-d transform-domain collaborative filtering. *IEEE Transactions on image processing*, 16(8):2080–2095, 2007. 2
- [9] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. 2
- [10] Weisheng Dong, Lei Zhang, Guangming Shi, and Xin Li. Nonlocally centralized sparse representation for image restoration. *IEEE transactions on Image Processing*, 22(4): 1620–1630, 2012. 2
- [11] Lu Fang, Mengqi Ji, Xiaoyun Yuan, Jing He, Jianing Zhang, Yinheng Zhu, Tian Zheng, Leyao Liu, Bin Wang, and Qionghai Dai. Engram-driven videography. *Engineering*, 25:101–109, 2023. 2
- [12] Ying Fu, Tao Zhang, Yinjiang Zheng, Debing Zhang, and Hua Huang. Hyperspectral image super-resolution with optimized rgb guidance. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11661–11670, 2019. 2
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 3, 4
- [15] Tao Huang, Yuan Zhang, Mingkai Zheng, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge diffusion for distillation. *Advances in Neural Information Processing Systems*, 36:65299–65316, 2023. 3
- [16] Xinyu Jia, Chuang Zhu, Minzhen Li, Wenqi Tang, and Wenli Zhou. Llivip: A visible-infrared paired dataset for low-light vision. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3496–3504, 2021. 2, 5
- [17] Aditya Kasliwal, Pratinav Seth, Sriya Rallabandi, and San-chit Singhal. Corefusion: Contrastive regularized fusion for guided thermal super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 507–514, 2023. 2, 3, 5, 7
- [18] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2
- [20] Lingshun Kong, Jiangxin Dong, Jianjun Ge, Mingqiang Li, and Jinshan Pan. Efficient frequency domain-based transformers for high-quality image deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5886–5895, 2023. 2
- [21] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yuetong Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022. 2
- [22] Shuang Li, Bingfeng Han, Zhenjie Yu, Chi Harold Liu, Kai Chen, and Shuigen Wang. I2v-gan: Unpaired infrared-to-visible video translation. In *Proceedings of the 29th ACM international conference on multimedia*, pages 3061–3069, 2021. 5
- [23] Xingyuan Li, Zirui Wang, Yang Zou, Zhixin Chen, Jun Ma, Zhiying Jiang, Long Ma, and Jinyuan Liu. Difiisr: A diffusion model with gradient guidance for infrared image super-resolution. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7534–7544, 2025. 3
- [24] Yawei Li, Yuchen Fan, Xiaoyu Xiang, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. Efficient and explicit modelling of image hierarchies for image restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18278–18289, 2023. 2
- [25] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1833–1844, 2021. 2
- [26] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11661–11670, 2019. 2

- 594        *ence on computer vision and pattern recognition workshops*,  
595        pages 136–144, 2017. 2
- 596 [27] Hanzhou Liu, Binghan Li, Chengkai Liu, and Mi Lu. De-  
597 blurdinat: A lightweight and effective transformer for image  
598 deblurring. *CoRR*, 2024. 2
- 599 [28] Jinyuan Liu, Xin Fan, Zhanbo Huang, Guanyao Wu, Risheng  
600 Liu, Wei Zhong, and Zhongxuan Luo. Target-aware dual  
601 adversarial learning and a multi-scenario multi-modality  
602 benchmark to fuse infrared and visible for object detection.  
603 In *Proceedings of the IEEE/CVF conference on computer vi-  
604 sion and pattern recognition*, pages 5802–5811, 2022. 2, 5
- 605 [29] Jinyuan Liu, Bowei Zhang, Qingyun Mei, Xingyuan Li,  
606 Yang Zou, Zhiying Jiang, Long Ma, Risheng Liu, and Xin  
607 Fan. Dcevo: Discriminative cross-dimensional evolutionary  
608 learning for infrared and visible image fusion. In *Proceed-  
609 ings of the Computer Vision and Pattern Recognition Con-  
610 ference*, pages 2226–2235, 2025. 2
- 611 [30] Jinyuan Liu, Bowei Zhang, Qingyun Mei, Xingyuan Li,  
612 Yang Zou, Zhiying Jiang, Long Ma, Risheng Liu, and Xin  
613 Fan. Dcevo: Discriminative cross-dimensional evolutionary  
614 learning for infrared and visible image fusion. In *Proceed-  
615 ings of the Computer Vision and Pattern Recognition Con-  
616 ference*, pages 2226–2235, 2025. 1
- 617 [31] Jiayi Ma, Yong Ma, and Chang Li. Infrared and visible im-  
618 age fusion methods and applications: A survey. *Information  
619 fusion*, 45:153–178, 2019. 1
- 620 [32] Julien Mairal, Michael Elad, and Guillermo Sapiro. Sparse  
621 representation for color image restoration. *IEEE Transac-  
622 tions on image processing*, 17(1):53–69, 2007. 2
- 623 [33] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik  
624 Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans.  
625 On distillation of guided diffusion models. In *Proceed-  
626 ings of the IEEE/CVF conference on computer vision and pattern  
627 recognition*, pages 14297–14306, 2023. 3
- 628 [34] Gaurav Parmar, Taesung Park, Srinivasa Narasimhan, and  
629 Jun-Yan Zhu. One-step image translation with text-to-image  
630 models. *arXiv preprint arXiv:2403.12036*, 2024. 3
- 631 [35] Raghunath Sai Puttagunta, Birendra Kathariya, Zhu Li, and  
632 George York. Multi-scale feature fusion using channel trans-  
633 formers for guided thermal image super resolution. In *Pro-  
634 ceedings of the IEEE/CVF Conference on Computer Vision  
635 and Pattern Recognition*, pages 3086–3095, 2024. 3
- 636 [36] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang  
637 Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng  
638 Yan, et al. Grounded sam: Assembling open-world models  
639 for diverse visual tasks. *arXiv preprint arXiv:2401.14159*,  
640 2024. 6
- 641 [37] Rafael E Rivadeneira, Angel D Sappa, Riad Hammoud,  
642 Jiyong Rao, Hang Zhong, Yu Wang, Shengjie Zhao, Zhi-  
643 wei Zhong, Yung-Hui Li, Shiqi Wang, Qiangqiang Shen,  
644 Han Zhang, and Xuanqi Zhang. Thermal image super-  
645 resolution challenge results-pbvs 2025. In *Proceedings of  
646 the IEEE/CVF conference on computer vision and pattern  
647 recognition*, pages 4630–4639, 2025. 5
- 648 [38] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear  
649 total variation based noise removal algorithms. *Physica D:  
650 nonlinear phenomena*, 60(1-4):259–268, 1992. 2
- [39] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin  
Rombach. Adversarial diffusion distillation. In *European  
Conference on Computer Vision*, pages 87–103. Springer,  
2024. 2, 3
- [40] Mark Sheinin, Aswin C Sankaranarayanan, and Srinivasa G  
Narasimhan. Projecting trackable thermal patterns for dy-  
namic computer vision. In *Proceedings of the IEEE/CVF  
Conference on Computer Vision and Pattern Recognition*,  
pages 25223–25232, 2024. 1
- [41] Ukcheol Shin, Jinsun Park, and In So Kweon. Deep  
depth estimation from thermal image. In *Proceedings of  
the IEEE/CVF Conference on Computer Vision and Pattern  
Recognition*, pages 1043–1053, 2023. 1
- [42] Yang Tan, Haitian Zheng, Yinheng Zhu, Xiaoyun Yuan, Xing  
Lin, David Brady, and Lu Fang. Crossnet++: Cross-scale  
large-parallax warping for reference-based super-resolution.  
*IEEE Transactions on Pattern Analysis and Machine Intelli-  
gence*, 43(12):4291–4305, 2020. 2
- [43] Linfeng Tang, Xinyu Xiang, Hao Zhang, Meiqi Gong, and  
Jiayi Ma. Divfusion: Darkness-free infrared and visible im-  
age fusion. *Information Fusion*, 91:477–493, 2023. 1
- [44] Zitian Tang, Wenjie Ye, Wei-Chiu Ma, and Hang Zhao. What  
happened 3 seconds ago? inferring the past with thermal  
imaging. In *Proceedings of the IEEE/CVF Conference on  
Computer Vision and Pattern Recognition*, pages 17111–  
17120, 2023. 1
- [45] Zhijun Tu, Kunpeng Du, Hanting Chen, Hailing Wang, Wei  
Li, Jie Hu, and Yunhe Wang. Ipt-v2: Efficient image process-  
ing transformer using hierarchical attentions. *arXiv preprint  
arXiv:2404.00633*, 2024. 2
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,  
Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia  
Polosukhin. Attention is all you need. *Advances in neural  
information processing systems*, 30, 2017. 3
- [47] Yufei Wang, Wenhan Yang, Xinyuan Chen, Yaohui Wang,  
Lanqing Guo, Lap-Pui Chau, Ziwei Liu, Yu Qiao, Alex C  
Kot, and Bihan Wen. Sinsr: diffusion-based image super-  
resolution in a single step. In *Proceedings of the IEEE/CVF  
conference on computer vision and pattern recognition*,  
pages 25796–25805, 2024. 3
- [48] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-  
moncelli. Image quality assessment: from error visibility to  
structural similarity. *IEEE transactions on image processing*,  
13(4):600–612, 2004. 5
- [49] Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learn-  
ing for image super-resolution: A survey. *IEEE transactions  
on pattern analysis and machine intelligence*, 43(10):3365–  
3387, 2020. 2
- [50] Jay Zhangjie Wu, Yuxuan Zhang, Haithem Turki, Xuanchi  
Ren, Jun Gao, Mike Zheng Shou, Sanja Fidler, Zan Goj-  
cic, and Huan Ling. Difix3d+: Improving 3d reconstruc-  
tions with single-step diffusion models. In *Proceedings of  
the Computer Vision and Pattern Recognition Conference*,  
pages 26024–26035, 2025. 4
- [51] Rongyuan Wu, Lingchen Sun, Zhiyuan Ma, and Lei Zhang.  
One-step effective diffusion network for real-world image  
super-resolution. *Advances in Neural Information Process-  
ing Systems*, 37:92529–92553, 2024. 3, 5, 7

- 709 [52] Rongyuan Wu, Tao Yang, Lingchen Sun, Zhengqiang Zhang,  
710 Shuai Li, and Lei Zhang. Seesr: Towards semantics-  
711 aware real-world image super-resolution. In *Proceedings of*  
712 *the IEEE/CVF conference on computer vision and pattern*  
713 *recognition*, pages 25456–25467, 2024. 3, 5, 7
- 714 [53] Zongsheng Yue, Jianyi Wang, and Chen Change Loy.  
715 Resshift: Efficient diffusion model for image super-  
716 resolution by residual shifting. *Advances in Neural Infor-*  
717 *mation Processing Systems*, 36:13294–13307, 2023. 3
- 718 [54] Jianing Zhang, Jiayi Zhu, Feiyu Ji, Xiaokang Yang, and  
719 Xiaoyun Yuan. Degradation-modeled multipath diffu-  
720 sion for tunable metalens photography. *arXiv preprint*  
721 *arXiv:2506.22753*, 2025. 3
- 722 [55] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shecht-  
723 man, and Oliver Wang. The unreasonable effectiveness of  
724 deep features as a perceptual metric. In *Proceedings of the*  
725 *IEEE conference on computer vision and pattern recogni-*  
726 *tion*, pages 586–595, 2018. 5
- 727 [56] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li,  
728 Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo,  
729 Yaqian Li, Shilong Liu, et al. Recognize anything: A strong  
730 image tagging model. In *Proceedings of the IEEE/CVF Con-*  
731 *ference on Computer Vision and Pattern Recognition*, pages  
732 1724–1732, 2024. 3
- 733 [57] Yingkai Zhang, Zeqiang Lai, Tao Zhang, Ying Fu, and  
734 Chenghu Zhou. Unaligned rgb guided hyperspectral image  
735 super-resolution with spatial-spectral concordance: Y. zhang  
736 et al. *International Journal of Computer Vision*, pages 1–21,  
737 2025. 2
- 738 [58] Zhifei Zhang, Zhaowen Wang, Zhe Lin, and Hairong Qi. Im-  
739 age super-resolution by neural texture transfer. In *Proceed-*  
740 *ings of the IEEE/CVF conference on computer vision and*  
741 *pattern recognition*, pages 7982–7991, 2019. 2
- 742 [59] Wenda Zhao, Shigeng Xie, Fan Zhao, You He, and Huchuan  
743 Lu. Metafusion: Infrared and visible image fusion via meta-  
744 feature embedding from object detection. In *Proceedings of*  
745 *the IEEE/CVF Conference on Computer Vision and Pattern*  
746 *Recognition*, pages 13955–13965, 2023. 2
- 747 [60] Zhicheng Zhao, Yong Zhang, Chenglong Li, Yun Xiao, and  
748 Jin Tang. Thermal uav image super-resolution guided by  
749 multiple visible cues. *IEEE Transactions on Geoscience and*  
750 *Remote Sensing*, 61:1–14, 2023. 3
- 751 [61] Haitian Zheng, Mengqi Ji, Haoqian Wang, Yebin Liu, and Lu  
752 Fang. Crossnet: An end-to-end reference-based super reso-  
753 lution network using cross-scale warping. In *Proceedings of*  
754 *the European conference on computer vision (ECCV)*, pages  
755 88–104, 2018. 2
- 756 [62] Naishan Zheng, Man Zhou, Jie Huang, Junming Hou, Haoy-  
757 ing Li, Yuan Xu, and Feng Zhao. Probing synergistic high-  
758 order interaction in infrared and visible image fusion. In  
759 *Proceedings of the IEEE/CVF conference on computer vi-*  
760 *sion and pattern recognition*, pages 26384–26395, 2024. 2
- 761 [63] Hang Zhong, Yu Wang, and Shengjie Zhao. Swinpaste: A  
762 swin transformer-based framework for rgb-guided thermal  
763 image super-resolution. In *Proceedings of the Computer Vi-*  
764 *sion and Pattern Recognition Conference*, pages 4589–4594,  
765 2025. 3, 5, 7