

000 SUPPLEMENTARY MATERIAL:  
 001 FLOW-MATCHING GUIDED DEEP UNFOLDING FOR  
 002 HYPERSPECTRAL IMAGE RECONSTRUCTION  
 003  
 004  
 005

006 **Anonymous authors**

007 Paper under double-blind review  
 008  
 009  
 010  
 011

012 1 DERIVATION OF CONDITIONAL OT OBJECTIVE  
 013

014 In the main paper, we directly introduced the regression objective for flow matching, parameterizing  
 015 the conditional vector field  $v_\theta(x_t, t, y)$  and learning it by least-squares regression. For completeness,  
 016 here we provide a detailed derivation showing how the simplified objective arises from the condi-  
 017 tional optimal transport (OT) formulation under a constant-velocity ODE. Specifically, we start from  
 018 the dynamic OT problem (Benamou & Brenier, 2000), connect it to displacement interpolation, and  
 019 show how this naturally leads to the constant velocity target  $x_1 - x_0$ . This derivation clarifies how  
 020 the ODE trajectory definition yields the training loss presented in Eqs. (10)–(11) of the main text.

021 1.1 CONDITIONAL OPTIMAL TRANSPORT VECTOR FIELD DERIVATION  
 022

023 Let  $y$  denote the conditioning signal (e.g., a CASSI measurement). We consider two conditional  
 024 distributions on  $\mathbb{R}^d$ : the source  $\mu_0^y$  with density  $p_0(x | y)$  and the target  $\mu_1^y$  with density  $p_1(x |$   
 025  $y)$ . A time-dependent conditional density  $\mu_t^y$  transported by a velocity field  $v_t^y(x)$  satisfies the  
 026 (conditional) continuity equation

$$027 \partial_t \mu_t^y(x) + \nabla \cdot (\mu_t^y(x) v_t^y(x)) = 0, \quad t \in [0, 1], \quad (1)$$

028 with boundary conditions  $\mu_0^y, \mu_1^y$ .  
 029

030 In the *dynamic OT* formulation (Benamou & Brenier, 2000), the 2-Wasserstein geodesic between  
 031  $\mu_0^y$  and  $\mu_1^y$  is obtained by minimizing the conditional kinetic energy  
 032

$$033 \inf_{\{\mu_t^y, v_t^y\}} \int_0^1 \int \frac{1}{2} \|v_t^y(x)\|_2^2 \mu_t^y(x) dx dt \quad \text{s.t.} \quad \begin{cases} \partial_t \mu_t^y + \nabla \cdot (\mu_t^y v_t^y) = 0, \\ \mu_0^y, \mu_1^y \text{ given.} \end{cases} \quad (2)$$

036 When a (conditional) OT map (Villani et al., 2008)  $T_y : \mathbb{R}^d \rightarrow \mathbb{R}^d$  exists, the geodesic curve is the  
 037 *displacement interpolation* (McCann, 1997)

$$038 X_t = (1 - t)X_0 + tT_y(X_0), \quad X_0 \sim \mu_0^y, \quad t \sim \text{Unif}[0, 1], \quad (3)$$

040 with interpolated law  $\mu_t^y = ((1 - t)\text{Id} + tT_y)_\# \mu_0^y$ . Along each sample path  $\gamma(t) = (1 - t)x_0 + tx_1$ ,  
 041 the geodesic velocity is constant:

$$042 \dot{\gamma}(t) = x_1 - x_0, \quad \forall t \in [0, 1]. \quad (4)$$

044 Thus the (conditional) ground-truth vector field is

$$045 v^*(x_t, t, y) = x_1 - x_0, \quad \text{whenever } x_t = (1 - t)x_0 + tx_1, \quad (x_0, x_1) \sim \pi_y, \quad (5)$$

046 where  $\pi_y$  is an (optimal) coupling between  $\mu_0^y$  and  $\mu_1^y$ .  
 047  
 048

049 1.2 FROM FLOW MATCHING REGRESSION TO THE CONSTANT-VELOCITY TARGET

050 Flow matching (FM) learns a parametric conditional field  $v_\theta(x_t, t, y)$  by least-squares regression:

$$052 \hat{\theta} = \arg \min_{\theta} \mathbb{E}_{t, x_t \sim \mu_t^y} [\|v^*(x_t, t, y) - v_\theta(x_t, t, y)\|_2^2]. \quad (6)$$

The sampling procedure consists of three main steps. First, a time variable  $t$  is uniformly sampled from the interval  $[0, 1]$ . Next, a pair  $(x_0, x_1)$  is drawn from the coupling distribution  $\pi_y$ . Finally, the intermediate state  $x_t$  is obtained by linearly interpolating between the two samples, that is,  $x_t = (1 - t)x_0 + tx_1$ . Under the constant-velocity ODE  $\dot{x}_t = x_1 - x_0$ , the target is  $x_1 - x_0$ . Hence the FM loss reduces to

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E}_{t, (x_0, x_1) \sim \pi_y} [\|(x_1 - x_0) - v_{\theta}((1 - t)x_0 + tx_1, t, y)\|_2^2], \quad (7)$$

which matches the simplified loss used in our main text.

### 1.3 INTERPRETATION

At optimum, the regressor satisfies

$$v_{\theta}^*(x_t, t, y) = \mathbb{E}[x_1 - x_0 \mid x_t, t, y], \quad (8)$$

i.e., the learned field is the conditional mean velocity of the OT geodesic. This connects flow matching training to estimating the OT vector field, ensuring theoretical consistency. Our proposed *mean-velocity loss* further enforces this property, improving global consistency when couplings are noisy.

## 2 EXTRA EXPERIMENTAL SETTINGS OF OUR FULL MODEL

**Stage design.** We follow an unfolding design with 10 stages. Parameter sharing across stages is enabled, consistent with previous LADE-DUN (Wu et al., 2024). To further enhance the learned prior, we integrate flow matching with a mean velocity regularization weight of  $\lambda_{\text{reg}} = 5$ .

**Training protocol.** We train 10-stage models with a batch size of 1 and patch size of  $256 \times 256$ . In both training phases, the maximum epoch is set to 300 and the learning rate is initialized at  $4 \times 10^{-4}$ . Gradient clipping is enabled to stabilize training. We employ the MultiStepLR scheduler from Paszke et al. (2019) with milestones at epochs  $\{50, 100, 150, 200, 250\}$  and decay factor  $\gamma = 0.5$ .

**Initialization and ODE solver.** We adopt Gaussian initialization with noise scale 1.0 and apply an ODE sampler based on the `rk45` solver (Dormand & Prince, 1980), with tolerance set to  $10^{-5}$ .

These additional settings highlight the key practical configurations in our pipeline and provide useful reference for reproducing and extending our work.

## 3 THEORETICAL INSIGHT EXPLANATION AND DETAILED MODEL ARCHITECTURE

### 3.1 ON THEORETICAL INSIGHT OF USING FLOW MATCHING PRIOR

We acknowledge that the manuscript emphasizes empirical evidence. The key insight is that flow matching enables high-quality latent priors without the heavy sampling overhead inherent to diffusion. In LADE-DUN, the number of diffusion steps must be aggressively reduced (e.g., 16 steps) to maintain feasible inference cost, which inevitably compromises the generated prior. In contrast, our flow-matching module supports RK-based sampling without step-count restrictions and therefore can freely increase sampling precision while keeping the runtime comparable to diffusion-based few-step priors. This capability directly improves prior quality and consequently reconstruction accuracy. The mean-velocity constraint further enhances this advantage by encouraging a globally consistent vector field, leading to more stable training and improved hyperspectral consistency.

### 3.2 LATENT ENCODER

The latent encoder (**LE**) follows the structure illustrated in Fig. 1a. Given the normalized measurement  $y_{\text{norm}}$  and input  $x$ , the encoder first applies a **Pixel Unshuffle** operation to rearrange spatial information into the channel dimension, enabling more compact feature extraction. This is followed by a sequence of stacked **Mblocks** (Howard et al., 2019), where local and global spatial dependencies are progressively modeled. To increase representational capacity while maintaining efficiency,

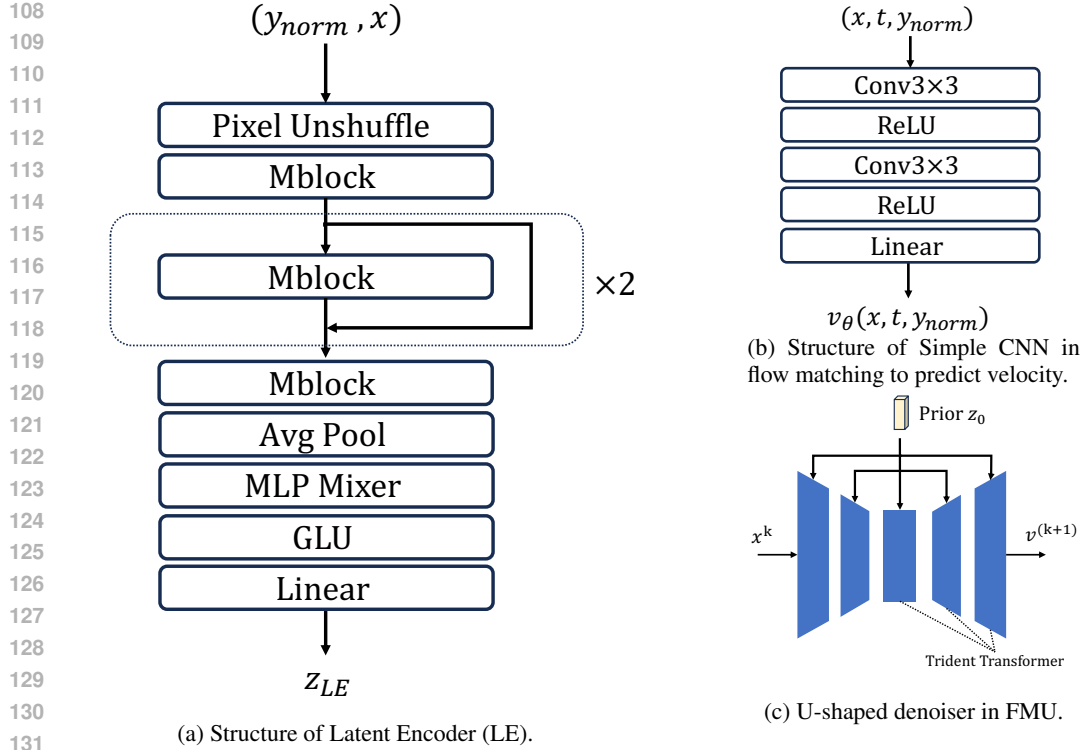


Figure 1: Detailed structures of LE, velocity predictor CNN in Flow Matching and denoiser in FMU.

two **Mblocks** are repeated in a loop as indicated in the figure. Subsequently, the features are further refined through an additional **Linear** layer, before being aggregated by an **average pooling**.

On top of the pooled features, the encoder employs an *MLP Mixer* and a *Gated Linear Unit (GLU)* to enhance channel interactions and improve non-linear feature fusion. Finally, a *linear projection* maps the encoded representation into the latent vector  $z$ , which serves as the compact and informative embedding for subsequent stages of our framework. This design ensures a balance between efficiency and expressive power, providing a strong latent representation tailored for HSI.

### 3.3 DENOISER IN FMU

We use a U-shaped Denoiser which consists of Trident Transformer (Wu et al., 2024). The structure of denoiser is shown in Fig. 1c, which incorporates different scales of Trident Transformer.

### 3.4 TRIDENT TRANSFORMER

we use Trident Transformer from Wu et al. (2024) to effectively aggregate high-quality degradation-free prior knowledge from flow matching for compensation. The overview of Trident Transformer assisted with flow-matching prior is in Fig. 2. The input of the transformer is  $U_i$ , in which  $i$  means that the Trident Transformer is of the  $i$ -th scales. The input is first divided to  $U_i^S$  and  $U_i^C$  by channels to extract different information. Then either is sent to Spectral Flow or Cross Spectral&Prior Flow. Flow-matching prior is sent to Cross Prior Flow to assist denoising. Inside each flow there are convolutions and QKV calculations. Refer to Wu et al. (2024) for more details.

### 3.5 VELOCITY PREDICTOR IN FLOW MATCHING

As for the network used to predict velocity field in flow matching, we adopt Simple CNN after ablations. Detailed structure is shown in Fig. 1b.

## 4 ADDITIONAL RESULTS AND TRAINING PSEUDOCODE

### 4.1 SIMULATION EXPERIMENT VISUALIZATIONS

We first show additional simulation experiment results to highlight the ability of our method to recover spatial details and spectral structures (see Fig. 3). These visualizations complement the main

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

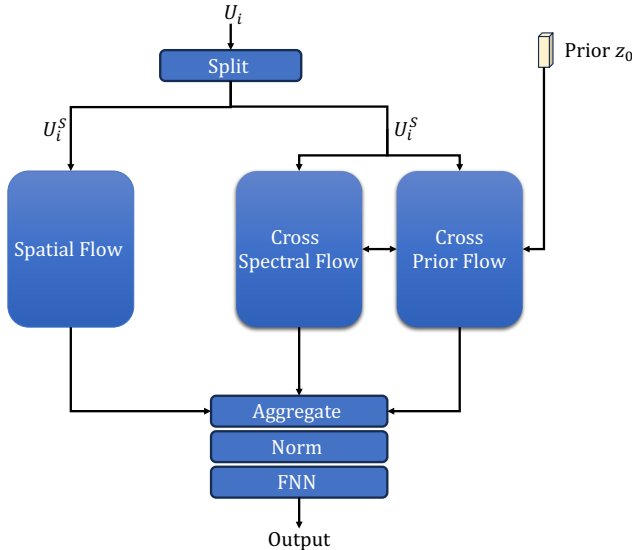


Figure 2: Structure of Trident Transformer.

text and provide a more comprehensive view of the reconstruction quality achieved by our approach. They confirm that our method produces sharper edges and more consistent spectral responses compared with competing with previous works.

#### 4.2 TWO-PHASE TRAINING PROCEDURE

Besides visual results, we present the two-phase training procedure of our model. In Phase 1, the latent encoder and unfolding module are jointly trained for accurate recovery from measurements. In Phase 2, flow matching is introduced to learn a generative prior, further enhancing reconstruction quality and robustness. The pseudocode for these two stages is summarized in Algorithm 1, 2. During inference, we set  $I_E = (y_{norm}, 0)$  instead of  $(y_{norm}, x_{gt})$  to ensure that no ground-truth information is used.

#### 4.3 FULL QUANTITATIVE CASSI RESULTS

In addition to the results reported in the main paper, we also provide full quantitative evaluations on the CASSI benchmark and additional visual comparisons across more scenes. (See table 3)

#### 4.4 EVALUATION METRICS, FAIRNESS, AND EFFICIENCY OF COMPARISON

To further assess generalization beyond PSNR and SSIM, we additionally report results on three widely used hyperspectral quality metrics: SAM and ERGAS, which are specifically designed to evaluate spectral fidelity of HSI reconstructions, and MANIQA, which measures perceptual quality in the RGB space. In particular, SAM quantifies the average spectral angle between reconstructed and ground-truth spectra, and ERGAS evaluates the global relative reconstruction error across all bands. MANIQA is computed by projecting the hyperspectral cube into the RGB space and applying a no-reference IQA model to the resulting pseudo-RGB image. As shown in Table 1, FMU achieves the best performance on four out of five metrics, significantly improving PSNR, SSIM, SAM, and ERGAS over all baselines. Although LADE-DUN is marginally better in MANIQA (0.6294 vs. 0.6293), this difference is negligible, whereas our gains on hyperspectral-specific metrics (SAM, ERGAS) are substantial. These results demonstrate that our evaluation goes beyond PSNR/SSIM and that FMU consistently offers better spectral fidelity and overall reconstruction quality under a fair and comprehensive comparison protocol. Beyond reconstruction quality, we also compare the efficiency of FMU with the diffusion-prior baseline LADE-DUN under a matched unfolding depth. Table 2 reports PSNR, FLOPs, and wall-clock inference time. Under a very similar FLOPs budget (96.69G vs. 98.84G), FMU improves PSNR from 40.97 dB to 42.13 dB (+1.16 dB) while reducing inference time from 0.521 s to 0.390 s (about  $1.3\times$  faster). This indicates that, despite slightly higher

Method	PSNR (dB) $\uparrow$	SSIM $\uparrow$	SAM (rad) $\downarrow$	ERGAS $\downarrow$	MANIQA $\uparrow$
BIRNAT	38.76	0.9792	0.1064	16.79	0.6119
BiSRNet	35.84	0.9524	0.2126	23.12	0.5714
CST-L-Plus	37.73	0.9762	0.1380	19.21	0.6180
DAUHST-9stg	38.81	0.9815	0.1120	17.05	0.6273
DGSMP	32.99	0.9471	0.1688	33.57	0.5664
MST++	39.21	0.9825	0.0989	16.25	0.6256
TSA-Net	37.68	0.9759	0.1103	18.82	0.6124
SPECAT	40.37	0.9860	0.0877	14.04	0.6292
LADE-DUN-10stg	40.97	0.9882	0.0708	12.70	<b>0.6294</b>
<b>FMU (Ours)</b>	<b>42.13</b>	<b>0.9900</b>	<b>0.0651</b>	<b>11.26</b>	0.6293

Table 1: Multi-metric comparison on the optical-filter-based system. FMU achieves the best performance on four out of five metrics.

Method	PSNR (dB) $\uparrow$	FLOPs (G)	Inference Time (s) $\downarrow$
LADE-DUN-10stg	40.97	<b>96.69</b>	0.521
<b>FMU (Ours)</b>	<b>42.13</b>	98.84	<b>0.390</b>

Table 2: Efficiency comparison with the diffusion-based LADE-DUN under a similar computational budget. FMU achieves both higher PSNR and lower inference latency.

theoretical FLOPs, the flow-matching prior yields a strictly better quality–efficiency trade-off than the diffusion prior in LADE-DUN.

#### 4.5 ADDITIONAL EXPERIMENTS ON KAIST AND ICVL (ZERO-SHOT GENERALIZATION)

To further evaluate robustness beyond the training distribution, we additionally conduct experiments on two unseen real-world datasets: KAIST and ICVL. Unlike the main benchmark, which follows the standard optical-filter-based setting, these two datasets contain different spectral characteristics and scene distributions, making them suitable for testing zero-shot generalization. **KAIST.** FMU preserves its advantage on previously unseen KAIST scenes, improving PSNR (43.44 vs. 43.04) and reducing ERGAS (13.52 vs. 14.09), indicating more accurate hyperspectral reconstruction and stronger spectral consistency. **ICVL.** The performance gap becomes more pronounced on ICVL. FMU improves PSNR by +1.66 dB and also yields clearly lower SAM and ERGAS. This suggests that FMU is more robust when the degradation process and spectral statistics differ from those observed during training. Overall, these results show that FMU consistently maintains its superiority without any fine-tuning on new datasets. This supports that flow-matching priors generalize more reliably than diffusion-based priors within physics-guided unfolding frameworks.

## REFERENCES

- Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the monge-kantorovich mass transfer problem. *Numerische Mathematik*, 2000.
- John R Dormand and Peter J Prince. A family of embedded runge-kutta formulae. *Journal of computational and applied mathematics*, 1980.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *ICCV*, 2019.
- Robert J McCann. A convexity principle for interacting gases. *Advances in mathematics*, 1997.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019.

Method	PSNR (dB)	SSIM
TwIST	23.12	0.669
DNU	30.74	0.863
MST++	35.99	0.951
BIRNAT	37.58	0.960
LRSDN	36.08	0.938
DAUHST-9stg	38.36	0.967
DADF-Plus-3	39.28	0.972
PADUT-5stg	37.84	0.967
RDLUF-MixS2-3stg	37.56	0.963
LADE-DUN-3stg	38.31	0.972
PADUT-12stg	38.89	0.974
RDLUF-MixS2-9stg	39.57	0.974
LADE-DUN-10stg	40.09	0.979
<b>FMU</b>	<b>41.02</b>	<b>0.988</b>

Table 3: CASSI reconstruction results of different methods.

More scenes from KAIST and ICVL				
Method	PSNR (dB) $\uparrow$	SSIM $\uparrow$	SAM (rad) $\downarrow$	ERGAS $\downarrow$
More scenes from KAIST				
FMU (Ours)	<b>43.44</b>	<b>0.9942</b>	0.0645	<b>13.52</b>
LADE-DUN	43.04	0.9939	<b>0.0638</b>	14.09
Scenes from ICVL				
FMU (Ours)	<b>33.58</b>	<b>0.9492</b>	<b>0.1046</b>	<b>73.20</b>
LADE-DUN	31.92	0.9396	0.1264	81.80

Table 4: Zero-shot reconstruction results on additional KAIST and ICVL scenes.

**Algorithm 1:** Phase 1: Learn Prior and Reconstruction**Data:** Ground-truth HSI  $\mathbf{x}$ , Measurement  $\mathbf{y}$ , Sensing matrix  $\mathbf{A}$ , Unfolding Network  $\phi_{FMU}$ **Result:** Trained Latent Encoder LE, Trained Unfolding Network  $\phi_{FMU}$ 

```

1  $\mathbf{y}_{norm} \leftarrow \Phi^T (\Phi \Phi^T)^{-1} \mathbf{y}$ ;
2  $I_E \leftarrow \text{concat}(\mathbf{y}_{norm}, \mathbf{x})$ ;
3 for  $epoch \leftarrow 1$  to  $E$  do
4   Shuffle dataset;
5   for each batch  $B$  in training data do
6      $I_E \leftarrow \text{concat}(\mathbf{y}_{norm}, \mathbf{x})$ ;
7      $\mathbf{z}_{LE} \leftarrow LE(I_E)$ ;
8      $\hat{\mathbf{x}} \leftarrow \phi_{FMU}(\mathbf{y}, \Phi, \mathbf{z}_{LE})$ ;
9      $\mathcal{L}_1 \leftarrow \|\hat{\mathbf{x}} - \mathbf{x}\|_1$ ;
10    Update parameters  $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_1$ ;
11  end
12 end

```

Cédric Villani et al. *Optimal transport: old and new*. Springer, 2008.Zongliang Wu, Ruiying Lu, Ying Fu, and Xin Yuan. Latent diffusion prior enhanced deep unfolding for snapshot spectral compressive imaging. In *ECCV*, 2024.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

---

**Algorithm 2:** Phase 2: Flow Matching for Prior Generation
 

---

**Data:** Measurement  $\mathbf{y}$ , Sensing matrix  $\mathbf{A}$ , Encoder LE, Unfolding Network  $\phi_{FMU}$ , Flow

Matching denoise network  $\phi_{FM}$

**Result:** Trained  $\phi_{FMU}$ , Trained  $\phi_{FM}$

```

1  $\mathbf{y}_{norm} \leftarrow \Phi^\top (\Phi \Phi^\top)^{-1} \mathbf{y};$ 
2  $I_E \leftarrow \text{concat}(\mathbf{y}_{norm}, \mathbf{x});$ 
3 Freeze encoder LE from Phase 1;
4 for  $epoch \leftarrow 1$  to  $E$  do
5   Shuffle dataset;
6   for each batch  $B$  in training data do
7      $I_E \leftarrow \text{concat}(\mathbf{y}_{norm}, \mathbf{x});$ 
8      $\mathbf{z}_{LE} \leftarrow LE(I_E);$ 
9      $\hat{\mathbf{z}}_{FM} \leftarrow \phi_{FM}(I_E);$ 
10     $\hat{\mathbf{x}} \leftarrow \phi_{FMU}(\mathbf{y}, \Phi, \mathbf{z}_E);$ 
11     $\mathcal{L}_2 \leftarrow \|\hat{\mathbf{z}}_{FM} - \mathbf{z}_{LE}\|_2^2 + \mathcal{L}_{rec};$ 
12    Update parameters  $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_2;$ 
13  end
14 end

```

---

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

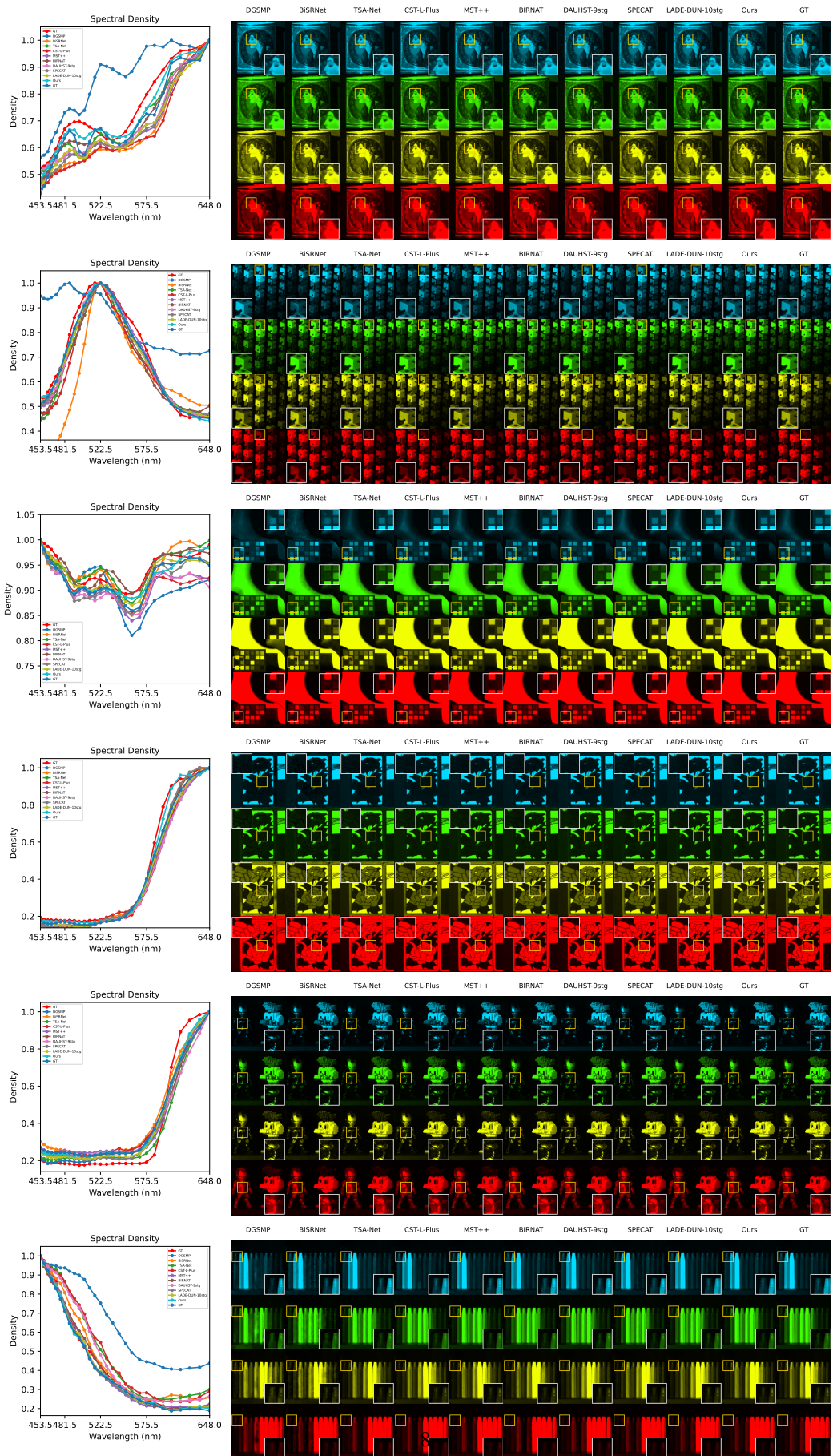


Figure 3: Qualitative comparisons across test cases. Each row shows results on a different example.