# Learning Image-Adaptive Codebooks for Class-Agnostic Image Restoration (Supplementary material)

## A. Semantic Grouped Classes

As described in the paper, we aggregate the 150 classes in ADE20K dataset [4] into five super-classes to train our basis codebooks. The overall details of how we divide the five sub-datasets are shown in Fig. 1. Although this is not a rigorous categorization, the codebook visualization in Section 3.1 empirically demonstrates that the grouping is meaningful to some extent.

## B. Network Architectures

In conjunction with the encoder-decoder network used in our AdaCode, we elaborate on the detailed architectures of the encoder $E$, the decoder $G$, and the discriminator $D$. For $E$ and $G$, we adopt the same autoencoder as VQGAN [2] in stage I and the same structure as FeMaSR [1] in stage II&III. For $D$, we adopt the same U-Net discriminator with spectral normalization as Real-ESRGAN [3].

## C. More Results

### C.1. Ablation on Codebook Size

We conduct experiments to empirically select the number of code entries in each basis codebook, as shown in Table. 1. Considering the tradeoff between performance and computation cost, we set the basis codebooks' size as $\{512, 256, 512, 256, 256\} \times 256$.

### C.2. Codebook Visualization

We visualize all the codes in our five basis codebooks in Fig. 2. As we discussed in Section 5, it is yet unclear how many basis codebooks and how many code entries in each codebook we need. It is also reflected by the visualization that there might be some redundancies in the codebooks.

### C.3. Qualitative Results

We show more results and comparisons for image reconstruction, super-resolution, and image inpainting in Fig. 3, Fig. 4, and Fig. 5, which empirically demonstrate the effectiveness of AdaCode.

Table 1: Stage I reconstruction performance on each super-class with different number of code entries. The chosen size is marked in red.

| Super-Class | Codebook Size | PSNR | SSIM | LPIPS |
|---|---|---|---|---|
| Architectures | $256 \times 256$ | 24.096 | 0.667 | 0.149 |
| | $512 \times 256$ | 24.260 | 0.684 | 0.144 |
| Indoor Objects | $256 \times 256$ | 26.565 | 0.788 | 0.110 |
| | $512 \times 256$ | 26.630 | 0.789 | 0.110 |
| Natural Scenes | $256 \times 256$ | 27.014 | 0.723 | 0.124 |
| | $512 \times 256$ | 27.693 | 0.743 | 0.110 |
| Street Views | $256 \times 256$ | 26.677 | 0.748 | 0.126 |
| | $512 \times 256$ | 26.937 | 0.755 | 0.127 |
| Portraits | $256 \times 256$ | 29.914 | 0.838 | 0.097 |
| | $512 \times 256$ | 29.662 | 0.837 | 0.098 |

## References

[1] Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo. Real-World Blind Super-Resolution via Feature Matching with Implicit High-Resolution Priors. In *Proceedings of ACM International Conference on Multimedia (MM)*, 2022.

[2] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming Transformers for High-Resolution Image Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[3] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-ESRGAN: Training Real-World Blind Super-Resolution With Pure Synthetic Data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[4] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene Parsing through ADE20K Dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2017.
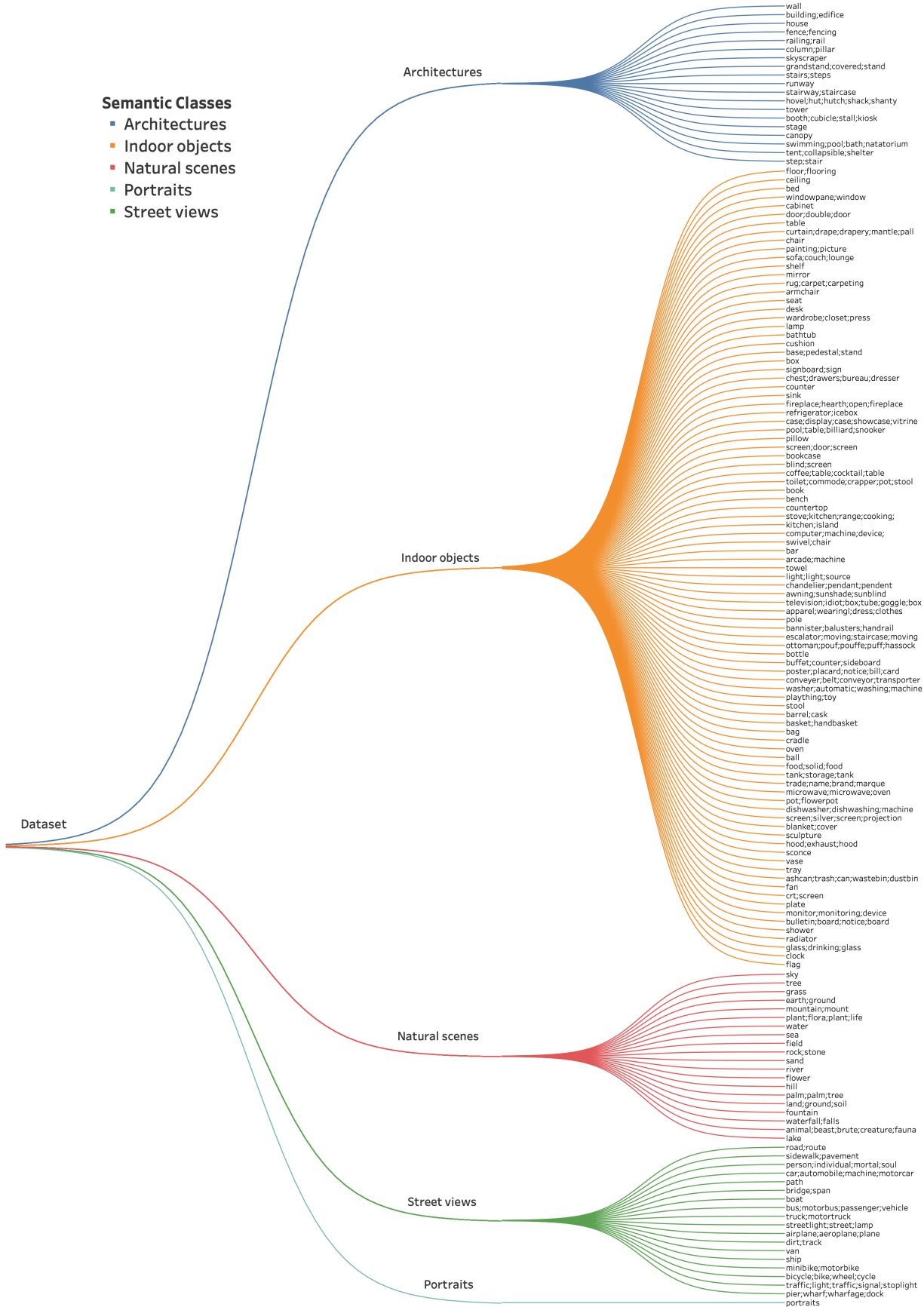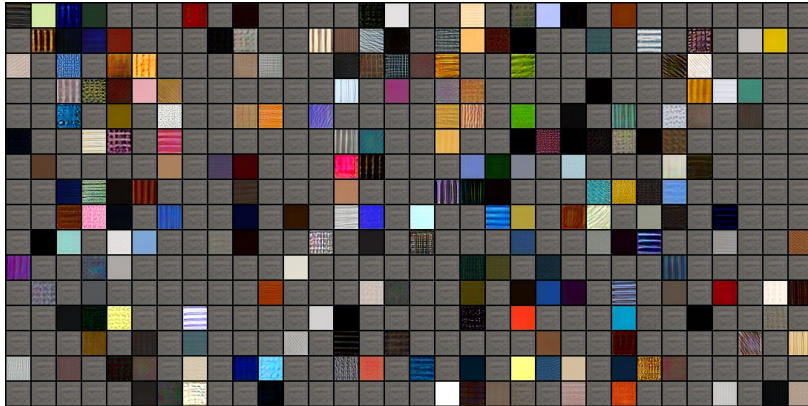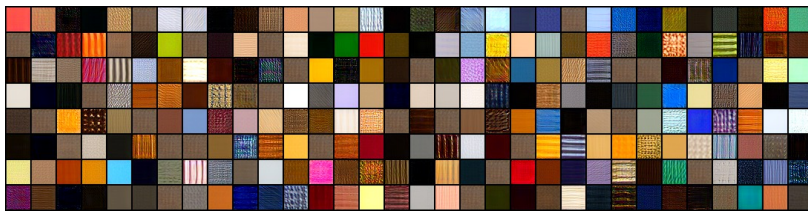
**Semantic Classes**
- Architectures
- Indoor objects
- Natural scenes
- Portraits
- Street views

**Architectures**
wall
building;edifice
house
fence;fencing
railing;rail
column;pillar
skyscraper
grandstand;covered;stand
stairs;steps
runway
stairway;staircase
hovel;hut;hutch;shack;shanty
tower
booth;cubicle;stall;kiosk
stage
canopy
swimming;pool;bath;natatorium
tent;collapsible;shelter
step;stair

**Indoor objects**
floor;flooring
ceiling
bed
windowpane;window
cabinet
door;double;door
table
curtain;drape;drapery;mantle;pall
chair
painting;picture
sofa;couch;lounge
shelf
mirror
rug;carpet;carpeting
armchair
seat
desk
wardrobe;closet;press
lamp
bathtub
cushion
base;pedestal;stand
box
signboard;sign
chest;drawers;bureau;dresser
counter
sink
fireplace;hearth;open;fireplace
refrigerator;icebox
case;display;case;showcase;vitrine
pool;table;billiard;snooker
pillow
screen;door;screen
bookcase
blind;screen
coffee;table;cocktail;table
toilet;commode;crapper;pot;stool
book
bench
countertop
stove;kitchen;range;cooking;
kitchen;island
computer;machine;device;
swivel;chair
bar
arcade;machine
towel
light;light;source
chandelier;pendant;pendent
awning;sunshade;sunblind
television;idiot;box;tube;goggle;box
apparel;wearingl;dress;clothes
pole
bannister;balusters;handrail
escalator;moving;staircase;moving
ottoman;pouf;pouffe;puff;hassock
bottle
buffet;counter;sideboard
poster;placard;notice;bill;card
conveyer;belt;conveyor;transporter
washer;automatic;washing;machine
plaything;toy
stool
barrel;cask
basket;handbasket
bag
cradle
oven
ball
food;solid;food
tank;storage;tank
trade;name;brand;marque
microwave;microwave;oven
pot;flowerpot
dishwasher;dishwashing;machine
screen;silver;screen;projection
blanket;cover
sculpture
hood;exhaust;hood
sconce
vase
tray
ashcan;trash;can;wastebin;dustbin
fan
crt;screen
plate
monitor;monitoring;device
bulletin;board;notice;board
shower
radiator
glass;drinking;glass
clock
flag

**Natural scenes**
sky
tree
grass
earth;ground
mountain;mount
plant;flora;plant;life
water
sea
field
rock;stone
sand
river
flower
hill
palm;palm;tree
land;ground;soil
fountain
waterfall;falls
animal;beast;brute;creature;fauna
lake

**Street views**
road;route
sidewalk;pavement
person;individual;mortal;soul
car;automobile;machine;motorcar
path
bridge;span
boat
bus;motorbus;passenger;vehicle
truck;motortruck
streetlight;street;lamp
airplane;aeroplane;plane
dirt;track
van
ship
minibike;motorbike
bicycle;bike;wheel;cycle
traffic;light;traffic;signal;stoplight
pier;wharf;wharfage;dock

**Portraits**
portraits

Figure 1: **Groups of 150 classes in ADE20K dataset [4].**
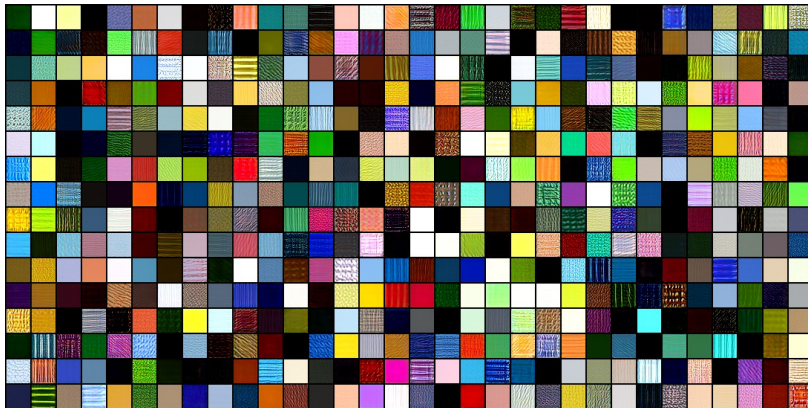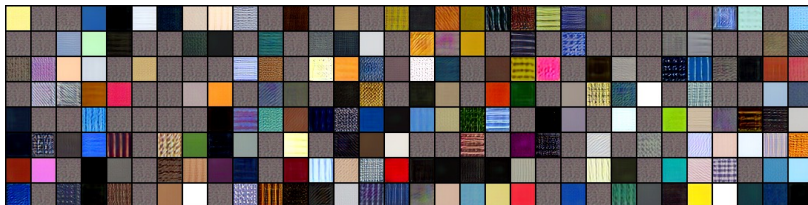
$Z_1$: Architecture

$Z_2$: Indoor Objects

$Z_3$: Natural Scenes
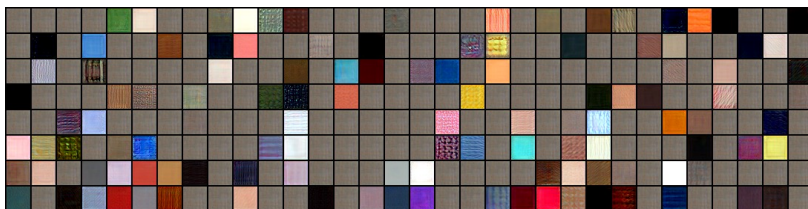
$Z_4$: Street Views

$Z_5$: Portraits

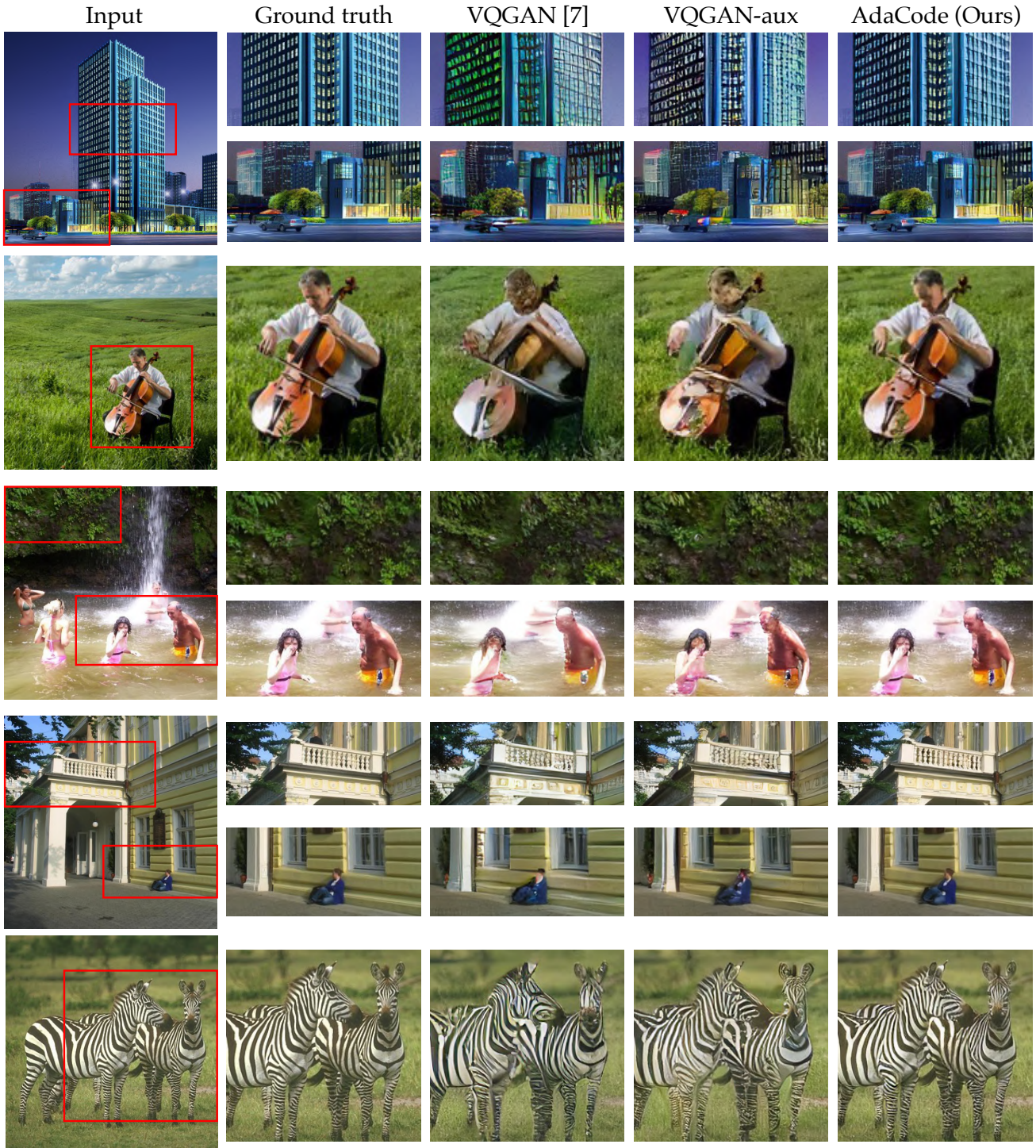Figure 2: **Visualization of all the basis codebooks.**
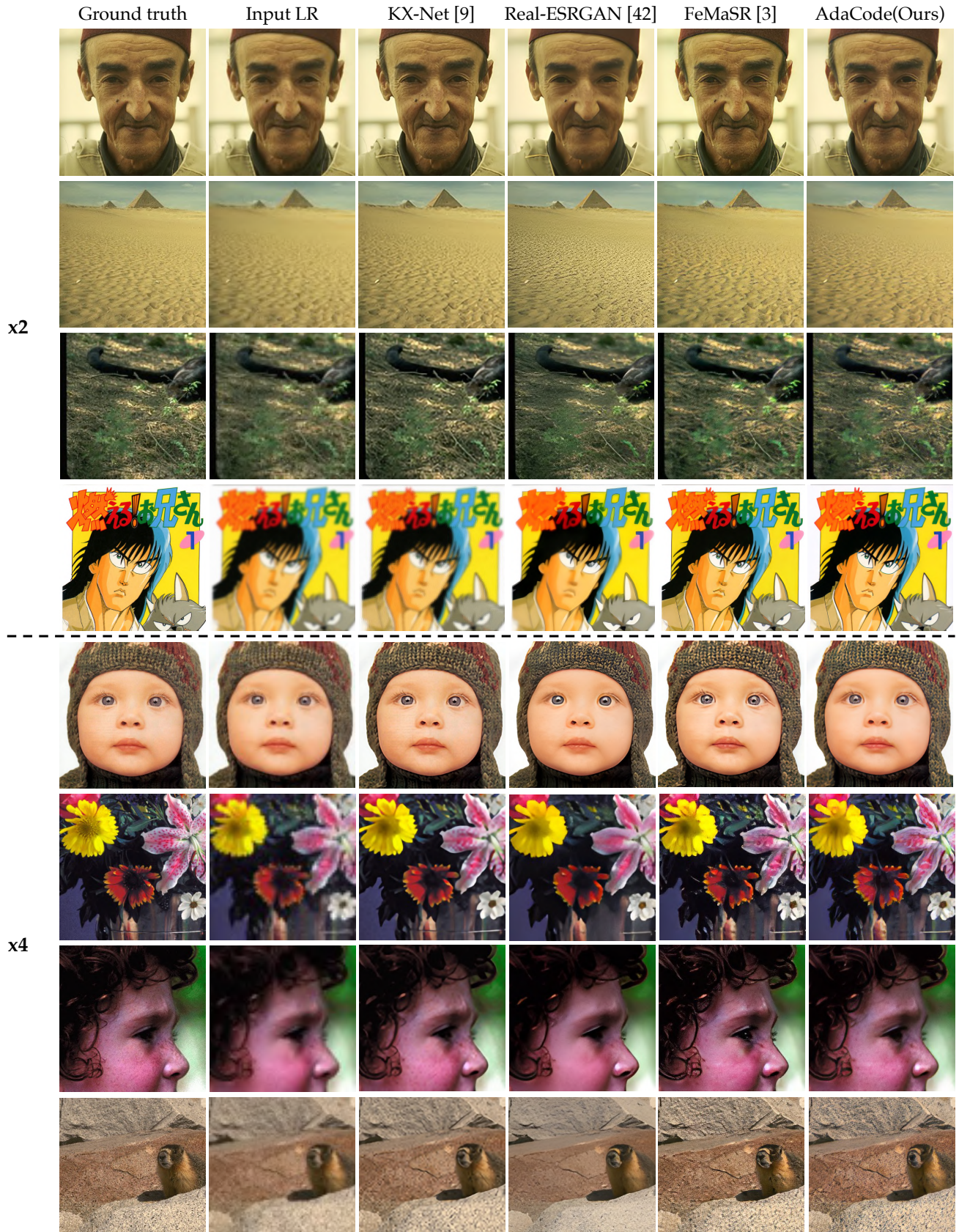
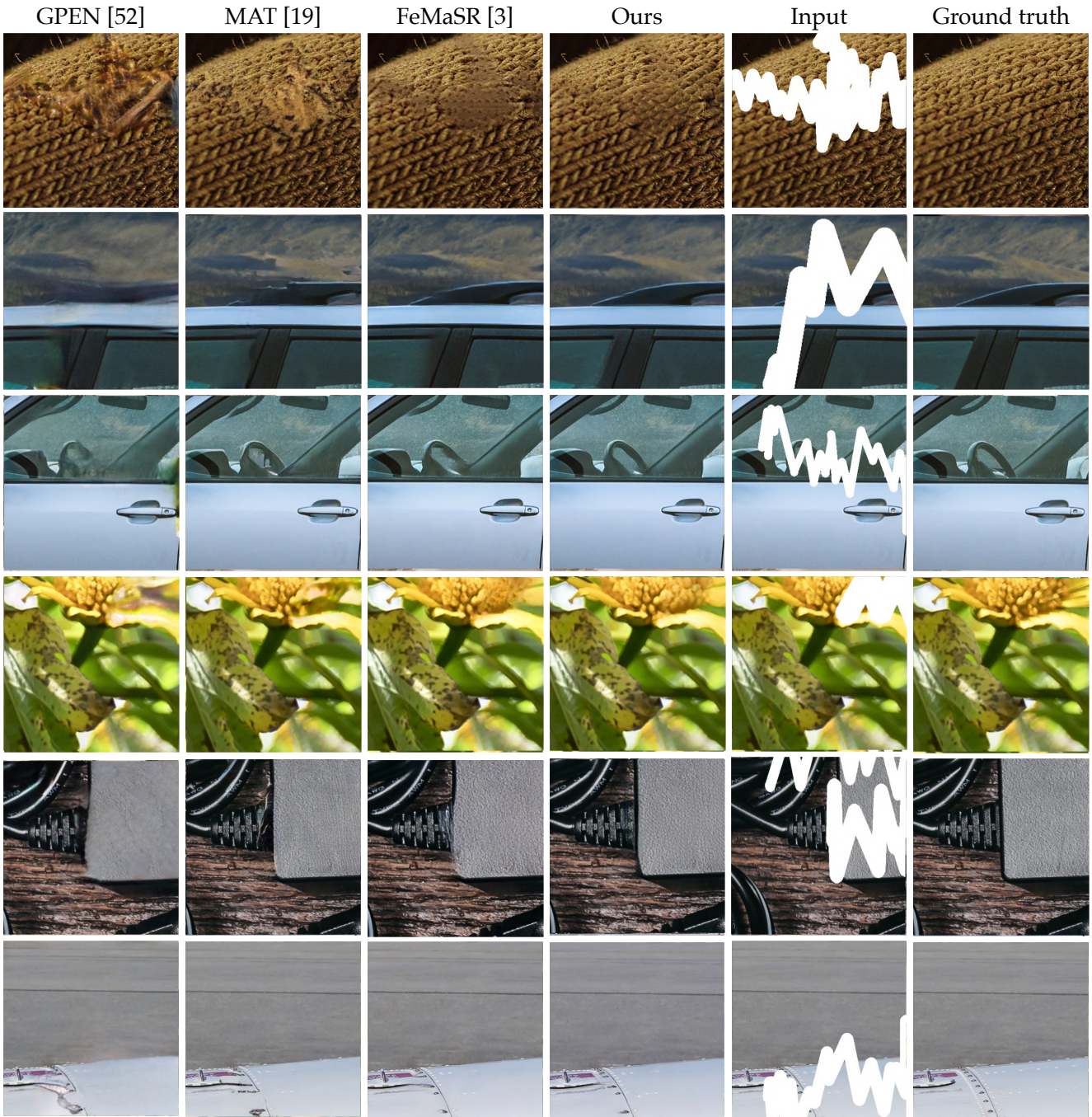Figure 3: **More results on Image Reconstruction.**

Figure 4: **More results on Super-Resolution.**

Figure 5: **More results on Image Inpainting.**