

Segmentation-Based Parametric Painting

Supplementary Material



Figure 10. Styles: (a) Expressionist, consisting of thick, textured, loose and low fidelity representation where subject matter is not a priority. (b) Painterly, think, shaped strokes with a loose touch, enough to balance the scene recognition and not yield a high-fidelity painting style, (c) Abstract, scattered, fewer strokes that capture the essence of the scene.

	Realistic	Painterly	Abstract Expressionist	Abstract
Num. Passes (P)	4	4	1	2
Num. layers (K)	1	<i>var</i>	<i>var</i>	<i>var</i>
Num. Strokes (T)	[36, 49, 64, 81]	[9, 16, 16, 9]	[9]	[9, 16]
Num. DAMs (V)	$(H//128 * W//128)$	[25, 25, 50]	$(H//128 * W//128)$	[25, 25]
DAMs mode	<i>U</i>	<i>U</i> if $p = 1$ else <i>S</i>	<i>U</i>	<i>S</i>

Table 5. Style hyperparameters by default. Users can fine-tune these parameters for further control over the desired styles.

6.0. Paintings at Large Resolution

More paintings at large resolution are shown at the end of this document.

6.1. Optimization Method

This section provides further implementation details of our method. The pseudocode is shown in Algorithm 2. We establish the nomenclature as follows:

- Input image X , canvas C
- Number of painting passes P (int), from 1 to P
- Number of painting layers K (int), from 0 to $K - 1$
- Number of Dynamic Attention Maps (DAMs), V (int). If uniform, DAMs follow an array of 128×128 patches uniformly distributed on image X and canvas C . $V = (H//128 * W//128)$
- Number of stroke parameters, T (int), per DAM, different per p , and can be different per k
- Semantic Segmentation Network, W , determines K

6.1. Implementation Details

For all styles, we use Adam optimizer [18] with a learning rate of 0.0002 and betas 0.5 and 0.99. We optimize all painting layers for 300 iterations, and set the canvas background color to black. For our semantic segmentation network, we use the DETR model [2], available at Hugging Face, which consists of a CNN (ResNet) backbone followed by an encoder-decoder Transformer, trained on the COCO Panoptic Segmentation task.

6.2. Style Parameters

Our method takes in an image X and a style input as a string between “realistic”, “painterly”, “abstract”, and “expressionist”. Note that our optimization method does not apply transfer style techniques. Rather, we adjust optimization hyper-parameters such as the amount of dynamic attention maps (DAMs), stroke budget, or painting passes, to achieve different style variations. We found the combination of such parameters empirically, and we keep the parameters shown at Table 5 as default. However, users can change it to fine-tune the paintings to further adapt the painting. All styles follow a coarse-to-fine strategy, halving the stroke thickness following the formula: $a_p = 2^{1-p}a_1$, where $a_1 = 0.8 * 128$, that is, the 80% of the length of a 128×128 patch. The set of segmentation layers $\{k\}_0^{K-1}$ is generated by a segmentation network W , such that $\{k\} = W(X)$.

642 **Realistic Style** Our realistic style is generated by setting
643 $P = 4$, $K = 1$, setting DAMs to uniform, and thus $V =$
644 $(H//128 * W//128)$. The stroke budget T on each patch
645 follows $s_p = (p + 5)^2$.

646 **Painterly Style** The painterly style is generated by setting
647 $P = 4$, applying a first foundational painting pass consist-
648 ing of uniform DAMs ($V = (H//128 * W//128)$) and no
649 segmentation distinction ($K_p = 1$ if $p = 1$). For $p > 1$, we
650 activate segmentation layers, that is, $K_p = |W(X)|$, where
651 X is the input image and $|\cdot|$ is the number of the segmen-
652 tation layers; and switch to selective DAMs. The stroke
653 budget T on each patch is limited to 16 except for the first
654 and last passes which is 9.

655 **Abstract Styles** Within the abstract style, our method can
656 yield different degrees of abstraction. For simplicity, we
657 grouped them together in the main paper. However, there
658 are two fundamental differences in the way they work. Fig-
659 ure 10 shows the two abstract styles. The expressionist style
660 (a) is captured by the use of loose, efficiently optimized
661 thick strokes with heavy texture. This produces a painting
662 that prioritizes the effect and ambient of the scene by not
663 caring too much about accurate shapes. The abstract style
664 (c) is a looser representation than the painterly style, yet
665 capturing shapes and subject matter better than the expres-
666 sionistic abstract style.

667 Given a segmentation network W , the expressionist style
668 is generated by setting DAMs to uniform, stroke budget
669 $T = 9$, $P = 1$, $K = |W(X)|$, where X is the input image
670 and $|\cdot|$ is the number of the segmentation layers. The ab-
671 stract style is achieved by setting DAMs to selective, $P = 2$,
672 stroke budget $T = [9, 16]$, $V = [25, 25]$, and $K = |W(X)|$.

673 6.3. Stroke Parameterization and Differentiable 674 Renderer

675 We use the differentiable renderer provided by [15], which
676 consists of 3 fully connected layers of dimensionality 512,
677 followed by convolutional layers, and a sigmoid nonlinear-
678 ity function to obtain a $1 \times 128 \times 128$ patch. A diagram of
679 the architecture is shown in Figure 11.

680 Each stroke is formally represented as $s_t =$
681 $(x_0, y_0, x_1, y_1, x_2, y_2, r_0, r_2, t_0, t_1, R, G, B)$, where (x, y)
682 are the Cartesian coordinates of the curve’s control points,
683 (r) denotes the radii at the endpoints, and (t) indicates the
684 transparency levels. Typically, the transparency variables
685 are set to unity, rendering all strokes opaque.

686 In the rendering pipeline, given a stroke budget of T ,
687 the parameters are either randomly or semi-randomly ini-
688 tialized in the configuration space of dimensions $T \times 13$.
689 Subsequently, these parameters are fed into the differen-
690 tiable rendering function G , which produces an array of T
691 rasterized alpha strokes with dimensions 128×128 . The

Algorithm 2 Pseudo-code of our method.

```

1: procedure PAINT(imagePath, style)
2:   painter  $\leftarrow$  INITIALIZEPAINTER(imagePath, style)
3:   Segment.LayersK  $\leftarrow$  W(imagePath)
4:   for  $p \leftarrow 1$  to  $P$  do ▷ Painting passes
5:     for  $k \leftarrow 0$  to  $K - 1$  do ▷ Painting layers
6:       DAMs  $\leftarrow$  GETDAMS(mode,  $p$ ,  $|V|$ )
7:       canvasp  $\leftarrow$  GETCANVAS(DAMs,  $p$ )
8:       strokes  $\leftarrow$  INITSTROKES(DAMs,  $T$ ,  $p$ )
9:       strokes  $\leftarrow$  OPTIMIZE(strokes, canvas, style,  $k$ )
10:      canvasp  $\leftarrow$  UPDATECANVAS(canvasp, strokes)
11:      Append canvasp to allCanvasesL
12:    end for
13:  end for
14:  COMPOSEFINALPAINTING(canvas)
15:  return final painting
16: end procedure
17: procedure INITIALIZEPAINTER(imagePath,
   parameters)
18:   Process input image and set up canvas
19:   Define renderer, segmentation network, and percep-
   tual network
20:   return painter object
21: end procedure
22: procedure GETDAMS(mode,  $p$ ,  $V$ )
23:   Computes DAMs based on mode and max number
    $V$ 
24:   return image  $X$  crops and DAMs coordinates
25: end procedure
26: procedure GETCANVAS(DAMs,  $p$ )
27:   Crops canvasp based on DAMs
28:   return canvas crops
29: end procedure
30: procedure INITSTROKES(DAMs,  $T$ ,  $p$ )
31:   Initializes random stroke parameters based on
   DAMS,  $T$  and  $p$ 
32:   return Stroke parameters
33: end procedure
34: procedure OPTIMIZE(strokes, canvas, style,  $k$ )
35:   Optimizes strokes based on style parameters
36:   Perform rendering to compute loss
37:   return strokes
38: end procedure
39: procedure COMPOSEFINALPAINTING(canvas)
40:   Merge layers and apply final adjustments to canvas
41:   Save final painting
42: end procedure

```

color channels are obtained by element-wise multiplication
of the alpha stroke with the respective RGB values. Each

692
693

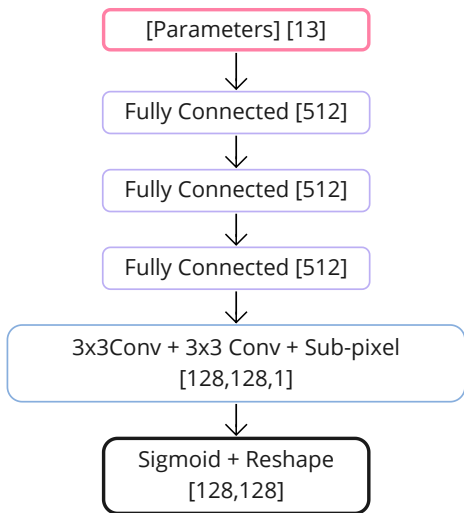


Figure 11. Neural differentiable renderer architecture provided by [15].

694 rasterized alpha stroke is sequentially composited onto the
695 canvas in accordance with the blending equation:

$$696 \quad \text{canvas}_t = \text{canvas}_{t-1} \odot \alpha_t + \text{stroke}_t \quad (3)$$

697 Here, t represents the temporal index within the stroke
698 budget T , α_t refers to the alpha channel generated by
699 the differentiable renderer G , and canvas_{t-1} denotes the
700 canvas state prior to the incorporation of the new stroke.
701 The initial canvas background is configured to be black
702 across all experimental setups. The loss function \mathcal{L} is com-
703 puted between the reference image I and the finalized art-
704 work C_T , followed by the backpropagation of the gradient
705 for parameter optimization. This cycle iterates until the op-
706 timization algorithm converges to a stable state.

707 6.4. Impact of Stroke Initialization

708 Stroke initialization critically affects the algorithm’s perfor-
709 mance. In its most rudimentary form, the algorithm ran-
710 domly initializes stroke parameters, uniformly distributed
711 within the range $[0, 1]$. In a more advanced form, we can
712 evenly distribute the total number of strokes along the two
713 dimensions of the canvas in the form of a grid. This is done
714 by anchoring the middle brushstroke coordinates to such
715 grid, and then control the length of the stroke and the po-
716 sition of the start and end brushstroke coordinates using a
717 Gaussian distribution. To illustrate, Figure 13 presents a se-
718 ries of examples showcasing paintings rendered at 128×128

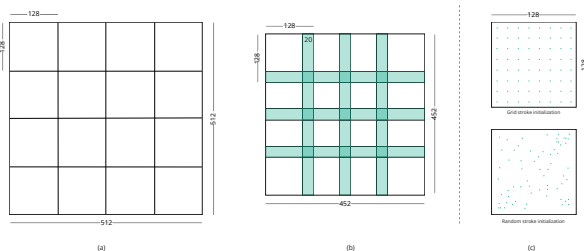


Figure 12. Canvas composition by patches. (a) Patches are organized into a structured grid without overlaps. (b) Example of an organization of patches in a structured grid with overlaps of 20 pixels. (c) Initialization of strokes before being optimized: strokes are evenly distributed across the x and y axis of the patch (upper part). Strokes are randomly sampled from a uniform distribution (bottom part).

pixels with different stroke initializations and stroke widths. These instances demonstrate that grid-based initialization, combined with Gaussian-controlled stroke lengths, substantially simplifies the optimization process. Intriguingly, although random stroke initialization may yield lower L_1 losses, it frequently results in less visually appealing outputs.

6.5. Patch Strategy and Stroke Initialization

Our method optimizes the stroke parameters of all N patches in the general canvas in batch. Once the strokes are optimized and all patches are painted using soft blending, we compose them back together with overlaps. For all experiments we use overlap of 20 pixels, as shown in Figure 12 (b).

We also perform an analysis on the computational cost in the number of patches and number of strokes. Figure 14 shows an analysis of the cost of increasing the number of patches (left) and increasing the stroke budget (right). Increasing resolution to 4x results in a cost of 1.56x. The bottleneck of SBR algorithms, regardless of whether we use patch-based strategies or not, is the number of strokes because of its linear relationship. That is, doubling the number of strokes results in double time to compute. This is because, generally, a blending operation is performed serially, instead of in batch. We perform this analysis on a single RTX A6000 Nvidia GPU.

6.6. Loss Ablation

We perform an ablation studies in loss functions. Generally, loss function, stroke model and number of strokes are responsible for the outcome of a painting. In this section, we analyze the impact of different loss function, see Figure 15. While L_1 works fine for reconstruction (top row), it

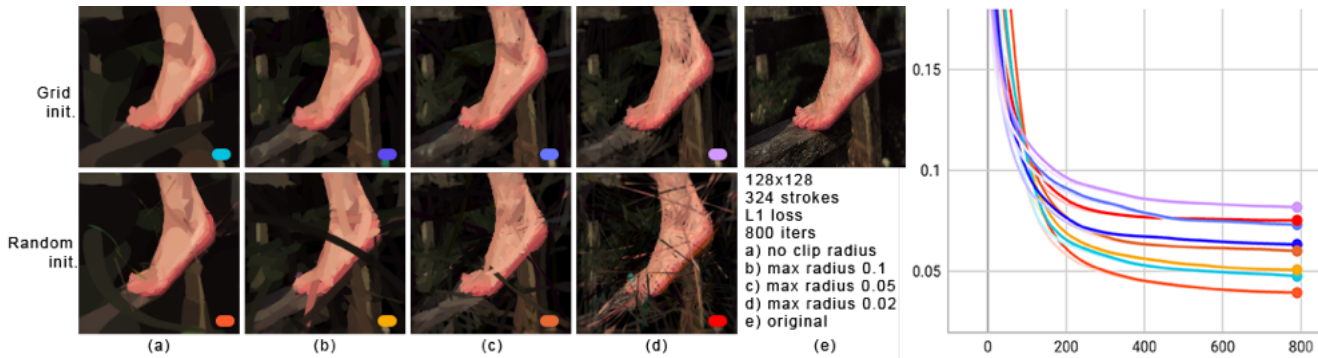


Figure 13. Difference of quality of paintings at 128x128 between random and grid stroke initialization. Top row corresponds to a grid stroke initialization, bottom row corresponds to random stroke initialization. Columns correspond to clipping the maximum stroke width the model is allowed to paint with: (a) no constraint, (b) maximum width $0.2 \times \text{canvas size}$, (c) $0.1 \times \text{canvas size}$ and (d) $0.04 \times \text{canvas size}$. All paintings use 324 strokes, and L1 as pixel loss function. We let the optimization run for 800 iterations.

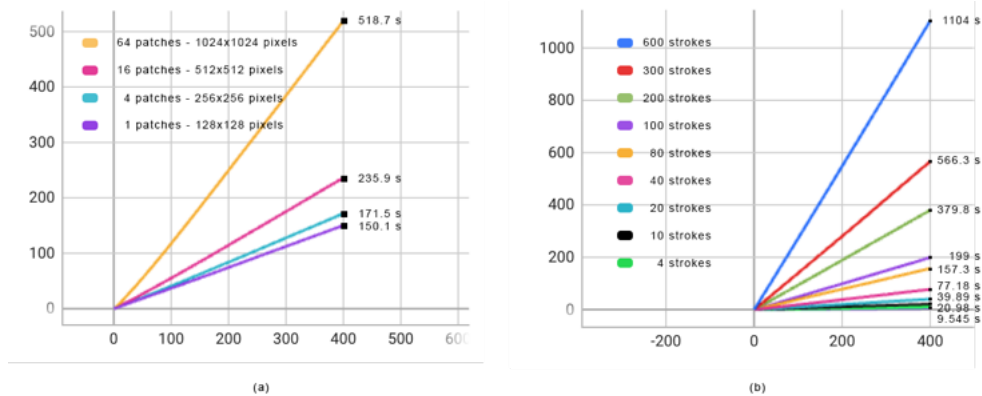


Figure 14. Patch and stroke cost. (a) Computational cost of painting by patches approach with a budget of 100 strokes per patch and 400 optimization steps. Each patch is 128x128. Increasing image size by double does not translate to a linear cost. (b) The computational cost of increasing stroke budget is approximately linear. Computed on a single RTX A6000 Nvidia GPU.

751 seems unable to capture finer details. Adding a perceptual
 752 loss helps alleviates this issue. We try using CLIP [29] as
 753 a type of perceptual loss, but it seems to add some noise to
 754 the painting. Second and fourth rows use a small weight for
 755 perceptual and CLIP losses, resulting in imperceptible contribu-
 756 tions. However scaling up the weight to 0.1 seems to
 757 work better for reconstructions, being the perceptual loss in
 758 the third row the one that achieves best reconstructions. For
 759 realistic paintings, we choose a combination of L1 + percep-
 760 tual loss, computed as shown in Section 3.3 in the main
 761 paper. CLIP loss is calculated as follows:

$$762 \quad \mathcal{L} = (\text{CLIP}(I) - \text{CLIP}(C))^2 \quad (4)$$

763 where CLIP is the image encoder of the CLIP ViT-B/32
 764 model.

7. Comparison with SOTA

765 We show more stroke distributions in Figure 16 across dif-
 766 ferent subject matters. SNP and L2P present a more uni-
 767 form distribution of strokes, barely differentiating between
 768 semantic areas or objects in the scene. While L2P presents
 769 a clear uniform distribution across the entire canvas, SNP
 770 is able to scatterly capture some regions of interest, but still
 771 suffer to bias the painting towards the main objects in the
 772 scene. PT does the absolute opposite, that is, it fixates a
 773 disproportionate amount of strokes to a specific area, which
 774 does not match semantic regions, while leaving the rest of
 775 the canvas with very few and big strokes. This is clearly
 776 visible in the picture at the middle row. The distribution is
 777 concentrated in a blob spanning the building and the pave-
 778 ment equally. The same event happens in the last row where
 779 the lake and the lavender flowers that are closer to the water
 780

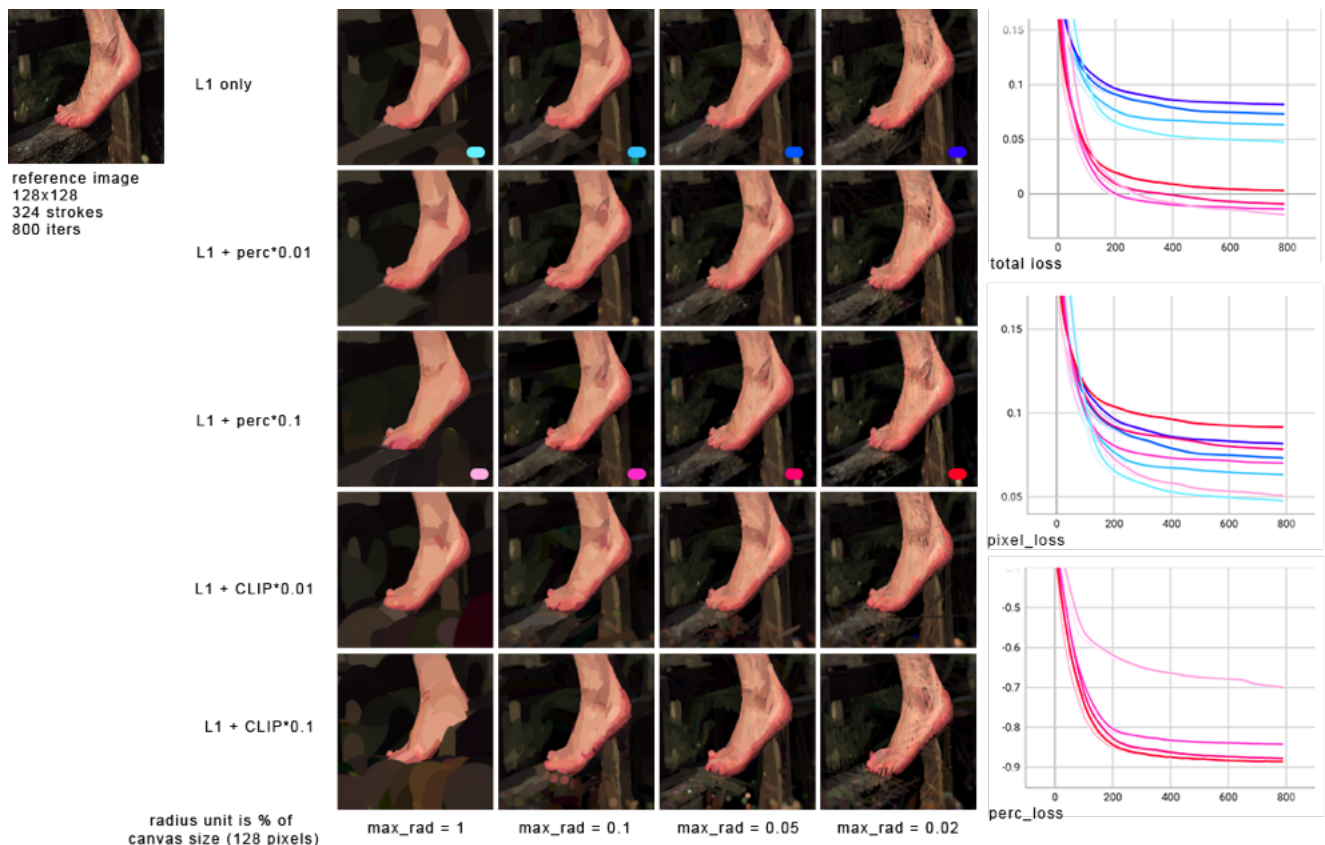


Figure 15. Difference in loss functions. We show the difference in the outcome when using L1, perceptual loss and CLIP loss. A weight of 0.01 seems very low for either perceptual or CLIP to make an impact. However, a weight of 0.1 makes the perceptual loss define some fine details than L1 loss alone cannot (see definition in the toes). CLIP loss introduces some noise in the painting.

781 receive a higher quantity of strokes.

782 Ours, however, sits in between L2P / SNP and PT. While
783 focusing on a specific object in the picture, with semantic
784 coherence, it is able to also distribute strokes on other im-
785 portant areas of the picture. For instance, the bottom row
786 shows how flowers, lake and mountains are better captured,
787 while leaving sky with fewer strokes. In the case of the
788 building, strokes are concentrated around it, since this is
789 clearly the most important semantic part of the image. How-
790 ever, our method is able to redirect its stroke distribution
791 according to the user’s desire as shown in Figure 17. In this
792 figure, we show a painting dissection based on our seman-
793 tic segmentation pipeline. Except for last row, each column
794 shows emphasis on a different k -th layer. The second row
795 shows a distribution of all the strokes in the painting, while
796 the fifth shows a distribution of only the strokes pertaining
797 to each k -th layer.

8. User Study

798 The user study is divided in two sections with 6 images
799 each. The first section evaluates realistic paintings, and asks
800 which painting is a better painting of an input image. The
801 second section evaluates visual appeal, and asks participants
802 to choose the painting that looks more natural and less com-
803 puter made. In total, we have 12 two-way side by side com-
804 parisons. The participants see first an input photograph, and
805 below the photograph we place a pair of paintings. We ask
806 to choose left or right painting. We compare our paintings
807 with three recent methods: an optimization-based method
808 “Stylized Neural Painting” (SNP) [41], a learning-based
809 method with Transformer “Paint Transformer” (PT) [23],
810 and a RL method “Learning to Paint” (L2P) [15]. We ran-
811 domize the order in which we show the pairings.
812

813 The complete user study is shown in Figures 24 to 26.
814 All four pairs from figure Figure 24 and the first two rows
815 in Figure 25 correspond to the first section of the user study.
816 These pairings correspond to the left side of Table 2 in the
817 main paper. IDs are Figure 24 top row: A5, second row:

818 A3, third row: A4, fourth row: A6. In Figure 25, top row:
819 A1, second row: A2. The last two rows of Figure 25 and
820 Figure 26 belong to the second section of the study. They
821 have the following IDs; Figure 25 third row: B1, fourth
822 row: B5. Figure 26, from top to bottom: B3, B6, B2 and
823 B4, respectively.

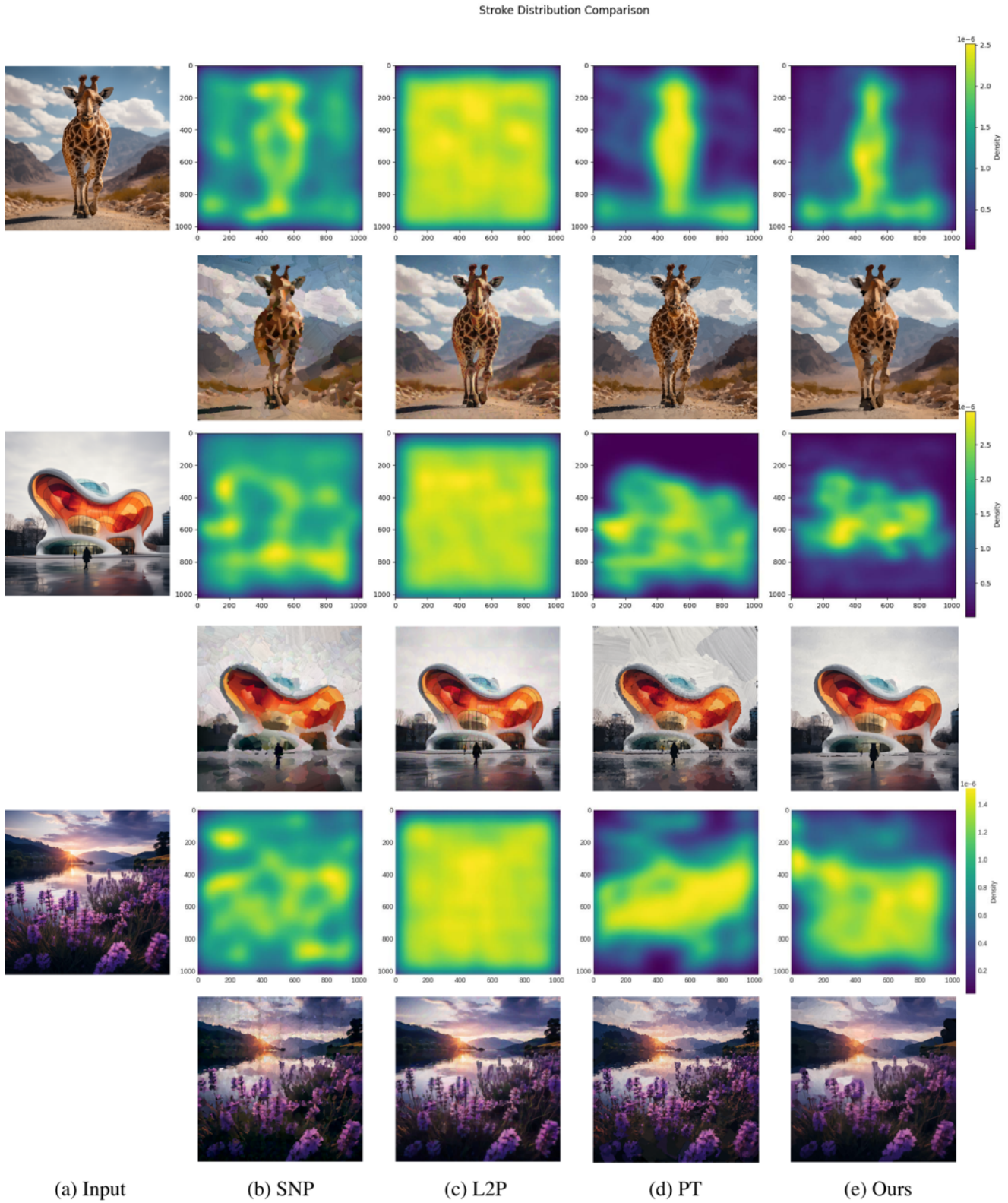


Figure 16. Stroke distribution analysis: SNP and L2P present a uniform distribution across the canvas. L2p is completely scene-agnostic, and SNP cannot distinguish between semantic areas. PT presents a severe contrast between objects and background, and struggles to differentiate between semantic areas of the image. Ours distribute strokes based on semantic areas.

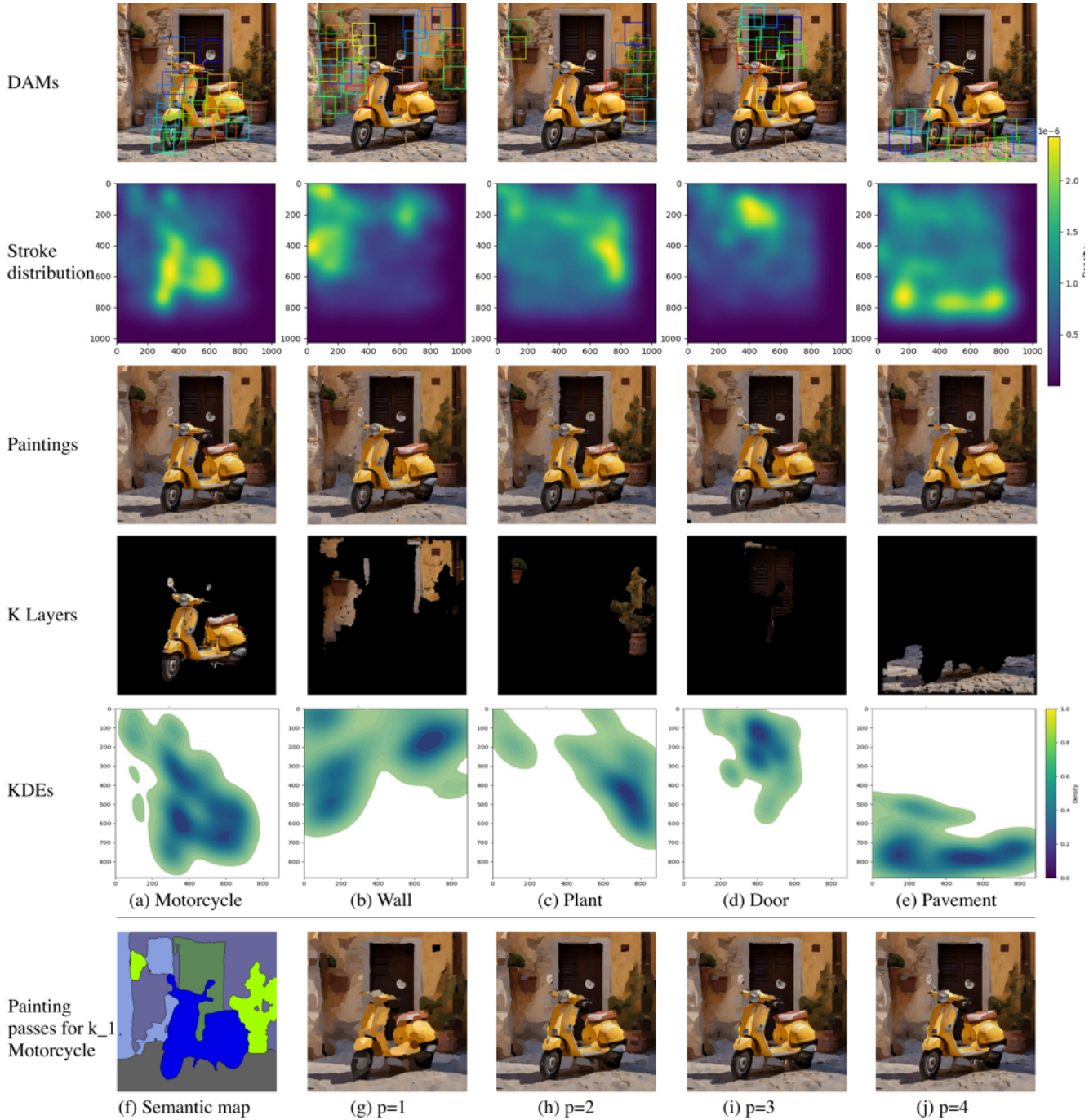


Figure 17. Painting analysis based on semantic layers and dynamic attention maps (DAMs). Top row shows DAMs per semantic layer k . Second row shows the overall stroke distribution over the canvas based on biasing the painting towards different semantic areas, and the third row shows the corresponding paintings. The fourth row shows each layer k independent from the rest of the painting. The fifth row shows the Kernel Density Estimation of the strokes that correspond to only such layer k . The bottom row shows the semantic map, and the four painting passes, from coarse to fine, with a bias on the vehicle over the rest of the scene.



(a) Ours - Realism

(b) L2P

Figure 18. Comparison of our method (a) against L2P [15] (b) on realistic images. Note how L2P paintings have visible seams and blurry areas.



(a) Ours - Painterly

(b) PT

(c) SNP

Figure 19. Comparison of our method (a) against PT [23] (b) and SNP [41] (c) on painterly images.



(a)
Expressionist



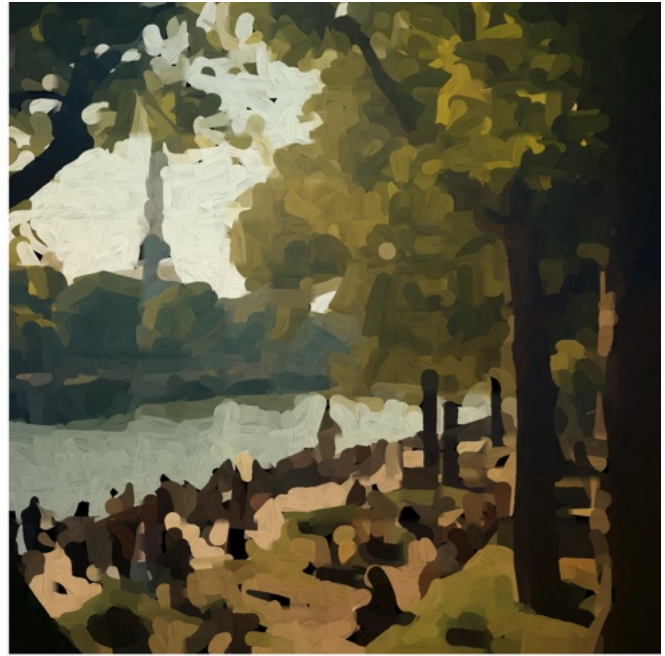
(b) Painterly



(c) Realism



(a) Input



(b) Expressionist



(c) Painterly



(d) Realism



(a) Input



(b) Expressionist



(c) Painterly



(d) Realism



(a)
Expressionist



(b) Painterly



(c) Realism



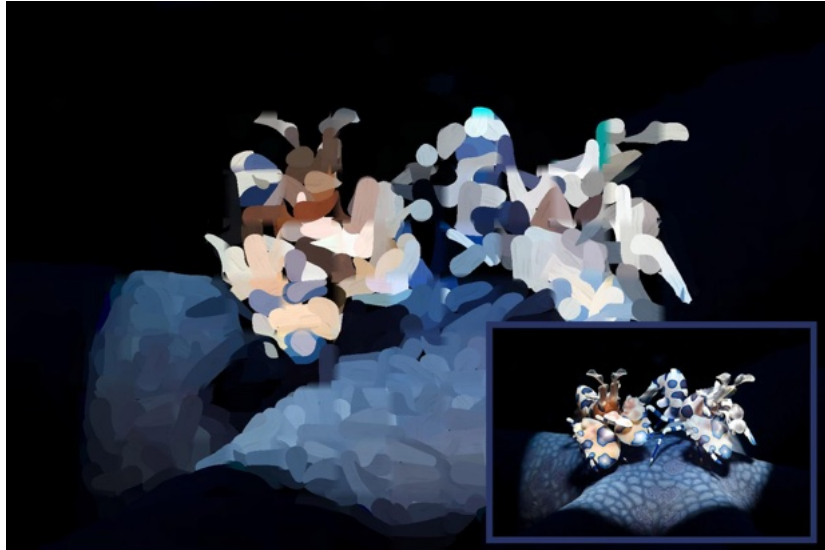
(a)
Expressionist



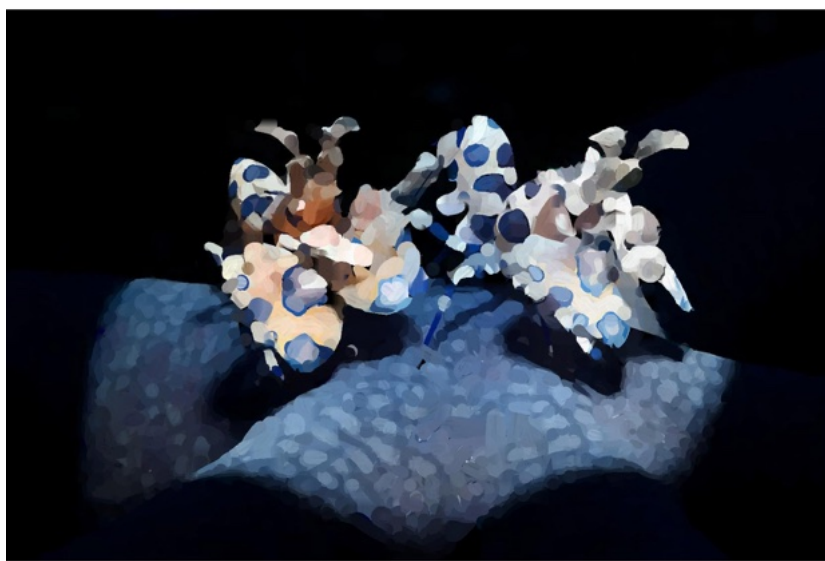
(b) Painterly



(c) Realism



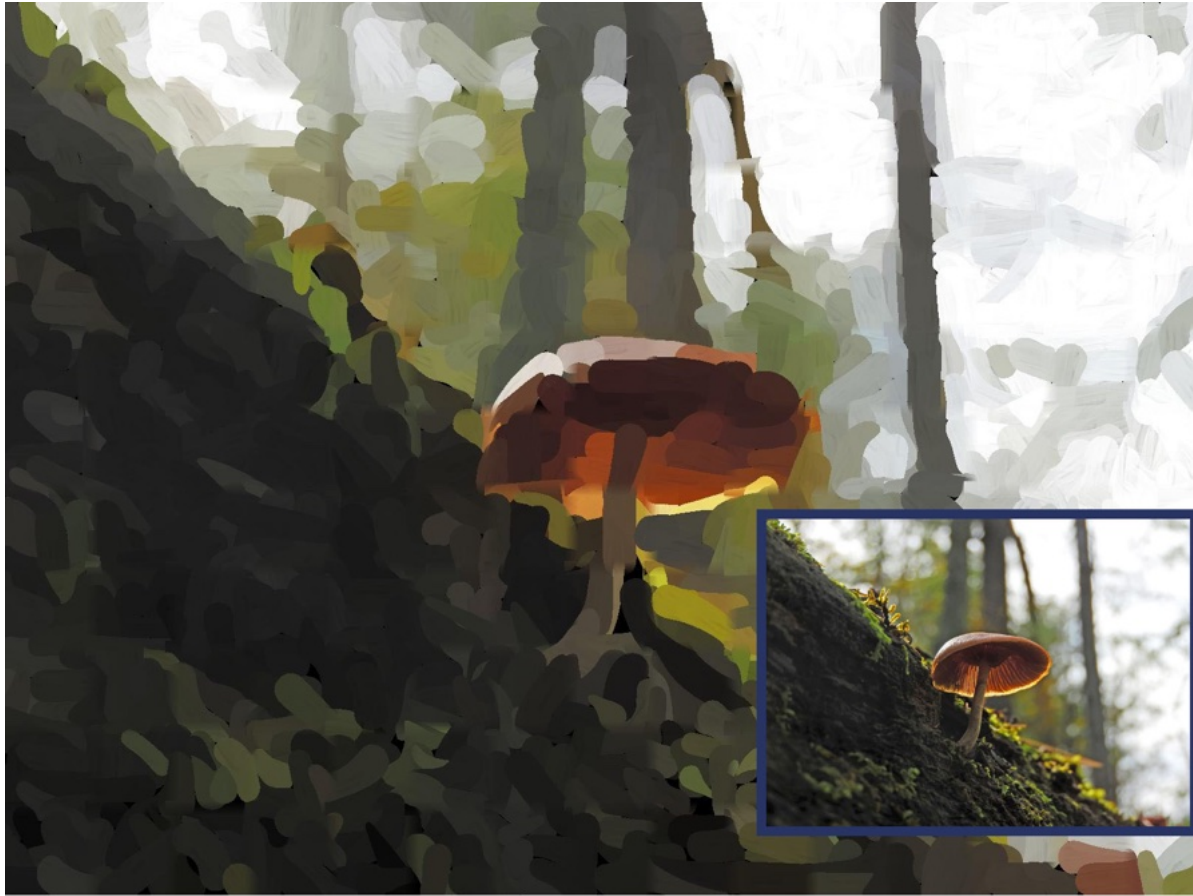
(a) Abstract



(b) Painterly



(c) Realism



(a) Abstract



(b) Painterly



Coarse-to-fine (from left to right and top to bottom)



Figure 20. Painterly style with dynamic attention maps corresponding to the last 3 painting layers are shown at the bottom.



Figure 21. Painterly style with dynamic attention maps corresponding to the last 3 painting layers are shown at the left. Painting process is shown at the right.



Figure 22. Realistic Style (left) and painterly style (right). Dynamic attention maps for the last painting pass are shown as an inset.



Figure 23. Realistic painting process. Painted using uniform dynamic attention maps. From coarse (top left) to fine (bottom right)



Figure 24. Pairings shown in user study from section 1. IDs from top to bottom: A5, A3, A4, and A6



Figure 25. Pairings shown in user study from section 1 and 2. IDs from top to bottom: A1, A2, B1, and B5

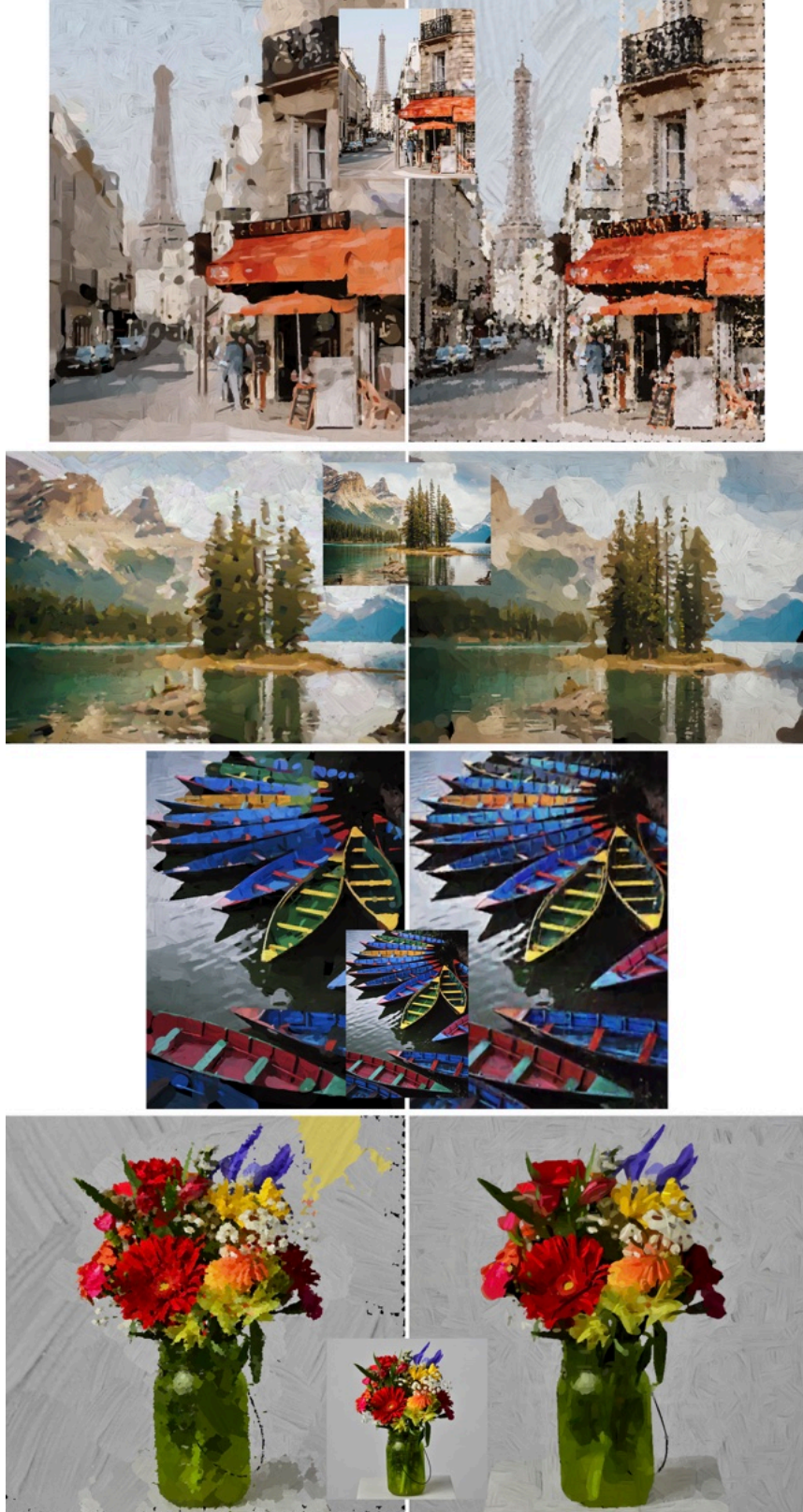


Figure 26. Pairings shown in user study from section 2. IDs from top to bottom: B3, B6, B2 and B4