## 1 Image similarity metrics

<sup>2</sup> An interface for a metric will take two images, x and  $\hat{x}$ , to compute a similarity metric,  $\sin(x, \hat{x})$  that

3 is symmetric. We take the view that unsuccessful rendering should be counted as absolute failures.

4 As such, our image similarity metrics computed are conditioned on the successful rendering of the

5 code into rendered images (i.e., if the rendering fails, the score will be 0 by default). In Image2Struct,

<sup>6</sup> we consider normalize metrics within the unit range so that they can be interpreted easily; a score

7 of zero implies complete dissimilarity whereas a score of 1 implies that the images are identical.

8 Without loss of generality, we assume both x and  $\hat{x}$  are of dimensions  $W \times H$ .

Earth Mover Similarity (EMS) We introduce an image metric named Earth Mover's Similarity
(EMS). It is inspired by the Earth Mover's distance (EMD)[1], which is a measure of dissimilarity
between two frequency distributions.

To compute the EMD between two images x and  $\hat{x}$ , we first transforms the images into signatures  $\mathcal{S}(x)$  and  $\mathcal{S}(\hat{x})$ , which are discrete distributions of features of Q elements.

$$\mathcal{S}(x) = \{(g_k, w_k^g) : 0 \le k < Q\} \qquad \& \qquad \mathcal{S}(\hat{x}) = \{(h_l, w_l^h) : 0 \le l < Q\}$$
(1)

We define a cost matrix  $C \in \mathbb{R}^{Q \times Q}$  wherein each element C[k, l] represents the cost of moving probability mass between  $g_k$  and  $h_l$ . We further denote the movement of probability mass between  $g_k$ and  $h_l$  by  $f_{k,l}$ . The optimal flow is the set of  $\{f_{k,l}^*\}$  that satisfies the following optimization problem:

$$\min \sum_{k} \sum_{l} f_{k,l} C[k,l] \qquad \qquad \text{subject to} \qquad (2a)$$

$$f_{k,l} \ge 0 \qquad \qquad \forall k, \forall l \qquad (2b)$$

$$\sum_{l} f_{k,l} \le w_k^g \qquad \qquad \forall k \qquad (2c)$$

$$\sum_{k} f_{k,l} \le w_k^h \qquad \qquad \forall l \qquad (2d)$$

$$\sum_{l} \sum_{k} f_{k,l} = \min\{\sum_{k} w_k^g, \sum_{l} w_l^h\}$$
(2e)

<sup>18</sup> The EMD can then be computed with Equation (3).

$$\operatorname{EMD}(x,\hat{x}) = \frac{\sum_{k} \sum_{l} f_{k,l}^{*} C[k,l]}{\sum_{k} \sum_{l} f_{k,l}^{*}}$$
(3)

<sup>19</sup> The signature is typically defined as the distribution of the grayscale values of an image when one <sup>20</sup> wants to compare images. In other words, S(x) is the probability mass function where the random <sup>21</sup> variable (i.e.,  $g_k$  or  $h_l$ ) is one of the possible pixel values (0 to 255) and the mass (i.e.,  $w_k^g$  or  $w_l^h$ ) is the <sup>22</sup> normalized count of the number of pixels in x with that value. In this formulation, spatial information <sup>23</sup> is lost and the metric is invariant to translation, reflection, and other pixel rearrangements. We solve <sup>24</sup> this problem by defining multidimensional signatures that consider the pixels' x- and y-coordinates <sup>25</sup> in addition to their values.

The support of our multidimensional signature,  $S_p$ , is all the possible combinations of the xcoordinates (x-pos), y-coordinates (y-pos), and the N possible pixel values (pix):

$$S_p(x) = \{ ((x \text{-pos}, y \text{-pos}, \text{pix})_k, w_k) : k \in \{0, 1, \cdots, WHN\} \}$$
(4)

The probability mass,  $w_k$ , takes the value of either  $\frac{1}{HW}$  or 0. The complexity of computing the cost matrix over  $S_p$  is  $O(W^2H^2)$ , making it difficult to compute for high resolution images. We therefore compute an approximated patch version of it, which we denoted as EMD<sub>block</sub>.

In EMD<sub>block</sub>, we first split two images, x and  $\hat{x}$ , each into K patches of dimensions  $r \times s$ :  $P_x^0, \dots, P_x^{K-1}$ . Our implementation sets r and s individually for every image such that there are  $8 \times 8$  patches in every image. To compare two patches  $P_x^t$  and  $P_{\hat{x}}^u$ , we treat each patch as separate images and compute the EMD using the multidimensional signature defined in Equation (4), which we will denote as  $\text{EMD}_p(P_x^t, P_{\hat{x}}^u)$ . Next, we define a separate cost matrix,  $C_p$ , such that the cost of moving one patch to another is the sum of the EMD between the patches and the Euclidean distance between them:

$$C_p[i,j] = EMD(P_x^i, P_{\hat{x}}^j) + ||(x \text{-} \text{pos}_i, y \text{-} \text{pos}_i), (x \text{-} \text{pos}_j, y \text{-} \text{pos}_j)||_2$$
(5)

EMD<sub>block</sub> attempts to minimize the cost of moving patches by solving the optimization problem 38 defined in Equation (2), but with the new cost function,  $C_p$ . By considering both the positions 39 and weights of the pixels within patches (through the multidimensional signature) and the distance 40 between patches, EMD<sub>block</sub> heavily penalizes random shuffling of pixels and assigns a lower score 41 (implying greater similarity) to a rendered image that contain blocks of similar but translated elements 42 as the input (see illustration in Figure 1). This property is useful for discerning between pairs 43 of images that contain similar elements -even if the elements are translated- and pairs where 44 distribution of colors in the rendered image is similar to the input image. 45

<sup>46</sup> Finally, we define the Earth Mover Similarity (EMS), a normalized similarity version of EMD<sub>block</sub>. It

47 compares  $\text{EMD}_{\text{block}}(x, \hat{x})$  against  $\text{EMD}_{\text{block}}(x, c(x))$ , the EMD between the reference image and a

constant black or white image (c(x)), whichever is the most dissimilar to the reference image x. An EMS of 0 indicates the least similarity and an EMS of 1 indicates the identity.

$$\mathrm{EMS}(x,\hat{x}) = 1 - \frac{\mathrm{EMD}_{\mathrm{block}}(x,\hat{x})}{\mathrm{EMD}_{\mathrm{block}}(x,c(x))}$$
(6)



Figure 1: An illustration of the two scales at which  $\text{EMD}_{\text{block}}$  operates. The left image is an altered copy of the right one in that 4 patches are manipulated.  $\text{EMD}_{\text{block}}$  computes an optimal flow where 3 of these patches (in red) are moved completely without modification. For the blue patch, it decides that it incurs a lower cost to move some pixels within the patch (the zoomed version on the right). On top of moving blocks or pixels,  $\text{EMD}_{\text{block}}$  can change the pixel colors at a cost (we do not illustrate color modification in this example for simplicity).

## 50 **References**

- 51 [1] Y. Rubner, C. Tomasi, and L.J. Guibas. A metric for distributions with applications to image
- databases. In Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271),
- <sup>53</sup> pages 59–66, 1998.