

## Supplementary Material

### A. Complementary definitions and details

**Ground-truth heatmap generation.** Note that we generate the ground-truth heatmap from ground-truth bounding boxes with the same method as CenterNet [29]. In particular, we apply a Gaussian kernel on the center of bounding boxes, where the kernel size is calculated according to the box size.

**Auxiliary task.** During training, our model will predict both faces and bodies, while in the evaluation, we only report the metrics regarding the faces.

**Center sampling.** Recall that in Equation (6) we define the average pooling positional encoding, now we introduce center sampling

$$cs(s_i)(u, v) = \phi((\bar{u}, \bar{v})), \quad (10)$$

where it first average the coordinate of an  $s_i \times s_i$  window and then encode it. When using center sampling, we replace it with  $\text{avgpool}(\mathbf{p}, s_i)(u, v)$  with it in (6). Since it discards the scale information, we adopt average pooling in our method.

**Evaluation details.** For the predicted bounding boxes, we apply a Non-Maximum-Suppression (NMS) IoU no more than 0.7, and retain the top-1000 predicted boxes based on confidence score. When evaluating the outpainting pipeline, we first apply the same NMS and top-1000 filter on the result of each image, then we aggregate all the remaining boxes and apply the NMS and filtering again. For the heatmap, we average the heatmaps over all images.

### B. Further discussion on the outpainting baseline

Tab. 7 reports the performance of the outpainting pipeline with varying numbers of samples. Increasing the number of samples improves metrics for outside faces, but degrades CE and AP on truncated faces, revealing a trade-off inherent to this approach. A further limitation is that the pipeline is not end-to-end trainable, making each component a potential bottleneck (Figure 7). Moreover, even with strong generative models, accessing the ideal conditional distribution remains an open challenge.

### C. Potential negative Societal impacts

We also note the potential for more troubling applications (dual use). Successfully detecting objects like humans faces beyond what is directly observable could serve opposing ends. Instead of directing the camera to avoid that area, extreme amodal face detection could be used to pursue

unseen-but-inferred objects. The existence of such applications does not negate the ethical case for extreme amodal face detection, though, which is based on its safety, privacy, and accessibility-enhancing potential.

Top- $\mu^2$	AP $\uparrow$	AP $_t$ $\uparrow$	AP $_o$ $\uparrow$	AP $_{o+}$ $\uparrow$	AP $_{o-}$ $\uparrow$	MAE $\downarrow$	MAE $_t$ $\downarrow$	MAE $_o$ $\downarrow$	MAE $_{o+}$ $\downarrow$	MAE $_{o-}$ $\downarrow$	mIoU $\uparrow$	Recall $\uparrow$	CE $\downarrow$	SE $\downarrow$
15	21.37	<u>62.51</u>	0.80	1.40	0.20	23.99	2.33	37.08	5.64	48.53	<u>18.08</u>	25.51	<b>93.27</b>	<b>88.34</b>
20	21.21	61.65	<b>0.99</b>	<b>1.74</b>	<b>0.24</b>	18.59	<u>2.15</u>	28.5	4.91	37.10	17.56	26.35	<u>93.64</u>	88.87
25	<b>21.49</b>	<b>62.73</b>	0.86	1.48	<b>0.24</b>	19.69	2.19	30.28	4.82	39.55	17.84	26.64	93.78	<u>88.69</u>
30	20.34	59.34	0.85	1.46	<u>0.23</u>	<b>16.31</b>	2.19	<b>24.94</b>	<u>4.51</u>	<b>32.38</b>	<b>18.09</b>	<b>26.85</b>	94.48	88.80
35	20.66	60.32	0.83	1.47	0.20	<u>17.29</u>	<b>2.07</b>	<u>26.53</u>	<b>4.33</b>	<u>34.61</u>	17.83	25.92	94.52	89.05
40	<u>21.38</u>	62.35	<u>0.89</u>	<u>1.56</u>	<u>0.23</u>	18.79	2.17	28.89	4.92	37.61	17.75	25.78	94.86	89.25

Table 5. Complete result of analysis on top- $\mu$  at scale  $\mathcal{S} = (2)$ .

Scale	AP $\uparrow$	AP $_t$ $\uparrow$	AP $_o$ $\uparrow$	AP $_{o+}$ $\uparrow$	AP $_{o-}$ $\uparrow$	MAE $\downarrow$	MAE $_t$ $\downarrow$	MAE $_o$ $\downarrow$	MAE $_{o+}$ $\downarrow$	MAE $_{o-}$ $\downarrow$	mIoU $\uparrow$	Recall $\uparrow$	CE $\downarrow$	SE $\downarrow$
(1)	<u>21.02</u>	<u>61.28</u>	<b>0.89</b>	<b>1.59</b>	0.19	20.31	2.25	31.32	5.14	40.85	<b>18.34</b>	25.25	<u>96.36</u>	89.57
(2)	19.06	55.74	0.72	1.27	0.17	21.28	<u>2.11</u>	32.81	5.32	42.82	17.60	23.67	96.62	90.34
(4)	19.93	58.62	0.58	1.03	0.12	21.87	2.22	33.77	5.15	44.19	17.41	23.09	96.60	90.41
(4, 2)	19.20	56.27	0.67	1.11	<u>0.22</u>	<b>13.64</b>	2.27	<b>20.68</b>	<u>4.82</u>	<u>39.55</u>	16.88	<u>25.34</u>	98.66	94.06
(2, 1)	<b>21.49</b>	<b>62.73</b>	<u>0.86</u>	<u>1.48</u>	<b>0.24</b>	19.69	2.19	30.28	<u>4.82</u>	<u>39.55</u>	17.84	<u>26.64</u>	<b>93.78</b>	<b>88.69</b>
(4, 2, 1)	19.96	58.30	0.79	1.42	0.16	<b>14.94</b>	<b>2.05</b>	<u>22.86</u>	<b>4.63</b>	<b>29.49</b>	<u>18.26</u>	23.91	98.25	93.06

Table 6. Complete result of analysis on multiple-scale.

Num of Samples	AP $\uparrow$	AP $_t$ $\uparrow$	AP $_o$ $\uparrow$	AP $_{o+}$ $\uparrow$	AP $_{o-}$ $\uparrow$	MAE $\downarrow$	MAE $_t$ $\downarrow$	MAE $_o$ $\downarrow$	MAE $_{o+}$ $\downarrow$	MAE $_{o-}$ $\downarrow$	mIoU $\uparrow$	Recall $\uparrow$	CE $\downarrow$	SE $\downarrow$
1	<b>9.07</b>	<b>24.01</b>	1.59	2.01	<b>1.17</b>	24.03	3.25	36.25	6.75	46.99	18.41	24.91	<b>93.68</b>	92.56
2	<u>7.75</u>	<u>20.13</u>	1.56	<u>2.23</u>	0.89	<u>14.16</u>	2.50	<u>20.95</u>	4.57	<u>26.92</u>	19.98	<b>25.27</b>	<u>95.43</u>	91.06
5	5.89	<u>15.02</u>	1.32	1.86	0.78	<b>13.75</b>	2.17	<b>20.45</b>	3.71	<b>26.54</b>	<u>20.47</u>	<u>25.15</u>	96.18	90.39
8	5.51	12.64	<b>1.94</b>	2.01	<b>1.17</b>	<u>14.16</u>	<u>2.08</u>	21.11	<u>3.58</u>	27.50	<b>20.53</b>	25.07	96.35	<u>90.24</u>
10	4.93	11.54	<u>1.62</u>	<b>2.47</b>	0.76	14.69	<u>2.07</u>	21.94	<b>3.48</b>	28.67	<b>20.53</b>	25.03	96.41	<b>90.18</b>

Table 7. Analysis of the number of outpainting samples.

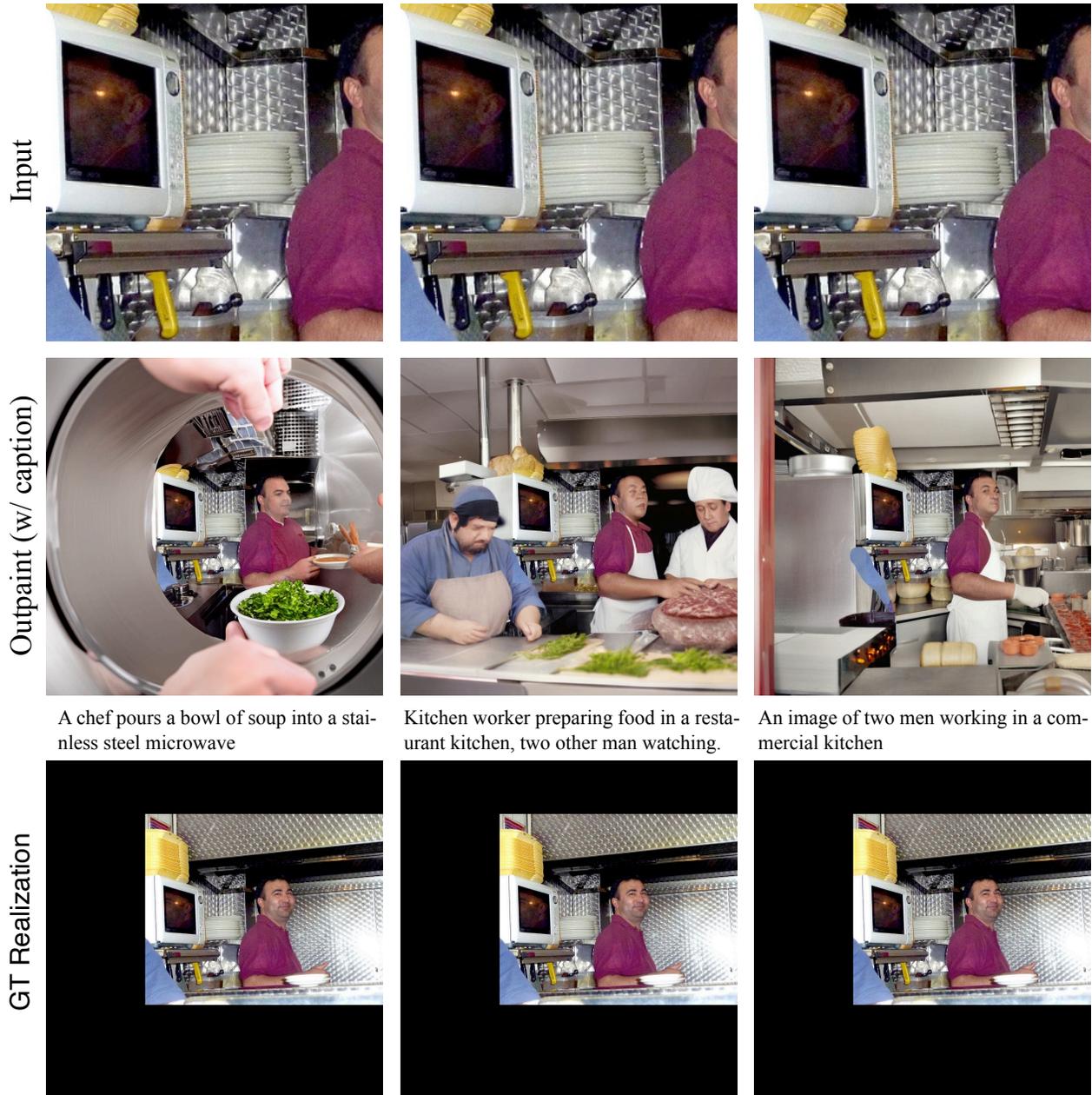
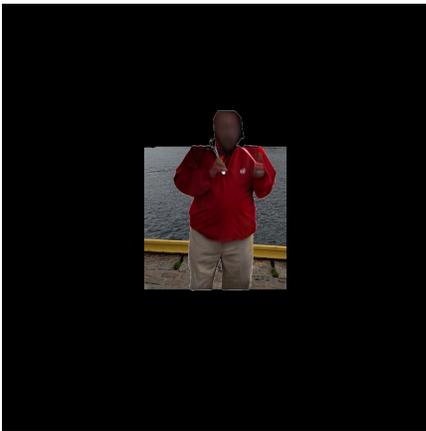


Figure 7. Outpainted example from SDXL [17] + BLIP2 [12]. These three examples show that the outpainted example can be bottlenecked by any one component, and the randomness of the outpainted result. The middle example demonstrates that when both components collaborate well, the left and right example shows the bottleneck made by either VLM or the outpainting model.

Input



Pix2Gestalt



GT Realization

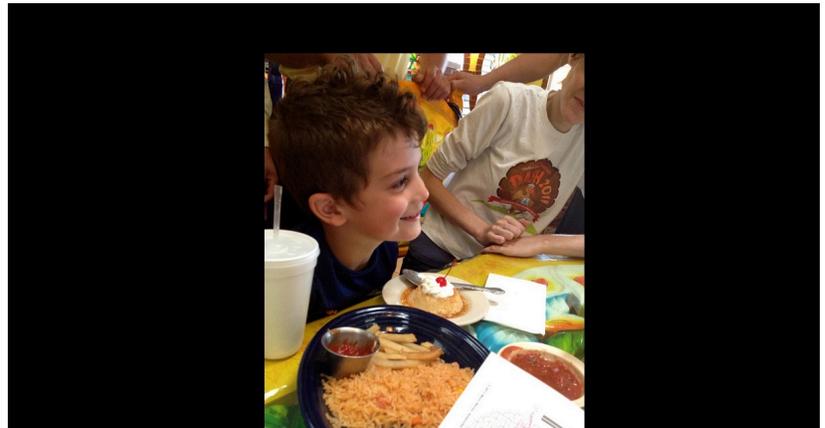
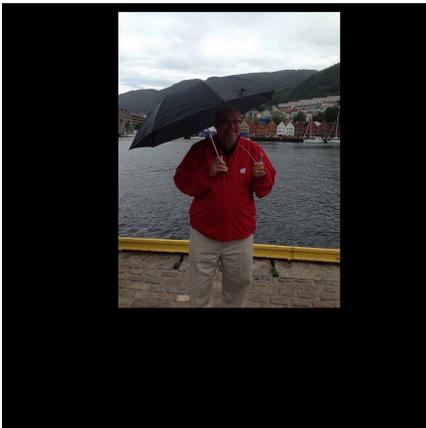


Figure 8. Completion examples with Pix2Gestalt [15]. The first example shows that the model struggles to complete out-of-frame regions despite strong visual evidence, while the second demonstrates effective in-frame occluder removal. Together, these cases highlight the distinction between in-frame completion and out-of-frame completion: strong performance on the former does not necessarily transfer to the latter.