

Generating Multiple Objects at Spatially Distinct Locations

Tobias Hinz, Stefan Heinrich, Stefan Wermter

We extend the GAN architecture in order to control object locations and identities in generated images.

INTRO

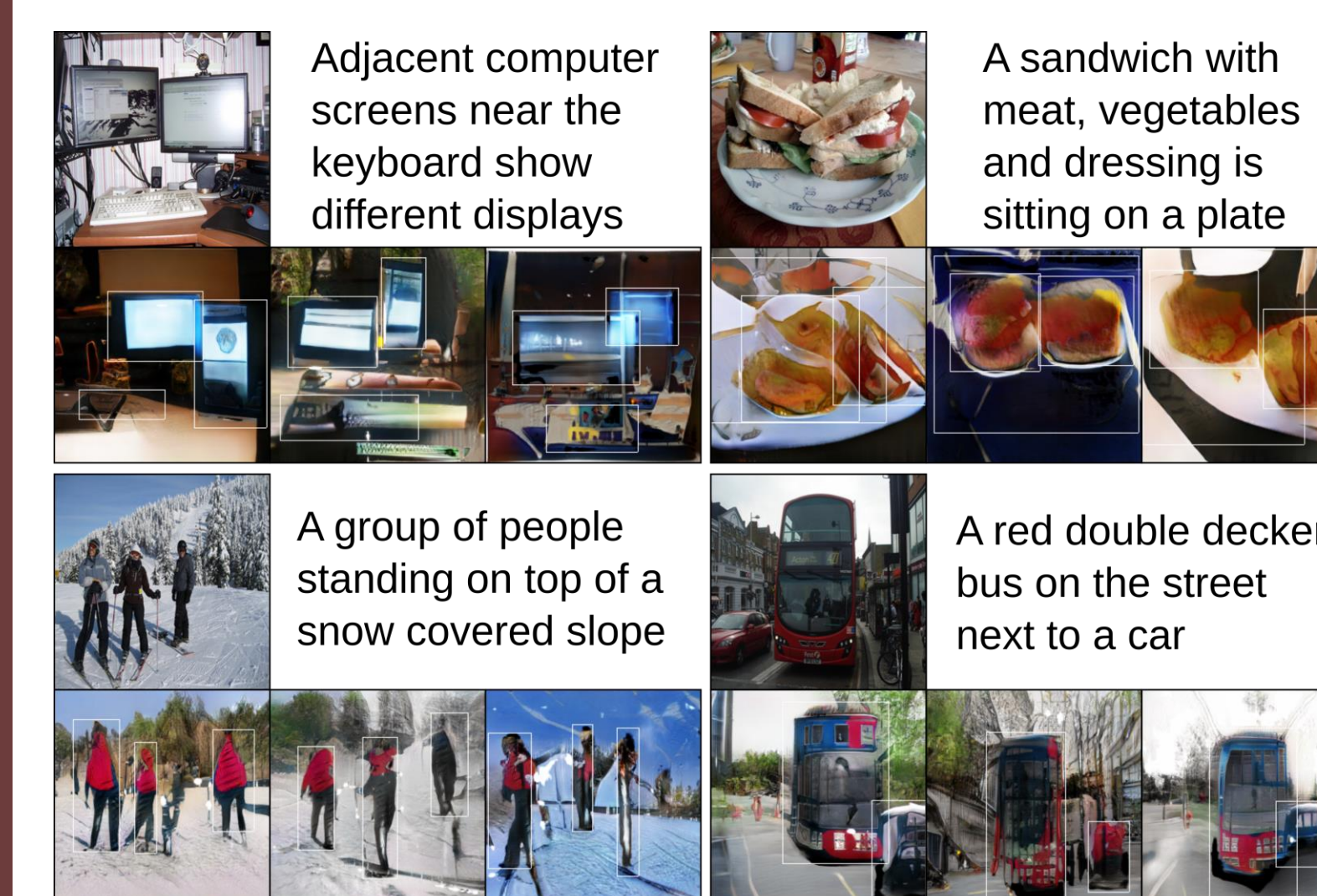
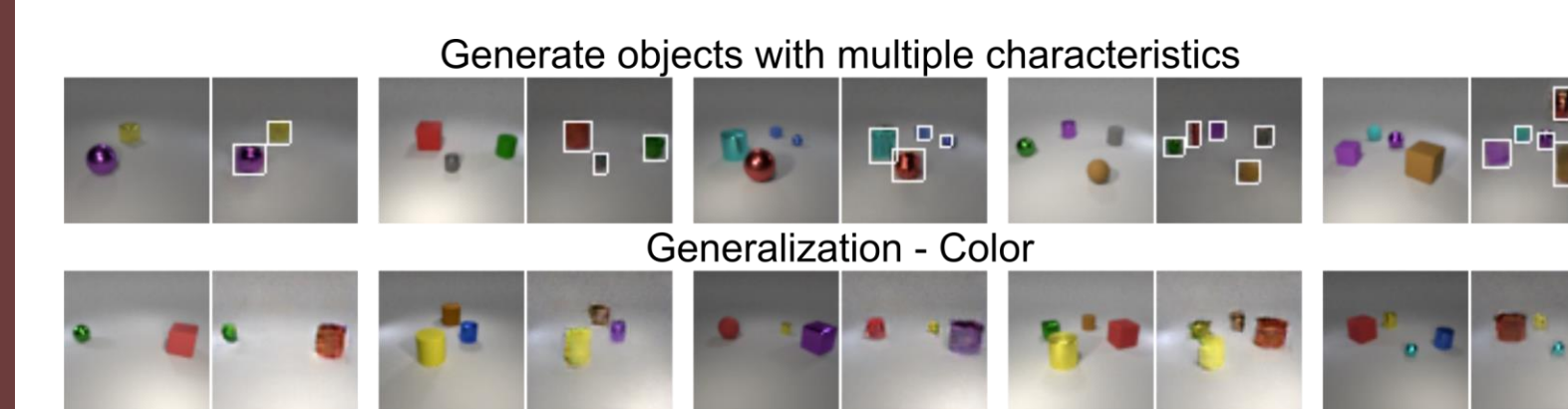
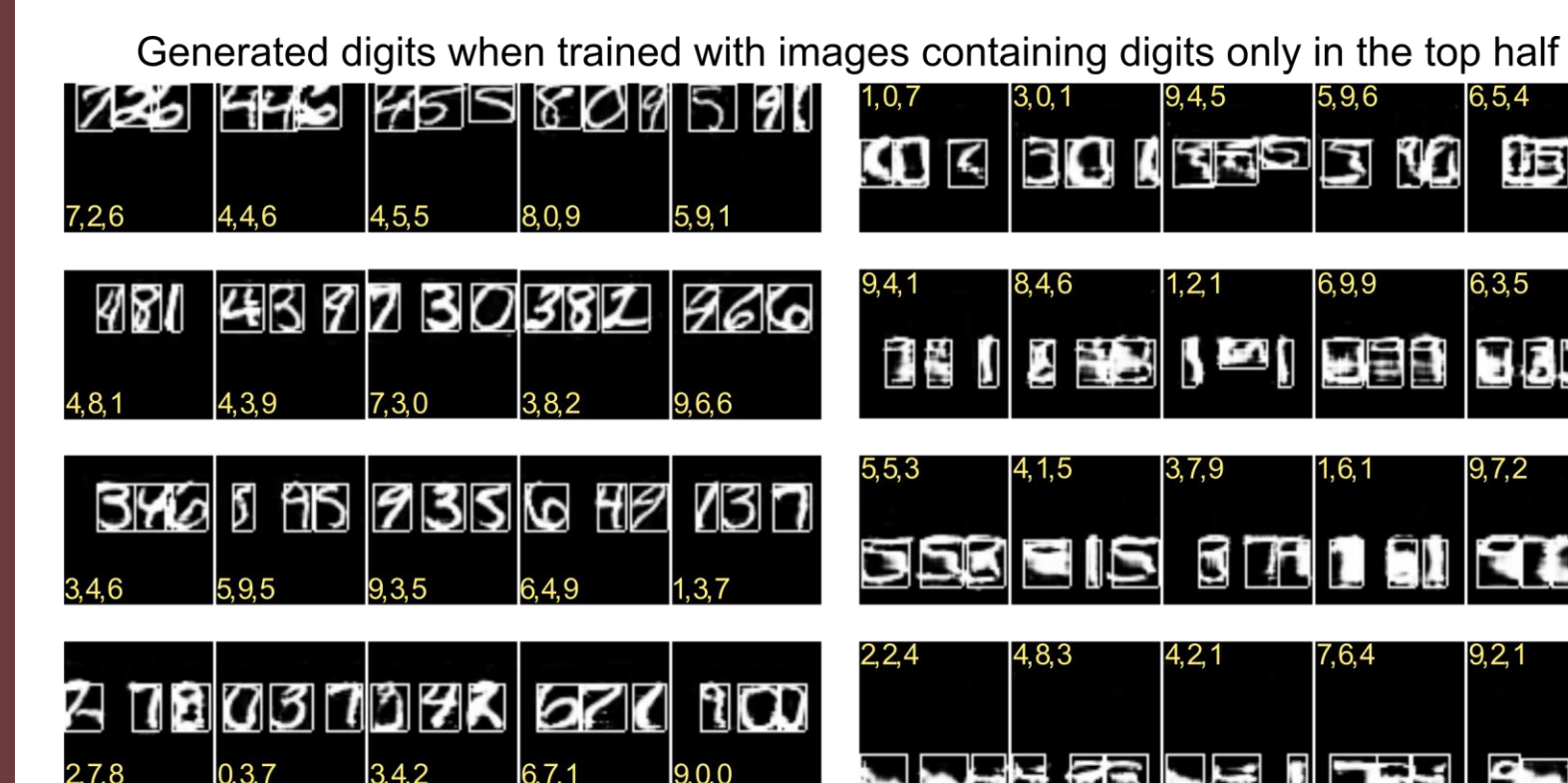
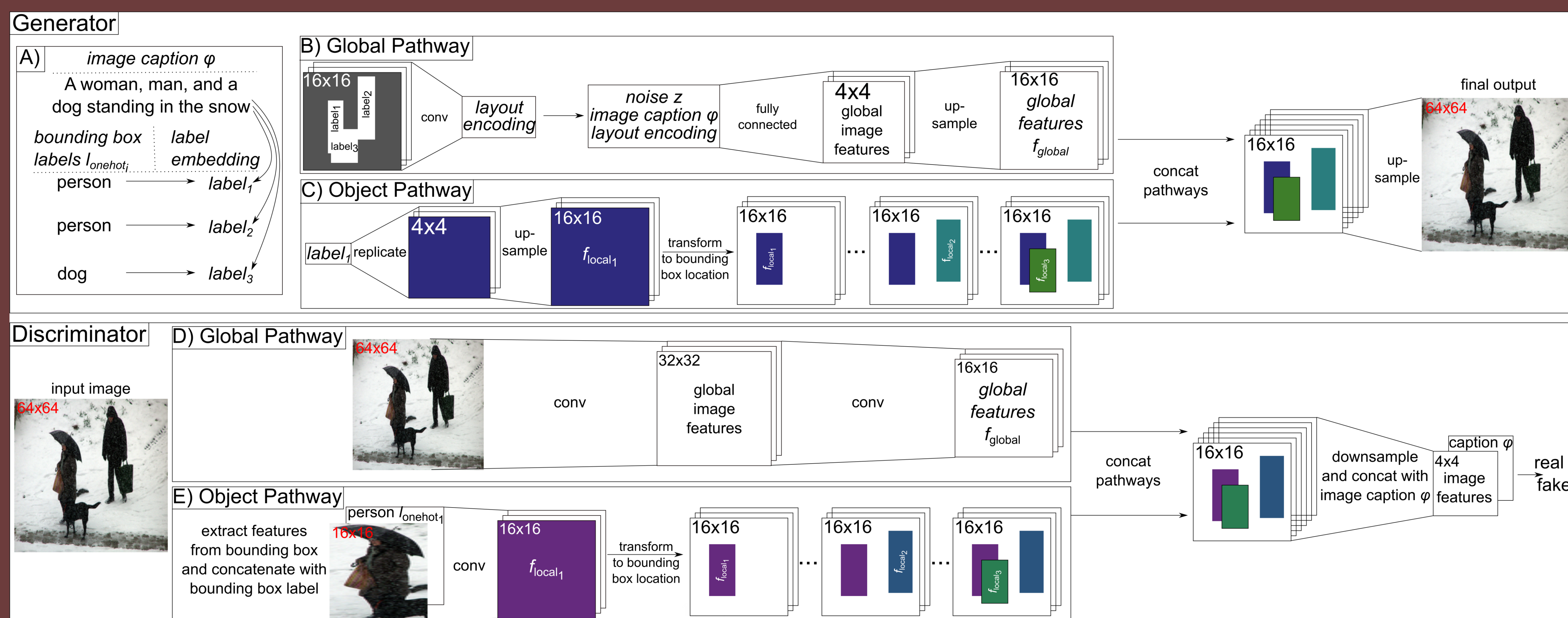
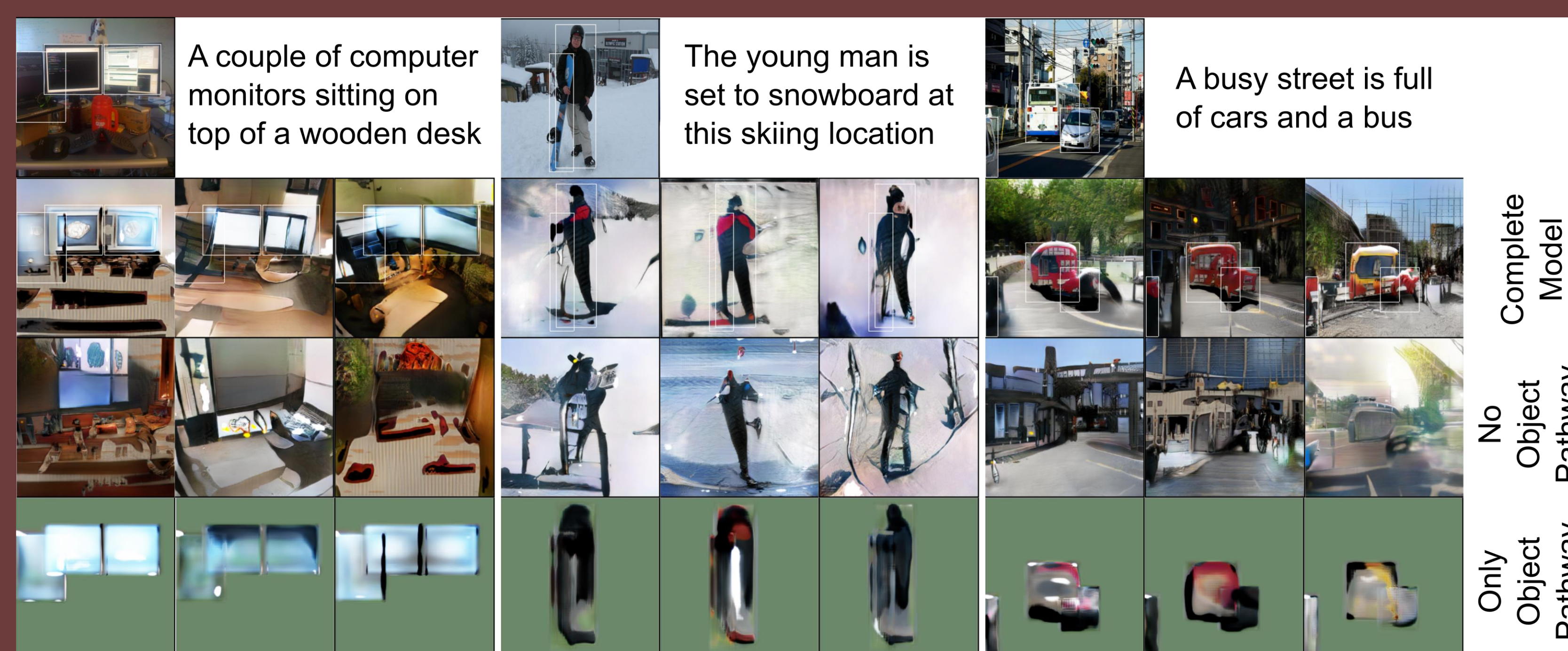
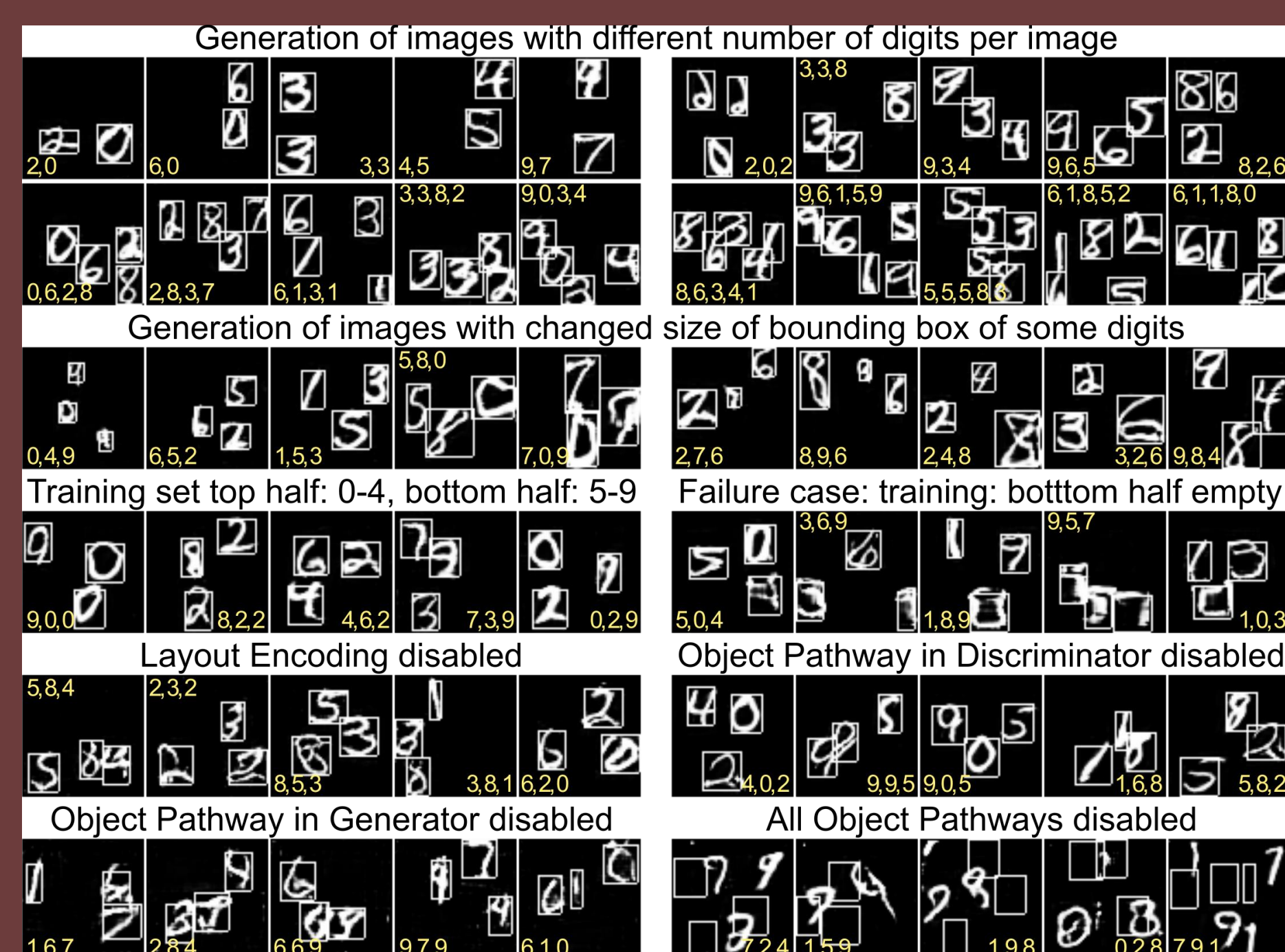
- one challenge with Generative Adversarial Nets is that it is hard to control the layout of the generated scenes
- a popular approach for increased control is to use scene layouts as additional input, however, this requires a lot of labeling
- our approach only needs a bounding box (thereby specifying location and shape) and a class label for each of the foreground objects we want to be in the image

METHODS

- we add an object pathway to the generator and the discriminator of a normal GAN
- the generator applies the object pathway iteratively at each bounding box location, conditioned by the label
- the discriminator's object pathway checks each bounding box location iteratively and evaluates whether the specified object is actually recognizable

RESULTS

- experiments show that the architecture does lead to increased control over the image generation process
- experiments on the Multi-MNIST and CLEVR datasets show that the architecture can generalize to novel object characteristics and locations, as well as to different numbers of generated objects per scene
- experiments on the MS-COCO dataset show that the object pathway learns features for the individual foreground objects and can lead to an overall higher quality of the generated images (based on IS and FID)



Model	Resolution	IS \uparrow	FID \downarrow
GAN-INT-CLS	64 x 64	7.88 \pm 0.07	60.62
StackGAN-V2	256 x 256	8.30 \pm 0.10	81.59
StackGAN	256 x 256	8.45 \pm 0.03 ¹	74.05
PPGN	227 x 227	9.58 \pm 0.21	
ChatPainter (StackGAN)	256 x 256	9.74 \pm 0.02	
Semantic Layout	128 x 128	11.46 \pm 0.09 ²	
HDgan	256 x 256	11.86 \pm 0.18	71.27 \pm 0.12 ³
AttnGAN	256 x 256	23.61 \pm 0.21 ⁴	33.10 \pm 0.11 ⁵
StackGAN + Object Pathways (Ours) ⁵	256 x 256	12.12 \pm 0.31	55.30 \pm 1.78
AttnGAN + Object Pathways (Ours)	256 x 256	24.76 \pm 0.43	33.35 \pm 1.15

¹ Recently updated to 10.62 \pm 0.19 in its source code.
² When using the ground truth bounding boxes at test time (as we do) the IS increases to 11.94 \pm 0.09.
³ FID score was calculated with samples generated with the pretrained model provided by the authors.
⁴ The authors report a "best" value of 25.89 \pm 0.47, but when calculating the IS with the pretrained model provided by the authors we only obtain an IS of 23.61. Other researchers on the authors' Github website report a similar value for the pretrained model.
⁵ We use the updated source code (IS of 10.62) as our baseline model.

Table 1: Comparison of the Inception Score (IS) and Fréchet Inception Distance (FID) on the MS-COCO data set for different models. Note: the IS and FID values of our models are not necessarily directly comparable to the other models, since our model gets, in addition to the image caption, up to three bounding boxes and their respective object labels as input at test time.

