1 We thank all reviewers for their positive and constructive comments, such as the application important, the results impressive,
2 and the method generalizable. Below we first address the common questions, and then questions by individual reviewers.
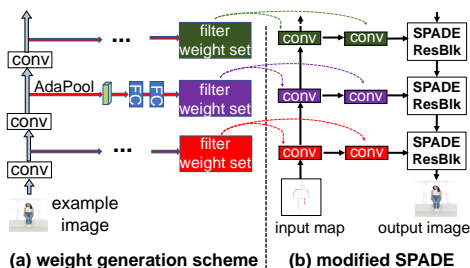


**(a) weight generation scheme** | **(b) modified SPADE**
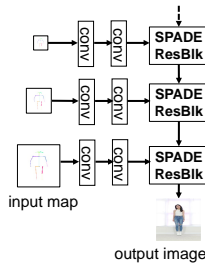**Fig A: Our network architecture**
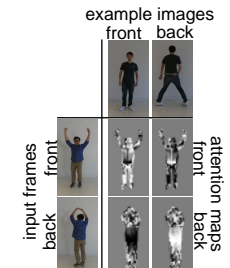
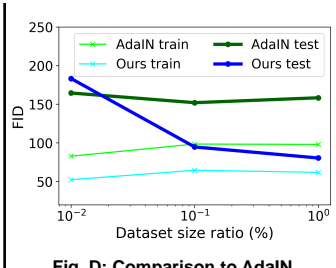**Fig B: Original SPADE**

**Fig. C: Example attn maps**

**Fig. D: Comparison to AdaIN with varying dataset sizes.**

3 **Details of network architecture.** We are sorry there were not enough
4 details provided in the paper due to the page limit. In Fig. A, we show
5 the network architecture when only one example image is used (i.e.
6 without the attention mechanism). Compared to the original SPADE
7 generator (Fig. B), we made two major modifications. First, instead of

| | AdaIN (avg) | AdaIN (attn) | Ours (attn) |
|---|---|---|---|
| FID | 129.90 | 113.83 | *76.53* |

**Table A: Comparison to AdaIN** (5 example images).

| | Ours (1 exp) | Ours (50 exps) | Ours (50 exps + finetune) | vid2vid |
|---|---|---|---|---|
| FID | 56.99 | 51.10 | *43.04* | 47.19 |

**Table B. Comparison to vid2vid**

8 generating the affine parameters at each layer independently, we reuse features from the previous layer when generating
9 parameters in the next layer (Fig. A(b)). Second, the weights in SPADE are determined on the fly (Fig. A(a)). After we
10 encode the example image, we apply adaptive pooling (AdaPool) to make it fixed spatial size. The results are fed to two
11 fully connected layers to generate the weights, which are used in corresponding convolutions (denoted by different colors).

12 **Attention mechanism.** We combine features from multiple example images before feeding them to the AdaPool layer in
13 Fig. A(a). The attention maps are soft spatial maps of size $\mathbb{R}^{K \times N \times N}$, where $K$ is the number of examples and $N$ is the
14 spatial dimension. The map determines that for each spatial location in the output $(x, y)$, which spatial location in which
15 example image $(k, x_k, y_k)$ carries most relevant information. For example, when example images include both front and
16 back of the target person, the attention maps can help capture corresponding body parts during synthesis (Fig. C).

17 **Comparison to AdaIN.** In AdaIN, information from the example image is represented as a scaling vector and a biased
18 vector. This operation could be considered as a 1x1 convolution with a group size equal to the channel size. From this
19 perspective, AdaIN is a constrained case of the proposed weight generation scheme, since our scheme can generate a
20 convolutional kernel with group size equal to 1 and kernel size larger than 1x1. Moreover, the proposed scheme can be
21 easily combined with the SPADE module. Specifically, we use the proposed generation scheme to generate weights for the
22 SPADE layers, which in turn generate spatially adaptive de-modulation parameters. To justify the importance of weight
23 generation, we compare with AdaIN both using weighted average and our attention module (Table A). We also compare
24 with AdaIN when different dataset sizes are used. The assumption is that when the dataset is small, both methods are able
25 to catch the diversity in the dataset. However, as the dataset size grows larger, AdaIN starts to fail since the expressibility is
26 limited, as shown in Fig. D.

27 **Limitations.** Although our network can, in principal, generalize to unseen domains, when the test domain is too different
28 from the training domains it will not perform well. For example, when testing on CG characters which look very different
29 from real-world people, the network struggles. We will include several failure examples in the revised version.

30 **R1: How many examples (the K parameter) are needed?** Although we demonstrated that more example images are
31 helpful (Fig. 7b in the paper), our method can work well when $K = 1$. We note that comparisons and examples shown in
32 the paper are using $K = 1$ (i.e., without the attention mechanism).

33 **R1: Comparison to the baselines and importance of weight generation.** As explained above, for all comparisons our
34 method only uses a single example image, which is the same as the baselines, so the comparisons are fair. This shows that
35 weight generation ($E_c$) is useful for generating good quality results. We will make this clear in the revised version.

36 **R2: Novelty.** To our best knowledge, we are the first to show that weight generation can be used to solve the adaptive
37 video-to-video synthesis problem. As recognized by R3, this will likely have a wide impact since the topic is of wide
38 interest and our method is generalizable. We will also fix the notations in the revised version.

39 **R3: Comparisons to Zhou et al. 2019 and [7].** We will cite and discuss these papers. Note that both these methods still
40 require *different* models for different persons, which has the same drawback as `vid2vid`. For example, they typically need
41 minutes of training data and days of training time, while our method only needs one image and negligible time for weight
42 generation. Moreover, since code of these methods are not released, we show comparisons to `vid2vid` in Table B for a
43 specific person. We find that our model renders comparable results even when $K = 1$. Moreover, if we further finetune our
44 model based on the example images, we can achieve comparable or even better performance.

45 **R3: Dataset.** The dataset is collected by the authors. Different videos are randomly divided into training/test sets. Each
46 video is further divided into clips to ensure each clip only contains one person and does not contain any scene transition,
47 which in general results in 30-1000 frames per clip. We will release the dataset for facilitating research in the field.

48 **R3: Human evaluation.** For each pair of comparisons, we generate 100 clips, each of them viewed by 60 workers. Orders
49 are randomized. Hence, each user preference score is computed based on the 6000 evaluations. As reported in the main
50 paper, our preference scores are significantly higher than $0.5$. We will also add the citation to the segmentation network.