

StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation.

20170264 Seungmin Park, 20190006 Sana Kang,
20190545 Seungjae Lim, 20200673 Hyunmin Choi

Abstract

Research before StarGAN had excellent results in image-to-image translation, but these approaches have limitations, which leads to independently building different models for every image domain pair, which requires processing three or more multi-domains. When doing so, there was a problem with limited scalability and robustness. To solve this problem, a stargan that can perform image-to-image translation for multiple domains using a single model was presented, and in this project, we want to reproduce and improve this stargan.

The unified model architecture of StarGAN enables simultaneous training of data sets from different domains within a single network, showing superior translation quality compared to the existing model and scalability and robustness by flexibly executing changes between domains. In this project, we reproduce and play StarGAN using pytorch. Then, after training using the CelebA dataset and RaFD, we would like to see the transformation result.

Based on this project, we found a section where vanilla StarGAN's performance decreased, analyzed the cause, and tried to come up with a solution. Our model is struggling for some specific class generation. Therefore, it can be seen that the current model has a low performance in finding geometric or structural patterns in faces. Therefore, we proposed self-attention as a solution to this.

1 Introduction

Image-to-Image translation, a topic we will cover in the paper we will reproduce, is the problem of changing one specific feature of a given image into another. This has been greatly advanced with the introduction of the Generative Adversarial Network (GAN). In particular, the problem solved in this paper is a multi-domain problem that converts two or more domains at the same time. Although multi-domain translation was possible in previous studies, there were several problems. First, the model needs a $k(k-1)$ generator to learn the mapping to the k domain, which is expensive to compute. Second, post-training models are inefficient because each generator does not share global features, so it cannot fully utilize all the training data. This is because, if the training data is not fully utilized, the translation quality will be lowered. Third, since each dataset is labeled separately, it is difficult to combine it into the learning domains of other datasets. Therefore, there were inefficiencies due to the lack of scalability for training data or for adding or changing features. So, we build a novel generative adversarial network that learns mappings between multiple domains using a single generator and discriminator to reproduce the paper that solved these problems. Successfully learn multi-domain image transformations between multiple data sets by utilizing the mask vector method. In this reproduction, we plan to reproduce it using CelebA and RaFD. After training the reproduced model with CelebA and RaFD data, we plan to show that it is superior to the basic model through qualitative and quantitative comparison of the transformation results of other trained models. Therefore, in this paper reproduction, it is meaningful to provide a model solution using one

generator for multiple domains, and to show ones. However, the limitation of Pix2pix is that, simultaneous learning and scalability between multiple domains. This is also meant to secure the problems of previous papers and to understand why the limitations of the previous papers were overcome.

127

2 Theoretical Background

Since the concept of Generative Adversarial Network (GAN) was first introduced to academia, various variations of GAN have emerged. To figure out common limitations several GAN models had shared, we conducted a preliminary study on the following models: Pix2pix, CycleGAN, and WGAN.

91

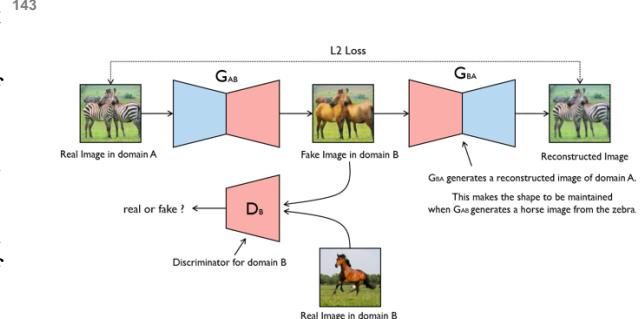
1.1 Pix2pix

Pix2Pix refers to an image-to-image translation GAN that learns a mapping from an image X with a random noise Z to the output image Y. It was designed in the early stages of i2i translation right after the advent of cGAN. To learn this model, we need to prepare the paired data set that consists of two different domains – X (the edge information) and Y (the original data). This edge information X is considered as the condition itself as in the conditional GAN. Therefore, for the Pix2pix model, we need to make the edge information of each data on our own.

Especially, Pix2Pix uses both generator and discriminator with the module of the form convolution BatchNorm ReLu. In this network, the input goes through a series of multiple layers in two main directions: towards a bottleneck layer and backward from it. This mechanism is to deal with a great amount of low-level information of the input and output across the network during image translation. Details of the learning process of Pix2pix are as follows. Firstly, Pix2pix uses Generator G to create a fake image $G(x)$ that meets condition x (edge information). After that, the fake image $G(x)$ and condition x are put together in Discriminator D. Based on it, the model carries out the learning process to make D figure out generated image $G(x)$ as a fake one. Conversely, for the G, the model tries to make it generate fake images which can deceive D as real

1.2 CycleGAN

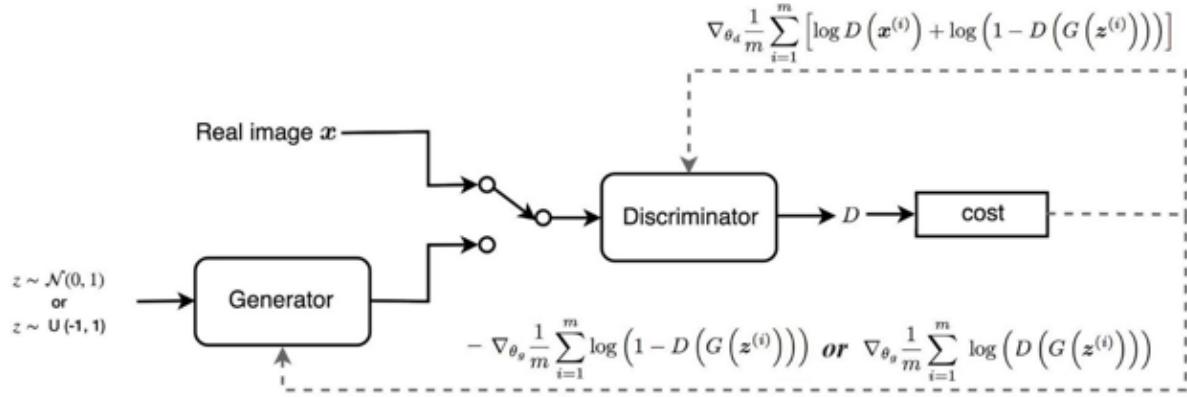
The Cycle Generative Adversarial Network, or CycleGAN, aims to train a deep convolutional neural network (CNN) for i2i translation tasks. In the learning process of this model, the image $G(x)$ made by the generator becomes reconstructed back into the original image x. This whole process with the cycle-consistency loss enables the cycle GAN to change domain-related features effectively while preserving the content of the original image. For this purpose, cycleGAN uses two translators, G and F. Translator G is responsible for the conversion from X to Y; while Translator F is responsible for the conversion from Y to X in reverse.



143
144 Figure 1: The process of CycleGAN.

However, there exist some cases where CycleGAN falls in a vicious loop – repeatedly presenting only one image corresponding to a particular domain regardless of input. For example, suppose you created a new zebra picture by adding a new attribute to the horse picture. In this case, for the corresponding fake image, the discriminator does not find the image as a fake one because it is very likely to exist. In this case, the generator eventually continues to generate that exact image that can deceive the discriminator over and over.

1.3 WGAN



160 Figure 2: The process of WGAN.

195

161
162 WGAN (Wasserstein GAN) refers to a type of
163 GAN that minimizes an approximation of the
164 Earth-Mover's distance rather than the vanilla
165 GAN with the Jensen-Shannon divergence.
166 Likewise, it aims to solve the problems of the
167 existing vanilla GAN. Those problems were [1]
168 the difficulty of balancing between the
169 discriminator and generator and [2] the
170 occurrence of mode dropping even after learning
171 is completed. The reason for these problems is
172 that the vanilla GAN model was not able to learn
173 itself to the optimum level because of the lack of
174 ability of the discriminator. To solve it, WGAN
175 uses the newly defined critic instead of the
176 discriminator. While the discriminator uses a
177 sigmoid function to determine whether the
178 generated image is true or fake, the critic uses
179 scalar values from the Earth Mover distance, the
180 measure to calculate the distance between
181 probability distributions. However, there exists a
182 weight clipping issue in the WGAN, so WGAN-
183 gp introduces the gradient penalty to solve it.

184

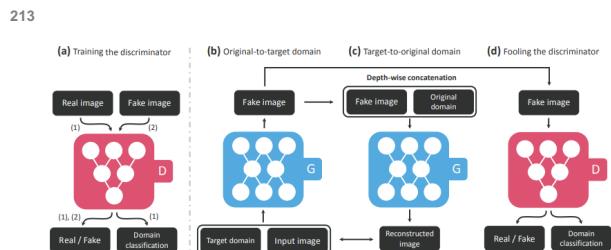
185 3 Challenges

186 After conducting the preliminary study, we
187 summarized the following two issues as
188 limitations that the previous GAN models had in
189 common. The first limitation is that it is difficult
190 to maintain the balance between the discriminator
191 and generator when training a model. The second
192 limitation is that no model has proposed an
193 efficient method for efficient image-to-image
194 translation in multiple domains.

196 4 StarGAN Design

197 So far, we describe about I2I-translation related
198 works and limitations of them. From now, we will
199 discuss about how StarGAN (Star Generative
200 Adversarial Networks) architecture and loss
201 functions are composed. Also, we will introduce
202 how StarGAN manipulates multiple datasets
203 containing different labels to perform image
204 translation using them.

205 Like other image translation networks, StarGAN
206 consists two modules, a discriminator D and a
207 generator G . Main purpose of discriminator is to
208 distinguish between real and fake image and
209 classify the image which is judged to be real into
210 its corresponding domain. Generator takes two
211 inputs, image and target domain label, and
212 generates a fake image.



213
214 Figure 3: Overview of StarGAN training process.

215 4.1 Network architecture: Loss functions and 216 Implementation

217 As mentioned before, we train G to translate an
218 input image x into an output image y on target
219 domain label c , $G(x, c) \rightarrow y$. In training phase, we
220 will randomly generate domain label c so that G
221 learns flexible translation. In case of discriminator
222 D , it will produce probability distribution over

224 source images and domain labels, $D : x \rightarrow$
 225 $\{D_{src}(x), D_{cls}(x)\}$.

$$L_G = L_{adv} + \lambda_{cls} L_{cls}^f + \lambda_{rec} L_{rec}$$

226
 227 **Adversarial Loss.** To make the generated images
 228 indistinguishable from real images, we adopt an
 229 adversarial loss something like this. G tries to
 230 minimize this objective, while D tries to
 231 maximize it.

$$232 \quad L_{adv} = E_x[\log D_{src}(x)] \\ 233 \quad + E_{x,c}[\log(1 - D_{src}(G(x, c)))]$$

234
 235 **Domain Classification Loss.** Given input x and
 236 target domain label c, starGAN should generate
 237 output image y which is properly classified into
 238 target domain. To achieve this, we need additional
 239 loss functions to optimize both D and G. In detail,
 240 domain classification loss for discriminator is
 241 defined below.

$$242 \quad L_{cls}^r = E_{x,c'}[-\log D_{cls}(c'|x)]$$

243 By minimizing this term, D could learn to classify
 244 a real image x to its original domain c_prime. On
 245 the other hand, the loss function for generator is
 246 defined below.

$$247 \quad L_{cls}^f = E_{x,c}[-\log D_{cls}(c|G(x, c))]$$

248 This term will help for distinguishing domain of
 249 fake images. By minimizing this term, G could
 250 learn to generate images that can be classified to
 251 target domain c.

252 **Reconstruction Loss.** To guarantee that
 253 translated images preserve the content of its
 254 original image while translation proceeds, we
 255 apply a cycle consistency loss to generator which
 256 is defined below.

$$257 \quad L_{rec} = E_{x,c,c'}[|x - G(G(x, c), c')|_1]$$

258 We adopt the L1 norm as our reconstruction loss.

259 Note that this is not for separate generators, but
 260 for single generator which translate in bi-direction.

261 **Full Objective.** Combining all loss functions we
 262 discuss about so far, the final object functions are
 263 defined as,

$$264 \quad L_D = -L_{adv} + \lambda_{cls} L_{cls}^r$$

265 where two lambdas are hyper-parameters that
 266 control the relative importance of domain
 267 classification and reconstruction losses. We use
 268 $\lambda_{cls} = 1$, $\lambda_{rec} = 10$ here.

269
 270
 271 **Improved GAN Training.** For higher quality
 272 image generation, we modified adversarial loss
 273 function to Wasserstein GAN objective with
 274 gradient penalty, which is defined below.

$$275 \quad L_{adv} = E_x[D_{src}(x)] - E_{x,c}[D_{src}(G(x, c))] \\ 276 \quad - \lambda_{gp} E_{x'}[(|\nabla_{x'} D_{src}(x')|_2 \\ 277 \quad - 1)^2]$$

278 Here, x' is sampled uniformly along a straight line
 279 between a pair of a real and a generated image.
 280 We use $\lambda_{gp} = 10$ for all experiments.

281
 282 **Network Architecture.** Adapted from
 283 CycleGAN, StarGAN has generator network
 284 composed of two convolutional layers with stride
 285 size of two for down-sampling, six residual blocks,
 286 and two transposed convolutional layers with
 287 stride size of two for up-sampling. For generator,
 288 instance normalization is applied.

289 4.2 Training with Multiple Datasets

290 Strength of StarGAN is that it can incorporate
 291 various kind of datasets containing different types
 292 of labels, so it could control all the labels at the
 293 test phase. But, label information is only partially
 294 known to each dataset.

295
 296 **Mask Vector.** To resolve this problem, we use a
 297 mask vector that allows StarGAN to ignore
 298 unspecified labels and focus on the explicitly
 299 known label.

300 **Training Strategy.** At training phase of StarGAN
 301 with multiple datasets, we use domain label
 302 concatenated with masked vector as input to the
 303 generator. By doing so, generator learns to ignore
 304 the unspecified labels, and focus on the explicitly
 305 given label. In discriminator's perspective,
 306 discriminator will try to minimize only the
 307 classification error associated to the known label.

308 Under multi-dataset training, discriminator will 357 latent vector z and the conditional vector c . In
309 learn all of the discriminative features for all 358 addition, IcGAN introduces an encoder to learn the
310 datasets. 359 inverse mapping of cGAN, $E_z : x \rightarrow z$ and $E_c : 360 x \rightarrow c$. This allows IcGAN to synthesis images by
361 only changing the conditional vector and
362 preserving the latent vector.

311 5 Experimental Results

312 In this section, we first compare our replicated
313 StarGAN against recent methods including original
314 StarGAN on facial attribute transfer by conducting
315 user-studies. Next, we prepare a new dataset named
316 AffectNet, to check that our model works well on
317 another dataset rather than RaFD and CelebA.
318 Lastly, we demonstrate empirical results that
319 StarGAN can learn image-to-image translation
320 from multiple datasets.

321

322 5.1 Baseline Models

323 As our baseline models, we adopt DIAT and
324 CycleGAN, both of which performs image-to-
325 image translation between two different domains.
326 We also adopt IcGAN as a baseline which can
327 perform attribute transfer using a cGAN.

328 For comparison, we have to implement all of
329 these models and train models multiple times for
330 every pair of two different domains. However, due
331 to lack of time to implement all of baseline models
332 and long training time, the output of baseline
333 models and original StarGAN is conducted from
334 original paper.

335

336 **DIAT** uses an adversarial loss to learn the mapping
337 from $x \in X$ to $y \in Y$, where x and y are face
338 images in two different domains X and Y ,
339 respectively. This method has a regularization term
340 on the mapping as $\|x - F(G(x))\|_1$ to preserve
341 identity features of the source image, where F is a
342 feature extractor pretrained on a face recognition
343 task.

344

345 **CycleGAN** also uses an adversarial loss to learn to
346 mapping between two different domains X and Y .
347 This method regularizes the mapping via cycle
348 consistency losses,

$$349 \|x - G_{YX}(G_{XY}(x))\|_1 \text{ and } \|y - G_{XY}(G_{YX}(y))\|_1.$$

350 This method requires two generators and
351 discriminators for each pair of two different
352 domains.

353

354 **IcGAN** combines an encoder with a cGAN model.
355 cGAN learns the mapping $G : \{z, c\} \rightarrow x$ that
356 generates an image x conditioned on both the

363
364 5.2 Datasets
365
366
367
368
369
370
371
372
373
374
375

369 **CelebA.** The CelebFaces Attributes (CelebA)
370 dataset contains 202,599 face images of celebrities,
371 each annotated with 40 binary attributes. We crop
372 the initial (178, 218) size images to (178, 178) then
373 resize them as (128, 128). We randomly select
374 2,000 images as test set and use all remaining
375 images for training data. We construct seven
376 domains using the following attributes: hair color
377 (black, blond, brown), gender (male / female), and
378 age (young / old).

379 **RaFD.** The Radboud Faces Database. This is the
380 original dataset that paper used. However, our
381 request to the RaFD website had been denied.
382 Hence, we prepare another dataset that also
383 contains large number of facial expression images.

384 **AffectNet.** is a large facial expression dataset with
385 around 0.4 million images manually labeled for the
386 presence of eight (neutral, happy, angry, sad, fear,
387 surprise, disgust, contempt) facial expressions
388 along with the intensity of valence and arousal. We
389 resize and crop images to (256, 256).

390 5.3 Training

391 As we mentioned above, due to implementation
392 time of baseline models and huge training time we
393 only trained our replicated StarGAN model. Since
394 we've trained with same parameter of original
395 paper, comparing with the result of original paper
396 and our model is reasonable enough. In AffectNet
397 dataset case, we adopt same parameter of RaFD
398 training of original paper. Each training takes about
one day with KCloud VPN GPU server.

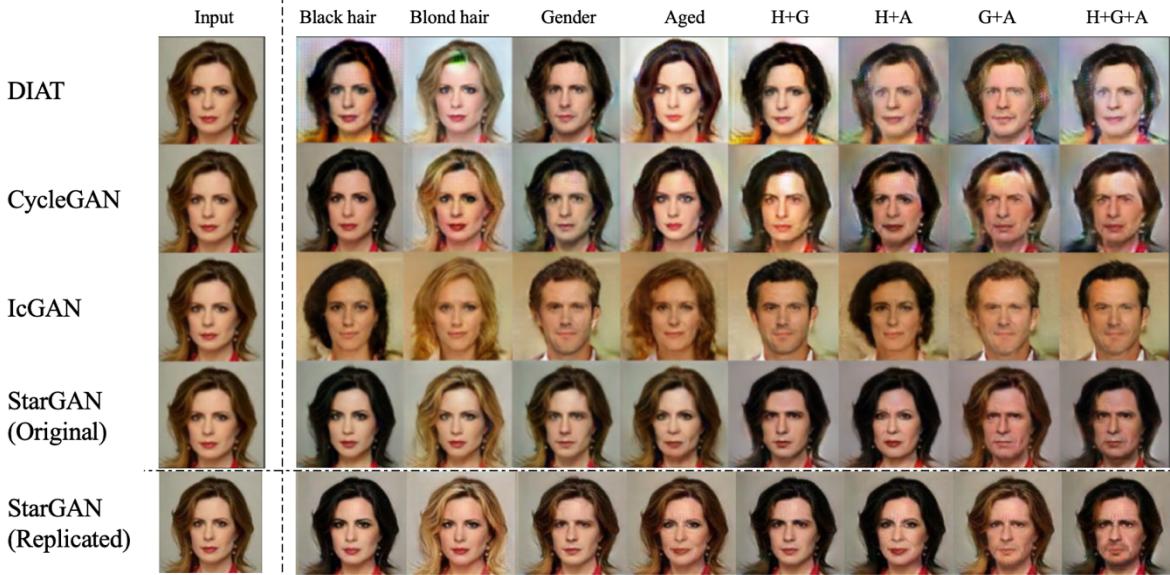


Figure 4: Facial attribute transfer results on the CelebA dataset. The first column shows the input images, next four columns show the single attribute transfer results, and rightmost columns show the multi-attribute transfer results. H: Hair color, G: Gender, A: Aged.

399

400 5.4 Experimental Results on CelebA

401 We first compare our proposed method to the
 402 baseline models on a single and multi-attribute
 403 transfer tasks. Since the implementation of the
 404 StarGAN in public github repository doesn't have
 405 multi-attribute transfer tasks, we implemented
 406 additional function that can support them.

407 As we mentioned above, all of result images are
 408 from original paper except our model.

409
 410 **Qualitative evaluation.** Fig. 4 shows the facial
 411 attribute transfer results on CelebA. We observed
 412 that our model provides a higher visual quality of
 413 translation results on test data compared to the
 414 cross-domain models and similar quality compared
 415 to the original StarGAN model.

416
 417 **Quantitative evaluation protocol.** For
 418 quantitative evaluations, the original paper used
 419 Amazon Mechanical Turk. Instead, for simplicity
 420 we performed two user studies in a survey format
 421 using google form to assess single and multiple
 422 attribute transfer tasks. Given an input image, each
 423 person were instructed to choose the best generated
 424 image based on perceptual realism, quality of
 425 transfer in attribute, and preservation of a figure's
 426 original identity.

427 First survey has a single attribute transfer in either
 428 hair color (black, brown, blond), gender, or age. In

429 another study, the generated images involve a
 430 combination of attribute transfers. Each person
 431 asked 10 questions and total 98 KAIST students
 432 were participated.

433

434 **Quantitative results.** Table1 and 2 show the
 435 results of our google form survey on single and
 436 multi-attribute transfer tasks, respectively. Our
 437 replicated StarGAN or original StarGAN obtained
 438 the majority of votes for best transferring attributes
 439 in all cases. Even though the difference of DIAT
 440 and StarGAN (original) in 'Gender' case is 26.5%
 441 and 35.7% in Table 1, the difference of DIAT and
 442 StarGAN (original) in 'G+A' case in Table 2
 443 becomes significant (8.2% and 46.9%). This result
 444 clearly showing the advantages of StarGAN in
 445 more complicated, multi-attribute transfer tasks.
 446 This is because unlike the other methods, StarGAN
 447 can handle image translation involving multiple
 448 attribute changes by randomly generating a target
 449 domain label in the training phase. Moreover, this
 450 result is from only one training of StarGAN and
 451 20~ training of other models(only support two
 452 domains), we checked that StarGAN is more
 453 efficient than any other baseline models enough.

454

Method	Hair Color	Gender	Aged
DIAT	5.1%	26.5%	0%

CycleGAN	1%	4.1%	2%
IcGAN	4.1%	29.6%	5.1%
StarGAN (original)	27.6%	35.7%	50%
StarGAN (our)	62.2%	4.1%	42.9%

455 Table 1. Google form survey evaluation for ranking
456 different models on a single attribute transfer task.
457 Each column sums to 100%.

Method	H+G	H+A	G+A	H+G+A
DIAT	4.1%	1%	8.2%	3.1%
CycleGAN	2%	12.2%	3.1%	2%
IcGAN	23.5%	1%	29.6%	36.7%
StarGAN (original)	29.6%	38.8%	46.9 %	44.9%
StarGAN (our)	40.8 %	46.9 %	12.2%	13.3%

459 Table 2. Google form survey evaluation for ranking
460 different models on a multi-attribute transfer task. H:
461 Hair color; G: Gender; A: Aged.

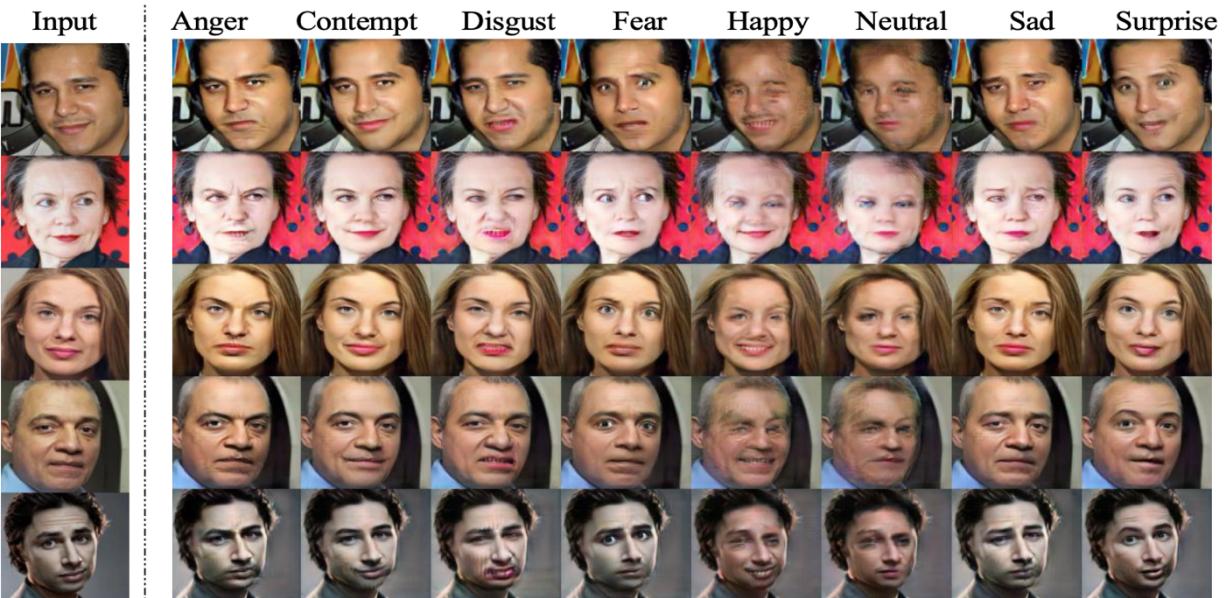


Fig 5. Facial expression synthesis results on the AffectNet dataset of our model.

463 5.5 Experimental Results on AffectNet

464 We next train our model on the AffectNet dataset
465 to learn the task of synthesizing facial expressions.
466 Since there are no any results about facial
467 expression synthesis on the AffectNet of baseline
468 models(due to lack of time to implement baseline
469 model and training AffectNet dataset for all
470 domain-domain pairs), we decide to only look at
471 the transformed images' quality.

472

473 Qualitative evaluation.

474 As seen Fig. 5, StarGAN clearly generates the
475 natural-looking facial expressions. However, our
476 results were not as natural as the results of RaFD in
477 original paper, which is presumed to be due to
478 differences in datasets. Since all of the RaFD
479 images are face-centered, looking straight ahead
480 and same background color, it'll be easy to catch
481 and extract features of each facial expression
482 domain. However, most of AffectNet images are
483 not face-centered and each person looking random
484 direction. Moreover, their background color,
485 chromaticity and brightness are very diverse.
486 Therefore, it would have been difficult to extract
487 common features in convolutional layer because
488 the direction of the face was different. In addition,
489 because each photo has a different color, the
490 transformation result did not follow the overall
491 color, and only the converted part showed
492 unnaturally different colors. For this reason, a low-
493 quality conversion seems to have appeared against
494 using RaFD dataset.

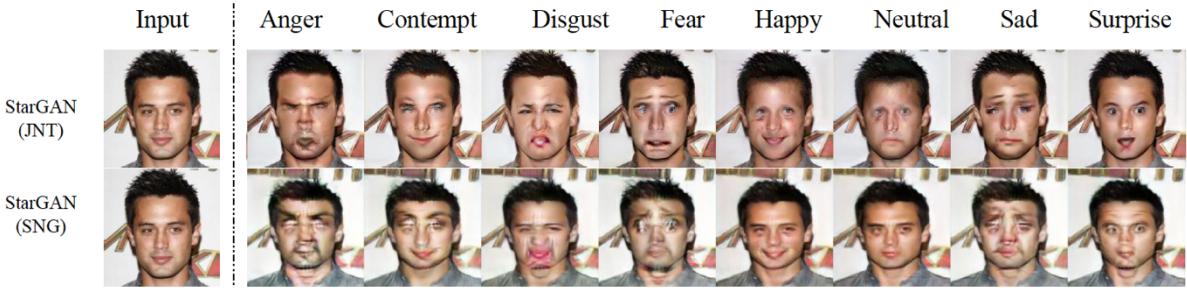


Fig 6. Facial expression synthesis results of StarGAN-SNG and StarGAN-JNT on CelebA dataset.



Fig 7. Learned role of the mask vector. All images are generated by StarGAN-JNT. The first row shows the result of applying the wrong mask vector, and the last row show the result of applying the proper mask vector.

495

496 5.6 Experimental Results on CelebA + 497 AffectNet

498 Finally, we empirically demonstrate that our
499 model can learn not only from multiple domains
500 within a single dataset, but also from multiple
501 datasets. We train our model jointly on the CelebA
502 and AffectNet datasets using the mask vector. To
503 distinguish between the model trained only on
504 AffectNet and the model trained on both CelebA
505 and AffectNet, we denote the former as StarGAN-
506 SNG (single) and the latter as StarGAN-JNT (joint).
507

508 **Effects of joint training.** Fig. 6 shows qualitative
509 comparisons between StarGAN-SNG and
510 StarGAN-JNT, where the task is to synthesize
511 facial expressions of images in CelebA. As the
512 similar result of original paper, StarGAN-JNT
513 exhibits emotional expressions with high visual
514 quality than StarGAN-SNG. Both results contain
515 some crushed image due to characteristics of
516 AffectNet dataset mentioned above, but the amount
517 of crush was greater at the StarGAN-SNG. The
518 high visual quality of StarGAN-JNT is due to the
519 fact that it learns to translate CelebA images during
520 training but not StarGAN-SNG. In other words,
521 StarGAN-JNT can leverage both datasets to

522 improve shared low-level tasks such as facial
523 keypoint detection and segmentation. By utilizing
524 both CelebA and AffectNet, StarGAN-JNT can
525 improve these low-level tasks, which is beneficial
526 to learning facial expression synthesis.

527
528 **Learned role of mask vector.** In this experiment,
529 we gave a one-hot vector c by setting the
530 dimension of a particular facial expression
531 (available from the AffectNet dataset) to one. In
532 this case, since the label associated with the
533 second dataset is explicitly given, the proper mask
534 vector would be $[0,1]$. Fig. 7 shows the case where
535 this proper mask vector was given and the
536 opposite case where a wrong mask vector of $[1,0]$
537 was given. When the wrong mask vector was used,
538 StarGAN-JNT fails to synthesize facial
539 expressions, and it manipulates the age of the
540 input image. This is because the model ignores the
541 facial expression label as *unknown* and treats the
542 facial attribute label as *valid* by the mask vector.
543 From this behavior, we can confirm that StarGAN
544 properly learned the intended role of a mask
545 vector in image-to-image translations when
546 involving all the labels from multiple datasets
547 together.

548 **6 Model Improvement Suggestions**

549

550 **6.1 Problems**

551 Even though we derive acceptable results in vanilla
 552 StarGAN, there are several parts that can be
 553 improved. From generated images by StarGAN,
 554 we've found out that our model is struggling for
 555 some specific class generation. (Fig 8.) In specific,
 556 generated images for facial expression domain was
 557 unrealistic relatively than hair, age, gender domain.
 558 We could conclude that current model has low
 559 performance at recognizing geometric or structural
 560 patterns in face.

561

576 generally. But it can't take advantage of structural,
 577 computational efficiency coming from local
 578 convolution.

579 **6.2 Our solution: Self-Attention Mechanism**

580 To resolve this problem, we suggest to use attention
 581 mechanism in StarGAN. Attention mechanism will
 582 fully utilize the long-range dependency in internal
 583 representation of image. Specifically, self-attention
 584 mechanism which translates convolution feature
 585 maps to attention feature maps is represented
 586 below. (Fig 9.) At the beginning, Input x is fed into
 587 3 different matrix function f, g, h which their
 588 weights pass through 1×1 convolution layer. Each of
 589 them has their own weights. To get attention map,
 590 matrix multiply result of $f(x)$ transpose with result
 591 of $g(x)$, then apply the softmax operation to result



592 Fig 8. Result of image translations with same source, different target domain.

593 Most of GAN based models containing StarGAN of multiplication.

594 use convolutional layers for image generation.

595 Because convolution operator has local receptive

596 field, convolution layer processes the information

597 in local neighborhood. Long-range dependency in

598 sample image will be managed only after several

599 convolution layers. Therefore, using only

600 convolution layers in GAN is inefficient for

601 modeling long-rang dependency, which leads to

602 lack of recognition of structural patterns. Simple

603 solution for this problem is increasing kernel size

604 of convolution layer. If size of kernel increases,

605 representational capacity of network will be

606 increase, and network can represent image more

$$594 \quad \beta_{j,i} = \frac{\exp(s_{ij})}{\sum_{i=1}^N \exp(s_{ij})},$$

$$595 \quad \text{where } s_{ij} = f(x_i)^T g(x_j)$$

596 Beta of above equation represents amount of
 597 relation between x_i and x_j . In other words, how
 598 much x_i is concerned in processing of x_j . After
 599 processing attention map, matrix multiply attention
 600 map with result of $h(x)$, then fed it into $v(x)$ which
 601 also passes its weight into 1×1 convolution layer to

602 get attention feature map o . Each element in output
603 feature map is defined below.

$$604 \quad o = v \left(\sum_{i=1}^N \beta_{j,i} h(x_i) \right), h(x_i) = W_h x_i, v(x_i) \\ 605 \quad = W_v x_i$$

606 The final output is defined below. In final output,
607 we add result of attention layer multiplied with
608 learnable scale parameter with input feature maps.
609 Introducing scale parameter γ , we can imply model
610 to consider local neighborhood first, then gradually
611 focus on non-local neighborhood information.

$$612 \quad y_i = \gamma o_i + x_i$$

613 Finally, loss function of attention applied generator
614 and discriminator is defined below. But we should
615 combine them with original loss functions of
616 StarGAN generator and discriminator properly to
617 utilize the attention mechanism to StarGAN.

$$619 \quad L_D = -E_{x,y \sim P_{data}} [\min(0, -1 + D(x, y))] \\ 620 \quad - E_{z \sim P_z, y \sim P_{data}} [\min(0, -1 \\ 621 \quad - D(G(z), y))] \\ 622 \quad L_G = -E_{z \sim P_z, y \sim P_{data}} [D(G(z), y)]$$

630 This unified model architecture of StarGAN allows
631 us to train multiple datasets simultaneously with
632 each different domain even within only one
633 network. Even more than that, we additionally
634 modified the adversarial loss function to the
635 objective function of WGAN (Wasserstein GAN)
636 along with the gradient penalty. We refer to our
637 model as Vanilla GAN including all jobs until this
638 point. Then, our second objective was to develop
639 this Vanilla StarGAN model with the application of
640 self-attention mechanism. This is our own original
641 work, and it aims to increase the performance of
642 StarGAN at recognizing geometric or structural
643 patterns in face. We proposed the self-attention
644 mechanism as a solution to this.

645
646 We also validated our theoretical approach
647 with experiments. We used CelebA and AffectNet
648 instead of RaFD. Due to excessively long
649 implementation time of baseline models on the
650 given GPU, we only trained our replicated
651 StarGAN model as mentioned before. In the
652 CelebA dataset, it was obvious that our StarGAN
653 model outperforms the baseline model. On the
654 other hand, we observed a low-quality conversion
655 issue when training the AffectNet dataset, and we
656 gave a logical reason for this phenomenon. We also
657 proceeded with joint training with CelebA and
658 AffectNet and further proposed the StarGAN-JNT.

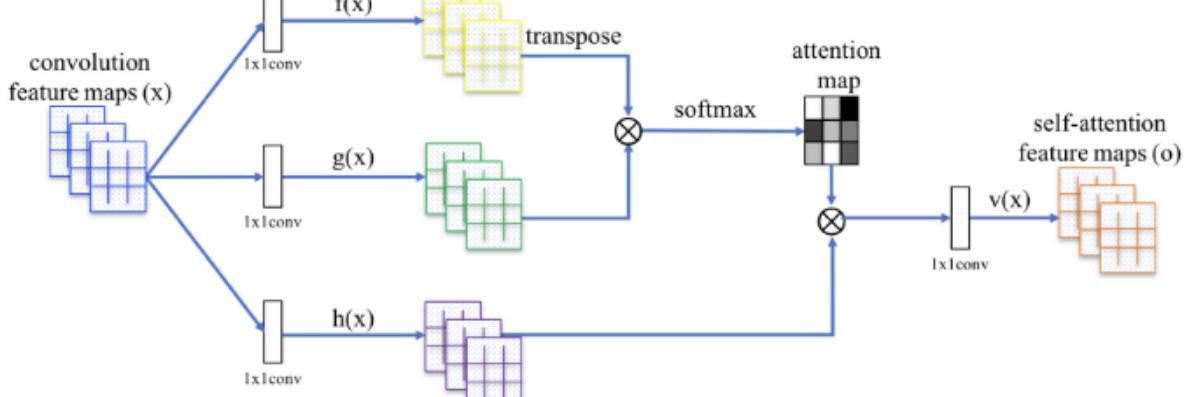


Fig 9. Overall structure of self-attention mechanism.

624 7 Discussion

625 In this paper, we proposed two specific research
626 questions. The first objective was to solve the
627 limited scalability and robustness issues of
628 previous I2I models when handling the multi-
629 domain. For this, we proposed the StarGAN model.

659 8 Conclusion

660 In this paper, we proposed the single-network
661 multi-domain I2I translation model StarGAN and a
662 possible improvement on this Vanilla StarGAN.
663 The generated images from the StarGAN model

664 present higher quality compared to previous
665 models. It's because the StarGAN possesses a high
666 generalization capability with the multi-task
667 learning setting. Furthermore, the mask vector of
668 StarGAN enables it to easily manage multiple
669 datasets with numerous domain labels sets. Also,
670 we observed that this Vanilla StarGAN shows a
671 poor performance when there exist geometric or
672 structural patterns in faces of the dataset. We also
673 gave a possible solution to this problem by
674 applying the concept of self-attention mechanism.
675 In the future investigation, we want to further study
676 how to maintain the balance between the
677 discriminator and generator when training a model
678 in the starGAN, which is one of the issues that most
679 I2I translation models need to deal with.

680 9 References

- 681 Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo
682 Ha, Sunghun Kim, Jaegul Choo. CVPR (2018, Oral).
683 StarGAN: Unified Generative Adversarial
684 Networks for Multi-Domain Image-to-Image
685 Translation <https://arxiv.org/abs/1711.09020>
- 686 Han Zhang Ian Goodfellow Dimitris Metaxa Augustus
687 Odena (2018) Self-Attention Generative
688 Adversarial Networks
689 <https://arxiv.org/abs/1805.08318>
- 690 Ali Mollahosseini, Behzad Hasani, and Mohammad H.
691 Mahoor, "AffectNet: A New Database for Facial
692 Expression, Valence, and Arousal Computation in
693 the Wild", IEEE Transactions on Affective
694 Computing, 2017
- 695 M. Li, W. Zuo, and D. Zhang. Deep identity-aware
696 transfer of facial attributes. arXiv preprint
697 arXiv:1610.05586, 2016. 2, 5, 8
- 698 Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning
699 face attributes in the wild. In Proceedings of the
700 IEEE International Conference on Computer
701 Vision (ICCV), 2015. 2, 4, 6
- 702 A. Odena, C. Olah, and J. Shlens. Conditional image
703 synthesis with auxiliary classifier gans. arXiv
704 preprint arXiv:1610.09585, 2016. 3, 5
- 705 G. Perarnau, J. van de Weijer, B. Raducanu, and J. M.
706 Álvarez. Invertibleconditionalgansforimageediting.
707 arXiv preprint arXiv:1611.06355, 2016. 5, 8
- 708 J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired
709 image- to-image translation using cycle-consistent
710 adversarial net- works. In Proceedings of the IEEE
711 International Conference on Computer Vision
712 (ICCV), 2017. 1, 2, 3, 4, 5, 8

713