# Adversarial defenses via a mixture of generators

**Maciej Żelaszczyk**
Faculty of Mathematics and Information Science
Warsaw University of Technology
00-662, Warsaw
`m.zelaszczyk@mini.pw.edu.pl`

**Jacek Mańdziuk**
Faculty of Mathematics and Information Science
Warsaw University of Technology
00-662, Warsaw
`mandziuk@mini.pw.edu.pl`

## Abstract

In spite of the enormous success of neural networks, adversarial examples remain a relatively weakly understood feature of deep learning systems. There is a considerable effort in both building more powerful adversarial attacks and designing methods to counter the effects of adversarial examples. We propose a method to transform the adversarial input data through a mixture of generators in order to recover the correct class obfuscated by the adversarial attack. A canonical set of images is used to generate adversarial examples through potentially multiple attacks. Such transformed images are processed by a set of generators, which are trained adversarially as a whole to compete in inverting the initial transformations. To our knowledge, this is the first use of a mixture-based adversarially trained system as a defense mechanism. We show that it is possible to train such a system without supervision, simultaneously on multiple adversarial attacks. Our system is able to recover class information for previously-unseen examples with neither attack nor data labels on the MNIST dataset. The results demonstrate that this multi-attack approach is competitive with adversarial defenses tested in single-attack settings.

## 1 Introduction

The discovery of adversarial examples in the image domain has fueled research to better understand the nature of such examples and the extent to which neural networks can be fooled. Work on new adversarial methods has produced various approaches to the generation of adversarial images. These *adversarial attacks* can be either geared toward misclassification in general or toward inducing the classifier assign the adversarial example to a specific class. Conversely, *adversarial defenses* are approaches designed with the goal of nullifying the effect of adversarial effects on a classifier.

In this work, we aim to show the possibility of defending against adversarial examples by utilizing a set of neural networks competing to reverse the effects of adversarial attacks.

**Our contributions:**

- We show that it is possible to use a GAN-inspired architecture to transform adversarial images to recover correct class labels.

- We show that such an architecture can be extended to handle multiple adversarial attacks at once by making the generators compete to recover the original images on which the adversarial examples are based.

- We demonstrate that such a system can be trained with no supervision, i.e. with no knowledge of: 1) the type of the adversarial attacks used, 2) the original images, 3) the class labels of the original images.

- A system trained in this way is able to recover correct class labels on previously-unseen adversarial examples for a number of adversarial attacks.

- The accuracy of a pretrained classifier on the images transformed by the system is meaningfully higher than that of a random classifier. For a number of our setups, including multi-attack ones, it is on par with the results of modern adversarial defenses, which are typically tested in single-attack settings.

- For a range of adversarial attacks, the system is able to invert the adversarial attack and produce images close to the clean images that were attacked.

- We highlight the fact that the components of such an unsupervised system specialize to some extent, revealing the distinction between attacks which produce images interpretable to humans and attacks which severely distort the visual aspects of the original images.

## 2 Related work

This work is related to a number of lines of research linked to adversarial examples. Szegedy et al. (2014) show that it is possible to construct adversarial images by adding perturbations to original images with the goal of the perturbation causing misclassification. This is done via the L-BFGS method. The resulting changes in the image may be imperceptible to the human observer but they still successfully cause the classifier to assign an incorrect class to the image relative to the unperturbed one. Goodfellow et al. (2015) show that the linear behavior of neural networks in high-dimensional spaces suffices to produce adversarial examples. They introduce the *fast gradient sign method* (FGSM) - a method of obtaining adversarial examples computationally less intensive than Szegedy et al. (2014). Ilyas et al. (2019) tie the existence of adversarial examples to the presence of statistical, but not human-interpretable, features in the data. They show that such features are pervasive in popular datasets and while predictive of the class, they are not necessarily robust.

The L-BFGS and FGSM methods were followed by an array of diverse adversarial attack methods. Papernot et al. (2016a) present the *Jacobian-based Saliency MapAttack* (JSMA). Madry et al. (2018) suggest that *projected gradient descent* (PGD) can be used to generate adversarial examples and that PGD is an universal adversarial method. Moosavi-Dezfooli et al. (2016) propose the DeepFool method to obtain adversarial examples and compare the robustness of classifiers to such examples. Carlini and Wagner (2017) show that *defensive distillation* (one of the adversarial defense methods) is still susceptible to
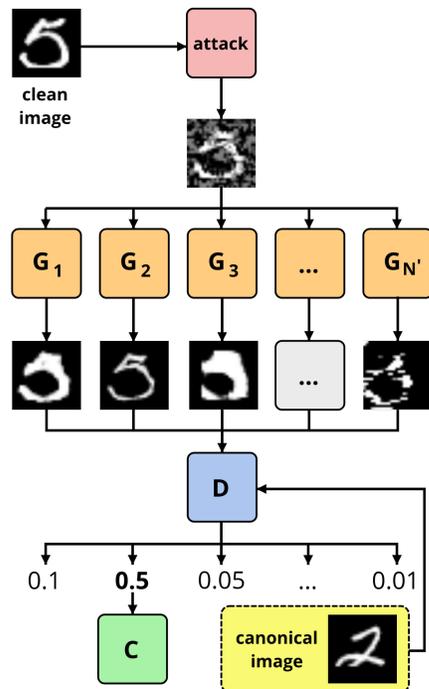


Figure 1: Overview of the system. A clean image is attacked and fed to the generators. The output of the generators is scored by the discriminator and the best-performing generator ($G_2$) is updated. The discriminator is trained on the output of all the generators and on canonical images. At test time, the output of the best-performing generator is passed to the pretrained classifier ($C$).

2

adversarial attacks. Brendel et al. (2018) introduce *Boundary Attack* - a decision-based adversarial attack gradually decreasing the size of the perturbation while ensuring that it, in fact, stays adversarial.

While initial research on adversarial examples focused on intentionally perturbing the whole image, further works have shown that attacks may be aimed at smaller areas - Su et al. (2019) show that one-pixel attacks can produce successful adversarial examples. Additionally, it has been shown that the adversarial attack do not necessarily need to be tied to a specific image. It is possible to generate general patches which can then be superimposed on original images to produce adversarial examples as demonstrated by Brown et al. (2017). These patches do not need to be produced in a strictly digital way - it is possible to physically print them and attach to other images and still get adversarial examples.

The emergence of adversarial examples has led to research on the possibility to defend from the adversarial attacks, a process known as an *adversarial defense*. Papernot et al. (2016b) introduce the *distillation defense* as a mechanism to counter adversarial attacks. Guo et al. (2018) focus on image transformations as adversarial defenses. In particular, they show that *total variance minimization* and *image quilting* can be effective adversarial defenses. Buckman et al. (2018) present the possibility of using image discretization as an adversarial defense. They show that a specific form of discretization, *thermometer encoding* increases the robustness of neural networks to adversarial attacks. Mustafa et al. (2019) propose to disentangle the intermediate learned features and show that this increases the robustness of the classifiers without deteriorating the performance on unperturbed images. Pang et al. (2019) posit that the robustness of neural networks to adversarial attacks can be improved by ensuring more diversity in ensemble models.

*Generative adversarial networks* (GANs) Goodfellow et al. (2014) are a network architecture directly stemming from the existence of adversarial examples. They are linked to our work not only via adversarial examples but also by the fact that we use a GAN-inspired architecture to handle image transformations. A more direct inspiration comes from the idea that it may be possible to train ensembles of conditional generators to identify and reverse transformations in the image domain - a concept proposed by Parascandolo et al. (2018).

## 3   Recovering relevant features through a set of generators

We approach the problem of finding adversarial defenses at the data preprocessing step. We aim to link ensemble methods with GAN-based training to invert the adversarial transformations.

### 3.1   Adversarial attacks as mechanisms

Following Parascandolo et al. (2018), let us consider a canonical distribution $P$ on $\mathbb{R}^d$ and $N$ *mechanisms* $A_1, \ldots, A_N$, which are functions. These mechanisms give rise to $N$ distributions $Q_1, \ldots, Q_N$, where $Q_j = A_j(P)$. We specifically consider a given mechanism $A_j$ to be an adversarial attack built with respect to a classifier $C$ on a given dataset. An important assumption is that at training time we receive two datasets: (1) a canonical, unperturbed dataset $\mathcal{D}_P$ drawn from the canonical distribution $P$ and (2) a transformed dataset $\mathcal{D}_Q = (x_i)_{i=1}^n$ sampled independently from a mixture of $Q_1, \ldots, Q_N$.

The aim is to learn, if possible, approximate inverse mappings from the transformed examples from $\mathcal{D}_Q$ to the base unperturbed images from $P$. In general, the task of inverting adversarial images may be a difficult one as adversarial attacks can be complex transformations, possibly severely distorting the canonical images they are based upon. Keeping this in mind, the learned mappings may not necessarily preserve the visual features but instead preserve features necessary for a neural network classifier to recover the correct label unseen by the inverse mappings.

An important point is that $\mathcal{D}_Q$ is constructed by randomly applying one of the attacks $A_1, \ldots, A_N$ to images from $P$. $\mathcal{D}_P$ contains images from $P$, other than the ones used to construct $\mathcal{D}_Q$.

### 3.2   Competitive generators

We consider a set of $N'$ functions $G_1, \ldots, G_{N'}$, which will be referred to as *generators*, parameterized by $\theta_1, \ldots, \theta_{N'}$. In principle, we do not require $N = N'$. We additionally consider a function $D : \mathbb{R}^d \to [0, 1]$, a *discriminator*, parameterized by $\theta_D$ which is required to take values close to 0 for

input outside of the support of the canonical distribution and values close to 1 for input from within this support.

The training is an adversarial procedure inspired by the training process of GANs (Goodfellow et al., 2014). Algorithm 1 outlines the training procedure. For each input example $x' \in \mathcal{D}_Q$, each generator $G_j$ is conditioned on this input. Based on the evaluation of the discriminator $D(G_j(x'))$, the generator with the highest score receives an update to its parameters $\theta_{j*}$, where $j^* = \operatorname{argmax} D(G_j(x'))$. Following that, $D$ receives an update to its parameters $\theta_D$ based on the output of all the generators and on the samples from the canonical distribution. The optimization problem can be stated as:

$$\theta_1^*, \ldots, \theta_{N'}^* = \underset{\theta_1, \ldots, \theta_{N'}}{\operatorname{argmax}} \mathbb{E}_{x' \sim Q} \left( \max_{j \in \{1, \ldots, N'\}} D(G_{\theta_j}(x')) \right). \tag{1}$$

The discriminator is trained to maximize the objective:

$$\max_{\theta_D} \left( \mathbb{E}_{x \sim P} \log \left( D_{\theta_D} \right) + \frac{1}{N'} \sum_{j=1}^{N'} \mathbb{E}_{x' \sim Q} \left( \log \left( 1 - D_{\theta_D}(G_{\theta_j}(x')) \right) \right) \right) \tag{2}$$

---

**Algorithm 1** Mixture of generators adversarial defense training algorithm

---

Choose $A_1, \ldots, A_N$ as adversarial attacks.
    **for** fixed number of initialization epochs **do**
        Train $G_1, \ldots, G_{N'}$ to approximate the identity transformation on adversarial images.
    **end for**
    **for** fixed number of training epochs **do**
        Select batch of canonical images from train set.
        Use attack $A_i$ selected at random to generate adversarial images from the batch of canonical images.
        Feed the adversarial images to $G_1, \ldots, G_{N'}$ to produce new images.
        Provide the produced images to $D$ to produce scores for each image.
        For the highest-scored generator $G_j$, update its weights against $D$.
        Update the weights of $D$ based on the images produced by all the generators $G_1, \ldots, G_{N'}$ as well as on canonical images.
    **end for**

---

## 4 Experiments

We evaluate the potential of the described system to recover the features necessary for correctly classifying adversarial images from the MNIST dataset. This is done by choosing the generators and the discriminator to be deep neural networks. The generators are fully-convolutional networks (Long et al., 2015) designed to preserve the size of the input data. The discriminator is a convolutional neural network with an increasing number of filters. Both network types use ELU activations (Clevert et al., 2016) and batch normalization (Ioffe and Szegedy, 2015). The architectural details of the networks are presented in the supplementary material.

Raw $28 \times 28$ pixel MNIST images are scaled to the $(0, 1)$ range. A separate

Table 1: Types of adversarial attacks. $\epsilon$ denotes a hyperparameter controlling the strength of the attack. An attack is considered successful if the LeNet5 classification of the attack image is incorrect

| Attack | $\epsilon$ | Success |
|---|---|---|
| fast gradient sign method (FGSM) | 0.5 | 89.8% |
| projected gradient descent (PGD) | 0.5 | 100.0% |
| DeepFool (DF) | 0.5 | 100.0% |
| additive uniform noise (AUN) | 3.5 | 90.6% |
| basic iterative method (BIM) | 0.2 | 90.6% |
| additive Gaussian noise (AGN) | 100 | 90.6% |
| repeated AGN (RAGN) | 15 | 94.5% |
| salt and pepper noise (SAPN) | 10 | 90.6% |
| sparse L1-descent (SLIDE) | 25 | 88.3% |

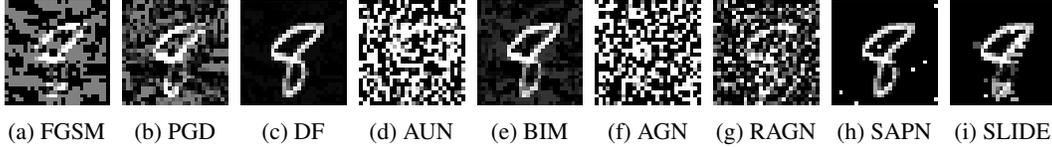| (a) FGSM | (b) PGD | (c) DF | (d) AUN | (e) BIM | (f) AGN | (g) RAGN | (h) SAPN | (i) SLIDE |

Figure 2: Sample attack images

LeNet5-inspired classifier (Lecun et al., 1998) is pretrained on the MNIST train set for 100 epochs with Adam, with an initial learning rate of $10^{-3}$. This architecture is chosen for its simplicity and relatively high accuracy. The network achieves a $98.7\%$ accuracy on the MNIST test set. The MNIST train set is divided equally into two disjoint train subsets. The first one serves as the canonical train set, which is not transformed further. The images from the other one are subject to adversarial attacks and form the transformed train set. It is worth stressing that there is no overlap between the canonical and transformed datasets, i.e. no images from the canonical train set are the basis for the images in the transformed dataset. Additionally, the system is unsupervised in that it does not use labels at train or test time.

The full list of potential attacks is given in Table 1 along with $\epsilon$ hyperparameters and the success rate of an attack on a random batch of 128 clean images from the MNIST train set. The attacks employed result in distinct visual distortions to the original images, while retaining relatively high success rates. Figure 2 shows the effects of applying the attacks to a random image from the MNIST train set. While most of the images are still recognizable to a human observer, the additive uniform noise and additive Gaussian noise attacks are not. Based on this, it can be expected that any system trying to recover original images from adversarial ones would have a hard time on these two attacks.

### 4.1 Training

We train our system by first initializing the generators over 10 epochs on the transformed images to approximate an identity transformation. This *identity initialization* procedure helps to stabilize training and partially mitigates the problem of one generator dominating the other ones (Parascandolo et al., 2018). Once identity initialization is complete, the system is trained for 100 epochs following the procedure outlined in Section 3.2. Adversarial images are fed into all the generators, which produce new images based on the input. The output of all the generators is then scored by the discriminator. The generator which achieves the highest score is then updated accordingly, as would be the case in a standard GAN training procedure. The rest of the generators remain unchanged (1). The discriminator additionally receives canonical images and scores them. It is then updated against the output of all the generators and based on the scores of the canonical images (2). We use Adam with an initial learning rate of $10^{-3}$ for all networks in both identity initialization and actual training.

### 4.2 Single-attack settings

In the first experiment, we train the system in single-attack settings with one generator. We repeat this for each attack separately. This means that the generator only receives images transformed by one attack in each setting. We evaluate the system on the previously-unseen MNIST test set by feeding transformed images to the generators, scoring the output of the generators by the discriminator and passing the highest-scored output on to the pretrained LeNet5 classifier. We measure the accuracy of the predictions in terms of the percentage of the correctly classified images. The results are presented in Table 2. For most attacks, the accuracy of the classification after the defense is above $90\%$, which highlights the possibility to use this approach as a successful adversarial defense mechanism. We also see that the success of the method seems to be related to the degree of visual distortions in the transformed images. The attacks for which our method generated less satisfactory results are also the ones which generate the most distorted images (the additive uniform noise and additive Gaussian noise attacks). Interestingly, for the DeepFool attack our method was able to marginally improve the original accuracy of the pretrained LeNet5 classifier on the test set. This suggests that for this particular attack our method comes close to completely removing the impact of the adversarial attack on classification.

Table 2: LeNet5 accuracy after defense in single attack/generator setting

| Attack | FGSM | PGD | DF | AUN | BIM | AGN | RAGN | SAPN | SLIDE |
|---|---|---|---|---|---|---|---|---|---|
| Accuracy | 97.5% | 93.5% | 98.9% | 18.9% | 97.5% | 11.4% | 92.9% | 98.1% | 95.3% |

## 4.3 Multi-attack settings

Given the results in the single-attack setting, we evaluate our system in settings with multiple attacks and generators. In the following experiments, we have chosen the number of attacks to be equal to the number of generators. This is not a hard and fast rule, and we could conceivably include a larger number of generators. In preliminary tests, we observed that limiting the number of generators can have adverse effects on the performance of the system but using a number of generators slightly lower than the number of attacks generally does not lead to a breakdown in performance. For simplicity, we proceed with a number of generators equal to the number of attacks.

Table 3: LeNet5 accuracy after defense with multiple attacks/generators. Additional results are presented in the supplementary material

| Number of attacks | 3 | 5 | 9 |
|---|---|---|---|
| Trained separately | 94.9% | 73.9% | 73.9% |
| Trained jointly | 84.4% | 64.2% | 63.3% |
| Faster initialization | 81.2% | 69.3% | 62.5% |

We train and evaluate our system in three multi-attack settings: **(a)** 3 generators, 3 attacks: FGSM, PGD, DeepFool, **(b)** 5 generators, 5 attacks: FGSM, PGD, DeepFool, AUN, BIM, **(c)** 9 generators, 9 attacks: all the attacks listed in Table 1.

For each of these settings, it is possible to adopt two approaches to training the architecture. In the first one, the system is trained jointly, end-to-end. This means that at train time neither of the generators has access to the information about the type of an attack applied to the images. In the other approach, we simply reuse the generators trained in single-attack settings and the discriminator trained in the first multi-attack approach. The main difference is that in the latter setting the generators do have information about the particular attack at train time. Such an adversarial defense is inherently an easier task - a reality reflected in the results presented in Table 3. A system trained separately (top row) consistently outperforms the system trained jointly (middle row) to the tune of a 10% difference in accuracy after the defense. However, our main interest is in the jointly-trained system, which does not have access to the type of attack at train time and can be trained in a fully unsupervised manner.

The results for the 3-attack version of the jointly-trained system are comparable to other adversarial defense methods (Schott et al., 2019) evaluated in single-attack settings. The advantage of our method is that it can be applied to multiple attacks at once with no supervision. Additionally, it can be used in a modular manner with pretrained classifiers. For the 5- and 9-attack versions of this system, the accuracy drops as the defense becomes increasingly difficult, but the negative impact of increasing the number of attacks from 5 to 9 is visibly less pronounced than in the case of moving from 3 to 5 attacks.

## 4.4 Faster initialization of the generators

The modular nature of the set of generators allows to potentially limit the computation time spent in the identity initialization phase. Since the goal of this phase is to bring all generators to an equal footing, it is possible to only initialize one generator, perturb its weights randomly and treat such random perturbations as new generators in the actual training procedure. We train such variants of the system where all the generators at the start of the actual training are copies of one identity-initialized generator with 5% of their weights perturbed by randomly setting them to 0. The rest of the training procedure remains unchanged. The results of this approach are presented in Table 3 (bottom row). The use of randomly-perturbed generators leads to qualitatively similar results, but allows to only initialize one generator instead of a full set of generators before the actual training commences.

The defenses generated by our system tend to reverse the effects of the adversarial attack, at least as far as a simple eyeball test is able to determine. The attack images based on one sample from the MNIST train set are presented in Figure 3 (first and third row) along with the images generated by our 9-attack system trained jointly (second and fourth row). Each defense is the image generated by

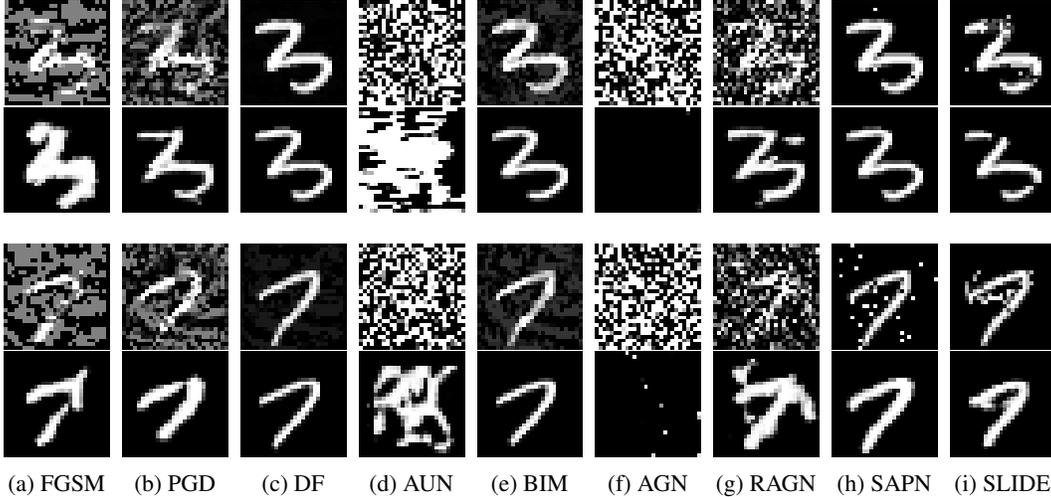(a) FGSM  (b) PGD  (c) DF  (d) AUN  (e) BIM  (f) AGN  (g) RAGN  (h) SAPN  (i) SLIDE

Figure 3: Sample attacks (first and third row) with the respective defenses (second and fourth row)

the best-scored generator. We see that for less distorted attack images our system is generally able to recover features corresponding to the clean image. In such cases, it does seem that the system is able to approximate the inversion of a particular attack. For more distorted images, for instance those generated by the AUN and AGN attacks, the system does not seem to approximate the inversion closely.

### 4.5 Large generator approach

One question that can be asked is how a mixture system fares against a single generator with similar capacity. Each of our generators has $28,609$ trainable weights. For a system consisting of 9 generators, this gives $257,481$ trainable weights. We construct a system with one generator with $333,697$ trainable weights to match the capacity of 9 standard generators. The details of the architecture are presented in the supplementary material.

We train a system with this large generator consistently with our standard training procedure for 9 attacks. This slightly more powerful system achieves an accuracy of $66.7\%$ - similar to the accuracy achieved by a mixture of generators. However, the advantage of training a mixture vs. a single large net is that the elements of the system can be inspected more easily. On top of that, this allows for more flexibility in the setup of the defense and it opens the possibility to reuse elements of the system.

### 4.6 Generator *expertise*

An interesting aspect of the mixture of generators is the extent to which particular generators specialize in either specific attacks or image classes. We perform a breakdown of the accuracy measure by digit class and by attack for the 9-attack jointly-trained system. In Figure 4, we present the results for 3 out of the 9 generators to highlight the potential types of generators trained.

The first generator (top row) is a *generalist* - it wins on a large number of sample images and achieves relatively high accuracy on the majority of the attacks. Notably, the two hard attacks we have identified so far (additive uniform noise and additive Gaussian noise attacks) are precisely the attacks this generator does not cover - it does not win on a single sample image from these attacks. *Generalists* are the most ubiquitous among the 9 trained generators.

The second generator only ever wins on samples from the two hard attacks. This generator is a *specialist* in that it focuses on specific kinds of attacks at the expense of other ones. The accuracy of the *specialist* is very high in only a handful of specific cases.

Finally, the third generator is a *marginalist* - much like the *specialist*, it focuses on two hard attacks but fails to win on a significant number of samples and has questionable accuracy levels.
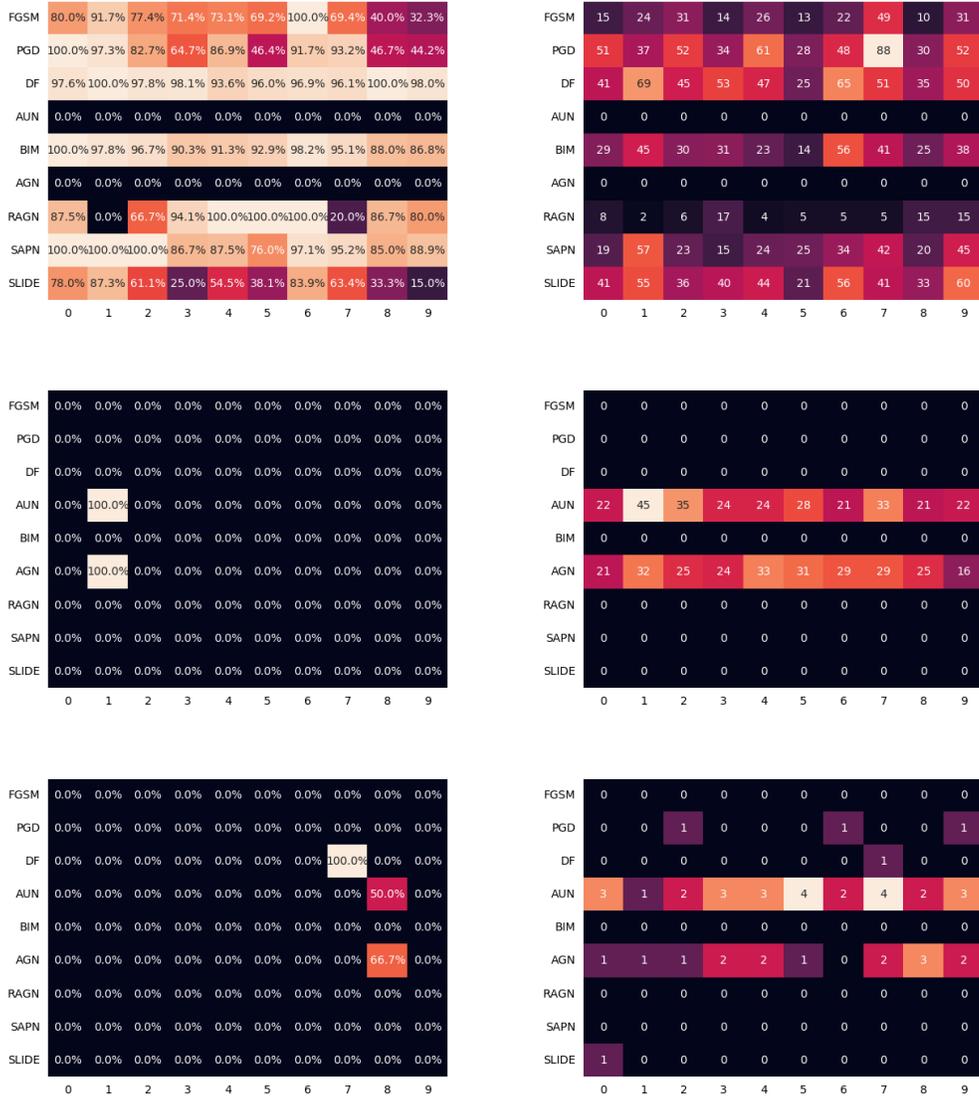
Figure 4: Examples of generator specialization: *generalist* (top row), *specialist* (middle row) and *marginalist* (bottom row). The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right)

# 5   Conclusions

In this work we show that it is possible to train a GAN-based mixture of generators to defend against adversarial examples as a data preprocessing step. The system is fully unsupervised in that it does not use attack or digit labels at all at train time. This system, trained adversarially with competing generators, is able to extract from adversarially transformed images features relevant for classification in both single-attack and multi-attacks settings. A system with 3 generators is able to defend against 3 adversarial attacks with results comparable with other modern single-attack adversarial defenses. At first, increasing the number of attacks to defend against, lowers the post-defense classification accuracy but the results suggest that extending the system further from 5 to 9 attacks/generators does not significantly decrease the accuracy.

The results show that for a majority of tested attacks, the system is not only able to recover statistical data features important for classification but that the generators also approximate the inverses of the

adversarial transformations, at least on a visual level. The images recovered by our system retain internal visual coherence and are immediately interpretable by a human observer, except for those generated by two attacks which severely distort the original images (AUN, AGN).

The system is able to internally distinguish between these hard attacks and other attacks, which produce adversarial images that are still recognizable to humans. The distinction between these types of attacks is visible in the process of generator specialization. Various generators specialize in various parts of the data space and the training process is able to produce more general generators, which mainly operate on easier attacks, and more specialized generators, which in turn tend to focus on harder attacks with different levels of success.

We show that infusing the system with knowledge as to the type of the attack by training generators separately increases the post-defense accuracy by about $10\%$ in 3-, 5- and 9-attack settings. It is also possible to limit the computational intensity of the system by copying and randomly perturbing one initialized generator without materially impacting the accuracy. The results for a more compartmentalized 9-attack/generator system are on par with the results for a similar system designed with one large generator. Given the increased interpretability and modularity of the mixture approach, it does seem that a modular approach might be better suited to the problem of inverting adversarial attacks.

Our research links adversarial defenses with ensemble methods and GAN-based training. From a more general point of view, it provides a fully-unsupervised approach to defending against adversarial examples without modifying the underlying classifier. This could provide for a practically-relevant method to approach the problem of adversarial examples in a real-life setting, particularly when the type of attack is not known beforehand and it is not feasible to modify the classifier.

This work also highlights interesting areas for further research. In particular, it hints at the possibility of linking adversarial defenses with research on causality. It would be particularly interesting to see whether it is possible to further specialize the components of ensemble systems to reveal the underlying causal structure behind adversarial attacks.

## Acknowledgment

## References

Brendel, W., Rauber, J., and Bethge, M. (2018). Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *International Conference on Learning Representations*.

Brown, T., Mane, D., Roy, A., Abadi, M., and Gilmer, J. (2017). Adversarial patch. In *NIPS Workshop*.

Buckman, J., Roy, A., Raffel, C., and Goodfellow, I. (2018). Thermometer encoding: One hot way to resist adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.

Carlini, N. and Wagner, D. (2017). Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57.

Clevert, D., Unterthiner, T., and Hochreiter, S. (2016). Fast and accurate deep network learning by exponential linear units (elus). In Bengio, Y. and LeCun, Y., editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc.

Goodfellow, I., Shlens, J., and Szegedy, C. (2015). Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*.

Guo, C., Rana, M., Cissé, M., and van der Maaten, L. (2018). Countering adversarial images using input transformations. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.

Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., and Madry, A. (2019). Adversarial examples are not bugs, they are features. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32*, pages 125–136. Curran Associates, Inc.

Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Bach, F. and Blei, D., editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–456, Lille, France. PMLR.

Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324.

Long, J., Shelhamer, E., and Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.

Moosavi-Dezfooli, S., Fawzi, A., and Frossard, P. (2016). Deepfool: A simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2574–2582.

Mustafa, A., Khan, S., Hayat, M., Goecke, R., Shen, J., and Shao, L. (2019). Adversarial defense by restricting the hidden space of deep neural networks. In *The IEEE International Conference on Computer Vision (ICCV)*.

Pang, T., Xu, K., Du, C., Chen, N., and Zhu, J. (2019). Improving adversarial robustness via promoting ensemble diversity. In *International Conference on Machine Learning*.

Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., and Swami, A. (2016a). The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS P)*, pages 372–387.

Papernot, N., McDaniel, P., Wu, X., Jha, S., and Swami, A. (2016b). Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE Symposium on Security and Privacy (SP)*, pages 582–597.

Parascandolo, G., Kilbertus, N., Rojas-Carulla, M., and Schölkopf, B. (2018). Learning independent causal mechanisms. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80 of *Proceedings of Machine Learning Research*, pages 4033–4041. PMLR.

Schott, L., Rauber, J., Bethge, M., and Brendel, W. (2019). Towards the first adversarially robust neural network model on MNIST. In *International Conference on Learning Representations*.

Su, J., Vargas, D. V., and Sakurai, K. (2019). One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2014). Intriguing properties of neural networks. In *International Conference on Learning Representations*.

# Supplementary material

## A  Architectural details

For brevity, we omitted the details of the system architecture in the main paper. The architectures of the generators and the discriminator are presented in Table 4. The generators are fully-convolutional deep neural networks with exponential linear units (ELUs) and batch normalization. They take an image of a given size as input and produce an image of the same size as output. In principle, each generator follows the logic of a conditional GAN generator, where the network is conditioned on the input. In this concrete case, the generator is conditioned on an image to produce an image of the same size. The goal is to invert the adversarial attacks and generate clean images.

Additionally, in Table 5 we present the architecture of the large generator discussed in Section 4.5 of the main paper. The large generator preserves the general characteristics of the individual generators by employing the same kernel size, padding, activation functions and batch normalization. The main difference between the large network and the individual generators is the number of layers and the fact that for the large generator the number of channels increases toward the middle layers and then decreases toward the output layer. This architecture is constructed so as to contain a number of parameters comparable to 9 individual generators. The aim is to test whether separate generators are able to learn comparably well as one large network while offering more interpretability.

Table 4: Generator/discriminator architectures

| CNN Generators | CNN Discriminator |
| --- | --- |
| Input 28x28 | Input 28x28 |
| Conv 3x3, 32, padding=1, ELU | 3 x Conv 3x3, 16, padding=1, ELU |
| BatchNorm | AvgPool 2x2 |
| Conv 3x3, 32, padding=1, ELU | 2 x Conv 3x3, 32, padding=1, ELU |
| BatchNorm | AvgPool 2x2 |
| Conv 3x3, 32, padding=1, ELU | 2 x Conv 3x3, 64, padding=1, ELU |
| BatchNorm | AvgPool 2x2 |
| Conv 3x3, 32, padding=1, ELU | FC 1024, ELU |
| BatchNorm | FC 1, Sigmoid |
| Conv 3x3, 1, padding=1, Sigmoid | |

Table 5: Large generator architecture

| CNN Generator |
| --- |
| Input 28x28 |
| Conv 3x3, 32, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 64, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 128, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 128, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 64, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 32, padding=1, ELU |
| BatchNorm |
| Conv 3x3, 1, padding=1, Sigmoid |

### A.1  Classification accuracy after defenses in repeated experiments

In the interest of verifying the stability of the obtained results, we repeat the experiments for multi-attack settings. The results of these repeated experiments are presented in Table 6. In particular,

we train the basic joint setup twice for 3, 5 and 9 attacks. Given that the results adhere to what we obtained previously, we further focus on the *faster initialization* setting. We repeat the experiments with faster initialization for 3, 5 and 9 attacks 10 times. Results for each individual experiment are presented along with means and standard errors. The outcomes show that the stability of the experiments increases with the number of generators. While we see some variability for the 3-attack settings, it decreases for the 5- and 9-attack settings, which suggests that the results are relatively stable, in particular for the settings with a larger number of attacks.

Table 6: LeNet5 accuracy after defense, multiple attacks/generators. Means for the *faster initialization* method are presented with standard errors.

| Number of attacks | 3 | 5 | 9 |
|---|---|---|---|
| Trained jointly | | | |
| Experiment 1 | 74.1% | 61.1% | 65.1% |
| Experiment 2 | 82.7% | 69.8% | 57.7% |
| Trained jointly - faster initialization | | | |
| Experiment 1 | 83.6% | 63.6% | 65.1% |
| Experiment 2 | 74.7% | 68.0% | 66.4% |
| Experiment 3 | 81.6% | 66.3% | 60.9% |
| Experiment 4 | 60.6% | 66.4% | 66.3% |
| Experiment 5 | 81.8% | 69.0% | 60.9% |
| Experiment 6 | 82.6% | 69.0% | 63.3% |
| Experiment 7 | 72.0% | 71.2% | 62.0% |
| Experiment 8 | 67.0% | 65.7% | 63.5% |
| Experiment 9 | 67.4% | 67.1% | 61.6% |
| Experiment 10 | 70.3% | 70.2% | 61.8% |
| Mean | $74.2 \pm 8.0\%$ | $67.7 \pm 2.3\%$ | $63.2 \pm 2.1\%$ |

## A.2 Generator analysis for single attacks

In reference to Table 2 of the main paper, Figures 5 and 6 provide an analysis of the LeNet5-based classification on images transformed by a system trained in single-attack settings. The results are reported for all 9 attacks considered in the paper. The heatmaps show that in single-attack settings, the system is able to recover features relevant for classification for most attacked images. Two attacks, the additive uniform noise attack (AUN) and the additive Gaussian noise attack (AGN) stand out as particularly hard in terms of recovering the original labels - which confirms the visual inspection of the samples generated by these attacks (discussed in the main paper). Interestingly, we still see some specialization. For instance, the system trained on AGN allows the LeNet5 classifier to achieve a perfect classification accuracy post-defense on images labeled as ones.

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 99.0% | 99.2% | 98.5% | 97.6% | 97.5% | 95.6% | 97.7% | 98.2% | 96.5% | 94.8% |
| PGD | 98.0% | 99.2% | 95.3% | 88.2% | 94.1% | 91.5% | 96.7% | 95.1% | 90.1% | 86.4% |
| DF | 99.6% | 99.5% | 99.0% | 99.1% | 99.0% | 98.0% | 98.6% | 99.1% | 98.7% | 98.3% |
| AUN | 28.3% | 8.1% | 34.8% | 10.2% | 40.2% | 10.7% | 14.2% | 12.6% | 25.9% | 5.1% |
| BIM | 98.6% | 98.9% | 98.0% | 97.5% | 98.2% | 97.2% | 97.2% | 96.8% | 97.1% | 95.4% |

Figure 5: Class accuracy breakdown for single-attack settings, part 1 (generators 1 to 5). Each heatmap corresponds to one generator. The heatmaps present the accuracy for a given class.

Figure 6: Class accuracy breakdown for single-attack settings, part 2 (generators 6 to 9). Each heatmap corresponds to one generator. The heatmaps present the accuracy for a given class.

## A.3 Generator analysis for multiple attacks - trained jointly

Similarly to Section A.2, we report the results for multi-attack/multi-generator settings discussed in the main paper. The setup with 3 attacks/generators is presented in Figure 7. Analogous results for the setup with 5 attacks/generators are presented in Figures 8 and 9. The largest setup with 9 attacks/generators is reported in Figures 10, 11 and 12. The figures for the 9 generator setup supplement the examples presented in Figure 4 in the main paper. Due to space limits we could not accommodate more heatmap examples in the body of the main paper.

The results confirm the earlier remarks that AUN and AGN, due to their high level of visual distortion of the images, are harder to defend against. We also see that there are roughly three types of generators: the ones covering a wide variety of attacks (*generalists*), generators covering only a subset of attacks/classes and recovering the label data relatively well (*specialists*), and generators covering a subset of attacks/classes without substantial success in recovering label data (*marginalists*). This is particularly visible for the 9-attack setting (Figures 10, 11 and 12).

Figure 7: Class/attack accuracy breakdown for the 3 attack/generator setting. Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
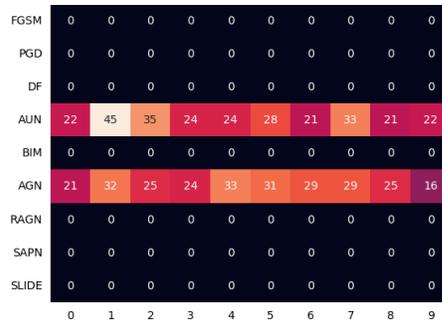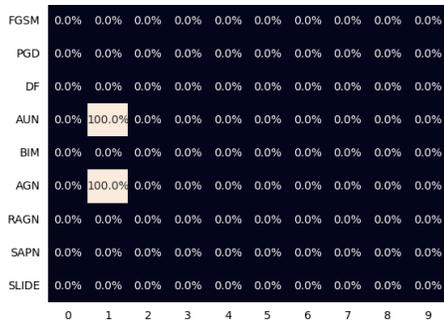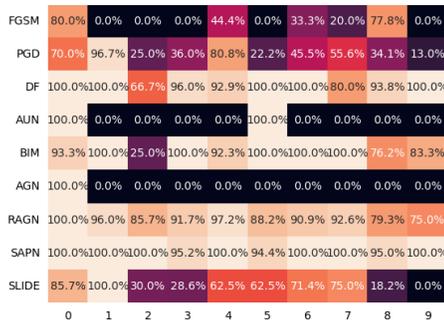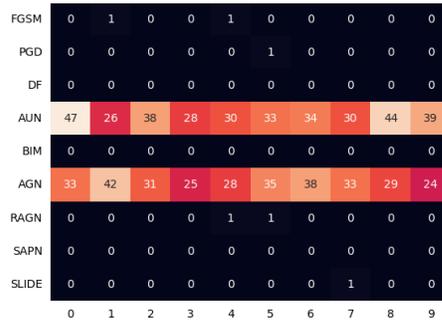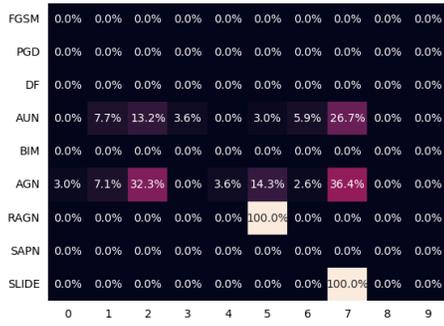
Figure 8: Class/attack accuracy breakdown for the 5 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

Figure 9: Class/attack accuracy breakdown for the 5 attack/generator setting, part 2 (generators 4 and 5). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

**Generator 1 — accuracy (left)**

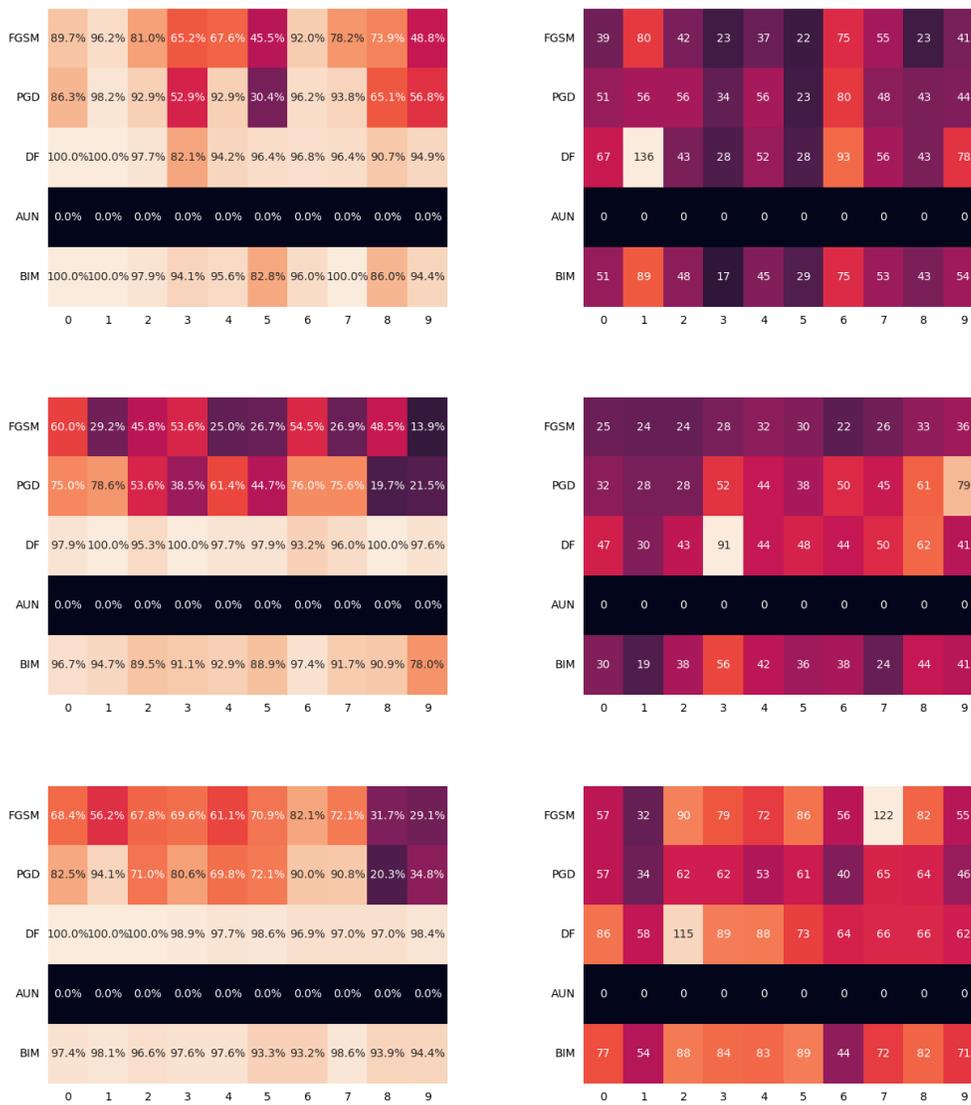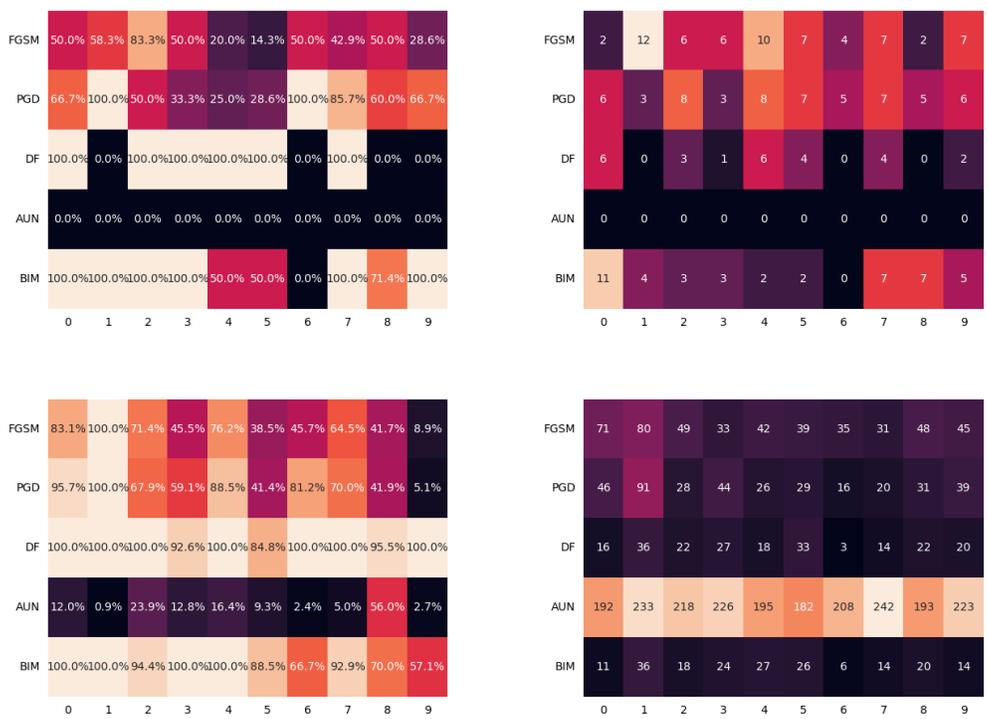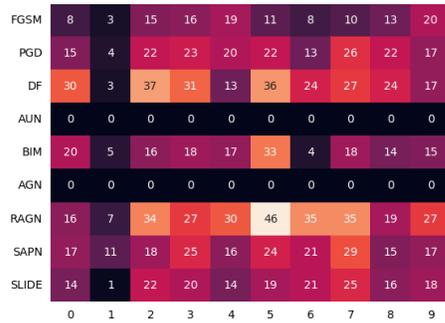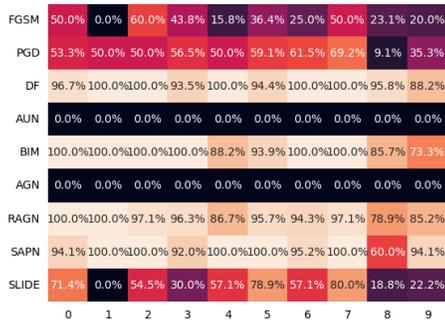| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 80.0% | 91.7% | 77.4% | 71.4% | 73.1% | 69.2% | 100.0% | 69.4% | 40.0% | 32.3% |
| PGD | 100.0% | 97.3% | 82.7% | 64.7% | 86.9% | 46.4% | 91.7% | 93.2% | 46.7% | 44.2% |
| DF | 97.6% | 100.0% | 97.8% | 98.1% | 93.6% | 96.0% | 96.9% | 96.1% | 100.0% | 98.0% |
| AUN | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| BIM | 100.0% | 97.8% | 96.7% | 90.3% | 91.3% | 92.9% | 98.2% | 95.1% | 88.0% | 86.8% |
| AGN | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| RAGN | 87.5% | 0.0% | 66.7% | 94.1% | 100.0% | 100.0% | 100.0% | 20.0% | 86.7% | 80.0% |
| SAPN | 100.0% | 100.0% | 100.0% | 86.7% | 87.5% | 76.0% | 97.1% | 95.2% | 85.0% | 88.9% |
| SLIDE | 78.0% | 87.3% | 61.1% | 25.0% | 54.5% | 38.1% | 83.9% | 63.4% | 33.3% | 15.0% |

**Generator 1 — samples won (right)**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 15 | 24 | 31 | 14 | 26 | 13 | 22 | 49 | 10 | 31 |
| PGD | 51 | 37 | 52 | 34 | 61 | 28 | 48 | 88 | 30 | 52 |
| DF | 41 | 69 | 45 | 53 | 47 | 25 | 65 | 51 | 35 | 50 |
| AUN | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| BIM | 29 | 45 | 30 | 31 | 23 | 14 | 56 | 41 | 25 | 38 |
| AGN | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| RAGN | 8 | 2 | 6 | 17 | 4 | 5 | 5 | 5 | 15 | 15 |
| SAPN | 19 | 57 | 23 | 15 | 24 | 25 | 34 | 42 | 20 | 45 |
| SLIDE | 41 | 55 | 36 | 40 | 44 | 21 | 56 | 41 | 33 | 60 |

**Generator 2 — accuracy (left)**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 45.5% | 91.7% | 42.9% | 50.0% | 18.5% | 52.5% | 50.0% | 56.4% | 26.5% | 9.7% |
| PGD | 41.2% | 80.0% | 48.7% | 34.4% | 50.0% | 55.2% | 68.4% | 63.6% | 18.2% | 15.6% |
| DF | 100.0% | 80.0% | 90.5% | 94.7% | 100.0% | 95.2% | 97.2% | 97.4% | 100.0% | 91.3% |
| AUN | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 33.3% | 66.7% | 0.0% | 33.3% | 0.0% |
| BIM | 87.5% | 100.0% | 85.7% | 100.0% | 100.0% | 94.1% | 100.0% | 100.0% | 92.3% | 86.7% |
| AGN | 0.0% | 0.0% | 50.0% | 50.0% | 33.3% | 0.0% | 0.0% | 20.0% | 33.3% | 0.0% |
| RAGN | 90.0% | 100.0% | 100.0% | 82.8% | 100.0% | 96.6% | 96.2% | 89.2% | 80.5% | 86.0% |
| SAPN | 100.0% | 83.3% | 100.0% | 86.2% | 91.3% | 88.9% | 97.6% | 95.1% | 90.7% | 100.0% |
| SLIDE | 57.1% | 66.7% | 50.0% | 24.2% | 48.1% | 33.3% | 76.9% | 61.4% | 22.2% | 10.8% |

**Generator 2 — samples won (right)**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 22 | 12 | 28 | 38 | 27 | 40 | 44 | 39 | 34 | 31 |
| PGD | 17 | 5 | 39 | 32 | 26 | 29 | 38 | 33 | 33 | 32 |
| DF | 8 | 5 | 21 | 19 | 26 | 21 | 36 | 39 | 36 | 23 |
| AUN | 1 | 2 | 4 | 4 | 3 | 3 | 3 | 4 | 3 | 1 |
| BIM | 8 | 2 | 14 | 22 | 12 | 17 | 21 | 15 | 26 | 30 |
| AGN | 4 | 3 | 2 | 2 | 3 | 0 | 0 | 5 | 3 | 2 |
| RAGN | 10 | 9 | 31 | 29 | 28 | 29 | 52 | 37 | 41 | 43 |
| SAPN | 12 | 6 | 23 | 29 | 23 | 27 | 42 | 41 | 43 | 39 |
| SLIDE | 14 | 6 | 20 | 33 | 27 | 21 | 26 | 44 | 36 | 37 |

**Generator 3 — accuracy (left)**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 100.0% | 50.0% | 83.3% | 66.7% | 0.0% | 50.0% | 0.0% | 0.0% | 25.0% | 0.0% |
| PGD | 50.0% | 0.0% | 50.0% | 38.5% | 100.0% | 0.0% | 40.0% | 0.0% | 0.0% | 0.0% |
| DF | 100.0% | 0.0% | 100.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| AUN | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| BIM | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| AGN | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% |
| RAGN | 100.0% | 100.0% | 75.0% | 16.7% | 80.0% | 0.0% | 57.1% | 50.0% | 0.0% | 0.0% |
| SAPN | 100.0% | 0.0% | 100.0% | 50.0% | 0.0% | 100.0% | 66.7% | 0.0% | 33.3% | 0.0% |
| SLIDE | 100.0% | 0.0% | 100.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 100.0% |

**Generator 3 — samples won (right)**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| FGSM | 2 | 4 | 6 | 3 | 1 | 2 | 2 | 8 | 4 | 1 |
| PGD | 4 | 0 | 2 | 13 | 1 | 5 | 5 | 2 | 0 | 0 |
| DF | 2 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| AUN | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| BIM | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| AGN | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| RAGN | 3 | 7 | 4 | 6 | 5 | 0 | 7 | 2 | 1 | 0 |
| SAPN | 1 | 0 | 1 | 2 | 0 | 1 | 3 | 0 | 3 | 0 |
| SLIDE | 2 | 0 | 3 | 0 | 2 | 1 | 0 | 0 | 1 | 1 |

Figure 10: Class/attack accuracy breakdown for the 9 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
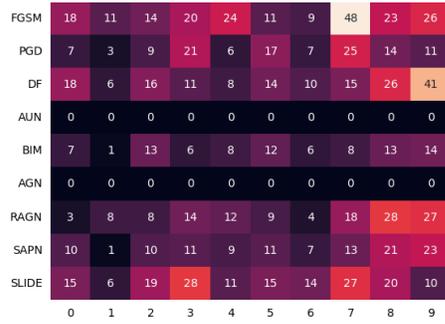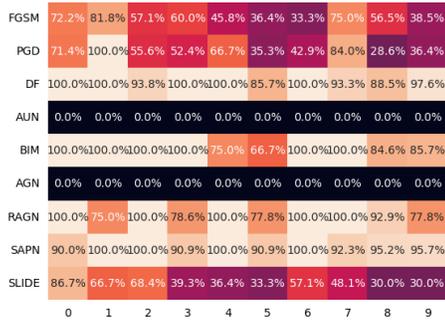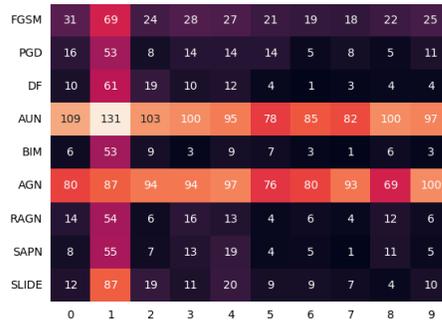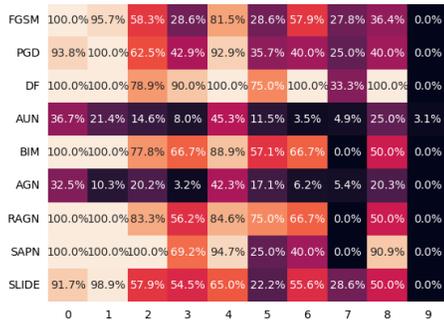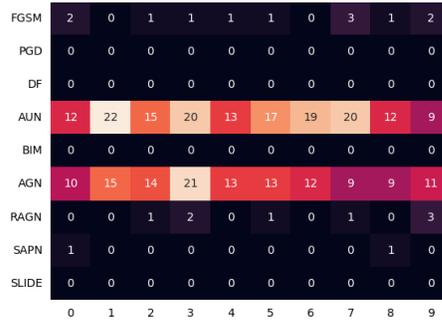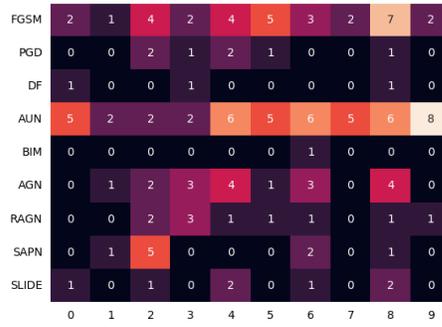
Figure 11: Class/attack accuracy breakdown for the 9 attack/generator setting, part 2 (generators 4 to 6). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

Figure 12: Class/attack accuracy breakdown for the 9 attack/generator setting, part 3 (generators 7 to 9). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
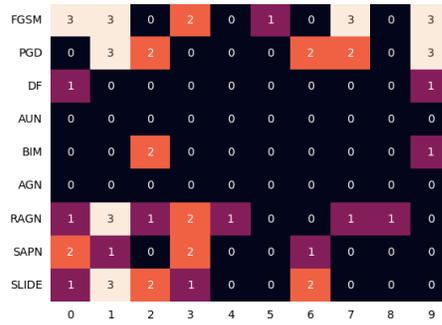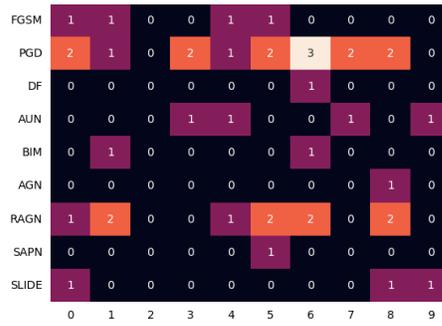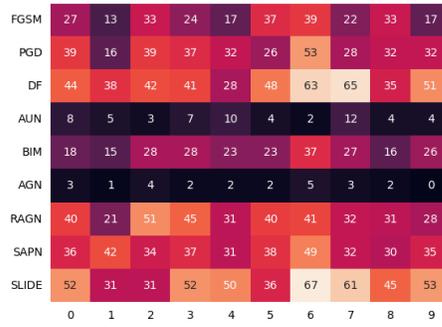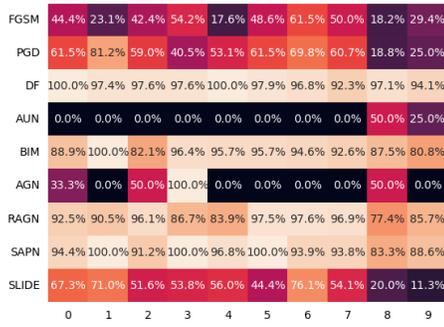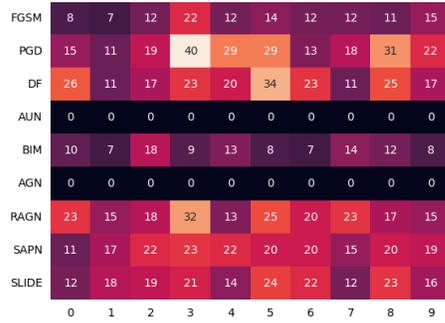
## A.4  Generator analysis in repeated experiments trained jointly with *faster initialization*

The results obtained using the faster initialization method do not materially differ (both in terms of post-defense accuracy and in terms of the qualitative assessment of the heatmaps) from the results of standard joint training without faster initialization shown in Section A.3. Due to the comparable outcomes of the two approaches, we present further results for the faster initialization method as it affords us the possibility to use less computational resources and speed up training.

The heatmaps for 3 sample sets of repeated multi-attack experiments with faster initialization are shown below in Figures 13-30 (the experiment numbers correspond to those presented in Table 6). Each set of experiments encompasses 3 settings: 3, 5 and 9 attacks/generators. The main observation is that the specialization of generators becomes more visible in settings with more attacks/generators. For the 3-attack setting, the generators are of the *generalist* type, with potential variation in the effectiveness of the defense on different attacks. For the 5-attack setting, specialization in the additive uniform noise attack (AUN) becomes more visible, occasionally producing *specialists*. Finally, for the 9-attack setting, there is further specialization as some of the generators become *specialists* or *marginalists*.

## A.4.1 Experiment 3, trained jointly - faster initialization



Figure 13: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 3 attack/generator setting. Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
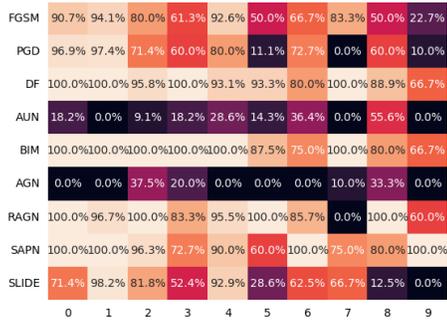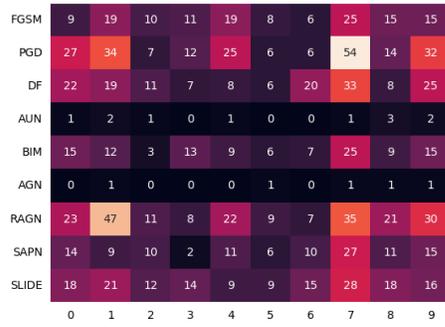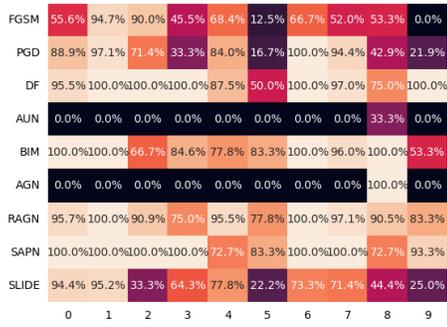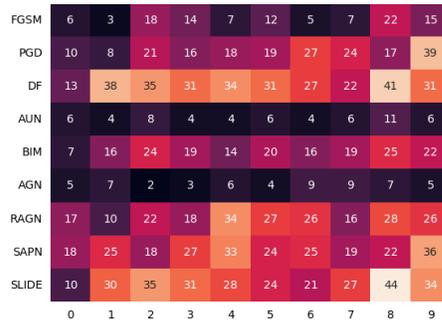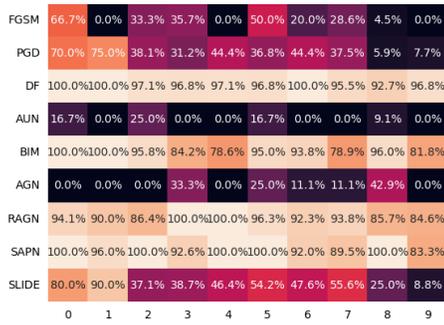
Figure 14: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
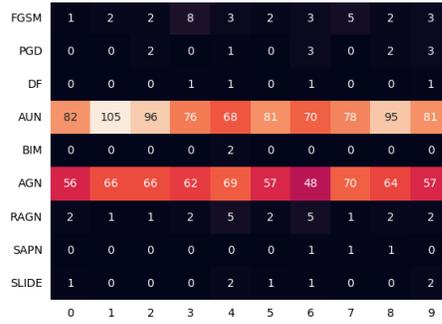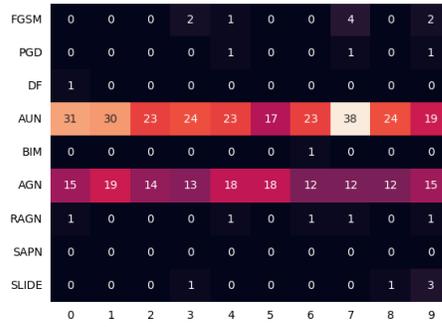
Figure 15: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 2 (generators 4 and 5). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

Figure 16: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

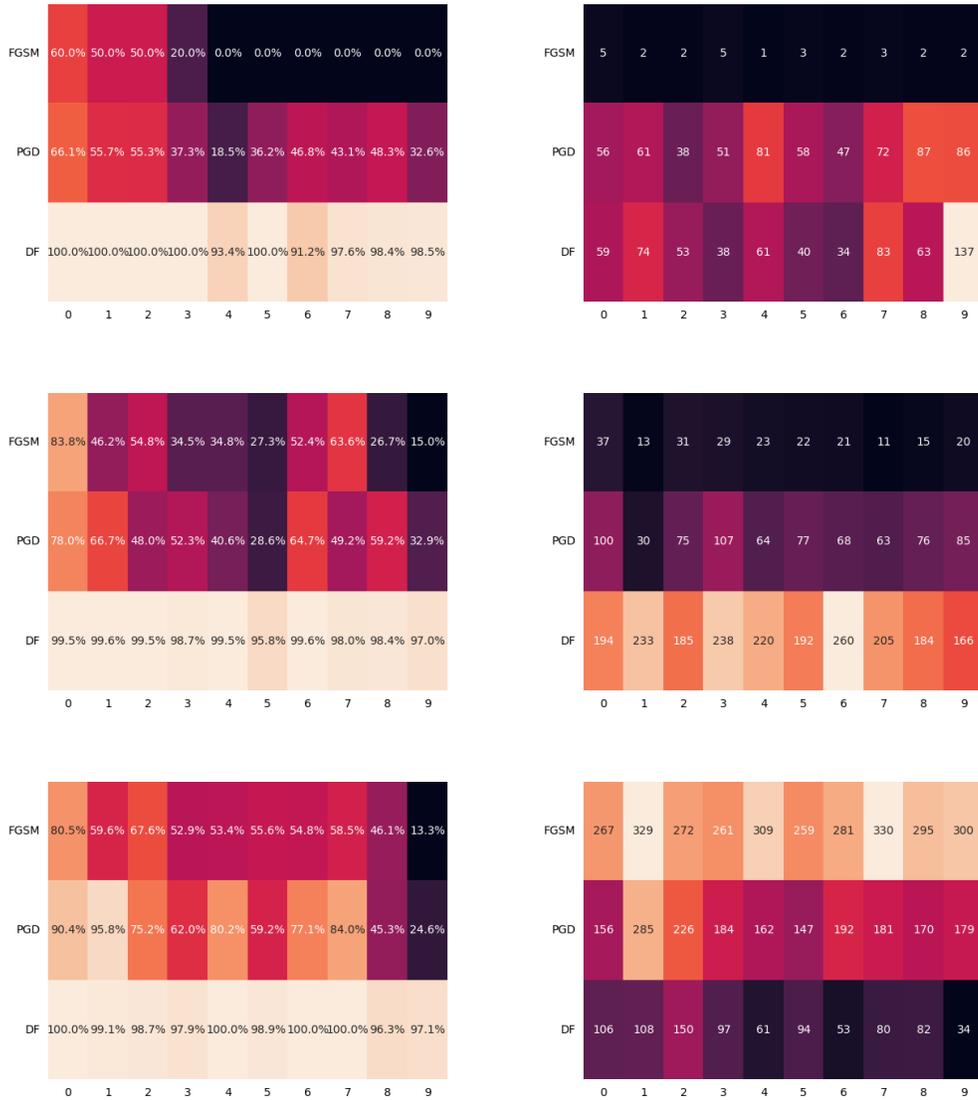Figure 17: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 2 (generators 4 to 6). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
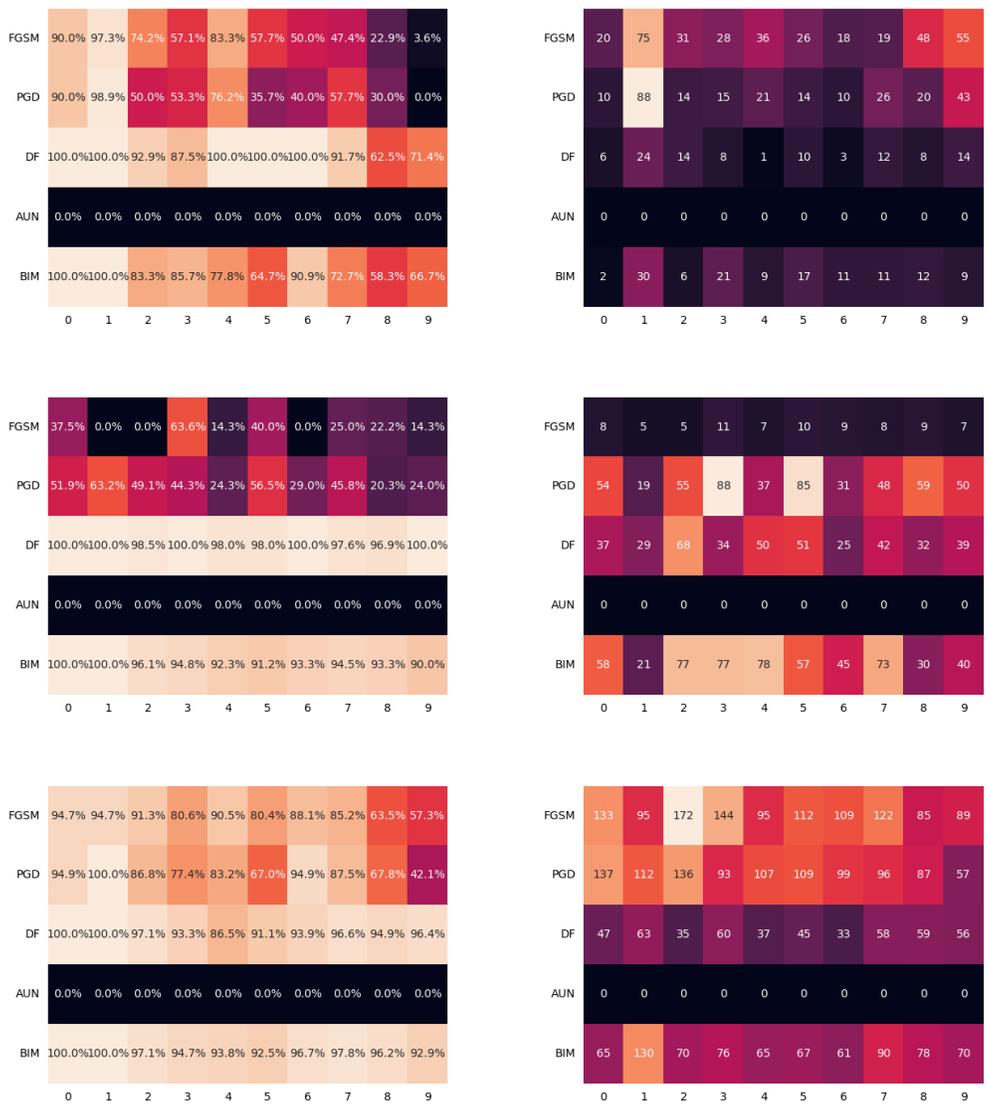
Figure 18: Experiment 3 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 3 (generators 7 to 9). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

Figure 19: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 3 attack/generator setting. Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
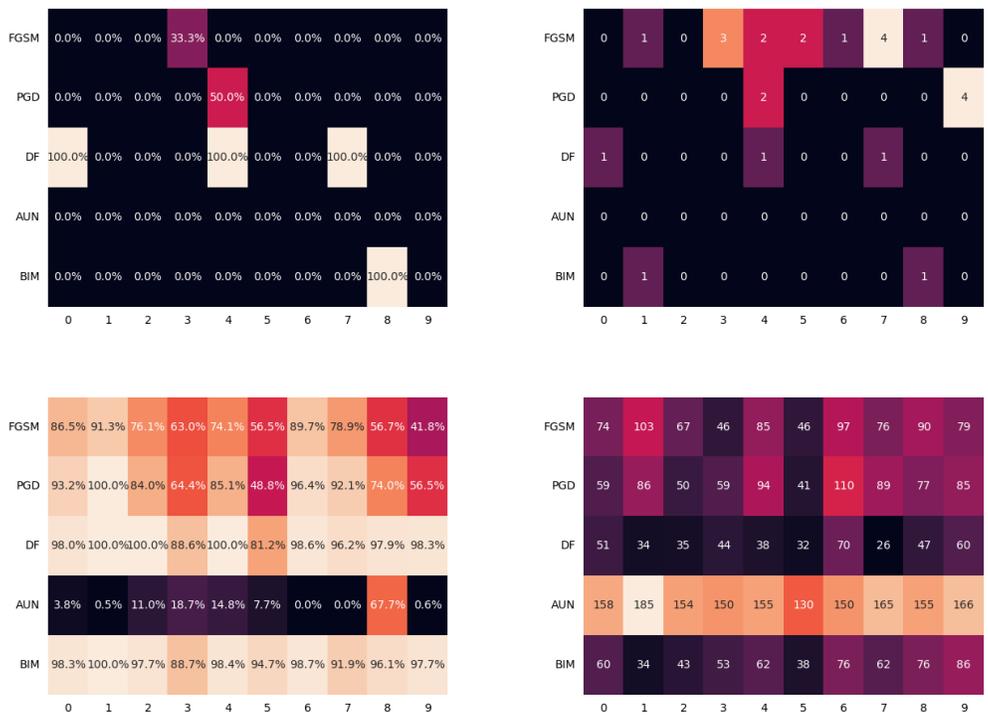
Figure 20: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
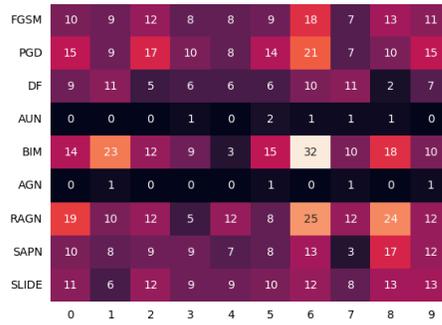
Figure 21: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 2 (generators 4 and 5). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
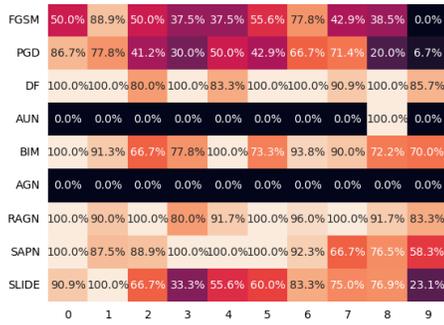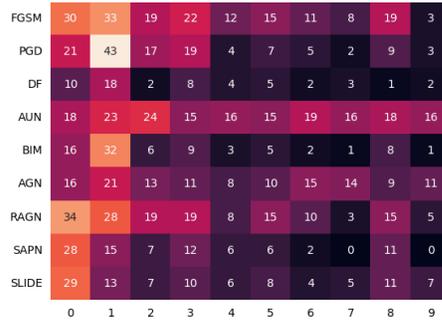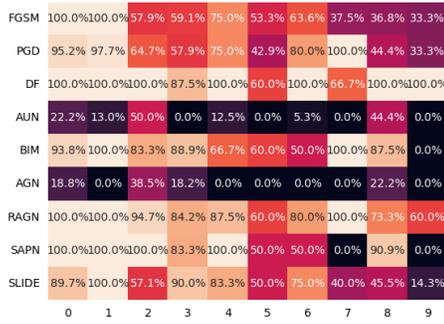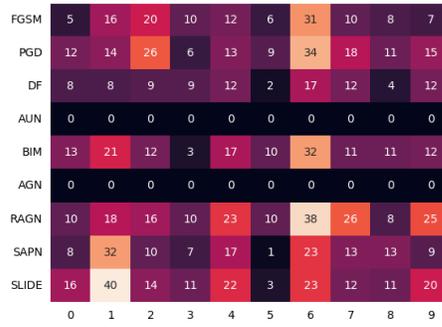
Figure 22: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
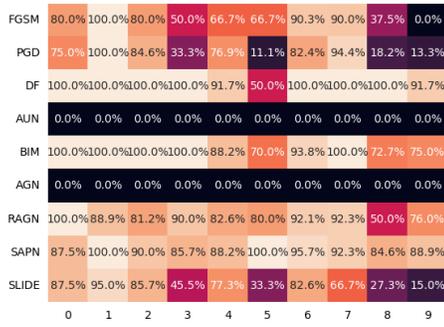
Figure 23: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 2 (generators 4 to 6). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).

Figure 24: Experiment 5 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 3 (generators 7 to 9). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
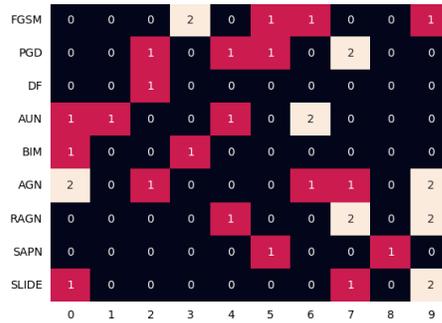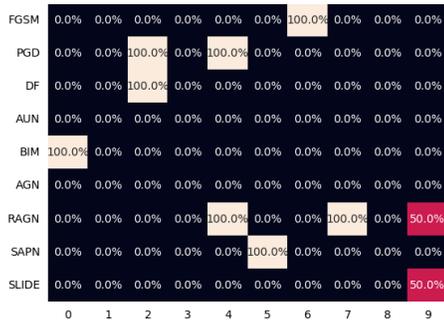
### A.4.3 Experiment 7, trained jointly - faster initialization



Figure 25: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 3 attack/generator setting. Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
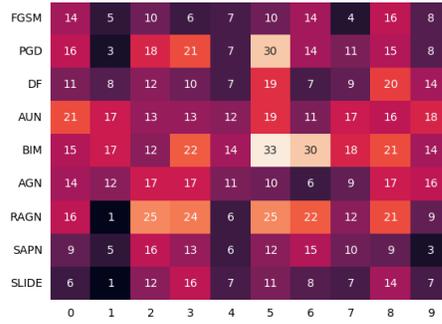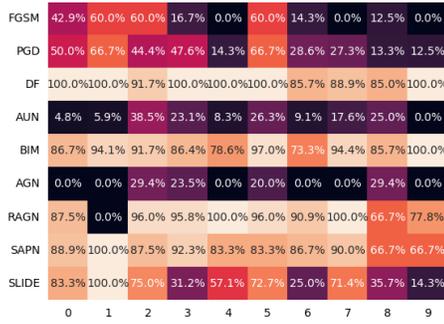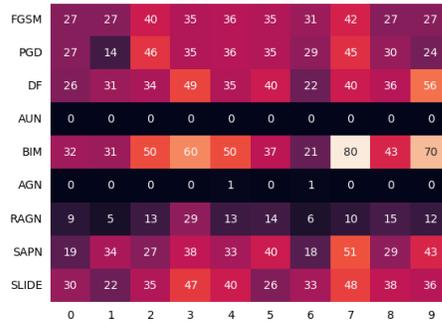
Figure 26: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
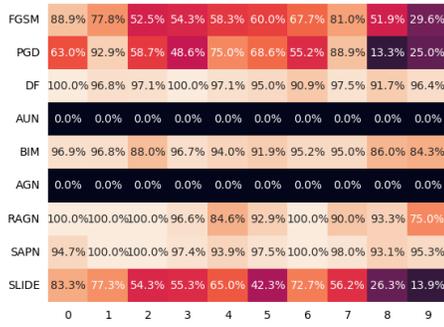
Figure 27: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 5 attack/generator setting, part 2 (generators 4 and 5). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
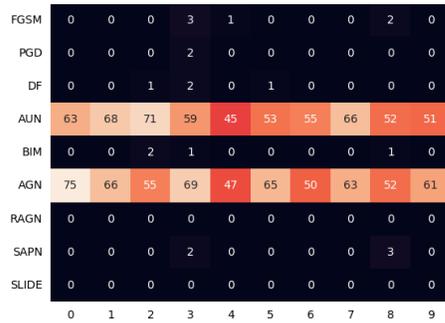
Figure 28: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 1 (generators 1 to 3). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
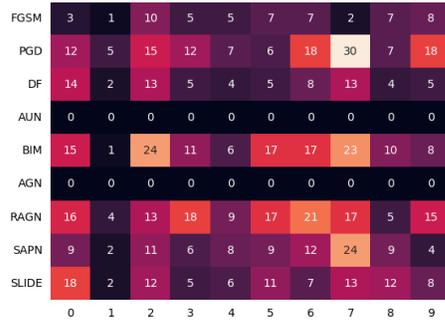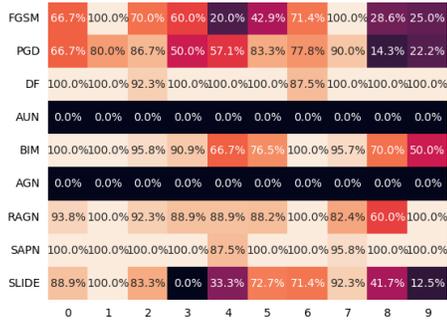
Figure 29: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 2 (generators 4 to 6). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).
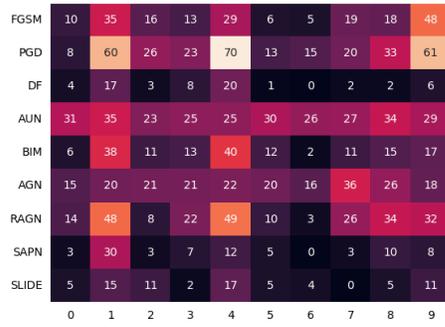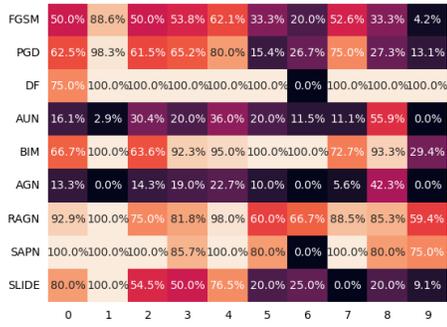
Figure 30: Experiment 7 (trained jointly - faster initialization). Class/attack accuracy breakdown for the 9 attack/generator setting, part 3 (generators 7 to 9). Each row corresponds to one generator. The heatmaps present the accuracy for a given class/attack (left) and the number of samples on which a given generator won (right).