

GradViT: Gradient Inversion of Vision Transformers

Ali Hatamizadeh*, Hongxu Yin*, Holger Roth, Wenqi Li,
Jan Kautz, Daguang Xu†, and Pavlo Molchanov†

NVIDIA

{ahatamizadeh, danny, hroth, wenqi, jkautz, daguangx, pmolchanov}@nvidia.com

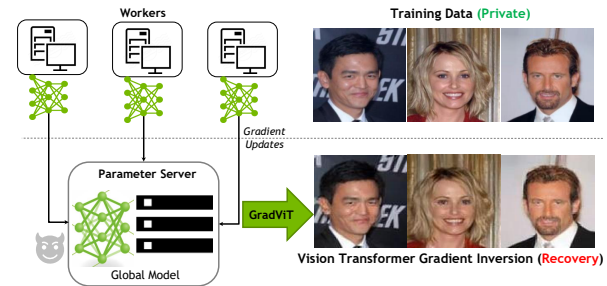
Abstract

In this work we demonstrate the vulnerability of vision transformers (ViTs) to gradient-based inversion attacks. During this attack, the original data batch is reconstructed given model weights and the corresponding gradients. We introduce a method, named GradViT, that optimizes random noise into naturally looking images via an iterative process. The optimization objective consists of (i) a loss on matching the gradients, (ii) image prior in the form of distance to batch-normalization statistics of a pretrained CNN model, and (iii) a total variation regularization on patches to guide correct recovery locations. We propose a unique loss scheduling function to overcome local minima during optimization. We evaluate GradViT on ImageNet1K and MS-Celeb-1M datasets, and observe unprecedentedly high fidelity and closeness to the original (hidden) data. During the analysis we find that vision transformers are significantly more vulnerable than previously studied CNNs due to the presence of the attention mechanism. Our method demonstrates new state-of-the-art results for gradient inversion in both qualitative and quantitative metrics. Project page at <https://gradvit.github.io/>.

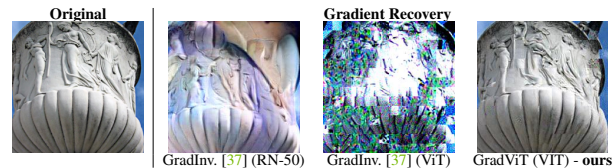
1. Introduction

Vision Transformers (ViTs) [8] have achieved state-of-the-art performance in a number of vision tasks such as image classification [39], object detection [6] and semantic segmentation [5]. In ViT-based models, visual features are split into patches and projected into an embedding space. A series of repeating transformer encoder layers, consisting of alternating Multi-head Self-Attention (MSA) and Multi-Layer Perceptron (MLP) blocks extract feature representation from the embedded tokens for downstream tasks (e.g., classification). Recent studies have demonstrated the effectiveness of ViTs in learning uniform local and global spatial dependencies [31]. In addition, ViTs have a great

*Equal contribution. †Equal advising.



(a) Recovering data from vision transformer gradient unveils intricate details.



(b) GradViT improves noticeably over prior art. Example within a batch of size 8.

Figure 1. Inverting gradients for image recovery. We show vision transformer gradients encode a surprising amount of information such that high-fidelity original image batches of high resolution can be **recovered**, see 112×112 pixel MS-Celeb-1M and 224×224 pixel ImageNet1K sample recovery above and more in experiments. Our method, GradViT, yields the **first** successful attempt to invert ViT gradients, not achievable by previous state-of-the-art methods. We demonstrate that ViTs, despite lacking batchnorm layers, suffer even **more data leakage** compared to CNNs. As insights we show that ViT gradients (i) encode **uneven** original information across layers, and (ii) **attention** is all that reveals.

capability in learning pre-text tasks and can be scaled for distributed, collaborative, or federated learning scenarios. In this work, we study vulnerability of sharing ViT's gradients in the above mentioned settings.

Recent efforts [9, 37, 43] have demonstrated the vulnerability of convolutional neural networks (CNN) to gradient-based inversion attacks. In such attacks, a malicious party can intercept local model gradients and reconstruct private training data in an optimization-based scheme via matching the compromised gradients. Most methods are limited to small image resolutions or non-linearity constraints amid the

hardness of the problem. Among these, GradInversion [37] demonstrated the first successful scaling of gradient inversion to deep networks on large datasets over large batches. In addition to gradient matching, GradInversion [37] is constrained to models with Batch Normalization layers to match feature distribution and bring naturality to the reconstructed images. However, vision transformers lack BN layers and are less vulnerable to previously proposed inversion methods. Naively applying CNN-Based gradient matching [9, 37] techniques for ViT inversion results in sub-optimal solutions due to inherent differences in architectures. Fig. 1 compares reconstruction results obtained by applying current state-of-the-art method GradInversion [37] on the CNN and ViT models. We clearly see significantly degraded visual quality when inverting the ViT gradients.

Since ViT-based models have a different architecture, operate on image patches, and contain no BNs as in CNN counterparts, it might be assumed as if they are more *secure* to gradient-based inversion attacks. On the contrary to this assumption, in this work we quantitatively and qualitatively demonstrate that ViT-based models are even *more vulnerable* than CNNs. To show that, we first study the challenges introduced by ViT’s architectural difference, then propose a novel method, named GradViT, which addresses them and obtains unprecedented high-fidelity and closeness to the original (hidden) data (Fig. 1). Specifically, in GradViT, we tackle the absence of BN statistics by using an independently trained CNN to match the feature distributions of natural images and the images under optimization. We use a ResNet-50 model trained with contrastive loss and its associated BN statistic as an image prior. That is, another model can serve as an image prior instead of the exact BN statistics and their corresponding updates. Moreover, we discover that the proposed image prior generalizes to unseen domain (*e.g.*, faces) which makes it universal.

In addition, while a gradient-based optimization attack can lead to a legitimate reconstruction of patches, their relative location will most likely be incorrect. This happens due to the lack of inductive image bias and permutation invariance in ViTs. To address this problem, we propose a patch prior loss that minimizes the total pixel distances of edges between patches. In other words, we enforce spatial constraints on shared borders (*i.e.*, vertical and horizontal) across neighboring patches as we expect no significant visual discontinuities between them. Minimizing all three losses simultaneously leads to sub-optimal solutions. Therefore, we propose a tailored scheduler to balance the contribution of each loss during training, which is observed to be critical to achieve a valid image recovery.

We validate the effectiveness of GradViT across a wide range of ViT-based models over changing datasets. We start with batch reconstruction of training images from ImageNet1K dataset [7] given the widely used ViT net-

works (*e.g.*, ViT-B/16,32, ViT-S, ViT-T, DeiT, etc.) as the base networks. Our results demonstrate new state-of-the-art benchmarks in terms of image reconstruction metrics. Furthermore, we demonstrate the possibility of detailed recovery of facial images by gradient inversion of a ViT-based model [42] from MS-Celeb-1M dataset [11]. Our findings demonstrate the vulnerabilities of ViT-based models to gradient inversion attacks and specifically for sensitive domains with human training data. With these concerns, we perform extensive studies to analyze the source of vulnerability in ViTs by investigating both layer-wise and component-wise contributions. Our findings provide insights for the development of protection mechanisms against such attacks, which can be beneficial for securing distributed training of ViTs in applications such as multi-node training or federated learning [14, 25].

Our main contributions are summarized as follows:

- We present GradViT, a first successful attempt at ViT gradient inversion, in which random noise is optimized to match shared gradients.
- We introduce an image prior based on CNNs trained with contrastive loss and show scalability across domains.
- We articulate a loss scheduling scheme to guide optimization out of sub-optimal solutions.
- We formulate a patch prior loss function tailored to ViT inversion that mitigates the issue of patch permutation invariance in the reconstructed image.
- We set a state-of-the-art benchmark for ViT gradient inversion across multiple ViT-based networks on ImageNet1K [7] and MS-Celeb-1M [11] datasets. Our method recovers high-resolution facial features with the most intricate details.
- We study the vulnerability of ViT components by performing layer-wise and component-wise analysis. Our findings show that gradients of deeper layers are more informative, and MSA gradients yield near-perfect input recovery.

2. Related Work

Image synthesis. Synthesizing images from neural networks have been a long-lasting important topic for vision, with generative models [10, 17, 27, 29, 40] being at the forefront and yielding state-of-the-art fidelity. However not all networks are equipped with image synthesis capacity as in GANs when pretrained on their target domains, and thus urge alternative methods to generate natural images from normally trained networks. To this end, one stream of work visualizes pretrained networks by analyzing intermediate representations [23, 24, 29, 30], while a more recent stream of method synthesizes natural images from a trained network through auxiliary generative networks [3, 19, 21, 22] or network inversion [2, 28, 32, 38]. The rapid progress of the

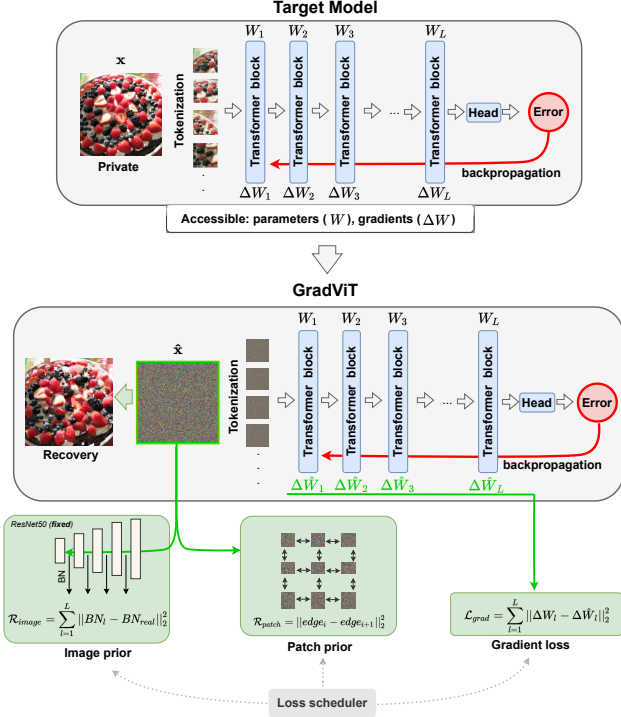


Figure 2. GradViT reconstructs original training data by gradient matching and an image prior network. The batch-wise statistics of synthesized images are matched with BN running statistics of the prior network to enhance fidelity. A loss scheduler balances training of the prior network with gradient matching to avoid local minima during optimization. We also propose a patch prior total variation loss to regularize the position of patches. GradViT is capable of large batch gradient inversion of up to 30 images.

field has shown complete viability to reverse out high fidelity images from deep nets on large-scale datasets with high image resolution. Yet all aforementioned methods reveal only dataset-level distributional prior, as opposed to image-level private visual features that impose privacy concerns.

Gradient inversion. Early efforts [26, 34] investigated the possibility of membership attacks and inferring properties of private training data by exploiting shared gradients. Beyond these membership attacks, Wang *et al.* [36] attempted to reconstruct one image from a client pool of private data using a GAN-based reconstruction model. This work was only evaluated for low-resolution images and a very shallow attack network. Furthermore, Zhu *et al.* [43] demonstrated successful joint image and label restoration by matching the gradients of trainable inputs. As opposed to previous efforts, this work used a relatively deeper CNN architecture [20], however it was still limited to low-resolution images (*e.g.*, CIFAR10) with single training mini-batches and incapable of handling non-continuous activation functions (*e.g.*, ReLU). Geiping *et al.* [9] mitigated this issue using a cosine similarity loss function to match the gradients sign. As a result,

this enabled reconstruction of input training data from more commonly-used networks such as ResNet-18 using higher resolution images (*e.g.*, ImageNet), but it only produces a single image. Yin *et al.* [37] introduced the GradInversion model that scales the attack to larger mini-batches with high-resolution ImageNet samples from a deep ResNet-50 network. In addition to gradient matching, GradInversion proposes to match the distribution of running mean and variance of batch normalization layers that are produced from a synthesized input image, augmented by multi-agent group consistency. Considering the prevalence of batch normalization in CNN-based architectures and the associated strong prior in running statistics, GradInversion significantly improves the fidelity of reconstructed images. Despite recent generative prior augmentation [15] and theoretical insights [16], gradient inversion attacks remain valid only for CNNs, with key assumptions nonexistent in ViTs.

3. GradViT

We next describe our proposed methodology in details. Fig. 2 illustrates an overview of the GradViT framework. Our inversion task is formulated as an optimization problem. Given randomly initialized input tensor $\hat{\mathbf{x}} \in \mathbb{R}^{N \times H \times W \times C}$ (N, H, W, C being batch size, height, width and number of channels) and a target network with weights \mathbf{W} and gradient updates $\Delta \mathbf{W}$ averaged over a mini-batch, GradViT recovers original image batch via the following optimization:

$$\hat{\mathbf{x}}^* = \underset{\hat{\mathbf{x}}}{\operatorname{argmin}} \Gamma(\mathbf{t}) \mathcal{L}_{\text{grad}}(\hat{\mathbf{x}}; \Delta \mathbf{W}) + \Upsilon(\mathbf{t}) \mathcal{R}_{\text{image}}(\hat{\mathbf{x}}) + \mathcal{R}_{\text{aux}}(\hat{\mathbf{x}}), \quad (1)$$

in which $\mathcal{L}_{\text{grad}}$ is a gradient matching loss, $\mathcal{R}_{\text{image}}$ and \mathcal{R}_{aux} are an image prior and auxiliary regularization. $\Gamma(\mathbf{t})$ and $\Upsilon(\mathbf{t})$ denote loss scheduler functions that balance contributions on the total loss at each training iteration \mathbf{t} . We solve the proposed optimization problem in an iterative manner. $\mathcal{L}_{\text{grad}}$ acts as a main force to reduce the error between the shared model’s gradients and the computed one, while other losses improve fidelity of the recovered images.

3.1. Gradient Matching

Gradient matching relies on valid target labels for simulation of gradients given synthesized inputs. Akin to [37], we first recover the labels through the negative sign traces of the gradient in the classification head, resulting in a label set $\hat{\mathbf{y}}$ for a batch size of N as

$$\hat{\mathbf{y}} = \operatorname{arg sort}(\min_i \nabla_{\mathbf{W}_{i,j}^{(\text{CLS})}} \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*))[: N], \quad (2)$$

in which $\Delta \mathbf{W}^{(\text{CLS})}$ denotes the gradient of the classification head of ViT, and \mathbf{x}^* and \mathbf{y}^* represent original training images and labels, respectively. Once the labels are restored, the ℓ_2 norm between the gradients from the synthesized

inputs and shared gradients are minimized according to

$$\mathcal{L}_{\text{grad}}(\hat{\mathbf{x}}; \Delta \mathbf{W}) = \sum_l \left\| \nabla_{\mathbf{W}^{(l)}} \mathcal{L}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) - \underbrace{\nabla_{\mathbf{W}^{(l)}} \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*)}_{\text{given batch gradient}} \right\|_2, \quad (3)$$

where $\nabla_{\mathbf{W}^{(l)}} \mathcal{L}(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is calculated based on the synthesized inputs $\hat{\mathbf{x}}$ and restored label $\hat{\mathbf{y}}$ in each layer l of the network.

3.2. Image Prior

As an image prior counterpart to guide the optimization process towards image naturalness, we look into auxiliary networks, such as CNN, to impose an image prior. In this paper, we use a self-supervised trained MOCO V2 ResNet-50 via contrastive loss [4, 12] for this task, which we observe scales across varying domains. More specifically, we use the stored BN statistics of the feature maps as a target distribution for the estimated per-layer statistics when passing the synthesized inputs through the network.

Given the batch-wise mean $\mu_l(\hat{\mathbf{x}})$ and variance $\sigma_l^2(\hat{\mathbf{x}})$ of synthesized inputs at layer l , the following image prior loss is minimized

$$\mathcal{R}_{\text{image}}(\hat{\mathbf{x}}) = \sum_l \left\| \mu_l(\hat{\mathbf{x}}) - \mu_{l,\text{BN}} \right\|_2 + \sum_l \left\| \sigma_l^2(\hat{\mathbf{x}}) - \sigma_{l,\text{BN}}^2 \right\|_2, \quad (4)$$

where $\mu_{l,\text{BN}}$ and $\sigma_{l,\text{BN}}^2$ denote the running mean and variance of the CNN prior across layers $l = 1, 2, \dots, L$. By aligning batch-wise and running statistics, the loss significantly enhances the image fidelity and visual realism, as we show later.

3.3. Loss Scheduler

Balancing losses is vital for ViT gradient inversion to yield valid recovery. In the early training stages, the gradient loss is very sensitive to abrupt changes in the pixel-wise values of the synthesized inputs. As a result, we observe an early stage minimization of both the gradient and image prior losses results in convergence to sub-optimal solutions. As a remedy, we activate the image prior loss only after the first half of training where the synthesized inputs are close to optimum for the gradient loss. Then, we reduce the contribution of the gradient loss to half for the rest of the training to allow for more effective prior extraction. For a total of \mathbf{T} training iterations, loss schedulers $\Gamma(\mathbf{t})$ and $\Upsilon(\mathbf{t})$ at iteration \mathbf{t} are defined as

$$\Gamma(\mathbf{t}) = \begin{cases} \alpha_{\text{grad}} & 0 < \mathbf{t} \leq \frac{\mathbf{T}}{2} \\ \frac{1}{2} \alpha_{\text{grad}} & \frac{\mathbf{T}}{2} < \mathbf{t} \leq \mathbf{T} \end{cases}, \quad (5)$$

$$\Upsilon(\mathbf{t}) = \begin{cases} 0 & 0 < \mathbf{t} \leq \frac{\mathbf{T}}{2} \\ \alpha_{\text{image}} & \frac{\mathbf{T}}{2} < \mathbf{t} \leq \mathbf{T} \end{cases}. \quad (6)$$

α_{grad} and α_{image} denote gradient and BN matching scale factors, respectively. We observe this scheduling is key to valid recovery, as shown in the ablations later.

3.4. Auxiliary Regularization

We also explore an extensive set of auxiliary image priors to govern image fidelity. Our auxiliary regularization loss consists of (i) a novel patch prior loss to regularize the permutation ordering of reconstructed patches, (ii) a registration loss to ensure consistency among final reconstructions of different optimization seeds and (iii) an image prior loss to improve the image quality:

$$\mathcal{R}_{\text{aux}}(\hat{\mathbf{x}}) = \alpha_1 \mathcal{R}_{\text{patch}}(\hat{\mathbf{x}}) + \alpha_2 \mathcal{R}_{\text{reg}}(\hat{\mathbf{x}}) + \alpha_3 \mathcal{R}_{\text{prior}}(\hat{\mathbf{x}}). \quad (7)$$

We next elaborate on each of the loss terms.

3.4.1 Patch Prior

As opposed to typical CNN-based networks, ViT-based models are permutation-invariant and lack inherent inductive image biases. The patch-based strategy that is used for feature extraction in ViTs greatly manifests itself during the inversion process in our GradViT, as several permutations of the same group of reconstructed patches can equally satisfy the minimization process. Hence, the reconstructed images suffer from an incorrect order of patches.

To mitigate this issue, we propose a new patch prior loss that enforces similarity between horizontal and vertical joints of adjacent patches. The main idea is that even though image tokens are regarded as separate entities when fed into transformers to learn attention, their associated patches are bonded spatially by nature - they have to form one single image when put next to each other. As a result, pixel values among adjacent patch edges shall be in similar ranges, and abrupt changes shall be penalized.

By assuming a patch size of $P \times P$ from an image of $H \times W$, our patch prior regularizes spatial positioning of neighboring patches by enforcing

$$\mathcal{R}_{\text{patch}}(\hat{\mathbf{x}}) = \sum_{k=1}^{\frac{H}{P}-1} \left\| \hat{\mathbf{x}}[:, P \cdot k, :, :] - \hat{\mathbf{x}}[:, P \cdot k - 1, :, :] \right\|_2 + \sum_{k=1}^{\frac{W}{P}-1} \left\| \hat{\mathbf{x}}[:, :, P \cdot k, :] - \hat{\mathbf{x}}[:, :, P \cdot k - 1, :] \right\|_2. \quad (8)$$

Our ablation studies demonstrate the effectiveness of the patch prior loss in enhancing the ordering of reconstructed patches. In other words, forcing patch boundaries to be smooth in color indirectly forces the optimizer to re-distribute the patches, such that the loss can be further reduced.

3.4.2 Registration

In the proposed framework, the final solution of reconstructed images depend on the optimization initialization (*i.e.*, randomly selected seeds). As a result, reconstructions with different image semantics and viewpoints may be produced. Inspired by Yin *et al.* [37], we also regularize the reconstruction of different seeds by aligning them with a consensus solution across all optimizations. Considering $\hat{\mathbf{x}}_{\mathcal{S}} = [\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_s]$ to represent all viable solutions for each optimization round, we first compute a consensus solution $\hat{\mathbf{x}}_{\mathcal{M}} = \frac{1}{|\hat{\mathbf{x}}_{\mathcal{S}}|} \sum_{\mathbf{s}} \hat{\mathbf{x}}_{\mathbf{s}}$ by pixel-wise averaging of all solutions. We perform an initial coarse alignment, by using a RANSAC-Flow based image alignment strategy [33], for each solution with respect to $\hat{\mathbf{x}}_{\mathcal{M}}$ as the target and obtain the final consensus solution $\hat{\mathbf{x}}_{\mathcal{C}}$ by averaging all registered inputs as in

$$\hat{\mathbf{x}}_{\mathcal{C}} = \frac{1}{|\hat{\mathbf{x}}_{\mathcal{S}}|} \sum_{\mathbf{s}} \mathbf{F}_{\hat{\mathbf{x}}_{\mathbf{s}} \rightarrow \hat{\mathbf{x}}_{\mathcal{M}}}(\hat{\mathbf{x}}_{\mathbf{s}}), \quad (9)$$

in which $\mathbf{F}_{\hat{\mathbf{x}}_{\mathbf{s}} \rightarrow \hat{\mathbf{x}}_{\mathcal{M}}}$ is a flow function for mapping the source candidate $\hat{\mathbf{x}}_{\mathbf{s}}$ to target $\hat{\mathbf{x}}_{\mathcal{M}}$. We minimize the ℓ_2 distance of all solutions with respect to the final consensus solution:

$$\mathcal{R}_{\text{reg}}(\hat{\mathbf{x}}) = \|\hat{\mathbf{x}} - \hat{\mathbf{x}}_{\mathcal{C}}\|_2. \quad (10)$$

3.4.3 Extra Priors

As a final step, we leverage extra conventional image prior losses [38] including ℓ_2 and total variation to improve the quality of reconstructions losses as:

$$\mathcal{R}_{\text{prior}}(\hat{\mathbf{x}}) = \mathcal{R}_{\ell_2}(\hat{\mathbf{x}}) + \mathcal{R}_{\text{TV}}(\hat{\mathbf{x}}). \quad (11)$$

At this stage, all three image regularization terms are balanced using scaling constants $\alpha_{1,2,3}$ in Eqn. 7, and then summed into gradient matching and image prior losses for input updates.

4. Experiments

4.1. Datasets

We next validate the effectiveness of our approach on the ImageNet1K [7] and MS-Celeb-1M datasets [11] for the task of image classification and face recognition, respectively. In addition to ImageNet1K as a widely adopted benchmarking task, the latter was chosen to demonstrate the risks of gradient inversion data leakage from a sensitive domain with considerable security concerns. For ImageNet1K experiments, we use images of resolution 224×224 px, whereas we resize MS-Celeb-1M images to 112×112 px amid network input requirement of [42].

4.2. Evaluation Metrics

To make our comparisons comprehensive, we report quantitative measurements in addition to qualitative results throughout our experiments. We adopt the commonly-used image quality metrics including (i) Peak Signal-to-Noise Ratio (PSNR), (ii) Learned Perceptual Image Patch Similarity (LPIPS) [41] and (iii) cosine similarity in the Fourier space (FFT_{2D}) to measure the similarity between the image recovery and original counterparts.

4.3. Implementation Details

We explore different variations of the ViT [8] and DeiT [35] models. MOCO V2-pretrained ResNet-50 model [4, 12] is used for all CNN experiments as a base for the image prior in GradViT. For the MS-Celeb-1M dataset, we use the FaceTransformer [42] that is a modified ViT model. We use an Adam optimizer [18] with an initial learning rate of 0.1 for 120K iterations and with cosine learning rate decay. For all experiments, we use an NVIDIA DGX-1 server and reconstruct the training images by only exploiting the shared gradients, using a mini batch size of 8 unless specified otherwise. We use $\alpha_{\text{grad}} = 4 \times 10^{-3}$, $\alpha_{\text{image}} = 2 \times 10^{-1}$, $\alpha_1 = 10^{-4}$, $\alpha_2 = 10^{-2}$ and $\alpha_3 = 10^{-4}$ as the scaling coefficients in the loss functions. According to the proposed loss scheduler as described in Sec. 3.3, we first start the optimization process with only the gradient matching loss for 60K iterations, and then decrease α_{grad} to 2×10^{-3} , jointly with adding the image prior loss.

5. Results

5.1. ImageNet1K

Table 1 presents quantitative comparisons between our method and the state-of-the-art benchmarks for batch gradient inversion on ImageNet1K, with Fig. 3 depicting our main qualitative results. GradViT is used for inversion of variants of ViT and DeiT models towards a target batch of size 8. Gradient inversion reconstructions of ViT-B/16 using GradViT outperform the previous state-of-the-art benchmarks (*i.e.*, ResNet-50 with GradInversion) by a large margin in terms of all image quality metrics. Applying GradInversion to ViTs results in unsatisfactory results. GradViT, for the first time, enables a viable, complete recovery of original images. More surprisingly, it yields unprecedented image realism and intricate original details that surpass even the best recovery from ResNet-50 using CNN-tailored GradInversion. This sets a new benchmark for gradient inversion on ImageNet1K.

5.2. MS-Celeb-1M

Fig. 4 shows the performance of GradViT on FaceTransformer [42]. We observe that GradViT recovers a substantial amount of original information, including face, hair, cloth-

| Gradient Inversion Method | Network | Image Reconstruction Metrics | | | Considerations | | |
|----------------------------|----------------|------------------------------|--------------------------------|--------------------|----------------|----------------------|-----------|
| | | PSNR \uparrow | FFT _{2D} \downarrow | LPIPS \downarrow | Type | Need Original Labels | GAN-Based |
| Random Noise | - | 1.351 | 0.706 | 9.964 | - | No | No |
| Latent projection [17] | BigGAN [1] | 10.149 | 0.275 | 0.722 | CNN | Yes | Yes |
| DeepInversion [38] | ResNet-50 [13] | 10.131 | 0.238 | 0.728 | CNN | Yes | No |
| Deep Gradient Leakage [43] | ResNet-50 [13] | 10.252 | 1.319 | 0.602 | CNN | No | No |
| Inverting Gradients [9] | ResNet-50 [13] | 11.703 | 0.355 | 0.749 | CNN | Yes | No |
| GradInversion [37] | ResNet-50 [13] | 12.929 | 0.175 | 0.484 | CNN | No | No |
| GradInversion [37] | ViT-B/16 [8] | 10.824 | 0.116 | 0.708 | ViT | No | No |
| GradViT | ResNet-50 [13] | 11.635 | 0.076 | 0.454 | CNN | No | No |
| GradViT | DeiT-B/16 [8] | 13.252 | 0.058 | 0.413 | ViT | No | No |
| GradViT | ViT-B/16 [8] | 15.515 | 0.032 | 0.295 | ViT | No | No |

Table 1. Quantitative comparisons of image reconstruction quality from batch of 8 images in ImageNet1K dataset. CNN-based networks use ResNet-50 for gradient inversion in line with prior work. GradViT outperforms all prior approaches by a large margin across image quality metrics.

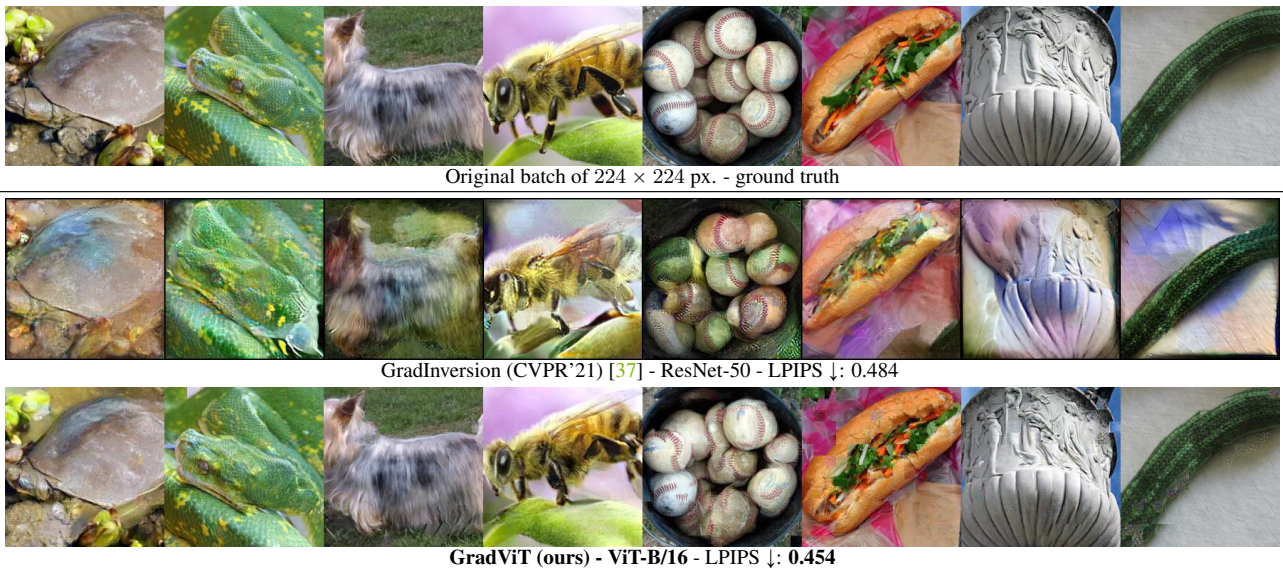


Figure 3. Qualitative comparisons of reconstructed images from batch of 8 images in ImageNet1K dataset using the proposed GradViT and state-of-the-art GradInversion [37]. GradViT outperforms GradInversion both qualitatively and quantitatively. It recovers the most intricate details, of very high image fidelity and naturalness, not only for the target objects, but also all the background scenes. Best viewed in color.

ing, and even background details close to the original images. These results demonstrate the vulnerabilities of ViTs under gradient inversion attacks, in a sensitive domain such as face recognition. Here a leakage of private data can lead to significant security concerns.

6. Analysis

6.1. Ablation Study

Table 2 provides both (i) quantitative comparisons to ablate the effectiveness of each training loss term on recovery quality and (ii) the associated qualitative comparisons. We observe that optimizing $\mathcal{L}_{\text{grad}} + \mathcal{R}_{\text{reg}}$ as in DeepInversion [37] restores certain features of the original training images. However, the reconstructions suffer from poor image fidelity and

loss of detailed semantics. Furthermore, naively optimizing the image prior loss \mathcal{R}_{D} results in sub-optimal solutions. Adding the scheduler alleviates this issue and results in substantially improved reconstructions. Adding the patch prior loss $\mathcal{R}_{\text{patch}}$ guides the location of recovered patches and significantly enhances the image quality. Please see supplementary materials for visualizations of synthesized images in various stages of training.

6.2. Varying Architecture & Patch Size

Table 3 shows the performance of GradViT given varying architectures and changing patch sizes. We observe that transformers with (i) a smaller patch size, (ii) more parameters, and (iii) stronger training recipe with distillation, reveal more original information and hence are more vulnerable in

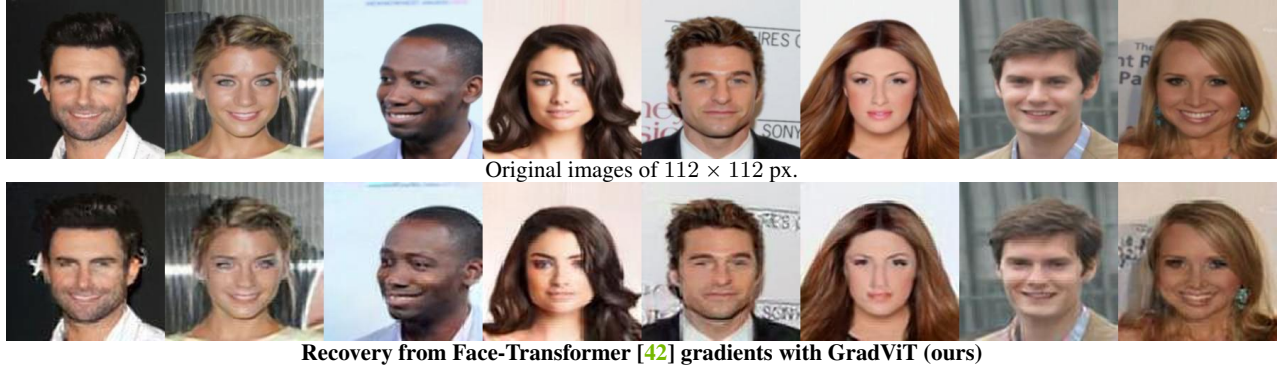


Figure 4. Qualitative comparison of reconstructed images from MS-Celeb-1M dataset using batch gradient inversion of Face-Transformer [42]. GradViT is able to recover detailed and facial features identical to the original. Recovery at batch size 4. Best viewed in color.

| Loss Function | $\mathcal{L}_{\text{grad}}$ | Image Reconstruction Metric | | |
|---|-----------------------------|-----------------------------|--------------------------------|--------------------|
| | | PSNR \uparrow | FFT _{2D} \downarrow | LPIPS \downarrow |
| Random | 8.143 | 0.706 | 9.964 | 1.351 |
| $\mathcal{L}_{\text{grad}} + \mathcal{R}_{\text{reg}}$ [37] | 4.190 | 11.431 | 0.071 | 0.498 |
| + $\mathcal{R}_{\text{image}}$ | 3.127 | 11.291 | 0.078 | 0.504 |
| + $\Gamma(\cdot), \Upsilon(\cdot)$ | 3.047 | 13.404 | 0.049 | 0.412 |
| + $\mathcal{R}_{\text{patch}}$ | 2.326 | 15.515 | 0.032 | 0.295 |

Table 2. Effect of each loss term on reconstruction quality of final synthesized images. Results presented among a batch of 8 images with total variation and ℓ_2 priors included by default in all runs.

| Network | Distilled | Image Reconstruction Metric | | |
|----------------|-----------|-----------------------------|--------------------------------|--------------------|
| | | PSNR \uparrow | FFT _{2D} \downarrow | LPIPS \downarrow |
| DeiT-T/16 [35] | No | 12.243 | 0.079 | 0.489 |
| DeiT-T/16 [35] | Yes | 13.212 | 0.076 | 0.454 |
| DeiT-S/16 [35] | No | 12.664 | 0.059 | 0.461 |
| DeiT-S/16 [35] | Yes | 13.092 | 0.055 | 0.419 |
| DeiT-B/16 [35] | No | 13.252 | 0.058 | 0.413 |
| DeiT-B/16 [35] | Yes | 13.708 | 0.041 | 0.407 |
| ViT-T/16 [8] | - | 12.521 | 0.062 | 0.483 |
| ViT-S/32 [8] | - | 12.365 | 0.063 | 0.505 |
| ViT-S/16 [8] | - | 13.658 | 0.042 | 0.412 |
| ViT-B/32 [8] | - | 13.599 | 0.048 | 0.436 |
| ViT-B/16 [8] | - | 15.515 | 0.032 | 0.295 |

Table 3. Quantitative comparisons of image reconstruction quality from gradient inversion of various ViT and DeiT models on ImageNet1K.

gradient inversion attacks. In addition, we observe more vulnerabilities in ViTs in terms of revealing more information than their counterpart DeITs.



Figure 5. Effect of increasing batch size on the quality of image recovery. ImageNet and MS-Celeb-1M images are reconstructed in 224×224 px and 112×112 px respectively. Representative sample reconstructions are presented for batch sizes of 8, 16, 30 and 48. The maximum number of batch sizes is limited to 30 for ImageNet dataset amid GPU memory constraint.

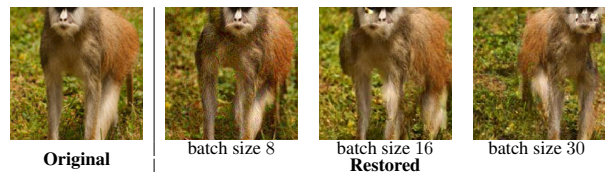


Figure 6. Visual comparison of reconstruction quality with different batch sizes on ImageNet. Although GradViT recovers major visual features, the quality decreases with increasing batch size.

6.3. Increasing the Batch Size

In Fig. 5, we study the effect of batch size on reconstruction image quality as gradients are averaged over a larger number of images. Considering GPU memory constraints, we experimented with maximum batch sizes of 30 and 64 for ImageNet1K and MS-Celeb-1M datasets, respectively. In both datasets, we observe that image quality degrades, as expected, at a larger batch size. For facial recovery, Grad-

ViT is still able to recover identifiable images even at the batch size of 30 (see examples in Fig. 5). In the Appendix we will also study the likelihood of person identification as a function of the batch size, and the potential of auxiliary GANs to improve fidelity. We also observe a similar trend on ImageNet1K, as shown in Fig. 6. Reconstruction at a batch size of 30 still reveals major visual features.

6.4. Tracing the Source

To give guidance on future defense regimes, we delve deep into tracing the source of information leakage – among all shared gradients, (i) *where* among all layers and (ii) *what* exact components leak the most original information?

Answering these questions are key to targeted protections for enhanced security. As an attempt, we ablate the contributions of gradients from varying ViT architecture sections to input recovery. More specifically, we conduct two streams of analysis. *Layer-wise*, we study the changing effects of removing gradient contributions from transformer layers of different depths. This hints at the possibility to share gradients separately as a remedy to prevent an overall inversion. *Component-wise*, we retrain by using gradients from either MSA or MLPs across all the layers in the target model, and analyze the strength of their links to original images. This gives insights on what exact transformation retains the most information. We base both analysis on ViT-B/16 and present our findings next.

6.4.1 Later Stages Reveal More

More specifically, we remove gradients from initial, middle, and later stages to ablate the impacts on recovery efficacy. To this end, we reconstruct images without including the gradients of layers 1 – 4, 5 – 8 and 9 – 12. Table 4 shows that reconstructions by excluding the gradients of earlier layers are more accurate than those of deeper layers, whereas dropping the later stage alters the recovery the most. In other words, gradients of deeper layers are more informative for inversion - see Fig. 7(a) for qualitative comparisons.

6.4.2 Attention is All That Reveals

We next perform two component-wise data leakage studies on the ViT-B/16 model by only utilizing the gradients of MLP or MSA blocks for the inversion attacks. We present results with Table 5 and Fig. 7(b). Table 5 demonstrates the importance of MSA gradients, as its reconstructions have significantly better image quality than images synthesized by MLP gradients. As illustrated in Fig. 7(b), reconstructions from MLP gradients lack important details, whereas utilizing the gradients of MSA layers alone can already yield high-quality reconstructions.

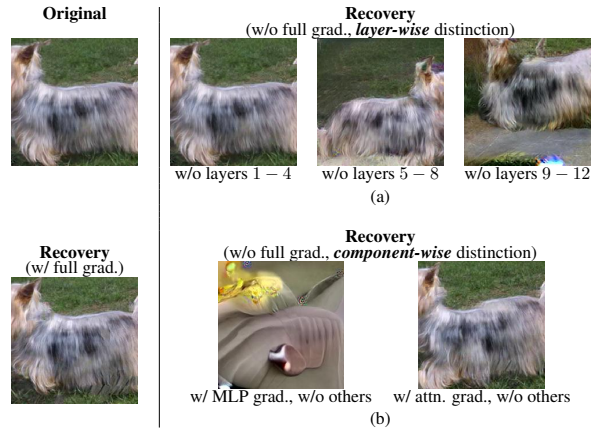


Figure 7. Reconstructed images from layer-wise and component-wise ablation studies using a batch size of 8. Later layers (9-12) contain the most critical information that leads to data leakage. The component-wise studies show that gradients of MSA blocks have more critical information than those of MLP blocks. See supplementary materials for more visualizations.

| Layer-wise Gradients | Image Reconstruction Metric | | |
|----------------------|-----------------------------|--------------------------------|--------------------|
| | PSNR \uparrow | FFT _{2D} \downarrow | LPIPS \downarrow |
| All (baseline) | 15.515 | 0.032 | 0.295 |
| w/o Layers 1-4 | 13.982 | 0.047 | 0.412 |
| w/o Layers 5-8 | 11.086 | 0.086 | 0.555 |
| w/o Layers 9-12 | 10.284 | 0.091 | 0.598 |

Table 4. Effect of layer-wise gradients on ViT-B/16 reconstructions.

| Component-wise Gradients | Image Reconstruction Metric | | |
|--------------------------|-----------------------------|--------------------------------|--------------------|
| | PSNR \uparrow | FFT _{2D} \downarrow | LPIPS \downarrow |
| All (baseline) | 15.515 | 0.032 | 0.295 |
| w/ MLP, w/o others | 12.256 | 0.066 | 0.568 |
| w/ MSA, w/o others | 13.559 | 0.047 | 0.408 |

Table 5. Effect of component-wise gradients on ViT-B/16 reconstructions.

7. Conclusion

In this work, we have introduced a methodology for gradient inversion of ViT-based models via (i) enforcing the matching of gradients to the shared target (ii) leveraging an image prior, (iii) and utilizing a novel patch prior loss to guide patch recovery locations. Through extensive analysis on ImageNet1K and MS-Celeb-1M datasets, we have shown state-of-the-art benchmarks for gradient inversion of deep neural networks. We have also conducted additional analysis to offer insights to the community and guide designs of ViT security mechanisms to prevent inversions that are shown even stronger than on CNNs. Homomorphic encryption and differential privacy have been shown effective against CNN-based gradient inversion attacks. However, future work is needed to study protection mechanisms against GradViT.

References

- [1] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019.
- [2] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. ZeroQ: A novel zero shot quantization framework. In *CVPR*, 2020.
- [3] Hanting Chen, Yunhe Wang, Chang Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chunjing Xu, Chao Xu, and Qi Tian. Data-free learning of student networks. In *ICCV*, 2019.
- [4] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [5] Bowen Cheng, Alexander G Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *arXiv preprint arXiv:2107.06278*, 2021.
- [6] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *CVPR*, pages 1601–1610, 2021.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [9] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients—How easy is it to break privacy in federated learning? In *NeurIPS*, 2020.
- [10] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of Wasserstein GANs. In *NeurIPS*, 2017.
- [11] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European conference on computer vision*, pages 87–102. Springer, 2016.
- [12] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [14] Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, and Sanjeev Arora. Evaluating gradient inversion attacks and defenses in federated learning. *NeurIPS*, 2021.
- [15] Jinwoo Jeon, Jaechang Kim, Kangwook Lee, Sewoong Oh, and Jungseul Ok. Gradient inversion with generative image prior. *arXiv preprint arXiv:2110.14962*, 2021.
- [16] Xiao Jin, Pin-Yu Chen, Chia-Yi Hsu, Chia-Mu Yu, and Tianyi Chen. CAFE: Catastrophic data leakage in vertical federated learning. *arXiv preprint arXiv:2110.15122*, 2021.
- [17] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Vinod K Kurmi, Venkatesh K Subramanian, and Vinay P Namboodiri. Domain impression: A source data free domain adaptation method. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 615–625, 2021.
- [20] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [21] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. In *CVPR*, pages 1215–1224, 2021.
- [22] Liangchen Luo, Mark Sandler, Zi Lin, Andrey Zhmoginov, and Andrew Howard. Large-scale generative data-free distillation. *arXiv preprint arXiv:2012.05578*, 2020.
- [23] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, 2015.
- [24] Aravindh Mahendran and Andrea Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *IJCV*, 2016.
- [25] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.
- [26] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *IEEE Symp. Security and Privacy (SP)*, 2019.
- [27] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.
- [28] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Inceptionism: Going deeper into neural networks. <https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>, 2015.
- [29] Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. In *CVPR*, 2017.
- [30] Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *NeurIPS*, 2016.
- [31] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *arXiv preprint arXiv:2108.08810*, 2021.
- [32] Shibani Santurkar, Andrew Ilyas, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Image synthesis with a single (robust) classifier. In *NeurIPS*, 2019.
- [33] Xi Shen, François Darmon, Alexei A Efros, and Mathieu Aubry. RANSAC-Flow: Generic two-stage image alignment. In *ECCV*, 2020.

- [34] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE Symp. Security and Privacy (SP)*, 2017.
- [35] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [36] Zhibo Wang, Mengkai Song, Zhifei Zhang, Yang Song, Qian Wang, and Hairong Qi. Beyond inferring class representatives: User-level privacy leakage from federated learning. In *INFOCOM*, 2019.
- [37] Hongxu Yin, Arun Mallya, Arash Vahdat, Jose M Alvarez, Jan Kautz, and Pavlo Molchanov. See through gradients: Image batch recovery via gradinversion. In *CVPR*, 2021.
- [38] Hongxu Yin, Pavlo Molchanov, Jose M. Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via Deep-Inversion. In *CVPR*, 2020.
- [39] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. *arXiv preprint arXiv:2106.04560*, 2021.
- [40] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, 2019.
- [41] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [42] Yaoyao Zhong and Weihong Deng. Face transformer for recognition. *arXiv preprint arXiv:2103.14803*, 2021.
- [43] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *NeurIPS*, 2019.