

Are we pruning the correct channels in image-to-image RIKEN translation models? (No, but we correct it.)

Yiyong Li¹, Zhun Sun^{1*}, Lichao² ^{*}corresponding author, zhunsun@gmail.com ¹BIGO Ltd, ²Tohoku University, ³AIP, RIKEN



(6)

TL;DR

- The commonly employed lasso-based channel regularization approach prunes channels with large visual effects incorrectly.
- We build a novel perturbation model to analyze what channels should be pruned for the instance normalization (IN)-based models.
- Using the perturbation bound, we achieve better compression performance in both on-training and zero-shot scenarios.





Step 2: Update Re-scalar as Pruning.

We conduct the pruning by updating $\gamma_{[I],i}^{(t)} = \beta_{[I],i}^{(t)} = 0$, if one of the following four conditions is satisfied, we increase the threshold ρ_1, ρ_2 after the numbers of survival channels stop decreasing during training.

$$\begin{array}{ll} \text{(i)} & \beta_{[I],i}^{(t-1)} \leq -\tau_{[I],i}^{(t-1)}, & \text{(ii)} & \gamma_{[I],i}^{(t-1)} = \mathbf{0}, \\ \\ \text{(iii)} & \frac{\sum_{j \in [D_I]} F_{i,j} \left(\mathcal{W}_{[I]}^{(t-1)}, \gamma_{[I]}^{(t-1)}, \mathbf{0} \right)}{\sum_{i \in [C_I]} P_{I,i} \left(\mathcal{W}_{[I]}^{(t-1)}, \gamma_{[I]}^{(t-1)}, \beta_{[I]}^{(t-1)} \right)} < \rho_1, \\ \\ \text{(iv)} & \frac{\sum_{j \in [D_I]} F_{i,j} \left(\mathcal{W}_{[I]}^{(t-1)}, \gamma_{[I]}^{(t-1)}, \beta_{[I]}^{(t-1)} \right)}{\sum_{j \in [D_I]} F_{i,j} \left(\mathcal{W}_{[I]}^{(t-1)}, \gamma_{[I]}^{(t-1)}, \beta_{[I]}^{(t-1)} \right)} \leq \rho_2 \end{array}$$

(a) Input (b) un-pruned (c) eb=0.194 (d) eb=1.604 (e) eb=2.430

Fig. 1: Examples of horse2zebra generation results using a pre-trained Cycle-GAN with specific channels pruned. It is conspicuous that the appearances alter more severely as the perturbation error bound (eb) of the channel grows.

Background

- On edge devices and mobile applications, image-to-image translation GAN are commonly employed for visual effects.
- Obtaining compact architectures with pruning and distillation is demanded, due to the limited computational resources.
- We focus on architectures with stacked Conv-IN-ReLU layers that are edge-friendly.
- Pruning methods that are proposed for general neural networks do not fit the image-toimage translation GAN, since the visual ef- The basic building block to



$\sum_{i \in [C_{l}]} P_{l,i} \left(\mathcal{W}_{[l]}^{(t-1)}, \gamma_{[l]}^{(t-1)}, \beta_{[l]}^{(t-1)} \right) < \rho_{2}.$

Results



Fig. 2: Examples of CycleGAN [1] results on horse2zebra **train** set. The methods and their FLOPs (in G) are annotated above the images. Each row shows a different sample. Our pruned model generates samples that have similar stripes to the one generated by CycleGAN.

horse2zebra			
Model	FLOPs	FID	# Pruned

summer2winter								
Model	FLOPs	FID	# Pruned					

fects are more sensitive than other tasks. construct GANs.

Perturbation Analysis

Proposition (Perturbation error bound) 1 Assume $\mathcal{Z} \in \mathbb{R}^{D \times W \times H}$ to be the output of conv, and $\mathcal{Z}^{\overline{i}}$ to denote the result by pruning the *i*th channel of the input. then the norm of perturbation $\Delta_i = \mathcal{Z} - \mathcal{Z}^{\overline{i}}$ is bounded by the following conditions: if $\gamma_i = 0$, then $\|\Delta_i\|_{\ell_1} = 0$; otherwise, we have

$$\|\Delta_{i}\|_{\ell_{1}} \leq \begin{cases} WH \sum_{j \in [D]} F_{i,j}(\mathcal{W}, \gamma, \mathbf{0}), & \beta_{i} \geq \tau_{i} \\ WH \sum_{j \in [D]} F_{i,j}(\mathcal{W}, \gamma, \beta), & |\beta_{i}| < \tau_{i} \\ 0, & otherwise \end{cases}$$
(1)

where $\mathbf{0} \in \mathbb{R}^{C}$ and $\tau_{i} = \sqrt{WH}|\gamma_{i}|$ for all $i \in [C]$, $F_{i,j}$ is the Sensitivity Measurement defined as

$$F_{i,j}(\mathcal{W},\gamma,\beta) = \sqrt{WH} |\gamma_i| \sqrt{\sum_{g,p\in[K]} \mathcal{W}(i,j,g,p)^2} + |\beta_i| \left| \sum_{g,p\in[K]} \mathcal{W}(i,j,g,p) \right|.$$
(2)

Step 1: Perturbation Error Bound as Loss.

Proposition 1 implies the worst-case scenario when pruning a certain

GAN-Comp	2.67G	64.95	—	Auto-GAN	4.34G	78.33	—
DMAD	2.41G	62.96	—	GAN-KD	3.20G	80.10	—
CAT	2.55G	60.18	—	SP-KD	3.20G	76.59	—
GCC	2.40G	59.19	—	DMAD	3.18G	78.24	—
OMGD [2]	1.408G [†]	51.97	_	 OMGD [2]	1.408G [†]	73.79	_
+ZSP	1.406G	51.70	2	+ZSP	1.404G	73.70	6
$+L_{BIG}$	1.408G*	46.72	—	+L _{BIG}	1.408G*	73.12	—
$+L_{BIG}+ZSP$	1.397G	47.03	10	 + <i>L_{BIG}</i> +ZSP	1.398G	73.13	9

Table 1: Performance of knowledge distillation models combining the BIG loss pre-training and/or zero-shot pruning. † stands for the officially released models, * stands for our pre-trained models. Even in the currently reported best performance distillation models, there still exist channels that can be pruned without influencing the generation results.



channel. When applied as loss, we enforce a part of the channels have smaller bounds than the rest. Minimizing the loss will push the model to learn more distinguish channels, which are selected to be pruned or not.

$$L_{all} = L_{GAN} + \lambda_1 L_{dist} + \lambda_2 L_{BIG},$$

$$L_{BIG}(\{\mathcal{W}_{[I]}\}, \{\gamma_{[I]}\}, \{\beta_{[beta]}\}) = \sum_{I \in [L]} \sum_{i \in [C_I]} P_{I,i} \left(\mathcal{W}_{[I]}, \gamma_{[I]}, \beta_{[I]}\right).$$

In the equation, $P_{I,i} = 0$ if $\gamma_{[I],i} = 0$, while if $\gamma_{[I],i} \neq 0$, then

$$\mathsf{P}_{I,i} = \begin{cases} \sum_{j \in [D]} \mathcal{F}_{i,j} \left(\mathcal{W}_{[I]}, \gamma_{[I]}, \mathbf{0} \right), & \beta_i \geq \tau_{[I],i} \\ \sum_{j \in [D]} \mathcal{F}_{i,j} \left(\mathcal{W}_{[I]}, \gamma_{[I]}, \beta_{[I]} \right), & |\beta_i| < \tau_{[I],i} \\ \mathbf{0}, & \text{otherwise} \end{cases}$$

o 25 50 75 100 125 150 175 epoch (a) (b)

Fig. 3: (a) Learning curves of FID score and the number of channels along with the epoch of a complete pruning run. (b) The generated zebra images from the models are marked as red dots in (a), and the methods and the model ids are annotated above the images.

References

(3)

(4)

(5)

- [1] Zhu, Jun-Yan and Park, Taesung and Isola, Phillip and Efros, Alexei A: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, 2017
- [2], Ren, Yuxi and Wu, Jie and Xiao, Xuefeng and Yang, Jianchao: Online multi-granularity distillation for gan compression, 2021