

---

# Anycost GAN Optimization

---

**Mahmoud Khalifa, Anthony Ou**  
mahmoudk@mit.edu, aou@mit.edu

## Abstract

Efficient inference of machine learning models have become relevant in the past year because of massive improvements in the quality of large models. The Anycost GAN is a model that has been optimized for practical use cases. In particular, it features a capability to handle a different computational cost loads, and input image qualities. However, we believe that the Anycost GAN can be made even more efficient with the ideas that we have learned in class. In this work, we will apply ideas of magnitude based weight pruning and parameter quantization to this model in an effort to further improve performance. While not all of our ideas were successful, we have shown through our exploration that it is possible to significantly improve the model performance with minimal loss in generated image quality.

## 1 Introduction

With recent developments in generative models, there has been a renewed interest in applications of machine learning models. However, many generative models such as Large Language models or diffusion models are too large and run at speeds too slow to be of use on edge devices. As such, there has been much interest in both reducing number of parameters and decreasing MACs. One such model is the Anycost GAN, a GAN made efficient for the purpose of quick generation while still maintaining the quality of the generated images (1).

The Anycost GAN is capable of generating images at a range of computational costs while producing visually consistent outputs. These varied costs include different image resolutions and generation times. It achieves this by using a version of neural architecture search. The model is trained so that it is flexible to channel reduction, as such it is capable of producing low-cost sub-generators by taking a subset of the trained weights without additional fine-tuning. While the Anycost GAN is very powerful for its consistency across different computational costs, it is not yet optimal. We believe that we can achieve more efficient performance of this model through the application of techniques that we have learned in this class. Namely we applied techniques of fine grained pruning, channel pruning, K-means quantization, and linear quantization to further optimize the Anycost GAN.

## 2 Related Works

### 2.1 Pruning

Considerable work has been done in structured weight-pruning, or the systematic removal of a subset of a model's weights to improve performance (2). Recent examples include applications to Large-Language models (3) and deep convolutional neural networks (4). As described in lectures, fine-grained pruning or magnitude based pruning has shown to be capable of drastically reducing model sizes with minimal loss of accuracy (5). It has specifically been applied to both vision networks (6), and language networks (7) and is central to the lottery ticket hypothesis (8). There have been many other efforts to both introduce and exploit sparsity in model parameters (9). Channel pruning is capable of giving drastic model speed-ups in addition to drastically reducing model size (10) (11). Due to its ability to reduce inference costs, channel pruning has been widely used to prepare models

38 for operation on edge devices (12). While these methods are generally coupled with fine-tuning to  
39 recover lost performance, it is possible to have high performance without fine-tuning (13).

40 **2.2 Quantization**

41 K-means quantization was first explored as Huffman coding of model parameters and has been shown  
42 to drastically improve model speed and reduce model size (14). A variation of this method has  
43 been used in many use cases including image compression (15) and weather forecasting (16). These  
44 methods have since been built upon to develop more sophisticated methods of quantization (17).  
45 Linear quantization casts all floating point arithmetic to integer arithmetic to drastically speed-up  
46 inference (18). These methods have been expanded to also make training more efficient and includes  
47 other efforts to (19). Other variations of this method have shown improvements in accuracy on  
48 edge devices (20). These methods have recently been specifically applied to GANs for improved  
49 performance in edge devices (21) and for improved handling of compressed images (22).

50 **3 Pruning**

51 In optimizing the Anycost GANs model, two primary techniques were employed: Fine Grained  
52 Pruning and Channel Pruning. These techniques aimed to reduce the model size and increase  
53 inference speed, respectively. The outcomes of these optimizations are discussed in the following  
54 subsections.

55 **3.1 Fine Grained Pruning**

56 Fine Grained Pruning, particularly Magnitude Based Fine Pruning, was employed to reduce the size  
57 of the neural network model while minimizing the impact on performance. This pruning technique  
58 targets individual weights in the network, removing those with the smallest absolute values under the  
59 assumption that they contribute less to the network's output.

60 The pruning process can be mathematically described as follows. Let  $W$  represent the weight matrix  
61 in a neural network layer. We define a pruning threshold  $\theta$ , typically based on the desired pruning  
62 percentage  $p$ . The pruned weight matrix  $W'$  is then derived by setting weights below this threshold  
63 to zero. The expression for this operation is:

$$W'_{ij} = \begin{cases} 0, & \text{if } |W_{ij}| \leq \theta \\ W_{ij}, & \text{otherwise} \end{cases}$$

64 In this process,  $W'_{ij}$  denotes the weight in the pruned matrix  $W'$ , and  $W_{ij}$  is the corresponding weight  
65 in the original matrix  $W$ . The threshold  $\theta$  is selected such that a predetermined percentage  $p$  of the  
66 smallest absolute weights in  $W$  are eliminated.

67 Our approach first involved examining the distribution of parameters across layers, as illustrated  
68 in Figure 8. We observed a concentration of parameters in the Convolutional Layers, leading us to  
69 focus our fine-grained pruning efforts there. The distribution of weights prior to pruning resembled a  
70 normal distribution (Figure 6). After applying Magnitude Based Fine Pruning, the resulting weight  
71 distribution is depicted in Figure 7.

72 Implementing Magnitude Based Fine Pruning resulted in a 20% reduction in the overall size of the  
73 model. We didn't see any meaningful change in inference speed look Figure 9. This significant  
74 reduction was achieved without noticeable degradation in the quality of the generated images, as  
75 evidenced in Figure 1. Attempts to further prune the network at 30% and 40% pruning ratios are  
76 shown in Figures 2 and 3, respectively. Figure 2 indicates a slight loss in image quality, suggesting  
77 that even a minor reduction in model size can adversely affect our specific task of generating faces. A  
78 similar trend is observed in Figure 3, where excessive pruning leads to poor results.

79 This experiment portrays the efficacy of Magnitude Based Fine Pruning in reducing model size while  
80 preserving output quality. However, it also highlights the need for careful calibration of the pruning  
81 process, as excessive pruning can lead to a noticeable decline in performance.

82 **3.2 Channel Pruning**

83 Channel Pruning focuses on accelerating the inference process of the neural network. This method  
84 involves reducing the number of channels in the convolutional layers, thereby decreasing the compu-  
85 tational load during inference and leading to faster processing times.

86 In the context of AnyCost GANs, which are designed to be adaptable across various computational  
87 constraints, applying Channel Pruning requires careful consideration. These models often have  
88 built-in mechanisms to adjust their complexity based on available resources, and Channel Pruning  
89 could potentially conflict with these dynamic adjustments.

90 Despite these considerations, we implemented Channel Pruning in our AnyCost GAN model. The  
91 process involved systematically reducing the number of channels in selected convolutional layers,  
92 with the goal of finding a balance between efficiency and performance.

93 The implementation of Channel Pruning resulted in a 18% speedup in inference time when we  
94 tested out different images with different settings and averaged the times in Figure 9. However, this  
95 increase in efficiency came at the cost of image quality. The observed severe decrease in image  
96 quality, as shown in Figure 4, highlights the trade-off inherent in Channel Pruning: while it enhances  
97 computational efficiency, it can also lead to a significant compromise in model performance if not  
98 carefully managed.

99 In our experiments, we observed that the reduction of channels led to a noticeable decline in the  
100 fidelity of the generated images. This degradation is particularly evident in complex image regions  
101 where higher representational capacity is required. It suggests that while Channel Pruning can  
102 be effective for reducing computational requirements, it must be applied judiciously, especially in  
103 generative models where image quality is paramount.

104 Moving forward, it is essential to consider adaptive pruning strategies that align with the dynamic  
105 nature of AnyCost GANs. This might include developing methods that can intelligently prune  
106 channels based on the specific requirements of the generated image or the available computational  
107 resources at inference time.

108 **4 Quantization**

109 **4.1 K-Means Quantization**

110 In an effort to optimize the AnyCost GAN, we employed K-means quantization, a technique aimed at  
111 reducing the model size and potentially speeding up inference by quantizing the model’s weights.  
112 The quantization process was conducted with varying bitwidths of 2, 4, and 8, setting the number of  
113 clusters  $k$  to  $2^{bitwidth}$ .

114 Mathematically, K-means quantization can be described as follows. For a given set of weights  $W$  in  
115 the neural network, the K-means algorithm partitions  $W$  into  $k$  clusters. Each weight  $w \in W$  is then  
116 replaced by the nearest centroid of the cluster it belongs to. The quantized weights  $W'$  are given by:

$$W' = \{c_i \mid c_i = \operatorname{argmin}_{c \in C} \|w - c\|^2, w \in W\}$$

117 where  $C$  is the set of centroids for the  $k$  clusters.

118 Our experiments revealed that for bitwidths of 2 and 4, the image quality deteriorated significantly, as  
119 evidenced by the results in Figure 5. This quality degradation is attributed to the coarse quantization,  
120 which leads to a loss of detail in the weight representation.

121 Furthermore, we observed a substantial increase in the initial setup time for the quantized models.  
122 The setup times were approximately one minute for a bitwidth of 2, five minutes for a bitwidth of 4,  
123 and an hour for a bitwidth of 8. This suggests that the complexity of the K-means clustering process  
124 increases with the number of clusters, impacting the initial preparation phase of the model.

125 Despite these drawbacks in initial setup time and image quality at lower bitwidths, the quantization  
126 process did not affect the operational efficiency during inference. Specifically, the time required for  
127 manipulating generated images (such as changing facial expressions) remained constant across all  
128 tested bitwidths.

129 These findings highlight the delicate balance in K-means quantization applied to AnyCost GANs,  
130 emphasizing the critical role of bitwidth selection. While lower bitwidths offer model size reduction,  
131 they come at the cost of significant image quality degradation and increased setup time, making them  
132 less viable for high-quality image generation and time-sensitive applications.

133 **4.2 Linear Quantization**

134 Based on what we learned in this class, we expected that implementing linear quantization into this  
135 model would result in drastic improvements in performance as linear quantization is independent  
136 from pruning and subnetwork search AnyCost performs. However, the implementation of AnyCost that  
137 we were using does not allow for easy integration of well used linear quantization packages and so  
138 we were forced to implement it manually. Given the prevalent usage of floating point noise, weight  
139 normalization, and adaptive channel inference it was infeasible to implement linear quantization  
140 end-to-end. Instead, we opted to quantize the convolution, as per figure 8 the convolutions are the  
141 computationally heaviest portion of this model. As such, for every convolution module, we quantized  
142 the input  $q$ , the weights  $w$  and the bias  $b$  then performed the convolution as integers to produce a  
143 quantized output  $y$ .

$$\begin{aligned}x &= S_x(q_x - Z_x) \\w &= S_w(q_w - Z_w) \\o &= S_o(q_o - Z_o) \\b &= S_b(q_b - Z_b)\end{aligned}$$

144 Where  $S$  and  $Z$  are the scale and zero-point respectively and  $q$  is the integer representation. Due to  
145 how the model applied scaling and normalization after selecting with its channel ratios, we decided  
146 not to implement per-channel quantization and instead used a single scale and zero-point for each  
147 convolution module. We used a global bitwidth of 8 bits to select our scale and zero-point, and used  
148 8-bit integers to represent  $q$ . The model was run normally once and the inputs and outputs for each  
149 module were used to calibrate the selection of the scales and zero-points. As we covered in Lab 2,  
150 the result for the output is:

$$q_o = (\text{CONV}[q_x, q_w] + q_b - \text{CONV}[Z_x, q_w]) \cdot (S_x S_w / S_o) + Z_o$$

151 Since all operations are now integer operations, we expected to see considerable speed-ups from this  
152 implementation. However, to the contrary we experienced very large decreases in performance; we  
153 determined that the reason for these slow-downs is that the PyTorch implementation for convolutions  
154 is far slower on integers than on floats. This is an issue that does not seem to have much support for.  
155 Given that the entirety of the code base is in PyTorch, we decided that it is infeasible to continue  
156 down this line of investigation as to do so would require a drastic overhaul of the codebase and all  
157 alternatives we explored seemed to either contain errors or be not applicable to our use case.

161 **5 Fine Tuning**

162 In our study, we have employed techniques such as Fine Grained Pruning and K-means Quantization  
163 to optimize the AnyCost GAN. Future work in this research is to explore the effects of fine-tuning  
164 on these optimized models. Fine-tuning, a process where a pre-trained model is further trained on a  
165 specific dataset or task, can lead to notable improvements in performance, especially in models that  
166 have undergone significant structural changes.

167 We anticipate that fine-tuning the pruned and quantized models could result in considerable enhance-  
168 ments. The idea is that fine-tuning allows the model to adjust its parameters more effectively to the  
169 dataset's specific features. This adjustment could help in offsetting performance losses incurred due  
170 to pruning or quantization. It might even make it feasible to prune the models further or to use lower  
171 bitwidths in quantization, thereby increasing the model's speed without a substantial compromise in  
172 image quality.

173 However, our ability to pursue fine-tuning is currently constrained by two significant factors:

- 174 1. Large Dataset Size: Our models were originally trained on the Flickr-Faces-HQ Dataset  
175 (FFHQ), which contains 70,000 high-quality images. The sheer size of this dataset presents  
176 a challenge for fine-tuning, especially considering the computational resources required to  
177 process such a large amount of data effectively.

178     2. Lack of GPU Resources: Another critical limitation is the absence of adequate GPU  
179       resources. Fine-tuning, particularly on large datasets like FFHQ, demands substantial  
180       computational power, typically provided by GPUs. Without access to such resources, it's  
181       impractical to undertake the fine-tuning process, which is essential to realize the full potential  
182       of the optimized models.

183     **6 Conclusion**

184     This study explored enhancing the Anycost GAN's efficiency, focusing on fine-grained pruning,  
185       channel pruning, and K-means quantization, alongside an attempt at linear quantization. Our results  
186       illuminated the complex trade-offs in model optimization: while fine-grained pruning effectively  
187       reduced size with minimal quality loss, channel pruning highlighted the balance between speed  
188       and fidelity. K-means quantization's impact on image quality was dependent on bitwidth selection,  
189       and our linear quantization efforts revealed integration challenges in complex models. Although  
190       constrained by resource limitations and unable to explore fine-tuning, our findings offer insights into  
191       the balance between model efficiency, size, and output quality. This research contributes to the further  
192       optimization of Anycost GAN and hopefully will help run it on smaller devices.

193 **Figures**

Figure 1: 20% Pruned Model

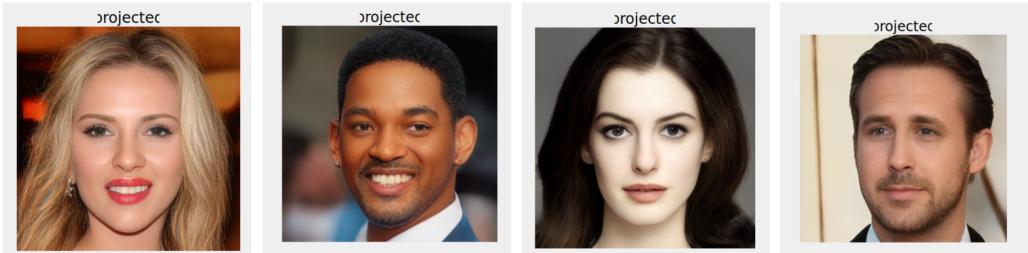


Figure 2: 30% Pruned Model

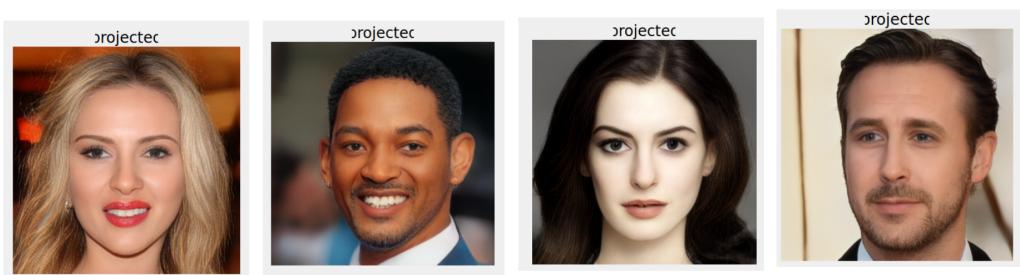


Figure 3: 40% Pruned Model

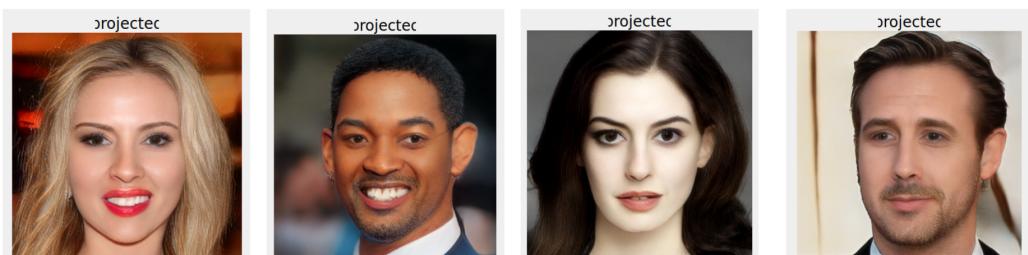


Figure 4: Channel Pruned Model

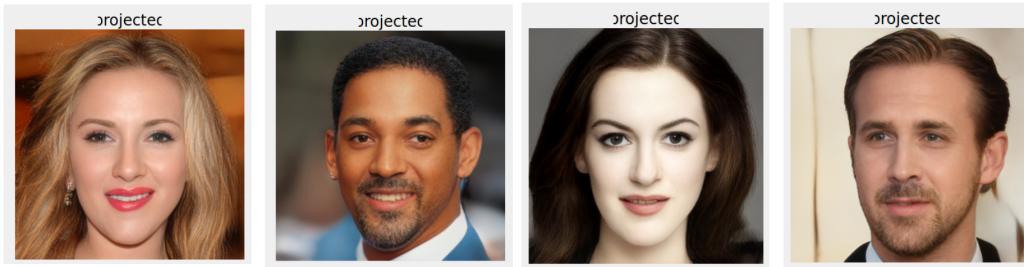


Figure 5: K-Means Quantized Model



Bitwidth = 2

Bitwidth = 4

Figure 6: Layer Weights

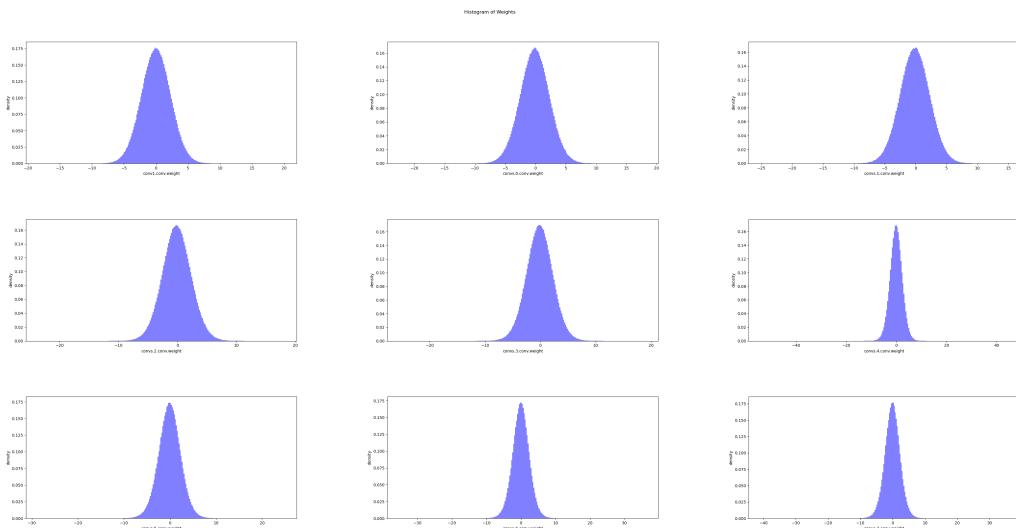


Figure 7: Layer Weights

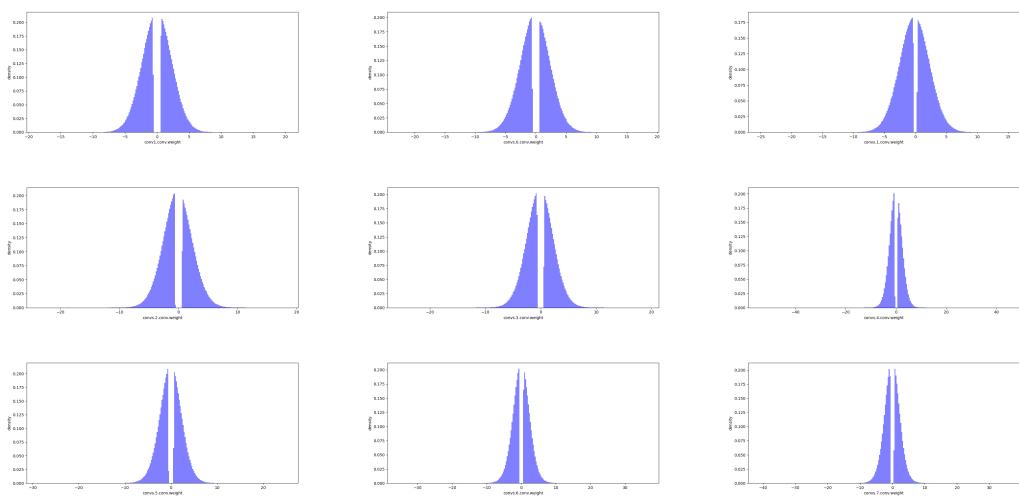


Figure 8: Layers Number of Parameters

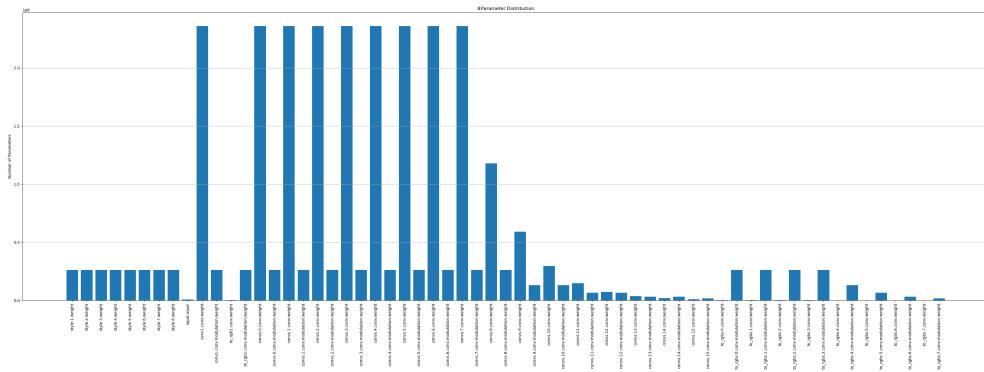


Figure 9: Average seconds to generate image

Image	Original	Fine Grained	Channel Pruning
Ryan	0.7s	0.72s	0.57s
Anne	0.7s	0.67s	0.58s
Will	0.68s	0.69s	0.56s
Scarlett	0.69s	0.67s	0.54s
Average	0.6925s	0.6875s	0.5625s

194 **References**

- 195 [1] J. Lin, R. Zhang, F. Ganz, S. Han, and J.-Y. Zhu, “Anycost gans for interactive image synthesis  
196 and editing,” *arXiv preprint arXiv:2103.03243*, 2021.
- 197 [2] G. Fang, X. Ma, M. Song, M. B. Mi, and X. Wang, “Depgraph: Towards any structural pruning,”  
198 2023.
- 199 [3] M. Xia, T. Gao, Z. Zeng, and D. Chen, “Sheared llama: Accelerating language model pre-  
200 training via structured pruning,” 2023.
- 201 [4] Y. He and L. Xiao, “Structured pruning for deep convolutional neural networks: A survey,”  
202 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, p. 1–20, 2023. [Online].  
203 Available: <http://dx.doi.org/10.1109/TPAMI.2023.3334614>
- 204 [5] S. Han, J. Pool, J. Tran, and W. J. Dally, “Learning both weights and connections for efficient  
205 neural networks,” 2015.
- 206 [6] Y. Guo, A. Yao, and Y. Chen, “Dynamic network surgery for efficient dnns,” 2016.
- 207 [7] T. Gale, E. Elsen, and S. Hooker, “The state of sparsity in deep neural networks,” 2019.
- 208 [8] J. Frankle, G. K. Dziugaite, D. M. Roy, and M. Carbin, “Linear mode connectivity and the  
209 lottery ticket hypothesis,” 2020.
- 210 [9] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste, “Sparsity in deep learning:  
211 Pruning and growth for efficient inference and training in neural networks,” 2021.
- 212 [10] Y. He, X. Zhang, and J. Sun, “Channel pruning for accelerating very deep neural networks,”  
213 2017.
- 214 [11] Z. Zhuang, M. Tan, B. Zhuang, J. Liu, Y. Guo, Q. Wu, J. Huang, and J. Zhu, “Discrimination-  
215 aware channel pruning for deep neural networks,” 2019.
- 216 [12] L. Xue, Z. Wang, X. Wang, and G. Li, “The road to on-board change detection: A lightweight  
217 patch-level change detection network via exploring the potential of pruning and pooling,” 2023.
- 218 [13] S. Bai, J. Chen, X. Shen, Y. Qian, and Y. Liu, “Unified data-free compression: Pruning and  
219 quantization without fine-tuning,” 2023.
- 220 [14] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with  
221 pruning, trained quantization and huffman coding,” 2016.
- 222 [15] C. Reich, B. Debnath, D. Patel, and S. Chakradhar, “Differentiable jpeg: The devil is in the  
223 details,” 2023.
- 224 [16] C. Chen, T. Zhou, Y. Zhao, H. Liu, L. Sun, and R. Jin, “Svq: Sparse vector quantization for  
225 spatiotemporal forecasting,” 2023.
- 226 [17] P. Stock, A. Joulin, R. Gribonval, B. Graham, and H. Jégou, “And the bit goes down: Revisiting  
227 the quantization of neural networks,” 2020.
- 228 [18] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko,  
229 “Quantization and training of neural networks for efficient integer-arithmetic-only inference,”  
230 2017.
- 231 [19] Y. Yang, X. Dai, J. Wang, P. Zhang, and H. Zhang, “Efficient quantization strategies for latent  
232 diffusion models,” 2023.
- 233 [20] C. Zhu, S. Han, H. Mao, and W. J. Dally, “Trained ternary quantization,” 2017.
- 234 [21] P. Wang, D. Wang, Y. Ji, X. Xie, H. Song, X. Liu, Y. Lyu, and Y. Xie, “Qgan: Quantized  
235 generative adversarial networks,” 2019.
- 236 [22] Q. Mao, T. Yang, Y. Zhang, S. Pan, M. Wang, S. Wang, and S. Ma, “Extreme image compression  
237 using fine-tuned vqgan models,” 2023.