

**Department of Physics and Astronomy
Heidelberg University**

Master Thesis in Physics
submitted by

Sebastian Heid

born in Gelnhausen (Germany)

2026

Diffusion Model Distillation

This Master Thesis has been carried out by Sebastian Heid at the
Computer Vision and Learning Lab in Heidelberg
under the supervision of
Prof. Dr. Carsten Rother
and
...

Abstract

Table of Contents

Table of Contents	IV
List of Abbreviations	IX
List of Figures	X
List of Tables	XII
1 Introduction	1
2 Background	4
2.1 Generative Models	4
2.1.1 Variational Autoencoder	4
2.1.2 Normalizing Flow	6
2.1.3 Generative Adversarial Network	7
2.2 Diffusion and Flow Matching Models	8
2.2.1 Denoising Diffusion Probabilistic Models	8
2.2.2 Denoising Diffusion Implicit Models	13
2.2.3 Score-Based Generative Models	14
2.2.4 Stochastic Differential Equations	15
2.2.5 Flow Matching	15
2.2.6 Latent Diffusion Models	17
2.2.7 Conditioning and Classifier-Free Guidance	18
2.3 Fundamentals of Neural Network Architectures	19
2.3.1 Transformer	19
2.3.2 Attention Mechanism	20
2.3.3 Positional Encoding	22
2.3.4 Vision Transformer	23
2.3.5 Adaptive Layer Normalization	23
2.4 Model Distillation	25
2.4.1 Knowledge Distillation	25
2.5 Evaluation Metrics	27
2.5.1 Fréchet Inception Distance	27
2.5.2 CLIP-Maximum Mean Discrepancy	27
2.5.3 Human Preference Score v2 Benchmark	28
2.5.4 GenEval Benchmark	28
2.5.5 Dense Prompt Graph Benchmark	29
3 Related Work	30

TABLE OF CONTENTS

4 Methods	34
4.1 Image Generation Models	34
4.1.1 Pixart- Σ	34
4.1.2 Flux.1-dev	37
4.2 Block Importance Analysis	42
4.2.1 Block Selection Criterias	42
4.2.2 Block Selection Algorithms	43
4.3 Training Data and Image Synthesis	44
4.3.1 Training Datasets	44
4.3.2 Evaluation Datasets	45
4.4 Model Compression Strategies	47
4.4.1 Block Removal	47
4.4.2 SVD Compression	47
4.4.3 Iteratively Repeated SVD Compression	48
4.4.4 GRASP	48
4.5 Knowledge Distillation Loss Functions	49
4.5.1 Normalization	49
4.5.2 Feature Mapping for Loss Computation	49
4.6 Pruning Schedule	50
4.6.1 One-Shot Model Pruning	50
4.6.2 Progressive Model Pruning	51
4.6.3 Linear Progressive Model Pruning	52
5 Experiments	55
5.1 PixArt- Σ Pruning	55
5.1.1 Experimental Default Setup	55
5.1.2 Block Importance Analysis	56
5.1.3 Block Selection Algorithms	61
5.1.4 One-Shot vs Progressive Model Pruning	63
5.1.5 Knowledge Distillation Loss	65
5.1.6 fine-tuning Protocol	66
5.1.7 Compression Strategy	68
5.1.8 Sensitivity Analysis of MLP and Attention Layers under SVD Compression	71
5.1.9 SVD Compression Scheduling	74
5.1.10 Summary of Best Settings for Structural Model Pruning	74
5.2 Flux.1-dev Pruning	76
5.2.1 Training Details	76
5.2.2 Qualitative Block Analysis	77
5.2.3 Training-Free Pruning	80
5.2.4 Linear Progressive SVD Compression	86
6 Appendix	90
6.1 Appendix A	90
6.1.1 Block Analysis	90
6.2 Appendix B - Detailed Overview of Selected Blocks Experiments PixArt- Σ	90
6.3 Experiments Flux.1-dev	91
6.3.1 GRASP	91

TABLE OF CONTENTS

Bibliography	92
--------------	----

List of Abbreviations

adaLN Adaptive Layer Normalization

ADD Adversarial Diffusion Distillation

AI Artificial Intelligence

ALD Annealed Langevin Dynamic

CFG Classifier-Free Guidance

CKA Central Kernal Alignment

CLIP Contrastive Language- Image Pre-training

CLS Classification

CMMID Fréchet Inception Distance

CNF Continous Normalizing Flow

CNN Convolutional Neural Network

DDIM Denoising Diffusion Implicit Model

DDPM Denoising Diffusion Probablistic Model

DiT Diffusion Transformer

DM Diffusion Model

DPG Dense Prompt Graph

DSG Dynamic Scene Graph

ELBO Evidence Lower Bound

FID Fréchet Inception Distance

GAN Generative Adversarial Network

GenAI Generative Artificial Intelligence

GRASP Gradient-based Retention of Adaptive Singular Parameters

HPSv2 Human Preference Score v2

Abbreviations

HSIC Hilbert-Schmidt Independence Criterion

iid independent and identically distributed

KL Kullback-Leibler

LADD Latent Adversarial Diffusion Distillation

LDM Latent Diffusion Model

LLM Large Language Model

LoRA Low-Rank Adaptation

LPIPS Learned Perceptual Image Patch Similarity

MC Markov Chain

MCMC Markov Chain Monte Carlo

MLE Maximum Likelihood Estimation

MLP Multi Layer Perceptron

MMD Maximum Mean Discrepancy

MMDiT Multi-Modal Diffusion Transformer

MoE Mixture-of-Expert

MSE Mean Squared Error

NCSN Noise-Conditional Score Network

NF Normalizing Flow

NNL Negative Log-Likelihood

ODE Ordinary Differential Equation

PTQ Post-Training Quantization

PTQD Post-Training Quantization for Diffusion Models

QAT Quantizaiton-Aware Training

RMS Root Mean Square

RMSNorm Root Mean Square Normalization

RNN Recurrent Neural Network

RoPE Rotary Positional Embedding

SDE Stochastic Differential Equation

SGM Score-Based Generative Model

Abbreviations

SOTA State-Of-The-Art

STE Straight-Through Estimator

SVD Singular Value Decomposition

VAE Variational Autoencoder

VI Variational Inference

ViT Vision Transformer

VRAM Video Random Access Memory

List of Figures

2.1	Visualization of DDPM	8
2.2	DDPM Training Scheme	12
2.3	DDPM Sampling Scheme	13
2.4	Rectified Flow	17
2.5	Latent Diffusion Model	18
2.6	Transformer Architecture	19
2.7	Self-Attention Mechanism	21
2.8	Cross-Attention Mechanism	21
2.9	Multi-Head Attention	22
2.10	Vision Transformer Architecture	23
2.11	GenEval Benchmark	29
4.1	PixArt- Σ Architecture	35
4.2	Transformer Block from PixArt- Σ	36
4.3	Flux.1-dev Architecture	39
4.4	Double- and Single-Stream Transformer Blocks from Flux.1-dev	40
4.5	LAION Dataset	44
4.6	LAION-PixArt- Σ dataset	45
4.7	LAION-Flux.1-dev dataset	45
4.8	Mapillary Dataset	46
4.9	Cityscapes Dataset	46
4.10	MHQ-30k Dataset	46
4.11	Model Compression Strategies	48
4.12	Feature Mapping Strategy	50
4.13	One-Shot Compression Scheme	51
4.14	One-Shot Compression Scheme	52
4.15	Linear Progressive Compression Scheme	54
5.1	Qualitative analyse of the block sensitivity of PixArt- Σ for block removal	57
5.2	Magnitude-base pruning versus CLIP-score	58
5.3	Magnitude-base pruning versus CMMD	58
5.4	CKA transformation intensity versus CLIP-score	59
5.5	CKA transformation intensity versus CMMD	59
5.6	Qualitative analyse of the block sensitivity of PixArt- Σ of block compression	61
5.7	Distribution of Selected Blocks: Optuna versus Greedy Algorithm	62
5.8	Benchmark Performance: Optuna versus Greedy Algorithm	63
5.9	Comparison of Pruning Strategies: Optuna versus Greedy Algorithm	63

LIST OF FIGURES

5.10	Benchmark Performance: Progessive versus One-Shot Model Pruning	64
5.11	Example images for one-shot versus progressive pruned models by 50%.	64
5.12	Benchmark Performance: Different Knowledge Distialltion Loss Configurations	66
5.13	Magnitude of Intermediate Features in PixArt- Σ	66
5.14	Benchmark Performance: fine-tuning Protocols	67
5.15	Benchmark Performance: Compression Strategies	69
5.16	Qualitative Comparison of Compression Strategies	70
5.17	Benchmark Performance: Compression of Individual Transformer Block Components	72
5.18	Example Images Sensitivity Analysis of MLP and Attention Layers	73
5.19	Flux.1-dev Analysis Block Removal - Single-Stream Blocks	78
5.20	Flux.1-dev Qualitative Analysis Block Removal - Double-Stream Blocks	79
5.21	Flux.1-dev Qualitative Analysis Block Compression - Single-Stream Blocks	79
5.22	Flux.1-dev Qualitative Analysis Block Compression - Double-Stream Blocks	80
5.23	Flux.1-dev Parameter Scaled CMMMD Values	80
5.24	Magnitude versus Importance Analysis of Singular Values	81
5.25	Benchmark Performance: SVD versus GRASP Compression Method	82
5.26	Flux.1-dev Images: SVD versus GRASP Compression	83
5.27	Flux.1-dev Text Generation Capabilities: SVD versus GRASP versus EcoDiff Compression . . .	85
5.28	Flux.1-dev Images: SVD versus GRASP versus EcoDiff Compression	85
5.29	Benchmark Performance: Random versus Metric-based Block Selection	86
5.30	Example images of the 49% models (5.8B parameters) obtained by random and importance-based block selection.	87
5.31	Benchmark Performance: Random versus Metric-based Block Selection	88
6.2	Training-Free: Flux.1-dev Block Importance Ranking	91

List of Tables

4.1	PixArt- Σ Parameter Distribution	36
4.2	Flux.1-dev Parameter Distribution	41
4.3	Parameter LAION-PixArt- Σ Dataset	44
4.4	Parameter LAION-Flux.1-dev Dataset	45
5.1	Hyperparameter Setting for PixArt- Σ fine-tuning	56
5.2	Training and Compression Settings PixArt- Σ for Block Removal	64
5.3	Training and Compression Settings PixArt- Σ for Block Removal and SVD	68
5.4	Flux fine-tuning Hyperparameters	77
5.5	Benchmark Performance: Training-Free Compression Methods	84
5.6	Overview of training details similar for both experiments, random and metric-based block selection.	87
5.7	Overview of compression metrics per iteration.	88
5.8	Benchmark Performance: Training-Free Compression Methods	89
6.1	Trainig-Free Flux Compression: SVD Ranks	91
6.2	Training Free Flux Compression: Model Compression Configuration	91

Chapter 1

Introduction

Over the past few years, generative artificial intelligence (GenAI) has gained immense importance and relevance in society, as well as in business and politics. In addition to large language models (LLMs) [1, 5, 172, 176], which have brought GenAI into the mainstream, generative image models [9, 15, 20, 93, 133, 144] are receiving more and more attention. The ever-increasing number of image generator models and simplified access present both opportunities and risks. On the one hand, they open up a range of new possibilities in areas such as design and marketing as well as gaming and entertainment. On the other hand, they enable the generation of highly realistic synthetic images and videos, often referred to as deepfakes, raising severe ethical and privacy concerns, which is why policymakers are striving to establish rule-based guidelines for their use (European AI Act).

The advent of deep learning based image generators is often marked by the introduction of variational autoencoders (VAEs) in 2013 [84], as they demonstrated that neural networks can learn complex data distribution and generate novel samples. Subsequently, generative adversarial networks (GANs) [53, 81, 152] became the dominant model type in the image generation field as they produce images of significantly higher quality than VAEs. However, GANs are notoriously difficult to train, as they often exhibit mode collapse and suffer from exploding or vanishing gradients [106, 149, 174].

In recent years, GANs were largely surpassed by diffusion models [68] deploying an iterative image generation denoising process. The breakthrough of diffusion models as the new state-of-the-art (SOTA) in image generation occurred mainly in 2022 when the latent diffusion models, most notably Stable Diffusion [145], were released. At this time, they excelled in image quality and image generation diversity, but came with the cost of high inference time due to the iterative nature of their image generation process. The transition from pixel space to latent space massively reduced the computing time of the models and enabled images with higher resolution to be generated.

The first stable diffusion models [133, 144] leveraged a U-Net architecture [146] and already required a huge amount of training data and computational resources. Their image generation process is typically controlled by textual prompts provided via cross-attention mechanisms [178]. To expand these capabilities, several new control and finetuning methods like ControlNet [194] and DreamBooth [147] were developed to enable precise spatial conditioning and personalization. Since the U-Net-based SD models already included transformer layers [178] for text-conditioning, the next logical evolution was to transition to purely transformer-based models known as diffusion transformers (DiTs) [131], which increased significantly in parameter size and required even more training data. In addition, a new formulation, so-called flow-matching, became predominant due to improved training stability and sampling efficiency which is leveraged by the current SOTA image generation models like Flux.1-dev [93] being a 11.9 billion parameter model.

Possible areas of application for diffusion and flow models include edge devices such as mobile phones [27, 103, 198]. In these cases, it is advantageous to run these models directly on the edge devices in order to avoid server costs, to make the associated functionality available offline and to protect the privacy of users so that, for

example, images do not have to be uploaded to external servers for processing. However, the current flow-based models like Flux.1-dev face two distinct challenges for real-time applications on edge devices. On the one hand, as mentioned before, they suffer from a long inference time due to the sequential application of the model. On the other hand, these models also require a substantial amount of VRAM during inference, e.g. Flux.1-dev requires about 32GB, which makes them unsuitable on edge devices.

Reducing VRAM and faster execution plays also a pivotal role in democratising these models and making them more environmentally friendly, as they can run on consumer-grade hardware and do not require power-intensive data centres reducing carbon footprint.

The main research stream for speeding up the inference time is to enable the model to generate high-quality images with fewer inference steps. While at the beginning, diffusion models required up to 1000 inference steps to generate one single image, using more advanced ODE-solvers for the denoising trajectory has led to significant reduction in inference steps.

Another line of research is step- and guidance distillation. Distillation methods generally attempt to transfer the capabilities of a teacher model to a student model. In step distillation, the same diffusion model is often used as both the teacher and student model, and the goal is to teach the student model to generate images of the same or similar quality with as few as four or two, sometimes even only a single inference step leveraging its original capabilities [117, 150, 151, 153, 160, 188]. This technique is for example applied to Flux.1-dev resulting in Flux.1-schnell, reducing the required inference steps to only four.

In guidance distillation, the goal is to incorporate the concept of classifier free guidance (CFG), which normally requires the model to be applied twice per inference step, into the model itself, thereby saving half of the model evaluations [122]. This procedure was utilized to distill Flux.1-pro into Flux.1-dev.

However, while step distillation and advanced solver effectively reduce the latency (inference time), they do not necessarily address the second bottleneck, namely high VRAM consumption. Proposed methods for VRAM reduction range from quantization [59, 185] and more efficient implementation of model components, e.g. Flash Attention [34], to structural or unstructural model pruning [17, 33, 82, 92, 96, 118, 173, 193, 196]. While quantization compresses the model size by reducing the precision of the model parameters, pruning directly addresses the phenomenon that modern models, which have a high number of parameters, exhibit high redundancy in and between the individual transformer blocks. For this reason, identification followed by compression of redundant model parts is a promising orthogonal approach to reducing inference steps for further efficiency gains.

The first part of this Master’s thesis systematically examines how the number of parameters in diffusion transformers can be effectively reduced while minimizing the degradation of image generation quality, using PixArt- Σ [20] as a representative baseline enabling the execution of numerous experiments due to its significantly smaller parameter count compared to Flux.1-dev. Specifically, the identification of architectural components with redundant parameters alongside the best pruning strategies, e.g., complete layer removal versus low-rank compression, and pruning schedules, e.g., one-shot versus progressive pruning are investigated. Furthermore, various post-compression retraining frameworks are evaluated to optimize performance recovery. The empirical results demonstrate that progressively compressing model blocks via low-rank approximation leads to the best results. Hereby, a linear distribution of compression ratios over transformer blocks meaning compressing the least important structures most aggressively and more critical structures less, is leveraged.

In the second part of this Master’s thesis, the previously identified best pruning strategy is applied to the Flux.1-dev model and compared with other existing compression methods.

Furthermore, a training-free compression strategy originally proposed for LLMs, namely GRASP, that is also based on singular value decomposition (SVD) is transferred to Flux.1-dev and tested in comparison with current alternative approach [196].

The remainder of this Master’s thesis is structured as follows: Chapter two provides a comprehensive overview of image generator models, emphasizing diffusion and flow models, alongside relevant metrics used for measuring

CHAPTER 1. INTRODUCTION

the performance. Building up on this context, chapter three contextualizes this work within the current literature by reviewing alternative model compression techniques. With the background established, chapter four introduces the specific methodologies investigated in this work. Their effectiveness is then rigorously evaluated through experimental results in chapter five. Finally, chapter six concludes the thesis with a comprehensive discussion of the main findings and their implications.

Chapter 2

Background

2.1 Generative Models

Generative models are a class of machine learning models that are trained to learn the true data distribution to generate new synthetic data elements. They are of high relevance for text and image generation tasks. In particular, famous models for text generation, also known as large language models (LLMs) include ChatGPT [1], Claude [5], Gemini [172] or Lama [176] whereas Stable Diffusion [133, 144], Flux [93] and Pixart [20] are widely used for image generation tasks. Both are being used more and more in the professional world and in daily life. Several different generator methods have been developed, such as variational autoencoders (VAEs) [84], generative adversarial networks (GANs) [53, 136], normalizing flows (NFs) [143], and diffusion models (DMs) [157].

The principal idea of a generative model f_θ is to learn a mapping from a simple distribution, e.g., a normal distribution $p_{\text{latent}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, often called the latent distribution, to a highly complex data distribution $p_{\text{data}}(\mathbf{x})$. A new data element \mathbf{x} is obtained by sampling from the latent distribution $\mathbf{z} \sim p(\mathbf{z})$ and applying the model

$$\mathbf{x} = f_\theta(\mathbf{z}) \sim p_{\text{model}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x}) \quad (2.1)$$

where the model's output distribution $p_{\text{model}}(\mathbf{x})$ should approximate the true data distribution as closely as possible.

In the following, the main concepts of VAEs, NFs and GANs are briefly presented to provide a comprehensive overview followed by an extensive discussion of diffusion and flow-based generative models in chapter 2.2 that are the most promising models today.

2.1.1 Variational Autoencoder

Autoencoder

An autoencoder [67] is a neural network commonly used for dimensionality reduction tasks. It aims to compress high-dimensional data to lower-dimensional space by preserving important information. The architecture consists of two main parts, namely the encoder and the decoder. The encoder maps high-dimensional input data to a latent space, effectively compressing the data, and the decoder learns to reconstruct the original data from its latent space representation. Typically, as learning objective, a simple reconstruction loss, e.g. mean-squared error, is applied comparing the original input data and the output of the autoencoder.

Evidence Lower Bound Loss (ELBO)

This paragraph and its mathematical derivations are mainly based on [13, 84]. The VAE [84] belongs to the class of probabilistic models. A common learning objective for probabilistic models is maximum likelihood estimation (MLE). Given a set of observations $\{\mathbf{x}_i\}_{i=1}^N$ the likelihood is given under the assumption of identical and independent distributed (iid) samples as

$$\mathcal{L}_{\text{MLP}} = \prod_{i=1}^n p_{\text{model}}(\mathbf{x}_i) . \quad (2.2)$$

However, in practice, minimizing the negative log-likelihood is often preferred over maximizing the likelihood due to possible instabilities caused by the product of likelihoods

$$\mathcal{L}_{\text{NLL}} = - \sum_{i=1}^n \log p_{\text{model}}(\mathbf{x}_i) . \quad (2.3)$$

This approach works well for simple distributions. However, for models that include latent variables $\mathbf{z} = \mathbf{z}_{1:m}$, such as VAEs, the likelihood becomes an intractable marginalization over the latent space

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (2.4)$$

leading to the so-called evidence lower bound (ELBO) loss [13] being used as training objective for VAEs. In general, untractable probability distributions are a core challenge in modern statistics. This is especially the case for posterior distributions in Bayesian statistics. Assuming the joint probability distribution

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z}) \cdot p(\mathbf{x}|\mathbf{z}) \quad (2.5)$$

is given with data $\{\mathbf{x}_i\}_{i=1}^N$ and latent variables $\{\mathbf{z}_j\}_{j=1}^M$. During inference, the in general untractable posterior $p(\mathbf{z}|\mathbf{x})$ is the quantity that needs to be computed. In the past, the dominant approach to tackle this challenge was using Monte Carlo Markov Chain (MCMC) simulations [175], which approximate the true posterior distribution by sampling. The drawback of these methods is that the sampling procedures are very slow. Therefore, in recent years, a new method has gained popularity, called variational inference (VI) [14]. In contrast to MCMC VI reformulates the problem as an optimization task

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} \text{KL}[q(\mathbf{z}) \| p(\mathbf{z}|\mathbf{x})] \quad (2.6)$$

where \mathcal{Q} is a family of distributions. The goal is to find the proposal density $q(\mathbf{z}) \in \mathcal{Q}$ that best approximates the true posterior distribution $p(\mathbf{z}|\mathbf{x})$. To compare the proposal density and the posterior the Kullback-Leibler divergence (KL) [91] is used having the general form of

$$\text{KL}[p_0(\mathbf{x}) \| p_1(\mathbf{x})] = \int p_0(\mathbf{x}) \log \frac{p_0(\mathbf{x})}{p_1(\mathbf{x})} d\mathbf{x} \quad (2.7)$$

comparing two distributions p_0 and p_1 . The key for a good choice of \mathcal{Q} is to choose a family that is complex enough so that it contains a $q^*(\mathbf{z})$ close to the true posterior but simple enough to enable efficient optimization. Applying VI to the problem of finding the true posterior yields

$$\text{KL} [q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})] = \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \quad (2.8)$$

$$= \int q(\mathbf{z}) \log q(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \log p(\mathbf{z}|\mathbf{x}) d\mathbf{z} \quad (2.9)$$

$$= \int q(\mathbf{z}) \log q(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \log \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})} d\mathbf{z} \quad (2.10)$$

$$= \mathbb{E} [\log q(\mathbf{z})] - \mathbb{E} [\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}) \quad (2.11)$$

From here, the ELBO loss is obtained by

$$\log p(\mathbf{x}) \geq \log p(\mathbf{x}) - \text{KL} [q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})] \quad (2.12)$$

$$= \mathbb{E} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E} [\log q(\mathbf{z})] \quad (2.13)$$

$$= \mathbb{E} [\log p(\mathbf{x}|\mathbf{z})] - \text{KL} [q(\mathbf{z}) \parallel p(\mathbf{z})] \quad (2.14)$$

$$= \text{ELBO} \quad . \quad (2.15)$$

Maximizing the ELBO loss is equivalent to minimizing the KL divergence from 2.6 and most importantly, it is tractable.

For training VAEs the ELBO objective is used. The difference compared to the MSE loss usually used for training autoencoders is that it contains an additional term from the KL divergence that enforces a specific structure on the latent space. Let θ be the decoder parameters and ϕ the encoder parameters then the objective of the VAE is given by

$$\mathcal{L}_{\text{VAE}} = \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \underbrace{\text{KL} [q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})]}_{\text{Regularization}} \quad . \quad (2.16)$$

In practice, the latent distribution is often set to a standard normal distribution $p_{\text{latent}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{1})$ and the reparameterization trick [84] enabling gradient computation is applied. After training, a new data point can be generated by sampling from the latent distribution $\mathbf{z} \sim p_{\text{latent}}(\mathbf{z})$ and applying the decoder of the trained model.

2.1.2 Normalizing Flow

Normalizing flows (NFs) [143] are another class of probabilistic models, similar to VAEs, which learn a mapping between a latent distribution and the data distribution. However, NFs have some key differences compared to VAEs. First, NFs are fully invertible, meaning the network F_θ consists of a sequence of (simpler) bijective operations so that an inverse transformation \bar{F}_θ mapping back from data to latent distribution exists. As a consequence, the latent space and the data space need to have the same dimension.

A second decisive advantage is that NFs allow the exact and tractable computation of the likelihood using the change of variables formula [130]

$$p_{\text{latent}}(\mathbf{z}) = p_{\text{model}}(\mathbf{x}) \left| \frac{\partial F_\theta(\mathbf{z})}{\partial \mathbf{z}} \right| = p_{\text{model}}(F_\theta(\mathbf{z})) \left| \frac{\partial F_\theta(\mathbf{z})}{\partial \mathbf{z}} \right| \quad (2.17)$$

for the latent distribution and

$$p_{\text{model}}(\mathbf{x}) = p_{\text{latent}}(\mathbf{z}) \left| \frac{\partial F_\theta(\mathbf{z})}{\partial \mathbf{z}} \right|^{-1} = p_{\text{latent}}(\bar{F}_\theta(\mathbf{x})) \left| \frac{\partial \bar{F}_\theta(\mathbf{x})}{\partial \mathbf{x}} \right| \quad . \quad (2.18)$$

It is worth noting that computing the likelihood (2.18) includes computing the determinant of a jacobian matrix. Therefore, NFs rely on carefully designed invertible layers (e.g., affine coupling layers [39]) where the determinant of the Jacobian can be easily computed. This enables training via exact maximum likelihood estimation and does not rely on the ELBO loss like VAEs.

2.1.3 Generative Adversarial Network

In contrast to VAEs, NFs or DMs, GANs [53] belong to the class of non-probabilistic generators which do not require an explicit likelihood function. The model, also called the generator G , is trained in an adversarial training setup that contains, on the one hand, the generator that should generate new data points being as realistic as possible and, on the other hand, a discriminator that should be able to distinguish if a data point is generated by the generator (“fake”) or from the original dataset (“real”). The generator and the discriminator are trained in an alternating fashion, where the generator tries to fool the discriminator by generating more realistic data points, and the discriminator D becomes better at distinguishing “fake” from “real”. This is formalized in a MinMax condition

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (2.19)$$

where $G(\mathbf{z})$ is a sample generated by the generator and \mathbf{x} is an element from the original dataset. The prediction of the discriminator is close to one if it identifies its input as “real” and close to zero if it is identified as “fake”. GANs were the dominant class of network architectures for image generation [81, 108, 152] but were overtaken by diffusion models over the last few years. The GANs training process turned out to be often unstable, and extensive hyperparameter tuning and engineering is required [106, 149]. One challenge is to carefully balance the learning dynamics of the generator and the discriminator. If one overpowers the other, e.g., the discriminator always perfectly distinguishes “real” from “fake” samples, the generator stops learning. Another difficulty is known as mode collapse [174], meaning that the generator just produces samples from a part of the original data distribution but fails to completely capture it. A third challenge is the problem of exploding and vanishing gradients [6] which prevents effective learning of the generator.

2.2 Diffusion and Flow Matching Models

Diffusion models [68, 145, 158, 159, 161, 165] are probabilistic generative models that have achieved remarkable success in image generation tasks [2, 93], demonstrating exceptional performance in terms of image quality and image diversity. While they require significant computational resources and have slower inference compared to some alternatives [53, 84, 143], they have become the dominant generative modeling paradigm in the research community. The main principles of diffusion models are reflected in two processes motivated by non-equilibrium thermodynamics, namely the forward process and the reverse process. In the forward process, the original image is gradually perturbed by systematically adding Gaussian noise until only pure noise is left. The reverse process describes the iterative noise removal by the diffusion model until a clean image is obtained. During image generation with a diffusion model, the starting point is randomly sampled noise to which the diffusion model is applied iteratively. This repeated application gradually produces a clean, high-quality image. To reduce the computational load, instead of executing the diffusion process in image space, [145] proposed to first convert the image/noise into a compressed latent space using an autoencoder, where the diffusion process then took place. This class of diffusion models is therefore called latent diffusion models and are widely applied used. Over time, several different formulations of diffusion models were proposed. However, there are three dominant formulations recognized in the literature [32]: Denoising diffusion probabilistic models (DDPMs) [68, 158] and its further development denoising diffusion implicit models (DDIMs, score-based generative models (SGMs) [161] and stochastic differential equations (SDEs) [165].

Recently, a generalization of DMs was proposed, known as flow matching or Flow-based models [107, 114]. Unlike classical diffusion, flow matching models directly learn a vector field that defines the path between the noise distribution and the data distribution. By enforcing a straight path trajectory, the sampling process becomes computationally more efficient. This formulation was successfully applied in several recent models like Flux-dev [93] or Sana [184].

In the following, the different formulations of diffusion models, DDPM, DDIM, SGM and SDE, are presented followed by an introduction to the concept of flow matching.

2.2.1 Denoising Diffusion Probabilistic Models

Unless stated otherwise, the discussion and mathematical formulations of DDPM and DDIM presented in this chapter are based on the lecture notes by Inga Strümke and Helge Langseth [167].

In DDPMs [68] the forward and reverse process are modeled as Markov Chain (MC), meaning that each step only depends on the immediate previous step. This key concept is illustrated in fig. 2.1. Let $\mathbf{x}_0, \dots, \mathbf{x}_T$ be the states of the MC with \mathbf{x}_0 representing the clean image and \mathbf{x}_T representing pure noise, then one step of the reverse process is given by $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$. This shows that the slightly denoised state \mathbf{x}_{t-1} just depends on state \mathbf{x}_t and θ denotes the parameters of the trained diffusion model being used to predict \mathbf{x}_{t-1} . The forward process, which progressively adds noise proceeds from right to left in fig. 2.1 and is given by $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ where no learned parameters are required.

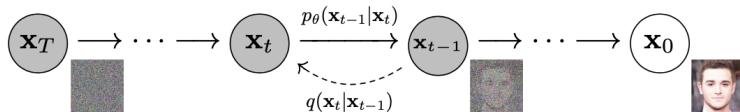


Figure 2.1: Visualization of the forward and backward process of DDPM taken from [68].

Forward Process

Formally, the forward process is defined as a MC that progressively injects noise into the clean image $\mathbf{x}_0 \sim p_{\text{data}}$ which is transformed such that for a sufficiently large number of steps T the final latent variable \mathbf{x}_T approximately follows a new often simpler prior distribution p_{prior} . This boundary condition is necessary so that in the reverse process, which starts with a sample of p_{prior} , the model is able to transfer it back to the image domain p_{data} . The transformation is governed by a variance schedule $\{\beta_t\}_{t=1}^T$, e.g., a linear [68], a cosine [126] or a learnable [83] schedule, to ensure smooth interpolation between p_{prior} and p_{data} . The schedule dictates the incremental increase of variance β_t .

Concretely, the intermediate state \mathbf{x}_t within the MC follows the distribution

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = K(\mathbf{x}_t; \mathbf{x}_{t-1}, \beta_t) \quad (2.20)$$

with the Markov kernel K , which is chosen as a Gaussian Markov diffusion kernel in DDPM leading to the explicit form

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) := \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}\right) \quad (2.21)$$

and the prior distribution $p_{\text{prior}}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$.

A key property of this Gaussian formulation is the existence of a closed-form expression for sampling any \mathbf{x}_t directly. At every timestep t the intermediate latent variable is computed by

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \mathbf{z}_t \quad (2.22)$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t(1 - \alpha_{t-1})} \mathbf{z}_{t-1} + \sqrt{1 - \alpha_t} \mathbf{z}_t \quad (2.23)$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \mathbf{z}_{t-1:t} \quad (2.24)$$

$$= \dots \quad (2.25)$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathbf{z}_{0:t}. \quad (2.26)$$

where $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ which shows that $q(\mathbf{x}_t | \mathbf{x}_0)$ follows a normal distribution

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right). \quad (2.27)$$

This underlines that during the forward process the mean of \mathbf{x}_t is gradually moved towards zero and the variance increases towards one. Therefore, the closed-form solution for \mathbf{x}_t is

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \boldsymbol{\epsilon} \quad (2.28)$$

with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This possibility to directly sample every \mathbf{x}_t is crucial for efficient training.

Reverse Process

While the forward process was all about adding noise to the clean image the reverse process aims recover the clear image by iteratively removing noise. Feller [46] showed that in the limit of infinitesimal steps the reverse process retains the same distributional form as the forward process, a Gaussian. Particularly, in the reverse process $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ is the quantity of interest and given by

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t) = q(\mathbf{x}_t | \mathbf{x}_{t-1}) \frac{q(\mathbf{x}_{t-1})}{q(\mathbf{x}_t)}. \quad (2.29)$$

However, it is not tractable because computing the marginal distributions $q(\mathbf{x}_{t-1})$ and $q(\mathbf{x}_t)$ requires the integration over the whole distribution $q(\mathbf{x}_0)$. Therefore, it is approximated by a neural network parameterized by θ that should learn the Gaussian distribution for each step. Although $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is not tractable $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is which is leveraged during training. The joint distribution learned by the model is given by

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (2.30)$$

where $p(\mathbf{x}_T)$ is pure noise and $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is a Gaussian with learned mean and covariance

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \quad . \quad (2.31)$$

As seen in the forward process, the variance of the noise follows a predefined schedule which means that the neural network just has to learn the mean $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ and the variance is fixed by $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$. Commonly the network does not directly predict the mean $\boldsymbol{\mu}_\theta$ but the noise $\boldsymbol{\epsilon}_\theta$ from which the mean can be computed.

Objective Function

Diffusion models belong to the class of probabilistic generative models. A natural choice for the objective of these models is the negative log-likelihood $-\log p_\theta(\mathbf{x}_0)$. As discussed in sec. 2.1.1 the likelihood is not always tractable. Therefore, the ELBO loss is used as learning objective. In the following the ELBO loss is derived for diffusion models. This derivation is based on [130] and [68]. Starting with the expectation of the negative log-likelihood under the data distribution $q(\mathbf{x}_0)$ typically equivalent to the empirical data distribution p_{prior}

$$-\mathbb{E}_{q(\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0)] \quad (2.32)$$

$$= -\mathbb{E}_{q(\mathbf{x}_0)} \left[\log \left(\int d\mathbf{x}_1 \dots d\mathbf{x}_T p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) \right) \right] \quad (2.33)$$

$$= -\mathbb{E}_{q(\mathbf{x}_0)} \left[\log \left(\int d\mathbf{x}_1 \dots d\mathbf{x}_T p(\mathbf{x}_T) q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) \prod_{t=1}^T \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right) \right] \quad (2.34)$$

$$= -\mathbb{E}_{q(\mathbf{x}_0)} \left[\log \mathbb{E}_{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \left[p(\mathbf{x}_T) \prod_{t=1}^T \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \right] \quad (2.35)$$

using Jensen's inequality

$$-\mathbb{E}_{q(\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0)] \leq -\mathbb{E}_{q(\mathbf{x}_0)} \left[\mathbb{E}_{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \left[\log \left(p(\mathbf{x}_T) \prod_{t=1}^T \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right) \right] \right] \quad (2.36)$$

$$= -\mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[\log \left(p(\mathbf{x}_T) \prod_{t=1}^T \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right) \right] \quad (2.37)$$

$$= \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[-\log p(\mathbf{x}_T) - \sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \quad (2.38)$$

$$=: \mathcal{L}_{\text{DDPM}} \quad (2.39)$$

the loss for training a diffusion model is obtained. This can be further transformed to

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[-\log p(\mathbf{x}_T) - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} - \log \frac{p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}{q(\mathbf{x}_1 | \mathbf{x}_0)} \right] \quad (2.40)$$

$$= \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[-\log p(\mathbf{x}_T) - \sum_{t=2}^T \log \left(\frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)} \cdot \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \right) - \log \frac{p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}{q(\mathbf{x}_1 | \mathbf{x}_0)} \right] \quad (2.41)$$

$$= \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[-\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T | \mathbf{x}_0)} - \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)} - \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1) \right] \quad (2.42)$$

$$= \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[\text{KL}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T)) + \sum_{t=2}^T \text{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)) - \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1) \right] . \quad (2.43)$$

For training, only non-constant terms are relevant. That is why the loss further simplifies to

$$\mathcal{L}_{\text{DDPM}} = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t)} [\text{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))] - \mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} [\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)] + \text{const} \quad (2.44)$$

$$\approx \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t)} [\text{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))] + \text{const} \quad (2.45)$$

assuming that the term $\mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} [\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)]$ is negligible. Note that in eq. 2.45 only Gaussians, $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$ with $\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$ and $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \hat{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$ are compared in the KL-divergence meaning there is a closed-form solution. The loss for the diffusion model boils down to

$$\mathcal{L}_{\text{DDPM}} = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t)} \left[\frac{1}{2\sigma^2} \|\hat{\mu}(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] . \quad (2.46)$$

The mean and the variance of the forward process $\hat{\mu}$ is obtained by

$$\hat{\mu}(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \quad (2.47)$$

$$\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t . \quad (2.48)$$

However, instead of predicting the mean, diffusion models are more often trained on predicting the noise ϵ_θ because [68] empirically found that this lead to more stable training

$$\mathcal{L}_{\text{DDPM}} = \sum_{t=2}^T \frac{1}{2\tilde{\beta}_t} \frac{(1 - \alpha_t)^2}{\alpha_t (1 - \bar{\alpha}_t)} \cdot \mathbb{E}_{q(x_t | x_0)} [\|\epsilon_\theta(x_t, t) - \epsilon_t\|_2^2] \quad (2.49)$$

which is achieved by reparameterization of eq. 2.46 via

$$\hat{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right) . \quad (2.50)$$

In addition, they found empirically that ignoring the weighting term of eq. 2.49 worked better

$$\mathcal{L}_{\text{DDPM}} = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}_t\|_2^2] . \quad (2.51)$$

DDPM Sampler

Algorithm 1 and fig. 2.2 depict the training process of the DDPM approach. In each training iteration, an image $\mathbf{x}_0 \sim p_{\text{data}}$ is sampled from the data distribution. Next, a timestep t is uniformly sampled from a predefined interval $[0, \dots, T]$ and Gaussian noise is drawn according to $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Using eq. 2.28 the noisy image \mathbf{x}_t at timestep t is computed. The diffusion model is applied to predict the noise from this noisy input, and the loss function (eq. 2.51) measures the difference between the predicted and true noise. Finally, the model parameters are updated using gradient descent or a more advanced iterative optimization algorithms.

Algorithm 1 Training of DDPM using learning rate η [167]

```

repeat
     $\mathbf{x}_0 \sim p_{\text{data}}$ 
     $t \sim \text{Uniform}(\{1, \dots, T\})$ 
     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
     $\mathbf{x}_t \leftarrow \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}$ 
     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \cdot \nabla_{\boldsymbol{\theta}} \|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|_2^2$ 
until converged

```

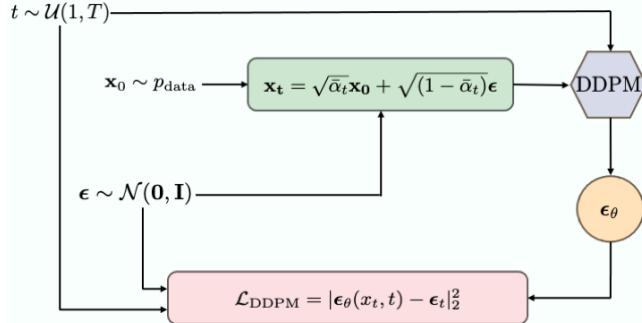


Figure 2.2: Visualization of the training procedure for DDPM. Taken from [16].

Algorithm 2 and fig. 2.3 show schematically the sampling process which reverses the diffusion process to generate images from noise. Starting with pure Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the process iteratively denoises the image. At each timestep t , the slightly denoised image \mathbf{x}_{t-1} is computed using

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z} \quad (2.52)$$

until $t = 2$. In the last step, the final clear images is obtained by

$$\mathbf{x}_0 = \frac{1}{\sqrt{\alpha_1}} \left(\mathbf{x}_1 - \frac{1 - \alpha_1}{\sqrt{1 - \bar{\alpha}_1}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_1, 1) \right) + \sigma_1 \mathbf{z} \quad (2.53)$$

with $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Here, it becomes clear that the model has to predict the noise T times which makes inference of such a diffusion model slow.

Algorithm 2 Sampling of DDPM [167]

```

 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
for  $t = T, \dots, 1$  do
    if  $t > 1$  then
         $\lambda \leftarrow 1$ 
    else
         $\lambda \leftarrow 0$ 
    end if
     $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
     $\mathbf{x}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \lambda \cdot \sigma_t \cdot \mathbf{z}$ 
end for
return  $\mathbf{x}_0$ 

```

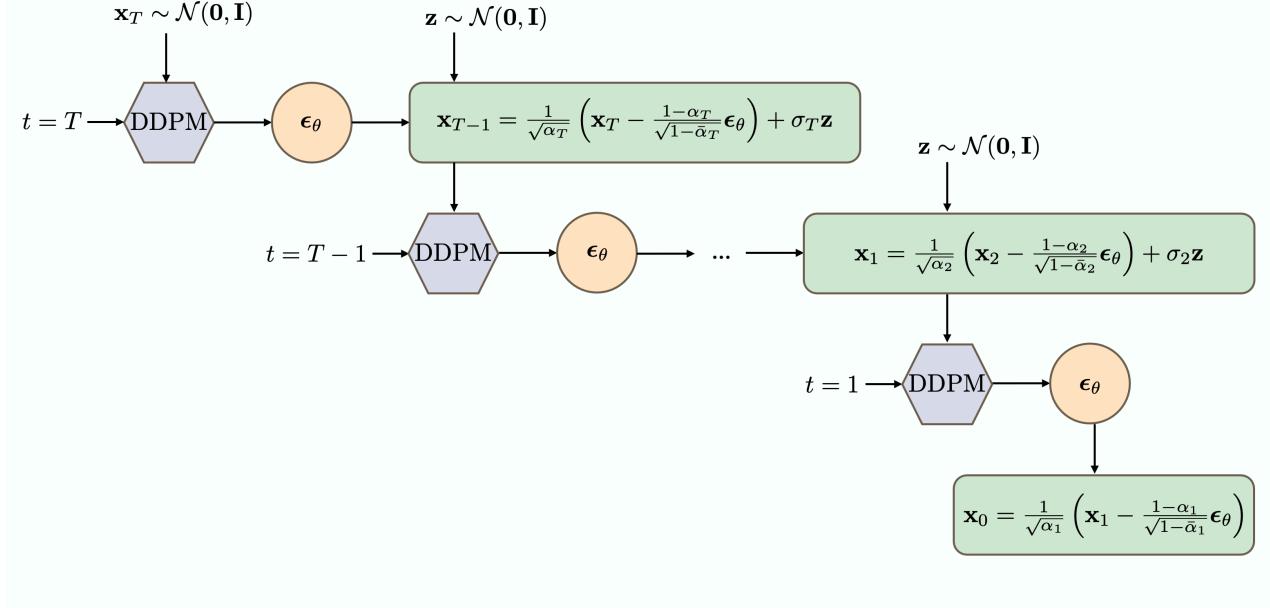


Figure 2.3: Visualization of the iterative sampling process of DDPM. Taken from [16].

2.2.2 Denoising Diffusion Implicit Models

The DDPM sampling process can be accelerated by replacing the Markovian diffusion process with a non-Markovian alternative resulting in denoising diffusion implicit models (DDIMs) [159]. Unlike DDPM, DDIM defines the joint distribution as

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = q(\mathbf{x}_T \mid \mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \quad . \quad (2.54)$$

It becomes clear that the forward process $q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)$ is not Markovian anymore due to its explicit dependence on \mathbf{x}_0 . The authors of [159] parameterize the reverse process as

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \cdot \mathbf{I} \right) \quad . \quad (2.55)$$

The corresponding sampling equation becomes

$$\mathbf{x}_{t-1} = \underbrace{\frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_\theta(\mathbf{x}_t, t))}_{\text{Predicted } \mathbf{x}_0} + \underbrace{\sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\mathbf{x}_t, t)}_{\text{Direction of } \mathbf{x}_t} + \underbrace{\sigma_t \cdot \mathbf{z}}_{\text{Noise}} \quad (2.56)$$

where σ_t is a parameter to choose. \mathbf{x}_{t-1} consists of three parts. The first part is a prediction of \mathbf{x}_0 based on ϵ_θ . The second part is a directional term describing the transition between \mathbf{x}_t and \mathbf{x}_0 and the last part adds Gaussian noise scaled by σ_t . The key advantage of DDIM is the flexibility in choosing σ_t , which does not need to follow a predefined variance schedule. Different values of σ_t yield different diffusion processes while using the same trained model. In particular, for the choice

$$\sigma_t = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} \sqrt{1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}}} \quad (2.57)$$

the DDIM approach becomes Markovian again and here, the deep connection between DDPM and DDIM can be seen. Another aspect to note is that if setting $\sigma = 0$ the whole process becomes deterministic.

2.2.3 Score-Based Generative Models

This discussion is mainly based on [32, 186]. Score-Based Generative Models (SGMs) [162, 163] do not learn directly the data distribution $p_{\text{data}}(\mathbf{x})$ but the score of the probability function $\nabla_x \log p_{\text{data}}(\mathbf{x})$. The score represents a vector field that points to the steepest increase in log probability density, thereby, guiding the data generation process to regions of higher data density. There are several approaches on how the score function can be learned [54, 164, 166]. Here, the one closest to the DDPM approach is presented [162]. During training, the data is perturbed by Gaussian noise at different noise levels and a neural network is trained to estimate the score function based on the noise level $s_\theta(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q(x_t)$. Perturbing the data is necessary because the score estimation is inaccurate in low-density regions and to ensure the process does not get stuck. These networks are called Noise-Conditional Score Networks (NCSNs). The learning objective boils down to

$$\mathbb{E}_{T \sim \mathcal{U}[1, T], \mathbf{x}_0 \sim q(\mathbf{x}_0), \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} [\lambda(t) \sigma_t^2 \| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t) \|^2] \quad (2.58)$$

$$= \mathbb{E}_{T \sim \mathcal{U}[1, T], \mathbf{x}_0 \sim q(\mathbf{x}_0), \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} [\lambda(t) \sigma_t^2 \| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) - s_\theta(\mathbf{x}_t, t) \|^2] + \text{const} \quad (2.59)$$

$$= \mathbb{E}_{T \sim \mathcal{U}[1, T], \mathbf{x}_0 \sim q(\mathbf{x}_0), \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} \left[\lambda(t) \left\| -\frac{\mathbf{x}_t - \mathbf{x}_0}{\sigma_t} - \sigma_t s_\theta(\mathbf{x}_t, t) \right\|^2 \right] + \text{const} \quad (2.60)$$

$$= \mathbb{E}_{T \sim \mathcal{U}[1, T], \mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\lambda(t) \|\epsilon + \sigma_t s_\theta(\mathbf{x}_t, t)\|^2] + \text{const} \quad . \quad (2.61)$$

In eq. 2.61 the relation to DDPM loss becomes visible by identifying $\epsilon(\mathbf{x}_t, t) = -\sigma_t s_\theta(\mathbf{x}_t, t)$ [186], showing that learning to predict noise is equivalent to learning the score function.

For sampling a new element from the data distribution [162] proposed Annealed Langevin Dynamics (ALD) which produces samples by only using the score function being approximated by the learned neural network

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \alpha \nabla_{\mathbf{x}} \log p(\mathbf{x}_{t-1}) + \sqrt{2\alpha} \mathbf{z}_t, \quad 1 \leq t \leq T \quad (2.62)$$

with $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and α regulates the step size of the update in the direction of the score. By applying eq. 2.62 iteratively to initial Gaussian noise, a new element of the data distribution is generated.

2.2.4 Stochastic Differential Equations

This discussion is mainly based on [32, 186]. Stochastic Differential Equations (SDEs) [166] are a generalization of SGMs and DDPMs to continuous time where the forward and the reverse process are modeled by the solution of stochastic differential equations. The forward process is described by the SDE

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w} \quad (2.63)$$

where $\mathbf{f}(\mathbf{x}, t)$ is the drift coefficient, $g(t)$ the diffusion coefficient and \mathbf{w} a standard Wiener process (aka Brownian Motion). [4] proved that the diffusion process in eq. 2.63 can be reversed leading to the reversed-time SDE given by

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q(\mathbf{x}_t)] dt + g(t)d\tilde{\mathbf{w}} \quad (2.64)$$

with a standard Wiener process $\tilde{\mathbf{w}}$ where the time flies backward. Both, forward and reverse SDE, share the same marginals. A score neural network s_θ is trained similarly to SGMs by generalizing the objective function (see eq. 2.58) to continuous time

$$\mathbb{E}_{t \sim \mathcal{U}[0, T], \mathbf{x}_0 \sim q(\mathbf{x}_0), \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} \left[\lambda(t) \|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)\|^2 \right] \quad (2.65)$$

such as $s_\theta(\mathbf{x}, t) \approx \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$. For converting noise to data, a solution of eq. 2.64 has to be computed using the score neural network and a SDE solver. Another finding of [166] is that for eq. 2.64 an ODE exists having the same marginal as the SDE being known as probability flow ODE in the literature. It's given by the deterministic version of eq. 2.64

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log q(\mathbf{x}_t) \right] dt \quad (2.66)$$

Both reversed-SDE or probability flow ODE can be used to generate new samples using SDE solvers [78, 166] or respectively ODE solvers [80, 116, 166, 195].

2.2.5 Flow Matching

The concept of flow matching [107] builds upon continuous normalizing flows (CNFs) [25]. The key elements of CNFs are the probability path $p_t(x)$ connecting two distributions p_0 and p_1 , its corresponding vector field $\mathbf{u}_t(\mathbf{x})$ and the flow $\phi_t(\mathbf{x})$ that is generated by the vector field and given by the ODE

$$\frac{d}{dt} \phi_t(\mathbf{x}) = \mathbf{u}_t(\phi_t(\mathbf{x})), \phi_0(\mathbf{x}) = \mathbf{x} \quad . \quad (2.67)$$

In flow matching, the vector field $\mathbf{u}_t(\mathbf{x})$ is approximated by a neural network $\mathbf{v}_{\theta,t}(\mathbf{x})$ to learn the probability path and enable mapping between both probability distributions. The corresponding flow matching loss function is

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t \sim \mathcal{U}(0, 1), \mathbf{x} \sim p_t(\mathbf{x})} \left[\|\mathbf{v}_{\theta,t}(\mathbf{x}) - \mathbf{u}_t(\mathbf{x})\|^2 \right] \quad . \quad (2.68)$$

However, the vector field $\mathbf{u}_t(\mathbf{x})$ is unknown which is why [107] constructed a sample-based vector field $\mathbf{u}_t(\mathbf{x} | \mathbf{x}_1)$ with $\mathbf{x}_1 \sim q(\mathbf{x})$, where $q(\mathbf{x})$ is the real, unknown data distribution from which only individual samples are known. They showed that this approach leads to identical gradients and is tractable.

Concretely, first the marginal probability path $p_t(\mathbf{x})$ is constructed by a mixture of simpler sample-based conditional probability paths $p_t(\mathbf{x} | \mathbf{x}_1)$

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} | \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1 \quad (2.69)$$

and its corresponding marginalized vector field is defined as

$$\mathbf{u}_t(\mathbf{x}) = \int \mathbf{u}_t(\mathbf{x} | \mathbf{x}_1) \frac{p_t(\mathbf{x} | \mathbf{x}_1) q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1 . \quad (2.70)$$

The trick is not to use the marginal distributions, because the integrals are still intractable. Instead, the loss function (eq. 2.68) is reformulated using the conditional vector field

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), q(\mathbf{x}_1), \mathbf{x} \sim p_t(\mathbf{x} | \mathbf{x}_1)} \left[\|\mathbf{v}_{\theta,t}(\mathbf{x}) - \mathbf{u}_t(\mathbf{x} | \mathbf{x}_1)\|^2 \right] . \quad (2.71)$$

The key insight is that \mathcal{L}_{CFM} leads to the same gradients w.r.t θ as \mathcal{L}_{FM} and can actually be computed due to the tractable conditional vector field.

There are many different possibilities to define the conditional vector field $\mathbf{u}_t(\mathbf{x} | \mathbf{x}_1)$. One choice which demonstrates the advantage of this approach over the diffusion process like defined in sec. 2.2.1 is to use optimal transport displacement interpolation. In this case, the mean $\mu_t(\mathbf{x}) = t\mathbf{x}_1$ and the standard deviation $\sigma_t(\mathbf{x}) = 1 - (1 - \sigma_{\min})t$ for the probability path are defined linearly in time resulting in the conditional vector field

$$\mathbf{u}_t(\mathbf{x} | \mathbf{x}_1) = \frac{\mathbf{x}_1 - (1 - \sigma_{\min})\mathbf{x}}{1 - (1 - \sigma_{\min})t} \quad (2.72)$$

which leads to more straight trajectories of the probability mass. These lines are computationally more efficient compared to the diffusion process with its complex, curved trajectories. The resulting flow can be solved with significantly fewer steps during inference while maintaining high image quality.

Rectified Flow

While standard flow matching connects noise and data via the probability path p_t , it typically relies solely on independent coupling, where noise-image pairs are matched randomly. This often leads to intersection of trajectories during the training phase. The consequence is that the neural network $\mathbf{v}_{\theta,t}(\mathbf{x})$ is presented with two different, contradictory target velocities at the specific location of the crossing. Therefore, [114] proposed an iterative algorithm called Reflow to untangle and straighten the trajectories. This is achieved by iteratively refining the learned rectified flow ODE. This approach offers significant computational advantages, specifically, the straightened paths create a more direct mapping between distributions, making the velocity field easier to learn and enabling efficient sampling with as few as 1-2 Euler integration steps, compared to the 50-1000 steps typically required by diffusion models.

The rectified flow is defined via the ODE

$$d\mathbf{x}_t = \mathbf{v}_{\theta,t}(\mathbf{x}_t) dt \quad (2.73)$$

where $\mathbf{v}_{\theta,t}$ is the vector field approximated by a neural network. The flow should learn to follow the straight trajectory $\mathbf{x}_1 - \mathbf{x}_0$. This is enforced by the learning objective

$$\min_{\mathbf{v}} \int_0^1 \mathbb{E} \left[\|(\mathbf{x}_1 - \mathbf{x}_0) - \mathbf{v}_{\theta,t}(\mathbf{x}_t)\|^2 \right] dt \quad (2.74)$$

where $\mathbf{x}_t = t\mathbf{x}_1 + (1 - t)\mathbf{x}_0$ is the linear interpolation between \mathbf{x}_0 and \mathbf{x}_1 . As depicted in fig. 2.4 (a) crossing trajectories occur when random coupling between samples from p_0 and p_1 is applied. Fig. 2.4 (b) shows that the connections of samples from p_0 and p_1 using the trained model are unwired due to the fact that the vector field of an ODE is uniquely defined at every point. However, the trajectories are not necessarily straight yet. To obtain straight trajectories, new couples $(\mathbf{x}_0^1, \mathbf{x}_1^1)$ are generated where $\mathbf{x}_0^1 \sim p_0$ and \mathbf{x}_1^1 is the solution obtained by applying the learned ODE. These newly acquired pairs are used for retraining the model (see fig. 2.4 (c) and (d)). The advantage is that the trajectories between the new pairs $(\mathbf{x}_0^1, \mathbf{x}_1^1)$ that the network encounters during training are straighter and non-overlapping. By iteratively training and sampling new pairs $(\mathbf{x}_0^k, \mathbf{x}_1^k)$, the

learned trajectories become increasingly straight and have a lower transportation cost. This iterative approach is called Reflow (see alg. 8)

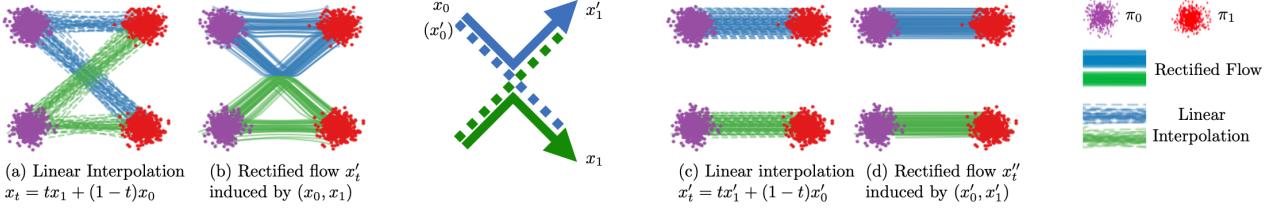


Figure 2.4: Visualization of the path unwiring and straightening of the first iteration in the Reflow algorithm. Taken from [114].

Algorithm 3 Rectified Flow: Main Algorithm [114]

Procedure: $x' = \text{RectFlow}((x_0, x_1))$

Inputs: Draws from a coupling (x_0, x_1) of π_0 and π_1 ; velocity model $v_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ with parameter θ .

Training: $\hat{\theta} = \arg \min_{\theta} \mathbb{E} \left[\|x_1 - x_0 - v(tx_1 + (1-t)x_0, t)\|^2 \right]$, with $t \sim \text{Uniform}([0, 1])$.

Sampling: Draw (x'_0, x'_1) following $dx_t = v_{\hat{\theta}}(x_t, t)dt$ starting from $x'_0 \sim \pi_0$ (or backwardly $x'_1 \sim \pi_1$).

Return: $x = \{x_t : t \in [0, 1]\}$.

Reflow (optional): $x^{k+1} = \text{RectFlow}((x_0^k, x_1^k))$, starting from $(x_0^0, x_1^0) = (x_0, x_1)$.

Distill (optional): Learn a neural network \hat{T} to distill the k -rectified flow, such that $x_1^k \approx \hat{T}(x_0^k)$.

2.2.6 Latent Diffusion Models

Diffusion models operating in pixel space produce images of high quality and diversity. However, a major problem is the high computational effort required during training and inference. This becomes particularly relevant for the generation of high-resolution images. For this reason, latent diffusion models (LDMs) were introduced by [145] and are widely used nowadays [20, 93, 133, 144]. With LDMs, the diffusion process is no longer performed in pixel space but in a compressed latent space (see fig. 2.5), which significantly reduces the computational effort. To perform the compression, an autoencoder is used, whereby additional regularizations are often applied such as KL-regularization [84] or VQ-regularization [177]. The training process of LDMs is divided into two steps. First, the autoencoder $(\mathcal{E}, \mathcal{D})$ is trained to reconstruct images learning to preserve local details and pixel-level information. In a second step, the diffusion model is trained on the compressed latent representations to learn structure and semantics. This procedure has the advantage that the perceptual task handled by the autoencoder and the generative task handled by the DM are learned separately, in contrast to DMs in pixel space where both tasks must be learned simultaneously. The training loss for a LDM is similar to the loss for DMs in pixel space. Following the DDPM formulation the loss for LDMs has the same form like for DMs (see eq. 2.51)

$$\mathcal{L}_{\text{LDM}} = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_t | \mathbf{z}_0)} \left[\|\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t) - \boldsymbol{\epsilon}_t\|_2^2 \right] \quad (2.75)$$

with \mathbf{z} denoting the latents, specifically, \mathbf{z}_0 being the latent representation of the clean image obtained by $\mathbf{z}_0 = \mathcal{E}(\mathbf{x}_0)$, where \mathcal{E} is the encoder of the autoencoder. During inference, random noise is sampled in the latent space and serves as input to the LDM during the iterative denoising process. The completely denoised latent is then fed to the decoder \mathcal{D} , which transforms it back to pixel space. Fig. 2.5 also shows the conditioning mechanism of the LDM via text, semantic maps, etc.. This important aspect based on cross-attention will be discussed in the following section 2.2.7.

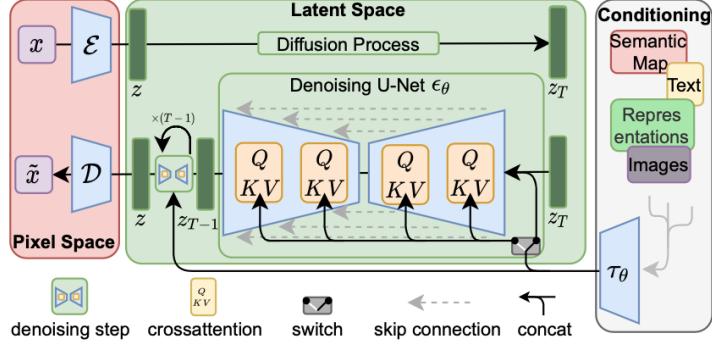


Figure 2.5: Visualizaiton of the structure of a latent diffusion model. Taken from [145].

2.2.7 Conditioning and Classifier-Free Guidance

In DMs, conditioning plays a pivot role in steering the diffusion process through incorporating extra inputs such as text, semantic maps [194], or images [170] users can precisely control the content of the final image. Formally, the diffusion model is trained to learn the conditioned distribution $p(\mathbf{x}|\mathbf{y})$ where \mathbf{y} represents the additional conditioning signal. Consequently, the noise prediction function is modified from $\epsilon_\theta(\mathbf{x}_t, t)$ to $\epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)$. Cross-attention layers (see section 2.3.2) [178] are an effective mechanism for integrating the additional guidance signal into the model. These layers are embedded at various resolutions within the U-net architecture, as shown in fig. 2.5, or within the transformer blocks of modern models like Pixart [20]. To process text prompts, a tokenizer first breaks down the text into discrete subword units, known as tokens, and assigns them a unique numerical ID. Next, the tokens are embedded in a higher-dimensional vector space by a text encoder, and positional encoding (see section 2.3.3) is applied to preserve the sequence order. These embeddings are then used as input for the cross-attention mechanism.

Classifier-Free Guidance

A frequently used method to improve the alignment between the generated image and the text prompt is classifier-free guidance (CFG) [69]. This method builds upon classifier guidance [38], which was introduced for class conditioned diffusion models. In the classifier guidance framework, during training an external classifier predicts the class label and its gradients are used to steer the diffusion model to produce images that can be easily identified by the classifier. However this necessitate the training of an additional classifier which is able to operate on noisy input images as standard classifier are unsuitable for this task. The goal is to steer the diffusion through the additional learning signal to producing images with higher fidelity to the conditioning signal. To facilitate the training process and get rid of the extra classifier CFG was introduced. In the CFG framework, the diffsuion model is jointly trained on conditional $\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y})$ and unconditional $\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y} = \emptyset)$ objectives. For that during training the condition is randomly dropped with the probability $p_{\text{unconditonal}}$. During inference, the same model is used twice to predict the conditional noise and the unconditional noise which are linearly combined to

$$\tilde{\epsilon}_\theta = (1 + w)\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y}) - w\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y} = \emptyset) \quad . \quad (2.76)$$

with w being a hyperparameter called guidance scale.

A high value for w pushes the generated images away from the unconditional distribution and towards the conditional distribution. However, this comes at the cost of higher computational inference time because the model has to be executed twice per timestep. Nevertheless, CFG lead to higher text coherencen and yields higher-contrast and sharper images.

2.3 Fundamentals of Neural Network Architectures

In this section, the focus lies on the key fundamental building blocks of the neural networks.

2.3.1 Transformer

The introduction of the transformer architecture [178] was a decisive moment for the whole AI community. It significantly pushed the development in natural language and image processing forward. Therefore, it is no surprise that the latest models are based on this architecture type [2, 5, 93, 172].

The transformer was introduced as an alternative to recurrent neural networks [148] for processing sequential data like text. In the original paper, the transformer architecture leverages an encoder-decoder structure, however, there are many variants like decoder-only [1, 176] or encoder-only transformer [37].

Transformers, as introduced in the original paper (Fig. 2.6), take text as input, which must first be converted into numerical form through a process called tokenization. This produces a sequence of tokens that are a numerical representation of words or subwords of the text. To provide the model with information about the position of each token in the sequence, positional encodings are added. These preprocessing steps occur before the input is fed into the transformer architecture. Each encoder block consists of two sub-layers. First, a multi-head self-attention mechanism is followed by a residual connection [61] and layer normalization [7]. Second, a feedforward network processes the output further. The decoder blocks contain three sub-layers. First, a masked multi-head self-attention prevents the model from attending to future tokens during the training process, second, a cross-attention mechanism introduces the output of the encoder in the decoder and third an additional feedforward network processes the output of the cross-attention layer. All three sub-layers are followed by a residual connection and layer normalization. The final output is projected to vocabulary-sized logits representing a probability distribution over possible next tokens.

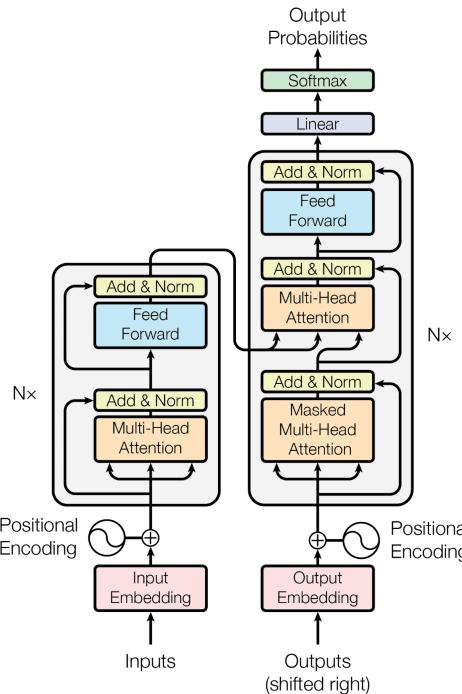


Figure 2.6: Visualization of the transformer architecture. Taken from [178].

2.3.2 Attention Mechanism

The attention mechanism is an integral part of the transformer architecture. It enables the model to focus on the input elements that are most important for the task at hand and ignore the less important ones. This mechanism was first introduced for RNNs [8] to overcome the information bottleneck, as traditionally in RNNs only the last hidden state is given as input for the next prediction. In [8], all previous hidden states were presented to the model as additional input, and the task of the attention mechanism was to filter out those that were relevant for the prediction of the next state. The authors of [178] leveraged the formalism and reformulated it for the transformer architecture that will be discussed in the following.

In the attention mechanism, so-called queries, keys and values are computed. The attention mechanism can be understood as a soft lookup table where queries search for relevant keys, and the corresponding values are retrieved and weighted by similarity. A distinction is made between self-attention (fig. 2.7) and cross-attention (fig. 2.8). In self-attention, the interaction of the tokens within a sequence is calculated. In contrast, in cross-attention you have two input sequences with the goal that the model aligns and relates information between the two different sequences.

Self-Attention

Let $X \in \mathbb{R}^{n \times d}$ be the embedding matrix of the tokens. Then the queries Q , keys K and values V are computed by linear projection of the embedding matrix

$$Q = X \cdot W_q^T \quad (2.77)$$

$$K = X \cdot W_k^T \quad (2.78)$$

$$V = X \cdot W_v^T \quad (2.79)$$

with $W_q^T \in \mathbb{R}^{d \times d_k}$, $W_k^T \in \mathbb{R}^{d \times d_k}$ and $W_v^T \in \mathbb{R}^{d \times d_v}$. Now, so-called scaled-dot product attention is used where the attention matrix A that assigns how much focus is put on the different parts of the token sequence is computed by the softmax of the scaled dot product between the queries and the keys

$$A = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) \quad (2.80)$$

The softmax function is defined as the following for a vector $\mathbf{z} = (z_1, z_2, \dots, z_N)$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (2.81)$$

The scaling of the dot product QK^T with the factor $\frac{1}{\sqrt{d_k}}$ is done in order to prevent vanishing gradients in the softmax function which occurs if the dot product results in very large values. Finally, the attention matrix is applied with the values. This can be summarized in the following formulation

$$Z = \text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (2.82)$$

One main challenge of the discussed attention mechanism is that it has quadratic complexity $\mathcal{O}(n^2)$ in computation and memory consumption which is tackled by several approaches [28, 30, 34, 87].

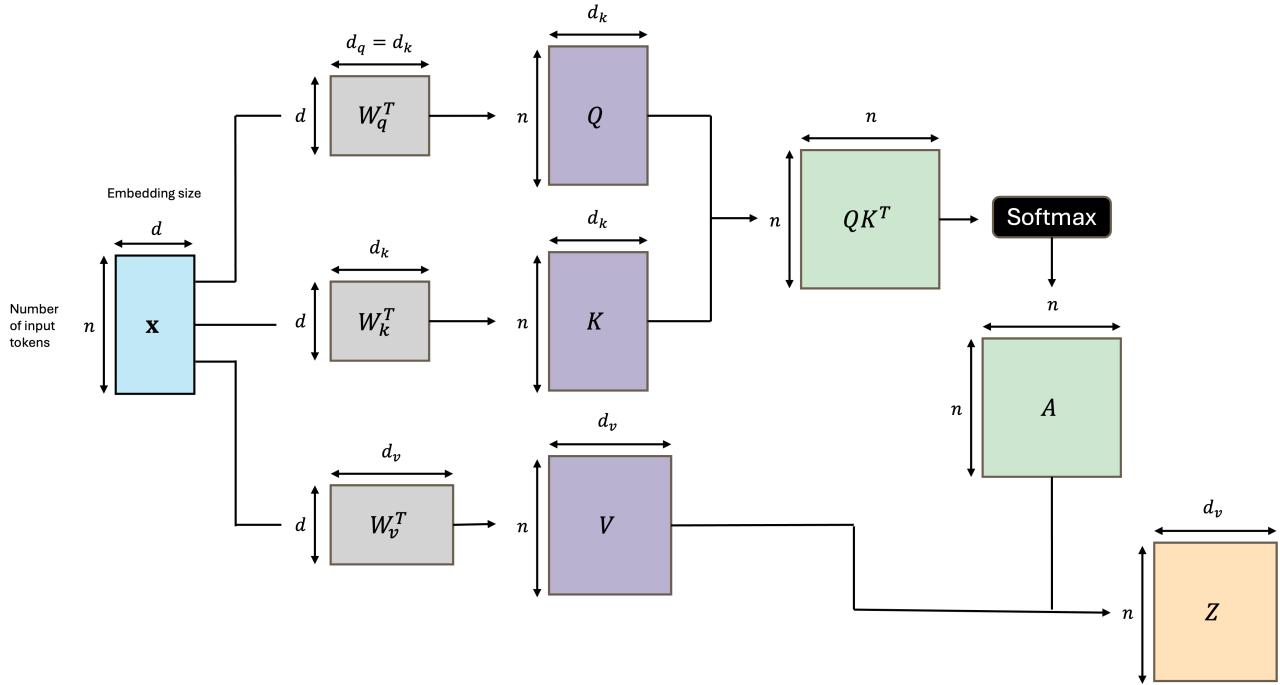


Figure 2.7: Visualization of the self-attention mechanism. Based on [141].

Cross-Attention

The goal of cross-attention is to relate two different sequences of tokens \mathbf{x}_1 and \mathbf{x}_2 . Commonly, the queries come from the target sequence ,e.g., the decoder of the transformer, while the keys and values come from the source sequence, e.g., encoder output of transformer. Apart from this difference, the computation is identical to self-attention.

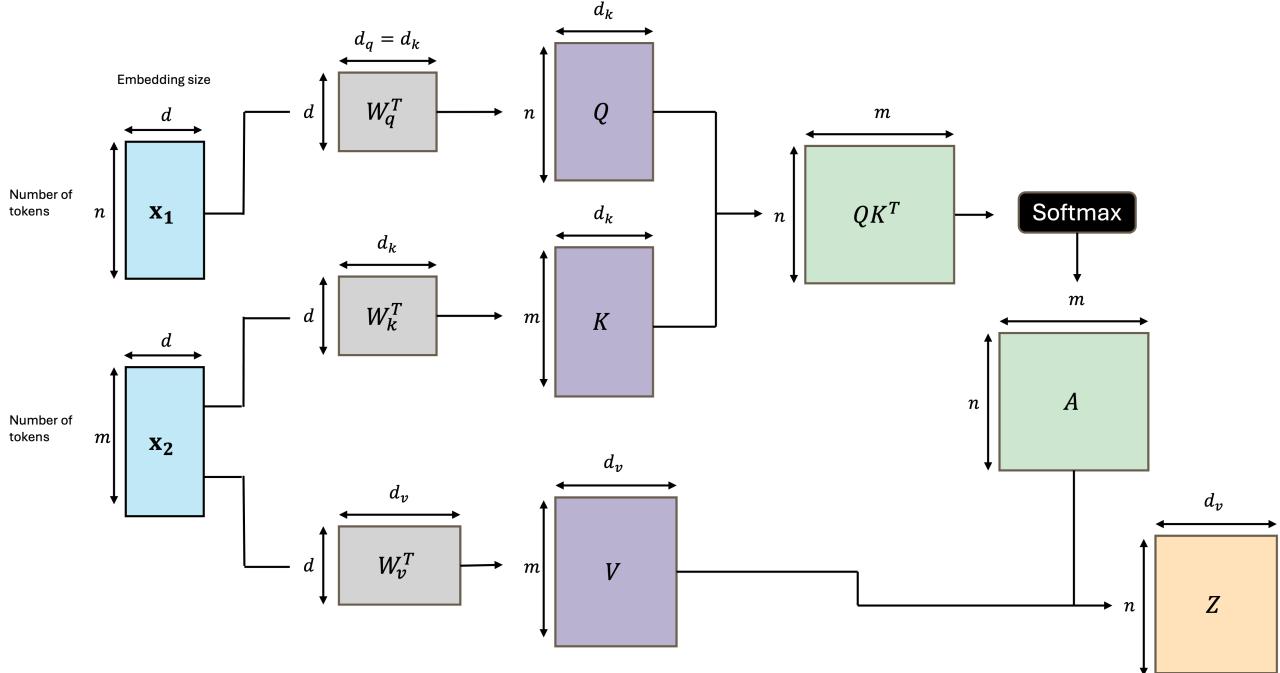


Figure 2.8: Visualization of the cross-attention mechanism. Based on [141].

Multi-Head Attention

Multi-head attention is another trick used in transformers to make the results of the attention operations even more expressive. Instead of only performing individual self- or cross-attentions in a transformer block, so-called multi-head attentions are used. Here, based on the input embeddings, not just one but N times keys, queries and values are calculated, each with different weight matrices $W_{q_i}^T \in \mathbb{R}^{d \times d_k}$, $W_{k_i}^T \in \mathbb{R}^{d \times d_k}$ and $W_{v_i}^T \in \mathbb{R}^{d \times d_v}$ with $i = 1, \dots, N$, and the attention mechanism is executed. Each attention that is carried out is designated as a head _{i} . The results of the N attentions are concatenated and combined with another linear projection $W^O \in \mathbb{R}^{Nd_v \times d}$.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_N) W^O \quad (2.83)$$

$$\text{head}_i = \text{Attention}(XW_{q_i}^T, XW_{k_i}^T, XW_{v_i}^T) \quad (2.84)$$

This is advantageous because the individual heads can concentrate on different tasks, thus reducing the complexity for each head compared to the case where only a single attention has to capture everything.

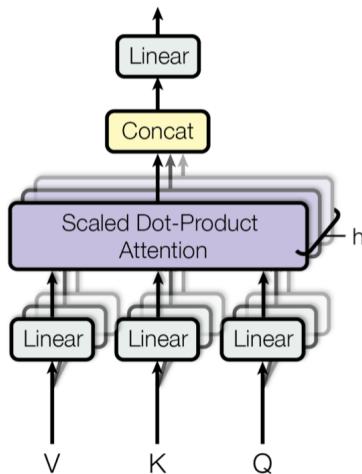


Figure 2.9: Visualization of the multi-head attention mechanism. Taken from [178].

2.3.3 Positional Encoding

Before the embedded token sequence is fed as input to the transformer, a positional information is added to each token to provide the network with information about the order of the tokens. There are different methods, which can be divided into the categories relative [101, 128, 168] versus absolute [88, 104, 115, 178] positional embedding or deterministic [88, 104, 128, 168, 178] versus learnable [101, 115] positional embedding.

Sinusoidal Positional Encoding

Sinusoidal positional encoding is used in the original transformer paper [178] which leverages the sinus and cosine function to inject positional information into the model. It belongs to the category of absolute and deterministic positional encoding.

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right) \quad (2.85)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right) \quad (2.86)$$

where d is the embedding dimension, pos is the position of the token in the sequence and i is the dimension index which is used to compute different frequencies for different parts of the embedding vector. The computed positional embedding has the same dimension like the token embedding vector and is added to it.

2.3.4 Vision Transformer

Original transformers were introduced for language processing tasks. There were some attempts to transfer the methodology of the transformer architecture to computer vision problems [73, 129, 139], but for a long time, convolutional neural networks (CNNs) [140] remained the dominant architecture. CNNs have an inherent advantage over transformers because they possess an image-related inductive bias. On the one hand, CNNs are translationally equivariant by construction and, on the other hand, the concept of locality is embedded in the way they are built. In contrast, these two concepts must be learned by a transformer, as they are not pre-defined by its architecture, which makes a sufficiently large dataset essential to achieve competitive performance.

The introduction of vision transformer (ViT) [41] for a classification task showed that given a sufficiently large dataset transformers can match or even surpass CNNs for vision tasks. The first key aspect in ViT is that not every single pixel is related to all other pixels because this would be computationally infeasible due to the quadratic complexity of the attention mechanism. Instead, the image is split $\mathbf{x} \in R^{H \times W \times C}$ into a sequence of smaller patches $\mathbf{x}_p \in R^{N \times (P^2 \cdot C)}$ where C is the number of channels, H the height of the image, W the width of the image, (P, P) the resolution of a patch and $N = \frac{HW}{P^2}$ the number of patches. At the beginning of the token sequence an additional classification (CLS) token is prepended. The patches are flattened and linearly projected to the encoder dimension. After that a learned positional encoding is added before they are fed into the encoder-only transformer architecture. On top of the transformer encoder a MLP head processes the embedded version of the classification token and provides probability for the different classes from the classification problem.

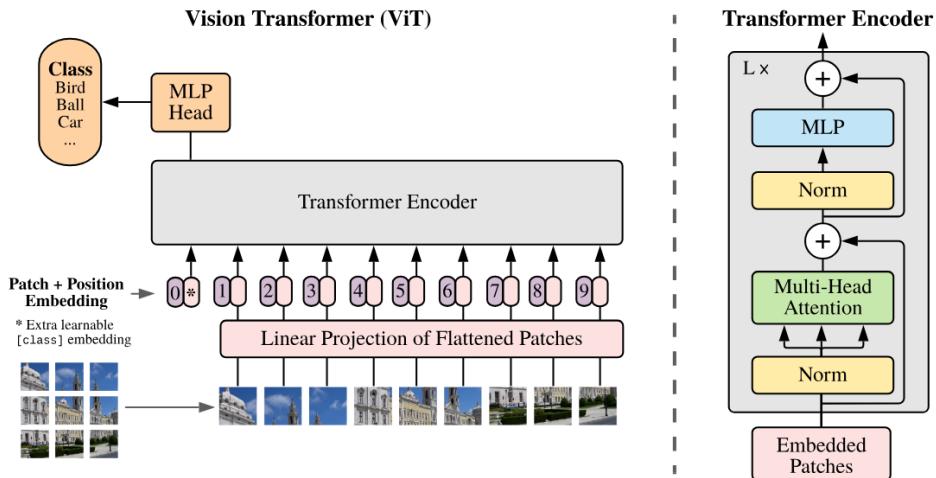


Figure 2.10: Visualization of vision transformer architecture. Taken from [41].

2.3.5 Adaptive Layer Normalization

Layer normalization [7] was introduced as an alternative to batch normalization [74] as a normalization method to reduce the internal covariate shift in deep neural networks independently of the batch size. The covariate shift describes the phenomenon whereby the range of the input to a layer which is the output of the previous layer changes permanently, which hinders effective learning because the layer needs to constantly adapt to a new input distribution. To stabilize and accelerate the training process normalizations are applied in the neural

network. In layer normalization, the activations of each layer \mathbf{z}_l are normalized to a zero mean and unit variance. In addition, a scale γ and a shift η parameter are learned to maintain the flexibility and expressiveness of the network and to ease the constrained range due to the normalization. The normalized layer output is given by

$$\hat{\mathbf{z}}_l = \gamma_l \frac{\mathbf{z}_l - \mu_l(\mathbf{z}_l)}{\sigma_l(\mathbf{z}_l)} + \eta_l \quad . \quad (2.87)$$

A further development of layer normalization occurred when [132] showed that the scale and shift parameters could be used to incorporate additional information into the model by parameterizing them as functions of an additional input c such that $\gamma \rightarrow \gamma(c)$ and shift $\eta \rightarrow \eta(c)$. This method is known as adaptive layer normalization (adaLN) and can be leveraged in diffusion models to incorporate the timestep or class conditioning [131].

2.4 Model Distillation

Diffusion models excel in generating high-quality and diverse images, showcasing remarkable capabilities in image synthesis. However, their main disadvantage is their prolonged inference time which stems primarily from two factors. First, diffusion models rely on an iterative generation process, removing noise step-by-step. To obtain high-quality images often 50-100 steps are required. Second, recent diffusion models like Flux-dev [93] are remarkably large, with billions of parameters, making each step computationally and memory-intensive. These issues limit their use in real-time applications or in memory-constraint environments, such as edge devices.

Diffusion model distillation addresses these challenges. The principal idea of most distillation methods is to transfer the knowledge and image generation capacity from a fully trained diffusion model, often referred to as teacher model, to a target model, often referred to as student model. The research landscape reveals two complementary strategies. First, accelerating inference by reducing the number of denoising steps, and second, creating more compact models through model size reduction. The first strategy, reducing the number of inference steps, is predominant in the literature and several different approaches such as adversarial distillation [151, 153], progressive distillation [105, 150], distribution matching distillation [187, 188], consistency models [117, 160] and guided distillation [122] were proposed. The second strategy, focusing on structural-based knowledge distillation [66] is very suitable for reducing the model complexity and parameter count.

All these distillation methods have in common that there is an efficiency-quality trade-off. The smaller the student model gets and the less inference steps are used the lower is the image quality.

This thesis focuses on reducing the model size of the models of Pixart- Σ and Flux-dev using the concept of knowledge distillation which is presented in the following section in detail.

2.4.1 Knowledge Distillation

The concept of knowledge distillation [66] was originally introduced for training a small model, referred to as the student, on classification tasks such as MNIST [94]. In this framework, a larger model with strong generalization capabilities, referred to as the teacher, serves as guidance. It was shown, that computing the loss between the probability outputs of the teacher and the student enables much more efficient training than relying solely on the original class labels. In such a setting, the small model can be trained on less data or with a higher learning rate. As a possible explanation, the authors hypothesized that even the very small probabilities the teacher assigns to incorrect classes, and particularly their relative magnitudes, carry valuable information for the student. For example, if a "2" from the MNIST dataset is to be classified, it is valuable information to know whether it resembles a "3" or a "5". When only hard targets, e.g., one-hot labels, are used, this relational information is lost. To make this information accessible to the student, one must deviate from the standard temperature $T = 1$ of the softmax function to soften the probabilities. Specifically, both teacher and student use temperature-scaling softmax($\frac{z_i}{T}$) where $T > 1$ making the probability distribution flatter and revealing more information about class similarities. In general, it has been shown that a combination of the distillation loss \mathcal{L}_{KD} , calculated using the outputs from the teacher, and the task loss $\mathcal{L}_{\text{Task}}$, calculated with the hard labels, is advantageous

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda_{\text{KD}} \mathcal{L}_{\text{KD}} \quad . \quad (2.88)$$

This idea has been adapted to train smaller versions of diffusion [82, 96, 193] and flow [17, 33, 45, 119, 173, 199] models. Here, the original model with its pretrained weights acts as the teacher, while the student model is obtained by removing individual parts of the teacher's architecture. Often the student is initialized with the weights of the teacher. The general idea in these approaches is to use the intermediate features as well as the final output as learning signal from the teacher to the student via MSE loss. The mapping between teacher and student layers is established by selecting a subset of the teacher's blocks that corresponds to the reduced depth of the student architecture. If the architectures and therefore the dimensions of the intermediate features differ

between teacher and student model, linear projection matrices ϕ_i are used to project the student's features to the same dimension as the teacher's features. Furthermore, due to different magnitudes in different layers, a weighting factor w_i for each layer i is applied such that the contribution of the different layers are in the same range.

$$\mathcal{L}_{\text{FeatKD}} = \sum_{i \in \text{Layers}} w_i \|F_i^T - \phi_i(F_i^S)\|^2 . \quad (2.89)$$

Additionally, the original diffusion loss $\mathcal{L}_{\text{Task}}$ is taken into account too such that the final loss is given by

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda_{\text{OutKD}} \mathcal{L}_{\text{OutKD}} + \lambda_{\text{FeatKD}} \mathcal{L}_{\text{FeatKD}} \quad (2.90)$$

where $\mathcal{L}_{\text{Task}}$ is the standard diffusion loss, $\mathcal{L}_{\text{FeatKD}}$ is the loss computed between the intermediate features and $\mathcal{L}_{\text{OutKD}}$ is the loss computed between the final predictions of the student and the teacher model.

2.5 Evaluation Metrics

For DMs, several different metrics were introduced to measure the overall image quality, the distance of the generated distribution with the true distribution or the image prompt alignment. In the following, the metrics used in this work are presented: Fréchet Inception Distance (FID) [65], CLIP-Maximum Mean Discrepancy (CMMMD) [77], Human Preference Score v2 (HPSv2) [183]

2.5.1 Fréchet Inception Distance

The FID [65] measures the similarity between the model distribution p_{model} with the true data distribution p_{data} . To compare the distributions a set of images is generated by the diffusion model and another set of images from the original data is needed. An inception model, e.g., Inception-v3 [169], extracts features for every image of the sets. Under the assumption that the features from both sets follow a multivariate normal distribution the mean and covariance of the generated image features ($\mu_{\text{model}}, \Sigma_{\text{model}}$) and the features from the original data ($\mu_{\text{data}}, \Sigma_{\text{data}}$) are computed and compared by the Fréchet Distance [42]

$$\text{FID} = \|\mu_{\text{data}} - \mu_{\text{model}}\|_2^2 + \text{Tr} \left(\Sigma_{\text{model}} + \Sigma_{\text{data}} - 2 (\Sigma_{\text{model}} \Sigma_{\text{data}})^{\frac{1}{2}} \right) \quad (2.91)$$

meaning lower FID values indicate a higher degree of similarity between the generated images and the original images.

Although FID is widely used [77] identified several shortcomings. First, the underlying Inception-v3 model is trained only on images from 1,000 classes which leads to embeddings that often fail to capture the rich and varied content of the images generated by current diffusion models. Secondly, the normality assumption for the feature distributions is frequently violated in reality which can lead to FID scores which do not align to human ratings. A third drawback is that the FID score does not capture complex distortion and depends highly sensitive on the sample size and exhibit bias. For these reasons [77] proposed an alternative score which is discussed next.

2.5.2 CLIP-Maximum Mean Discrepancy

Similar to the FID, the CMMMD [77] score measures the discrepancy between the data and the model distribution. It replaces the Inception-v3 embeddings with Contrastive Language- Image Pre-training (CLIP) [135] embeddings. Since CLIP is trained on more images and leverages natural language supervision, its embeddings have better representation able to capture more semantics details of the images. The features are compared via the maximum mean discrepancy (MMD) [55, 56]

$$\begin{aligned} \text{MMD}^2 &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(\mathbf{x}_{\text{data},i}, \mathbf{x}_{\text{data},j}) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\mathbf{x}_{\text{model},i}, \mathbf{x}_{\text{model},j}) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(\mathbf{x}_{\text{data},i}, \mathbf{x}_{\text{model},j}) . \end{aligned} \quad (2.92)$$

with $\mathbf{x}_{\text{data},i} \sim p_{\text{data}}$, $\mathbf{x}_{\text{model},i} \sim p_{\text{model}}$ and a positive semi-definite Gaussian RBF kernel $k = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$. Through the choice of this kernel no distributional assumption over the features is needed and in contrast to FID score it is an unbiased estimator. [77] showed in several ablation studies that the CMMMD score aligns more closely with human preferences and can detect distortion that FID struggles to identify.

2.5.3 Human Preference Score v2 Benchmark

The HPSv2 benchmark [183] is designed to rank generated images based on human preferences. To model the human preference score, a CLIP model[135] is finetuned on the HPDv2 dataset [183] which contains 434,000 images, generated by nine different text-to-image models or taken from the COCO dataset with over 798,000 human preference choices. A human preference choice always includes two images \mathbf{x}_1 and \mathbf{x}_2 and a binary ranking \mathbf{y} indicating whether image \mathbf{x}_1 is preferred over image \mathbf{x}_2 , e.g. $\mathbf{y} = [1, 0]$ means that \mathbf{x}_1 is preferred by the annotator. The human preference score is determined by the finetuned CLIP model via

$$s_\theta(p, \mathbf{x}) = \mathcal{E}_{\text{txt}}(p) \cdot \mathcal{E}_{\text{img}}(\mathbf{x}) \cdot 100 \quad (2.93)$$

where p is the text prompt, \mathbf{x} is the image, and \mathcal{E}_{img} and \mathcal{E}_{txt} are the image and text encoder of the CLIP model, respectively. The HPSv2 benchmark contains four categories, namely "Animation", "Concept-Art", "Painting", and "Photo" with 800 prompts each. The prompts are derived on COCO captions [26] and DiffusionDB [181] but cleaned and modified by ChatGPT. A higher score computed with eq. 2.93 correspond to better image quality.

2.5.4 GenEval Benchmark

Unlike global metrics such as CLIP and FID, which struggle to capture the compositional correctness of images, the GenEval Benchmark [52] is designed to evaluate text-to-image alignment based on an object-focused framework. It leverages existing models, such as object detection model, semantic model, CLIP, to analyze the content of the image and determine whether it follows the prompt closely. The benchmark is divided into six categories of varying difficulty.

1. Single object: A single object is described by the prompt.
2. Two objects: Two different objects are described by the prompt.
3. Counting: One object should appear several times.
4. Colors: An object should have a specific color.
5. Position: Two or more objects should be placed in a certain order relative to each other.
6. Attribute binding: Correctly associating attributes (e.g., colors) to specific objects (e.g., a red cube and a blue sphere).

For each category, the prompt follows a specific structure. For the benchmark, 553 different prompts are used and for each prompt four images are generated. The scoring is based on a binary classification determining whether the object, number, and color of objects, etc. are correctly present in the image (see fig. 2.11). This returns a ratio representing the proportion of images that contain all described aspects of the prompt.

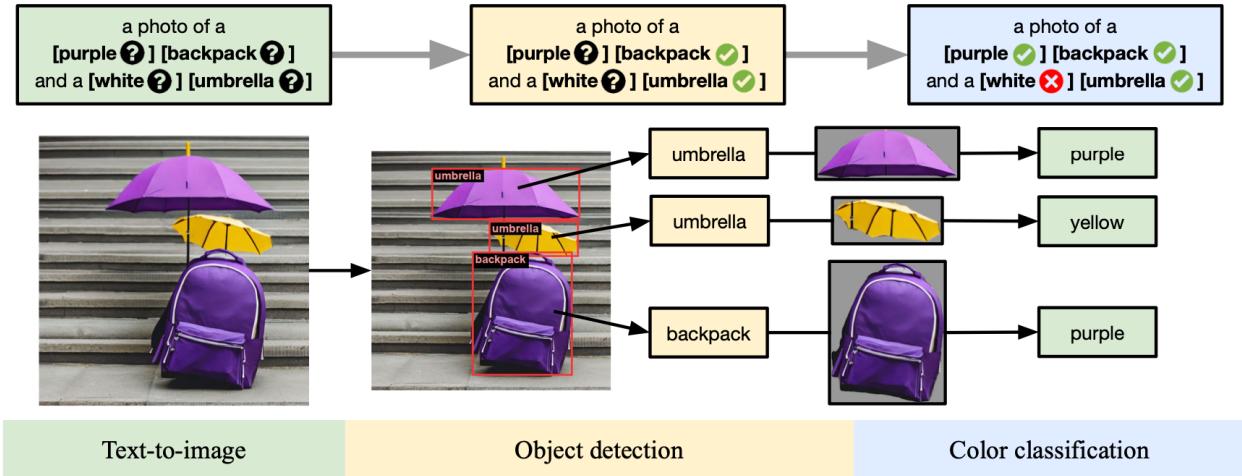


Figure 2.11: Visualization of the evaluation procedure for the GenEval benchmark. Taken from [52].

2.5.5 Dense Prompt Graph Benchmark

The goal of the Dense Prompt Graph (DPG) benchmark [71] is to evaluate the performance of image generation models on long and detailed prompts. It leverages existing datasets, such as COCO [26], PartiPrompts [189] and DSG-1k [29] and enriches their prompts with GPT-4 [1] by adding details, attributes, and spatial relations between objects. A fourth dataset, Object365 [156] is leveraged which is a dataset originally developed for object detection models. The various objects in this dataset are organized into super-categories. Typically, one to four objects from the same super-category are randomly sampled, after which GPT-4 generates a dense prompt based on these objects. In total, 1,065 different prompts were collected.

For evaluation, GPT-4 generates binary questions following the Dynamic Scene Graph (DSG) framework [29] which organizes the questions in a hierarchical tree where questions in child nodes can only be marked as correct if the questions in their parent nodes are also correct. This structure ensures that the answers follow the laws of logic. For example, the parent node normally contains an existence question, e.g. "Is there a motorcycle?" while the child nodes contain questions about its attributes, e.g. "Is the motorcycle blue". Therefore, if there is no motorcycle also the question regarding the color is rendered obsolete. The vision language model mPLUG-large [98] analyzes the generated images with respect to the questions. The evaluation follows a hierarchical aggregation process where the prompt score is defined as the arithmetic mean of all individual question results associated with a specific prompt, while the final DPG-Bench score is calculated by averaging these individual prompt scores across the entire dataset.

Chapter 3

Related Work

Advances in model performance are achieved not only through new algorithms and architectures, but also through the continuous scaling of parameter counts. This scaling paradigm is facilitated due to the continuous advancement in computation hardware and the availability of internet scale training datasets. While massive architecture yields superior generative capabilities, there is a growing need to achieve comparable performance with smaller models. The hardware required for such large models is extremely expensive and power-intensive. The latter aspects leads to a substantial carbon footprint and causes environmental concerns. Furthermore, deploying these models is often infeasible in a ressource constraint setup, e.g. on edge devices. Consequently, a diverse body of literature has emerged around model compression techniques, aiming to minimize the parameter count and the memory footprint during inference and, ideally, also reducing the inference time. This chapter provides a comprehensive overview of various model compression methods, with a particular focus on methodologies adapted for diffusion and flow models as well as methods to accelerate inference. Both is critical for edge device deployment.

Step Distillation

Diffusion and flow models suffer from long inference times due to their iterative image generation process, which requires repeated application of the model. Step distillation methods achieve a significant speed-up of the image generation process by reducing the number of required inference steps down to four or even a single step [117, 150, 151, 153, 160, 188]. In step distillation frameworks, the teacher and the student model often share the same architecture. The teacher is leveraged to retrain the student to generate images in a few inference steps. To this end, several different approaches have been proposed.

One foundational approach is progressive distillation [150] where the number of inference steps is iteratively halved and the previous student model is utilized as the teacher model in the next iteration. Consistency models [117, 160] even achieve a one-step image generation process by enforcing the self-consistency property which dictates that all points on a given diffusion trajectory must be mapped to the same clean image. Alternatively, rather than strictly following the teacher’s trajectory, recent methods rely on adversarial or statistical matching. Latent adversarial diffusion distillation (LADD) [151, 153] incorporates a discriminator in a GAN-like training setup to provide an additional learning signal. Similarly, moving away from exact sample regression, distribution matching distillation [188] focuses on matching the entire image distribution. Concretely, it optimizes the student model to produce samples that collectively match the statistical properties of real data, and not just the output of the teacher on specific noise-image pairs.

Importantly, these step distillation approaches are orthogonal to structural model compression. While this Master’s thesis focuses on architectural model reduction via pruning and SVD, step distillation can be seamlessly

integrated with compressed models to achieve even higher inference efficiency which is demonstrated at the example of Flux.1-schnell in Sec. [?].

Quantization

Quantization is a widely adapted technique to reduce model size and speed up the inference process [18, 59, 75]. In quantization, the high-precision weight values e.g., FP32, are mapped to discrete or lower-precision values e.g., INT8, thereby significantly reducing the memory footprint required to store the weights. Quantization methods can be categorized in mainly two different approaches.

First, post-training quantization (PTQ) [63, 75, 124] describes all methods which operate on fully trained models. They only rely on a small calibration dataset to determine the scaling and zero-point parameters for mapping the high-precision values to their discrete low-precision counterparts. These methods are computationally highly efficient as no training of the complete, potentially very large model is necessary.

Second, quantization-aware training (QAT) [11, 44, 62] integrates the quantization process into the training or finetuning process of the model. During the forward pass the weights are approximated with the chosen low-precision format. However, the rounding function is not differentiable. That is why in the backward pass a straight-through estimator (STE) is employed to approximate the gradients. Although the quantized weights are simulated during the forward pass, the high-precision weights are updated during training. Upon completion of the training, the high-precision weights are permanently replaced with the quantized weights. While generally QAT leads to superior performance compared to PTQ the full retraining of the model makes it computationally demanding.

The application of quantization to diffusion and flow-matching models presents unique challenges due to the iterative denoising process, wherein activations are time dependent. Moreover, even minor quantization errors can exert a significant influence because of their accumulation over several inference steps. Consequently, diffusion-specific quantization methods have been developed to address these specific challenges. Methods such as post-training quantization for diffusion models (PTQD) [64] or Q-Diffusion [102] leveraging a timestep aware calibration dataset sampling features across the entire reverse diffusion process to compute optimal scaling factors have been proposed. Building on this [23] successfully transferred these advanced quantization methods to DiTs. A critical challenge inherited from transformer-based architectures is the emergence of extreme activation outliers which is addressed by [100]. They absorb these anomalous values into high-precision, low-rank components successfully compressing models like Flux.1-dev to 4-bit precision. Most recently, [185] demonstrated that the weights of the Flux.1-dev can be even further compressed to 1.58-bit.

Ultimately, quantization serves as an orthogonal technique to model pruning. Once a model’s architecture is structurally pruned, post-training quantization can be applied to further reduce its memory footprint.

Non-Architectural VRAM Reduction

Alongside methods that reduce the model parameter count, several strategies have been developed to minimize VRAM consumption during inference.

An integral part of modern model architectures is the attention mechanism, which scales quadratically with the sequence length, leading to severe memory bottlenecks. Consequently, hardware-aware algorithms such as FlashAttention [34] or xFormers [97] were introduced to optimize memory access patterns and realize more efficient computation for large attention matrices.

Furthermore, modern models such as Flux.1-dev are standardly deployed using 16-bit half-precision formats. A primary limitation of FP16 is its restricted dynamic range such that gradient or activation outliers can easily exceed the maximum representable value, leading to overflow that crashes both inference and training. For this reason, modern GPUs, e.g. the NVIDIA A100, H100, and H200, natively support the BFloat16 format (BF16).

BF16 utilizes the same 16-bit memory footprint as FP16 but preserves the dynamic range of FP32 by sacrificing decimal precision [79]. Fortunately, neural networks have proven to be highly robust against slight precision loss.

A third option for reducing VRAM consumption is to load only the currently active layers into the VRAM (GPU), leaving the rest of the model in the RAM (CPU). This results in slower inference and training, as there is overhead involved in moving the model components to the VRAM and back to the RAM. Moreover, during training the optimizer states and gradients can be offloaded to CPU saving memory. These memory amnagement strategies are natively supported by modern training frameworks from huggingface [58] or deepspeed [138, 142].

Model Pruning

Model pruning [76, 113, 123] focuses on reducing the actual parameter count of models. It can be dissected in mainly two categories: structural and unstrcultural model pruning. In unstructured model pruning [50, 60, 127] individual, less important weights across the model’s architecture are masked out such that they do not contribute to the inference process anymore. The process to identify parameters to prune can be based on several criteria, e.g. random weight selection [112] or magnitude-based pruning [49, 59]. While this allows removing many parameters while maintaining high performance and theoretically fewer operations (FLOPs) are executed per inference, it does not necessarily speed up inference or reduce the VRAM consumption in practice [51, 120, 182]. The reason is that to achieve a reduction in VRAM consumption the pruned, sparse matrices need to be converted into a sparse tensor format, e.g. compressed sparse row (CSR) [10], which stores the non-zero entries and their coordinates leading to a memory overhead that contradicts the efforts to reduce VRAM. Therefore, a high degree of sparsity is required to compensate for the overhead and ensure low memory requirements. Moreover, modern GPUs are highly optimized for dense matrix multiplication which is why running inference with sparse tensors on standard hardware may be slower than using dense matrix multiplication. Nevertheless there are emerging custom kernels and specialized hardware to accelerate unstructured sparsity calculations [51].

This master thesis focuses on the second pruning method, namely structured pruning for transformer-based diffusion and flow models. Unlike unstructured pruning, this method targets coherent model components, such as attention heads, channels, or entire transformer blocks. To recover degraded performance and lost generation capabilities, most methods subsequently apply retraining, specifically knowledge distillation.

Based on the modification strategy the literature can be divided into two main categories, namely approaches that remove model components entirely and those which compresses them.

Most research on compressing Unet-based diffusion models, such as SDXL [133], adopts the approach of strucual component removal [82, 96, 150, 171, 193], yet differs significantly in pruning granularity. For instance, [82] removed predominantly entire middle transformer blocks from the SDXL architecture, identifying them as highly redundant. In contrast, [193] focused on pruning individual transformer layers within the blocks, whereas [96] employed a hybrid strategy that combines block removal with layer-wise pruning for higher compression ratios. Notably, [96] observed that the decoder plays a more critical role than the encoder of the U-Net for the generation process and should therefore be subject to less aggressive compression. To recover the model’s fidelity after pruning, most approaches leverage knowledge distillation [66] to retrain their pruned models. They mainly employ the unpruned model as a teacher to guide the recovery process of the pruned student model. Adressing optimization challenges in this phase, [193] demonstrated that using normalized feature loss terms is advantageous as the magnitude of individual features from different layers may vary significantly. Furthermore, [150] demonstrated that a progressive pruning schedule for SDXL leads to superior results compared to one-shot pruning.

Recently, the landscape of diffusion models has shifted from U-Net to purely transformer-based architectures where compression methods become increasingly critical due to the steep increase in model size. The primary

advantage of these architectures is their dimensional homogeneity. Unlike U-Net architectures which may require down- or upsampling and channel adjustments, entire transformer blocks can be removed without adjusting the feature dimension. Besides straightforward transformer block removal strategies [33, 45, 173], several compression techniques have been proposed using a hybrid depth- and width-pruning paradigm [92, 118]. Regarding depth reduction, complete transformer blocks are removed. To simultaneously address width [118] identified that especially the text-stream and the MLPs of both streams exhibit high redundancy and replaced them with lightweight linear layers. Similarly, [17] proposes replacing blocks with linear layers and instead of performing a full model finetuning only LoRAs on the preceding and subsequent blocks are trained to mitigate the performance loss.

In a distinct approach, [199] leverages a mixture-of-expert (MoE) framework to dynamically activate only a small fraction of the model’s parameters during each inference step. The core advantage of this approach is that it exploits the temporal variability of feature importance, recognizing that different transformer blocks are critical at different timesteps. Their routing mechanism adaptively selects the most important blocks and experts for the current time step, making the model more expressive than standard static pruning methods. However, while the reduction in activated parameters per inference step significantly accelerates inference, it does not necessarily alleviate the VRAM consumption bottleneck. Crucially, since the complete model (including both activated and inactivated parameters) typically remains in memory, the total memory footprint remains close to the original model size, which distinguishes this method significantly from memory-focused pruning methods. Finally, a distinct category of component compression is low-rank factorization, commonly implemented via singular value decomposition (SVD) [197]. This method approximates a dense weight matrix by decomposing it into products of two low-rank matrices, effectively reducing the parameter count while preserving the original architectural topology.

In the domain of LLM SVD has been extensively investigated as static post-training compression strategy [72, 179]. However, the straightforward application of SVD has proven to lead to significant performance drops even at low compression ratios, since LLMs exhibit massive activation outliers and SVD removes components critical for these activations as the singular value selection is solely based on their magnitude [110, 190].

In this work, SVD is revisited within the context of diffusion and flow models. Unlike the zero-shot compression often applied in LLMs, SVD is identified as a highly-effective initialization strategy for compressed model weight matrices before retraining. This work demonstrates that a progressive linear compression scheme utilizing SVD can effectively compress Flux.1-dev, outperforming SOTA pruning approaches on several benchmarks. In addition, it is demonstrated that SVD can be effectively deployed as a training-free compression approach for moderate compression ratios while maintaining acceptable performance and that [110] can be transferred from LLMs to flow models like Flux.1-dev.

Chapter 4

Methods

4.1 Image Generation Models

In the following section, the models, namely Pixart- Σ [20] and Flux.1-dev [93], which were used for the experiments in this work are described in detail.

4.1.1 Pixart- Σ

PixArt- Σ [20], the latest addition to the PixArt family [21, 22], was developed to demonstrate that high-fidelity image generation can be achieved with significantly reduced training costs and model parameters. Pixart- Σ is a further development of Pixart- α . Therefore, the shared architecture of Pixart- α and Pixart- Σ is introduced first. Next, the key points for efficient training of Pixart- α are described that are also of relevance for PixArt- Σ , followed by the specific enhancements implemented in PixArt- Σ .

Architecture

The shared architecture of PixArt- Σ and PixArt- α is based on the class-conditional diffusion transformer DiT-XL/2 [131] which contains 28 transformer blocks. The transformer blocks have primarily three inputs (see Fig. 4.1) which are described in the following.

The first input is the **timestep** t which is projected into a time embedding vector using a 256-frequency embedding, from which a MLP extracts the global scale and shift parameters. The global scale and shift parameters are used as input for the customized adaptive layer normalization called adaLN-single layer, which modulates the hidden states based on the time embedding.

The second input is the **text prompt** that describes the content of the image. It is tokenized into 120 tokens in Pixart- α , respectively 300 tokens in Pixart- Σ , and encoded by the T5-XXL text encoder [137]. A linear projection layer is then applied to reduce the high dimensional output from the T5-XXL encoder after which it is incorporated via cross-attention layers into the model.

The third input is the **latent image** which is encoded via a pre-trained autoencoder, e.g. SDXL-VAE [133] in Pixart- Σ , transforming the image $(3, H, W)$ into latent space $(4, H/8, W/8)$. Patches of size 2×2 are extracted from the latent representation of the image and flattened into a sequence of visual tokens. Finally, 2D sine-cosine positional encoding is applied before feeding them as input to the transformer blocks.

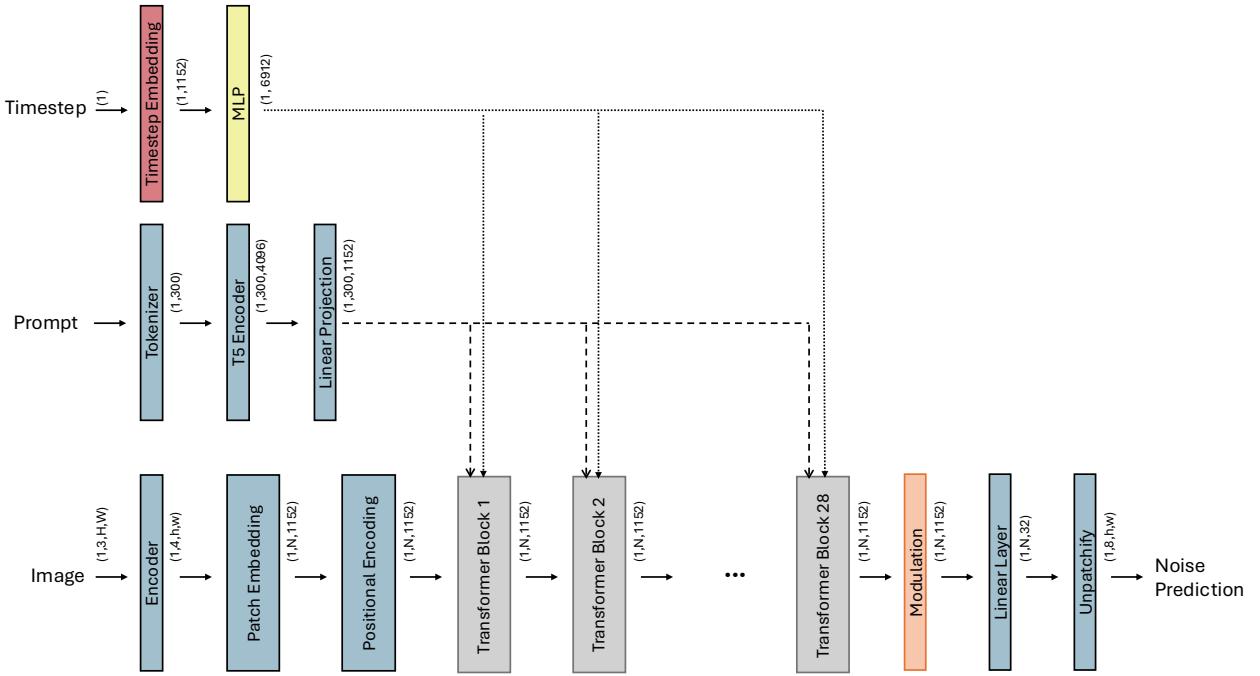


Figure 4.1: Overview of the Pixart- Σ architecture. H and W are the height and width of the image. $h = \frac{H}{8}$ and $w = \frac{W}{8}$ represent the dimensions in latent space and $N = \frac{H}{8 \cdot p} \frac{W}{8 \cdot p}$ with p being the patch size.

The main components of the diffusion model are the transformer blocks. They consist of four main parts, namely the multi-head self-attention layer, the multi-head cross-attention layer, the MLP and the adaLN-single mechanism (see Fig. 4.2).

The **multi-head self-attention** layer allows the tokens to attend to each other and capture spatial relationships. It is embedded between time-dependent modulation layers at the beginning of the transformer block.

The **multi-head cross-attention** layer incorporates text conditioning into the model to align the generated image with the text prompt and is placed between the self-attention and the MLP layer.

The **MLP** processes each token individually to refine the features. Similar to the self-attention layer, it is positioned between modulation layers at the end of the transformer block.

adaLN-single is a modification of adaLN method [132] (see Sec. 2.3.5) replacing the static parameters of the normalization layer with dynamic values that are predicted directly from conditioning signals such as the timestep. This enables the model to efficiently modulate the feature distribution in each block, thereby controlling the generation process globally. In [131] both class and time conditioning were handled via adaLN. In contrast, Pixart- Σ only uses it for incorporating the time embedding. To reduce the number of parameters global scale and shift parameters \bar{S} are extracted from one shared MLP and each block has additional learnable parameters E_i which are added to the global parameters via a summation function g $S_i = g(\bar{S}, E_i)$ to provide flexibility for every transformer block.

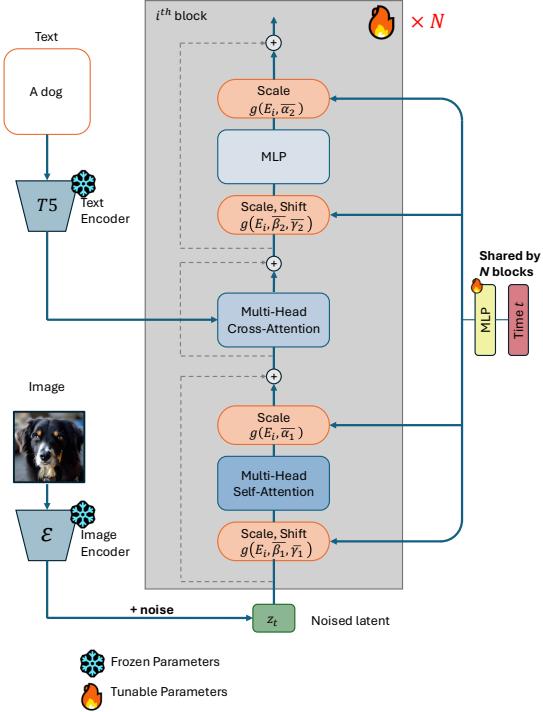


Figure 4.2: Visualization of the transformer block within PixArt- Σ . Taken from [22].

The output of the final transformer block is modulated and processed by a linear layer before the token sequence is reshaped back into spatial dimensions. The final output comprises eight channels as not only the noise but also the variance is predicted. The total number of parameters of individual parts of the transformer blocks and the total model can be found in Tab. 4.1.

Table 4.1: Number of parameters for the individual components of the Transformer block and the total Pixart- Σ model.

Component	Parameters (Mio)
Multi-Head Self-Attention	5.3
Multi-Head Cross-Attention	5.3
Feed-Forward Network (MLP)	10.6
Total (Single Block)	21.3
Total Model (PixArt-Σ)	610.9

Training Strategy

The training of PixArt models proceeds in three distinct stages, explicitly decoupling the learning of pixel dependencies from text-concept alignment to improve training efficiency.

First, the focus is on learning pixel dependencies to generate semantically meaningful images. To achieve this, the model is initialized with weights pre-trained on ImageNet [36] using class-conditional generation. This pre-training provides a strong initialization for the visual distribution at a low computational cost.

In a second step, the focus is on learning text alignment so that the model generates images that accurately reflect the text prompt. For this purpose, the SAM [85] dataset was utilized with corresponding high-density prompts generated by the vision-language model LLaVA [109].

Finally, the model undergoes finetuning on high-resolution and high-aesthetic quality images to refine visual details which was done on a data set that is not publicly accessible.

PixArt- Σ Enhancements

Pixart- Σ is a further development of Pixart- α which is able to generate images at higher resolution, up to 4k, and of higher quality. In this process, the weights from Pixart- α are used as initialization for Pixart- Σ introducing three main changes, namely a more powerful VAE, an improved dataset and a more efficient self-attention mechanism leveraging KV-compression.

First, Pixart- Σ utilizes a more powerful VAE encoder, specifically the VAE from SDXL [133], to obtain more expressive image features.

Second, the image dataset was extended with high resolution images up to 4k. Moreover, a more powerful model, Share-Captioner [24], was leveraged to create captions of higher quality and length. To accommodate this the number of tokens for the text embedding is increased from 120 to 300.

Third, a compression technique for the KV tokens in the self-attention layers is introduced to reduce the computational cost for high resolution images which otherwise scales quadratically with the number of tokens $O(N^2)$. The KV tokens are compressed via a convolutional 2x2 kernel operating in spatial space exploiting the redundancy of feature semantics within a $R \times R$ window. However, the Q tokens are kept uncompressed to mitigate the information loss. Consequently, the computational cost reduces from $O(N^2)$ to $O(\frac{N^2}{R^2})$. The attention operation (compare to Eq. 2.82) changes to

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{Q f_c(K)^T}{\sqrt{d_k}} \right) f_c(V) \quad (4.1)$$

with f_c representing the compression operation.

To conclude, Pixart- Σ shows that leveraging a weaker model as a base to introduce improvements helps to accelerate the training and development process enormously.

4.1.2 Flux.1-dev

Flux.1-dev [93] is a publicly available, guidance-distilled rectified flow model based on the commercial model Flux.1 [93]. Due to the absence of a paper or technical report describing the exact training process, the data set used and the specific design choices, the following section focuses on presenting the model's architecture. The architecture analysis is based on the following git repository [89].

Architecture

Flux.1-dev is based on a multi-modal diffusion transformer (MMDiT) [43] architecture containing 19 double-stream blocks followed by 38 single-stream blocks (see Fig. 4.3) adding up to a total of 11.9 billion parameters (see Tab. 4.2).

For prompt processing, it leverages two distinct text encoders, namely CLIP ViT-L/14 [135] and T5-XXL [137]. For image encoding a VAE is used compressing the height H and width W of the image by a factor of eight. Each double-stream block has mainly four inputs: a guidance vector, the positional information of the image patches as well as the text and the image hidden states. Hereby, the text and the image are processed separately in two distinct streams allowing the model to apply modality-specific weights. This dual-path approach is the defining characteristic of a MMDiT. In contrast, the single-stream blocks have only three inputs as they combine the text and image hidden states and process them together.

The **guidance vector** is a combination of the timestep, a guidance parameter and the pooled vector of the clip embedding. In guidance distillation, the model is trained to simulate classifier free guidance internally. This

speeds up the inference process because in contrast to classifier free guidance where the model is applied twice per inference step only one inference pass per step is needed. Hereby, the guidance is a scalar which corresponds to the guidance scale in standard classifier free guidance. First, the guidance parameter and the timestep are embedded via MLPs into high dimensional tensors which are summed together. Then the pooled clip vector for the text prompt is embedded via another MLP and added. The resulting vector, called guidance vector, is injected via modulation in every double-stream and every single-stream block.

Another input to both block types is the concatenated **positional encoding** for the text prompt and the image. First, based on the latent representation of the image, 3D coordinates for the position of every patch in the image are created. The 3D vector for the image represents the 2D positional IDs (x, y) of the patch in the image and a time dimension t for the possibility to adapt Flux.1-dev for videos. However, for image generation the time dimension is always set to zero. Text position IDs are included for architectural consistency but are initialized to zero, as text tokens lack meaningful 2D spatial coordinates in a latent image grid. In contrast to Pixart models, Flux.1-dev uses rotary positional embedding (RoPE) [168]. The embedded vector is then used in the attention mechanism to rotate the key and query vectors such as the spatial relationships within the image and the sequential structure of the text are preserved.

For the text embedding, the **text prompt** is tokenized into 512 tokens such that the model is capable of processing long and detailed prompts. The T5-XXL encoder embeds the tokens into high-dimensional vectors $(1, 512, 4096)$ and a linear projection is applied before it is fed into the first double-stream block. The subsequent double-stream blocks always receive the hidden states of the text from the previous block.

The **image** $(3, H, W)$ is encoded via a VAE to a latent representation $(16, h = \frac{H}{8}, w = \frac{W}{8})$. To prepare the latents for the transformer, 2×2 neighboring latent pixels are packed into a single patch and therefore increasing the dimension from 16 to 64 during the patchify operation while simultaneously reducing the sequence length N to $\frac{h}{2} \cdot \frac{w}{2}$ $(1, N = \frac{h}{2} \cdot \frac{w}{2}, 64)$. These patches are then linearly projected into a higher dimensional embedding space before being fed into the first double block.

The output, sepcifically the text and image hidden states, of the 19th double-stream block are concatenated and used as single input to the first single-stream block.

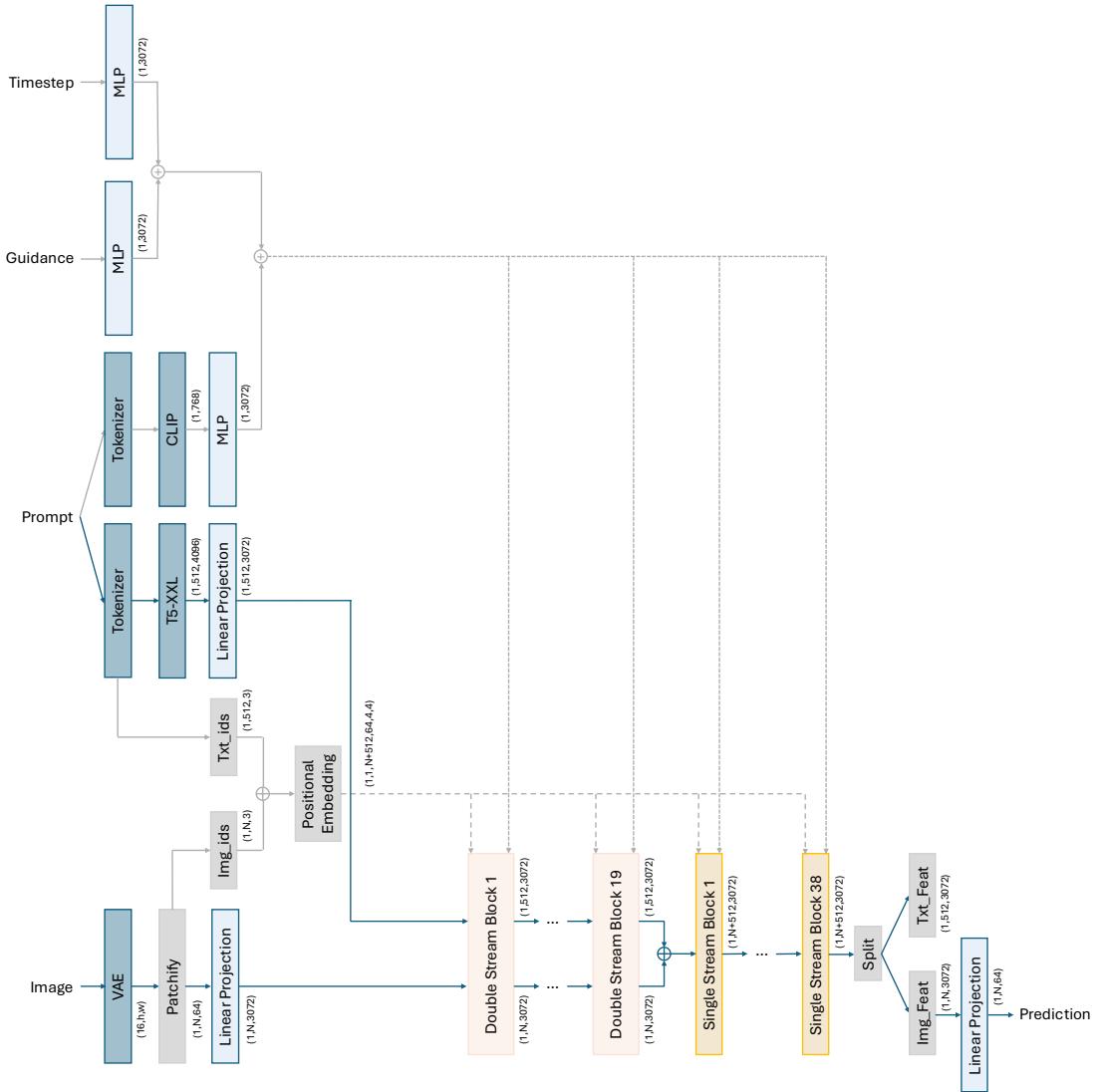


Figure 4.3: Overview of the Flux.1-dev architecture.

Double-Stream Blocks

The double-stream blocks (see Fig. 4.4 on the right) are the fundamental building blocks of the first stage of the Flux.1-dev architecture. In contrast to other diffusion-based models [20, 133, 144] the block consists of two separate streams, which process image and text information separately and only exchange information during the attention process. The blocks consist of three main components: the modulation, which incorporates the guidance vector encoding information about the timestep, the guidance, and the pooled CLIP vector, the attention mechanism, which bidirectionally transfers information from text to image features and vice versa and the MLPs, which further process the features.

The **modulation** of the text and image features is implemented via adaLN. As shown in Fig. 4.4, the guidance vector passes through two separate projection layers, which generate each two sets of scale and shift parameters for adaLN as well as two gating parameters. The scaling and shifting of the features to incorporate the time, guidance and pooled text conditioning is applied twice, once before the attention mechanism and once before the MLP. Subsequently, the gating parameters scale the updated features before they are added to the residual stream. This technique is often applied in very deep neural networks. The linear projection layers producing the gating values are initialized with zeros. Consequently, the gating paramters start at zero, such that the

block initially acts as an identity function, which stabilizes the training.

The **attention** mechanism enables the exchange of information between the text and the image features. First, the text and image features are processed by separate linear layers which project them into query, key and value representations and are normalized via RMSNorm [192]. Subsequently, they are concatenated into a joint sequence. At this stage the RoPEs are integrated into the model by rotating the joint query and key vectors before the attention mechanism is applied such that self- and cross-attention take place simultaneously. Afterwards, the updated text and image features are separated and processed by individual projection layers. The third main parts are the **MLPs** in each stream which process the updated feature vectors further to increase the representational capacity further. Like the attention mechanism, the MLP is also enclosed by a modulation and connected via a gated residual connection.

Single-Stream Blocks

The single-stream blocks (see Fig. 4.4 on the left) are the fundamental building blocks of the second stage of the Flux.1-dev architecture. Their primary components are the self-attention mechanism and the modulation layers. First, the concatenated image and text features X are modulated using adaLN where the shift and scale parameters are obtained by processing the guidance vector through a linear layer. Similar to the double-stream architecture, a gating mechanism is applied after the attention mechanism is executed. To optimize throughput, the model uses a fused linear projection where the hidden states for the attention mechanism and the parallel feed-forward path are computed simultaneously. The outputs of the attention mechanism and the feed-forward path are concatenated together and processed by a subsequent linear projection before the gating mechanism is applied. Finally, the updated features are combined with the residual path to form the final output of the block.

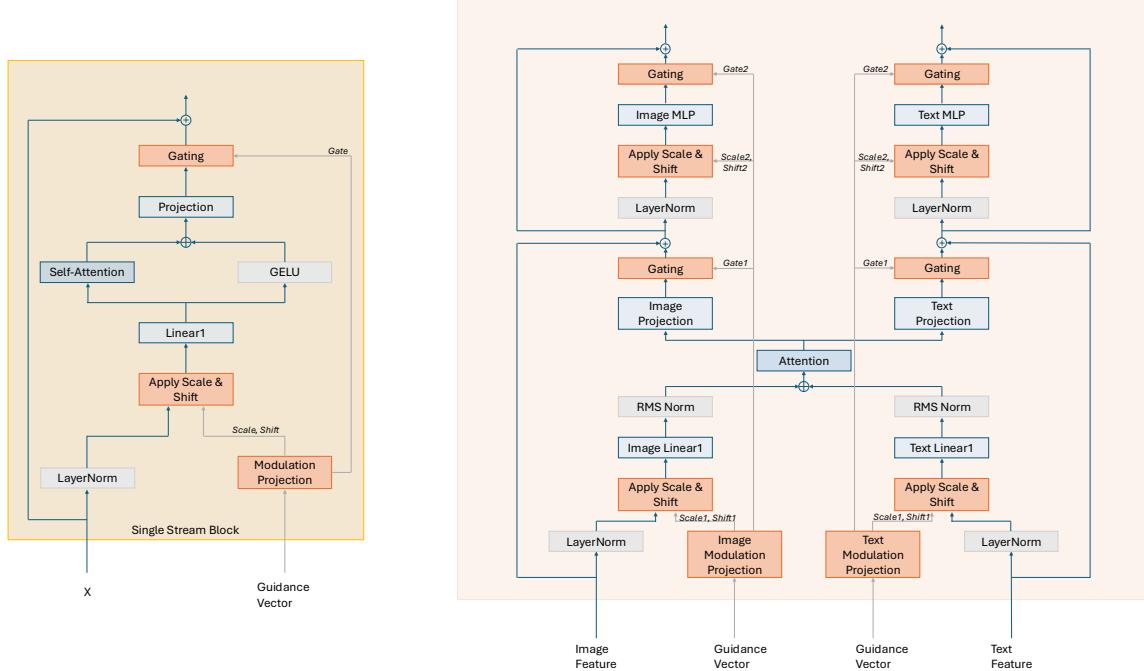


Figure 4.4: Visualization of single-stream block (left) and double-stream block (right) of Flux.1-dev. Graphics are adapted from [17].

The number of parameters of the individual components of a single- and a double-stream block as well as

the total parameter count can be found in Tab. 4.2.

Table 4.2: Number of parameters for the individual components of double-stream block from Flux.1-dev.

Component	Parameters (Mio)
Double-Stream Image Linear Layer (Attention)	28.3
Double-Stream Image MLP	75.5
Double-Stream Image Modulation	56.6
Double-Stream Image Projection	9.4
Double-Stream Text Linear Layer (Attention)	28.3
Double-Stream Text MLP	75.5
Double-Stream Text Modulation	56.6
Double-Stream Text Projection	9.4
Total Double-Stream Block	339.8
Single-Stream Fused Linear Layer (Linear1)	66.1
Single-Stream Projection	47.2
Single-Stream Modulation	28.3
Total Single-Stream Block	141.6
Total Model (Flux.1-dev)	11,901.4

4.2 Block Importance Analysis

Identifying the least important parts of the model which, can be compressed or removed with minimal impact on the final model performance, is crucial for effective structural model pruning. In the following, various block selection criterias are first presented, followed by selection algorithms tackling the challenge of identifying subsets of blocks for simultaneous compression.

4.2.1 Block Selection Criterias

Magnitude-Pruning

In magnitude pruning [47, 60, 95], the magnitudes of the weights within individual model components are computed to serve as a proxy for their relative importance. The underlying assumption is that small weight magnitudes correspond to a low overall impact on the model performance, implying that architectural components with lower weight magnitudes should be prioritized for compression. The aggregated magnitude M_i of a transformer block B_i is defined as the total Frobenius norm of its constituent weight matrices $W \in B_i$. To account for the multiple sub-layers within a block, we compute the block-wise magnitude as

$$M_i = \sqrt{\sum_{W \in B_i} \|W\|_F^2} \quad \text{with} \quad \|W\|_F = \sqrt{\sum_j \sum_k |w_{jk}|^2} . \quad (4.2)$$

Consequently, the block with the smallest magnitude $M_s = \min_i(M_i)$ is identified as the primary candidate for compression.

Representational Similarity Analysis

Another approach to identify redundant model components is to compare their input representations with their corresponding output representations [121]. If the transformation between the input and output is marginal, indicating that the specific component performs an approximate identity mapping, the model part is considered a candidate for removal. To quantify the similarity between the input and output features of a transformer block, the Central Kernel Alignment (CKA) metric [90, 134] is employed. Let $X \in \mathbb{R}^{n \times d}$ denote the input and $Y \in \mathbb{R}^{n \times d}$ the output features to be compared. First, the kernels K and L need to be chosen which are the inner products of the input $K = XX^T$ and the output $L = YY^T$ in the linear case. To ensure invariance to mean shifts, the centering matrix $H = I - \frac{1}{n}11^T$ is applied. The final CKA score is then given by

$$\text{CKA}(K, L) = \frac{\text{HSIC}(K, L)}{\sqrt{\text{HSIC}(K, K)\text{HSIC}(L, L)}} . \quad (4.3)$$

This calculation utilizes the Hilbert-Schmidt Independence Criterion (HSIC) [57] $\text{HSIC}(K, L) = \frac{1}{(n-1)^2} \text{tr}(KHLH)$ [35].

Metric Based Selection

A third approach for identifying the least important model components involves compressing all components, such as each transformer block, independently and evaluating the resulting configurations using image quality metrics on a small reference dataset. Specifically, let $\{B_1, \dots, B_L\}$ denote the individual blocks of a transformer-based model architecture. Each block is modified separately resulting in L distinct model configurations $\{\{\hat{B}_1, \dots, B_L\}, \dots, \{B_1, \dots, \hat{B}_i, \dots, B_L\}, \dots, \{B_1, \dots, \hat{B}_L\}\}$ where \hat{B}_i indicates that the i^{th} block is modified. For each model configuration, the metrics are computed. In this work, the CMMD and/or CLIP-score are employed to rank the importance of the individual model components. While the CLIP-score is used

to evaluate semantic prompt-image coherence, CMMD serves as a measure of general image quality and distributional fidelity. Since CMMD and the CLIP-score produce values in different ranges and have opposing optimization direction (minimization versus maximization), for both metrics a separate ranking of the blocks is created where $r_{\text{CMMD},i}$ and $r_{\text{CLIP},i}$ are the corresponding ranks of block B_i . The final importance score is obtained by computing the aggregated rank $r_{\text{final},i}$ as

$$r_{\text{final},i} = r_{\text{CMMD},i} + r_{\text{CLIP},i} \quad . \quad (4.4)$$

In the linear progressive compression approach (see Sec. 4.6.3) the compression strategy of the individual model parts leads to different numbers of removed parameters due to repeated compression and different compression ratios per block. Therefore, the CMMD value normalized by the number of removed parameters $\frac{\text{CMMD}_i}{\text{Removed Parameters}}$ is leveraged to quantify the degradation of image quality per removed parameter and to rank the blocks to ensure a fair comparison.

4.2.2 Block Selection Algorithms

The necessity of block selection algorithms arises from a specific challenge in metric-based structural pruning. Specifically, the importance of individual model components, e.g. transformer blocks, is typically determined in isolation, meaning only one block is modified, while the remaining parts of the model are held fixed during evaluation. However, in practical applications, several components are compressed simultaneously. It cannot necessarily be assumed that the simultaneous removal of components with the lowest individual impact will yield the optimal combination for preserving image quality. In the following, methods to select the best subset of model components for compression are described, and subsequently evaluated in Sec. 5.1.2.

Greedy Block Selection

The greedy algorithm is applied in conjunction with the metric-based selection criterion. It is an iterative procedure in which the importance of all currently uncompressed blocks is determined to identify and compress the least critical one. Following this step, the block importance is re-evaluated for all remaining uncompressed blocks. Concretely, let the model consist of L transformer blocks $B = \{B_1, \dots, B_L\}$. In the first iteration, the metrics are computed for all L possible model configurations $\mathcal{M}^{(1)} = \{(B \setminus \{B_i\}) \cup \{\hat{B}_i\} \mid 1 \leq i \leq L\}$ and the optimal block B_j is selected according to the ranking. Let \bar{B}_j denote the compressed version of block B_j . In the second iteration, the metrics are re-computed for all $L - 1$ remaining blocks $\mathcal{M}^{(2)} = \{(B \setminus \{B_i, B_j\}) \cup \{\hat{B}_i, \bar{B}_j\} \mid 1 \leq i \leq L, i \neq j\}$. For the n^{th} iteration $L - (n - 1)$ uncompressed blocks remain. Let B_{comp} denote the set of all previously compressed blocks then the potential candidates for compression are given by $\mathcal{M}^{(n)} = \{(B \setminus (B_{\text{comp}} \cup \{B_i\})) \cup \bar{B}_{\text{comp}} \cup \{\hat{B}_i\} \mid B_i \in B \setminus B_{\text{comp}}\}$. This process is repeated until the desired number of blocks to be compressed has been identified.

Black-Box Optimization with Optuna

Optuna [3] is an open-source hyperparameter optimization framework that primarily leverages a Tree-structured Parzen Estimator [12], a Bayesian optimization algorithm, to search the space of all possible hyperparameter settings which in this work corresponds to the discrete combinatorial space of block pruning configurations. It partitions the observed hyperparameter configurations into promising and non-promising candidates based on a user-defined objective function. By modeling these two groups using probability density functions, it draws a new set of hyperparameters that belongs to the promising group with high probability. Through this informed search, the global optimum can be located much faster than with non-adaptive methods such as random or grid search. For an in-depth description of the Optuna algorithm the reader is referred to the original paper [3] as

a more detailed explanation would go beyond the scope of this work.

4.3 Training Data and Image Synthesis

4.3.1 Training Datasets

In this section, the datasets leveraged for retraining PixArt- Σ and Flux.1-dev after pruning are briefly described. Both synthetic, self-generated data sets and pre-existing data sets were used.

LAION Aesthetic

LAION-5B [154] is a dataset containing 5.85 billion web images of varying sizes, of which 2.32 billion have English captions. All images are filtered using CLIP embeddings to prevent that harmful and inappropriate images appear in the dataset. LAION-5B is partitioned into several subsets. Specifically, an aesthetic subset [155] comprising 120 million images was created by training a linear model on top of the CLIP embeddings to predict an aesthetic score. Only images exceeding a certain threshold were added to this subset. In this work, 608k images from the aesthetic subset are utilized as the training dataset. Due to poor quality of the web-scraped captions which come along with the LAION-5B dataset, new synthetic prompts were generated using the visual language model JoyCaption [48]. Five example images of the dataset are displayed in Fig. 4.5.



Figure 4.5: Examples of LAION dataset.

LAION-Pixart- Σ Dataset

The pre-trained Pixart- Σ model is utilized to generate a synthetic dataset of 100k images conditioned on prompts generated by the JoyCaptions model from the aesthetic LAION subset. Tab. 4.3 presents the configuration parameters used in the generation process. Fig. 4.6 shows example images based on the JoyCaption prompts providing a direct visual comparison to the original LAION samples presented in Fig. 4.5.

Table 4.3: Generation Parameters for LAION-PixArt- Σ dataset.

Parameters	Value
Image Height	512
Image Width	512
Guidance Scale	3.5
No. Inference Steps	20
Max Sequence Length	300



Figure 4.6: Examples of PixArt- Σ generated dataset based on JoyCaptions prompts from LAION aesthetic dataset.

LAION-Flux.1-dev Dataset

Following the same methodology as for the LAION-PixArt- Σ dataset a LAION-Flux.1-dev dataset was created. Specifically, 190k prompts generated with JoyCaptions based on the aesthetic subset of LAION-5B were utilized to generate synthetic images using the Flux.1-dev model. The configuration used during the image generation process is detailed in Tab. 4.4. In addition, example images based on the same prompts as those used for the LAION- PixArt dataset are presented in Fig. 4.7.

Table 4.4: Generation Parameters for LAION-Flux.1-dev dataset.

Parameters	Value
Image Height	1024
Image Width	1024
Guidance Scale	3.5
No. Inference Steps	50
Max Sequence Length	512



Figure 4.7: Examples of Flux.1-dev generated dataset based on JoyCaptions prompts from LAION aesthetic dataset.

4.3.2 Evaluation Datasets

This section presents data sets on which the compressed models were evaluated.

Mapillary

The Mapillary [125] dataset consists of 25,000 street level images. Additionally, it contains fine-grained semantic masks for 66 different object classes. In order to achieve high diversity, the images originate from various geographical regions, e.g., Europe, North and South America, Africa, Wider Geographic Oceania, Asia, and were captured during different seasons and at varying times of day with diverse lighting conditions. In addition, no image sequences, which are often recorded by car cameras, were integrated, only individual and unique images to avoid redundancy. Most of the images were taken from the street or sidewalk. These were supplemented

by images taken on highways, in rural areas, and off-road environments. The use of different camera sensors and perspectives contribute to the diversity of the dataset. For quality assurance, only images with at least full HD resolution (1920×1080) that contain only a small amount of motion blur were selected. Due to privacy protection, faces and number plates are blurred.



Figure 4.8: Example images of Mapillary dataset.

Cityscapes

The Cityscapes [31] dataset contains a total of 25,000 images (1024×2048) extracted from stereo-video-sequences of urban street scenes. The images were captured at 50 different cities, in spring, summer or fall. However, there are only images taken during daylight hours and in good to moderate weather conditions. While 20,000 images provide coarse annotations, 5,000 images of the dataset are densely annotated at pixel level covering 30 different object classes. The dataset is highly relevant, e.g., for the automotive sector, because of autonomous driving. However, due to its high degree of standardization it is less diverse than the Mapillary dataset.



Figure 4.9: Example images of Cityscapes dataset.

MJHQ-30k

The MJHQ-30k [99] dataset consists of 30,000 synthetic images from ten different categories (animals, art, fashion, food, indoor, landscape, logo, people, plants, vehicles) each containing 3,000 images to ensure a balanced distribution. These images were generated using Midjourney 5.2 [99] and selected based on their text-alignmen (CLIP) and their aesthetic [86] scores.



Figure 4.10: Example images of MJHQ-30k dataset.

4.4 Model Compression Strategies

There are two primary strategies in model pruning. On the one hand, individual model components can be removed entirely, or, on the other hand, they can be compressed, e.g. via low-rank approximation. In transformer-based architectures like PixArt- Σ and Flux.1-dev, a natural choice is to partition the model according to individual transformer blocks. Therefore, any reference to "model parts" in the following sections specifically denotes these individual transformer blocks. Both strategies and their respective variants are described in the section below.

4.4.1 Block Removal

In structural model pruning, entire functional components of a model are eliminated [82, 96, 193] while others that are considered as highly relevant for the models performance are maintained. Concretely, if a transformer block B_i is removed, the output features of block B_{i-1} (denoted as x_{i-1}) are directly passed to block B_{i+1}

$$x_{i+1} = B_{i+1}(x_{i-1}) \quad . \quad (4.5)$$

Due to consistent dimensionality of the hidden features across the model's depth no additional dimension adaption needs to be done. A schematic sketch of this procedure is included in Fig. 4.11.

4.4.2 SVD Compression

Singular Value Decomposition (SVD) [197] factorizes a matrix $W \in \mathbb{R}^{n \times m}$ into three matrices U, Σ, V such that

$$W = U\Sigma V^T \quad (4.6)$$

where $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{m \times m}$ are orthonormal matrices and $\Sigma \in \mathbb{R}^{n \times m}$ is a diagonal matrix. The entries of Σ are referred to as singular values and are sorted in descending order $\sigma_1 \geq \sigma_2 \geq \sigma_3 \dots \geq \sigma_p$ with $p = \min(n, m)$. The magnitude of each singular value serves as a proxy for the importance of the corresponding rows and columns of U and V .

In model compression [180, 191], SVD is utilized to approximate weight matrices by retaining only the $k < p$ largest singular values, such that $U_{\text{red}} \in \mathbb{R}^{n \times k}$, $V_{\text{red}} \in \mathbb{R}^{m \times k}$ and $\Sigma_{\text{red}} \in \mathbb{R}^{k \times k}$, thereby reducing the total parameter count depending on the choice of k . Specifically, the matrix W is approximated by two matrices $A = U_{\text{red}}\sqrt{\Sigma_{\text{red}}}$ and $B = \sqrt{\Sigma_{\text{red}}}V_{\text{red}}^T$. Given a feature vector \mathbf{x} , the standard transformation $\mathbf{x}_{\text{next}} = W\mathbf{x} + \mathbf{b}$ is replaced by

$$\mathbf{x}_{\text{next}} = AB\mathbf{x} + \mathbf{b} \quad (4.7)$$

where \mathbf{b} is the bias term, which remains unchanged. Rather than defining a fixed rank k , it is more practical to specify the percentage c of parameters to be removed. For a weight matrix $W \in \mathbb{R}^{n \times m}$, the initial number of parameters is given by $P_{\text{orig}} = n \times m$. The target number of parameters is then given by $P_{\text{targ}} = P_{\text{orig}} \times (1.0 - c)$. Consequently, the target rank is determined as

$$k_{\text{targ}} = \max \left(1, \int \left(\frac{P_{\text{targ}}}{n + m} \right) \right) \quad . \quad (4.8)$$

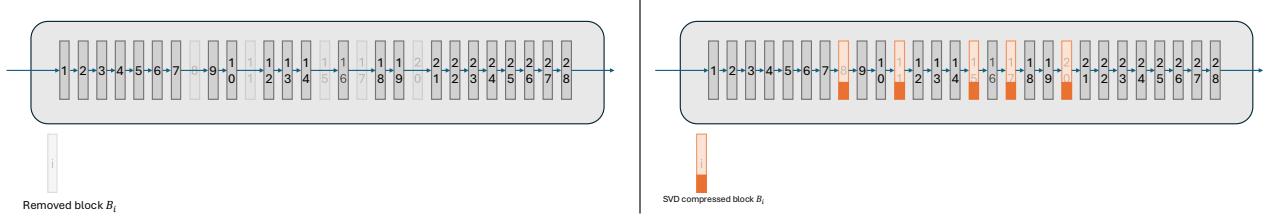


Figure 4.11: Example of removal (left) and compression (right) strategy based on a simplified representation of PixArt- Σ architecture. The depictions show the case that the blocks 8,11,15,17 and 20 are removed or compressed.

4.4.3 Iteratively Repeated SVD Compression

In the final compression setup (see Sec. 4.6.3), individual transformer blocks are compressed and subsequently retrained multiple times using SVD. Initially, an original weight matrix W is decomposed into matrices A and B such that $W \approx AB$. During the subsequent retraining, the matrices A and B are updated to A_t and B_t . Therefore, the relation $W \approx A_t B_t$ does not necessarily need to hold with respect to the initial matrix W , as the model adapts to the low rank constraint during training. Therefore, if the transformer block containing the weight matrix W is to be compressed again, the current effective weight matrix $W_t = A_t B_t$ is first reconstructed, followed by the application of the SVD algorithm $W_t = U^{(2)} \Sigma^{(2)} V^{(2)T}$. Similar to the initial step, after choosing a further reduced rank, the new smaller matrices $A^{(2)}$ and $B^{(2)}$ are constructed via $A^{(2)} = U_{\text{red}}^{(2)} \sqrt{\Sigma_{\text{red}}^{(2)}}$ and $B^{(2)} = \sqrt{\Sigma_{\text{red}}^{(2)}} V_{\text{red}}^{(2)T}$. By following this procedure, reconstructing $W_t^{(j)} = A_t^{(j)} B_t^{(j)}$, choosing a new, smaller rank and obtaining the new $A^{(j+1)} = U_{\text{red}}^{(j+1)} \sqrt{\Sigma_{\text{red}}^{(j+1)}}$ and $B^{(j+1)} = \sqrt{\Sigma_{\text{red}}^{(j+1)}} V_{\text{red}}^{(j+1)T}$, a transformer block can be gradually compressed across multiple iterations.

4.4.4 GRASP

The method named Gradient-Based Retention of Adaptive Singular Parameters GRASP was originally proposed for compressing LLMs [111]. It is based on the idea of exploiting the inherent low-rank structure of weight matrices within LLMs. Crucially, the authors hypothesized that the standard SVD approach where the singular values are selected solely based on their magnitude is suboptimal. They argue that even small singular values may be vital for the model’s performance and should be preserved. Consequently, they proposed a gradient-based criterion instead of a magnitude-based selection for the singular values.

In this work, the GRASP method is transferred to image generation models demonstrating its application on Flux.1-dev. First, the candidate blocks for pruning are identified using the method described in Sec. 4.2. Next, these blocks are compressed via SVD. Afterwards, all parameters are frozen and only the singular values are treated as trainable parameters for gradient computation. For computing the gradients, a calibration dataset is utilized consisting of 200 images. For each image, the gradients with respect to each singular value are computed and accumulated over all timesteps $t \in \{1, \dots, 1000\}$ whereby the normal flow matching loss is utilized. Finally, the gradients are averaged over the calibration dataset to obtain a final estimate of the importance of each singular value. Note that gradients are accumulated for scoring but the singular values themselves are not updated via optimization. Ultimately, instead of retaining the k largest singular values by magnitude, the k singular values with the largest accumulated gradients are selected for the low-rank approximation of the original weight matrix.

4.5 Knowledge Distillation Loss Functions

The general learning objective for knowledge distillation as described in Eq. 2.90 is

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda_{\text{OutKD}} \mathcal{L}_{\text{OutKD}} + \lambda_{\text{FeatKD}} \mathcal{L}_{\text{FeatKD}} . \quad (4.9)$$

In the experiment Sec. 5.1.5 the importance of the individual components of this loss function is investigated.

4.5.1 Normalization

It is important to note that the features of different blocks may have varying magnitudes. Consequently, [193] demonstrated that applying a normalization to the features is beneficial. Therefore, after computing the difference between the teacher and the student features a normalization is applied to ensure that the individual blocks contribute equally to $\mathcal{L}_{\text{FeatKD}}$. The normalization scheme adopted in this work is taken from [199].

Let T_l and S_l denote the features of the teacher and the student corresponding to transformer block $l \in \{1, \dots, L\}$. For a batch size of B , each feature tensor has the dimensions (B, D_l) where D_l is the total number of flattened features. The i -th sample in the batch for layer l is denoted as $T_{i,l}$ and $S_{i,l}$, respectively. First, the MSE is computed between the teacher and the student features

$$\mathcal{L}_{i,l} = \frac{1}{D_l} \|T_{i,l} - S_{i,l}\|_2^2 . \quad (4.10)$$

Next, the RMS of the teacher features is computed serving as a measure of feature magnitude

$$\eta_{i,l} = \sqrt{\frac{1}{D_l} \|T_{i,l}\|_2^2} . \quad (4.11)$$

Additionally, the loss is scaled by the magnitude of the task loss $\mathcal{L}_{i,\text{Task}}$. The normalization factor is then given by

$$w_{i,l} = \underbrace{\left(\frac{\mathcal{L}_{i,\text{Task}}}{\mathcal{L}_{i,l} + \epsilon} \right)}_{\text{Loss Balancing}} \cdot \underbrace{\left(\frac{\sum_{k=1}^L \eta_{i,k}}{L \cdot \eta_{i,l} + \epsilon} \right)}_{\text{Norm Normalization}} \quad (4.12)$$

where a small ϵ is added to the denominator to prevent division by zero. The final feature loss is computed by

$$\mathcal{L}_{\text{FeatKD}} = \frac{1}{B} \sum_{i=1}^B \left(\frac{1}{L} \sum_{l=1}^L w_{i,l} \cdot \mathcal{L}_{i,l} \right) . \quad (4.13)$$

Note that this normalization applies not only to $\mathcal{L}_{\text{FeatKD}}$ but also to $\mathcal{L}_{\text{OutKD}}$ which compares the final output of the teacher with the student's output.

4.5.2 Feature Mapping for Loss Computation

The assignment of learning signals based on intermediate features differs significantly between the complete block removal and the SVD block compression approach (see Fig. 4.12). In the SVD compression approach, the intermediate features of every student block S_l are directly compared with the features of the corresponding teacher block T_l , for all $l \in \{1, \dots, L\}$ assuming the model consists of L transformer blocks.

In contrast, the block removal approach necessitates a dynamic mapping strategy for selecting the correct intermediate features. If a block B_l is removed, the features of the preceding student block S_{l-1} are compared with the features of the teacher's block T_l . This effectively forces block B_{l-1} to approximate the removed block

B_l . Consequently, if several consecutive blocks B_{l-i} to B_l (where $i \geq 1$) are removed, the student features $S_{l-(i+1)}$ are compared to T_l .

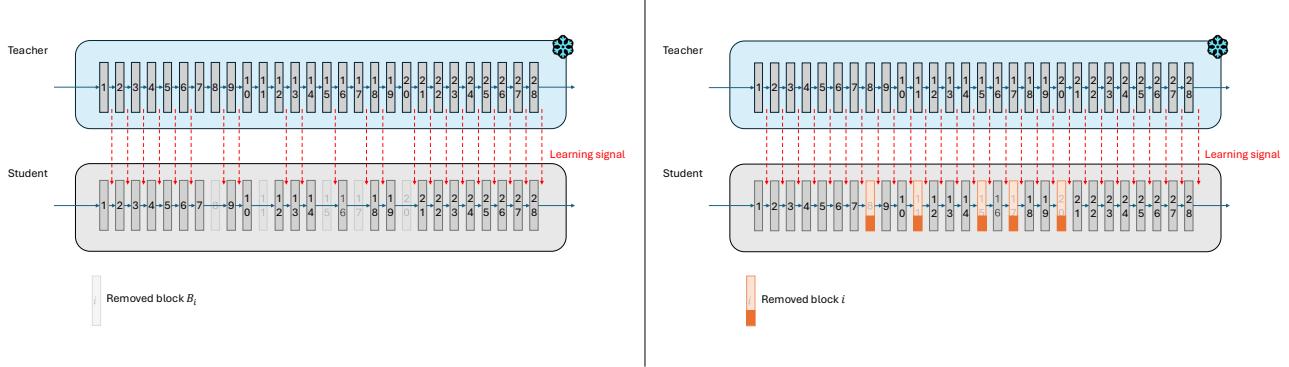


Figure 4.12: Visualization of the feature mapping between student and teacher model for loss computation.

4.6 Pruning Schedule

The pruning schedule for models dictates when and to what extent the compression of individual model components is performed. There are mainly two options for the temporal progression of parameter removal. Concretely, all parameters can be removed at once or iteratively. Regarding the structural distribution of removed parameters, one needs to decide how they are allocated across the model and whether model components can only be compressed once or multiple times. To explore these aspects, three distinct pruning schedules, namely one-shot, progressive, and linear progressive pruning, are investigated in this Master’s thesis and presented in the following subsections.

4.6.1 One-Shot Model Pruning

The standard approach for compressing transformer-based diffusion and flow models [33, 96, 173, 193] involves three main steps. First, the transformer blocks $\{B_1, \dots, B_L\}$ are sorted according to their importance, e.g. using CLIP-score and CMMMD as importance proxies (see Sec. 4.2.1). Let $\{B_1^*, \dots, B_L^*\}$ denote the sorted blocks, where B_1^* represents the least important one. Second, based on the target compression ratio of the total model, the number of blocks $k \leq L$ to be pruned is determined and the blocks $\{B_1^*, \dots, B_k^*\}$ are modified based on the chosen compression strategy (see Sec. 4.4). Third, the compressed model is retrained, e.g. using knowledge distillation (see Sec. 2.4.1), to recover the generative capabilities. This approach introduces a significant, abrupt structural alteration of the model, leading to a severe drop in performance. A scheme of this schedule is depicted in Fig. 4.13.

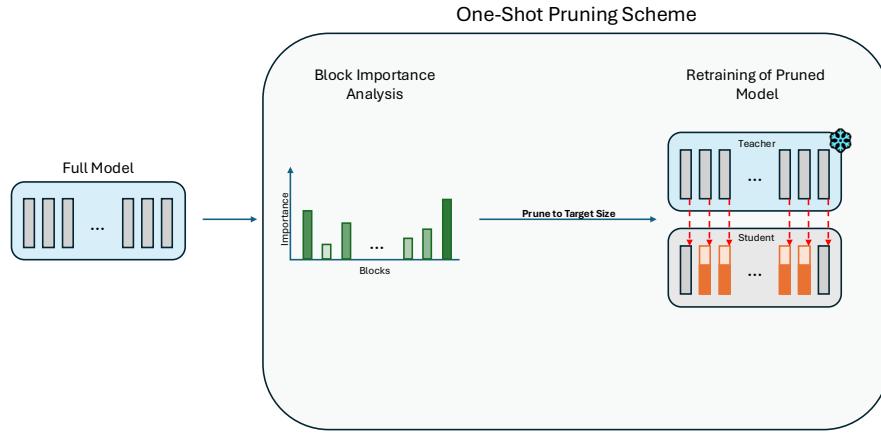


Figure 4.13: Visualization of the one-shot compression procedure which iteratively compresses the model.

4.6.2 Progressive Model Pruning

In contrast to one-shot model compression, progressive model compression, e.g. [150], iteratively applies the three-step procedure described above. In each iteration, the importance of previously unmodified blocks is evaluated, the number of blocks need to be modified is determined, and the compressed model is retrained. The crucial difference is that the number of parameters removed in a single iteration is significantly smaller than in one-shot model pruning. This leads to a gradual and less severe performance degradation in each iteration, which is easier to recover from. Note that each transformer block is only compressed once and in a subsequent iteration, the block importance is re-assessed only for all remaining unmodified blocks, as the relative block ranking may have changed as a result of the retraining. A scheme of this schedule is depicted in Fig. 4.14.

A design choice of the progressive model compression concerns the number of parameters that are removed per iteration. Either, the number of parameters can be set to a fixed value, or it could be chosen dynamically based on the importance metrics, e.g., by defining a threshold $\text{CMMD} = 0.1$ and initiating retraining as soon as it is exceeded by compressing an additional transformer block. The latter is utilized as standard setting for all experiments. This typically implies that in early iterations more parameters are removed than in later iterations.

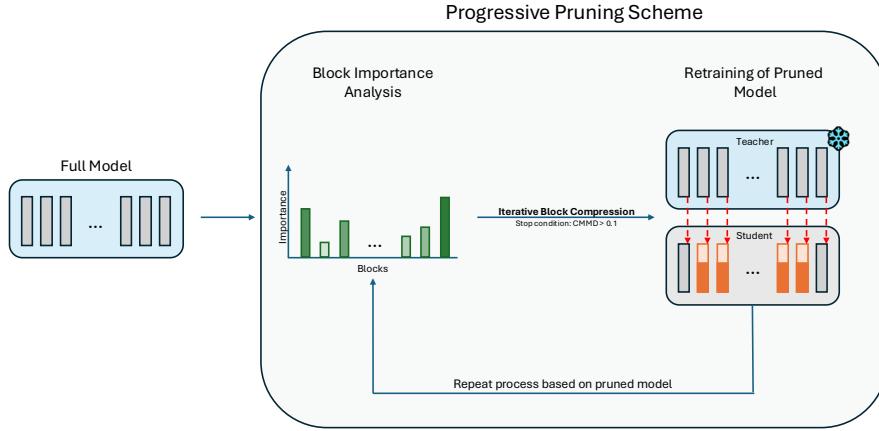


Figure 4.14: Visualization of the progressive compression procedure which iteratively compresses the model.

4.6.3 Linear Progressive Model Pruning

In this work, a new linear progressive compression schedule is proposed, acting as an advancement of progressive model compression, which introduces a fourth step, namely determining individual compression ratios of each block during each iteration. Instead of compressing all blocks uniformly and each block only once, as it was done in standard progressive pruning, the compression ratio of each block is determined only based on its importance, allowing blocks to be modified across multiple iterations. Specifically, the assigned compression ratio linearly decreases with increasing block importance.

The first step of linear progressive model pruning is similar to standard progressive pruning. The block importance is determined by leveraging the CMMD. However, as multiple compressions of the same transformer block are possible, the CMMD value is normalized by the number of removed parameters to ensure fair comparison, as reducing an already compressed block by $c_{\text{fixed}} = 60\%$ removes less parameters than compressing a full transformer block by the same percentage. Let P_i ($i \in \{1, \dots, L\}$) denote the current number of parameters of block B_i . Then, the final importance score for each block B_i is given by

$$s_i = \frac{\text{CMMD}_i}{P_i \cdot c_{\text{fixed}}} \quad (4.14)$$

where CMMD_i is the metric obtained from the dataset generated by the current model with the newly compressed block B_i . Based on these scores, the blocks are sorted into an order $\{B_1^*, \dots, B_L^*\}$ where B_1^* has the smallest importance score.

Second, in a crucial departure from standard progressive pruning, to determine the compression ratio for each block a linear compression schedule is applied where the least important block is compressed most, and blocks of higher importance less. Specifically, three parameters need to be specified by the user to compute the individual compression ratios c_i , namely the total number of parameters to be removed N , the number of blocks to be compressed during the current iteration k and the compression ratio of the least important block $c_1 = S$ to which is referred to as base compression ratio. To determine the compression ratios c_i of the remaining blocks $B_i^* \in \{B_1^*, \dots, B_k^*\}$, the base compression ratio S is decreased linearly by x according to the block's rank in the importance-sorted list. Therefore, the compression ratio for every block B_i^* is given by

$$c_i = S - ((i - 1) \cdot x) \quad (4.15)$$

with $i \in \{1, \dots, k\}$. The goal is to identify x based on the parameters S, N, k . The total number of parameters to be removed is computed by

$$P_{\text{total}} = \sum_{i=1}^k P_i^* \cdot c_i \quad (4.16)$$

where P_i^* is the current parameter count of block B_i^* . To solve for x , Eq. 4.15 and Eq. 4.16 are combined

$$\begin{aligned} & \sum_{i=1}^k P_i^*(S - (i-1) \cdot x) = N \\ \Leftrightarrow & \sum_{i=1}^k (P_i^* \cdot S) - \sum_{i=1}^k (P_i^* \cdot (i-1) \cdot x) = N \\ \Leftrightarrow & \sum_{i=1}^k (P_i^* \cdot S) - N = x \cdot \sum_{i=1}^k (P_i^* \cdot (i-1)) \end{aligned} \quad (4.17)$$

resulting in the final equation for the slope

$$x = \frac{\sum_{i=1}^k (P_i^* \cdot S) - N}{\sum_{i=1}^k (P_i^* \cdot (i-1))} \quad . \quad (4.18)$$

The input parameters S, N, k specified by the user have to meet two conditions. First, $N \leq S \cdot \sum_{i=1}^k P_i^*$ meaning the total number of parameters to be removed must not exceed the amount removed if all blocks were compressed by the base ratio S . Otherwise, the slope x would become negative, causing the compression ratio to increase for more important blocks which would contradict the ranking strategy. Second, the compression ratio of the last block c_k must not be negative which leads to the following condition

$$\begin{aligned} & 0 \leq c_k = S - (k-1) \cdot x \\ \Leftrightarrow & 0 \leq S - (k-1) \frac{S \sum_{i=1}^k P_i^* - N}{\sum_{i=1}^k (P_i^* \cdot (i-1))} \\ \Leftrightarrow & N \geq S \sum_{i=1}^k P_i^* - \frac{S \sum_{i=1}^k (P_i^* \cdot (i-1))}{k-1} \quad . \end{aligned} \quad (4.19)$$

Combining these constraints, a valid range for the total number of parameters to remove in one iteration based on the choice of S and k is obtained by

$$S \sum_{i=1}^k P_i^* - \frac{S \sum_{i=1}^k (P_i^* \cdot (i-1))}{k-1} \leq N \leq S \cdot \sum_{i=1}^k P_i^* \quad . \quad (4.20)$$

Third, after computing the compression ratios c_i , the model can be compressed based on the SVD compression scheme as described in Sec. 4.4.3 and is subsequently retrained. Finally, this procedure is repeated until the desired model size is reached.

A scheme of the whole procedure is depicted in Fig. 4.15. where step three and four, model compression and retraining, are summarized in the last element of the graphic.

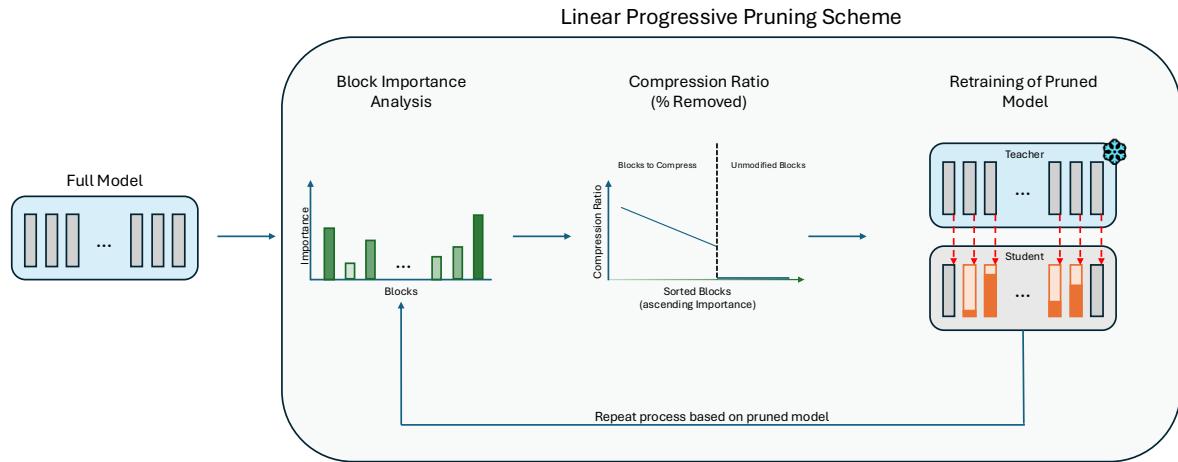


Figure 4.15: Visualization of the linear progressive compression procedure which iteratively compresses the model.

Chapter 5

Experiments

This master thesis investigates model pruning techniques, specifically focusing on the removal and SVD compression of transformer blocks within the PixArt- Σ and Flux.1-dev architectures. PixArt- Σ is employed as an experimental baseline to evaluate various design choices in the pruning framework due to its relatively low parameter count enabling rapid iteration and extensive ablation studies. Subsequently, the most effective strategies are applied to Flux.1-dev to assess their efficacy in high-parameter regimes.

5.1 PixArt- Σ Pruning

5.1.1 Experimental Default Setup

In the following section, individual components of the compression pipeline such as loss function, block selection algorithm, compression strategy, ..., are systematically investigated. To isolate the impact of each component, a default experimental setup is defined below. Unless stated otherwise, all subsequent experiments utilize the baseline configuration. It is important to note that this baseline configuration serves as consistent starting point for the experiments and does not necessarily represent the final, best setting which is derived later in this work. The following elements are constituents of the baseline configurations:

- **Learning Objective:** Knowledge distillation is utilized for fine-tuning the pruned models. Therefore, the learning objective is the sum of the unnormalized feature and output distillation loss $\mathcal{L}_{\text{OutKD}} + \mathcal{L}_{\text{FeatKD}}$.
- **Pruning Schedule:** The progressive compression strategy was applied. As soon as the CMMD value became larger than 0.1 in the block selection algorithm, the model was retrained to recover its generation capabilities.
- **Structural Compression Strategy:** Complete transformer block removal was used to reduce the model size.
- **Block Selection Algorithm:** The greedy algorithm based on CMMD and CLIP-score was utilized, computed on 100 images.
- **Training Datasets:** Training was performed sequentially starting on the LAION dataset followed by further training on the LAION-PixArt dataset. For each iteration, the models were trained for two epochs on each datasets.
- **Fine-tuning Protocol:** The recovery of the pruned model followed a structured three stage protocol. First, for local stabilization, optimization was restricted to the block B_{l-1} immediately preceding the removed block B_l for one epoch (38k steps) on the LAION dataset. In case where several blocks were

removed simultaneously, all corresponding preceding blocks were updated while the rest of the network was frozen. Second, the complete model was optimized for an additional epoch on the LAION dataset. Third, for domain-specific refinement, the model was finetuned for two epochs (12.5k steps) on the LAION-PixArt dataset.

- **Evaluation** The pruned models are evaluated on the HPSv2 and the GenEval benchmark. The HPSv2 benchmark investigates the overall image quality by predicting a human preference score and the GenEval benchmark analyze text-image alignment.
- **Training Hyperparameters:** All other training hyperparameters are summarized in tab. 5.1.

The implementation of all experiments build up on the official git repository for PixArt- Σ [19]. For evaluation of the model, the HPSv2 and the GenEval benchmarks were utilized in this section.

Table 5.1: Default hyperparameter settings for PixArt- Σ fine-tuning after compression.

Hyperparameter	Value
<i>Optimization</i>	
Optimizer Type	CAME
Learning Rate	2×10^{-5} (Block Removal)/ 2×10^{-6} (SVD Compression)
Weight Decay	0.03
Betas	(0.9, 0.999, 0.9999)
Epsilon	(10^{-30} , 10^{-16})
<i>Training Schedule</i>	
Warmup Steps	1000
Batch Size	16
Gradient Clipping	0.001
<i>Model & Data</i>	
Precision	bf16
Resolution	512×512
Class Dropout Prob.	0.1
Max Prompt Length	300

5.1.2 Block Importance Analysis

In structural pruning, the selection of components to be removed or compressed is essential to ensure that the compressed model retains the highest possible image generation quality. Qualitatively, Fig. 5.1 shows the effect of removing every transformer block from PixArt- Σ individually on its image generation capabilities without subsequent fine-tuning. It is clearly visible that removing the first and last transformer blocks (blocks 1-4 and blocks 27-28) exerts the most significant influence on the image quality. This confirms the observation from previous work [121] that redundancies are primarily find in the middle part of transformer based models. To identify which blocks are best to prune meaning which blocks have a minimal contribution on the total image generation process, several different metrics can be employed. In the following section, different approaches are investigated to determine the importance of the individual transformer blocks. First, the traditional magnitude-based pruning and representational dissimilarity by computing the central kernal alignment score (CKA) (see Sec. 4.2.1) between the input and output of each transformer block are evaluated. Moreover, they are verified against to standard performance metrics such as the CLIP-score and CMMMD (see Sec. 4.2.1).

Normally, several transformer blocks are removed or compressed simultaneously. Due to block interdependencies, two different options to identify the best grouping for blocks are investigated. First, a greedy-algorithm is applied, iteratively searching for the best next block to prune based on the already compressed model from

previous iterations. Second, the block selection process is formulated as hyperparameter search leveraging the Optuna framework.

Qualitative Block Analysis

Fig. 5.1 displays an example for the qualitative influence of the removal of individual transformer blocks from the PixArt- Σ architecture. For each image a different block was completely removed. The removal of the first four blocks and the last two blocks exhibits a strong degeneration in image quality which confirms the finding that the first and last layers are more critical for transformer-based models from [121]. In diffusion or flow models the first layers aggregate global context while the middle layers refine details. The last layers are important for high frequency details as seen in Fig. 5.1 where for removed block 27 or 28 the concept of a man is present in the image but very blurry.

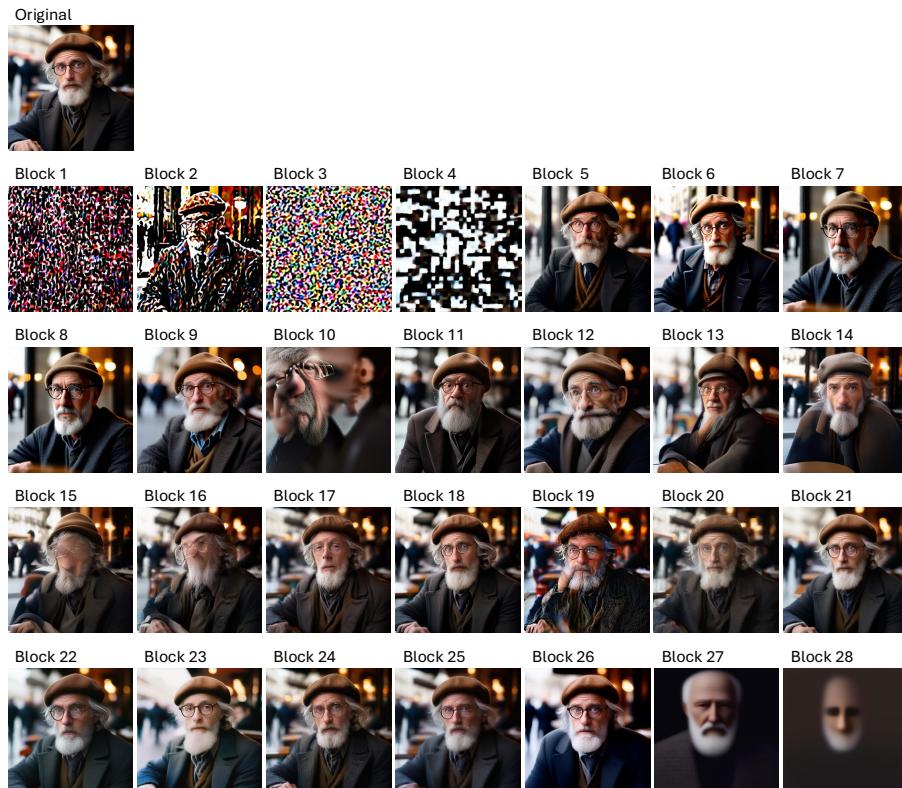


Figure 5.1: The figure shows the influence of selectively removing individual transformer blocks (1–28) on the final generation result.

Compression Criteria

In magnitude-based pruning strategies the model parts with the lowest averaged weight magnitude are considered to be less important for the overall model performance. Interestingly, in PixArt- Σ the weight magnitudes of the initial layers, which have a substantial influence on the final image synthesis, have the lowest weight magnitudes across the architecture (see Fig. 5.2 and 5.3 which both show the weight magnitudes for every block compared to the CLIP-score or CMMD). This observation contradicts with the assumption of low impact by low weights magnitudes. It seems that despite the small weight magnitude the first layers act as important initial feature translator having a massive influence on model performance. Furthermore, the weight magnitudes of the final layers are only marginally higher than those of the middle layers, while the middle layers' weight magnitudes exhibit a high degree of similarity. This makes a precise selection of blocks for removal based solely on this

heuristic difficult. Consequently, using the weight magnitude to identify redundant blocks within PixArt- Σ is not a suitable criterion for structural pruning in this context.

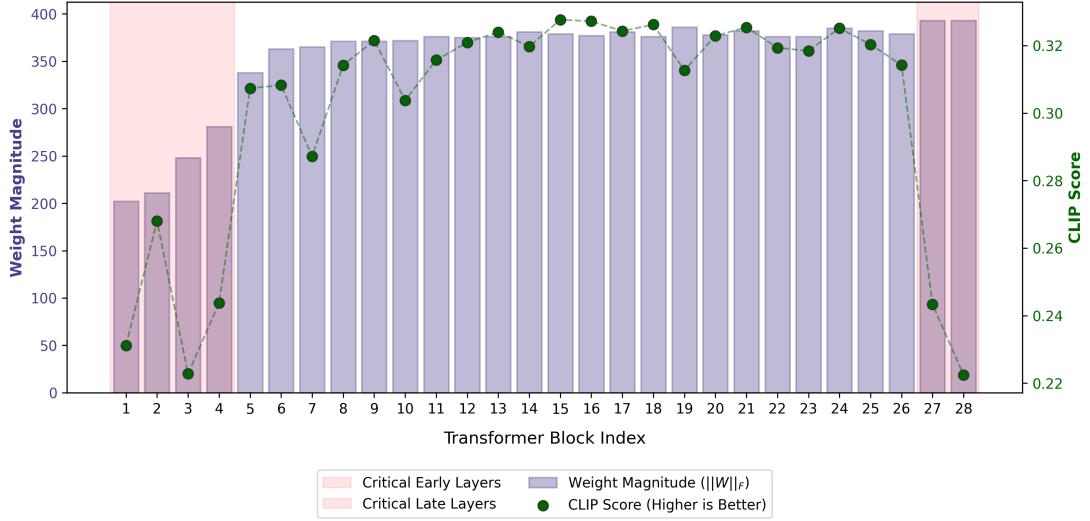


Figure 5.2: The histogram shows the magnitudes of the weights of the individual transformer blocks (1-28). The blocks which have shown qualitatively (see fig. 5.1) to be very important for the image generation capabilities of PixArt- Σ (see fig. 5.1) are marked with red. In addition, the CLIP-scores for each block is displayed. The CLIP-scores follow the qualitative observations in contrast to the magnitude-based analysis.

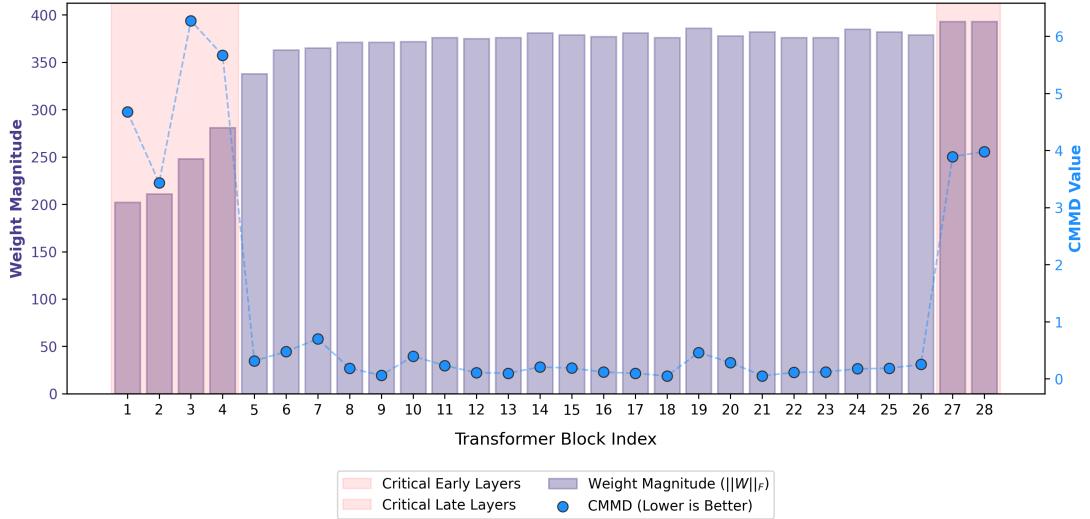


Figure 5.3: The histogram shows the magnitudes of the weights of the individual transformer blocks (1-28). The blocks which have shown qualitatively to be very important for the image generation capabilities of PixArt- Σ (see fig. 5.1) are marked with red. In addition, the CMMMD for each block is displayed. The CMMMD follow the qualitative observations in contrast to the magnitude-based analysis.

The CKA score is computed for every block and every sampling step across 100 example prompts taken from the LAION dataset. Subsequently, these scores are averaged across all samples and time steps resulting in a mean similarity score $\mu_{CKA,i}$ for each individual PixArt- Σ block. Averaging over the time steps should provide an estimate of the general importance of the block. However, the importance of a block at different time steps can vary drastically (see Appendix A). A low value indicates that the input and output of a specific block differs significantly, suggesting higher importance of the block.

Both Fig. 5.4 and 5.5 show the transformation intensity $1 - \mu_{CKA,i}$ for every block compared with CMMMD or

CLIP-score, a greater value signifies a more substantial feature transformation. While blocks 17 and 22 stand out with their high values, the implied importance of these blocks cannot be confirmed by qualitative analysis (see Fig. 5.1). Furthermore, blocks 1, 2, 4, 27, and 28 are not uniquely identified as critical blocks for the performance of the model, as the values for blocks 19, 20, and 21 are higher, which would mark them as more relevant. As this ranking also fails to align with the qualitative review of block influence on the final image, the use of representational similarity for block selection in PixArt- Σ must be viewed critically too.

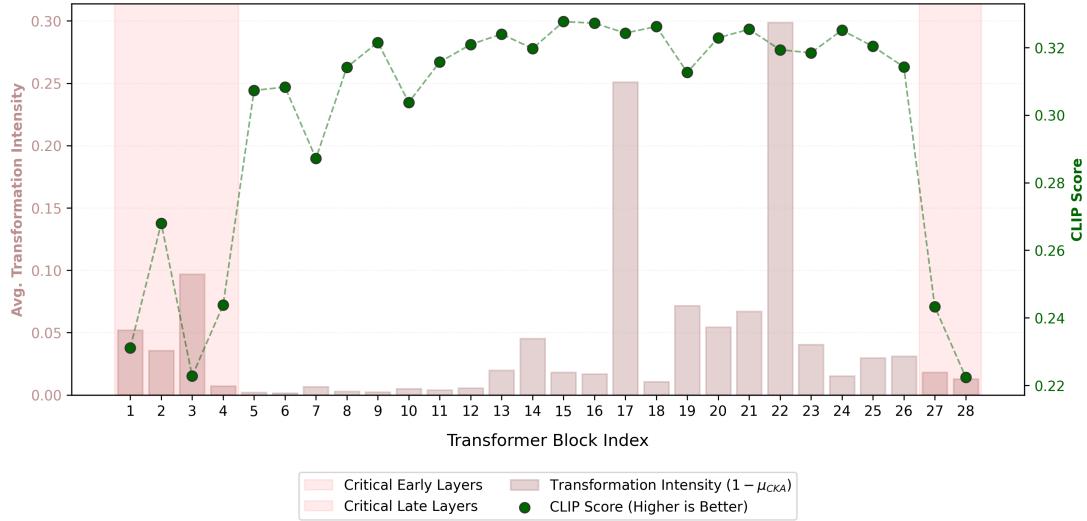


Figure 5.4: The histogram shows the transformation intensity based on CKA comparing the input and output features of the individual transformer blocks (1-28). The blocks which have shown qualitatively (see fig. 5.1) to be very important for the image generation capabilities of PixArt- Σ are marked with red. In addition, the CLIP-scores for each block is displayed. The CLIP-scores follow the qualitative observations in contrast to the transformation intensity analysis.

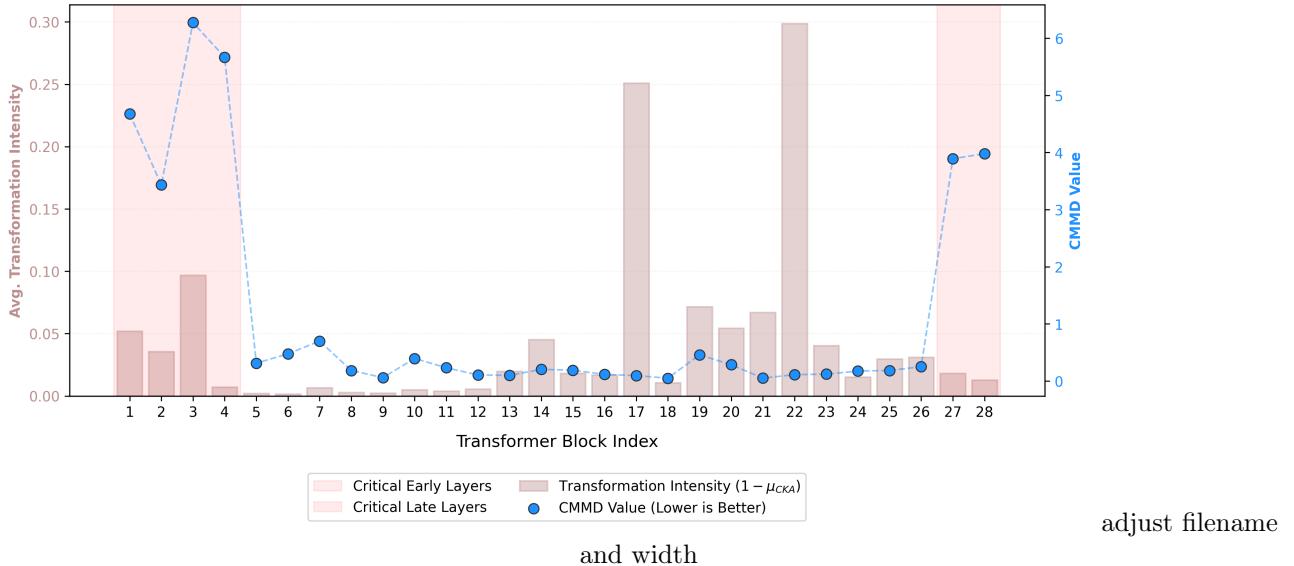


Figure 5.5: The histogram shows the transformation intensity based on CKA comparing the input and output features of the individual transformer blocks (1-28). The blocks which have shown qualitatively (see fig. 5.1) to be very important for the image generation capabilities of PixArt- Σ are marked with red. In addition, the CMMD for each block is displayed. The CMMD follow the qualitative observations in contrast to the transformation intensity analysis.

Due to the contradictory results of magnitude-based pruning and CKA-scores when compared with the

qualitative observations, a third methodology was investigated, namely quantifying the importance of individual blocks based on standard performance metrics [17]. Specifically, the CLIP-score and CMMD are computed on a small reference dataset consisting of 100 images from the LAION dataset. It should be noted that the standard protocols for computing statistically robust CMMD and CLIP-scores typically require larger datasets, e.g. upwards of 1,000. However, given the iterative nature of the greedy search (see Sec. 4.2.2) and the high computational overhead associated with repeated evaluations, using such a large reference set was prohibitively expensive. Therefore, a sample size of 100 images serves as a sufficient proxy to capture the relative importance ranking of the transformer blocks, providing the necessary directional guidance for the pruning process. As illustrated in the Fig. 5.4 and 5.5 (and likewise in Fig. 5.2 and 5.3) both metrics closely align with the qualitative evaluation, confirming that the first four and the last two blocks exert the greatest impact on the final image quality. Consequently, these metrics are selected as primary criteria for identifying redundant blocks in PixArt- Σ . This ranking serves as foundation for the subsequent, more sophisticated block selection algorithm.

Qualitative Block Analysis for SVD Compression

In the experiments in which the blocks were compressed via SVD, the scores for block selection were computed based on compressed blocks. In Fig. 5.6 the results of compressing every transformer block individually by 60% are displayed. One important point to note is that compressing the first blocks (1-4) which have a large impact when the blocks are completely removed (see Fig. 5.1) resulting only in minimal image changes compared to the original image generated with the uncompressed model. Also the compression of block 28 results in less severe image quality degeneration than completely removing it. This demonstrates that these blocks are critical for the overall performance but also contain high degree of redundancies. Notably, the compression of block 27 also exhibits a strong smoothing of details like it was already observed by removing it indicating its importance. Compressing the middle blocks does not lead to a big drop of image quality but they differ even more from the original image compared to the first blocks. A possible explanation is that for setting the global structure which is done by the first blocks a low rank structure is sufficient. The middle blocks which are responsible for the details, e.g. direction in which the eagle looks, are more sensitive for compression.

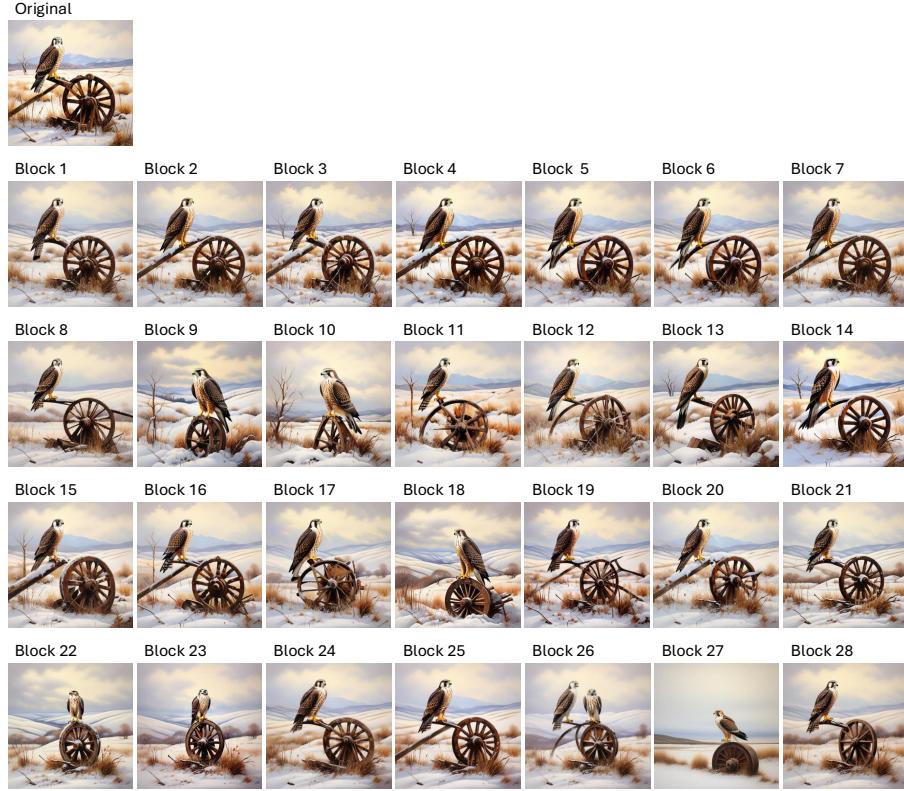


Figure 5.6: The figure shows the influence of selectively compressing individual transformer blocks (1–28) by 60% on the final generation result.

5.1.3 Block Selection Algorithms

A further key challenge after establishing the pruning criteria as the combination of CLIP-score and CMMD in model pruning is identifying the optimal combinations of blocks to remove simultaneously. When removing only a single block, there are merely 28 possibilities for PixArt- Σ . However, if the goal is to reduce the model size by 50% necessitating the removal of 14 blocks, the number of possible block combinations escalates to $\binom{28}{14}$. Therefore, it is computationally infeasible to evaluate the CLIP-score and CMMD for every possible combination.

Greedy

A natural alternative is to perform a greedy search algorithm. In this approach, the importance ranking for all blocks is initially computed, after which the least important block is identified and compressed. Next, the CLIP-scores and CMMD are re-evaluated for the remaining blocks within the newly reduced architecture. By iteratively evaluating and compressing the least important block the model architecture is progressively reduced while accounting for the interdependencies between blocks.

Optuna

Another approach is to reformulate the block selection problem as hyperparameter search using the hyperparameter framework Optuna. The number of blocks to compress N_c , the blocks that are compressed $\{B_i\}_{i=1,\dots,N_c}$ and the compression ratios of each block $\{\pi_i\}_{i=1,\dots,N_c}$ can be potential hyperparameters to optimize. Here, the focus is on identifying the blocks $\{B_i\}_{i=1,\dots,N_c}$ given N_c and $\{\pi_i\}_{i=1,\dots,N_c}$ to reduce the search space. As the objective function, the CMMD metric is computed on 100 reference images from the LAION dataset. To

further speed up the search process, the results of the greedy-algorithm are provided as a prior to the Optuna optimization framework, ensuring that the first trials are not picked randomly which would increase the search time drastically until a good setting is found.

The Optuna hyperparameter search was tested for the iterative SVD pruning framework with fixed compression ratios $\pi_i = 0.6 \forall i \in \{1, \dots, N_c\}$, where the number of compressed transformer blocks was chosen for each iteration individually. The procedure to determine the least important transformer blocks as follows. First, the greedy algorithm was leveraged to find a prior of potentially good candidates for compression. Afterwards, the Optuna hyperparameter framework was applied to refine the block selection. In total 300 different configurations were tried out by Optuna where the decision criterion was solely the CMMD score. After the best block combination was identified, the pruned model was finetuned. This procedure was executed repeatedly for four iterations. Fig. 5.7 shows the prior blocks from the greedy algorithm as well as the final block selection of the Optuna algorithm. On the right side the corresponding CMMD values which were computed on a small reference dataset of 100 images are displayed for each configuration. Note this is the CMMD value before fine-tuning to compare the prior with the final selection. The grey blocks, represent the consensus where the Optuna algorithm did not change the blocks from the prior whereas the yellow and green blocks mark the exclusive blocks selected by the greedy algorithm and the Optuna algorithm respectively. For iterations one to three, one observes that the majority of blocks from the greedy prior were maintained by the Optuna algorithm, only in iteration four two out of three blocks differ. The large overlap is a first indication that the greedy algorithm already performs well. However, because the greedy result was used as prior the optuna framework search is focused around it which might also explain the high similarity. Overall, one can conclude that the CMMD values obtained by the Optuna algorithm are consistently smaller (better) than the ones obtained by the greedy algorithm. This indicates that at least the pre-training performance is enhanced by the refinement step.

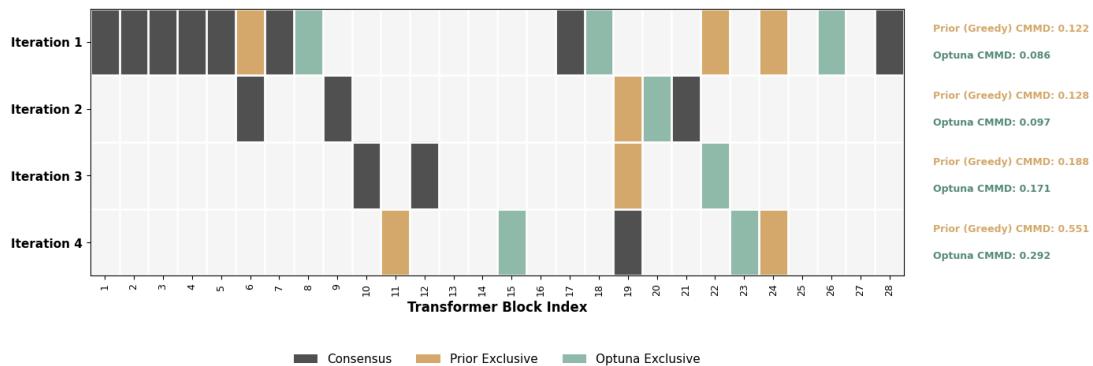


Figure 5.7: Comparison of removed blocks per iteration by the Greedy algorithm versus the Optuna framework. In each iteration, the Greedy algorithm's selection served as the initialization prior for the Optuna hyperparameter search. The heatmap visualizes the divergence in block selection: While gray blocks represent consensus, the colored blocks highlight where Optuna (green) found a superior structural reduction compared to the Greedy baseline (yellow). The model was fine-tuned between iterations based on the optimized configuration. The annotations demonstrate that Optuna consistently achieves lower (better) CMMD scores.

Above, it is observed that in each iteration, the Optuna framework was able to identify a block combination for compression which resulted in a smaller CMMD value compared to the prior greedy algorithm.

To finally evaluate both, Greedy and Optuna, approaches, the final post-fine-tuning performance was compared. In this step, PixArt- Σ is compressed twice, once based on the greedy algorithm and once based on the Optuna algorithm as described above, each time targeting the same number of parameters (same compression rates $\pi_i = 0.6$ and same number of blocks in each iteration). Fig. 5.9 shows the corresponding distribution of blocks being compressed in each iteration. Please note that here "Greedy" specifies the compression scheme where the greedy algorithm is strictly utilized as the block selection algorithm, in contrast to the previous

section, where "Greedy" just represented the prior for the Optuna algorithm. Looking at the block distribution, one observes that the early layers are compressed heavily in both schemes first, indicating a consensus on the high redundancy of these initial blocks. In the later layers, no clear structure is visible and the selections diverge. This lack of a consistent pattern suggests that the deeper blocks are largely interchangeable regarding their impact on the model's capacity. The evaluation of the fine-tuned compressed models (see Fig. 5.8) does not provide clear evidence that the Optuna block selection algorithm led to better performing models post-fine-tuning. One observes that for the smallest model, the HPSv2 score is even better for the greedy algorithm, and a similar trend is observed for the second and third iteration regarding the GenEval score. This demonstrates that the fine-tuning process acts as an equalizer, effectively eliminating the small primary advantage of the Optuna algorithm.

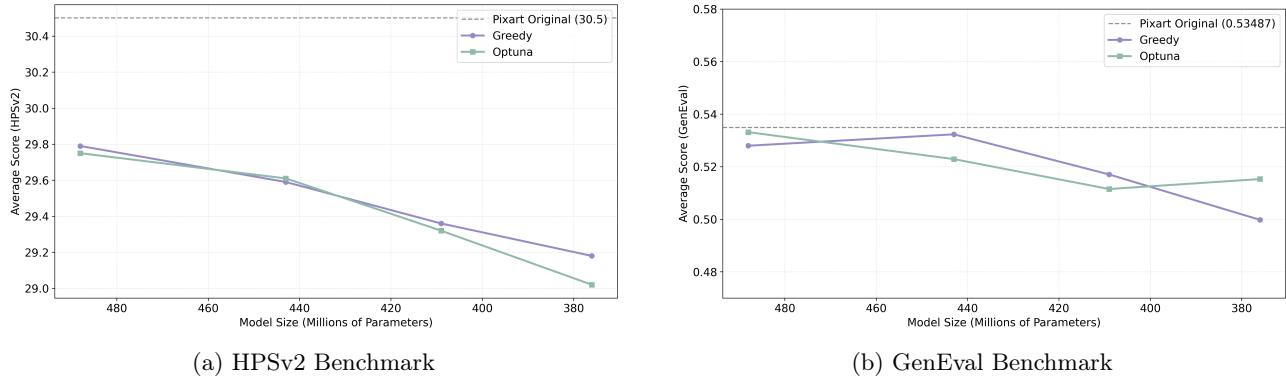


Figure 5.8: Compare the performance of pruned models based on the Optuna and Greedy block selection algorithm.

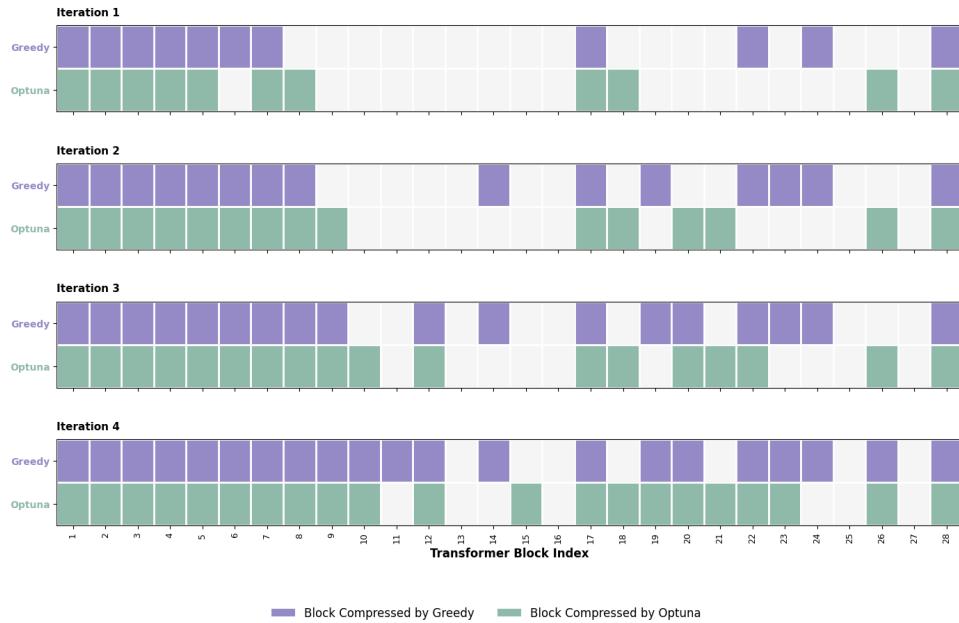


Figure 5.9: Here, all compressed blocks based on a greedy and optuna pruning scheme are displayed per iteration.

5.1.4 One-Shot vs Progressive Model Pruning

One-shot model pruning compresses the model to its target size in a single step, while progressive pruning employs an iterative scheme of incrementally removing parameters, followed by recovery fine-tuning phases.

Both frameworks are tested using the complete block removal approach. The number of blocks removed per model and the corresponding number of fine-tuning steps are summarized in Tab. 5.2.

Fig. 5.10 displays the performance of the one-shot versus progressive pruning on the HPSv2 and GenEval benchmarks across various compression ratios. To ensure a fair comparison, the models for both methods are fine-tuned for the same number of steps and on exactly the same data. Both benchmarks demonstrate that one-shot compression is a viable strategy in the low-compression regime (10% to 20% parameter reduction). However, as the compression ratio increases, the progressive approach becomes advantageous because the incremental removal of parameters allows the model to recover its generative capabilities more effectively compared to the abrupt removal of a large amount of parameters at once. Fig. 5.11 presents visual evidence that the model progressively pruned by 50% generates images of higher quality than those obtained using the one-shot approach. Nevertheless, both approaches exhibit noticeable visual degradation at a compression ratio of 50 %.

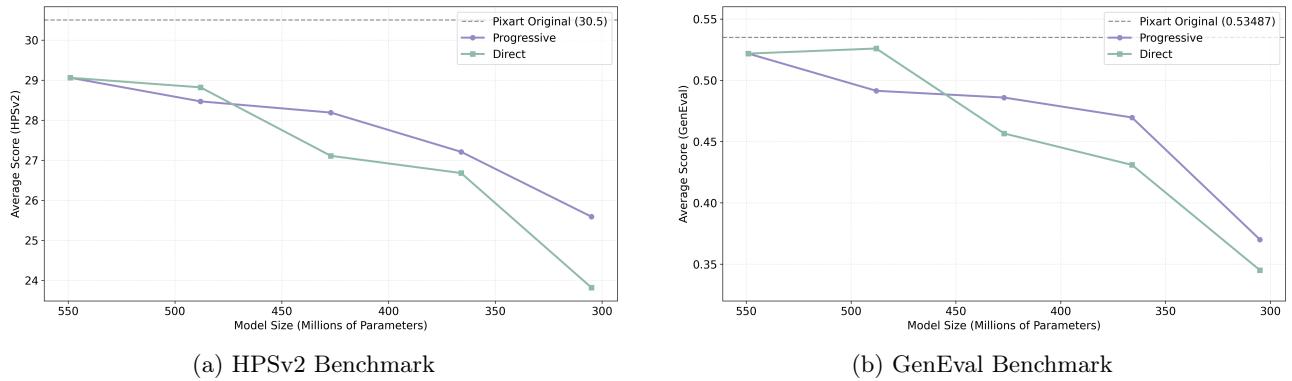


Figure 5.10: Performance comparison between progressive and one-shot pruning across varying degrees of model compression.



Figure 5.11: Example images for one-shot versus progressive pruned models by 50%.

Table 5.2: Training and compression settings for compressed PixArt- Σ models.

Model Variant (Remaining %)	Number of Removed Blocks	Training Steps (in Thousands)
Model 90%	3	88.5
Model 80%	6	265.5
Model 70%	8	531.0
Model 60%	11	708.0
Model 50%	14	973.5

5.1.5 Knowledge Distillation Loss

In this section, the effectiveness of the individual components of the knowledge distillation loss (see Eq. 2.90) and their combinations is investigated. In total, five different loss functions are tested:

- **Output-Only Distillation ($\mathcal{L}_{\text{OutKD}}$):**

The MSE between the final output of the student and the teacher serves as the sole learning signal.

- **Standard Distillation ($\mathcal{L}_{\text{OutKD}} + \mathcal{L}_{\text{FeatKD}}$):**

The unweighted sum of the final output loss and the intermediate feature loss is leveraged as the learning objective.

- **Hybrid Distillation ($\mathcal{L}_{\text{OutKD}} + \mathcal{L}_{\text{FeatKD}} + \mathcal{L}_{\text{Task}}$):**

The original diffusion loss is added to the standard distillation loss.

- **Normalized Distillation ($\mathcal{L}_{\text{OutKD,norm}} + \mathcal{L}_{\text{FeatKD,norm}}$):**

The final output and the intermediate feature loss are normalized (see Sec. 4.5.1) before being summed together to ensure that individual layers do not dominate the total loss.

- **Normalized Hybrid Distillation ($\mathcal{L}_{\text{OutKD,norm}} + \mathcal{L}_{\text{FeatKD,norm}} + \mathcal{L}_{\text{Task}}$):**

This configuration extends the normalized distillation objective by incorporating the original diffusion task loss.

PixArt- Σ is pruned to varying sizes between 90% and 50% utilizing the described loss configurations. The setup follows the exact same settings as described in Tab. 5.2 are utilized.

Fig. 5.12 presents the performance of the pruned models obtained by the different loss configurations. According to the HPSv2 benchmark, which measures overall image quality, the normalization of the intermediate features and the final output loss proves advantageous (see blue and red lines in Fig. 5.12). Both loss configurations which utilize normalization achieve a consistently higher HPSv2 score compared to their unnormalized counterparts. The normalization ensures that the contribution of each intermediate feature to the overall loss is balanced. This prevents the optimization process from prioritizing the later layers, where the errors are numerically larger, to the detriment of the earlier layers. The numerically larger discrepancies between the teacher and the student stem from the larger feature magnitudes in deeper layers, as displayed in Fig. 5.13 which presents the RMS of the intermediate features after every transformer block in PixArt- Σ averaged over 100 prompts. The increase in feature magnitude results from the signal accumulation inherent to the PixArt- Σ architecture, which follows the residual formulation $\mathbf{x}_{l+1} = \mathbf{x}_l + B_l(\mathbf{x}_l)$ where B_l represents the l^{th} transformer block.

Between the two loss configurations applying normalization, the normalized hybrid distillation performs best. This suggests that it is advantageous for the fine-tuning process of the pruned models that besides the teacher’s learning signal the original diffusion loss is provided as well. The original loss might become critical when the teacher is uncertain itself about a prediction as it anchors the student to the ground truth.

The benefit of the original diffusion loss is also confirmed by the unnormalized settings where the hybrid distillation consistently achieves a higher HPSv2 score compared to the standard distillation approach.

Interestingly, the output-only approach outperforms both standard and hybrid distillation settings for shallow compressed models. This suggests that the unbalanced features prevent effective learning. However, for strongly compressed models the output-only loss configuration leads to an abrupt drop in quality. This shows that the intermediate feature loss, despite being suboptimally distributed, serves as a critical stabilizer for the fine-tuning process.

The GenEval benchmark mainly confirms the observations from the HPSv2 benchmark, that the normalized hybrid distillation approach performs best overall and that adding the original diffusion loss boosts performance. Notably, the output-only configuration results are even more pronounced, performing comparably well or better

than to the normalized hybrid distillation approach in the low-compression regime (90% to 70%). This suggests that at low compression, removing the constraints of feature matching allows the student to optimize semantic alignment very effectively. However performance drops severely for larger compression rates.

Furthermore, it is worth noting that for the high-compression regime of 60% or 50% the unnormalized hybrid distillation approach slightly outperforms the normalized distillation approach. This indicates that in the high-compression regime, the original diffusion loss is critical for maintaining text-to-image alignment.

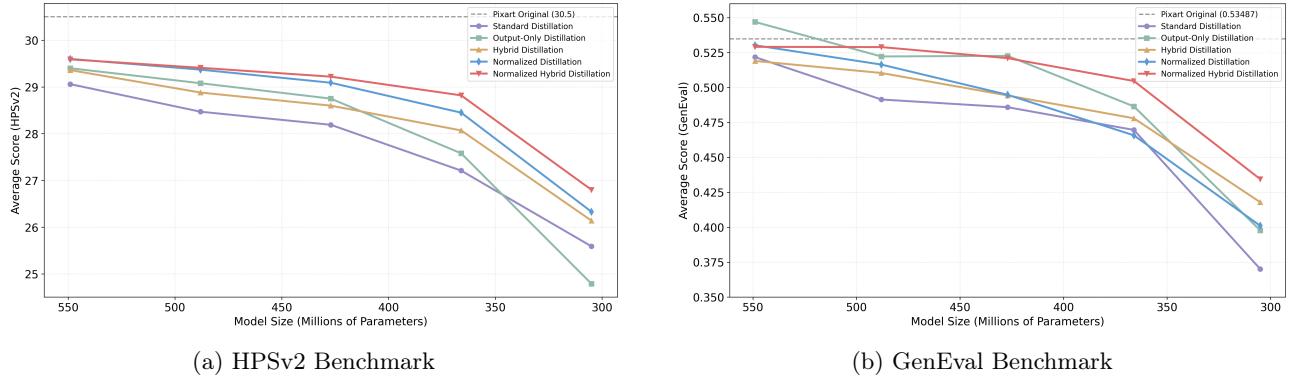


Figure 5.12: Impact of different loss configurations on performance across varying compression ratios.

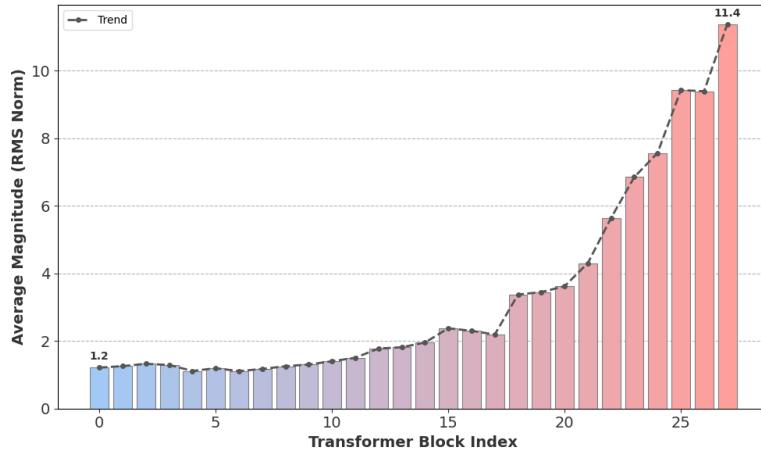


Figure 5.13: The RMS norm of the output features of the individual transformer blocks from PixArt- Σ varies significantly. Features from later layers have a larger magnitude than earlier layers. The RMS norm is computed on 100 prompts from LAION dataset and averaged.

5.1.6 fine-tuning Protocol

After removing blocks from the PixArt- Σ architecture, the reduced model is finetuned to recover its image generation capabilities. In total, three different fine-tuning protocols are tested to identify the most effective one. In all protocols, the model was trained for 88.5k optimization steps to enable fair comparison.

- **Three Stage fine-tuning Protocol**

This is the default fine-tuning protocol described in sec. 5.1.1. It is repeated for the reader's convenience.

1. Local Stabilization: The blocks immediately preceding the removed blocks are finetuned for 38k steps on the LAION dataset while all other parts of the model are frozen.
2. Global Optimization: The complete model is optimized for one epoch (38k steps) on the LAION dataset.

3. Domain Specific Refinement: The LAION-Pixart dataset is leveraged to finetune the pruned model for 12.5k steps on its teachers domain.

- **Two Stage Finetuning Protocol**

In the following protocol, instead of performing local fine-tuning, the entire model is directly finetuned for 2 epochs.

1. Global Optimization: The complete model is optimized for two epoch (76k steps) on the LAION dataset.
2. Domain Specific Refinement: The LAION-Pixart dataset is leveraged to finetune the pruned model for 12.5k steps on its teachers domain.

- **Two Stage LAION only fine-tuning Protocol**

The absence of the domain specific refinement is investigated in this protocol.

1. Local Stabilization: The blocks immediately preceding the removed blocks are finetuned for 38k steps on the LAION dataset while all other parts of the model are frozen.
2. Global Optimization: The complete model is optimized for two epoch (50.5k steps) on the LAION dataset.

The evaluation of the performance obtained for models of different pruning ratios (90% to 50%) can be found in fig. 5.14. For both the HPSv2 and GenEval benchmarks, fine-tuning solely on the LAION dataset (yellow line) leads to a sharp decline for strongly compressed models. This is a strong indication that the image quality being higher in the LAION-PixArt dataset compared to the LAION dataset is highly relevant for the recovery of the generation capability of a pruned transformer-based diffusion model. Moreover, it could be observed that for shallow compressed models (10% to 20%) the two stage fine-tuning protocol (green line) performs better than the three stage protocol (purple line). This indicates that in the low compression regime the model benefits from the extended global fine-tuning and no local repair is needed. For the middle compressed models (70%-60%), the three stage fine-tuning is consistently better or on par with the two stage protocol. This performance gap suggests that the local stabilization step provides necessary structural regularization when the model capacity is reduced to a critical threshold. For strong compression there are contradictory results. According to the HPSv2 benchmark the three stage approach is still marginally better than the two stage approach. However, the GenEval benchmark suggests that it is reverse. Given the high degree of degradation in both models at this stage, these marginal differences may also reflect the inherent instability of training extremely pruned networks rather than a definitive superiority of one method over the other.

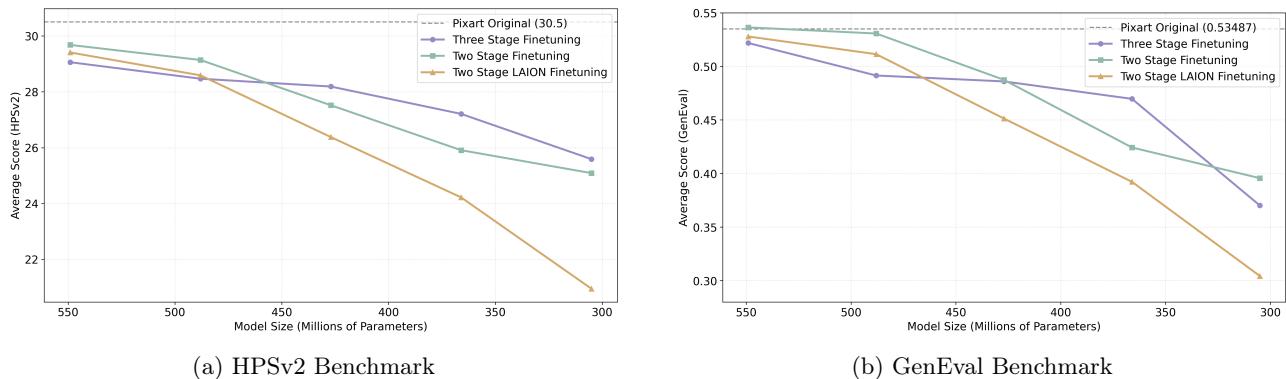


Figure 5.14: Impact of different fine-tuning protocols on performance across varying compression ratios.

5.1.7 Compression Strategy

As an alternative to the pruning strategy of removing entire transformer blocks from PixArt- Σ , a compression strategy utilizing SVD compression is investigated. This decomposition is applied to the MLP, self-attention and cross-attention matrices within a transformer block. To evaluate the efficiency of the SVD compression method against the structural block removal method, a representative operating point is chosen. The rank for the MLP is set to $r_{\text{MLP}} = 512$ corresponding to a compression of approximately 44% of its parameters and the ranks for the self- and cross-attention matrices are fixed to $r_{\text{self-attn}} = r_{\text{cross-attn}} = 128$ representing a parameter reduction of approximately 78%. This serves as a starting point. Later, the choice of compressing attention matrices more aggressively than the MLP matrices is further investigated (see sec. 5.1.8).

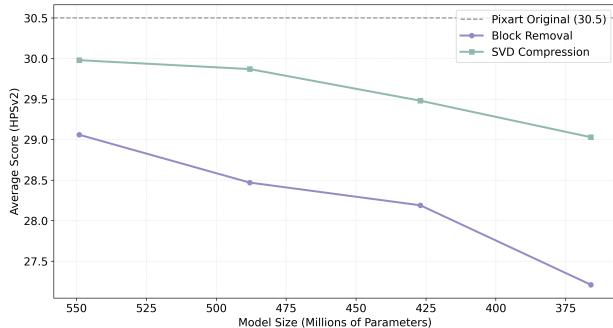
Tab. 5.3 compares the configurations of both methods, specifically the number of blocks removed or compressed and corresponding training duration for each model. Note that the number of training steps differs significantly, as the default criterion - a CMMRD value exceeding 0.1 - was employed to the pruning rate per iteration, consistent with the default settings. The observation that the SVD allows for more parameters to be removed with a less severe impact on the CMMRD value on the small reference dataset is an initial indication that compression is a more nuanced method than block removal. Fig. 5.15 presents the performance for both methods on the HPSv2 and GenEval benchmarks. For both benchmarks the SVD compression method outperforms the complete block removal method in low and high compression regimes. This suggests that distributing the parameter reduction across multiple blocks via compression is less adverse to model performance compared to removing individual blocks completely.

For a qualitative comparison of both methods, Fig. 5.16 provides a set of example images for compression ratios ranging from 90% to 50%. The advantage of the SVD compression method becomes particularly evident in the high-compression regimes (60% and 50% compression ratios). Its images contain more details, exhibit fewer artefacts and are more closely resemble to the original image. Specifically, column two demonstrates that the woman's face is accurately preserved up to 60% compression ratio when applying the SVD compression technique. Moreover, column one shows that in high-compression regimes, complete block removal can lead to severe deviations compared to the original image, in contrast to the SVD compression technique. Additionally, column five reveals that the images generated via complete block removal exhibit reduced sharpness and the plants seem to be blurry compared to those produced by SVD compression.

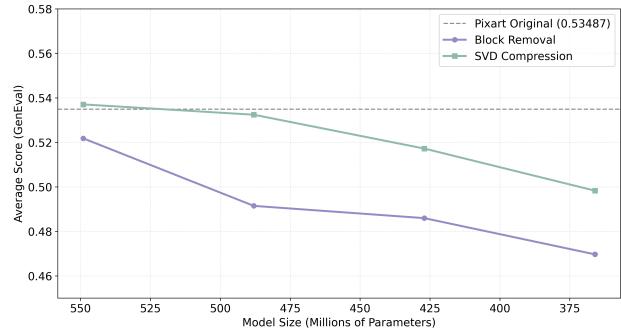
Interestingly, at a 90% and 80% compression ratio, the SVD compression technique leads to a valid but semantically altered image in column one, where the man's head is replaced by that of an alien. However, at different compression stages the model recovers and correctly generates a human head again. In conclusion, the qualitative results confirm the quantitative results of the benchmarks and demonstrate a clear advantage of SVD compression over entire block removal.

Table 5.3: Comparing training steps and modified blocks for block removal versus SVD compression method.

Model Variant (Remaining %)	Number of Removed Blocks	Number of SVD Compressed Blocks	Training Steps (in Thousands)
Model 90%	3	-	88.5
SVD Model 90%	-	6	88.5
Model 80%	6	-	265.5
SVD Model 80%	-	11	88.5
Model 70%	8	-	531.0
SVD Model 70%	-	17	265.5
Model 60%	11	-	708.0
SVD Model 60%	-	22	442.5



(a) HPSv2 Benchmark



(b) GenEval Benchmark

Figure 5.15: Impact of different compression strategies on performance across varying compression ratios.

CHAPTER 5. EXPERIMENTS

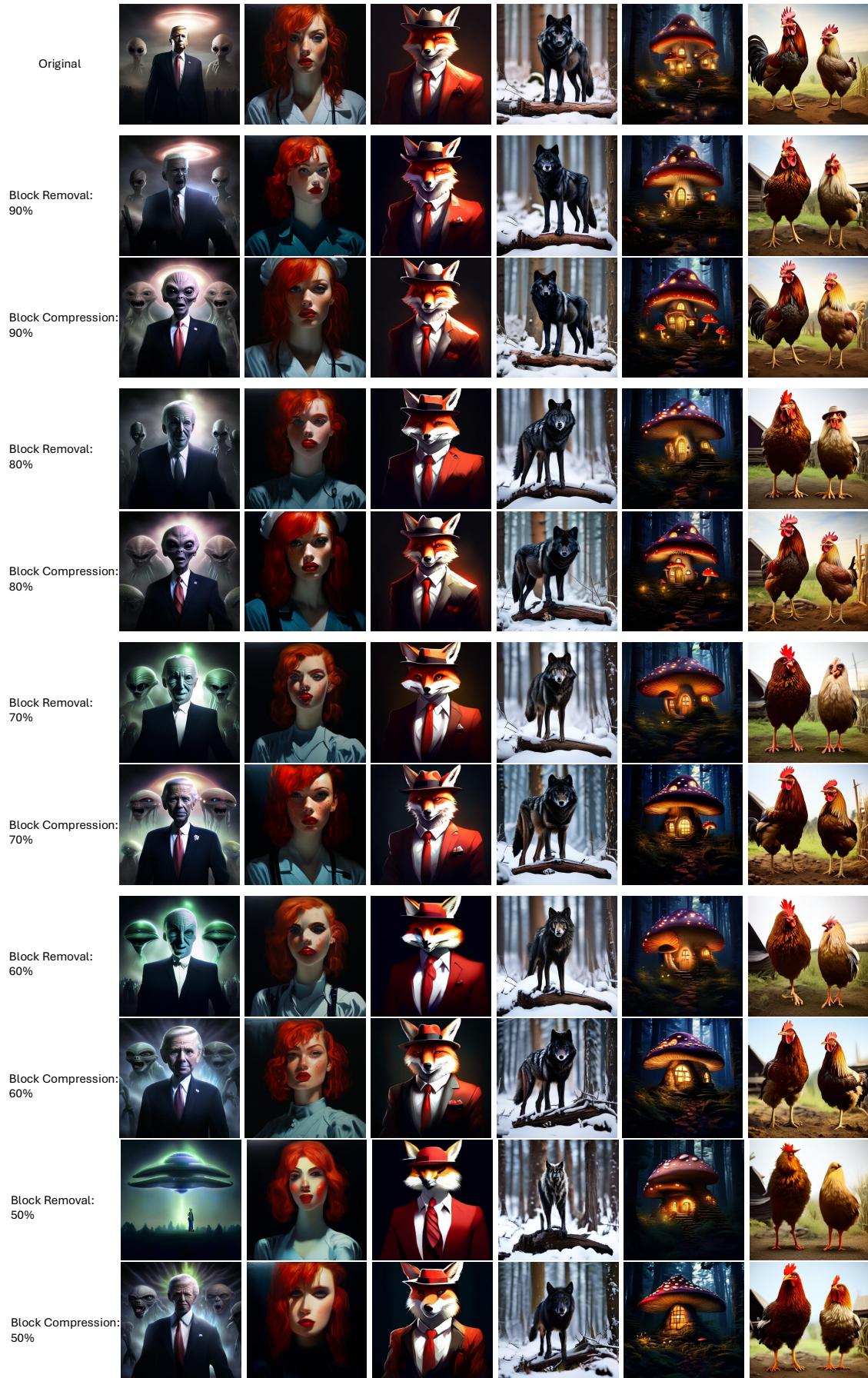


Figure 5.16: Qualitative comparison between both compression strategies: Complete block removal versus block compression with SVD

5.1.8 Sensitivity Analysis of MLP and Attention Layers under SVD Compression

In this section the SVD compression strategy is applied but otherwise the default settings remains. The goal is to investigate how to distribute the parameter removal of the individual components of one transformer blocks best. Therefore, individual compression ratios are applied to the MLP matrix and the attention matrices whereby both self- and cross-attention matrices are pruned by the same amount. In total three different scenarios are explored:

- **Equal Compression:**

Both MLP and attention matrices are compressed by the same percentage $p_{\text{MLP}} = p_{\text{Attn}} = 60\%$.

- **Attention-Dominant Compression:**

The attention matrices are more strongly compressed $p_{\text{Attn}} = 80\%$ compared to the MLP matrix $p_{\text{MLP}} = 40\%$.

- **MLP-Dominant Compression:**

The MLP matrix are more strongly compressed $p_{\text{Attn}} = 80\%$ compared to the attention matrices $p_{\text{MLP}} = 40\%$.

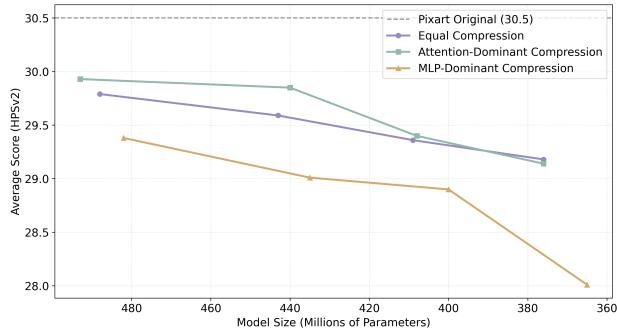
Note that the combined parameter count of self- and cross-attention is approximately equal to the parameter count of the MLP. Therefore, all configurations remove a similar amount of total parameters of one transformer block.

The scores in fig. 5.17 reveal two distinct trends. The HPSv2 score demonstrates that prioritizing the compression of the attention matrices consistently yields higher image quality than the method focused on compressing the MLP matrices. The GenEval benchmark, however, shows that although the attention-dominant approach is advantageous in the low-compression regime, it declines sharply under heavier compression, eventually falling below the MLP-dominant approach for models with approximately 400 million parameters.

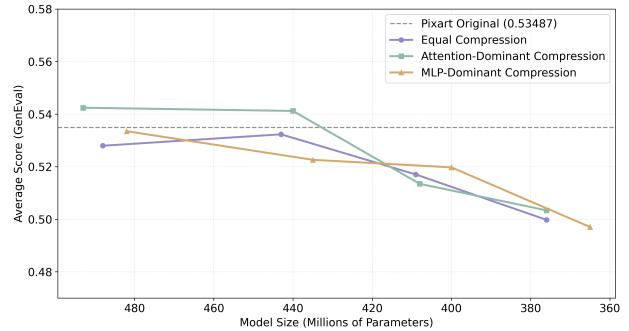
One possible explanation for the divergence between the two benchmarks is the distribution of tasks among the different components of a transformer block. A recent study suggests [40] that MLPs are crucial for knowledge storage and the processing of visual features. Consequently, strong compression of these layers leads to a significant decrease in image quality (see HPSv2 benchmark yellow line). Fig. 5.18 visually supports the benchmark results. The images clearly demonstrate that the MLP-dominant approach lacks detail, texture, and structural integrity in the high compression regime. At model sizes below 71%, the MLP-dominant approach hallucinates an unprompted headband in the samurai image (column 4) and the facial features become distorted at the highest compression ratio. Moreover, for the 60% model, the burger (column 2) is hovering in the air instead of being placed on the table, and the humanoid (column 6) is completely distorted. In contrast, the models of the attention-dominant approach, despite being approximately of the same size, do not indicate such severe structural errors, although a general degeneration of image quality at higher compression rate is also visible. In contrast, the GenEval benchmark measures text-image alignment. Since the cross-attention mechanism is responsible for integrating the text prompts into the visual latents its compression in many transformer blocks severely compromise this linking. While the image quality suffers less than with the strong MLP compression, the model's text comprehension deteriorates disproportionately. This explains why the attention-dominant approach loses its advantage for heavy compression, and is in some cases even worse than the MLP-dominant approach in the GenEval benchmark.

The equal compression strategy exhibits results similar to the attention-dominant compression approach in the high-compression regimes but is slightly worse in the low-compression regime.

CHAPTER 5. EXPERIMENTS



(a) HPSv2 Benchmark



(b) GenEval Benchmark

Figure 5.17: Impact of varying degrees of compression of individual transformer block components on performance across varying compression ratios.



Figure 5.18: Example images for differently strong compression of MLP and Attention matrices. Note that due to the slightly different matrix sizes of MLP and attention there is a different in model size by 1% between both methods.

5.1.9 SVD Compression Scheduling

In previous experiments, SVD compression was applied to each block at most once. However, restricting compression to untouched blocks may be suboptimal. Incrementally reducing the size of an already compressed block could potentially preserve model performance better than pruning a dense block. For this reason, the linear progressive compression scheme defined in Sec. 4.6.3 is investigated. Assuming that a fixed number of parameters should be removed per iteration, a design choice is how to distribute them across the model. One can either compress many blocks slightly or compress a few blocks heavily. To investigate which works better, the linear progressive compression scheme is performed for four different cases where the number of compressed blocks per iteration differs:

- **SVD_7**: Seven blocks are compressed per iteration.
- **SVD_14**: Fourteen blocks are compressed per iteration.
- **SVD_21**: Twenty-one blocks are compressed per iteration.
- **SVD_28**: All blocks are compressed per iteration.

As block selection criteria the parameter-weighted CMMD value is utilized as in previous experiments in this Master's thesis to SVD compression technique and within each iteration 10% of the original number of parameters are removed.

Untersuche ob immer neue Blöcke gepruned werden oder eher dieselben mehrfach

Trainiere bestes Model und Gebe ergebnisse auf HPSv2, GenEval und DPG an

fine-tuning auf Mapillary und Cityscapes

Evaluiere Bildqualität, Text-Generierungsfähigkeiten, Bilder auf Mapillary und Cityscapes

5.1.10 Summary of Best Settings for Structural Model Pruning

Considering the results of all experiments, they lead to the following best configuration for structural model pruning:

- **Block Importance Analysis**

The experiments (sec. 4.2) showed that using the CMMD value as primary block selection criteria is working well. Particularly in linear SVD compression scheduling, where the number of compressed parameters varies per block, the CMMD can be weighted by the number of removed parameters to ensure a fair comparison.

As block selection method, we found that the greedy algorithm already performs well enough such that more advanced methods like Optuna do not provide substantial benefits that would justify the additional compute.

- **Pruning Method**

The progressive model pruning has proven advantageous over the one-shot model pruning in the high-compression regime where models are compressed more than 70% of the original size (sec. 5.1.4). Therefore, this is the first choice.

- **Knowledge Distillation Loss**

The results indicates that the combination of the original diffusion loss together with the normalized feature and output loss leads to the strongest performance in all pruning regimes (sec. 5.1.5).

- **fine-tuning Protocol**

The finding suggests that incorporating the LAION-PixArt dataset into the training procedure is critical as it contains images of higher quality demonstrated in the higher HPSv2 and GenEval scores for the three and two stage fine-tuning protocols. However, the results regarding the superiority of three stage protocol over the two stage protocol are inconclusive. It is not clear if local fine-tuning is beneficial or not for high compression regimes (sec. 5.1.6). Therefore, the two stage approach is chosen for the final setting.

- **Structural Compressin Strategy**

The experiments demonstrated that leveraging SVD compression over complete block removal leads to superior performance even for less training steps (sec. 4.4).

- **Sensitivity Analysis of MLP and Attention Layers under SVD Compression**

The sensitivity analysis for MLP and attention layers revealed that strong compression of MLPs led to more severe deterioration than heavy compression of attention layers. However, compressing both MLP and attention layers equally yielded performance comparable to the strong attention compression experiment. Consequently, the equal compression scheme was adopted for the final configuration.

- **SVD Compression Scheduling**

5.2 Flux.1-dev Pruning

Flux.1-dev is a 11.9 billion parameter rectified flow model capable of generating detailed, high-quality, and high-aesthetic images. Due to its high parameter count, it is one primary target for model reduction techniques in the literature and should be used to demonstrate the effectiveness of the best settings found for PixArt- Σ by comparing it to various competitors. Furthermore, the method is also applied to Flux.1-schnell to demonstrate that it also works for step-distilled models. In addition, as training of Flux.1-dev is extremely computationally intensive, the training-free GRASP method (see Sec. 4.4.4) originally developed for LLM compression is transferred to Flux1.dev.

5.2.1 Training Details

In this section, all experiments are built upon the training and inference code for Flux.1-dev published in the following git repository [89]. The hyperparameters used for fine-tuning Flux.1-dev are specified in Tab. 5.4. For pruning Flux.1-dev the best settings found for PixArt- Σ are leveraged and specified here again for completeness:

- **Block Importance Analysis**

To determine the importance of individual blocks in Flux.1-dev, every block is individually compressed by 60% and a small reference dataset (100 JoyCaption prompts based on LAION dataset) is utilized to compute the CMMMD score. Afterwards, the CMMMD score is weighted by the number of removed parameters, e.g. a single-stream block contains less parameters than a double-stream block. In contrast to the experiments with PixArt- Σ , the scoring is computed in a single initial pass rather than leveraging a dynamic greedy algorithm that would iteratively reevaluate the remaining blocks after each sequential compression due to the high computational demands of Flux.1-dev.

- **Pruning Method**

Progressive structural model pruning is performed.

- **Knowledge Distillation Loss**

The original flow matching loss as well as the normalized intermediate feature and output losses are utilized whereby the normalization is computed for features from double-stream and single-stream block separately

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda_{\text{OutKD}} \mathcal{L}_{\text{OutKD}} + \lambda_{\text{SingleFeatKD}} \mathcal{L}_{\text{SingleFeatKD}} + \lambda_{\text{DoubleFeatKD}} \mathcal{L}_{\text{DoubleFeatKD}} \quad (5.1)$$

with $\lambda_{\text{OutKD}} = 10$, $\lambda_{\text{SingleFeatKD}} = 0.01$ and $\lambda_{\text{DoubleFeatKD}} = 0.1$.

- **fine-tuning Protocol**

The complete model is directly retrained on the LAION-Flux.1-dev dataset with image resolution of 1024×1024 .

- **Structural Compression Strategy**

The SVD is leveraged for compressing Flux.1-dev blocks. The individual components of double-stream and single-stream blocks and their corresponding parameter counts can be found in Tab. 4.2.

- **SVD Compression Scheduling**

The linear progressive pruning scheme is utilized to compress Flux.1-dev

A specific

Table 5.4: Configuration details for the compressed model retraining, including optimizer settings and flow matching specific parameters. These parameters represent the configuration setup for the training using the git repository [89].

Hyperparameter	Value
<i>Optimization & Scheduler</i>	
Optimizer Type	Adafactor
Optimizer Arguments	relative_step=False, scale_parameter=False, warmup_init=False
Learning Rate	1.8×10^{-6}
LR Scheduler	constant_with_warmup
LR Warmup Steps	4240
Gradient Accumulation	4
Batch Size	1
<i>Model Precision & Hardware</i>	
Mixed Precision	b1f16
Weight Precision (Base)	fp8
Full bfloat16 Training	True
Gradient Checkpointing	True
SDPA (Scaled Dot Product Attn)	True
GPU	1 × H200
<i>Flow Matching & Loss Settings</i>	
Discrete Flow Shift	3
Timestep Sampling	sigmoid
Max Timestep	1000
Loss Type	L2
Huber Parameter (c)	0.1
Huber Schedule	snr
Model Prediction Type	raw
Guidance Scale	1.0
Sample Sampler	euler_a
<i>Architecture & Data Settings</i>	
T5 Max Token Length	512
Clip Skip	1
Apply T5 Attn Mask	True
Data Loader Workers	6
Seed	42

5.2.2 Qualitative Block Analysis

To identify blocks of lower importance in the Flux.1-dev architecture, the CMMD values weighted by the removed parameters are calculated, as previously discussed in Sec. 4.2.1. To obtain a better understanding of the impact of the individual blocks, a qualitative block analysis is performed for Flux.1-dev as in Sec. 5.1.2 for PixArt- Σ .

Fig. 5.19 presents images generated after removing a single-stream block comprising approximately 142 million parameters. In most cases, removing a single block does not lead to observable artifacts. The direction in which the tiger head points may change, but visible artifacts or severe changes are only observed when removing single-stream block 35, which results in vertical stripes, or block 38, which leads to overly smoothed details.

Similarly, only the removal of a few double-stream blocks, containing approximately 340 million parameters, leads to severe degradation of the generated images (see Fig. 5.20). In particular, double-stream blocks 1, 2 and 3 seem to be critical, because the image generation capabilities of Flux.1-dev completely collapse after their removal. A possible reason is that early blocks are essential for the text-image alignment and the global layout structure. In addition, the removal of the last double-stream block leads to a notable lighting change, while the

content of the image remains intact.

Overall, the double-stream blocks seem to have a greater impact on the image generation quality, which may be explained by their higher parameter count. Interestingly, for the double-stream blocks, a similar U-shaped importance ranking emerges, as previously noted for PixArt- Σ indicating that the transformer blocks at the very beginning and end are the most critical.

Fig. 5.21 and Fig. 5.22 present visual results for compressing individual blocks by 60% rather than removing them entirely. Similar to PixArt- Σ , blocks whose removal leads to severe degradation in image quality, can be compressed without creating artifacts, e.g. single-stream blocks 35 and 38 or double-stream blocks 1, 2 and 19. The only double-stream block which appears to be highly critical and even compression exerts a substantial impact on the generated image quality is double-stream block number 3, whose compression results in visible stripes.

Ultimately, the visual robustness of the compressed double-stream blocks strongly aligns with the parameter-weighted CMMD predictions displayed in Fig. 5.23. It shows that all double-stream blocks except block number 3 are the first choice for compression. However, no clear importance structure within the double- or single-stream blocks can be identified.

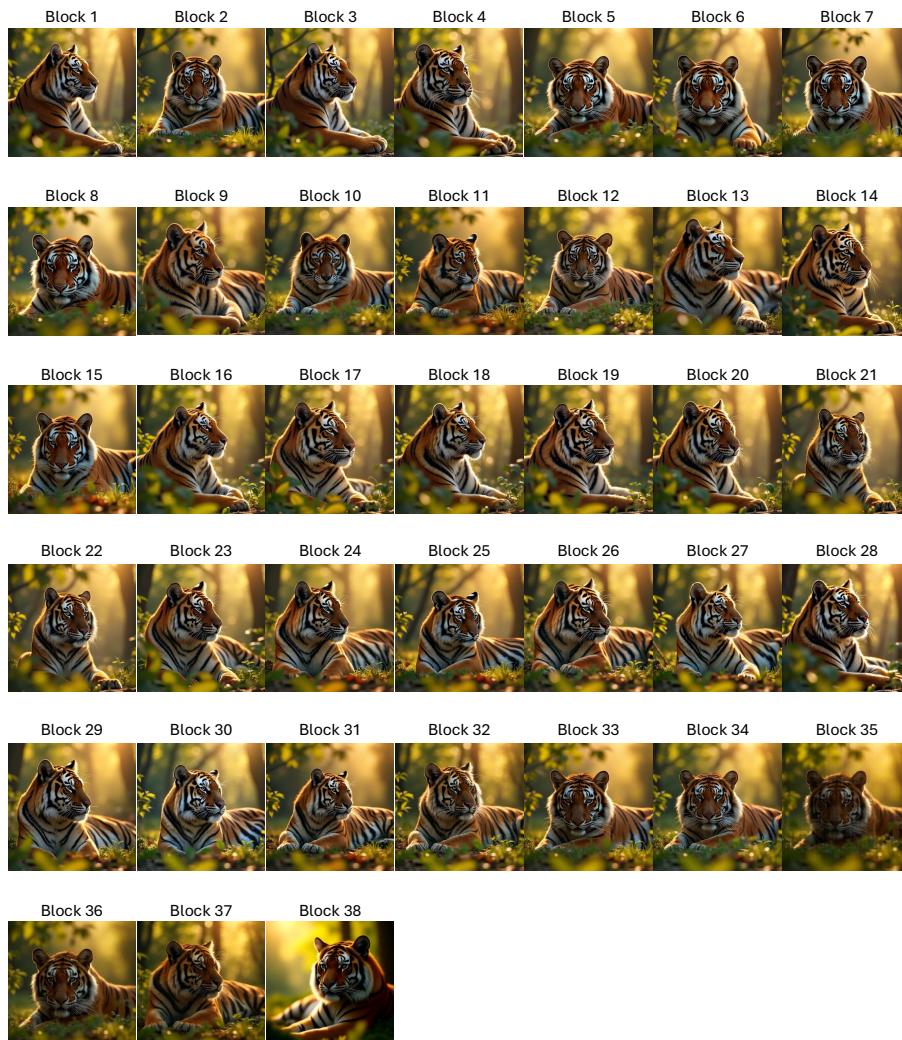


Figure 5.19: Qualitative analysis of the impact of removing individual single-stream blocks from the Flux.1-dev architecture.

CHAPTER 5. EXPERIMENTS

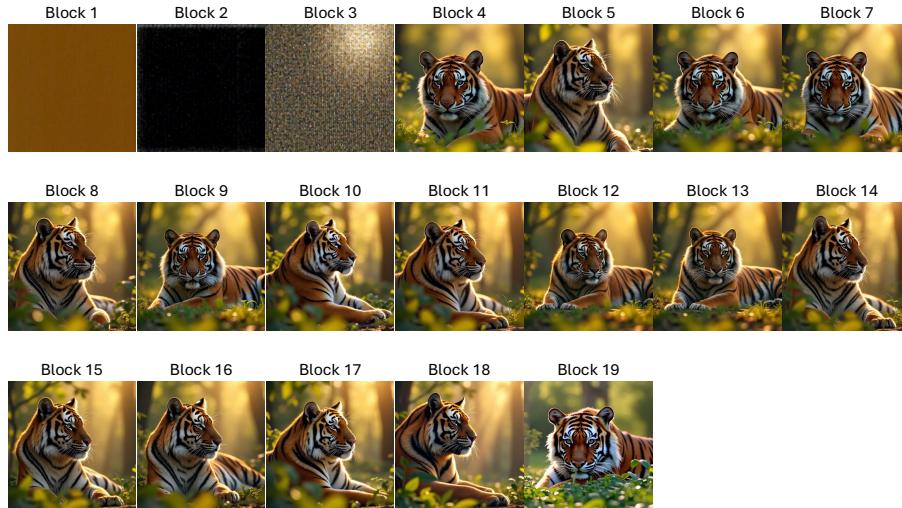


Figure 5.20: Qualitative analysis of the impact of removing individual double-stream blocks from the Flux.1-dev architecture.



Figure 5.21: Qualitative analysis of the impact of compressing individual single-stream blocks from the Flux.1-dev architecture by 60%.

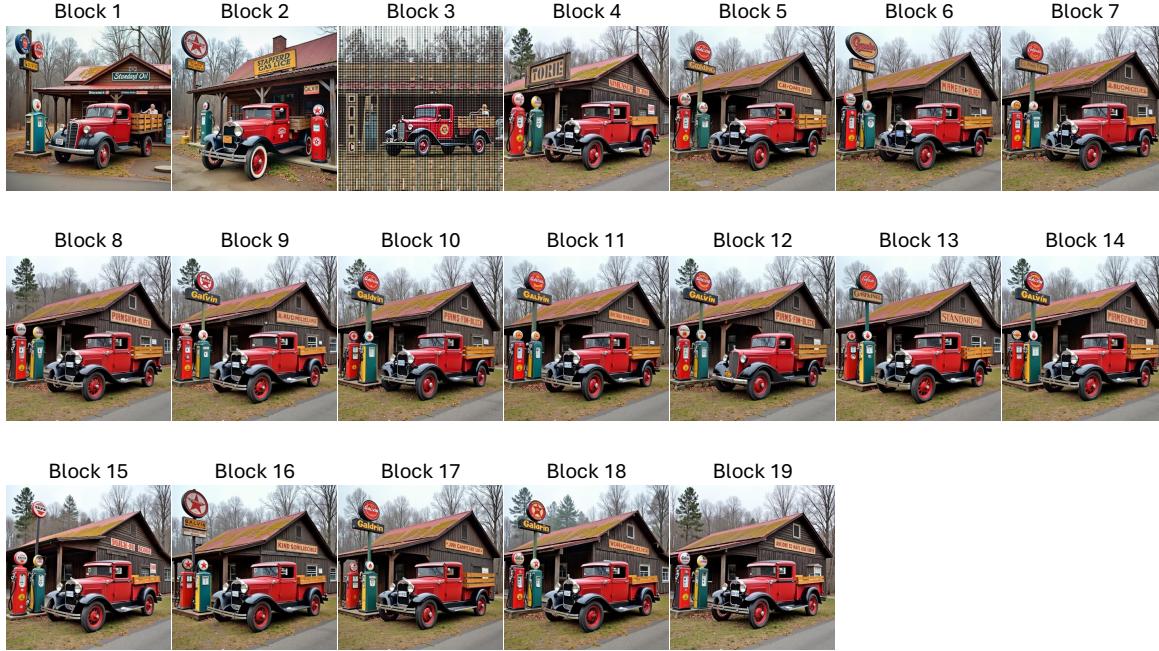


Figure 5.22: Qualitative analysis of the impact of compressing individual double-stream blocks from the Flux.1-dev architecture by 60%.

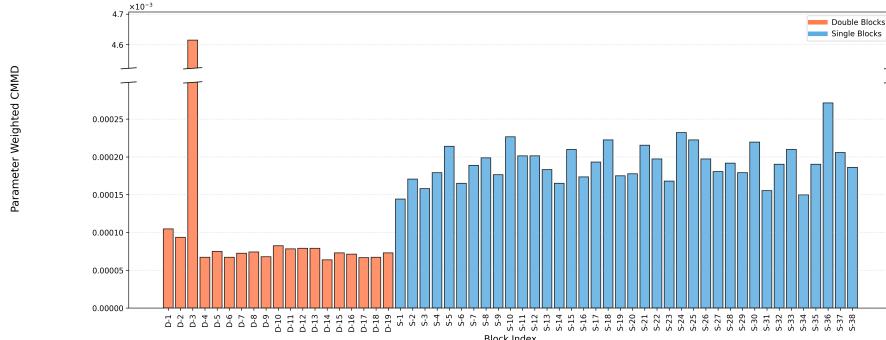


Figure 5.23: CMMD values scaled by the number of removed parameters for double- and single-stream blocks compressed by 60% of Flux.1-dev .

5.2.3 Training-Free Pruning

In this section, GRASP (see Sec. 4.4.4) is transferred from LLMs to Flux.1-dev. The goal is to compress Flux.1-dev via SVD but identify singular values based on their importance for the image generation task rather than solely based on their magnitude as is standard practise in ordinary SVD compression. The performance of GRASP is compared to Eco-Diff [196] another recently proposed training-free method to prune Flux.1-dev and to the simple application of SVD without retraining.

First, the extent to which the gradient-based method GRASP selects different singular values than the standard magnitude-based SVD selection is examined. This experiment focuses exclusively on the compression of the Flux.1-dev double-stream blocks, as they represent the model’s most parameter-intensive components. Exclusively for this setup a pre-existing importance ranking for the Flux.1-dev transformer blocks from [17] is utilized. In total, four different compression levels are tested. The compression rate for the double-stream blocks is fixed to 76% meaning the levels vary solely in the number of blocks subject to compression. For the exact ranks used during SVD for the individual components please refer to Tab. 6.1.

Fig. 5.24 illustrates exemplarily the magnitude and the importance estimated by GRASP for each singular value of the image modulation matrix from double-stream block number four. The total view (Fig. 5.24a) reveals a few significant outliers with very high magnitude and importance. Therefore, to investigate the correlation between magnitude and importance for the majority of singulare values Fig. 5.24b provides a zoomed-in version. Here, one can observe a general positive trend where higher magnitude correlates with higher importance. However, there are numerous singular values with low magnitude but high importance and high magnitude values with negligible important values. This discrepancy indicates that GRASP might effectively identify important features that magnitude-based methods would miss, potentially leading to an improvement over ordinary SVD.

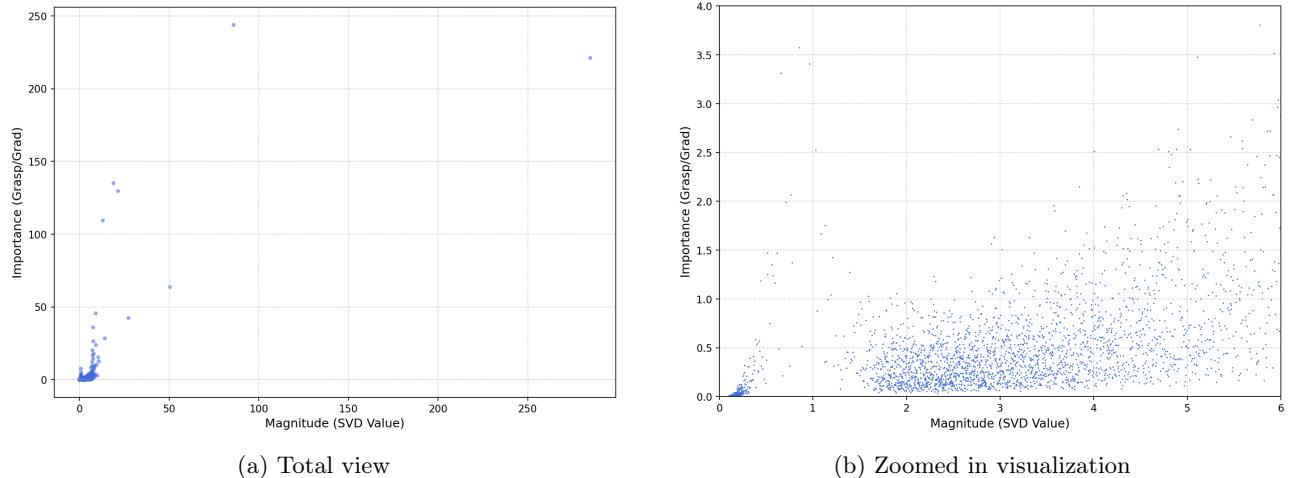
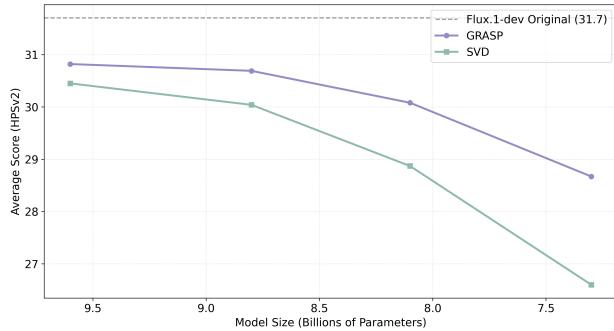
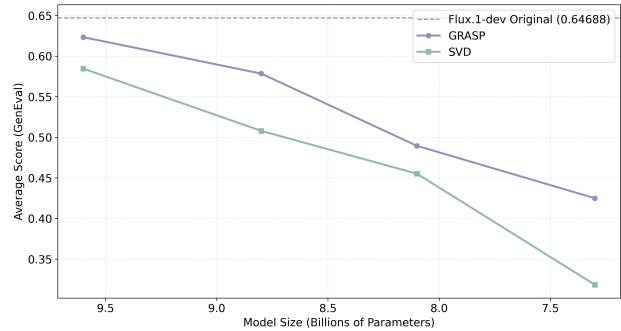


Figure 5.24: Comparing the magnitude (x-axis) with the gradient-based importance (y-axis) of singular values from SVD. The visualization is from the 4th block of the image modulation matrix.

The potential for improvements over SVD implied by the analysis of the correlation between magnitude and importance is confirmed by the evaluation of both methods on the HPSv2 and GenEval benchmarks (see Fig. 5.25). The results show that models obtained using the most important singular values (GRASP) instead of those with the highest magnitude outperform standard SVD across all compression ratios on both benchmarks. Fig. 5.26 displays example images, demonstrating that especially for high compression ratios (68% and 60%) the divergence between both methods is pronounced. SVD based models exhibit significant degradation in image quality and detail at these compression stages. Furthermore, their text-generation capabilities degrades more severely than those of GRASP based models as the kangaroo image still correctly displays "Welcome Friends" at the strongest compression for GRASP but not for SVD (see second column in Fig. 5.26). However, models compressed with GRASP also struggles to generate text consistently correctly at high compression (see fourth columns) and for the highest compression rate some severe artifacts (see third column) can be observed. Nevertheless, both image quality and text generation capabilities are superior with the GRASP method. As image details are typically high-frequency information the results suggest that the high-magnitude singular values primarily encode low-frequency information and the low-magnitude singular values often discarded by ordinary SVD store the high-frequency information. By prioritizing gradients, GRASP better preserves these subtle but important singular values and, thereby mainiting superior text-image alignment and preserving finer details.



(a) HPSv2 Benchmark



(b) GenEval Benchmark

Figure 5.25: Impact of singular value selection based on magnitude (SVD) or importance (GRASP) on performance across varying compression ratios.

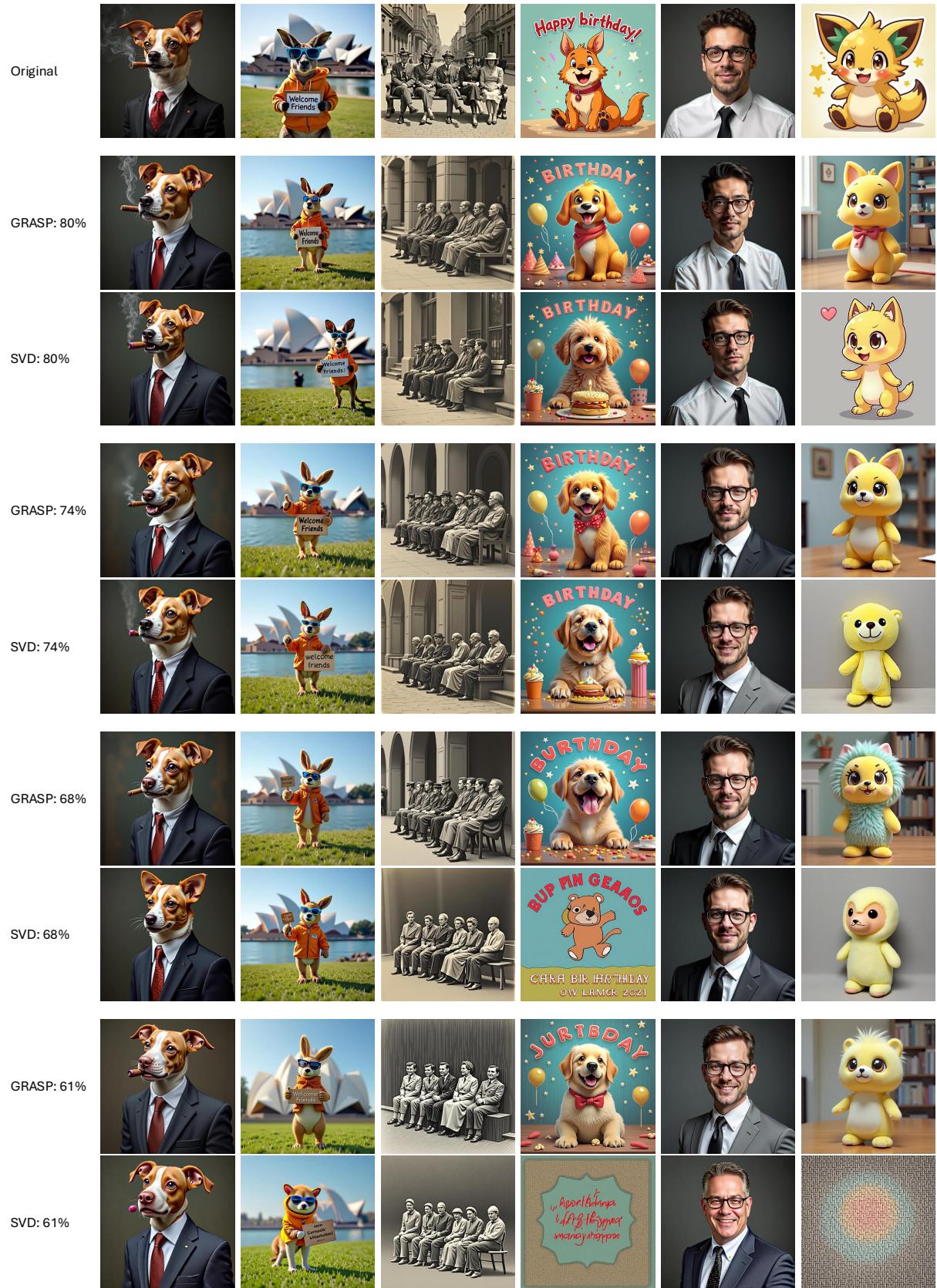


Figure 5.26: Images from compression Flux.1-dev with SVD or GRASP respectively.

Finally, GRASP and SVD are compared to EcoDiff (see Tab. 5.8) demonstrating that these approaches are

significantly superior in maintaining image quality for the same degree of parameter reduction. For a parameter reduction of 20%, the SVD-based methods obtain a GenEval score of 0.58 and 0.62 whereas EcoDiff drops to 0.18 indicating a loss of text-alignment. Fig. 5.28 presents example images from the GenEval Benchmark providing visual proof of the strong degeneration and poor text-alignment for the compressed Flux.1-dev model when reduced to less than 90% of its original size using EcoDiff. In addition, Fig. 5.27 presents images testing the text-generation capabilities of models compressed by 20% leveraging EcoDiff, SVD, or GRASP as the compression methods. It clearly demonstrates that EcoDiff-compressed models completely fail to generate text and visually appealing images. Overall, GRASP-compressed models are slightly superior to simple SVD-compressed models in the text generation task (columns 2 and 4) as the generated text contains fewer errors. However, column 4 also indicates that the GRASP-compressed model also exhibits minor flaws. Similarly, the HPSv2 score of EcoDiff measuring human preferences is only 23.6 while SVD and GRASP reach 30.45 and 30.82. Notably, when compressing Flux.1-dev by 30% utilizing the EcoDiff method, the model completely collapses. Visual inspection shows that the model ceases to generate meaningful images, producing grey images with undefined texture, in contrast to the performance observed with GRASP. This performance gap indicates that structurally removing entire neurons breaks the Markovian generation trajectory of Flux.1-dev and validates the effectiveness of the general concept of SVD that approximates the original information flow instead of relying on the hard removal of neurons as in EcoDiff.

Nevertheless, while all methods effectively reduce the VRAM consumption during inference, standard implementations of SVD and GRASP do not inherently reduce the inference time as EcoDiff. This is because low-rank decomposition replaces one large matrix by two smaller matrices which doubles the number of matrix multiplications and only leads to runtime savings after stronger compression (see Sec. ??).

Table 5.5: Performance Comparison of Compression Methods across HPSv2, GenEval, and DPG Benchmarks.

Ratio	Method	HPSv2 ↑	GenEval ↑	DPG ↑
90%	EcoDiff	29.67	0.452	77.98
	SVD	30.86	0.653	82.77
	GRASP	31.08	0.659	
80%	EcoDiff	23.633	0.187	49.12
	SVD	30.45	0.585	79.84
	GRASP	30.82	0.623	81.19
70%	EcoDiff	7.26	0.002	17.02
68%	SVD	28.87	0.455	73.13
68%	GRASP	30.08	0.490	75.50



Figure 5.27: The images provide examples for text generation capabilities of models compressed to 80% of the original parameters using EcoDiff, SVD or GRASP as compression method.

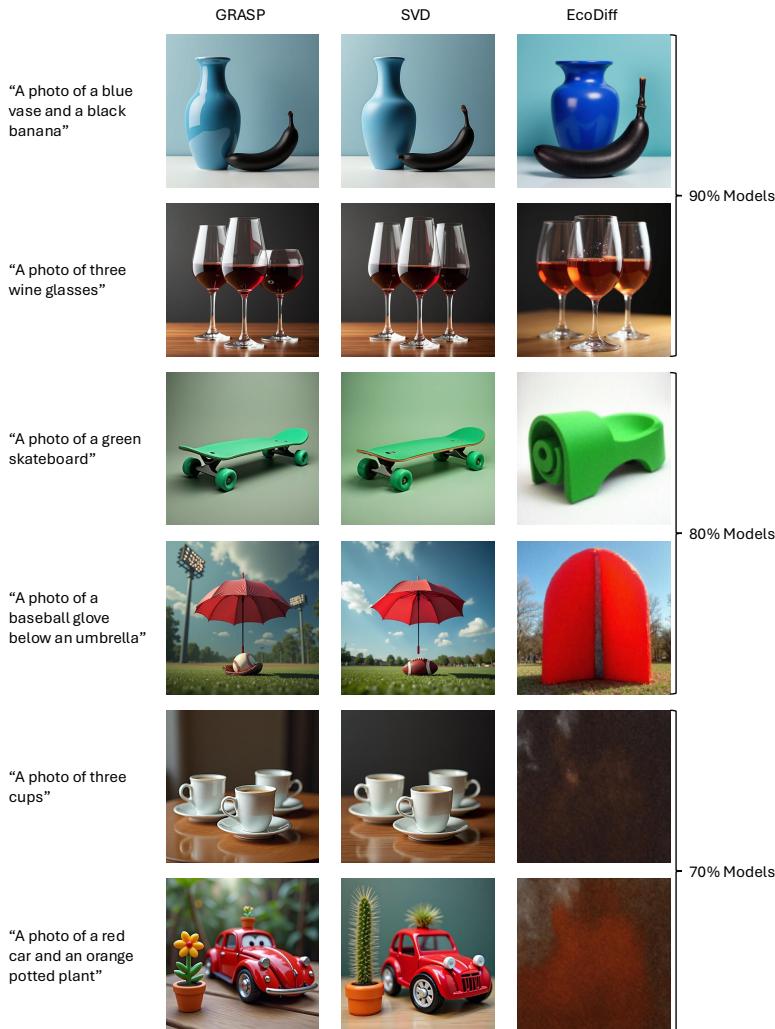


Figure 5.28: Examples from GenEval benchmark for Flux.1-dev compressed by 10%, 20% and 30% using EcoDiff, SVD and GRASP.

5.2.4 Linear Progressive SVD Compression

This section applies the optimal pruning strategies identified for PixArt- Σ to Flux.1-dev and Flux.1-schnell. Prior to presenting the final results, two preliminary experiments are conducted. First, it is demonstrated that block importance computation on a reference dataset comprising only 100 images proves advantageous over randomly block selection. Second, following the methodology in Sec. 4.6.3 the optimal distribution of removed parameters across the Flux.1-dev transformer blocks is analyzed.

Block Importance Analysis

The image generation process using Flux.1-dev is computationally demanding and time-consuming. Consequently, the reference dataset utilized for block importance determination is limited to 100 images. Furthermore, block importance is only computed once per iteration, rather than being re-evaluated iteratively as in the experiments of PixArt- Σ . Typically, image quality metrics are computed on several thousands of images to ensure high diversity and statistical reliability. However, [77] demonstrated that CMMD is effective even with small sample sizes and is considerably more reliable than FID. Thus, the following experiment demonstrates that despite the constrained reference dataset and the one-time calculation of the block importance per iteration, the analysis remains robust and yields superior results compared to random block selection.

In both experiments, the default training settings (see Sec. 5.2.1) are used applying four subsequent iterations within the linear progressive pruning framework. Details regarding the model sizes and number of training steps are provided in Tab. 5.6.

The performance over various benchmarks (see Fig. 5.29) demonstrates that importance-based block selection, even with a small calibration dataset, leads to superior results compared to random block selection. A possible reason is that the CMMD score based on 100 images already effectively indicates whether a block is in general suitable for compression or if it should be maintained. For example, as shown in Sec. 5.2.2, the third double-stream block is highly important and therefore is preserved by the metric-based approach, whereas random selection risks pruning it, potentially leading to a significant performance drop. Qualitatively, the performance difference between models obtained by random versus importance-based block selection can be observed in Fig. 5.30 and aligns with the quantitative benchmark results. The images generated by the model (5.83B parameters) based on random block selection exhibit artifacts, such as a bear with two heads or a chainsaw with two chain blades, while the importance-based approaches preserves the correct semantic structures much better. In addition, its text-generation capabilities are inferior (see the fifth column).

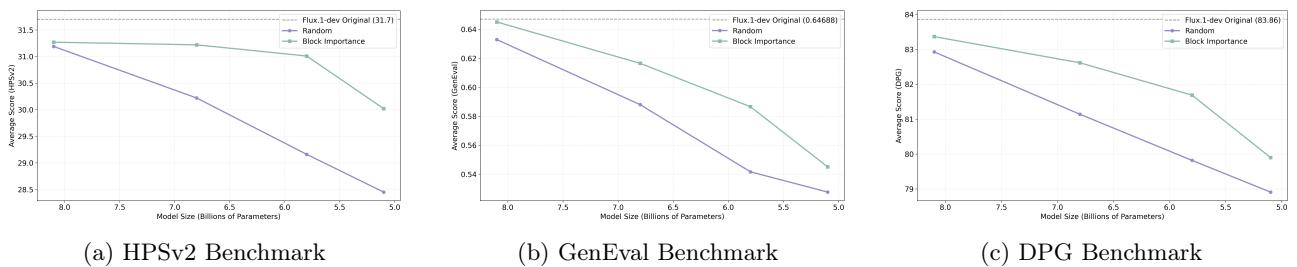


Figure 5.29: Demonstrating the advantage of metric-based block selection even with a small reference dataset.



Figure 5.30: Example images of the 49% models (5.8B parameters) obtained by random and importance-based block selection.

Table 5.6: Overview of training details similar for both experiments, random and metric-based block selection.

Iteration	Model Size (%)	Training Steps per Iteration
1	68%	40k
2	57%	20k
3	49%	20k
4	43%	20k

Sensitivity Analysis: The Role of Parameter Removal Distribution

The goal of this experiment is to investigate how the parameters to be removed should be optimally distributed across the model. Three different scenarios are examined where Flux.1-dev is pruned by 10% in each iteration with a varying number of modified blocks and consequently, varying base compression ratios:

- **20 Blocks:** In each iteration 20 transformer blocks are modified.
- **30 Blocks:** In each iteration 30 transformer blocks are modified.
- **40 Blocks:** In each iteration 40 transformer blocks are modified.

Note that the number of blocks does not specify whether double- or single-stream blocks are compressed. This depends entirely on the importance ranking of the blocks. The base compression ratio, defined as the compression ratio of the least important block, may need to be adapted across different iterations depending on which blocks are ranked as least important. For example if in the fifth iteration of the experimental setting **20 Blocks**, single-stream blocks and already heavily compressed double-stream blocks are further compressed, a higher base compression ratio is needed to remove the same number of parameters as in the first iteration. An overview of which value was chosen in each iteration can be found in Tab. 5.7. For all experiments, the models are retrained for 20k steps in each iteration.

Fig. 5.31 presents the benchmark performance of all three experiments. For all experiments, significant performance drop is observed when reducing the model size from 60% (7.14B parameters) to 50% (5.95B parameters) of the original model size across all benchmarks. Notably, according to the HPSv2 benchmark, the **20 Blocks** experiment performs worse at each compression stage than the other two settings. In contrast, its performance on the GenEval benchmark is superior at the 80% (9.52B parameters) and the 70% (8.33B parameters) pruning stage. Furthermore, on the DPG benchmark this setting shows superior performance in low-compression regime but drops below the other experiments in high-compression regime. Overall, there

appears to be a slight tendency suggesting that in highly compressed models, it might be beneficial to distribute the removed parameters more uniformly. However, a definitive conclusion cannot be drawn.

Furthermore, when comparing the **30 Blocks** and **40 Blocks** experiments, the results are very similar on the HPSv2 benchmark. For the GenEval and DPG benchmark the **30 Blocks** setting seems to be slightly superior in the low-compression regime and the **40 Blocks** setting performs better in the high-compression regime. However, a clearly superior setting cannot be determined.

To conclude, GenEval and DPG benchmark measure the text-image coherence while HPSv2 investigates human preference. It is notable that the **20 Blocks** and **30 Blocks** settings are performing better in the low-compression regime for GenEval (exception 90% **20 Blocks**) and DPG than the **40 Blocks** setting. This suggests that for prompt adherence it is advantageous when the removal of parameters is concentrated on less blocks. In contrast, the results suggest for the high-compression regime that the more uniform distribution of removed parameters is beneficial.

Regarding the HPSv2, the **40 Blocks** setting seems to perform for nearly all compression stages best which indicates that for overall image quality a uniform distribution is beneficial. In contrast, for the overall image quality a more uniform distribution of removed parameters seems beneficial.

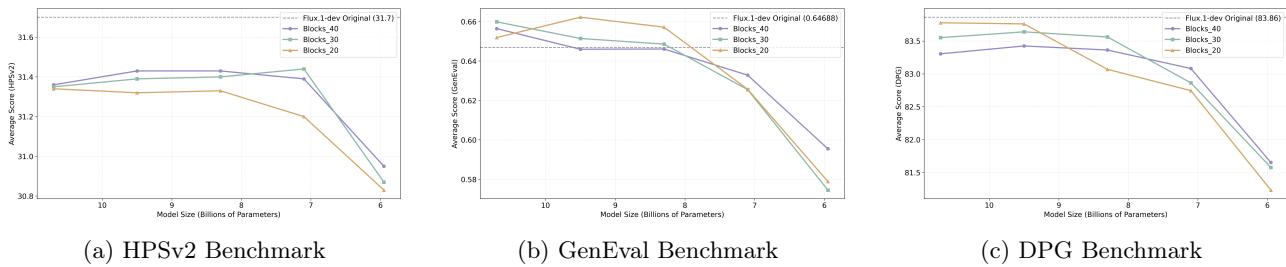


Figure 5.31: Demonstrating the advantage of metric-based block selection even with a small reference dataset.

Table 5.7: Overview of compression metrics per iteration.

Iteration	Model Size	# Blocks Modified	Base Compression Ratio
1	90%	20	0.3
		30	0.25
		40	0.2
2	80%	20	0.3
		30	0.25
		40	0.2
3	70%	20	0.3
		30	0.3
		40	0.2
4	60%	20	0.35
		30	0.35
		40	0.25
5	50%	20	0.6
		30	0.4
		40	0.35

Table 5.8: Performance Comparison of Compression Methods across HPSv2, GenEval, and DPG Benchmarks.

Ratio	Method	HPSv2 ↑	GenEval ↑	DPG ↑	GPU Memory (%) ↓
100 %	Flux.1-dev	31.70	0.647	83.86	100
68%	FluxLite [33]	31.31	0.523	79.31	-
	PPCL ¹ [118]	-	0.605	80.0	-
	EcoDiff-Lora [196]	25.99	0.399	-	-
	TinyFusion ¹ [45]	-	0.511	77.2	-
	HierarchicalPrune ¹ [92]	-	0.503	75.7	-
57%	Our	31.31	0.650	83.13	78
	Our	31.31	0.650	83.13	70
	Our	31.01	0.587	81.69	
43%	MoE [199] ²	-	0.570	81.63	
	Our	30.02	0.545	79.70	

¹ The numbers are taken from [118] and are not self computed.² The numbers are taken from [199] and are not self computed.

Chapter 6

Appendix

6.1 Appendix A

6.1.1 Block Analysis

Heatmap for cka: Block importance for different time steps

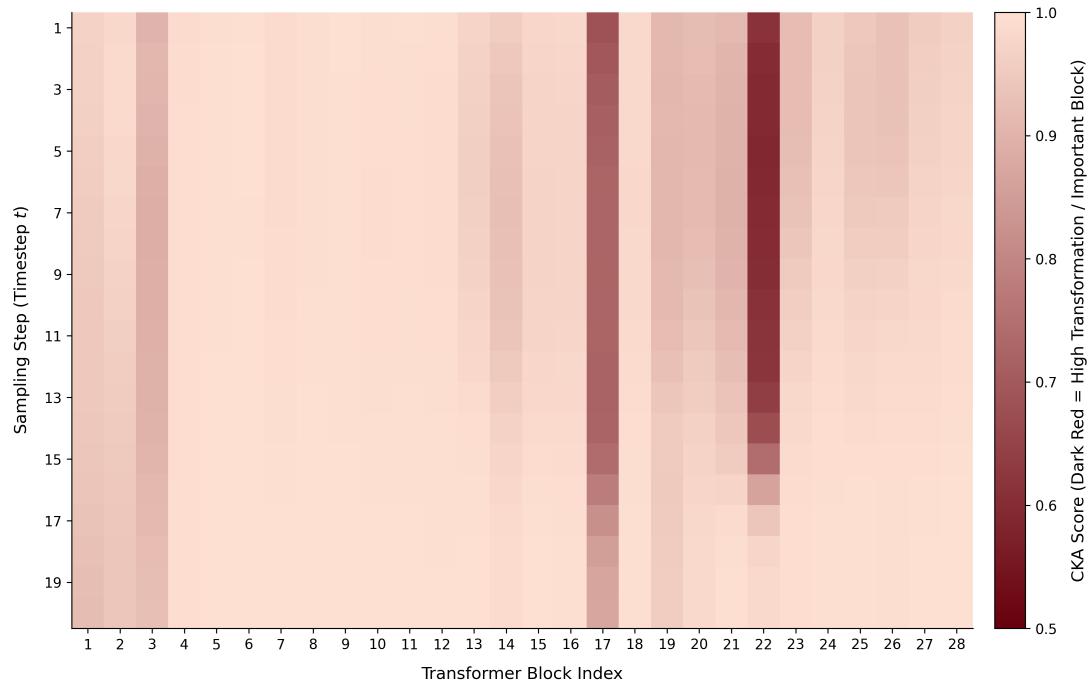


Figure 6.1: CKA-Kernel Clip

6.2 Appendix B - Detailed Overview of Selected Blocks Experiments PixArt- Σ

In this section, the removed or compressed blocks for each experiments are visualized. Different colors indicates during which iteration of the progressive pruning framework the specific block was selected.

6.3 Experiments Flux.1-dev

6.3.1 GRASP

Block Importance Ranking

For this experiment the importance of blocks determined in [17] was leveraged. Below the figure from the corresponding work is presented.

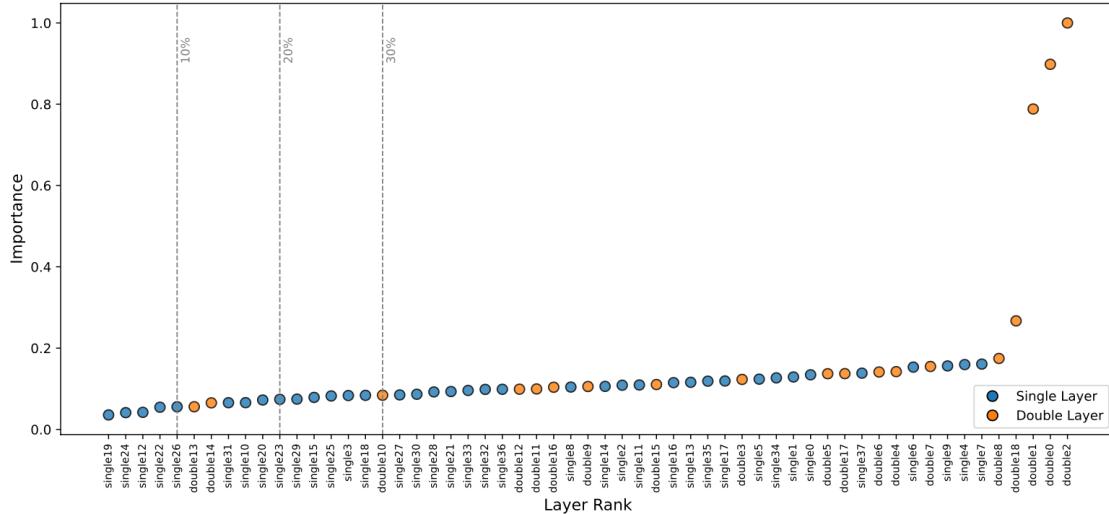


Figure 6.2: Block importance for Flux.1-dev taken from [17].

Compression Details

Table 6.1: Overview of compressed double-stream block components and the respective rank used in SVD.

Components Double Block	Rank for SVD
Double-Stream Image Linear Layer (Attention)	307
Double-Stream Image MLP	491
Double-Stream Image Modulation	526
Double-Stream Text Linear Layer (Attention)	307
Double-Stream Text MLP	491
Double-Stream Text Modulation	526

Table 6.2: Overview of Pruned Models and Block Configurations

Model Size	# Pruned	Pruned Block IDs
80%	9	13, 14, 10, 12, 11, 16, 9, 15, 3
74%	12	13, 14, 10, 12, 11, 16, 9, 15, 3, 5, 17, 6
68%	15	13, 14, 10, 12, 11, 16, 9, 15, 3, 5, 17, 6, 4, 7, 8
60%	18	13, 14, 10, 12, 11, 16, 9, 15, 3, 5, 17, 6, 4, 7, 8, 18, 1, 0

Bibliography

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, and Shyamal Anadkat. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Vivago AI. Hidream-l1. 2025. URL: <https://vivago.ai/home>.
- [3] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631, 2019.
- [4] Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- [5] Anthropic. Claude 3 haiku: our fastest model yet. 2024. URL: <https://www.anthropic.com/news/clause-3-haiku>.
- [6] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR.
- [7] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [8] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [9] Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Lluis Castrejon, Kelvin Chan, Yichang Chen, Sander Dieleman, Yuqing Du, et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024.
- [10] Nathan Bell and Michael Garland. Efficient sparse matrix-vector multiplication on cuda. Technical report, Nvidia Technical Report NVR-2008-004, Nvidia Corporation, 2008.
- [11] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [12] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. *Advances in neural information processing systems*, 24, 2011.
- [13] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [14] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.

BIBLIOGRAPHY

- [15] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Leo Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. *OpenAI Blog*, 1(8):1, 2024.
- [16] Anja Butter, Nathan Huetsch, Sofia Palacios Schweitzer, Tilman Plehn, Peter Sorrenson, and Jonas Spinner. Jet diffusion versus jetgpt—modern networks for the lhc. *SciPost Physics Core*, 8(1):026, 2025.
- [17] Fuhan Cai, Yong Guo, Jie Li, Wenbo Li, Xiangzhong Fang, and Jian Chen. Fastflux: Pruning flux with block-wise replacement and sandwich training. *arXiv preprint arXiv:2506.10035*, 2025.
- [18] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Zeroq: A novel zero shot quantization framework. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13169–13178, 2020.
- [19] Junsong Chen, Chongjian Ge, Enze Xie, et al. Pixart-sigma: Official pytorch implementation. <https://github.com/PixArt-alpha/PixArt-sigma>, 2024.
- [20] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In *European Conference on Computer Vision*, pages 74–91. Springer.
- [21] Junsong Chen, Yue Wu, Simian Luo, Enze Xie, Sayak Paul, Ping Luo, Hang Zhao, and Zhenguo Li. Pixart- δ : Fast and controllable image generation with latent consistency models. *arXiv preprint arXiv:2401.05252*, 2024.
- [22] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- [23] Lei Chen, Yuan Meng, Chen Tang, Xinzhu Ma, Jingyan Jiang, Xin Wang, Zhi Wang, and Wenwu Zhu. Q-bit: Accurate post-training quantization for diffusion transformers. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28306–28315, 2025.
- [24] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahu Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer, 2024.
- [25] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- [26] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [27] Yu-Hui Chen, Raman Sarokin, Juhyun Lee, Jiujiang Tang, Chuo-Ling Chang, Andrei Kulik, and Matthias Grundmann. Speed is all you need: On-device acceleration of large diffusion models via gpu-aware optimizations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4651–4655, 2023.
- [28] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.

BIBLIOGRAPHY

- [29] Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation. *arXiv preprint arXiv:2310.18235*, 2023.
- [30] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, and Lukasz Kaiser. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- [31] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- [32] Florinel-Alin Croitoru, Vlad Hondu, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10850–10869, 2023.
- [33] Javier Martín Daniel Verdú. Flux.1 lite: Distilling flux1.dev for efficient text-to-image generation. 2024.
- [34] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022.
- [35] MohammadReza Davari, Stefan Horoi, Amine Natik, Guillaume Lajoie, Guy Wolf, and Eugene Belilovsky. Reliability of cka as a similarity measure in deep learning. *arXiv preprint arXiv:2210.16156*, 2022.
- [36] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [37] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186.
- [38] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [39] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- [40] Yihe Dong, Lorenzo Noci, Mikhail Khodak, and Mufan Li. Attention retrieves, mlp memorizes: Disentangling trainable components in the transformer. *arXiv preprint arXiv:2506.01115*, 2025.
- [41] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, and Sylvain Gelly. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [42] DC Dowson and BV666017 Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982.
- [43] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.

BIBLIOGRAPHY

- [44] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.
- [45] Gongfan Fang, Kunjun Li, Xinyin Ma, and Xinchao Wang. Tinyfusion: Diffusion transformers learned shallow. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18144–18154, 2025.
- [46] William Feller. On the theory of stochastic processes, with particular reference to applications. In *Proceedings of the First Berkeley Symposium on Mathematical Statistics and Probability*, pages 403–432, Berkeley, CA, 1949. The Regents of the University of California.
- [47] Determine Filters’Importance. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [48] fpgaminer. JoyCaption: An open, uncensored image captioning VLM. <https://github.com/fpgaminer/joycaption>, 2024. Accessed: 2026-01-22.
- [49] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- [50] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International conference on machine learning*, pages 10323–10337. PMLR, 2023.
- [51] Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. Sparse gpu kernels for deep learning. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–14. IEEE, 2020.
- [52] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.
- [53] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [54] Jackson Gorham and Lester Mackey. Measuring sample quality with kernels. In *International Conference on Machine Learning*, pages 1292–1301. PMLR.
- [55] Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. A kernel method for the two-sample-problem. *Advances in neural information processing systems*, 19, 2006.
- [56] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *The journal of machine learning research*, 13(1):723–773, 2012.
- [57] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*, pages 63–77. Springer, 2005.
- [58] Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. Accelerate: Training and inference at scale made simple, efficient and adaptable. <https://github.com/huggingface/accelerate>, 2022.
- [59] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.

BIBLIOGRAPHY

- [60] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015.
- [61] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- [62] Yefei He, Jing Liu, Weijia Wu, Hong Zhou, and Bohan Zhuang. Efficientdm: Efficient quantization-aware fine-tuning of low-bit diffusion models. *arXiv preprint arXiv:2310.03270*, 2023.
- [63] Yefei He, Luping Liu, Jing Liu, Weijia Wu, Hong Zhou, and Bohan Zhuang. Ptqd: Accurate post-training quantization for diffusion models. *Advances in Neural Information Processing Systems*, 36:13237–13249, 2023.
- [64] Yefei He, Luping Liu, Jing Liu, Weijia Wu, Hong Zhou, and Bohan Zhuang. Ptqd: Accurate post-training quantization for diffusion models. *Advances in Neural Information Processing Systems*, 36:13237–13249, 2023.
- [65] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [66] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [67] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- [68] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, June 01, 2020 2020. URL: <https://ui.adsabs.harvard.edu/abs/2020arXiv200611239H>, doi:10.48550/arXiv.2006.11239.
- [69] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [70] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- [71] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.
- [72] Ting Hua, Xiao Li, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Dynamic low-rank estimation for transformer-based language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9275–9287, 2023.
- [73] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 603–612.
- [74] Sergey Ioffe. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [75] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2704–2713, 2018.

BIBLIOGRAPHY

- [76] Steven A Janowsky. Pruning versus clipping in neural networks. *Physical Review A*, 39(12):6600, 1989.
- [77] Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and Sanjiv Kumar. Rethinking fid: Towards a better evaluation metric for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9307–9315, 2024.
- [78] Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021.
- [79] Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellemundi, Dipankar Das, Kunal Banerjee, Sasikanth Avancha, Dharma Teja Vooturi, Nataraj Jammalamadaka, Jianyu Huang, Hector Yuen, et al. A study of bfloat16 for deep learning training. *arXiv preprint arXiv:1905.12322*, 2019.
- [80] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022.
- [81] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in neural information processing systems*, 34:852–863, 2021.
- [82] Bo-Kyeong Kim, Hyoung-Kyu Song, Thibault Castells, and Shinkook Choi. Bk-sdm: A lightweight, fast, and cheap version of stable diffusion. In *European Conference on Computer Vision*, pages 381–399. Springer.
- [83] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- [84] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013.
- [85] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023.
- [86] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in neural information processing systems*, 36:36652–36663, 2023.
- [87] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.
- [88] Shun Kiyono, Sosuke Kobayashi, Jun Suzuki, and Kentaro Inui. Shape: Shifted absolute position embedding for transformers. *arXiv preprint arXiv:2109.05644*, 2021.
- [89] kohya ss. sd-scripts. <https://github.com/kohya-ss/sd-scripts>, October 2024.
- [90] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pages 3519–3529. PMLR, 2019.
- [91] Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- [92] Young D Kwon, Rui Li, Sijia Li, Da Li, Sourav Bhattacharya, and Stylianos I Venieris. Hierarchicalprune: Position-aware compression for large-scale diffusion models. *arXiv preprint arXiv:2508.04663*, 2025.
- [93] Black Forest Labs. Flux model. 2024. URL: <https://bfl.ai>.

BIBLIOGRAPHY

- [94] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 2002.
- [95] Jaeho Lee, Sejun Park, Sangwoo Mo, Sungsoo Ahn, and Jinwoo Shin. Layer-adaptive sparsity for the magnitude-based pruning. *arXiv preprint arXiv:2010.07611*, 2020.
- [96] Youngwan Lee, Kwanyong Park, Yoorhim Cho, Yong-Ju Lee, and Sung Ju Hwang. Koala: Empirical lessons toward memory-efficient and fast diffusion models for text-to-image synthesis. *Advances in Neural Information Processing Systems*, 37:51597–51633, 2024.
- [97] B. Lefauveux, F. Massa, D. Liskovich, W. Xiong, V. Caggiano, S. Naren, M. Xu, J. Hu, M. Tintore, S. Zhang, P. Labatut, D. Haziza, L. Wehrstedt, J. Reizenstein, and G. Sizov. xformers: A modular and hackable transformer modelling library. <https://github.com/facebookresearch/xformers>, 2022.
- [98] C Li, H Xu, J Tian, W Wang, M Yan, B Bi, J Ye, H Chen, G Xu, Z Cao, et al. mplug: Effective and efficient vision-language learning by cross-modal skip-connections. arxiv 2022. *arXiv preprint arXiv:2205.12005*, 1, 2022.
- [99] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024. [arXiv:2402.17245](https://arxiv.org/abs/2402.17245).
- [100] Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdquant: Absorbing outliers by low-rank components for 4-bit diffusion models. *arXiv preprint arXiv:2411.05007*, 2024.
- [101] Shanda Li, Chong You, Guru Guruganesh, Joshua Ainslie, Santiago Ontanon, Manzil Zaheer, Sumit Sanghali, Yiming Yang, Sanjiv Kumar, and Srinadh Bhojanapalli. Functional interpolation for relative positions improves long context transformers. *arXiv preprint arXiv:2310.04418*, 2023.
- [102] Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17535–17545, 2023.
- [103] Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snapfusion: Text-to-image diffusion model on mobile devices within two seconds. *Advances in Neural Information Processing Systems*, 36:20662–20678, 2023.
- [104] Tatiana Likhomanenko, Qiantong Xu, Gabriel Synnaeve, Ronan Collobert, and Alex Rogozhnikov. Cape: Encoding relative positions with continuous augmented positional embeddings. *Advances in Neural Information Processing Systems*, 34:16079–16092, 2021.
- [105] Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion distillation. *arXiv preprint arXiv:2402.13929*, 2024.
- [106] Zinan Lin, Ashish Khetan, Giulia Fanti, and Sewoong Oh. Pacgan: The power of two samples in generative adversarial networks. *Advances in neural information processing systems*, 31, 2018.
- [107] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [108] Bingchen Liu, Yizhe Zhu, Kunpeng Song, and Ahmed Elgammal. Towards faster and stabilized gan training for high-fidelity few-shot image synthesis. In *iclr*.
- [109] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.

BIBLIOGRAPHY

- [110] Kainan Liu, Yong Zhang, Ning Cheng, Zhitao Li, Shaojun Wang, and Jing Xiao. GRASP: Replace redundant layers with adaptive singular parameters for efficient model compression. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 26344–26359, 2025.
- [111] Kainan Liu, Yong Zhang, Ning Cheng, Zhitao Li, Shaojun Wang, and Jing Xiao. Grasp: Replace redundant layers with adaptive singular parameters for efficient model compression. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 26344–26359, 2025.
- [112] Shiwei Liu, Tianlong Chen, Xiaohan Chen, Li Shen, Decebal Constantin Mocanu, Zhangyang Wang, and Mykola Pechenizkiy. The unreasonable effectiveness of random pruning: Return of the most naive baseline for sparse training. *arXiv preprint arXiv:2202.02643*, 2022.
- [113] Shiwei Liu, Tim Van der Lee, Anil Yaman, Zahra Atashgahi, Davide Ferraro, Ghada Sokar, Mykola Pechenizkiy, and Decebal Constantin Mocanu. Topological insights into sparse neural networks. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 279–294. Springer, 2020.
- [114] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- [115] Xuanqing Liu, Hsiang-Fu Yu, Inderjit Dhillon, and Cho-Jui Hsieh. Learning to encode position for transformer with continuous dynamical model. In *International conference on machine learning*, pages 6327–6335. PMLR.
- [116] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- [117] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- [118] Jian Ma, Qirong Peng, Xujie Zhu, Peixing Xie, Chen Chen, and Haonan Lu. Pluggable pruning with contiguous layer distillation for diffusion transformers. *arXiv preprint arXiv:2511.16156*, 2025.
- [119] Jian Ma, Qirong Peng, Xujie Zhu, Peixing Xie, Chen Chen, and Haonan Lu. Pluggable pruning with contiguous layer distillation for diffusion transformers. *arXiv preprint arXiv:2511.16156*, 2025.
- [120] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European conference on computer vision (ECCV)*, pages 116–131, 2018.
- [121] Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect, 2024. URL <https://arxiv.org/abs/2403.03853>, 2(3):4, 2024.
- [122] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306.
- [123] Michael C Mozer and Paul Smolensky. Using relevance to reduce network size automatically. *Connection Science*, 1(1):3–16, 1989.

BIBLIOGRAPHY

- [124] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In *International conference on machine learning*, pages 7197–7206. PMLR, 2020.
- [125] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kortschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Proceedings of the IEEE international conference on computer vision*, pages 4990–4999, 2017.
- [126] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR.
- [127] Oscar Nordström. Unstructured pruning of pre-trained language models tuned for sentiment classification., 2022.
- [128] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- [129] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International conference on machine learning*, pages 4055–4064. PMLR.
- [130] Jan Pawłowski and Tilman Plehn. Physics and machine learning. Lecture script, Institut für Theoretische Physik, Universität Heidelberg, January 2024. Version from January 30, 2024.
- [131] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.
- [132] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [133] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- [134] Ian Pons, Bruno Yamamoto, Anna H Real Costa, and Artur Jordao. Effective layer pruning through similarity metric perspective. In *International Conference on Pattern Recognition*, pages 423–438. Springer, 2024.
- [135] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.
- [136] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [137] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [138] Samyam Rajbhandari, Olatunji Ruwase, Jeff Rasley, Shaden Smith, and Yuxiong He. Zero-infinity: Breaking the gpu memory wall for extreme scale deep learning. In *Proceedings of the international conference for high performance computing, networking, storage and analysis*, pages 1–14, 2021.

BIBLIOGRAPHY

- [139] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in neural information processing systems*, 32, 2019.
- [140] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in neural information processing systems*, 32, 2019.
- [141] Sebastian Raschka. Understanding and coding the self-attention mechanism of large language models from scratch. 2023. URL: <https://sebastianraschka.com/blog/2023/self-attention-from-scratch.html>.
- [142] Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. {Zero-offload}: Democratizing {billion-scale} model training. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, pages 551–564, 2021.
- [143] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR.
- [144] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695.
- [145] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, December 01, 2021 2021. CVPR 2022. URL: <https://ui.adsabs.harvard.edu/abs/2021arXiv211210752R>, doi:10.48550/arXiv.2112.10752.
- [146] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, pages 234–241. Springer.
- [147] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation, 2023. URL: <https://arxiv.org/abs/2208.12242>, arXiv:2208.12242.
- [148] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [149] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- [150] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- [151] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11.
- [152] Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. In *International conference on machine learning*, pages 30105–30118. PMLR.
- [153] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pages 87–103. Springer.

BIBLIOGRAPHY

- [154] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
- [155] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-Aesthetic v1. <https://github.com/LAION-AI/laion-datasets/blob/main/laion-aesthetic.md>, 2022. Accessed: 2026-01-22.
- [156] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019.
- [157] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr.
- [158] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr.
- [159] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, October 01, 2020 2020. ICLR 2021; updated connections with ODEs at page 6, fixed some typos in the proof. URL: <https://ui.adsabs.harvard.edu/abs/2020arXiv201002502S>, doi:10.48550/arXiv.2010.02502.
- [160] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
- [161] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [162] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [163] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [164] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in artificial intelligence*, pages 574–584. PMLR.
- [165] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [166] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [167] Inga Strümke and Helge Langseth. Lecture notes in probabilistic diffusion models. *arXiv preprint arXiv:2312.10393*, 2023.
- [168] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

BIBLIOGRAPHY

- [169] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- [170] Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. Ominicontrol: Minimal and universal control for diffusion transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14940–14950, 2025.
- [171] Siao Tang, Xin Wang, Hong Chen, Chaoyu Guan, Yansong Tang, et al. Lightweight diffusion models with distillation-based block neural architecture search. *arXiv preprint arXiv:2311.04950*, 2023.
- [172] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Jo-han Schalkwyk, Andrew M Dai, Anja Hauth, and Katie Millican. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [173] TencentARC. flux-mini. <https://huggingface.co/TencentARC/flux-mini>, 2024. Hugging Face Model Hub.
- [174] Hoang Thanh-Tung and Truyen Tran. Catastrophic forgetting and mode collapse in gans. In *2020 international joint conference on neural networks (ijcnn)*, pages 1–10. IEEE.
- [175] Luke Tierney. Markov chains for exploring posterior distributions. *the Annals of Statistics*, pages 1701–1728, 1994.
- [176] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, and Shruti Bhosale. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [177] Aaron Van Den Oord and Oriol Vinyals. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [178] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [179] Xin Wang, Samiul Alam, Zhongwei Wan, Hui Shen, and Mi Zhang. SVD-LLM v2: Optimizing singular value truncation for large language model compression. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4287–4296, 2025.
- [180] Xin Wang, Samiul Alam, Zhongwei Wan, Hui Shen, and Mi Zhang. Svd-llm v2: Optimizing singular value truncation for large language model compression. *arXiv preprint arXiv:2503.12340*, 2025.
- [181] Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 893–911, 2023.
- [182] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. *Advances in neural information processing systems*, 29, 2016.
- [183] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.

BIBLIOGRAPHY

- [184] Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, and Yao Lu. Sana: Efficient high-resolution image synthesis with linear diffusion transformers. *arXiv preprint arXiv:2410.10629*, 2024.
- [185] Chenglin Yang, Celong Liu, Xueqing Deng, Dongwon Kim, Xing Mei, Xiaohui Shen, and Liang-Chieh Chen. 1.58-bit flux. *arXiv preprint arXiv:2412.18653*, 2024.
- [186] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- [187] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and Bill Freeman. Improved distribution matching distillation for fast image synthesis. *Advances in neural information processing systems*, 37:47455–47487, 2024.
- [188] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6613–6623.
- [189] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfai Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- [190] Zhihang Yuan, Yuzhang Shang, Yue Song, Dawei Yang, Qiang Wu, Yan Yan, and Guangyu Sun. ASVD: Activation-aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023.
- [191] Zhihang Yuan, Yuzhang Shang, Yue Song, Dawei Yang, Qiang Wu, Yan Yan, and Guangyu Sun. Asvd: Activation-aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023.
- [192] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in neural information processing systems*, 32, 2019.
- [193] Dingkun Zhang, Sijia Li, Chen Chen, Qingsong Xie, and Haonan Lu. Laptop-diff: Layer pruning and normalized distillation for compressing diffusion models. *arXiv preprint arXiv:2404.11098*, 2024.
- [194] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023.
- [195] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. *arXiv preprint arXiv:2204.13902*, 2022.
- [196] Yang Zhang, Er Jin, Wenzhong Liang, Yanfei Dong, Ashkan Khakzar, Philip Torr, Johannes Stegmaier, and Kenji Kawaguchi. Learnable sparsity for vision generative models. *arXiv preprint arXiv:2412.02852*, 2024.
- [197] Zhihua Zhang. The singular value decomposition, applications and beyond. *arXiv preprint arXiv:1510.08532*, 2015.
- [198] Yang Zhao, Yanwu Xu, Zhisheng Xiao, Haolin Jia, and Tingbo Hou. Mobilediffusion: Instant text-to-image generation on mobile devices. In *European Conference on Computer Vision*, pages 225–242. Springer, 2024.

BIBLIOGRAPHY

- [199] Youwei Zheng, Yuxi Ren, Xin Xia, Xuefeng Xiao, and Xiaohua Xie. Dense2moe: Restructuring diffusion transformer to moe for efficient text-to-image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18661–18670, 2025.