

The Hidden Cost of an Image: Quantifying the Energy Consumption of AI Image Generation

Giulia Bertazzini^a, Chiara Albisani^a, Daniele Baracchi^a, Dasara Shullani^a,
Roberto Verdecchia^a

^a*Department of Information Engineering, University of Florence, Florence, 50139, Italy*

Abstract

With the growing adoption of AI image generation, in conjunction with the ever-increasing environmental resources demanded by AI, we are urged to answer a fundamental question: What is the environmental impact hidden behind each image we generate? In this research, we present a comprehensive empirical experiment designed to assess the energy consumption of AI image generation. Our experiment compares 17 state-of-the-art image generation models by considering multiple factors that could affect their energy consumption, such as model quantization, image resolution, and prompt length. Additionally, we consider established image quality metrics to study potential trade-offs between energy consumption and generated image quality. Results show that image generation models vary drastically in terms of the energy they consume, with up to a 46x difference. Image resolution affects energy consumption inconsistently, ranging from a 1.3x to 4.7x increase when doubling resolution. U-Net-based models tend to consume less than Transformer-based one. Model quantization instead results to deteriorate the energy efficiency of most models, while prompt length and content have no statistically significant impact. Improving image quality does not always come at the cost of a higher energy consumption, with some of the models producing the highest quality images also being among the most energy efficient ones.

Keywords: Diffusion models, green ai, energy consumption, image quality assessment.

Email addresses: giulia.bertazzini@unifi.it (Giulia Bertazzini), chiara.albisani@unifi.it (Chiara Albisani), daniele.baracchi@unifi.it (Daniele Baracchi), dasara.shullani@unifi.it (Dasara Shullani), roberto.verdecchia@unifi.it (Roberto Verdecchia)

1. Introduction

Text-to-image models have seen rapid adoption in the most recent years and are, with high probability, here to stay. Albeit the numerous benefits image generation models offer, a looming question remains to date unanswered: *What is the environmental impact of image generation models?* The ever-increasing environmental impact of Artificial Intelligence (AI) is becoming an established concern in academic literature that can no longer be neglected [1]. Although past studies covered a wide range of topics related to the environmental sustainability of AI, to the best of our knowledge, no research effort so far considered the energy required by image generation models.

In this research, we aim to gain a comprehensive understanding of the energy consumption of image generation models by presenting a comprehensive empirical experiment comprising 17 distinct state-of-the-art diffusion models. The experiment is conducted considering various aspects that may influence model energy efficiency, such as model quantization, image resolution, and prompt length. To gain systematic and sound insights, our research method comprises multiple experimental re-runs and statistical analyses, leading to the execution of more than 9k distinct experimental measurements. In addition to model energy consumption, we complement our results by evaluating also the quality of the generated images, achieved *via* the use of consolidated metrics specifically designed to measure the quality of generated images.

On one hand, our research targets researchers interested in understanding and improving the environmental sustainability of image generation, by providing a comprehensive analysis on the factors influencing energy consumption of this process. On the other hand, our study is intended to reach a wider audience, by providing concrete evidence grounded in empirical data of the to date overlooked environmental cost hidden behind image generation.

To support scrutiny of our results, aligned with open science principles, we make a comprehensive supplementary package of our study available in our companion material.

1.1. Related Work

Recent research has highlighted the substantial carbon footprint and high energy consumption associated with training large deep learning models. However, relatively few studies focus on their efficiency during inference, despite its critical role in real-world deployment.

Gowda et al. [2] analyze the energy consumption of various deep learning models, providing a detailed trade-off between accuracy and efficiency. Their study quantifies the energy consumption for both training and testing in terms of backbone architecture, dataset, GFLOPs, and parameter count, by disregarding however image generation tasks.

A large-scale assessment on the environmental impact of machine learning is conducted by Castano et al. [3]. They evaluated the carbon footprint of 1,417 models hosted on Hugging Face, covering diverse domains such as multimodal models, audio, and computer vision. Notably, they find no statistically significant differences in carbon footprint between pre-trained and fine-tuned models.

As other pioneering study in this space, Budenny et al. [4] assesses the CO_2 emissions of text-to-image generative models. They focus only on two models, Malevich and Kandinsky [5], evaluating their consumption based on training epochs, loss, GPU/CPU usage, and batch size.

Among the models considered in this work, only PixArt- α [6] reported CO_2 emissions. We therefore assert that, while commonly computational cost of training is considered – typically in terms of GPU hours or Neural Function Evaluations (NFEs) – the environmental impact of models is most commonly overlooked.

Synthetic image generation models are usually evaluated using both quantitative metrics, such as FID (Fréchet Inception Distance) [7], Precision-Recall [8], IS (Inception Score) [9], and CLIPScore [10], as well as qualitative user studies. FID and Precision-Recall require a reference dataset, representing the real data distribution against which distances are measured. In contrast, IS operates without the need of a reference dataset. All these discussed metrics rely on features extracted by a pre-trained InceptionV3 [11] network. Unlike them, CLIPScore does not assess image quality directly but instead measures how well the generated images align with the given prompts.

Some research has explored the impact of image resolution on generation quality. PixArt- σ [12] finds that increasing resolution improves both FID and CLIP scores, testing images at 512 and 1024 pixels with aspect ratios between 1 and 9. Similarly, Lumina [13] provides qualitative comparisons of images upscaled from 512 to 2048 pixels, benchmarking against PixArt- α [6]. Stable Diffusion 1.5 [14] examines $4\times$ upscaling with FID and IS metrics.

Beyond resolution, studies investigate how architectural choices influence generative performance. Factors such as number of diffusion steps, guidance scale, and network type (e.g., UNet or LoRA) are frequently considered. Stable Diffusion XL Turbo [15] explores variations in loss functions, discriminator types, conditioning methods (image- or text-based), and initialization strategies. Stable

Diffusion 3 [16] studies the impact of model depth and guidance weight on FID, while Flash Stable Diffusion [17] and Flash Stable Diffusion XL [17] analyze the effects of timestamp sampling.

To the best of our knowledge, no prior work has systematically examined the energy consumption of image generation by analyzing the interplay between model architecture, quantization, resolution, and prompt length. In this study, we provide new insights into the environmental impact of diffusion models during inference, while also investigating the relationship between a model’s environmental footprint and the quality of the generated images.

2. Preliminaries

In this section, we provide an overview of the key concepts related to image generation with diffusion models, that are at the basis of this study. Specifically, we focus on two main classes of diffusion models: **U-Net-based models**, that leverage U-Net backbone, and **Transformer-based models**, that replace the U-Net with a Vision Transformer.

2.1. U-Net-based Diffusion Models

Diffusion models have emerged as a leading paradigm in image generation. Originally introduced by [18], Denoising Diffusion Probabilistic Models (DDPMs) [19] involve a forward diffusion process and a reverse denoising process. These early models are defined as *unconditional*, as they generate images without external input.

In the *forward diffusion process*, Gaussian noise $\epsilon \sim \mathcal{N}(0, I)$ is progressively added to an input image $x_0 \sim q(x_0)$, sampled from the data distribution, over T timesteps, following a noise schedule. By timestep T , the image is transformed into pure noise. This process is modeled as a Markov chain, defined by $q(x_1, \dots, x_T) := \prod_{t=1}^T q(x_t|x_{t-1})$ where $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$. In this formulation, t and β_t denote the timestep and noise schedule, respectively.

On the other side, the *reverse denoising process* consists of removing the noise and reconstructing the original image x_0 , starting from the completely noisy version, $x_T \sim \mathcal{N}(0, I)$. This is done by iteratively denoising each timestep, using a U-Net based neural network predicting the mean and the variance of the denoised image at each timestep: $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$, where $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ represent the mean and variance of the Gaussian distribution.

To improve the efficiency of diffusion model training on limited computational resources without compromising quality and flexibility, [14] introduced latent diffusion models (LDMs). LDMs use a Variational Autoencoder (VAE) to encode the high-resolution input image into a low-dimensional latent representation, allowing the U-Net to perform the diffusion process with reduced computational complexity. Additionally, a CLIP-based text encoder [20] converts input text prompts into embeddings that condition the U-Net and effectively guide the diffusion process. These kind of models are categorized as *conditional*.

2.2. Transformer-based Diffusion Models

Transformer-based diffusion models [21] replace the commonly-used U-Net backbone with a transformer [22] that operates on latent patches. Starting from Gaussian noise, a transformer network reverses the diffusion process to generate the target image, leveraging diffusion timesteps to generate distinct features at each stage.

These models support class conditioning, enabling to generate images that correspond to specific class labels. This feature allows the network to produce images from predefined categories by conditioning the transformer on both the class label and the timestep embedding, making it highly flexible for controlled image generation.

The architecture of these models follows the Vision Transformers (ViT) [23], combining multi-head self-attention (MHSA) and a multi-layer perceptron (MLP) blocks, described as $X \leftarrow X + \alpha \text{MHSA}(\gamma X + \beta)$, $X \leftarrow X + \alpha' \text{MLP}(\gamma' X + \beta')$, where $X \in \mathbb{R}^{N \times C}$ represents image tokens, with N as the number of tokens and C as the channel dimension. The parameters $\alpha, \gamma, \beta, \alpha', \gamma', \beta'$ come from adaptive layer normalization (adaLN) integrating class condition embedding E_{cls} and timestep embedding E_t .

3. Methodology

This section documents the methodology used to study the energy consumption of different diffusion models during the image generation process. To ensure in-depth analysis, a wide set of diffusion models was selected. The study investigated how changes of key factors, such as models, quantization, image resolution, and prompt length, impact energy consumption. In addition to the energy consumption analysis, we also conducted a quality assessment of the generated

images to explore the relationship between model energy consumption and output quality. The details of our experimental setup are presented in the following sections.

3.1. Experimental Variables

To evaluate the energy consumption of different diffusion models during the image generation process, we selected a diverse set of diffusion models, including both U-Net-based and Transformer-based ones. These categories were based on their different backbone architecture, resulting in distinct image generation processes and therefore potentially different energy consumption patterns.

In designing our experiments, we identified key variables to assess their influence on energy consumption while ensuring the fairest comparison possible between models. To analyze model behavior across multiple realistic usage scenarios, we selected **model quantization**, **image resolution**, and **prompt length** as the primary experimental variables. Conversely, we excluded the number of diffusion steps, as each diffusion model is typically designed to operate with a specific step count, which can range from as few as 2 steps to 50 or more. Running a model designed for 2 steps with 50 steps, or *vice versa*, would not only degrade image quality but also introduce significant biases in the comparison, as models are not intended to function under such conditions.

3.1.1. Diffusion Models Selection

For this study, we selected a total of **17 diffusion models**, based on their underlying architectures: U-Net-based and Transformer-based models. By including models from both groups, we aimed to capture a comprehensive range of diffusion-based generation techniques.

In selecting the models, we prioritized state-of-the-art models over older ones, with a focus on the models that are more widely used. We also included a few slightly outdated models that remain popular and continue to be highly downloaded by end users. In this way, we ensured a representation of both cutting-edge technologies and well-established ones for our experimentation. Furthermore, we selected models that are free for academic use and accessible offline, i.e., excluding subscription-based solutions such as DALL·E [24] and Midjourney [25], which are not publicly accessible and hence we could not experiment on. All chosen models are available on Hugging Face, allowing for their standardized use and minimizing implementation variability. The models selected for our experimentation are:

U-Net-based: Stable Diffusion 1.5 (SD_1.5) [14], Stable Diffusion XL (SDXL) [26], Stable Diffusion XL Turbo (SDXL_Turbo) [15], Stable Diffusion XL Lightning (SDXL_Lightning) [27], Hyper Stable Diffusion (Hyper_SD) [28], Segmind Stable Diffusion 1B (SSD_1B) [29], Latent Consistency Model Segmind Stable Diffusion 1B (LCM_SSD_1B) [30], Latent Consistency Model Stable Diffusion XL (LCM_SDXL) [30], Flash Stable Diffusion (Flash_SD) [17], Flash Stable Diffusion XL (Flash_SDXL) [17].

Transformer-based: PixArt- α (PixArt_Alpha) [6], PixArt- σ (PixArt_Sigma) [12], Flash PixArt (Flash_PixArt) [17], Stable Diffusion 3 (SD_3) [16], Flash Stable Diffusion 3 (Flash_SD3) [17], Lumina-Next-SFT (Lumina) [13], Flux.1 schnell (Flux_1) [31].

Many of the models considered in this study are derivatives of a base model that has been distilled or otherwise modified to enhance image quality and improve efficiency. For instance, Stable Diffusion XL serves as the foundation for multiple variants, including Turbo, Lightning, LCM, and Flash. While these versions all originate from the same base model, they integrate distinct optimizations and refinements that influence the image generation process, potentially leading to energy consumption variations.

3.1.2. Model Quantization

To address the substantial memory demands of diffusion models, quantization emerged as an effective technique for reducing memory usage while preserving high model performance. Given the widespread use of quantization, we study the energy consumption of each investigated model **in both quantized and unquantized forms**. For this purpose, we utilized Quanto¹, a quantization toolkit built on PyTorch provided through Hugging Face Optimum, a suite of tools for hardware optimization. Among the available quantization options, we selected `int8` quantization, as it is a widely used approach that significantly reduces memory usage while maintaining acceptable image quality.

3.1.3. Image Resolution

Image resolution is a critical factor that influences both computational complexity and visual quality of generated images. Higher resolutions generally produce higher visual quality, while being characterized by increased computational requirements, and hence higher memory consumption, processing times, and potentially energy consumption.

¹<https://github.com/huggingface/optimum-quanto>

In our study, we considered two widely used square resolutions, 512×512 and 1024×1024 , along with a non-square resolution, namely 768×1024 . This approach allowed us to explore the potential impact of aspect ratio variations on generation performance and visual quality of the images.

3.1.4. Prompt Length

Understanding whether longer prompts introduce unexpected computational overhead or efficiency variations in the image generation process is important for optimizing prompt engineering strategies for different applications. Subtle variations in processing time, memory usage, or energy consumption could arise when handling more complex or descriptive input prompts.

To evaluate the effect of prompt length on energy consumption, we considered three different lengths: short, medium, and long. For **short-length prompts**, we randomly selected 30 classes from CIFAR-100 [32] dataset and we used them without any adaptation as short prompts. Therefore, short prompts consist of one or two words as presented in the original dataset. The CIFAR-100 prompts were then expanded into **long-length prompts** by using the Phi-3-Mini-4K-Instruct [33] large language model (LLM), a lightweight and state-of-the-art open model, specifically designed for general-purpose AI systems. We used 13 examples from Stable Diffusion’s best-performing photographic prompts [34] as a reference, instructing the LLM to add additional details to the input short prompt. Therefore, long prompts consists of an average of 33 words. Finally, long prompts were manually truncated to obtain **medium-length prompts**, consisting of an average of 14 words.

3.2. Experiment Execution

We carried out a preliminary test to evaluate the influence of prompt semantic content on energy consumption. We selected as prompts the 100 distinct classes of the CIFAR-100 dataset, leveraging its broad class variety. Each model processed the prompts 10 times in randomized order to mitigate ordering biases. Then, to ensure consistency, all models were tested using their default settings, avoiding potential parameter-related influences.

The results of the preliminary test indicated no significant correlation between the semantic content of the prompts and energy consumption, confirming that semantic variability is not a critical factor for this study. Further details on the preliminary test are provided in are provided in the supplementary material of this work.

Building upon the findings of the preliminary test, we conducted the main experiments of this study, which focused on measuring energy consumption for each model while varying the independent variables outlined in Section 3.1. To this end, we randomly selected 30 prompts from the CIFAR-100 dataset classes to form the short-length prompts. Then we expanded and truncated these prompts to obtain the long-length and medium-length prompt, following the approach detailed in Section 3.1.4. Similar to the preliminary experiment, to prevent biases associated to prompt ordering and boundary effects, the prompts were randomly shuffled before each experiment.

Our setup involved generating images using 17 different diffusion model under their default settings while varying two quantization parameters, three image resolution levels, and three prompt lengths. As a results, each prompt underwent $17 \times 2 \times 3 \times 3 = 540$ experimental runs, leading to a total of $540 \times 30 = 9180$ experiments across all prompts.

Energy consumption was measured *via* CodeCarbon [35], a Python library allowing to measure the energy consumed during the generation process of a image in kilowatt-hours (kWh) through the formula $E = E_{CPU} + E_{GPU} + E_{RAM}$, i.e., the sum of CPU, GPU, and RAM energy consumption.

All experiments were performed on a dedicated workstation with an AMD Ryzen 9 7950X processor, 128 GB of RAM, and an NVIDIA GeForce RTX 4090.

3.3. Image Quality Assessment

Since no definitive reference set exists for text-to-image diffusion models, a common approach [6, 17] is to use the validation set of the dataset from which class labels were used as prompts. In our setting, instead of using the CIFAR-100 validation images, which have low resolution (32×32), we employed the ImageNet-1k validation set. However, since the class labels in ImageNet-1k do not fully overlap with our selected prompt classes, we applied a sampling strategy to guarantee that the same subjects are uniformly represented. Specifically, we utilized our prompt classes as queries to identify corresponding classes in ImageNet-1k. First, we selected all ImageNet-1k classes which had at least one match with a prompt. If a unique match was found, we retrieved all 50 images for that class. Otherwise, where multiple matches existed (e.g., prompt “clock” matched multiple classes such as “wall clock”, “analog clock” and “digital clock”), we performed uniform random sampling across the matched classes to maintain balance and diversity in the final dataset.

Via our filtering strategy we obtained a final dataset for quality assessment consisting of 1,100 images from the ImageNet-1k validation set and 396 images

for each diffusion model considered. FID computation was performed using IQA-PyTorch library [36], while Precision and Recall were computed using the implementation provided by [37].

4. Results

In this section, we present the results of our experiment by considering the various independent variables studied, i.e., models, quantization, image resolution, and prompt length. As concluding remark, we additionally explore the interplay between energy consumption and image quality.

4.1. Median Model Energy Consumption

An overview of the overall energy consumption of the models (i.e., the median energy consumption of each model across all independent variables) is depicted in Figure 1. As we can observe from Figure 1, the considered models display a notable difference in terms of median energy consumption, with Lumina being the models consuming the highest energy value (4.08×10^{-3} kWh) and LCM.SSD_1B the lowest one (8.6×10^{-5} kWh). While the difference of median energy consumption is less drastic for other models (e.g., PixArt_Alpha and PixArt_Sigma differ only for a negligible 4×10^{-5} factor), we conclude that choosing a model over another one can lead to drastic energy savings up to 46x. Overall, we can observe that generally U-Net-based models tend to consume less than Transformer-based ones, reflecting the greater computational complexity inherent to transformer architectures. As additional observation, we note that the variability in terms of energy consumption differs between models, with some displaying only a limited energy consumption fluctuation across settings (e.g., the Pixart-based models), while others a much higher one (e.g., Lumina, SD_1.5, and Flash_SD).

☞ Median Model Energy Consumption: Energy consumption varies drastically across models, showcasing up to a 46x increment. U-Net-based models tend to consume less than Transformer-based ones, while consumption fluctuations vary from model to model.

4.2. Model Quantization

Figure 2 illustrates, on a logarithmic scale, the energy consumption (in kWh) of the analyzed models across three image resolutions, comparing their performance when quantization using the *int8* data type and *no quantization* is applied. The energy consumed values for each model are the medians across the

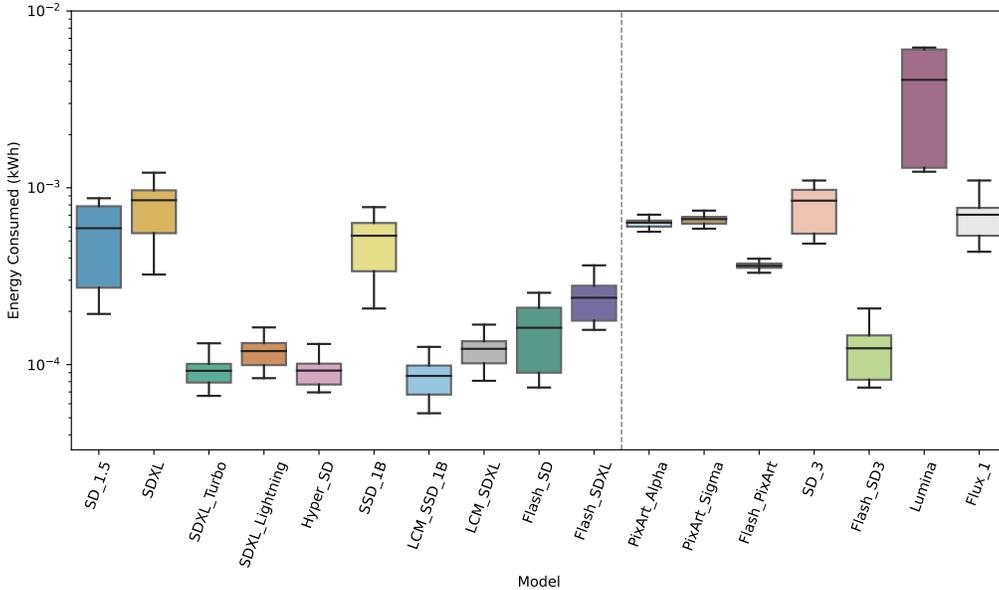


Figure 1: Boxplot of energy consumption across the investigated models. The dotted grey lines separate U-Net (left) from Transformer-based models (right).

30 prompts. The results indicate that most models exhibit increased energy consumption when quantized to the int8 data type, with the exception of four models—SDXL_Turbo, Hyper_SD, Flash_SD3, and Flux_1—which show a modest reduction in energy usage (up to 12.13%). This observation suggests that quantization primarily optimizes memory usage without optimizing the energy efficiency of the image generation process. Furthermore, as we can see from Table 1, the percentage difference in energy consumption between the quantized and non-quantized versions is generally negligible, except for few models, e.g., SD_1.5, SDXL, and SSD_1B.

🍃 Model Quantization: Counterintuitively, model quantization leads in the vast majority of cases to an energy consumption increase (up to 64.54%) and only in rare cases to a negligible energy saving (up to 12%).

4.3. Image Resolution

Figure 3 shows the impact of image resolution on the energy consumption. Without any surprise, most of the investigated models exhibit the expected trend of increased energy consumption as image resolution increases.

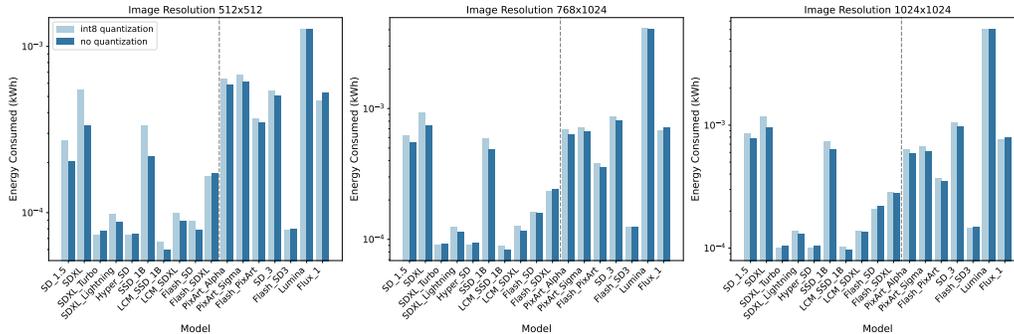


Figure 2: Energy consumption (kWh) for image generation process of diffusion models at varying quantization levels. Bar charts (log scale) compare int8-quantized vs. non-quantized models at fixed resolutions. The dotted grey lines separate U-Net (left) from Transformer-based models (right). Each bar shows the median over 30 prompts.

Table 2 highlights the varying energy consumption scaling of diffusion models as image resolution increases fourfold, from 512×512 to 1024×1024 . Energy consumption variation is calculated as the ratio of the energy used to generate a 1024×1024 resolution image to the energy used for 512×512 resolution image. Regarding U-Net-based models, some of them exhibit a relatively low scaling ratio (around 1.3), such as SDXL_Turbo and Hyper_SD, indicating a more efficient handling of higher resolutions. In contrast, models like SD_1.5 (3.85) and SSD_1B (2.90) show a much higher ratio, suggesting that their energy consumption is almost directly proportional to image resolution. For Transformer-based models, the PixArt family maintains a ratio close to 1.00, indicating that their energy consumption remains constant across resolutions. Other models, such as SD_3 and Flash_SD3, have slightly higher ratios (~ 1.9), implying a modest consumption increase. Lumina stands out with the highest ratio (4.75), exceeding even the expected quadratic scaling, which may suggest additional computational overhead or inefficiencies when processing high-resolution images. As peculiar result, the PixArt-based models showcased a higher energy consumption for non-square formats. Further experiments conducted by considering other 11 formats between the range of 256×256 to 1536×1536 confirmed this trend. We conjecture this trend is due to the training data used for the model, which could have been formatted only in square resolution. Additional details on the follow-up experiments are available in the supplementary material of this work.

Table 1: Percentage variation of energy consumption of diffusion models at varying quantization levels. Variation is computed as the difference between non-quantized and quantized model for each resolution. Negative values indicate lower energy consumption in non-quantized models.

Type	Model	Image Resolution		
		512×512	768×1024	1024×1024
U-Net-based	SD_1.5	-33.33%	-13.67%	-8.46%
	SDXL	-64.58%	-26.61%	-22.34%
	SDXL_Turbo	5.32%	2.13%	3.51%
	SDXL_Lightning	-11.80%	-8.08%	-6.09%
	Hyper_SD	1.49%	3.60%	2.59%
	SSD_1B	-54.42%	-22.22%	-17.84%
	LCM_SSD_1B	-11.42%	-6.43%	-5.44%
	LCM_SDXL	-12.74%	-7.85%	-3.47%
	Flash_SD	-13.38%	0.96%	5.30%
	Flash_SDXL	3.77%	2.66%	-1.73%
Transformer-based	PixArt_Alpha	-8.30%	-8.76%	-7.95%
	PixArt_Sigma	-8.88%	-7.11%	-9.75%
	Flash_PixArt	-5.25%	-6.40%	-6.02%
	SD_3	-8.35%	-7.13%	-7.64%
	Flash_SD3	1.05%	0.67%	0.75%
	Lumina	-0.15%	-0.16%	-0.20%
	Flux_1	12.13%	5.97%	2.88%

Image resolution: Image resolution affects differently the energy consumption of the models, with some showcasing low scaling ratio (e.g., 1.3x when image size is doubled) while other a considerable increase (e.g., 4.7x). Pixart family models remain close to constant across all square resolutions, albeit consuming considerably more for non-square ones.

4.4. Prompt Length

From a Quantile-Quantile plot analysis the energy consumption of each model does not result to follow a normal distribution. Therefore, we study the correlation between energy consumption and prompt length *via* the Kruskal-Wallis non-parametric test. Results showcase that for all models the p-values are notably above the conventional significance level of 0.05, implying that prompt length is not a determining factor in the energy consumption of the examined models. Further details regarding this experiment are provided in the supplementary material of this work.

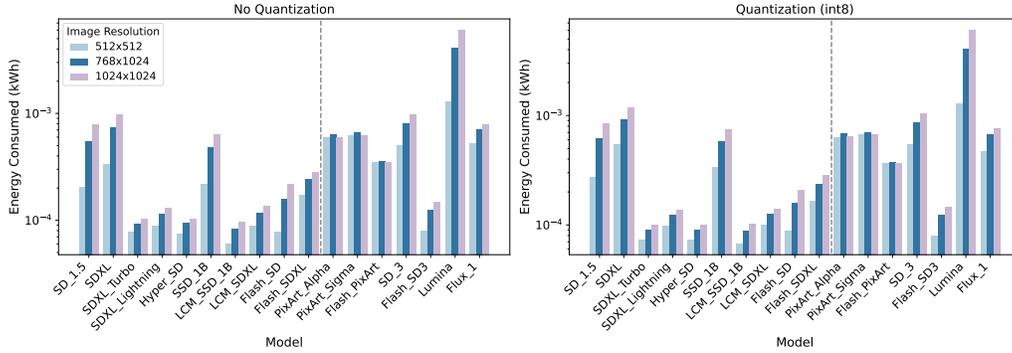


Figure 3: Energy consumption (kWh) for image generation process of diffusion models at varying image resolution. Bar charts (log scale) compare 512×512 , 768×1024 , and 1024×1024 image resolutions at fixed quantization setting. The dotted grey lines separate U-net (left) from Transformer-based models (right). Each bar shows the median over 30 prompts.

🍀 Prompt Length: Prompt length does not impact in a statistically significant manner the energy consumed by image generation.

4.5. Image Quality Assessment

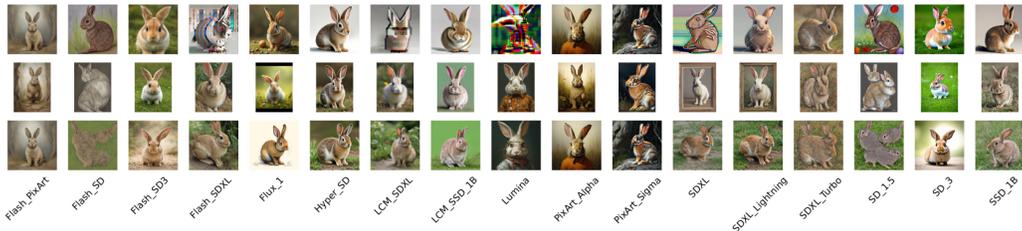


Figure 4: Example images generated by each model at different resolutions. Each row corresponds to a different resolution: 512×512 (top), 768×1024 (middle), and 1024×1024 (bottom). The images were generated using the short prompt “rabbit” with int8 quantization.

Figure 5 documents the performance of each diffusion model in terms of FID and energy consumption. Low values of FID correspond to better quality. Notably, Lumina emerges as the least efficient model, both from the quality and the energy perspective. On the other hand, Flash_SD3 offers the best compromise, achieving the lowest FID among all models while maintaining low energy consumption. As illustrated in Figure 4, where the same subject is represented at dif-

Table 2: Scaling ratios of energy consumption for image generation at 1024×1024 resolution compared to 512×512 resolution across different diffusion models.

Type	Model	No Quantization	Quantization (int8)
U-Net-based	SD_1.5	3.85	3.13
	SDXL	2.88	2.14
	SDXL_Turbo	1.33	1.35
	SDXL_Lightning	1.48	1.41
	Hyper_SD	1.38	1.37
	SSD_1B	2.90	2.21
	LCM_SSD_1B	1.62	1.54
	LCM_SDXL	1.51	1.39
	Flash_SD	2.80	2.35
	Flash_SDXL	1.61	1.70
Transformer-based	PixArt_Alpha	1.00	1.00
	PixArt_Sigma	0.99	1.00
	Flash_PixArt	1.00	1.01
	SD_3	1.92	1.91
	Flash_SD3	1.85	1.86
	Lumina	4.75	4.75
	Flux_1	1.50	1.63

ferent resolutions, the final output of a diffusion model can vary significantly depending on the resolution. In particular, Lumina’s high FID score is largely due to the poor quality of its 512×512 images, as the model is fine-tuned on resolutions of 1024×1024 and above, decreasing its ability to generate 512-resolution images effectively. Similar effects are observed for other models, such as Flash_SD, SD_1.5 and SDXL_Turbo, which exhibit intermediate FID scores despite producing semantically altered images at 768×1024 and 1024×1024 resolutions.

FID captures both fidelity (realism) and diversity (variations in the original data) of generative models but struggles to separate their contributions when FID scores are similar. Precision and Recall overcome this limitation by independently measuring fidelity and diversity, providing a more nuanced evaluation of image quality. In Figure 6 each diffusion model is represented as a dot in the Precision-Recall space, with size indicating energy consumption. As we can observe, Lumina appears as the least precise model. Differently, while SD_1.5 and Hyper_SD achieve similar scores for FID in Figure 5, their distinct positions in the Precision-Recall space highlight qualitative differences. SD_1.5 shows lower precision but higher recall with respect to Hyper_SD. This behavior can be seen as a drawback, as high recall often results from generating numerous diverse but unrealistic samples rather than effectively capturing subject heterogeneity [37]. Based on this observation, Hyper_SD, SD_3, and Flash_SD3 stand out as the highest-quality models, combining strong performance with good energy efficiency.

Additional discussions on quality metrics with resolution-specific plots are provided in the supplementary material of this work.

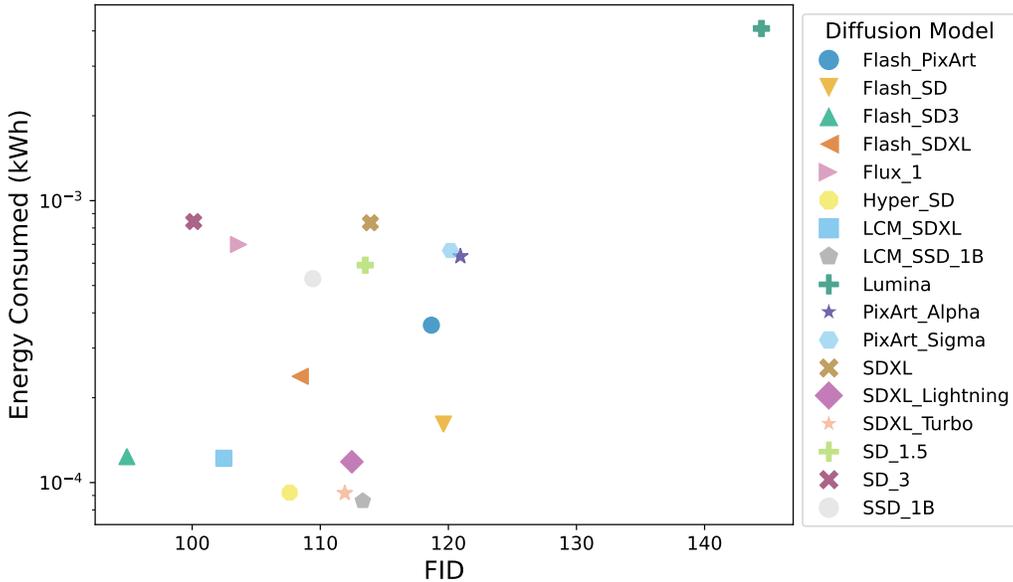


Figure 5: Energy consumption (kWh) vs. FID scores for diffusion models. Y-axis is on a logarithmic scale. Each dot represents the median energy consumption per model.

🍃 Image Quality: Higher image quality does not come at the cost of higher energy consumption, with some of the most energy efficient models also producing among the highest quality images.

5. Conclusions

In this study, we analyze the energy consumption of 17 diffusion models, considering multiple influencing factors at interplay. The measured energy varies significantly across models (up to a 47x increase), with U-Net-based models generally being more efficient than Transformer-based ones. Counterintuitively, model quantization often increases energy consumption (up to 64.54%), with only a few cases showing minor savings. Image resolution affects energy consumption inconsistently, ranging from a 1.3x to 4.7x increase when doubling resolution. Considering the quality-energy connection, Lumina ranks lowest, while Flash SD3 excels in both aspects. In conclusion, with this research, we present a fine-grained

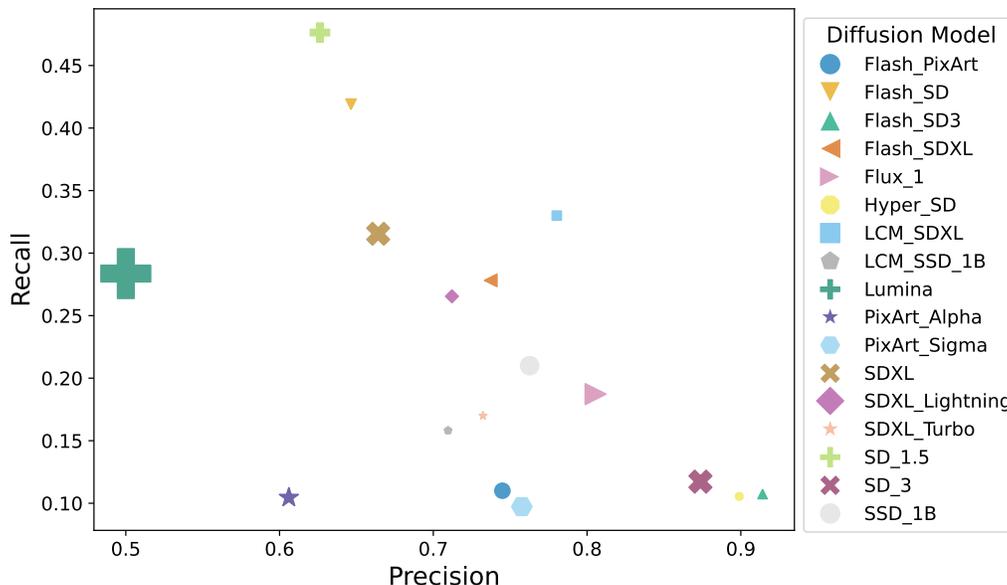


Figure 6: Precision vs. Recall for diffusion models. Shape size is proportional to the median energy consumption per model.

empirical perspective on the energy required to generate images. With our study, we hope to have set a stepping stone to progress image generation not only by producing more eye pleasing images, but also by improving the environmental impact hidden behind each image we generate.

References

- [1] R. Verdecchia, J. Sallou, L. Cruz, A systematic review of Green AI, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 13 (4) (2023) e1507.
- [2] S. N. Gowda, X. Hao, G. Li, S. N. Gowda, X. Jin, L. Sevilla-Lara, Watt for what: Rethinking deep learning’s energy-performance relationship, *arXiv preprint arXiv:2310.06522* (2023).
- [3] J. Castaño, S. Martínez-Fernández, X. Franch, J. Bogner, Exploring the carbon footprint of hugging face’s ml models: A repository mining study, in: *2023 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, IEEE, 2023, pp. 1–12.

- [4] S. A. Budenny, V. D. Lazarev, N. N. Zakharenko, A. N. Korovin, O. Plosskaya, D. V. Dimitrov, V. Akhripkin, I. Pavlov, I. V. Oseledets, I. S. Barsola, et al., Eco2ai: carbon emissions tracking of machine learning models as the first step towards sustainable ai, in: *Doklady Mathematics*, Vol. 106, Springer, 2022, pp. S118–S128.
- [5] A. Razzhigaev, A. Shakhmatov, A. Maltseva, V. Arkhipkin, I. Pavlov, I. Ryabov, A. Kuts, A. Panchenko, A. Kuznetsov, D. Dimitrov, Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion, in: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2023, pp. 286–295.
- [6] J. Chen, J. Yu, C. Ge, L. Yao, E. Xie, Y. Wu, Z. Wang, J. Kwok, P. Luo, H. Lu, et al., Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis, *arXiv preprint arXiv:2310.00426* (2023).
- [7] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, S. Hochreiter, Gans trained by a two time-scale update rule converge to a local nash equilibrium, in: *Neural Information Processing Systems*, 2017.
URL <https://api.semanticscholar.org/CorpusID:326772>
- [8] M. S. M. Sajjadi, O. Bachem, M. Lučić, O. Bousquet, S. Gelly, Assessing Generative Models via Precision and Recall, in: *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [9] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, Improved techniques for training gans, *Advances in neural information processing systems* 29 (2016).
- [10] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, Y. Choi, Clipscore: A reference-free evaluation metric for image captioning, *arXiv preprint arXiv:2104.08718* (2021).
- [11] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [12] J. Chen, C. Ge, E. Xie, Y. Wu, L. Yao, X. Ren, et al., Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation, *arXiv preprint* (2024).

- [13] P. Gao, L. Zhuo, Z. Lin, C. Liu, J. Chen, R. Du, E. Xie, X. Luo, L. Qiu, Y. Zhang, et al., Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers, arXiv preprint arXiv:2405.05945 (2024).
- [14] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 10684–10695.
- [15] A. Sauer, D. Lorenz, A. Blattmann, R. Rombach, Adversarial diffusion distillation, in: European Conference on Computer Vision, Springer, 2025, pp. 87–103.
- [16] P. Esser, S. Kulal, A. Blattmann, R. Entezari, J. Müller, H. Saini, Y. Levi, D. Lorenz, A. Sauer, F. Boesel, et al., Scaling rectified flow transformers for high-resolution image synthesis, in: Forty-first International Conference on Machine Learning, 2024.
- [17] C. Chadebec, O. Tasar, E. Benaroché, B. Aubin, Flash diffusion: Accelerating any conditional diffusion model for few steps image generation (2024). arXiv:2406.02347.
- [18] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, S. Ganguli, Deep unsupervised learning using nonequilibrium thermodynamics, in: F. Bach, D. Blei (Eds.), Proceedings of the 32nd International Conference on Machine Learning, Vol. 37 of Proceedings of Machine Learning Research, PMLR, Lille, France, 2015, pp. 2256–2265.
URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>
- [19] J. Ho, A. Jain, P. Abbeel, Denoising diffusion probabilistic models, Advances in neural information processing systems 33 (2020) 6840–6851.
- [20] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., Learning transferable visual models from natural language supervision, in: International conference on machine learning, PMLR, 2021, pp. 8748–8763.

- [21] W. Peebles, S. Xie, Scalable diffusion models with transformers, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 4195–4205.
- [22] A. Vaswani, Attention is all you need, Advances in Neural Information Processing Systems (2017).
- [23] D. Alexey, An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv: 2010.11929 (2020).
- [24] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, M. Chen, Hierarchical text-conditional image generation with clip latents, arXiv preprint arXiv:2204.06125 1 (2) (2022) 3.
- [25] Midjourney, <https://www.midjourney.com/> (2022).
- [26] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, R. Rombach, Sdxl: Improving latent diffusion models for high-resolution image synthesis, arXiv preprint arXiv:2307.01952 (2023).
- [27] S. Lin, A. Wang, X. Yang, Sdxl-lightning: Progressive adversarial diffusion distillation, arXiv preprint arXiv:2402.13929 (2024).
- [28] Y. Ren, X. Xia, Y. Lu, J. Zhang, J. Wu, P. Xie, X. Wang, X. Xiao, Hyper-sd: Trajectory segmented consistency model for efficient image synthesis, arXiv preprint arXiv:2404.13686 (2024).
- [29] Y. Gupta, V. V. Jaddipal, H. Prabhala, S. Paul, P. V. Platen, Progressive knowledge distillation of stable diffusion xl using layer level loss (2024). arXiv:2401.02677.
- [30] S. Luo, Y. Tan, L. Huang, J. Li, H. Zhao, Latent consistency models: Synthesizing high-resolution images with few-step inference, arXiv preprint arXiv:2310.04378 (2023).
- [31] B. F. Labs, Flux, <https://github.com/black-forest-labs/flux> (2023).
- [32] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny images (2009).

- [33] M. Abdin, J. Aneja, H. Awadalla, A. Awadallah, A. A. Awan, N. Bach, A. Bahree, A. Bakhtiari, J. Bao, H. Behl, et al., Phi-3 technical report: A highly capable language model locally on your phone, arXiv preprint arXiv:2404.14219 (2024).
- [34] D. Wind, Top 40 useful prompts for stable diffusion xl, <https://medium.com/phygital/top-40-useful-prompts-for-stable-diffusion-xl-008c03dd0557>, accessed: 2025-03-07 (Dec. 2023).
- [35] B. Courty, V. Schmidt, S. Luccioni, Goyal-Kamal, MarionCoutarel, B. Feld, J. Lecourt, LiamConnell, A. Saboni, Inimaz, supatomic, M. Léval, L. Blanche, A. Cruveiller, ouminasara, F. Zhao, A. Joshi, A. Bogroff, H. de Lavoreille, N. Laskaris, E. Abati, D. Blank, Z. Wang, A. Catovic, M. Alencon, M. Stechly, C. Bauer, L. O. N. de Araújo, JPW, MinervaBooks, mlco2/codecarbon: v2.4.1 (May 2024). doi:10.5281/zenodo.11171501.
URL <https://doi.org/10.5281/zenodo.11171501>
- [36] C. Chen, J. Mo, IQA-PyTorch: Pytorch toolbox for image quality assessment, [Online]. Available: <https://github.com/chaofengc/IQA-PyTorch> (2022).
- [37] M. F. Naeem, S. J. Oh, Y. Uh, Y. Choi, J. Yoo, Reliable fidelity and diversity metrics for generative models (2020).

Appendix A. Preliminary Experiment

The results of the correlation between the semantic content of prompts and the energy consumed during the image generation process are presented in Table A.3. Before performing the correlation analysis, we verified that the energy consumption for each model followed a normal distribution, by examining Quantile-Quantile (QQ) plots, which are illustrated in Figure A.7. If the data closely follows the normal distribution, the points on the Q-Q plot will move on a diagonal line (depicted in red), while deviations from the reference line indicate departures from the expected distribution. Figure A.7 shows that for all the analyzed models, data follows a normal distribution. This step ensured the validity of applying Pearson’s correlation coefficient, which relies on the assumption of normally distributed data.

From Table A.3, we can observe that the Pearson correlation coefficients (PCC) for all models ranged from -0.042 to 0.036, with all the corresponding p-values exceeding the conventional significance level of 0.05. These findings indicate that there is no statistically significant relationship between the prompt content and the energy consumed during the generation process, supporting the validity of the followed experimental setup.

Additionally, Figures A.8 to A.10 show the box plots which illustrate the distribution of energy consumption across different models for the 100 different prompts. As we can see, the median energy remains quite stable across all semantics, with negligible variations among prompts. This observation supports the finding from Pearson correlation test that prompt semantic does not impact on the energy consumption.

Appendix B. Image Resolution

Figure B.11 illustrates the additional experiments conducted on PixArt_Alpha, PixArt_Sigma, and Flash_PixArt models to further investigate the impact of image resolution on energy consumption. We can notice that for all models the energy consumption for non-square resolution is much higher compared to square resolutions.

Appendix C. Prompt Length

Before performing the correlation analysis between the length of prompts and the energy consumed during the image generation process, we verified if the en-

Table A.3: Pearson correlation coefficients (PCC) and corresponding p-values assessing the relationship between the semantic content of prompt and energy consumption during image generation process across different models.

Model	PCC	P-value
Flash_PixArt	0.033	0.329
Flash_SD	-0.035	0.296
Flash_SD3	0.035	0.307
Flash_SDXL	-0.038	0.263
Flux_1	-0.024	0.485
Hyper_SD	-0.007	0.841
LCM_SDXL	-0.029	0.392
LCM_SSD_1B	-0.009	0.795
Lumina	0.024	0.481
PixArt_Alpha	0.021	0.532
PixArt_Sigma	0.036	0.294
SD_1.5	-0.042	0.213
SD_3	-0.030	0.378
SDXL	0.027	0.426
SDXL_Lightning	-0.030	0.371
SDXL_Turbo	-0.011	0.751
SSD_1B	-0.003	0.934

ergy consumption for each model followed or not a normal distribution, by examining Quantile-Quantile (QQ) plots, which are illustrated in Figure C.12. If the data closely follows the normal distribution, the points on the Q-Q plot will move on a diagonal line (depicted in red), while deviations from the reference line indicate departures from the expected distribution. Figure C.12 shows that for all the analyzed models, data does not follow a normal distribution, as the points deviate significantly from the straight diagonal line. This step ensured the validity of applying Kruskal-Wallis test, which is a non-parametric statistical test, which does not assume a normally distributed data.

Moreover, Figure C.13 show the box plots which represent the distribution of energy consumption across different models for the three different prompt lengths. As we can notice, the median energy consumption remains stable across different prompt lengths for all models, with no clear increasing or decreasing trend. This observation suggests that prompt length does not have a significant influence on energy consumption, as no evident pattern is observed across different models. This finding is also supported by the Kruskal-Wallis statistical test.

Table C.4: Rank sum and p-values obtained from Kruskal-Wallis test to assess the relationship between the length of prompt and energy consumption during image generation process across different models.

Model	Rank Sum	P-value
Flash_PixArt	0.77	0.68
Flash_SD	3.66	0.16
Flash_SD3	0.08	0.96
Flash_SDXL	0.95	0.62
Flux_1	2.92	0.23
Hyper_SD	0.06	0.97
LCM_SDXL	0.04	0.98
LCM_SSD_1B	0.25	0.88
Lumina	0.26	0.88
PixArt_Alpha	0.75	0.69
PixArt_Sigma	0.51	0.77
SD_1.5	0.56	0.76
SD_3	0.25	0.88
SDXL	0.14	0.93
SDXL_Lightning	0.92	0.63
SDXL_Turbo	2.99	0.23
SSD_1B	0.20	0.91

Table C.4 presents the results of the Kruskal-Wallis test, which examines whether variations in prompt lengths lead to statistically significant differences in energy consumption.

Appendix D. Quality Assessment

As previously stated in Section 3.3, the validation set of ImageNet-1k was properly filtered to ensure that only the classes matching with the prompt labels were retained, with random uniform sampling in case of multiple matches. Unfortunately, for prompts *apples*, *girl*, *poppies*, *aquarium fish*, *oak*, *shrew*, *ray* and *tulips* were not possible to find a matching ImageNet-1k class. Consequently, we excluded all images associated with the aforementioned prompts from the calculation of the quality metrics.

Figure D.14 and Figure D.15 illustrate how the quality-energy relationship varies when considering images at specific resolutions separately. Flash_SD3 maintains consistent FID and Precision-Recall performance across different resolutions, whereas many other models exhibit significant variations. Among them, Flash_SD shoes low FID combined with high precision at 512×512 but shifts

to higher FID, lower precision and higher recall at 768×1024 and 1024×1024 . These trends are visually confirmed by the generated images in Figure D.16, Figure D.17 and Figure D.18.

Figure D.18 presents the same subject (prompt number 0) across different diffusion models at 1024×1024 resolution with int8 quantization, varying only the prompt length. As discussed in Appendix Appendix C, prompt length generally does not impact the energy consumption of the generation process. From a quality perspective, short-prompt images cannot be definitively classified as lower quality. However, it is interesting to note that the most significant differences in image content arise between short and long prompts. This aligns with the goal of prompt engineering for text-to-image generation, which is to produce more detailed and complex images. Other examples are shown in Figure D.16 and Figure D.17.

Figure D.19 presents the average CLIPScore values for all examined models across varying prompt lengths. The results indicate that shorter prompts consistently achieve the highest scores for every model, whereas longer prompts yield the lowest. This pattern suggests that when prompts contain only limited information, the generated images are less likely to deviate from the requested content, resulting in looser alignment with the given text. In contrast, longer prompts introduce more details, which the generation process may struggle to accurately capture, leading to greater misalignment, and consequently lower CLIPScore values.

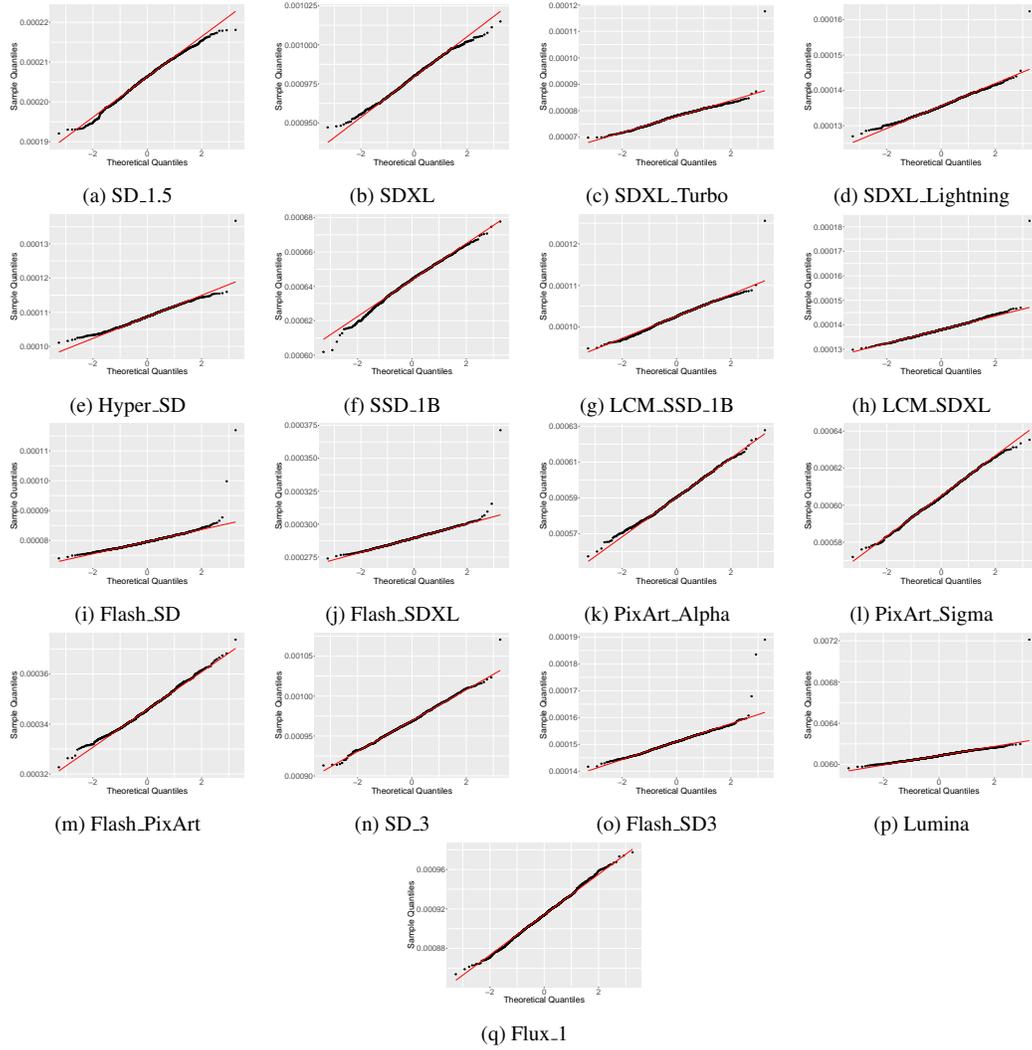
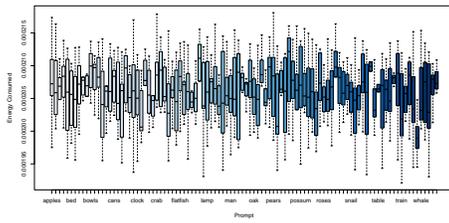
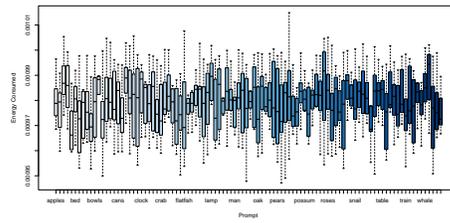


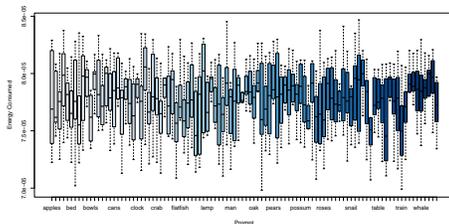
Figure A.7: Q-Q plots comparing the energy consumption of the analyzed models across various prompt contents to the normal distribution. Each plot visualizes the distribution of energy consumption values for a specific model, assessing how closely they follow a normal distribution.



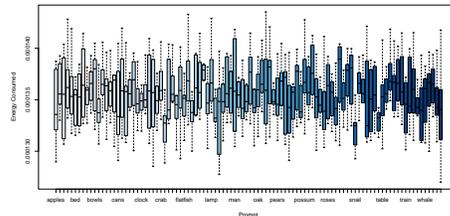
(a) SD_1.5



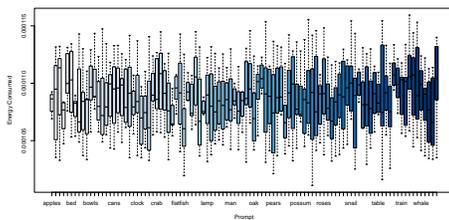
(b) SDXL



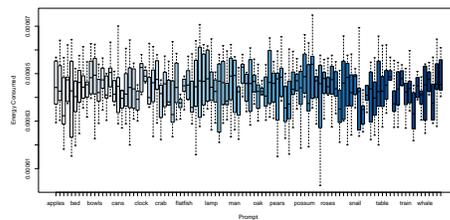
(c) SDXL_Turbo



(d) SDXL_Lightning

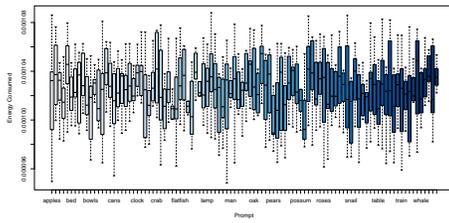


(e) Hyper_SD

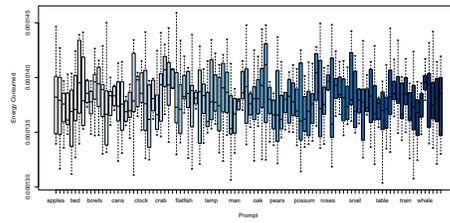


(f) SSD_1B

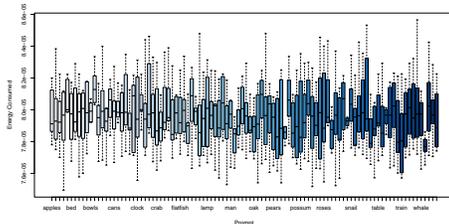
Figure A.8: Box plots of the energy consumption across various prompt lengths for the analyzed models. (Group 1).



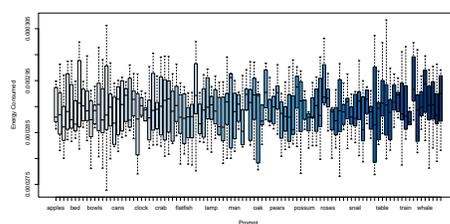
(a) LCM_SSD_1B



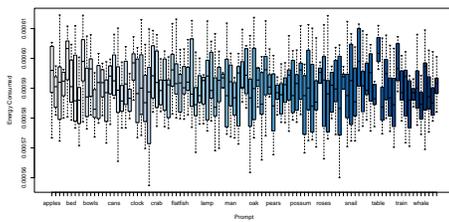
(b) LCM_SDXL



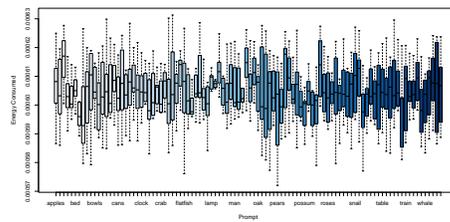
(c) Flash_SD



(d) Flash_SDXL

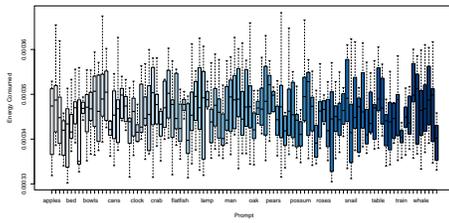


(e) PixArt_Alpha

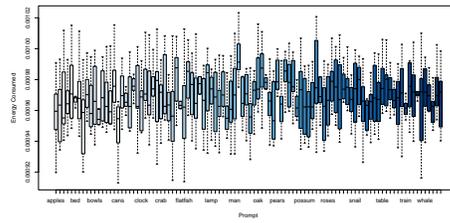


(f) PixArt_Sigma

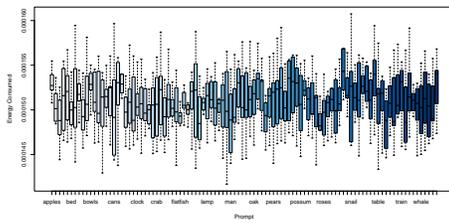
Figure A.9: Box plots of the energy consumption across various prompt lengths for the analyzed models. (Group 2).



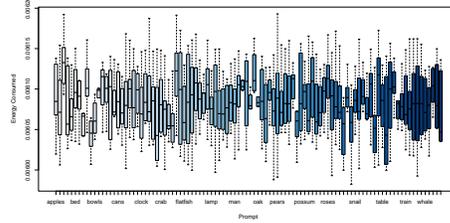
(a) Flash_PixArt



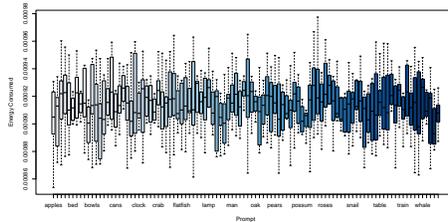
(b) SD_3



(c) Flash_SD3



(d) Lumina



(e) Flux_1

Figure A.10: Box plots of the energy consumption across various prompt lengths for the analyzed models. (Group 3).

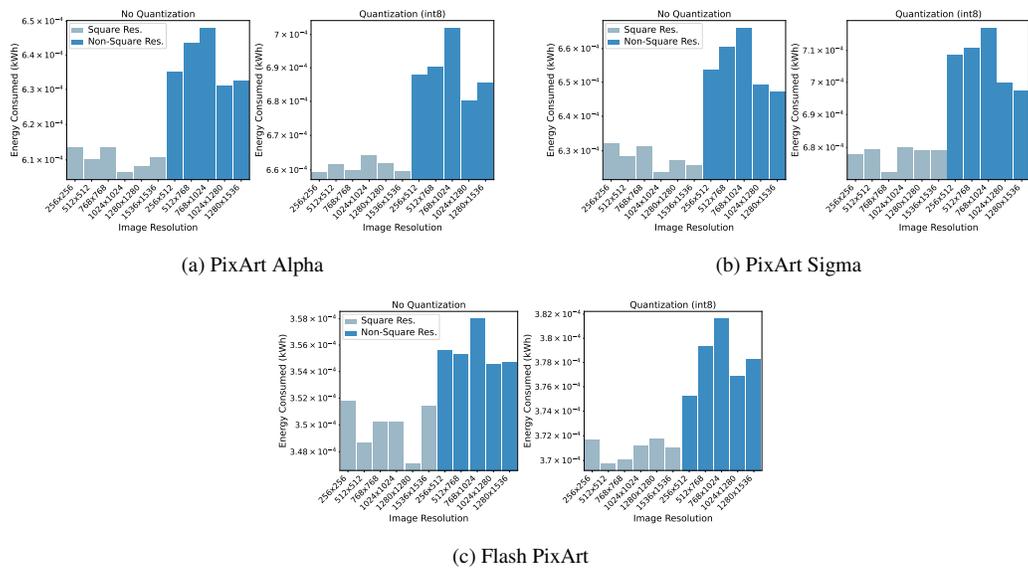


Figure B.11: Energy consumption (kWh) of PixArt_Alpha, PixArt_Sigma, and Flash_PixArt models at varying image resolution. Bar charts (log scale) show the energy consumed across different image resolutions, categorized as square and non-square resolutions, at fixed quantization setting. Each bar shows the median over 30 prompts.

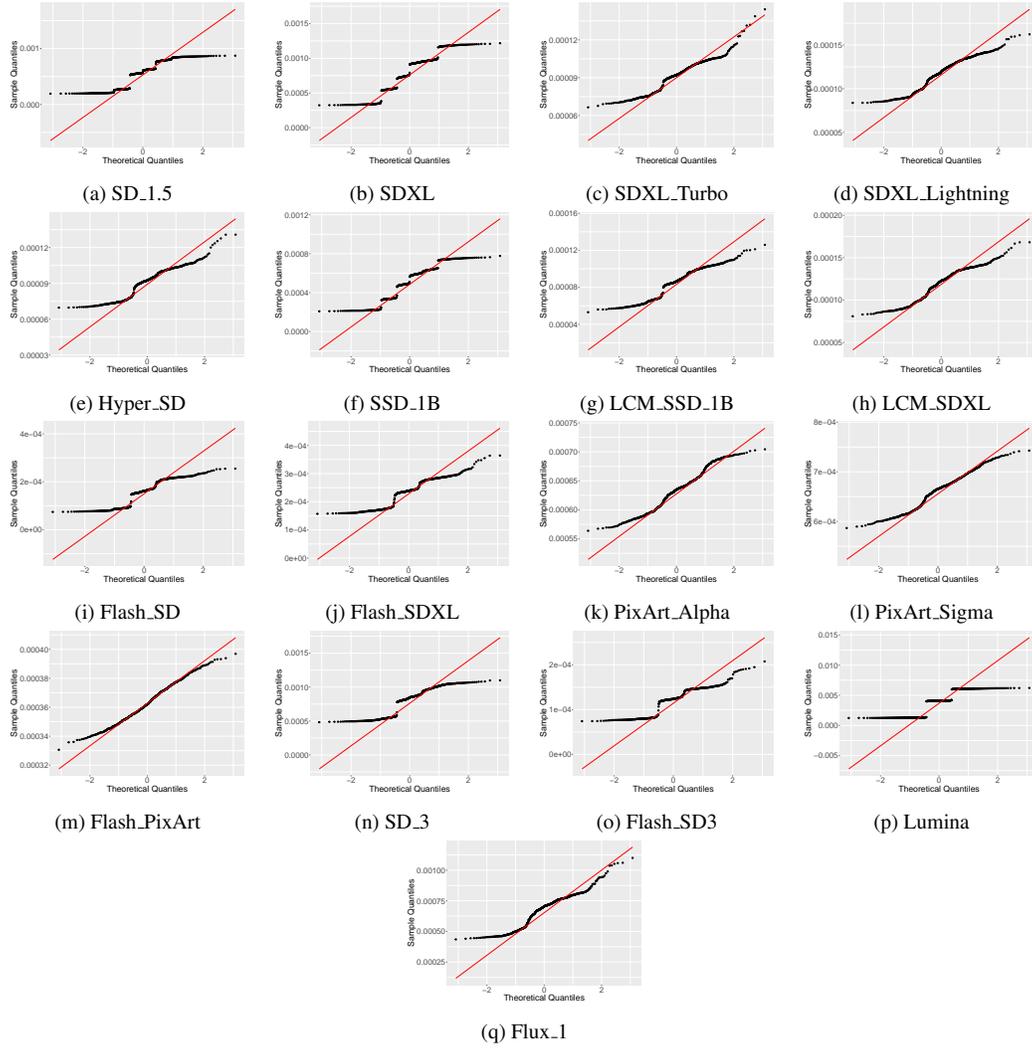


Figure C.12: Q-Q plots comparing the energy consumption of the analyzed models across various prompt length to the normal distribution. Each plot visualizes the distribution of energy consumption values for a specific model, assessing how closely they follow a normal distribution.

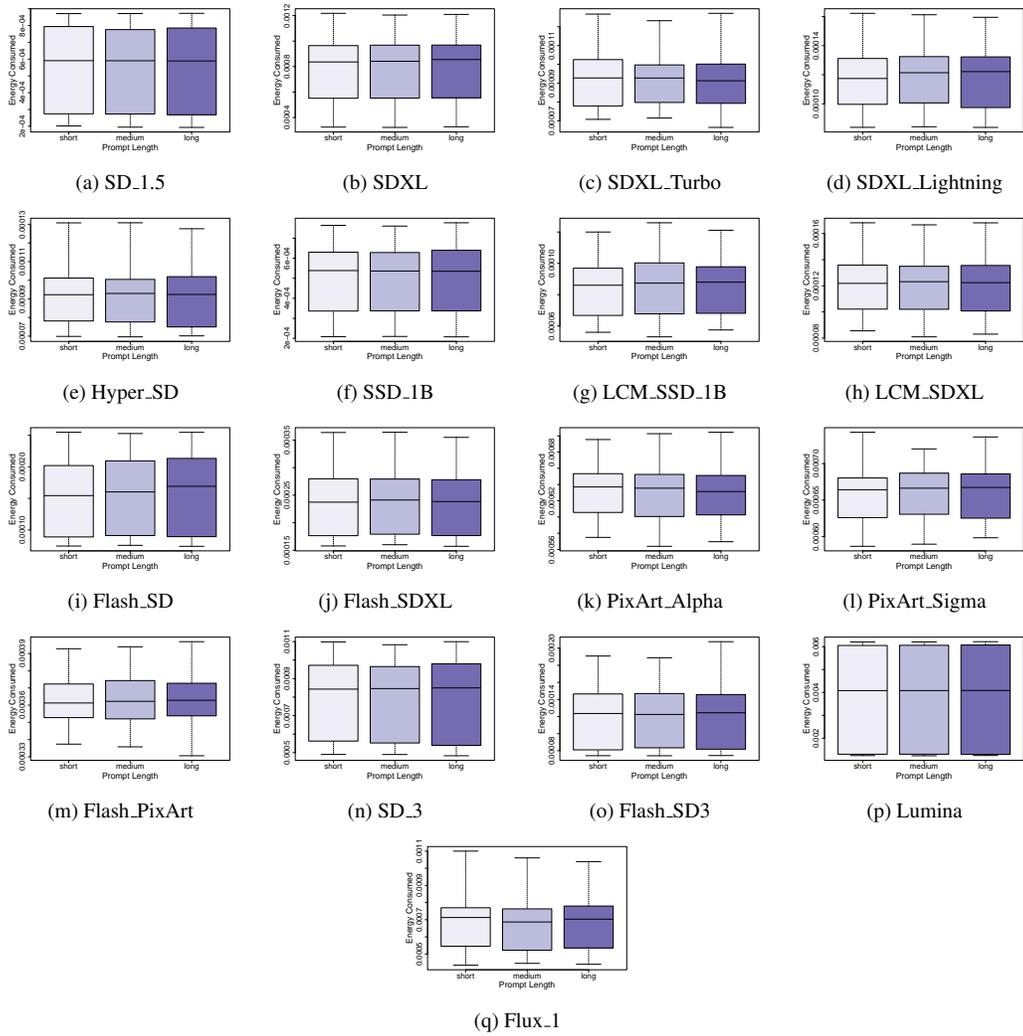


Figure C.13: Box plots of the energy consumption across various prompt length for the analyzed models.

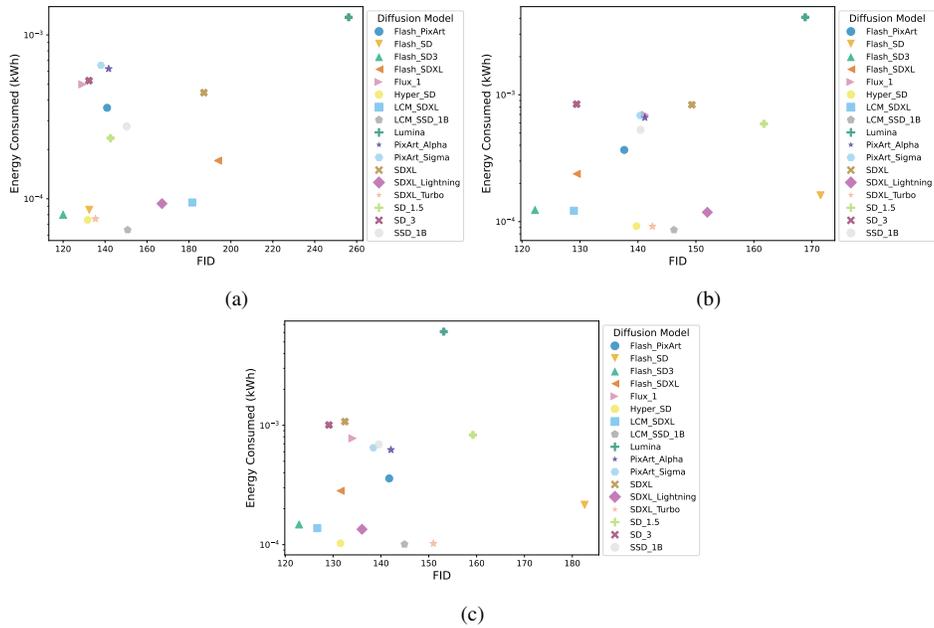


Figure D.14: Energy consumption (kWh) vs FID values for all the models at resolution (a) 512×512 (b) 768×1024 and (c) 1024×1024 .

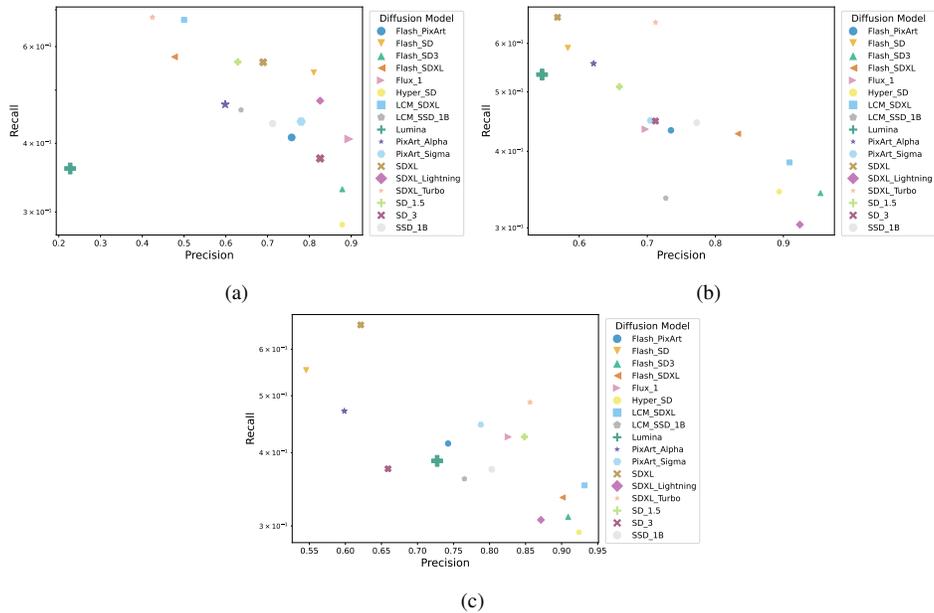


Figure D.15: Precision and Recall values for each model at resolution (a) 512×512 (b) 768×1024 and (c) 1024×1024 . Each symbol of the plot has different size based on the energy consumed.

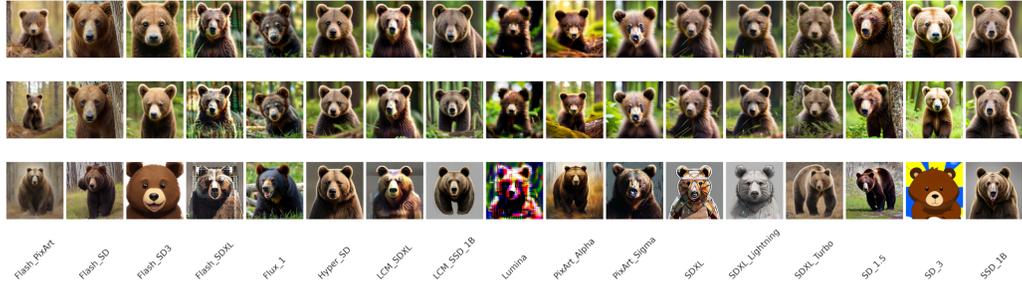


Figure D.16: Example images generated by each model at resolution 512×512 , quantization int8 and prompt “bear”. Each row corresponds to a different prompt length: long (top), medium (middle), and short (bottom).



Figure D.17: Example images generated by each model at resolution 768×1024 , quantization int8 and prompt “tree”. Each row corresponds to a different prompt length: long (top), medium (middle), and short (bottom).

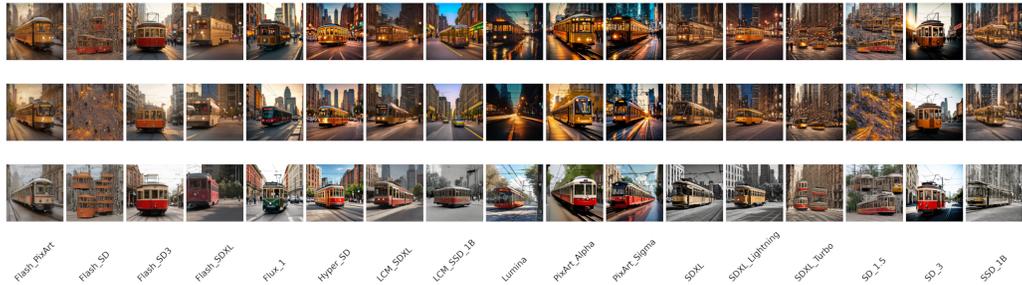


Figure D.18: Example images generated by each model at resolution 1024×1024 , quantization int8 and prompt “streetcar”. Each row corresponds to a different prompt length: long (top), medium (middle), and short (bottom).

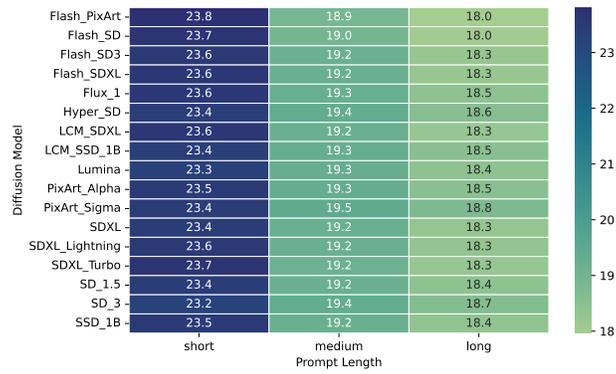


Figure D.19: Average CLIPScore values for each diffusion model with respect to different prompt lengths.