# A Comparative Analysis of Text-to-Image Model Architectures

The field of artificial intelligence has witnessed a remarkable transformation in recent years, with text-to-image models emerging as a particularly captivating area of innovation. These sophisticated machine learning systems possess the ability to translate natural language descriptions into visually compelling imagery, opening up unprecedented possibilities for creative expression, content generation, and technological advancement across various industries. From generating novel artistic concepts to streamlining design workflows, the impact of these models is rapidly expanding, making a thorough understanding of their underlying mechanisms crucial for developers, researchers, and end-users alike. This report aims to provide a comprehensive analysis of the architectures of several leading text-to-image models, including the recently introduced FLUX model, while considering the latest research and advancements in the field to illuminate the factors that contribute to their performance and to explore the complex question of what constitutes the "best" architecture. The analysis will delve into the technical intricacies of models such as FLUX, the Stable Diffusion family (including SDXL and SD3), the DALL-E family (with a focus on DALL-E 3), the Imagen family (particularly Imagen 3), as well as Midjourney, Recraft V3, Luma Photon, Ideogram, Amazon Titan, and Google's Gemini, highlighting their unique architectural characteristics and capabilities.

## Deconstructing the FLUX Model Architecture

FLUX, developed by Black Forest Labs and introduced in August 2024 [1], represents a notable addition to the landscape of high-performing text-to-image models [2]. Its architecture employs a sophisticated hybrid approach, strategically blending the strengths of multimodal diffusion and transformer blocks [4]. Multimodal diffusion in this context refers to the model's ability to process and integrate information from both textual and visual modalities during the image generation process, allowing for a more nuanced understanding of the relationship between the input prompt and the desired output image. Transformer blocks, a fundamental component of modern natural language processing and increasingly in computer vision, play a crucial role in processing sequential data, such as the text prompt and the latent representations of the image as it is being generated [4]. This hybrid design allows FLUX to effectively capture both the global context provided by the text and the fine-grained details required for high-quality image synthesis.

Beyond its core hybrid architecture, FLUX incorporates several key technical innovations that contribute to its performance [4]. One such innovation is the integration

of flow matching, a relatively recent technique in the field of generative modeling that simplifies the diffusion process [4]. Unlike traditional diffusion models that involve complex iterative denoising steps, flow matching aims to directly learn a vector field that maps the initial noise distribution to the desired data distribution, potentially leading to more efficient training and the generation of more coherent images [6]. The use of rotary positional embeddings is another distinguishing feature of FLUX's architecture [4]. These embeddings are specifically designed to encode information about the spatial relationships between different parts of the image, allowing the model to better understand and represent the arrangement of objects and details within the generated scene, which is crucial for accurate prompt adherence. Furthermore, FLUX utilizes parallel attention layers [4]. Attention mechanisms, a key aspect of transformer networks, enable the model to focus on the most relevant parts of the input data when generating the output. By processing these attention calculations in parallel, FLUX can potentially achieve faster generation times without compromising the quality or accuracy of the generated images. This efficient handling of complex spatial relationships and improved alignment between textual descriptions and visual outputs allows FLUX to excel in areas like image fidelity and prompt adherence.

Recognizing the diverse needs of users, FLUX is available in three distinct variants: Schnell, Dev, and Pro [3]. The Schnell variant is specifically optimized for speed, making it well-suited for applications where rapid image generation is paramount, even if it involves a slight trade-off in ultimate image quality [3]. The Dev variant stands out as an open-weight version, intended primarily for non-commercial use, such as research and personal development, and it offers the flexibility for users to fine-tune the model for specific tasks [3]. Finally, the Pro version represents the top-tier offering, providing the highest level of image detail, output diversity, and accuracy in following complex prompts, making it ideal for commercial projects that demand the utmost performance [3]. Access to the Pro version is typically provided through an API, allowing for seamless integration into various applications and workflows. The existence of these tailored variants underscores the practical considerations in the design of text-to-image models, where a balance must often be struck between factors such as speed, quality, accessibility, and the specific requirements of the intended use case.

## Exploring the Landscape of Text-to-Image Model Architectures

Beyond the intricacies of the FLUX model, a diverse range of architectures underpins the current generation of text-to-image models, each with its own set of strengths and limitations.

## A. Stable Diffusion Family (SDXL, SD3)

The Stable Diffusion family, developed by Stability AI, has been a significant force in the text-to-image generation landscape, known for its open-source nature and adaptability [8]. The evolution of this family has seen notable architectural advancements.

**SDXL:** Stable Diffusion XL represents a substantial upgrade, primarily characterized by its significantly larger UNet backbone, which boasts approximately three times more model parameters compared to its predecessors [9]. This increase in model capacity allows SDXL to potentially capture more complex relationships between the text prompt and the image features, leading to higher quality outputs. A key architectural innovation in SDXL is the employment of two distinct text encoders: OpenCLIP ViT-bigG and CLIP ViT-L [9]. The outputs from these two encoders are concatenated, effectively providing the model with two different perspectives on the input text prompt. This dual encoding strategy enables SDXL to capture a broader range of semantic information and potentially achieve more nuanced and accurate image generation, as different encoders may excel at understanding different aspects of the text [11]. Furthermore, SDXL incorporates a technique to condition the model on the original image resolution during training [10]. This is achieved by using Fourier feature encoding to represent the original height and width of the training images, which are then provided as additional parameters to the model. This approach addresses the limitations of earlier models in handling images of varying sizes and allows users to explicitly control the output resolution during inference. To further refine the image generation process, SDXL introduces a novel method for conditioning on image crops during training [10]. By randomly sampling crop coordinates during data loading and using these as conditioning parameters, the model learns to better handle objects that might have been cropped in the training data, mitigating issues with the generation of incomplete or poorly framed objects. Moreover, SDXL utilizes multi-aspect training as a fine-tuning stage [10]. This involves training the model on images with various aspect ratios, allowing it to generate images that are well-suited for different display formats without requiring specific adjustments during prompting. Finally, SDXL often incorporates a refinement model that operates on the latent vectors generated by the base model [10]. This specialized model focuses on enhancing the local quality and finer details of the images, resulting in a two-stage generation process that prioritizes both overall coherence and intricate visual accuracy.

**SD3:** The latest iteration in the family, Stable Diffusion 3, marks a significant architectural shift by leveraging a Multimodal Diffusion Transformer (MMDiT) architecture [6]. This departure from the U-Net architecture that characterized earlier

Stable Diffusion models reflects a growing trend in the field towards adopting transformer-based architectures for their ability to handle long-range dependencies and potentially offer improved scalability and performance [12]. Within the MMDiT framework, SD3 employs separate sets of weights for processing image and language representations [6]. This separation allows the model to develop specialized processing pathways for each modality, potentially leading to a better understanding of the text prompt and improved spelling and text rendering capabilities in the generated images. Building upon the foundational Diffusion Transformer (DiT) architecture, SD3 uses embeddings of the timestep (indicating the stage of the diffusion process) and the text conditioning vector as inputs to a modulation mechanism [6]. This enables conditional image generation, where the characteristics of the generated image are directly influenced by the provided text prompt and the current stage of the denoising process. The model constructs a sequence that includes embeddings of both the text and the image inputs, along with positional encodings [12]. This sequence-based representation allows the transformer network to effectively process the combined information and capture the complex relationships between the text and the visual features. In a move towards greater efficiency, SD3 optimizes memory usage during inference by removing the memory-intensive T5 text encoder [12]. This optimization results in significantly reduced memory requirements without a substantial loss in the overall performance of the model, making it more practical for deployment on a wider range of hardware.

## B. OpenAI's DALL-E Family (DALL-E 3)

OpenAI's DALL-E family has consistently pushed the boundaries of text-to-image generation, known for its ability to generate imaginative and often surreal imagery [14].

**DALL-E 3:** The latest iteration, DALL-E 3, introduces a significant architectural enhancement through its tight integration with ChatGPT [16]. This unique combination allows users to leverage the conversational capabilities of ChatGPT to brainstorm and refine their prompts iteratively. Users can start with a simple idea and then engage in a dialogue with ChatGPT to add details, specify styles, and make adjustments, which are then translated into detailed prompts for the DALL-E 3 image generation model. This seamless collaboration between a language model and an image generation model significantly improves prompt adherence and allows for a more intuitive and adaptive creative process. The underlying image generation in DALL-E 3 is powered by a generative adversarial network (GAN) architecture [17]. GANs consist of two neural networks, a generator and a discriminator, that are trained in a competitive manner. The generator aims to produce realistic images from the text prompt, while the discriminator tries to distinguish between real images and those generated by the

generator. This adversarial training process enables GANs to generate highly detailed and realistic images. The image generation process in DALL-E 3 involves iteratively refining an initial image based on the text description and the feedback received from ChatGPT [17]. This iterative approach, guided by the conversational AI, likely contributes to the improved quality and the model's ability to closely follow the nuances of the user's prompts. A crucial factor in the enhanced performance of DALL-E 3 is the use of improved synthetic captions for its training dataset [19]. Recognizing that the quality of training data significantly impacts the model's ability to understand and follow instructions, OpenAI developed a bespoke image captioning system to generate more detailed and accurate descriptions of the images in its training set. By training on this higher-quality data, DALL-E 3 demonstrates a notable improvement in its ability to translate complex textual descriptions into corresponding visual outputs.

## C. Google's Imagen Family (Imagen 3)

Google's Imagen family has established itself as a strong contender in the text-to-image arena, particularly known for its ability to generate photorealistic images with a deep understanding of language [1].

**Imagen 3:** The latest evolution, Imagen 3, continues this trend by building upon a latent diffusion model architecture [21]. Operating in the latent space, a compressed representation of the image, allows Imagen 3 to handle large-scale image generation tasks with greater computational efficiency while maintaining high-quality outputs. This approach enables faster training times and more efficient inference, making it practical for a wide range of applications. A key feature of Imagen 3 is its support for multi-stage upsampling [22]. The model natively generates images at a resolution of 1024x1024 pixels but can then upscale these images through multiple stages to achieve resolutions exceeding 8000x8000 pixels. This process ensures that even when images are significantly enlarged, they retain their sharpness and visual coherence, making Imagen 3 well-suited for creating large-format visuals or detailed print media. Imagen 3 also demonstrates advanced spatial and numerical reasoning capabilities [22]. It can accurately depict specific quantities of objects and understand complex spatial relationships described in the text prompt, a capability that is crucial for applications such as advertising and product design where precise visual configurations are required. Furthermore, Imagen 3 excels in generating photorealistic imagery and exhibits strong prompt-image alignment [21]. Its ability to closely adhere to the details provided in the text prompt, combined with its capacity to produce lifelike representations, makes it a preferred choice for applications that demand realistic visuals, such as marketing materials and simulation environments.

## D. Midjourney

Midjourney has gained significant popularity for its ability to generate highly artistic and aesthetically pleasing images, often with a unique and distinctive style [1]. Unlike some other leading models, the specific architectural details of Midjourney are not publicly disclosed [25]. Midjourney operates as an independent research lab focused on exploring new mediums of thought and expanding the imaginative powers of the human species [25]. User interaction with Midjourney primarily occurs through the Discord platform, where users provide text prompts along with various parameters to control the style, aesthetics, and other characteristics of the generated images [26]. Prompt engineering plays a crucial role in achieving the desired results with Midjourney, as the model is highly responsive to the specific wording and details included in the prompt [27]. Despite the lack of transparency regarding its internal architecture, Midjourney has consistently demonstrated its ability to generate highly detailed and visually appealing images, particularly in the realm of artistic creations [2]. Its unique aesthetic and capacity for creative interpretation have made it a favorite among artists, designers, and enthusiasts seeking visually striking and imaginative outputs.

## E. Recraft V3

Recraft V3, also known by its model name "red_panda" [32], is presented as a revolutionary AI model that "thinks in design language" [33]. This suggests that its architecture is specifically designed to understand and generate visuals based on principles of graphic design, making it particularly well-suited for design-oriented tasks. A standout feature of Recraft V3 is its exceptional ability to generate images with long and accurate text [32]. This capability distinguishes it from many other text-to-image models that often struggle with rendering extended text passages legibly and correctly within the generated image. Furthermore, Recraft V3 supports both raster and vector output formats [32]. The ability to generate vector graphics is particularly valuable for designers as it allows for the creation of scalable images that can be resized without loss of quality, making it ideal for logos, illustrations, and other design elements. The model also offers precise control over various aspects of image generation, including the size and positioning of text, style customization, improved inpainting (editing within an image), and new outpainting (extending beyond the image boundaries) capabilities [33]. This level of user control makes Recraft V3 a powerful tool for designers who require specific visual outcomes and the ability to fine-tune their creations. Notably, Recraft V3 has achieved the top ranking in Hugging Face's Text-to-Image Benchmark by Artificial Analysis, outperforming all competitor models in key quality metrics [33]. This benchmark performance provides strong

evidence of its state-of-the-art capabilities in image generation.

## F. Luma Photon

Luma Photon, announced by Luma Labs, is built on a new and groundbreaking architecture that is claimed to deliver ultra-high image quality with significantly improved cost efficiency, up to 10 times better than comparable models [7]. This suggests a substantial advancement in both the performance and resource utilization of the model. At the core of Luma Photon is a bespoke Universal Transformer architecture [39]. The specifics of this architecture are still emerging, but the use of a transformer-based approach aligns with the trend towards leveraging these networks for their capabilities in handling complex data relationships. A notable feature of Luma Photon is its powerful new image reference system [37]. This system allows users to provide multiple images as part of their prompt, enabling more intuitive and flexible control over the generated output by guiding the model with visual examples. Luma Photon also boasts the unique capability (currently in beta) of generating consistent characters from just a single input image [37]. This feature is particularly valuable for storytelling and campaign creation, as it allows users to maintain a consistent visual identity across multiple generated images. According to Luma Labs, Photon has outperformed every other model on the market in terms of quality, creativity, and understanding in large-scale, double-blind evaluations [37]. This claim, based on rigorous user preference testing, positions Luma Photon as a leading model in terms of subjective visual quality and its ability to interpret and fulfill user prompts effectively.

## G. Ideogram

Ideogram models, particularly Ideogram 2.0, are built upon a transformer-based architecture that is specifically designed and optimized for text comprehension, generation, and editing within images [42]. This architectural focus allows Ideogram to excel in tasks that require accurate and legible text rendering. Ideogram is recognized for its ability to create images with perfectly rendered text most of the time [30]. This makes it a valuable tool for generating content such as posters, social media graphics, and other visuals where integrated text is essential. The architecture incorporates a refined attention mechanism that enhances the model's ability to process and generate large volumes of text while maintaining high coherence and contextual accuracy [42]. This suggests an advancement over standard transformer attention mechanisms, specifically tailored for handling text-rich image generation tasks. According to its developers, Ideogram 2.0 significantly outperforms existing models in key metrics such as image-text alignment, subjective preference, and, most notably, text rendering accuracy [42]. This indicates a substantial improvement over its

predecessor and positions Ideogram 2.0 as a leader in generating images where the accurate integration of text is a primary requirement.

**H. Amazon Titan Image Generator**

The Amazon Titan Image Generator utilizes a text-conditioned diffusion model as its core architecture for image generation [44]. This is a common and effective approach in modern text-to-image models, where the input text prompt is encoded into numerical vectors that then guide the diffusion process, transforming random noise into a coherent image that aligns with the textual description. The Titan Image Generator offers a broad range of capabilities beyond basic text-to-image generation, including image conditioning (generating images with similar composition to a reference image), various image editing functionalities such as inpainting and outpainting, automatic background removal, and the ability to control the color palette of the generated images [7]. This comprehensive set of features makes it a versatile tool for a wide array of image manipulation tasks. Furthermore, the Amazon Titan Image Generator incorporates built-in support for responsible AI use [45]. This includes mechanisms for detecting and removing harmful content from the training data and for rejecting inappropriate user inputs. Additionally, all images generated by the model include an invisible watermark by default, which is designed to help reduce the spread of misinformation and encourage the responsible use of AI-generated content.

**I. Google's Gemini**

Google's Gemini model leverages Imagen as its key architecture for text-to-image generation [21]. Therefore, the underlying principles of its image generation capabilities are rooted in the latent diffusion model architecture, similar to Imagen 3. However, Gemini extends its capabilities far beyond just text-to-image generation [21]. It is designed as a highly multimodal model capable of understanding and generating a wide variety of inputs and outputs, including text, images, audio, video, and even code. This signifies a sophisticated and versatile architecture that can process and synthesize information across different modalities. Specifically regarding image generation, Gemini 2.0 Flash Experimental supports the ability to generate both text and inline images [7]. This enables conversational image editing, where users can interact with the model to make changes to an image through natural language prompts, and the model can respond with updated images, creating a more interactive and intuitive editing experience. The focus on multimodal capabilities and conversational interaction distinguishes Gemini from models primarily focused solely on text-to-image generation.

# Key Architectural Components and Techniques

The diverse architectures of text-to-image models rely on several fundamental components and techniques that are crucial to their functionality.

## A. Transformer vs. U-Net in Diffusion Models

The choice between transformer-based architectures and U-Net architectures has become a significant point of discussion in the field of diffusion models for image generation. Traditionally, the U-Net architecture has been the dominant choice [6]. Known for its encoder-decoder structure with skip connections, the U-Net excels at capturing local features and preserving the spatial resolution of the image, which were initially considered highly beneficial for the iterative denoising process in diffusion models [13]. However, with the remarkable success of transformer networks in natural language processing and increasingly in computer vision, transformer-based architectures, particularly Diffusion Transformers (DiT), are gaining prominence in text-to-image generation [6]. Transformers operate on latent patches of the image and are adept at capturing long-range dependencies and global context through their attention mechanisms [13]. While U-Nets might face challenges in modeling these long-range interactions [13], transformers inherently possess this capability, which is increasingly seen as advantageous for understanding complex prompts and generating coherent scenes. The shift towards transformer-based architectures in leading models like Stable Diffusion 3 and Luma Photon suggests a potential superiority in certain aspects of text-to-image synthesis, especially as the complexity and fidelity of desired outputs continue to increase.

## B. Diffusion Process Techniques

The core of many modern text-to-image models lies in the diffusion process, which involves several key techniques. Latent Diffusion Models (LDMs) have become a standard approach [1]. LDMs perform the computationally intensive diffusion and denoising steps in the latent space of a pre-trained autoencoder. This significantly reduces the computational resources required compared to operating directly in the pixel space, enabling the generation of high-resolution images more efficiently. The diffusion process itself typically involves two main phases: forward diffusion and reverse diffusion [6]. In forward diffusion, Gaussian noise is gradually added to a training image over multiple steps until it becomes indistinguishable from pure noise. The reverse diffusion, or denoising, process is where the model comes into play. It learns to predict and iteratively remove the added noise, step by step, conditioned on the input text prompt, to reconstruct the original image [6]. This process is often guided by conditioning, where the text prompt embedding is provided to the model to influence

the characteristics of the generated image [6]. To further enhance the quality and the model's adherence to the prompt, a technique called classifier-free guidance is often employed [52]. This involves training a conditioned model where the conditioning input (like the caption) is sometimes randomly dropped during training. At inference time, the model is run both with and without the conditioning, and the results are combined to improve the output. More recently, a technique called flow matching has emerged as a way to simplify the training of diffusion models by directly learning a mapping between the noise and data distributions [4].

### C. Text Encoding Strategies (Single vs. Dual)

The way in which the input text prompt is encoded into a numerical representation plays a crucial role in the performance of text-to-image models. Many early models and some current ones utilize a single text encoder, such as CLIP (Contrastive Language-Image Pre-training) [50]. CLIP has been a popular choice due to its ability to learn a joint embedding space for text and images, effectively aligning textual descriptions with their corresponding visual features. However, some more recent architectures, like SDXL, employ dual text encoders [9]. By using two separate encoders (e.g., OpenCLIP and CLIP), these models can capture different levels of semantic information from the prompt, potentially leading to a richer understanding and the generation of more detailed and accurate images. For instance, one encoder might be better at capturing the overall scene description, while the other focuses on finer details or stylistic elements [11]. Another active area of research involves exploring the use of Large Language Models (LLMs) as text encoders [54]. LLMs possess superior text representation capabilities, support multiple languages, and can handle longer and more complex contexts compared to models like CLIP. While using LLMs as text encoders holds significant promise for improving the language understanding in text-to-image generation, it also presents challenges in effectively aligning the text features learned by LLMs with the visual information required for image synthesis.

## Insights from Latest Advancements and Research

The architectures and capabilities of text-to-image models are constantly evolving, driven by ongoing research and the introduction of new techniques. The increasing adoption of transformer-based architectures, as seen in Stable Diffusion 3 and Luma Photon, suggests a growing recognition of their advantages in scalability, performance, and handling complex prompts [6]. The integration of flow matching in models like FLUX points towards potentially more efficient training methods and improved image coherence [4]. The success of DALL-E 3 highlights the critical role of high-quality training data, with the use of synthetic captions demonstrating a

significant impact on prompt following abilities [19]. Google's Gemini exemplifies the trend towards greater multimodality, with its capacity to handle and generate various types of data beyond just text and images [21]. Models like Ideogram and Recraft V3 showcase substantial advancements in text rendering accuracy within images, making them particularly valuable for applications where integrated text is important [30]. The continued exploration of using LLMs as text encoders indicates a strong interest in further enhancing the language understanding capabilities of these models [54]. Finally, the development of new evaluation metrics tailored for creative use cases, as seen with Luma Photon, underscores the need for more nuanced ways to assess the performance of these models beyond traditional metrics [37]. These advancements collectively demonstrate a rapid pace of innovation in the field, with a clear focus on improving efficiency, quality, prompt adherence, and expanding the range of applications for text-to-image technology.

## Towards the "Best" Architecture: Considerations and Trade-offs

Determining the "best" architecture for a text-to-image model is not a straightforward task, as the ideal choice is highly dependent on the specific use case, the desired characteristics of the output images, and the available computational resources. Several crucial trade-offs must be considered when evaluating different architectures. For instance, some models might prioritize generating extremely high-fidelity images with intricate details, but this often comes at the cost of slower generation times, as seen with models like Imagen 3. Conversely, other architectures might focus on achieving rapid generation speeds, which can be beneficial for iterative design processes, but may involve some compromises in the ultimate image quality, as exemplified by the FLUX Schnell variant [5]. The size and complexity of the model also play a significant role. Larger models with a greater number of parameters, such as SDXL and SD3, have the potential to generate more complex and detailed images, but they typically require more powerful hardware and computational resources for training and inference. On the other hand, more efficient architectures, like the one claimed by Luma Photon, aim to strike a better balance between performance and resource usage.

The level of prompt adherence versus the degree of creative interpretation is another important consideration. Some models, particularly those integrated with conversational AI like DALL-E 3, strive for a very strict adherence to the user's prompt, ensuring that the generated image closely matches the described scene [18]. In contrast, other models, such as Midjourney, are known for their more artistic and sometimes unexpected interpretations of prompts, which can be desirable for creative exploration but might not be suitable for applications requiring precise replication of a

description [24]. For specific use cases, such as graphic design or content creation that requires integrated text, the text rendering capabilities of a model become paramount. Models like Ideogram and Recraft V3 excel in this area [30], but this might not be a primary focus for all text-to-image architectures. The level of control and customization offered by a model is also a significant factor, especially for professional users. Models like Recraft V3 and the Stable Diffusion family provide users with more granular control over various aspects of the image generation process, allowing for fine-tuning of styles, compositions, and other parameters [8]. Finally, the choice between open-source and proprietary models presents its own set of trade-offs. Open-source models like Stable Diffusion offer greater flexibility, transparency, and the benefit of community contributions, but proprietary models like DALL-E and Midjourney might have access to more cutting-edge research and features, albeit with less visibility into their internal workings.

Considering the diverse needs of different applications further highlights the context-dependent nature of the "best" architecture. For graphic design, models with high text rendering accuracy and the ability to output vector graphics, such as Recraft V3, might be the most suitable [32]. For generating photorealistic imagery for applications like stock photos or marketing materials, models with strong prompt adherence and high resolution capabilities, such as Imagen 3 and DALL-E 3, would likely be preferred [18]. In the realm of artistic creation and visual exploration, models known for their aesthetic quality and creative interpretations, like Midjourney and Luma Photon, might be the top choices [24]. For applications requiring rapid prototyping or real-time generation, models with fast generation times, such as FLUX Schnell or potentially more lightweight and efficient architectures, would be advantageous [5]. Finally, for research and development purposes, open-source models with flexible architectures, such as the Stable Diffusion family and FLUX Dev, provide the necessary platform for experimentation and further innovation.

| Architecture Characteristic | Examples of Models Prioritizing This Characteristic | Potential Trade-offs |
| --- | --- | --- |
| Image Quality | Imagen 3, DALL-E 3, Luma Photon, Recraft V3 Pro | Slower generation speed, higher computational cost |

| | | |
|---|---|---|
| **Generation Speed** | FLUX Schnell, potentially lightweight models | Possible slight reduction in image quality |
| **Computational Resources** | Luma Photon (claimed efficiency), potentially smaller models | May limit complexity and detail of generated images |
| **Prompt Adherence** | DALL-E 3 (with ChatGPT), Imagen 3 | Potentially less creative or surprising outputs |
| **Text Rendering** | Ideogram, Recraft V3 | Might not be a primary focus in all architectures |
| **Control & Customization** | Recraft V3, Stable Diffusion family | Can require more expertise in prompt engineering |
| **Artistic Interpretation** | Midjourney, Luma Photon | May not always strictly follow the prompt |
| **Vector Output** | Recraft V3 | Not a common feature in most text-to-image models |

## Conclusion

The landscape of text-to-image model architectures is characterized by a remarkable diversity of approaches, each with its own set of strengths and trade-offs. From the hybrid architecture of FLUX to the transformer-based innovations in Stable Diffusion 3 and Luma Photon, and the GAN-powered creativity of DALL-E 3, the field is witnessing rapid innovation across various architectural paradigms. While diffusion models, often operating in the latent space for efficiency, remain a dominant approach, the increasing adoption of transformer networks highlights their growing importance in handling complex prompts and achieving high-quality image synthesis. Advancements in areas such as text rendering, prompt understanding, efficiency, and multimodality continue to expand the potential applications of these models. Ultimately, the "best" architecture is not a fixed entity but rather a dynamic concept that depends on the specific requirements and priorities of the user or application. The choice involves navigating a complex interplay between factors like image quality, generation speed, computational resources, prompt adherence, text rendering capabilities, user control, and the desired level of artistic interpretation. As research in this field progresses, further innovations in model architecture and training techniques are expected to push the boundaries of what is possible with

text-to-image generation, promising even more powerful and versatile tools for creative expression and technological advancement in the future. Continued exploration and experimentation with the diverse range of available models will be key for users to identify the architecture that best aligns with their specific needs and goals.

**Works cited**

1. Text-to-image model - Wikipedia, accessed March 19, 2025, https://en.wikipedia.org/wiki/Text-to-image_model
2. FLUX.1 vs Midjourney: Text to Image AI Showdown | getimg.ai Blog, accessed March 19, 2025, https://getimg.ai/blog/flux-1-vs-midjourney-ultimate-text-to-image-ai-showdown
3. FLUX.1: A Deep Dive - Superteams.ai, accessed March 19, 2025, https://www.superteams.ai/blog/flux-1-a-deep-dive
4. FLUX.1 Text-to-Image AI: Next-Gen Diffusion Model for Visual Fidelity, accessed March 19, 2025, https://www.ikomia.ai/blog/flux1-text-to-image-diffusion-model
5. Comparing Flux.1 and Stable Diffusion - E2E Networks, accessed March 19, 2025, https://www.e2enetworks.com/blog/comparing-flux-1-and-stable-diffusion---a-technical-deep-dive
6. Stable Diffusion 3: Guide to the Latest Text-to-Image Model by Stability AI - Analytics Vidhya, accessed March 19, 2025, https://www.analyticsvidhya.com/blog/2024/06/stable-diffusion-3/
7. Text to Image Models and Providers Leaderboard - Artificial Analysis, accessed March 19, 2025, https://artificialanalysis.ai/text-to-image
8. Comparing Text to Image models and providers - Recraft, accessed March 19, 2025, https://www.recraft.ai/blog/comparing-popular-and-high-performing-text-to-image-models-and-providers
9. Stable Diffusion XL - Hugging Face, accessed March 19, 2025, https://huggingface.co/docs/diffusers/training/sdxl
10. research-papers/Summaries/Diffusion/SDXL.md at main - GitHub, accessed March 19, 2025, https://github.com/garg-aayush/research-papers/blob/main/Summaries/Diffusion/SDXL.md
11. SDXL: Two text encoders, two text prompts : r/StableDiffusion - Reddit, accessed March 19, 2025, https://www.reddit.com/r/StableDiffusion/comments/15c2n0q/sdxl_two_text_encoders_two_text_prompts/
12. Stable Diffusion 3: Multimodal Diffusion Transformer Model Explained - Encord, accessed March 19, 2025, https://encord.com/blog/stable-diffusion-3-text-to-image-model/
13. Diffusion Transformer (DiT) Models: A Beginner's Guide - Encord, accessed March 19, 2025, https://encord.com/blog/diffusion-models-with-transformers/

14. Comparing Text-to-Image models: DALL-E and Stable Diffusion | by vTeam.ai | Data Science in your pocket | Medium, accessed March 19, 2025, https://medium.com/data-science-in-your-pocket/comparing-text-to-image-models-dall-e-and-stable-diffusion-131795e4f037

15. The best AI image generators of 2024: Tested and reviewed | ZDNET, accessed March 19, 2025, https://www.zdnet.com/article/best-ai-image-generator/

16. AlonzoLeeeooo/awesome-text-to-image-studies: A collection of awesome text-to-image generation studies. - GitHub, accessed March 19, 2025, https://github.com/AlonzoLeeeooo/awesome-text-to-image-studies

17. DALL-E 3: A Fusion of Imagination and Conversation | by Sukriti Mehrotra | Medium, accessed March 19, 2025, https://medium.com/@sukritimehrotra/dall-e-3-a-fusion-of-imagination-and-conversation-4c8ec8930442

18. DALL·E 3 | OpenAI, accessed March 19, 2025, https://openai.com/index/dall-e-3/

19. DALLE-3 research paper , Adept's Fuyu-8B open source , Waymo simulator - TLDR, accessed March 19, 2025, https://tldr.tech/ai/2023-10-20

20. Improving Image Generation with Better Captions - OpenAI, accessed March 19, 2025, https://cdn.openai.com/papers/dall-e-3.pdf

21. Text-to-image AI | Google Cloud, accessed March 19, 2025, https://cloud.google.com/use-cases/text-to-image-ai

22. What is Imagen 3? - Integrail, accessed March 19, 2025, https://integrail.ai/blog/what-is-imagen-3

23. A developer's guide to Imagen 3 on Vertex AI | Google Cloud Blog, accessed March 19, 2025, https://cloud.google.com/blog/products/ai-machine-learning/a-developers-guide-to-imagen-3-on-vertex-ai

24. Tested: The Best AI Image Generators for 2025 - PCMag, accessed March 19, 2025, https://www.pcmag.com/picks/the-best-ai-image-generators

25. Midjourney, accessed March 19, 2025, https://www.midjourney.com/

26. Generative AI Meets Architecture: Using Midjourney to Generate Innovative Ideas - Maket.ai, accessed March 19, 2025, https://www.maket.ai/post/generative-ai-meets-architecture-using-midjourney-to-generate-innovative-ideas

27. How to use Midjourney for architect? - Future Architecture, accessed March 19, 2025, https://futurearchi.blog/en/midjourney-architect/

28. Architects' Guide To Midjourney: An Adventure in AI-Generated Imagery for Concept Development - Architizer, accessed March 19, 2025, https://architizer.com/blog/practice/tools/an-architects-guide-to-midjourney-ai-generated-imagery/

29. Midjourney Prompts for Architecture - Freeflo, accessed March 19, 2025, https://freeflo.ai/midjourney-prompts/architecture

30. Best AI image generators of 2025 - Tom's Guide, accessed March 19, 2025, https://www.tomsguide.com/best-picks/best-ai-image-generators

31. The 7 best AI image generators in 2025 - Zapier, accessed March 19, 2025, https://zapier.com/blog/best-ai-image-generator/

32. Recraft V3 AI: Advanced Image Generator - PIXEL DOJO, accessed March 19, 2025, https://pixeldojo.ai/recraft-v3-ai-image-generator

33. Recraft introduces a revolutionary AI model that thinks in design language, accessed March 19, 2025, https://www.recraft.ai/blog/recraft-introduces-a-revolutionary-ai-model-that-thinks-in-design-language

34. Recraft Blog, accessed March 19, 2025, https://www.recraft.ai/blog

35. Best DALL-E alternatives - Recraft, accessed March 19, 2025, https://www.recraft.ai/blog/best-dall-e-alternatives

36. Recraft V3 | Text to Image | AI Playground - Fal.ai, accessed March 19, 2025, https://fal.ai/models/fal-ai/recraft-v3

37. Luma Photon, accessed March 19, 2025, https://www.lumalabs.ai/photon

38. Luma Photon Image Generation Model: Iteratively Generating Images with Natural Language Descriptions, Balancing Picture Quality and Creativity - Chief AI Sharing Circle, accessed March 19, 2025, https://www.aisharenet.com/en/luma-photon-tuxiangbeng/

39. AI Video Generator API - Luma Dream Machine, accessed March 19, 2025, https://lumalabs.ai/dream-machine/api

40. Consistent AI Image Generators : luma ai photon - Trend Hunter, accessed March 19, 2025, https://www.trendhunter.com/trends/luma-ai-photon

41. This new AI image generator lets you create reusable characters - ZDNet, accessed March 19, 2025, https://www.zdnet.com/article/this-new-ai-image-generator-lets-you-create-reusable-characters/

42. Ideogram 2.0: Generating Text on Images with Unmatched Accuracy, accessed March 19, 2025, https://neurohive.io/en/ai-apps/ideogram-2-0-generating-text-on-images-with-unmatched-accuracy/

43. Stable Diffusion 3: Research Paper - Stability AI, accessed March 19, 2025, https://stability.ai/news/stable-diffusion-3-research-paper

44. Amazon Titan Image Generator - AWS AI Service Cards, accessed March 19, 2025, https://docs.aws.amazon.com/ai/responsible-ai/titan-image-generator/overview.html

45. Amazon Titan Image Generator Demo - AWS, accessed March 19, 2025, https://aws.amazon.com/de/awstv/watch/5e67a1e3606/

46. Use Amazon Titan models for image generation, editing, and searching | AWS Machine Learning Blog, accessed March 19, 2025, https://aws.amazon.com/blogs/machine-learning/use-amazon-titan-models-for-image-generation-editing-and-searching/

47. Generate images | Gemini API | Google AI for Developers, accessed March 19, 2025, https://ai.google.dev/gemini-api/docs/image-generation

48. From CNN to Transformer: A Review of Medical Image Segmentation Models - PMC, accessed March 19, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC11300773/

49. Why U-Net instead of Transformers? : r/learnmachinelearning - Reddit, accessed March 19, 2025, https://www.reddit.com/r/learnmachinelearning/comments/1c2kbsx/why_unet_instead_of_transformers/

50. nicd.org.uk, accessed March 19, 2025, https://nicd.org.uk/knowledge-hub/image-to-text-latent-diffusion-models#:~:text=The%20diffusion%20process%20creates%20the,together%20text%20and%20image%20modalities.

51. Diffusion and Denoising: Explaining Text-to-Image Generative AI - Exxact Corporation, accessed March 19, 2025, https://www.exxactcorp.com/blog/deep-learning/diffusion-and-denoising-explaining-text-to-image-generative-ai

52. From DALL·E to Stable Diffusion: How Do Text-to-Image Generation Models Work?, accessed March 19, 2025, https://www.edge-ai-vision.com/2023/01/from-dall%C2%B7e-to-stable-diffusion-how-do-text-to-image-generation-models-work/

53. From Text to Image: How do image generators like DALL-E 2 work? - Mario Dias - Medium, accessed March 19, 2025, https://itsmariodias.medium.com/from-text-to-image-how-do-image-generators-like-dall-e-2-work-17703661944c

54. An Empirical Study and Analysis of Text-to-Image Generation Using Large Language Model-Powered Textual Representation - arXiv, accessed March 19, 2025, https://arxiv.org/html/2405.12914v2