# Padding Tone: A Mechanistic Analysis of Padding Tokens in T2I Models

**Michael Toker**[1]    **Ido Galil**[1,2]    **Hadas Orgad**[1]    **Rinon Gal**[2]    **Yoad Tewel**[2]
**Gal Chechik**[2,3]    **Yonatan Belinkov**[1]
[1]Technion – Israel Institute of Technology    [2]NVIDIA    [3]Bar-Ilan University

## Abstract

Text-to-image (T2I) diffusion models rely on encoded prompts to guide the image generation process. Typically, these prompts are extended to a fixed length by adding padding tokens before text encoding. Despite being a default practice, the influence of padding tokens on the image generation process has not been investigated. In this work, we conduct the first in-depth analysis of the role padding tokens play in T2I models. We develop two causal techniques to analyze how information is encoded in the representation of tokens across different components of the T2I pipeline. Using these techniques, we investigate when and how padding tokens impact the image generation process. Our findings reveal three distinct scenarios: padding tokens may affect the model's output during text encoding, during the diffusion process, or be effectively ignored. Moreover, we identify key relationships between these scenarios and the model's architecture (cross or self-attention) and its training process (frozen or trained text encoder). These insights contribute to a deeper understanding of the mechanisms of padding tokens, potentially informing future model design and training practices in T2I systems.

## 1 Introduction

Text-to-image (T2I) models consist of two main components: a text encoder, which generates representations of the user's prompt, and a diffusion model, which generates an image based on this representation. To standardize sequence lengths for efficient batch processing in training and inference, input prompts are padded to a fixed length with a special padding token. Unlike language models, where padding tokens are explicitly masked and thus ignored, the computation process of the T2I models can use these tokens as any other token. Despite their ubiquity, the potential impact of
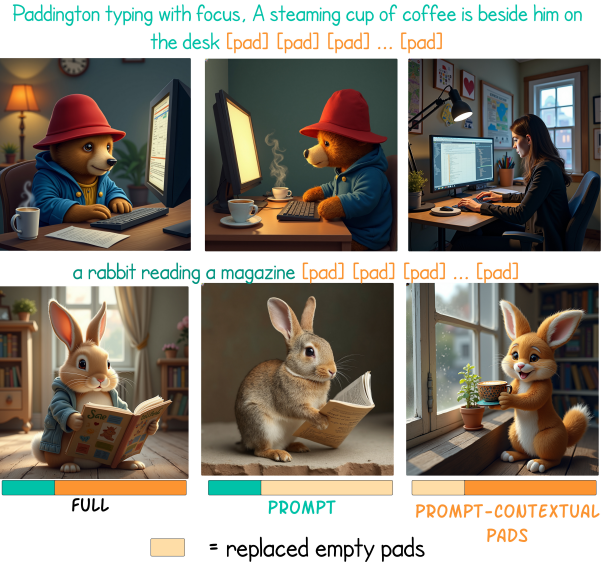


Figure 1: Images generated with FLUX from different segments of the input prompt. Description of each column, from left to right: (1) An image generated using the full prompt (both prompt tokens and padding tokens encoded together), (2) An image generated using only the prompt tokens and clean padding tokens, (3) An image generated using only the prompt-contextual pads encoded with the prompt, while the prompt tokens were replaced with clean pad tokens.

padding tokens on image generation outcomes has been overlooked.

We introduce two methods to evaluate the influence of tokens on different model components: (1) Intervention in the Text Encoder Output (ITE) and (2) Intervention in the Diffusion Process (IDP). Both methods build on causal mediation analysis, also known as activation patching (Imai et al., 2010; Vig et al., 2020; Zhang and Nanda, 2024). This technique involves perturbing specific inputs or intermediate representations to observe their effect on the output, helping to pinpoint the influential elements.

In ITE we selectively perturb specific segments of the text encoder's output representations to iso-

late the contributions of two key elements: prompt tokens and padding tokens. Next, we generate images using the modified prompt representations and analyze the results. The perturbation involves replacing selected token representations with those from a prompt that consists solely of padding tokens, referred to as *clean pads*. These clean pads differ from the original padding tokens, which contain contextual information from the prompt. The method is illustrated in Figure 2. If padding tokens carry meaningful information, we expect two outcomes: (a) replacing the prompt tokens with clean pads should still result in an image reflecting elements of the original prompt, while (b) replacing the padding tokens with clean pads should alter the image either semantically or stylistically.

In cases where our analysis with ITE indicates that padding tokens are not used by the text encoder, we further examine the role of padding tokens in the diffusion process. Particularly, we investigate whether significant information is written into the padding token representations throughout the diffusion process. Here we employ IDP, illustrated in Figure 8, to interpret the causal effect of the padding tokens during the diffusion process. We begin with a standard prompt padded to a fixed length, as well as an "only pads" prompt. However, in IDP, token replacement occurs before each attention block within the diffusion process and at every diffusion step. We repeat the procedure of selectively replacing either prompt tokens or padding tokens with clean pads, similarly to ITE. Figure 1 illustrates an example of images generated using this method.

We analyze six different T2I models and highlight two scenarios where padding tokens are utilized. First, when the text encoder was not frozen during training or fine-tuning, it learns to encode meaningful semantic information into these tokens. Second, in architectures with multi-modal attention mechanisms—such as Stable Diffusion 3 (Esser et al., 2024) and FLUX[1]—padding tokens carry meaningful information throughout the diffusion process, even if the text encoder itself does not directly encode it. Here, the padding tokens seem to act as "registers", with information written into their representations to store and recall, similarly to findings from both language models and vision-language models (Darcet et al., 2024; Burtsev et al., 2020).

---

[1]blackforestlabs.ai

To summarize, our main contributions are:

1. We propose two causal methods for analyzing the use of specific tokens in both the text encoder and diffusion model of the T2I pipeline, and apply them to investigate the role of padding tokens.

2. We find that T2I models with frozen text encoders ignore padding tokens. However, when the text encoder is trained or fine-tuned, padding tokens gain semantic significance.

3. We uncover that even when padding tokens are not utilized by the text encoder, for some architectures of the diffusion model, they can still function as "registers" and play a meaningful part in the diffusion process.

## 2 Analysis of Padding in Text Encoding

In the T2I pipeline, the text encoder processes the input prompt $P = [P_1, .., P_k]$, with a prompt length of $k$. To ensure a consistent input length, the prompt is usually padded to a fixed length, denoted as $N$. We denote this padded version of the prompt as $P_{\text{full}}$, which is a concatenation of the $k$ prompt tokens and the $N - k$ padding tokens:

$$P_{\text{full}} = [P_1, \ldots, P_k, \mathsf{pad}, \ldots, \mathsf{pad}]. \quad (1)$$

The text encoder then processes $P_{\text{full}}$, producing a constant-length encoded representation, which is subsequently used by the diffusion model for conditional image generation. We denote this encoded full prompt representation as $E_{\text{full}}$.

### 2.1 Method

Our goal is to evaluate the information encoded in the prompt-contextual padding tokens, and to measure their effect on the generated image. To do so, as illustrated in Figure 2, we generate images using partial representations of $E_{\text{full}}$ that isolate the effect of the padding tokens. We generate images based on modified representations of $E_{\text{full}}$ and compare them to images generated from the full prompt $E_{\text{full}}$. This enables us to visually express the information from different parts of the text input.

Specifically, to remove information coming from a subset of the tokens, we replace them with "clean" padding tokens that were not influenced by the user's prompt. To obtain these clean padding tokens, we encode $S_{\text{clean}} = [\mathsf{pad}, \mathsf{pad}, \ldots, \mathsf{pad}]$, a fixed-length sequence made entirely of padding tokens, and denote their embeddings as $E_{\text{clean}}$.

These encoded pad tokens are then used in constructing the final mixed representation, which com-
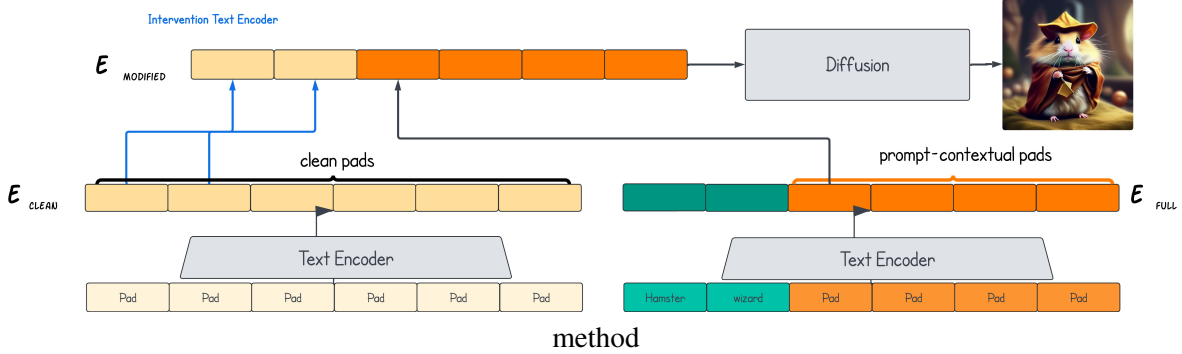
method

Figure 2: ITE: Interpreting information within pad tokens in the text encoder. We first encode the full prompt and an clean pads separately. Next, we keep the tokens we want to interpret and replace all other tokens with clean pad tokens. We then generate an image conditioned on this mixed representation. In the example shown here, we interpret the pad tokens in LLaMA-UNet, revealing semantic information embedded within the pad tokens.
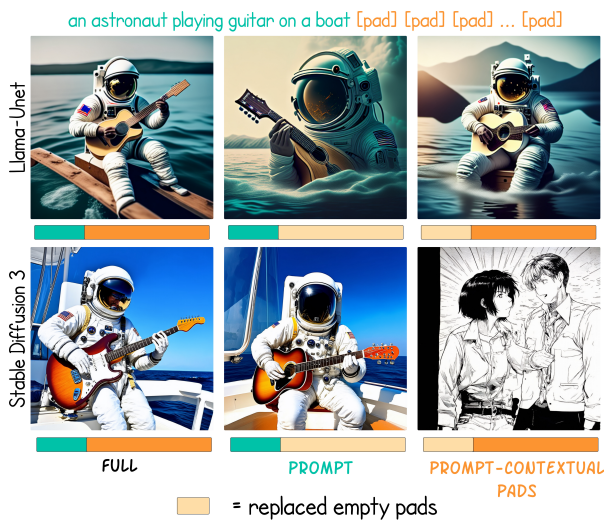


Figure 3: Images generated from different segments of the input prompt using ITE. Description of each column, from left to right: (1) An image generated using the full prompt (both prompt tokens and padding tokens encoded together), (2) An image generated using only the prompt tokens and clean padding tokens, (3) An image generated using only the prompt-contextual pads encoded with the prompt, while the prompt tokens were replaced with clean pad tokens.

bines both the encoded prompt and clean padding tokens. We use the encoded pad tokens since they contain no information related to the current prompt, while maintaining the same length and distribution of the text encoder's output. This allows us to effectively isolate the contribution of the padding tokens that are encoded alongside the full prompt tokens, helping us understand how much of the information in the final representation comes from the prompt itself versus the pads. Figure 3 demonstrated our method. First, we generate an image from the full prompt, which is how the image

is generated in the standard pipeline (left column). Then, we generate an image that demonstrates the information in the non-pad tokens, by replacing the pad tokens with clean pads (middle column). Lastly, we generate an image demonstrating the information within the pad tokens, by replacing the non-pad tokens with clean pads (left column).

More formally, the mixed representation for generating an image from the prompt tokens only (middle column) is:

$$E_{\text{prompt}} = \left[ E_{\text{full}}^{0:k}, E_{\text{clean}}^{k:N} \right], \qquad (2)$$

where $E_x^{i:j}$ represents the encoded tokens from index $i$ to $j$, and for a representation that generates an image from the prompt-contextual padding tokens only (right column):

$$E_{\text{pads}} = \left[ E_{\text{clean}}^{0:k}, E_{\text{full}}^{k:N} \right] \qquad (3)$$

## 2.2 Experimental Setup

**Models.** We use six T2I models. These models can be divided into two categories based on their training approach: those with pretrained frozen text encoders during the training: Stable Diffusion 3 (Esser et al., 2024), Stable Diffusion 2, Stable Diffusion XL (Podell et al., 2024), FLUX; and those with some learned weights as part of the text to image training: LDM (Rombach et al., 2022) and Lavi-Bridge (Zhao et al., 2024) (LLaMA-UNet version). The first group can be divided to two subgroups: models that use vision-language cross-attention with the text representations in the diffusion process (Stable Diffusion 2, Stable Diffusion XL) and models that use the text representations as part of vision-language self-attention, allowing
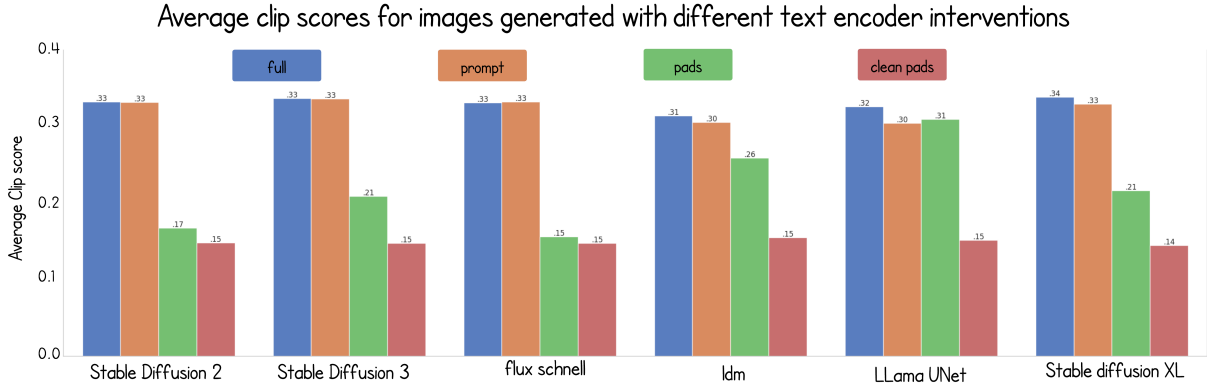
Figure 4: Average CLIP score over 5,000 images generated from the different representations: full prompt, only prompt, prompt-contextual pads and clean pads. LDM and LLaMA-UNet are the only models achieving high CLIP scores for images generated from padding tokens, indicating their use during text encoding. See Table 4 in the Appendix for standard deviations.

text representations to change throughout diffusion (FLUX, Stable Diffusion 3). Appendix C provides more information regarding each of the models.

**Data.** Our prompts are based on the Parti dataset (Yu et al., 2022), a benchmark containing over 1600 diverse and challenging prompts used to evaluate T2I models. To prevent using prompts that have leaked into the training corpus of the models, we select prompts from eight different challenge categories in Parti, and use GPT-4o[2] to generate an alternative set of prompts with similar style and complexity. We then manually review the prompts to ensure their coherence. This process results in 500 new prompts. The complete list of categories, along with the prompt used with GPT, can be found in Appendix A, and the full dataset is included in the supplementary material.

Each of the 500 prompts is used to generate 10 images from different random seeds, resulting in 5,000 images for each configuration of model and representation. We investigate three representations: $E_{full}$, $E_{prompt}$ (Eq. 2), $E_{pads}$ (Eq. 3), and $E_{clean}$ as a lower-bound control, with their corresponding images denoted as "full", "prompt", "prompt-contextual pads" and "clean", respectively.

**Metrics.** To evaluate the generated images, we employ two key metrics: CLIP score (Hessel et al., 2021), which measures how well the generated images align with the prompts, and KID (Kernel Inception Distance) (Bińkowski et al., 2018), to evaluate the quality of generated images. KID is used to measure the similarity between the feature

distributions of images generated from full representation and generated images after some causal intervention. Unlike FID (Heusel et al., 2017), which is based on Gaussian approximations, KID uses the maximum mean discrepancy (MMD) measure, making it more robust in practice, especially when dealing with smaller sample sizes.

## 2.3 Results

Figure 4 shows the average CLIP score over generations from different representations: "full", "prompt", "prompt-contextual pads" and "clean". Stable Diffusion (versions 2+3) and FLUX models appear to make little to no use of padding tokens: CLIP scores for the full and prompt representations are nearly identical, while the prompt-contextual pads—containing only padding tokens—yield significantly lower scores. In contrast, LLaMA UNet and LDM contain significant semantic information in padding, with a higher CLIP score for the "prompt-contextual pads", although the degradation in performance from "full" to "prompt" is small.

**Text encoder training objective and its influence on padding usage** Our results suggest that the training objective of the text encoder significantly impacts how padding tokens are utilized. Many current T2I models, such as Stable Diffusion and FLUX, employ a frozen text encoder, with the diffusion model being trained on its encoded outputs. It may be that because the text encoder is not explicitly trained to process padding tokens for image generation, it does not effectively incorporate them during the textual encoding. As shown in Figure 4 and Table 1, in models that use frozen text encoders,

---

[2]openai.com/index/hello-gpt-4o

|  | KID Score | |
| Model | Prompt | Pads |
|---|---|---|
| **Flux-schnell** | 0.01 | 14.52 |
| **LDM** | 0.88 | 4.53 |
| **LLaMA UNet** | 7.37 | 0.48 |
| **Stable Diffusion 2** | 0.02 | 31.09 |
| **Stable Diffusion 3** | 0.01 | 15.74 |

Table 1: KID scores between the images generated from the prompt-contextual pads vs. images generated only from prompt representations. Lower is better. The KID is calculated w.r.t images generated from the full representation.

| Pad Segment | CLIP Score |
|---|---|
| 1 | $0.30 \pm 0.018$ |
| 2 | $0.23 \pm 0.018$ |
| 3 | $0.17 \pm 0.022$ |

Table 2: Average CLIP scores for different prompt-contextual pad segments in LLaMA-UNet: the first 20% of the pads, the next 20%, and then the subsequent 20%. We observe that the semantic information degrades gradually, with most of it concentrated in the initial tokens.

images generated using the "prompt" representation yield the same CLIP score as those generated using the "full" representation, while images generated from "prompt-contextual padding" representations result in a very low CLIP score, almost as low as those generated from clean padding. In these models, the prompt KID is very high, meaning that the images are out of distribution. This suggests that in these models, the text encoder does not encode any meaningful information in the padding tokens, which makes them unnecessary for generating the final image.

Other models, like LDM and Lavi-Bridge, propose adapting the text encoder specifically for the image generation task. These methods train the text encoder, including the use of padding tokens, on the image generation objective, allowing it to effectively learn how to utilize padding. In these models, the results differ: images generated from full prompt tokens have lower scores compared to those generated using prompt representations, suggesting that the information encoded in the prompt tokens is insufficient to generate the correct images. Furthermore, images generated from the prompt-contextual padding tokens in these models yield much higher CLIP scores, even surpassing images generated from full prompt tokens in one of the models. KID of "prompt-contextual pads" in these models is comparatively low, indicating that the images generated from pads come from a closer distribution compared to the images generated from the full representation. Overall, this indicates that padding tokens play an important role in the text encoding process for image generation in these adapted models.

**How many padding tokens do text encoders use?** We focus on the LLaMA-UNet model and analyze padding behavior. We divide the padding tokens into five segments, each containing 20% of the total padding tokens in their natural order. For each segment, we mask both the prompt tokens and pad tokens in the other segments, then generate images from this mixed representation.

The CLIP scores can be found in Table 2. Our observations reveal that the information encoded in padding tokens varies based on their proximity to the prompt tokens, with those closer to the prompt carrying more significant information. We hypothesize that this behavior may be due to the text encoder's use of causal masking or the positional encoding scheme applied to the padding tokens. Only the padding tokens that are closer to the prompt tokens appear to be utilized effectively.

Since LLaMA is a language model adapted for image generation using LoRa training, we can load the LoRa with a scaling factor, $\alpha$, to observe how gradually removing LoRa affects the number of used pad tokens. Our results in Figure 5 show that as $\alpha$ decreases, fewer pad tokens are used. This indicates that part of what the LoRa learns involves encoding information in more pad tokens.

## 3 Analysis of Padding in the Diffusion Process

Even when padding tokens contain no meaningful information after text encoding, the diffusion model might still make use of them during the diffusion process. To generate images from text prompts, T2I models use an attention mechanism to condition the generation process, typically following two common approaches: cross-attention and MM-DiT (Esser et al., 2024) blocks. In cross-attention, used in models like Stable Diffusion 2/XL, the model converts image patches into query vectors and text tokens into key and value vec-
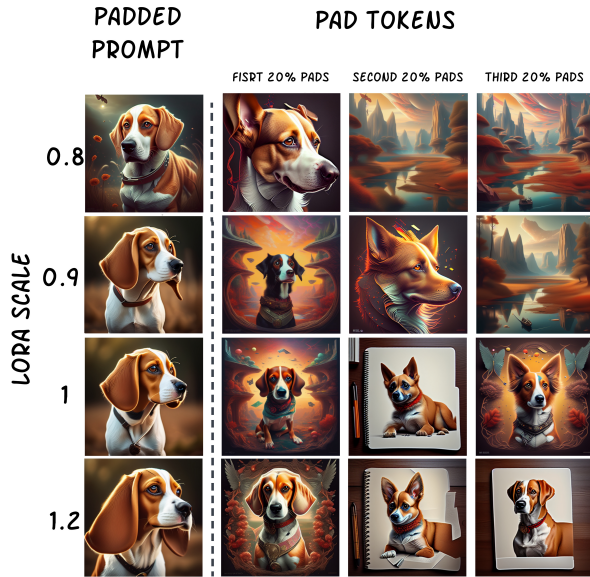
Figure 5: Images generated from Lavi-bridge with LoRa loaded with scaling factor $\alpha$ (y-axis). We analyze pad token segments: the first column shows the full image, and the next columns show three consecutive 20% of the pads. As $\alpha$ decreases, fewer pad tokens are used.
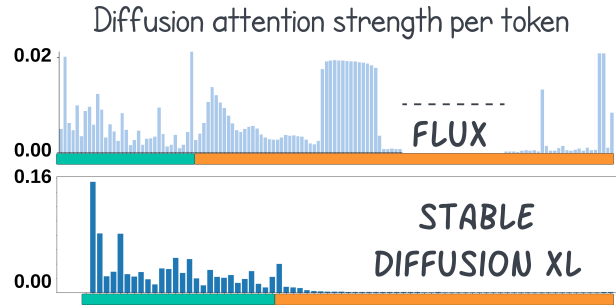


Figure 6: Attention histogram for Stable Diffusion XL and FLUX* for each token reveals that while both models exclude semantic information from padding tokens, FLUX utilizes these tokens, whereas Stable Diffusion does not. *In FLUX, we have removed the long middle part with low attention in order to improve visualization.



Figure 7: Attention maps for FLUX diffusion show strong alignment between prompt tokens and semantically relevant image tokens. These maps also reveal high attention for padding tokens with the main objects in the image.

tors. The image patches gather information from the text based on an attention map, but the text representation remains unchanged throughout the process. In contrast, MM-DiT blocks, found in models like FLUX and Stable Diffusion 3, implement a multi-modal self-attention, by projecting both image patches and text token representations into query, key, and value vectors. Thus, both the image and text representations update and influence each other during the attention process. We therefore expect that models implementing cross-attention where the pads are not used in the text encoder would also not use them in the diffusion process. However, models implementing *MM-DiT* blocks can potentially aggregate information into the padding tokens, even if initially they contain no information.

**Motivation: attention maps and qualitative examples.** To explore this, we examine the attention maps between image patches and text representations, resulting in an attention map for each token (see example in Figure 7). While in Stable Diffusion XL only the prompt (and the end-of-text) tokens significantly attend to main areas in the image, in FLUX not only prompt tokens, but also many pad tokens contribute much attention to image areas (Figure 6). Moreover, generating images with FLUX and Stable Diffusion XL, with and

without padding (Figure 11, App. E), reveals that FLUX without padding often misses key details, while Stable Diffusion XL remains consistent in its generations.

### 3.1 Method

To interpret the causal effect of tokens during the diffusion process, we develop IDP, illustrated in Figure 8. The diffusion process consists of several diffusion steps, where each step begins with the current latent image representation and the full encoded text representation. Since we look only at models where padding tokens do not carry meaningful information in the text encoder, we hypothesize that the diffusion model might be using these tokens as "registers" to store and recall information, subsequently passing it to the image patches, similar to the findings of Darcet et al. (2024) in their work on VLMs with image patches.

In this case, we conduct the intervention before each attention block to ensure that the attention mechanism fully incorporates the tokens we wish to interpret. We use a full prompt and a "clean pads" prompt, whose representations per diffusion layer are denoted as $E_{\text{full}}^{(l)}$ and $E_{\text{clean}}^{(l)}$, respectively. We explore two directions: first, we replace the
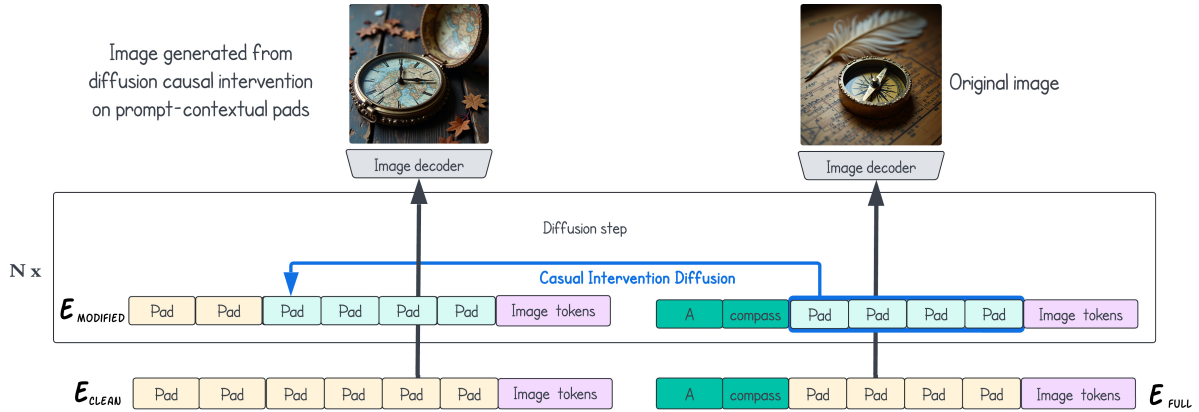
Figure 8: IDP: Interpreting information within pad tokens in the diffusion model. We perform a diffusion of two prompts simultaneously: the full prompt and an clean pads. During the diffusion, we keep the tokens we want to interpret (here: the prompt-contextual padding tokens) and replace all other tokens with clean pad tokens. We perform this intervention before each attention block in the diffusion model, through all diffusion steps. We then generate an image conditioned on this mixed representation. In the example shown here, we interpret the pad tokens in FLUX, revealing semantic information embedded within the pad tokens during diffusion.

| Representation | CLIP Reference | |
| | Image | Prompt |
| --- | --- | --- |
| Pads | 0.76 ±0.022 | 0.23 ±0.015 |
| Prompt | 0.90 ±0.036 | 0.33 ±0.028 |
| Clean | 0.46 ±0.018 | 0.10 ±0.009 |
| Full | 1.0 ±0.0 | 0.34 ±0.020 |

Table 3: Average CLIP scores between images generated (with FLUX) with different IDP interventions and either the full prompt or an image generated from the full prompt. 'Pad': prompt-contextual pads; 'Prompt': prompt tokens; 'Clean': a prompt full of pads, used for comparison; 'Full': a prompt with real tokens and pads.

prompt tokens with clean pads, using Equation 2. If the images generated from these representations still contain prompt-relevant information, it would suggest that the pads are being utilized. Second, we replace the prompt-contextual pads with the clean pads, as in Equation 3. If the resulting images remain unchanged, this would indicate that the pads are not used for encoding information.

## 3.2 Results

Table 3 shows the results of our intervention in the diffusion process. The table shows average CLIP scores of images generated with different IDP interventions, to assess the role of pad and prompt tokens. First, we compute CLIP scores vis-a-vis the full text prompt (Prompt column). As may be expected, images generated from the prompt tokens are similar to the prompt text to the same extent as images generated from the full prompt. Inter-

estingly, images generated from only the pads are much more similar to the text prompt than images generated from clean prompts, indicating that pad tokens are used by the diffusion model to produce images that relate to the prompt.

Next, we compute CLIP scores vis-a-vis images generated from the full prompt (Table 3, Image column). The CLIP score between images generated from full prompts and images generated when using only padding tokens in the diffusion is approximately 76—significantly higher than the score for randomly generated images from a 'clean' padding prompt. This is further evidence that the padding tokens contain visual information closely related to the content of the prompt tokens. However, the CLIP score when using images generated with IDP from the prompt tokens is still higher, suggesting that some information is lost when only using padding tokens in the diffusion model.

Finally, we provide qualitative example in Figure 9 and more examples in Figure 10 (Appendix E), which show that images generated from the prompt-contextual pads with IDP have meaningful semantic information. While images generated solely from prompt tokens typically align with the semantic meaning of the prompt, different visual features are often missing when padding tokens are excluded, while the same features are presented in the padding tokens. It appears that the diffusion model uses padding tokens to create additional visual information, while semantic content remains primarily in the prompt tokens.
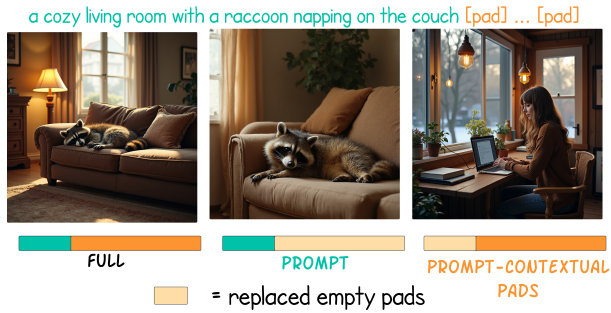
Figure 9: Images generated with FLUX from different prompt segments show distinct alignments: prompt tokens produce semantically accurate images, while the visual nuance like 'cozy' emerges only from the prompt-contextual pad tokens.

## 4 Related Work

**Special tokens and additional computation**
While padding tokens are generally used for efficient batch processing without fulfilling a functional role, other special tokens are known to carry various roles. In transformer language models, attention is often directed to special tokens, including punctuation marks ('.'), [SEP], or just the first token; this has been referred to as null or no-op attention (Vig and Belinkov, 2019; Kovaleva et al., 2019; Clark et al., 2019; Rogers et al., 2020). Some have added special tokens to enable additional processing, such as 'registers' in vision transformers (Darcet et al., 2024) or 'memory tokens' in language models (Burtsev et al., 2020). More generally, language models benefit from additional computation via chain-of-thought reasoning (Wei et al., 2024). Finally, several studies found it useful to *train* models to perform additional computation with custom tokens, including filler tokens like '.....' (Pfau et al., 2024), so-called 'pause tokens' (Goyal et al., 2024), or 'meta-tokens' for additional reasoning steps (Zelikman et al., 2024). This idea can be traced back to adaptive computation time techniques (Graves, 2016; Banino et al., 2021). Our work contributes to this literature by analyzing the role of padding tokens in T2I models.

**Interpreting vision-language models.** Compared to uni-modal models, VLMs have seen relatively few attempts at interpretation. CLIP (Radford et al., 2021) has been a focus of several studies: Goh et al. (2021) identified multimodal neurons responding to specific concepts, while Gandelsman et al. (2023) decomposed its image representations into text-based characteristics. In the realm of text-to-image models, Tang et al. (2023) introduced a method to interpret T2I pipelines by analyzing the influence of input words on generated images through cross-attention layers. Chefer et al. (2024) decomposed textual concepts, with a focus on the diffusion component. Basu et al. (2024) employed causal tracing to investigate the storage of knowledge in T2I models like Stable Diffusion. Toker et al. (2024) analyzed the text encoder in T2I pipelines, offering a view into intermediate representations rather than just its final output.

Our work takes a unique direction by focusing specifically on padding tokens, which have been largely overlooked in prior research. While previous research has illuminated how prompt tokens guide image generation, we show that padding tokens, often thought to be inert, can play a more active role—encoding semantic information or even functioning as "registers" that influence model computations. This adds a new dimension to the interpretation of T2I models, suggesting that even these seemingly unimportant tokens may hold valuable information or operational significance.

## 5 Discussion

This work addresses a design decision present in every T2I model that has remained largely unexplored: the choice to include padding tokens during both textual encoding and the diffusion process. As more studies begin integrating large language models (LLMs) into T2I pipelines using techniques like fine-tuning, LoRA, or adapters, the role of padding tokens becomes increasingly crucial. Training these models with padding tokens could influence a wide range of methods that assume subject information is encoded in specific tokens (Chefer et al., 2023; Rassin et al., 2023; Hertz et al., 2023; Gal et al., 2022), potentially altering their implementation when padding tokens carry significant semantic information. This factor should be carefully considered when deciding whether to train with or ignore padding tokens.

Furthermore, future research could explore how incorporating padding tokens into training might provide computational advantages in more integrated, end-to-end architectures, potentially allowing models to dynamically allocate resources by adjusting the use of padding tokens as needed.

## Limitations

While we have studied multiple T2I models representing several architectures, our work did not cover the vast space in this area. Our prompt selection offers some variety, but it may not capture all edge cases, potentially overlooking cases where padding tokens are used differently. Additionally, although we rely on widely used metrics like CLIP Score and KID for evaluation, these may not capture all nuances of image quality.

## Ethical Considerations

In developing our code, we used both Copilot and GPT-4o, but carefully reviewed each line to ensure it aligned with our intended implementation. For writing and rephrasing improvements, we used Wordtune and GPT-4o. Every generated suggestion was carefully reviewed and adjusted to ensure our original intent remained intact.

## References

Andrea Banino, Jan Balaguer, and Charles Blundell. 2021. Pondernet: Learning to ponder. In *8th ICML Workshop on Automated Machine Learning (AutoML)*.

Samyadeep Basu, Nanxuan Zhao, Vlad I Morariu, Soheil Feizi, and Varun Manjunatha. 2024. Localizing and editing knowledge in text-to-image generative models. In *The Twelfth International Conference on Learning Representations*.

Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. 2018. Demystifying mmd gans. In *International Conference on Learning Representations*.

Mikhail S Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V Sapunov. 2020. Memory transformer. *arXiv preprint arXiv:2006.11527*.

Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. 2023. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics*, 42(4):1–10.

Hila Chefer, Oran Lang, Mor Geva, Volodymyr Polosukhin, Assaf Shocher, Inbar Mosseri, Lior Wolf, et al. 2024. The hidden language of diffusion models. In *The Twelfth International Conference on Learning Representations*.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. 2024. Vision transformers need registers. In *The Twelfth International Conference on Learning Representations*.

Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. 2024. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*.

Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*.

Yossi Gandelsman, Alexei A Efros, and Jacob Steinhardt. 2023. Interpreting CLIP's image representation via text-based decomposition. *arXiv preprint arXiv:2310.05916*.

Gabriel Goh, Nick Cammarata, Chelsea Voss, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. 2021. Multimodal neurons in artificial neural networks. *Distill*, 6(3):e30.

Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2024. Think before you speak: Training language models with pause tokens. In *The Twelfth International Conference on Learning Representations*.

Alex Graves. 2016. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*.

Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2023. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations, 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Kosuke Imai, Luke Keele, and Dustin Tingley. 2010. A general approach to causal mediation analysis. *Psychological methods*, 15(4):309.

Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. Revealing the dark secrets of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4365–4374, Hong Kong, China. Association for Computational Linguistics.

Jacob Pfau, William Merrill, and Samuel R Bowman. 2024. Let's think dot by dot: Hidden computation in transformer language models. *arXiv preprint arXiv:2404.15758*.

Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2024. SDXL: improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. 2023. Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map alignment. *arXiv preprint arXiv:2306.08877*.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10674–10685. IEEE.

Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. 2023. What the DAAM: Interpreting stable diffusion using cross attention. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5644–5659, Toronto, Canada. Association for Computational Linguistics.

Michael Toker, Hadas Orgad, Mor Ventura, Dana Arad, and Yonatan Belinkov. 2024. Diffusion lens: Interpreting text encoders in text-to-image pipelines. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 9713–9728. Association for Computational Linguistics.

Jesse Vig and Yonatan Belinkov. 2019. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.

Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Simas Sakenis, Jason Huang, Yaron Singer, and Stuart Shieber. 2020. Causal mediation analysis for interpreting neural nlp: The case of gender bias. *arXiv preprint arXiv:2004.12265*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2024. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. 2022. Scaling autoregressive models for content-rich text-to-image generation. *Trans. Mach. Learn. Res.*, 2022.

Eric Zelikman, Georges Raif Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah Goodman. 2024. Quiet-STar: Language models can teach themselves to think before speaking. In *First Conference on Language Modeling*.

Fred Zhang and Neel Nanda. 2024. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*.

Shihao Zhao, Shaozhe Hao, Bojia Zi, Huaizhe Xu, and Kwan-Yee K Wong. 2024. Bridging different language models and generative vision models for text-to-image generation. *arXiv preprint arXiv:2403.07860*.

## A  Data Creation

- We randomly selected 50 samples from each of the following categories in the Parti dataset:
    - Fine-grained Detail
    - Imagination
    - Simple Detail
    - Style and Format
    - Complex
    - Linguistic Structures
    - Perspective
    - Quantity

- For each category, we used the following prompt with GPT-4o: Create an alternative CSV with different prompts of similar style and complexity.

- For the categories Style and Format, and Simple Detail, we repeated this process twice, generating a total of 100 examples for each.

- In the end, we obtained 500 prompts overall.

## B    Attention Between Image and Text in Different Architectures

To condition the generation process on a textual prompt, T2I models typically employ an attention mechanism. There are two popular methods for achieving this: through cross-attention mechanism, used in models like Stable Diffusion 2 and Stable Diffusion XL, and *MM-DiT* (Esser et al., 2024) blocks, found in models such as FLUX and Stable Diffusion 3. In the cross-attention mechanism, image patches are projected into query vectors $Q$ while text tokens are projected into key and value vectors $K$ and $V$. Essentially, each image patch draws information from the text tokens based on the attention map $A$:

$$A = softmax(QK^\top/\sqrt{d_k}), \qquad (4)$$

where $d_k$ represents the dimensionality of the key vectors. It is important to note that only the image patches extract information from the text tokens, while the text tokens remain constant throughout the computation process. Alternatively, the *MM-DiT* blocks implement a self-attention mechanism where both the image patches and text tokens are concatenated into a single set and then projected into $Q$, $K$ and $V$ vectors. In this formulation, both the image and text draw information from each other, using the following attention map:

$$A = softmax([Q_{txt}, Q_{img}][K_{txt}, K_{img}]^\top/\sqrt{d_k}), \qquad (5)$$

where $Q_{txt}$, $K_{txt}$ are the text query and key vectors, and $Q_{img}$, $K_{img}$ are the image query and key vectors. Here, both the image patches and text tokens are updated after the operation.

## C    Models

The models with frozen text encoders are:

1. *Stable Diffusion 2* employs a single frozen CLIP-based text encoder.

2. *Stable Diffusion 3* utilizes a combination of two frozen CLIP text encoders along with a frozen T5 encoder.

3. *FLUX* utilizes a frozen T5 text encoder and CLIP encoder. A key distinction between FLUX and Stable Diffusion models is that the latter incorporates a transformer architecture with self-attention to both the image and text latent representations in the diffusion process. This allows the diffusion model to modify text representations dynamically during the diffusion.

The models with trained text encoders:

1. *LDM* uses a BERT text encoder, which is trained jointly with the diffusion model on the image generation task.

2. *Lavi-Bridge* employs a LLaMA that is trained jointly with the diffusion model on the image generation task.

## D    Technical Details

All experiments were conducted using NVIDIA A100 GPUs with 8 cores, ensuring high computational performance and efficiency for our model evaluations. The total computational time across all experiments amounted to approximately 200 GPU hours.

## E    Qualitative Examples

The following figures provide visual examples illustrating the impact of padding tokens in the T2I pipeline, highlighting some key findings from our analysis.

Figure 10: Additional examples of images generated from different segments of the input prompt using IDP. Description of each column, from left to right: (1) An image generated using the full prompt (both prompt tokens and padding tokens encoded together), (2) An image generated using only the prompt tokens and clean padding tokens that were not encoded with the prompt, (3) An image generated using only the padding tokens encoded with the prompt, while the prompt tokens were replaced with clean pad tokens. See Figure 8 for further technical details.

Figure 11: Examples of images generated from the same prompts with maximum padding and without padding in Stable Diffusion XL and FLUX. Images generated by Stable Diffusion XL maintain consistent quality, while produced by FLUX without padding often miss key details. For example, given the prompt *"a compass beside a feather,"* images with padding typically include textured paper with text or a manuscript. In contrast, for the prompt *"a boy visiting a zoo,"* images generated without padding result in vague animal shapes (first column) or hybrids, such as a mix between a giraffe and a horse (third image). However, adding padding leads to more visually coherent animals.

# F Complementary results

| Model | Clean Pads | Full | Pads | Prompt |
|---|---|---|---|---|
| flux-schnell | 0.039 | 0.037 | 0.036 | 0.036 |
| ldm | 0.033 | 0.037 | 0.043 | 0.042 |
| LLaMA unet | 0.034 | 0.035 | 0.034 | 0.041 |
| stable diffusion 2 | 0.037 | 0.033 | 0.037 | 0.034 |
| stable diffusion 3 | 0.039 | 0.035 | 0.046 | 0.036 |
| stable diffusion XL | 0.023 | 0.036 | 0.043 | 0.039 |

Table 4: Calculated Standard Deviation of CLIP Scores for each Model and different text encoder interventions.