

# Follow the Flow: On Information Flow Across Textual Tokens in Text-to-Image Models

\*Guy Kaplan<sup>◇</sup>, \*Michael Toker<sup>♣</sup>, Yuval Reif<sup>◇</sup>, Yonatan Belinkov<sup>♣</sup>, Roy Schwartz<sup>◇</sup>

<sup>◇</sup>Hebrew University of Jerusalem  
{guy.kaplan3,yuval.reif,roy.schwartz1}@mail.huji.ac.il

<sup>♣</sup>Technion – Israel Institute of Technology  
{tok,belinkov}@campus.technion.ac.il

## Abstract

Text-to-Image (T2I) models often suffer from issues such as semantic leakage, incorrect feature binding, and omissions of key concepts in the generated image. This work studies these phenomena by looking into the role of information flow between textual token representations. To this end, we generate images by applying the diffusion component on a subset of contextual token representations in a given prompt and observe several interesting phenomena. First, in many cases, a word or multiword expression is fully represented by one or two tokens, while other tokens are redundant. For example, in “San Francisco’s Golden Gate Bridge”, the token “gate” alone captures the full expression. We demonstrate the redundancy of these tokens by removing them after textual encoding and generating an image from the resulting representation. Surprisingly, we find that this process not only maintains image generation performance, but also reduces errors by 21% compared to standard generation. We then show that information can also flow between *different expressions* in a sentence, which often leads to *semantic leakage*. Based on this observation, we propose a simple, training-free method to mitigate semantic leakage: replacing the leaked item’s representation after the textual encoding with its uncontextualized representation. Remarkably, this simple approach reduces semantic leakage by 85%. Overall, our work provides a comprehensive analysis of information flow across textual tokens in T2I models, offering both novel insights and practical benefits.<sup>1</sup>

## 1 Introduction

Text-to-image (T2I) models consist of two main components: a text encoder and a diffusion model (Ho et al., 2020; Song & Ermon, 2019). The former processes the user’s prompt, transforming it into a representation that guides the latter in generating the image. While T2I models are powerful and widely used, misalignment issues frequently arise. Generated images often fail to capture key concepts from the user’s prompt, leading to catastrophic neglect (Chefer et al., 2023a), semantic leakage of properties between entities (Rassin et al., 2022), or incorrectly bind different concepts (Huang et al., 2023; Rassin et al., 2022). Prior work has attempted to address these issues by modifying the cross-attention mechanism between the diffusion model and the encoded textual tokens (Rassin et al., 2023; Chefer et al., 2023a; Dahary et al., 2024), under the implicit assumption that the text encoder accurately captures the user’s intent and that each textual token primarily represents the item it is intended to convey.

In this work, we challenge these assumptions and investigate how information flows within tokens during the textual encoding process, and what information is encoded in

<sup>\*</sup>Equal contribution

<sup>1</sup>We release our code at <https://github.com/tokeron/lens>

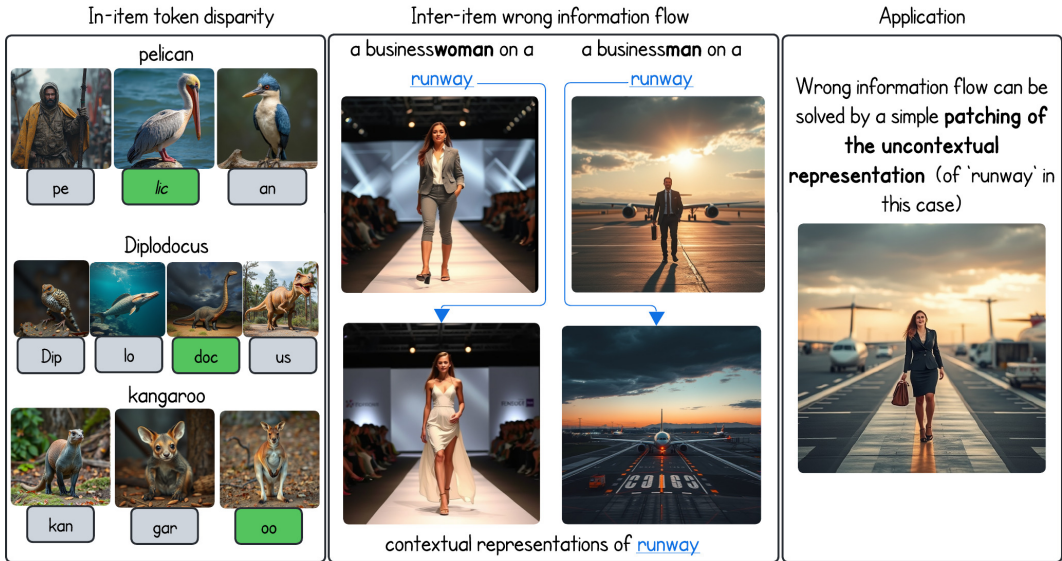


Figure 1: **Our main findings.** Left: Information within a lexical item is unevenly distributed across its token representations, with a few carrying most of the meaning. Middle: Context may cause incorrect information flow between items, distorting their representations and leading to unintended images. In this example, the representation of “runway” is distorted in the context of “a businesswoman”. Right: We propose a simple solution—encode the leaked item (“runway” in this case) separately and patch it into the representation of the prompt after textual encoding.

each textual token representation.<sup>2</sup> To this end, we analyze lexical items within prompts, where each lexical item is a single word or multiword expression that functions as a unit of meaning. Our analysis focuses on two complementary aspects. First, *in-item information flow*—how different sub-word token representations of the same lexical item encode and distribute information. For example, do all the token representations of the item “pelican” (“pe”, “lic”, and “an”) encode this lexical item (Fig. 1, top left). Second, *inter-item information flow*—how information flows between different lexical items (Feng et al., 2024; Feng & Steinhardt, 2024). For example, given the prompt “A businesswoman on a runway”, does the representation of “runway” change in the context of the word “businesswoman” (Fig. 1, bottom middle).

For *in-item information flow* (Section 4), we adapt a causal intervention method (Toker et al., 2025; see Fig. 2) to generate images from individual token representations. Specifically, we first encode the entire sentence, and then generate an image from each individual contextualized token representation. The resulting image reflects the information that the diffusion model can reconstruct when conditioned on that single token’s contextual representation. Our experiments with FLUX (black-forest labs, 2024) show that, in most cases, each lexical item is fully represented by only one or two tokens, while its remaining tokens might be considered redundant (Section 4.2). Finally, we demonstrate that our findings can be used to improve image–text alignment in image generation. To this end, we mask all the *non-representative tokens* (tokens that do not represent the lexical item) just before the diffusion, and generate images directly from the remaining representations. We find that this operation alone reduces image generation error by 21% compared to standard image generation.

For *inter-item information flow* (Section 5), we adapt the same causal intervention method—this time to assess the information encoded in one lexical item with respect to others in the prompt. For instance, in the sentence “A pool and a table”, we examine the information flow between the items “pool” and “table”. We find that in 89% of cases, each

<sup>2</sup>Throughout this work, we study token representations at the output of the text encoder. For brevity, we sometimes omit the word “representation”.

lexical item’s information is preserved and not influenced by others. However, in some cases, information flows between tokens. Interestingly, this flow can occur between lexical items that should not influence each other, for example, “table” influencing the representation of “pool” in the sentence above, changing pool’s representation to a “pool table”. We show that this leakage can be mitigated by patching the leaked lexical item’s representation with an uncontextualized representation of this item. Notably, this simple intervention prevents leakage in 85% of cases.

In summary, we analyze information flow in T2I models and uncover two key findings: (1) lexical items are often represented by just one or two tokens, and removing *non-representative tokens* can improve generation; (2) while different lexical items usually do not influence each other’s representations, when they do, it can happen between items that should not influence each other—leading to *semantic leakage*. Finally, based on our observation, we propose a simple method to mitigate semantic leakage, leading to a 85% reduction in leakage errors.

## 2 Related work

### 2.1 Interpretability

**Interpretability in multimodal models.** Previous work on interpretability in T2I models analyzed the text encoder. [Chefer et al. \(2023b\)](#) examined the concepts encoded in the latent space of CLIP ([Radford et al., 2021](#)) and how different concepts relate to each other. [Toker et al. \(2024\)](#) studied the progression of textual representations through the layers of the text encoder, using the diffusion model as a lens into the encoding process. Another line of work explored the connection between the text encoder and the image generator ([Tang et al., 2023](#)). In this work we investigate how information flows between tokens within a single layer. We focus on the final layer, as it directly conditions the diffusion process.

**Information flow in LLMs.** Interpretability research has examined how token information propagates via attention in LLMs ([Vig & Belinkov, 2019](#); [Clark et al., 2019](#)). However, raw attention weights can be misleading ([Pruthi et al., 2020](#); [Jain & Wallace, 2019](#)). To address this, methods such as Attention Rollout ([Abnar & Zuidema, 2020](#)) were developed to more reliably measure information flow through attention. Other studies have employed Patchscopes ([Ghandeharioun et al., 2024](#)) and the logit lens ([nostalgebraist, 2020](#)) to interpret the encoding processes in LLMs. These works demonstrate that modern language models go far beyond simple subword processing. For instance, [Kaplan et al. \(2025\)](#) showed that models perform an implicit detokenization process, effectively fusing subword tokens into coherent word-level representations. [Feucht et al. \(2024\)](#) reported an “erasure” effect, in which the contributions of earlier tokens are overwritten by later ones—suggesting the emergence of an internal lexicon. Another prominent direction is probing ([Adi et al., 2016](#); [Liu et al., 2019](#)). For example, [Zhang & Bowman \(2018\)](#); [Brunner et al. \(2019\)](#) trained probes to classify, from a token’s intermediate representation, the identity of tokens a fixed number of positions away. However, these methods are not fully reliable, as probes can exploit spurious correlations ([Belinkov, 2022](#)). Moreover, these works typically assess whether information is present or flows between tokens, but not whether the next step in the pipeline (e.g., the LM head or decoder) actually uses this information. In our analysis, we focus on intervention methods that test whether the next component in the image generation pipeline can retrieve and use the relevant information.

### 2.2 Challenges and solutions in multimodal models.

The recent rise of T2I models has introduced several challenges, prompting various approaches to address them. This work shows that a better understanding of information flow between tokens can help mitigate many of these challenges.

**Semantic leakage.** One notable challenge in T2I models is *semantic leakage* ([Rassin et al., 2022](#)). For example, given the prompt “a bat flying in a baseball stadium”, the generated image might include a baseball bat instead of a bat (the animal). A similar effect has been

observed in language models, where the context influences generation even when it conceptually should not (Gonen et al., 2024).

**Attribute binding.** Another issue is *attribute binding*, which concerns the correct association of attributes with entities mentioned in the prompt. For example, in the prompt “a yellow flamingo and a pink sunflower”, models must correctly bind “yellow” to “flamingo” and “pink” to “sunflower”. Several approaches have been proposed to address these challenges, including performing an optimization in the latent space to produce an attention pattern that is better aligned with the syntactic structure of the sentence (Rassin et al., 2023), and merging tokens to concentrate information into a single representation (Hu et al., 2024).

**Catastrophic neglect.** A further prominent challenge is *catastrophic neglect*, where entities mentioned in the prompt are completely omitted from the generated image. Many works have attempted to mitigate this, one of the first being *Attend and Excite* (Chefer et al., 2023a), which optimizes the image latents to increase attention on the neglected entity tokens.

### 3 Methodology

#### 3.1 Intervention on the text encoder

Our goal in this paper is to evaluate the effect of information flow between textual token representations, which later condition the diffusion process. We are interested in analyzing how information is distributed across tokens within a lexical item (e.g., do all the tokens in “San Francisco’s Golden Gate Bridge” represent that item, or is the information more diffusely spread?). To do so, we must be able to interpret the information encoded in individual tokens. We are also interested in measuring how different lexical items influence one another. To this end, we seek to isolate a subset of tokens and evaluate their joint representation. For these purposes, we adjust *text intervention in diffusion models* (Toker et al., 2025), allowing us to generate images from arbitrary subsets of tokens within the encoded user prompt.

Given a prompt with  $N$  tokens,  $t_1, t_2, \dots, t_N$ , our goal is to isolate and interpret the information encoded by a subset of these tokens. Let  $S \subset \{1, \dots, N\}$  be the index set of selected tokens, where  $1 \leq |S| \leq N$ . We begin by encoding the full prompt using the text encoder  $E$ , yielding the final hidden states  $e_1, \dots, e_N$ . Separately, we encode a sequence consisting entirely of pad tokens to obtain pad embeddings  $p_1, \dots, p_N$ . We then construct a *patched prompt* by replacing all hidden states outside  $S$  with the corresponding pad embeddings:

$$\tilde{t}_i = \begin{cases} e_i & \text{if } i \in S, \\ p_i & \text{otherwise} \end{cases} \quad \text{for } i = 1, \dots, N.$$

The patched sequence  $\tilde{t}_1, \dots, \tilde{t}_N$  is then used to guide the diffusion model. Generating an image from this patched representation allows us to visualize and isolate the individual contributions of the selected tokens as interpreted by the diffusion model. To evaluate the information in the generated images, we use a vision-language model (VLM): For the *inter-item information flow* experiments, we assess whether the generated image represents the full lexical item, while for the *in-item information flow* experiments, we measure whether information from other lexical items is present in the generated image. See Fig. 2 for an illustration of the method applied to interpret individual subtoken representations.

#### 3.2 Experimental setup

**Models.** We experiment with FLUX-Schnell (black-forest labs, 2024)—a recent state-of-the-art T2I model. We also report similar findings on FLUX-dev in Appendix A.5.1. FLUX models employ T5-XXL (Raffel et al., 2019) as their text encoder, enabling bidirectional information flow between tokens. In Appendix A.5.2 we also include another model, SDXL-Turbo (Sauer et al., 2023), which uses a different text encoder—CLIP (Radford et al., 2021). Unlike T5-xxl, CLIP is a causal language model. See Appendix A.5.2 for the analysis of SDXL-Turbo.

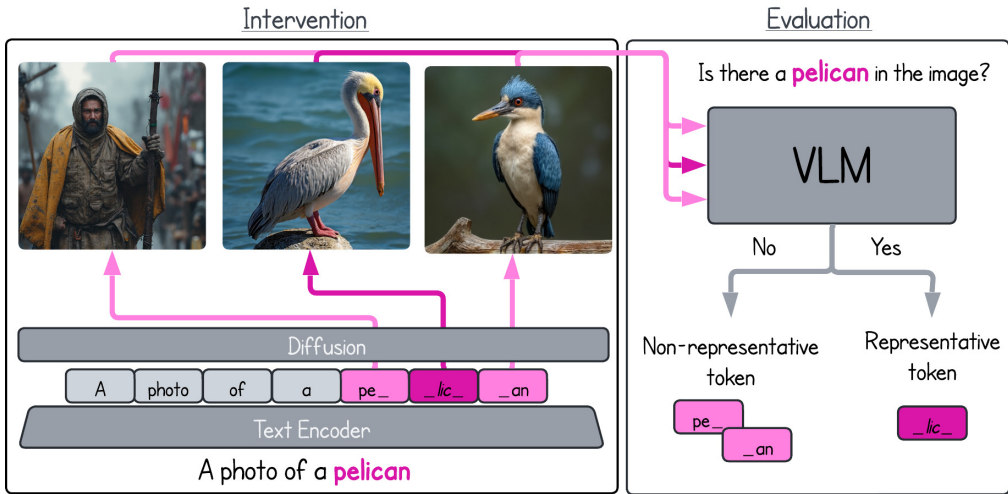


Figure 2: **Evaluating in-item information flow:** Our proposed framework interprets the information flow within a lexical item. We generate images from each token comprising the lexical item (left-hand side) and analyze the resulting images using a VLM (right-hand side). In this example, we interpret the different tokens composing the item “**pelican**” and find that only one of them, “**lic**”, represents the concept of a pelican, while the other two tokens (“**pe**”, and “**an**”) do not.

**Data.** We use a subset of prompts from DrawBench (Saharia et al., 2022) and Parti (Yu et al., 2022b). In total, we use 1,053 prompts. For each prompt, we generate five images using different random seeds. To conduct the analysis in our paper, we first find the lexical items in each prompt. We use GPT-4o (OpenAI et al., 2024) to extract the lexical items, resulting in a total of 4,864 unique lexical items. See Appendix A.1 for the exact prompt we used and further technical details. We then use spaCy (Honnibal et al., 2020) to determine the part of speech of each lexical item, and retain only nouns, proper nouns, and adjectives, as these are typically concrete and can be identified in their visual representation. We end up with 3,891 unique lexical items and use them to evaluate both *inter-item information flow* and *in-item information flow*.

**Evaluation.** To evaluate the content of generated images, we also employ Qwen2-VL-72B-Instruct (Wang et al., 2024), a model with strong general vision capabilities. We restrict our evaluation to binary (yes/no) questions. See Appendix A.1 for the exact prompt used and additional details.

## 4 In-item information flow

### 4.1 How is information distributed across tokens?

In this section, we explore how information is distributed across the token representations within a lexical item. We aim to answer two main questions: Do all tokens encode the same information? And is the information evenly distributed across tokens or concentrated in specific token representations? These questions are central, as many T2I applications (Chefer et al., 2023a; Rassin et al., 2023; Dahary et al., 2024) treat all tokens of a lexical item equally. Discovering a non-uniform distribution of information across tokens could improve application effectiveness by focusing on the most informative tokens.

We apply our intervention method (Section 3.1) to the hidden representation of each token in the prompt. We then generate an image conditioned on each token’s contextualized representation. We use Qwen2-VL to assess whether an image represents the lexical item it is part of. We repeat this analysis for each lexical item, for each prompt in our combined dataset (Section 3.2). The results show that in 89% of the cases, at least one of the tokens



Figure 3: Examples illustrating the effect of removing *non-representative tokens*. **Top row:** Images generated after removing *non-representative tokens* (Representative tokens are shown in **bold**; non-representative tokens are in gray.). **Bottom row:** Images generated from the full prompt. **Left:** In most cases, removal results in no noticeable effect on the generation. **Right:** In some cases, removal improves alignment with the prompt.

fully represents the lexical item. Out of the remaining 11%, about 85% cannot be generated even from the full lexical item, suggesting that the model is unfamiliar with the concept. See Appendix A.4 for a deeper analysis of these cases. The remaining (1.3% of all lexical items) can be generated from the full lexical item, indicating that in these cases, the lexical item’s representation is distributed across multiple tokens.

Focusing on the cases where at least one token represents the lexical item, we examine the number of *representative tokens* (those that result in an image that represents the item) and *non-representative tokens* (those that do not). We compare the proportion of *representative tokens* across different lexical items of different lengths. Our results (Fig. 7 in Appendix A.3) show that typically one or two tokens represent the concept, while the remaining tokens are *non-representative tokens*. As the lexical item becomes longer, the number of *non-representative tokens* increases. Focusing on these *non-representative tokens*, we find that they account for 52% of the tokens within the lexical items represented by least two tokens. In the next section, we examine the effect of removing *non-representative tokens*.

#### 4.2 Are all prompt tokens necessary?

Our analysis reveals that 52% of lexical item tokens are non-representative. In this section, we examine whether these tokens have any effect on image generation. To answer this question we apply our intervention method for each prompt, generating an image after masking all the *non-representative tokens*. We then use Qwen-VL to measure whether the image aligns with the prompt and compare it to the image generated without this intervention.

**Non-representative tokens are redundant.** Our results show that removing *non-representative tokens* generally does not harm the generation (see Fig. 3, left-hand side). Of the cases where the original generation was aligned with the prompt, in 98% of the cases, the generated image still aligns with the prompt, suggesting that these tokens are largely redundant. Moreover, in cases where the original image failed to align with the prompt, we observe a 21% reduction in error rate (see Fig. 3, right-hand side). We attribute this improvement to the model being guided more by the *representative tokens* after removing the *non-representative ones*. See Table 2 in Appendix A.3 for full results.

**Can we identify redundant tokens without generating their images?** Identifying which tokens are representative can improve image generation by removing redundant ones as shown in Fig. 3 (right). Moreover, many existing T2I methods operate at the token level, typically using all of a lexical item’s tokens to represent the item. (Chefer et al., 2023a;

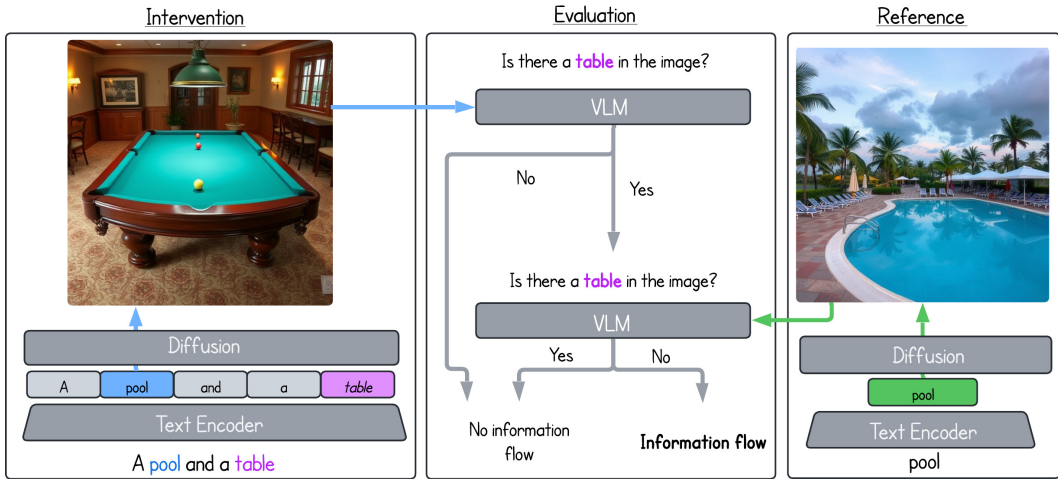


Figure 4: **Evaluating inter-item information flow:** Our proposed framework to interpret the information flow between lexical items in the prompt. For each lexical item, we generate an image from its contextual representations (left), and from its uncontextualized representation (right), and analyze the generated images using a VLM (middle). In this example, we interpret the item “pool” and assess whether it is influenced by the item “table”. To do so, we ask a VLM whether the image generated from the token “pool” contains a “table”. To ensure this is the result of information flow, and not a natural correlation between a pool and tables, we generate the item “pool” without context (right-hand side), and verify that “table” is not present in this image. If this is the case, we conclude that there was information flow from “table” to “pool”.

Dahary et al., 2024). Refining these methods to operate only on *representative tokens* could lead to improvements without requiring any additional changes.

However, generating an image for each token is computationally expensive. To this end, we propose two ways to estimate token importance without image generation. First, we calculate the edit distance between a lexical item (e.g., “giraffe”, which is tokenized as “gir”, “a” and “ffe”, and a given token (e.g., ‘a’). We find that higher edit distance indicates a lower likelihood that the token is representative (Pearson correlation of  $-0.44$ ). While this correlation is relatively strong, it is not accurate enough for practical use. Instead, we train a  $k$ -NN classifier ( $k=5$ , Euclidean distance) to predict token redundancy from encoded representations. We extract 6,966 tokens (4,864 unique lexical items) from our dataset, each labeled for whether the VLM labeled this token as representative. We find that this simple classifier achieves a precision of 92% in predicting token redundancy, highlighting its potential as an efficient tool for removing *non-representative tokens* while preserving the essential ones. Further details about the classifier can be found in Appendix A.1.

### 5 Inter-item information flow

User prompts are typically composed of multiple lexical items. Here we investigate the influence of lexical items on one another. For example, in the sentence “a pool next to a table”, is the representation of the item “pool” affected by the presence of the item “table”? More generally, given a prompt, made up of several lexical items, we aim to measure whether the presence in the prompt of one item affects the representation of another.

To this end, we conduct the following experiment: given a prompt, we isolate each lexical item one at a time. First, we encode the entire prompt. Next, for each item, we patch the representations of all other tokens using our method Section 3.1, and generate an image based on the modified sequence. Specifically, we define a subset  $S \subset \{1, \dots, N\}$  containing the token indices of the selected lexical item, construct the patched sequence  $\tilde{t}_1, \dots, \tilde{t}_N$ , and use it to guide the diffusion model. This process produces one image per lexical item

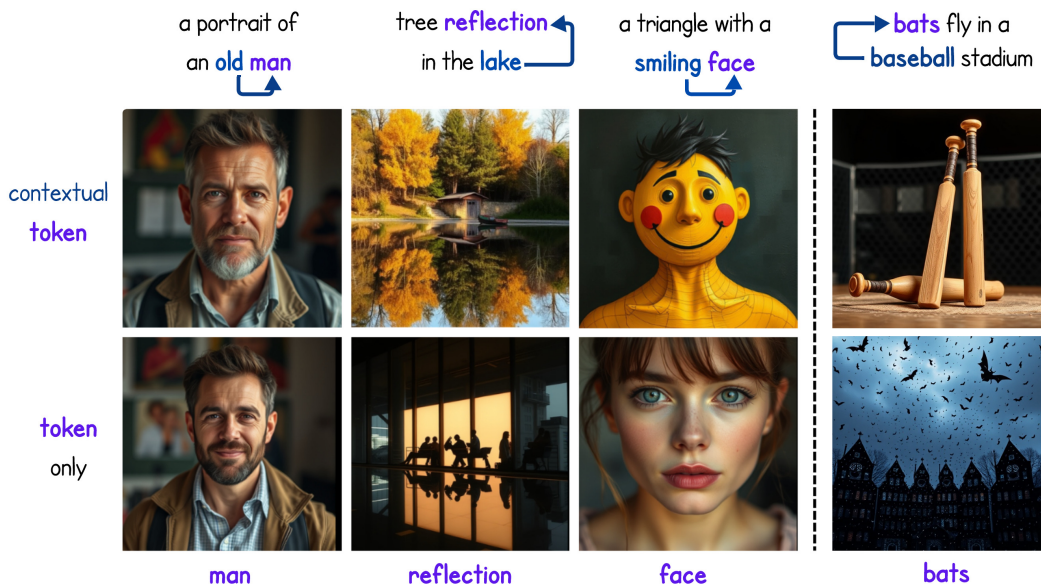


Figure 5: **Examples of information flow between items.** Top: Images generated from a **lexical item** encoded alongside **another item** that alters its representation. Bottom: Images generated from the uncontextualized representation of the same lexical item. The first three images (from the left) demonstrate correct information flow, while the last image (far right) demonstrates incorrect information flow.

in the prompt—each encoded in the full context of the prompt but with all other token representations masked. These images allow us to assess what information is encoded in each item’s representation.

We then generate an image for the uncontextualized version of each lexical item—encoding only that item in isolation, without the rest of the prompt. To assess whether a lexical item  $y$  influences another item  $x$ , we use Qwen2-VL to check if  $y$  is present in the image generated from  $x$ , both with and without context. We then compare the Qwen2-VL answers between the contextualized and uncontextualized versions. Refer to Appendix A.1 for more details regarding Qwen2-VL and the prompts we use. An influence is considered *True* if  $x$  did not exhibit  $y$  in the uncontextualized representation but did in the contextualized one. Confirming that  $y$  does not appear in the isolated representation of  $x$  is crucial, as we want to ensure that the presence of  $y$  in  $x$  is due to information flow that occurs during textual encoding, and not due to some underlying correlation between the items. This approach allows us to quantify the degree to which lexical items affect one another within a prompt. See Fig. 4 for a demonstration of this framework.

Our results show that in 89% of the cases, lexical items are not *influenced* by other items in the prompt. In the remaining 11%, one lexical item incorporates information about another item that appears in the same prompt; see Fig. 5 for examples. Interestingly, we find that when such influence occurs, it often involves items that frequently co-occur in natural language, and are therefore likely to appear together in the training data. For example, the item “pool” often shifts its representation from a swimming pool to a pool table when “table” is present in the context. This phenomenon occurs even when “table” is not used as a pool table, as in the sentence “A pool next to a table” (see Fig. 4).

To measure the extent of information flow between items that should not influence each other (as in the “pool” and “table” example), we divide the identified cases of flow between items into two categories. The first includes cases where the items are semantically linked, such as “black” and “bear” in the prompt “a black bear”. The second includes cases where the items are not semantically linked, such as “black” and “bear” in the sen-



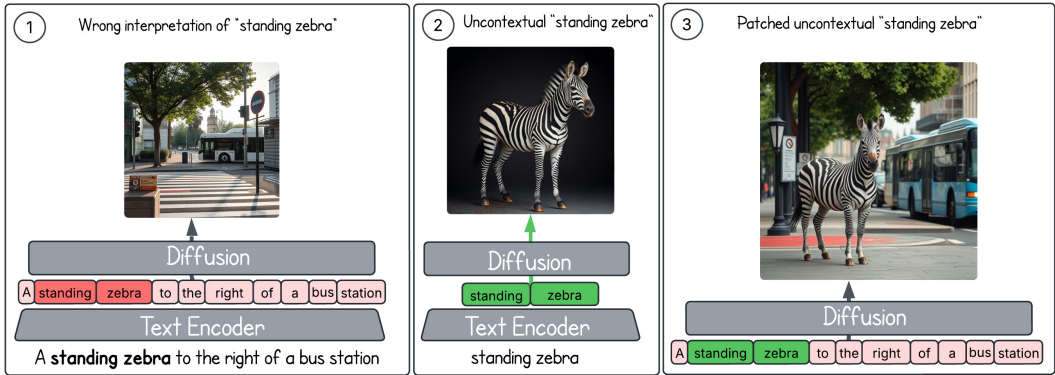


Figure 6: **Removing semantic leakage by replacing the contextually leaked concept representation.** (1) Regular generation produces an image showing a crosswalk to the right of a bus station. (2) Generation from the prompt “standing zebra”, without any context, results in the correct interpretation of the zebra as an animal. (3) Generation using the original prompt, but with the leaked concept “standing zebra” replaced by its **uncontextualized representation**, yields a correct image.

tence “a green bear by the black tree”. To achieve this, we use an LLM as a judge (details in Appendix A.1). Our analysis reveals that 31% out of the influenced instances exhibit unintentional information flow.

**Application.** Based on our findings, we present a method for mitigating semantic leakage. Our approach first encodes the prompt as usual, then separately encodes the item suspected of leakage. Finally, we replace the representation of the contextualized, leaked item in the original prompt’s encoding with its uncontextualized representation, and condition the diffusion process on this mixed representation (see Fig. 6).

To evaluate the effectiveness of our method, we adapt prompts from [Rassin et al. \(2022\)](#) as these prompts often induce semantic leakage. For example, given the prompt “a bat is flying over a baseball stadium”, models tend to generate a flying baseball bat rather than the animal, due to the strong association between “bat” and “baseball stadium”. Since the dataset contains only about 30 examples, we use GPT-4o to expand it to 110 prompts. See Appendix A.2 for details on the augmentation process, including the full list of prompts.

Out of the 110 prompts, 104 exhibit semantic leakage (as determined by manual human evaluation), where the image contains a concept not intended by the prompt. We then select the lexical items in each prompt that we suspect of having leaked, and apply our method on them. Remarkably, we achieve an 85% improvement on this set. See Fig. 8 in Appendix A.3 for additional qualitative examples.

## 6 Conclusions

In this study, we examined the information flow within lexical items and between different lexical items during the textual encoding of a prompt in T2I models. We introduced an intervention to evaluate how information is distributed across various token representations. For *in-item information flow*, our analysis revealed that information within a concept is not uniformly distributed; instead, it is primarily concentrated in a few specific tokens. For *inter-item information flow*, we found that information typically does not flow between different items. When it does, the flow is often unintended and may result in *semantic leakage*. Our experiments suggest that semantic leakage originates within the text encoder itself, challenging the common assumption that each token primarily captures its intended meaning. We showed that semantic leakage phenomenon can be largely mitigated by simply encoding the leaked item separately and patching it into the prompt representation before the diffusion.

While prior work has focused on modifying cross-attention to address challenges in text-to-image models, our results suggest that these methods should also account for information flow between tokens that occurs earlier, during the textual encoding.

## 7 Acknowledgments

This research was supported in part by the Israel Science Foundation (grant No. 448/20), NSF-BSF grant 2020793, an Azrieli Foundation Early Career Faculty Fellowship, an AI Alignment grant from Open Philanthropy, the European Union (ERC, Control-LM, 101165402), and by an Academic Gift by NVIDIA. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

## References

- Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4190–4197, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.385. URL <https://aclanthology.org/2020.acl-main.385/>.
- Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. *arXiv preprint arXiv:1608.04207*, 2016.
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, March 2022. doi: 10.1162/coli.a.00422. URL <https://aclanthology.org/2022.cl-1.7/>.
- black-forest labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. On identifiability in transformers. *arXiv preprint arXiv:1908.04211*, 2019.
- Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM transactions on Graphics (TOG)*, 42(4):1–10, 2023a.
- Hila Chefer, Oran Lang, Mor Geva, Volodymyr Polosukhin, Assaf Shocher, Michal Irani, Inbar Mosseri, and Lior Wolf. The hidden language of diffusion models. *arXiv preprint arXiv:2306.00966*, 2023b.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT’s attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes (eds.), *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 276–286, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828. URL <https://aclanthology.org/W19-4828/>.
- Omer Dahary, Or Patashnik, Kfir Aberman, and Daniel Cohen-Or. Be yourself: Bounded attention for multi-subject text-to-image generation. In *European Conference on Computer Vision*, pp. 432–448. Springer, 2024.
- Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context?, 2024. URL <https://arxiv.org/abs/2310.17191>.
- Jiahai Feng, Stuart Russell, and Jacob Steinhardt. Monitoring latent world states in language models with propositional probes, 2024. URL <https://arxiv.org/abs/2406.19501>.

- Sheridan Feucht, David Atkinson, Byron Wallace, and David Bau. Token erasure as a footprint of implicit vocabulary items in llms, 2024. URL <https://arxiv.org/abs/2406.20086>.
- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. Patchscopes: A unifying framework for inspecting hidden representations of language models. *ArXiv*, abs/2401.06102, 2024. URL <https://api.semanticscholar.org/CorpusID:266933130>.
- Hila Gonen, Terra Blevins, Alisa Liu, Luke Zettlemoyer, and Noah A Smith. Does liking yellow imply driving a school bus? semantic leakage in language models. *arXiv preprint arXiv:2408.06518*, 2024.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. *spacy: Industrial-strength natural language processing in python*, 2020. URL <https://spacy.io/>.
- Taihang Hu, Linxuan Li, Joost van de Weijer, Hongcheng Gao, Fahad Shahbaz Khan, Jian Yang, Ming-Ming Cheng, Kai Wang, and Yaxing Wang. Token merging for training-free semantic binding in text-to-image synthesis. *Advances in Neural Information Processing Systems*, 37:137646–137672, 2024.
- Kaiyi Huang, Chengqi Duan, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench++: An enhanced and comprehensive benchmark for compositional text-to-image generation. *IEEE transactions on pattern analysis and machine intelligence*, PP, 2023. URL <https://api.semanticscholar.org/CorpusID:259847295>.
- Sarthak Jain and Byron C. Wallace. Attention is not Explanation. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 3543–3556, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1357. URL <https://aclanthology.org/N19-1357/>.
- Guy Kaplan, Matanel Oren, Yuval Reif, and Roy Schwartz. From tokens to words: On the inner lexicon of llms, 2025. URL <https://arxiv.org/abs/2410.05864>.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. Linguistic knowledge and transferability of contextual representations. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 1073–1094, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1112. URL <https://aclanthology.org/N19-1112/>.
- nostalgebraist. Interpreting GPT: The logit lens. *lesswrong*, 2020., 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,

Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.

Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C. Lipton. Learning to deceive with attention-based explanations. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4782–4793, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.432. URL <https://aclanthology.org/2020.acl-main.432/>.

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

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2019. URL <https://arxiv.org/abs/1910.10683>.

Royi Rassin, Shauli Ravfogel, and Yoav Goldberg. DALLE-2 is seeing double: Flaws in word-to-concept mapping in Text2Image models. In Jasmijn Bastings, Yonatan Belinkov,

- Yanai Elazar, Dieuwke Hupkes, Naomi Saphra, and Sarah Wiegrefe (eds.), Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP, pp. 335–345, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.blackboxnlp-1.28. URL <https://aclanthology.org/2022.blackboxnlp-1.28/>.
- Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map alignment. Advances in Neural Information Processing Systems, 36:3536–3559, 2023.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems, 35:36479–36494, 2022.
- Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation, 2023. URL <https://arxiv.org/abs/2311.17042>.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Advances in neural information processing systems, 32, 2019.
- Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. What the DAAM: Interpreting stable diffusion using cross attention. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 5644–5659, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.310. URL <https://aclanthology.org/2023.acl-long.310/>.
- Michael Toker, Hadas Orgad, Mor Ventura, Dana Arad, and Yonatan Belinkov. Diffusion lens: Interpreting text encoders in text-to-image pipelines. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 9713–9728, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.524. URL <https://aclanthology.org/2024.acl-long.524/>.
- Michael Toker, Ido Galil, Hadas Orgad, Rinon Gal, Yoad Tewel, Gal Chechik, and Yonatan Belinkov. Padding tone: A mechanistic analysis of padding tokens in t2i models, 2025. URL <https://arxiv.org/abs/2501.06751>.
- Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes (eds.), Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pp. 63–76, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4808. URL <https://aclanthology.org/W19-4808/>.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024.
- 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. Scaling autoregressive models for content-rich text-to-image generation, 2022a. URL <https://arxiv.org/abs/2206.10789>.
- Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, et al. Scaling autoregressive models for content-rich text-to-image generation. arXiv preprint arXiv:2206.10789, 2(3):5, 2022b.

Kelly Zhang and Samuel Bowman. Language modeling teaches you more than translation does: Lessons learned through auxiliary syntactic task analysis. In Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pp. 359–361, 2018.

## A Appendix

### A.1 Technical details

**Lexical item classification** We define a *lexical item* as either a single word or a compound expression of multiple words that, in context, conveys a unified semantic meaning. A compound expression is treated as a single item when its words form a fixed lexical unit with cohesive semantics rather than merely exhibiting a modifier–head relationship. For example, while “broken mirror” describes a mirror’s state, expressions like “hot air balloon” or “teddy bear” denote entities with distinct identities. Similarly, although phrases such as “identical twins” or “baseball bat” might be interpreted as separate concepts, conventional usage supports their treatment as unified entities. We employ the reasoning model `03-mini-high` as a classifier to tag multi-word lexical items in both the target prompts and the dataset. The model returns a list of identified multi-word expressions, while the remaining untagged words are treated as individual lexical items.

**Redundant token classification** We propose a probing classifier to predict whether a token is redundant (i.e., *non-representative* of its lexical item) using solely its encoded representation, without the need to generate an image. For this purpose, we extracted the 6,966 tokens corresponding to 4,864 unique lexical items from our dataset. Each token was annotated with a binary label indicating whether it represents the lexical item by the VLM (see A.1 for more details).

We split the data into training and validation sets using an 80-20 ratio. A  $k$ -nearest neighbors (k-NN) classifier with  $k = 5$  and Euclidean distance as the similarity measure was then applied to predict token redundancy directly from the encoded representations. The results on the evaluation set are presented in Table 1. The high precision indicates that, in practical settings, one can remove tokens predicted as redundant with a high degree of certainty that *representative tokens* will not be inadvertently discarded.

Table 1: Performance of k-NN classifier ( $k = 5$ ) for predicting token redundancy

Metric	Accuracy	Precision	Recall	F1-score
Score	0.82	0.92	0.74	0.82

**Evaluation visual generations.** We evaluate whether an image matches a textual description using Qwen2-VL-72B-Instruct (Wang et al., 2024). The following prompt is used:

“In Yes, No and maybe. Does every image match one of those descriptions: (description string)? Answer Yes if all images match or relate to at least one description, Maybe if only some match, otherwise No.”

Here, the *textual description* can be either a single lexical item or a complete textual prompt.

**Evaluating generated textual descriptions.** We employ gpt4o (OpenAI et al., 2024) to evaluate the textual interpretations produced by Patchscopes. We use the following prompt:

“In Yes, No and Maybe. Does every image match the description: {Patchscopes.description} ? Answer Yes if all images match or relate to the description, Maybe if only some match, otherwise No.”

**Classifying the relations between items.** We enhance our leakage validation by distinguishing between cases where two lexical items exhibiting semantic leakage are perceptually bound together—for example, “old” and “man” in the prompt “a portrait of an old man”—and cases where they are not as “cone hat” and “eating” in the prompt “A person searing a cone hat is eating” (see Fig. 5). To achieve this, we use a large language model (LLM) as a judge. Specifically, we employ the following prompt:

“In Yes or No: in this prompt: {input\_prompt}, are {item\_1} and {item\_2} perceptually bound together?”

We then filter out all cases where the lexical items are perceptually bound together and find that only 6.5% instances exhibit unintentional leakage.

### A.2 Data

**DrawBench** (Saharia et al., 2022): We include all categories except for “misspelling”, “rare words”, and “text”. Overall we end up with 134 prompts from DrawBench.

**Parti Prompts** (Yu et al., 2022a): We include all categories except “Style & Format”, “Writing & Symbols”, and “Arts”. Overall we end up with 923 prompts from Parti Prompts.

In total, we end up with 1,056 prompts.

**Extended dataset of leakage prompts** Our augmentation process incorporates two components. First, we generate variations of existing prompts from Rassin et al. (2022) (e.g., modifying ‘a gentleman with a bow in the forest’ to ‘a man wearing a bow in the jungle’). Second, we introduce novel prompts with potential semantic leakage. For these prompts, we applied a one-lexical item change test by generating an image from a similar prompt that substitutes the affected or leaked item with an alternative term (e.g., replacing ‘bishops’ with ‘cardinals’ or ‘checkers’ in ‘chess in 2 bishops playing chess’). This test ensures that minimal lexical modifications do not alter the intended semantic meaning while producing a different image due to semantic leakage from another item in the prompt (see the first two columns in Fig. 8 for few visual examples) . Together, these methods enrich the dataset and provide a robust framework for analyzing semantic leakage.

The full list of prompts is available in the following anonymous Git repository: <https://anonymous.4open.science/r/TokenRole-3E19>.

### A.3 Additional results

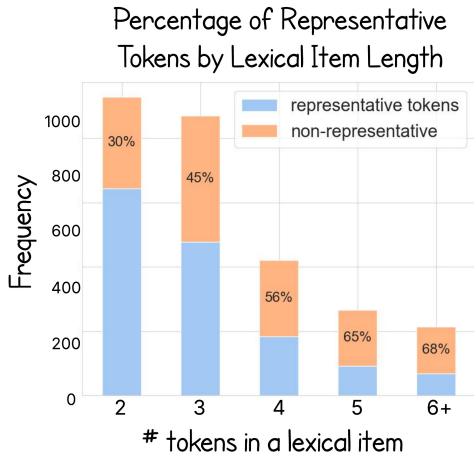


Figure 7: Distribution of representative vs. *non-representative tokens* across different lexical item lengths.

### A.4 Catastrophic negligence

*Catastrophic negligence* refers to cases where a concept or entity is completely absent from the generated image. Chefer et al. (2023a) attempted to address this issue by increasing attention to the neglected item, implicitly assuming that the model understands the concept.



Method	Success Rate
Regular Generation	84.1%
After Removing Redundant Tokens	87.6%

Table 2: Image-text alignment success rate (evaluated with Qwen-VL): comparison between standard generation and generation with redundant tokens removed

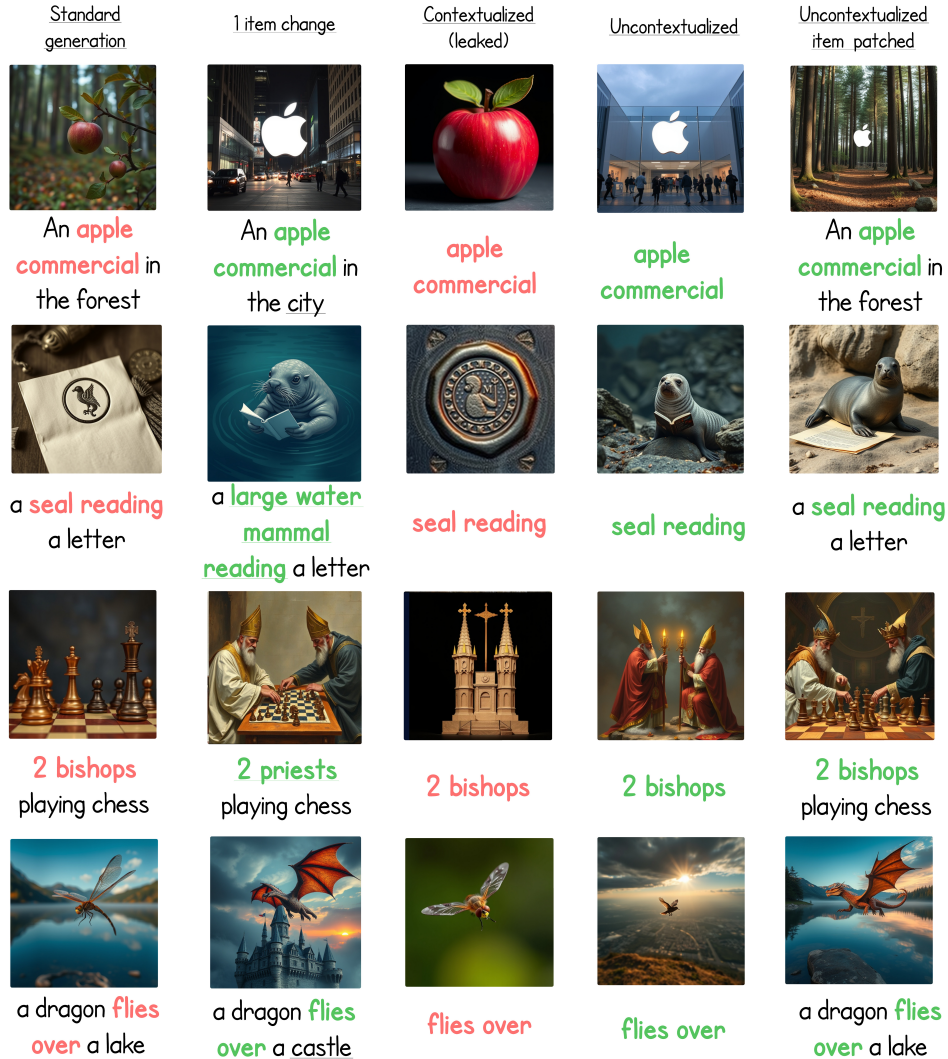


Figure 8: Examples from our semantic leakage method. **Left:** standard generation of leakage contained prompt. **Second:** generation using a one-lexical item change test as part of the dataset creation (a minimal substitution to verify that a slight lexical change yields a different image). **Third:** image from the contextual representation (misinterpreted item). **Forth:** image from the uncontextualized representation (correct interpretation). **Right:** final generation after patching the correct, uncontextualized representation into the prompt.

We find that in 7% of the cases, the item is not generated from any of its constituent tokens and is absent from the full prompt’s image. To better understand this phenomenon, we leverage Patchscopes (Ghandeharioun et al., 2024), a tool that interprets intermediate token representations by generating textual explanations using the decoder in encoder-decoder architectures. Using Patchscopes, we observe that in approximately 67% of these cases the tokens are explained coherently, indicating that the item is adequately represented in

the text encoder. For example, when we patch the encoded representation of the lexical item tuba into the prompt “describe X”, Patchscopes describes it as “a musical instrument played by blowing into the mouthpiece.” However, when we try to generate an image from its representation, the model fails (see Fig. 9), suggesting that the diffusion model may have encountered too few examples of the intended entity during training.

Additionally, we identify two abstract categories that consistently result in failure in image generation, while Patchscopes accurately describes the concept: negation and numbers. For instance, the lexical item “without” failed to generate the an image without the item in 100% of 100 observed instances. See Fig. 9.

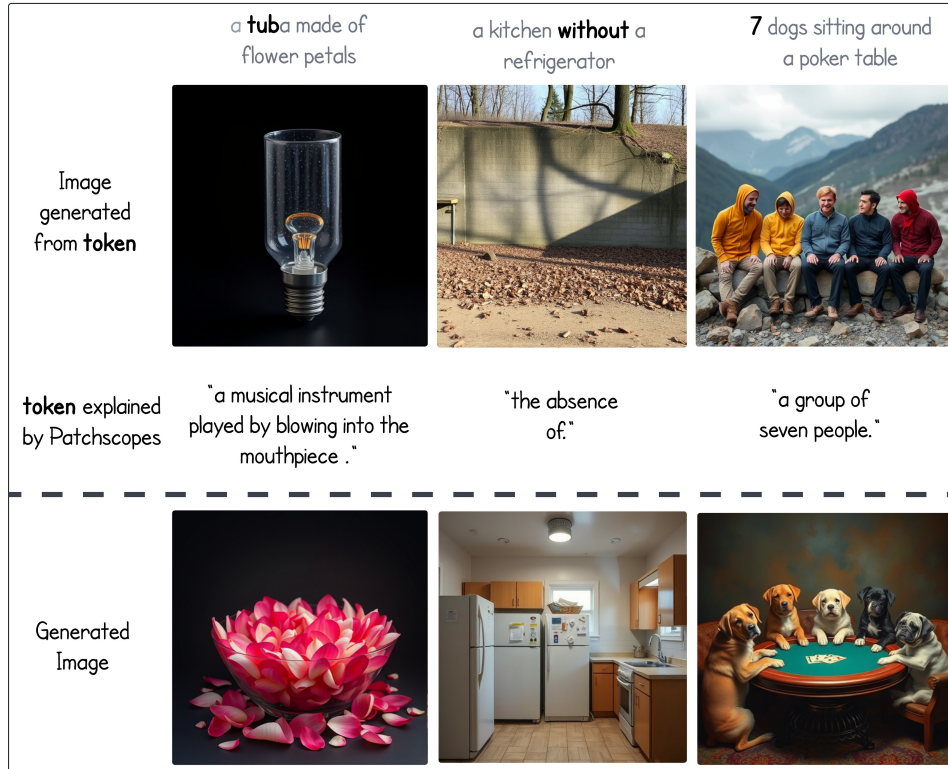


Figure 9: Cases where intervention-based interpretations fail, while *Patchscopes* accurately describe token meaning. Each column corresponds to a different prompt. The first row shows images generated from the contextual **bolded** token. The second row presents the same token explained by *Patchscopes*, and the third row shows the final image generated from the full prompt. **Left:** The image generated from the token “tub” (from “tuba”) depicts a lamp, whereas *Patchscopes* correctly identifies it as a musical instrument. **Middle:** The concept “without” fails to yield an image lacking a refrigerator, although *Patchscopes* correctly describes its semantics. **Right:** The concept “7” leads to an incorrect number of entities in the image, while *Patchscopes* correctly identifies the quantity.

## A.5 Additional models

### A.5.1 FLUX-Dev

In addition to our primary experiments with FLUX, we repeated all analyses using the Flux-dev variant. The redundant versus representative token experiments yielded similar trends, with 55% of tokens identified as representative and 45% as non-representative—values closely matching those observed with FLUX. Likewise, our inter-item flow experiments confirmed that information flow occurred in 11% of cases (and 3.1% miss intended leakage), reinforcing the overall patterns reported in the main text. Notably, while the aggregate trends are consistent across models, the specific lexical items

resolved can differ between FLUX-schnell and Flux-dev, indicating potential a slightly different inner-lexicon. These findings underscore the robustness of our approach while highlighting model-dependent nuances in token representation and information flow dynamics.

### A.5.2 SDXL-Turbo



Figure 10: Images generated from individual subtokens in SDXL-Turbo. We find that, in many cases, the representation of an item is not clearly reflected in any of its subtokens—for example, in the case of the token “skateboard.” Another interesting observation is that the last token of a lexical item often carries its representation, as seen in the “square” token of “times square.” We also observe that the EOT token incorporates information from the full prompt.

Our analysis reveals that SDXL-Turbo, which uses the CLIP text encoder, behaves markedly differently from FLUX, which relies on the encoder in the encoder-decoder T5-XXL. In SDXL-Turbo, the text encoder is a causal language model, meaning each token’s encoding is influenced only by its preceding tokens during the encoding process.

We repeated our *in-item information flow* experiments using SDXL-Turbo. Our first observation is that most generated images are either abstract or unrelated to the intended lexical items. According to our analysis, 55% of lexical items in SDXL-Turbo lack any representative token (compared to just 11% in FLUX). Moreover, when a representative token is present in CLIP, it is typically the final token of the lexical item (see Fig. 10). This is aligned with the unidirectional encoding of the model.

Another phenomenon we observe—consistent with CLIP’s training objective—is the unusually dominant role of the end-of-sequence (EOS) token. Images generated from the EOS token often encapsulate nearly the full semantic content of the prompt. In our evaluation, 62% of EOS-generated images matched the prompt ( compared to 73% when using the full prompt). We believe this also causes our intervention method to be less effective, since when we interpret a single token, we patch all other tokens, including the EOT token—which usually contains a lot of information—with tokens derived from an empty prompt (see Fig. 10).