

# AI's Blind Spots: Geographic Knowledge and Diversity Deficit in Generated Urban Scenario

Ciro Beneduce

Mobile and Social Computing Lab  
Bruno Kessler Foundation  
Trento, Italy  
cbeneduce@fbk.eu

Massimiliano Luca

Mobile and Social Computing Lab  
Bruno Kessler Foundation  
Trento, Italy  
mluca@fbk.eu

Bruno Lepri

Mobile and Social Computing Lab  
Bruno Kessler Foundation  
Trento, Italy  
lepri@fbk.eu

## Abstract

Image generation models are revolutionizing many domains, and urban analysis and design is no exception. While such models are widely adopted, there is a limited literature exploring their geographic knowledge, along with the biases they embed. In this work, we generated 150 synthetic images for each state in the USA and related capitals using FLUX 1 and Stable Diffusion 3.5, two state-of-the-art models for image generation. We embed each image using DINO-v2 ViT-S/14 and the Fréchet Inception Distances to measure the similarity between the generated images. We found that while these models have implicitly learned aspects of USA geography, if we prompt the models to generate an image for "United States" instead of specific cities or states, the models exhibit a strong representative bias toward metropolis-like areas, excluding rural states and smaller cities. In addition, we found that models systematically exhibit some entity-disambiguation issues with European-sounding names like Frankfort or Devon.

## CCS Concepts

- Computing methodologies → Image generation;
- Information systems → Geographic information systems;
- Social and professional topics → Bias, discrimination and fairness.

## Keywords

Diffusion models, geographic knowledge, geographic bias, urban imagery

## ACM Reference Format:

Ciro Beneduce, Massimiliano Luca, and Bruno Lepri. 2025. AI's Blind Spots: Geographic Knowledge and Diversity Deficit in Generated Urban Scenario . In . ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nmnnnnnnnnnnnnn>

## 1 Introduction

Generative Artificial Intelligence (AI) has rapidly emerged as a transformative technology, significantly impacting urban spatial analysis, scenario generation, movement prediction and other urban tasks [1, 16, 17]. When it comes to deal with images, diffusion

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference'17, Washington, DC, USA*

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM  
<https://doi.org/10.1145/nmnnnnnnnnnnnnn>



**Figure 1: Outputs from FLUX 1-schnell (top) and SD 3.5-L (bottom) for three prompts: "USA," "Boston," and "Juneau". The cities are correctly generated, showing the detailed geographic knowledge of the models. Despite that, the U.S. representations default to a stereotypical metropolitan scenario."**

models have become especially influential thanks to their ability to generate realistic and context-rich images [12] for various applications, including image editing [9], generation of artistic visuals [10], and even spatiotemporal text-to-video generation [6, 14]. While it is fundamental to examine the geographic knowledge and biases that models encode is crucial for a fair, inclusive, and accurate application of AI in the urban domain [2, 3], this area is remarkably unexplored.

For instance, some studies highlighted the tendency of vision models to reinforce stereotypes or misrepresent specific groups or regions due to training data imbalances [7, 15]. In spatial contexts, generative image models have proven to exhibit regional biases, for instance, generating less realistic or less diverse visuals for prompts about certain continents (e.g., Africa or West Asia) compared to others like Europe or North America [4].

However, most studies rely on textual input or structured metadata, rather than purely visual generation. As such, there is a lack of work investigating whether these models can generate geographically coherent imagery solely from visual representations, as this study explores.

This paper addresses these research gaps by systematically exploring both geographic representation biases and geographic knowledge embedded in contemporary generative AI models, contributing to the broader pursuit of equitable and inclusive technological solutions.

In particular, in this work we investigate the geographic awareness and presence of representation bias of two state-of-the-art image generation models: FLUX 1 schnell [8] and Stable Diffusion 3.5 [13]. To this end, each model generates 150 synthetic images for every U.S. state and the respective capital city. After, we employed the visual embedding DINO-v2 ViT-S/14 [11] to capture high-level semantic features and Fréchet Inception Distance (FID) [5] to measure the similarity between images.

We found that models encode an implicit geographic knowledge proven by the fact that regions of interest that are spatially close to each other are also close to each other in terms of FID.

Interestingly, when generically prompted to generate an image for "USA", both models demonstrated a significant representation bias favouring metropolitan areas with rural areas and smaller cities being underrepresented. Such behaviors suggest a critical diversity deficit, highlighting the tendency of such models to default to stereotypical urban-centric portrayals instead of reflecting comprehensive geographic diversity, of which they have, indeed, knowledge. Recognising and mitigating these biases is vital for developing generative models capable of accurately representing the diverse geographic realities present in the real world.

## 2 Methodology

### 2.1 Models and Prompt

To investigate the geographic representation capabilities and biases of state-of-the-art generative models, we selected two open text-to-image diffusion models: **FLUX 1-schnell** and **Stable Diffusion 3.5-Large** (SD-3.5-L). Both models are representative of recent advances in generative AI and are widely used in academic and applied settings. **FLUX 1-schnell** [8] is a 12-billion-parameter diffusion transformer model. It uses latent adversarial diffusion distillation to achieve high-quality image synthesis with few (2-4) sampling steps. We opted for the "schnell" variant because it allows fast inference while retaining high fidelity. **Stable Diffusion 3.5-large** [13] is an 8.1-billion-parameter model based on a multimodal diffusion transformer architecture. SD-3.5-L excels at producing photorealistic images and demonstrates advanced semantic understanding of complex prompts, making it a representative state-of-the-art open-source model.

To ensure consistency across all queries, we designed our experiment around a fixed prompt template:

"A photorealistic high-resolution street-view photo of {LOCATION}"

The placeholder {LOCATION} was systematically replaced with each of the 50 U.S. states and their respective capital cities, as well as a general reference to "USA". This resulted in 101 distinct prompts. For each prompt, we generated images, totalling 15,150 images per model. The generation process was kept uniform across all runs by using the default inference settings (guidance scale and number of diffusion steps) provided by each model's implementation.

We selected this prompt structure to balance between interpretability and geographic specificity. The street-view framing aims to invoke grounded, spatially contextualised outputs rather than abstract or symbolic imagery. This framing is particularly relevant to urban studies, visual culture, and perception research, where the

appearance of built environments serves as a key indicator of place identity.

Additionally, we intentionally avoided conditioning the prompt on aesthetic styles, time of day, or weather conditions to minimise confounding variables. Our focus was to probe how models represent place identity given minimal stylistic direction, and to assess the variability and consistency within and across locations.

The combination of geographically fine-grained inputs, a photorealistic framing, and repeated sampling enables us to evaluate both the models' internal geographic knowledge and any systematic biases that emerge in their visual outputs. The inclusion of a generic prompt for "USA" further allows us to test for national stereotypes and their alignment with actual geographic variation.

### 2.2 Embedding

To measure and compare the visual semantics encoded by each model's outputs, we employed feature embeddings generated by the DINO-v2 ViT-S/14 model [11]. DINO-v2 is a self-supervised vision transformer trained to learn semantically rich image representations without labelled data. Its ViT-S/14 variant produces 384-dimensional feature vectors from input images and has demonstrated robust performance in urban analysis, including visual place recognition and land-cover classification tasks [11].

In our setup, each image was resized to 256 pixels on the shortest side, centre-cropped to 224×224 pixels, and lastly normalised. Using the pretrained model, each image was passed through DINO-v2 to extract a single embedding vector. For every location (state, capital, and the USA reference), we computed the mean vector and the empirical covariance matrix across the set of image embeddings.

### 2.3 Distance Metrics

Building on the embedding representations, we quantified structural relationships between locations using the Fréchet Inception Distance (FID) [5]. This metric not only allows comparisons between each state or capital and the "USA" reference, but also enables a systematic analysis of all state-state and capital-capital similarities. FID measures the dissimilarity between two multivariate Gaussian distributions, each defined by the mean and covariance of image embeddings for a given location. For two distributions with means  $\mu_1, \mu_2$  and covariances  $\Sigma_1, \Sigma_2$ , FID is computed as follows:

$$\text{FID}(\mu_1, \Sigma_1; \mu_2, \Sigma_2) = \|\mu_1 - \mu_2\|_2^2 + \text{Tr}(\Sigma_1 + \Sigma_2 - 2(\Sigma_1 \Sigma_2)^{1/2}) \quad (1)$$

This equation comprises two key components that capture differences between distributions. The first is the mean term ( $\|\mu_1 - \mu_2\|_2^2$ ), representing the squared Euclidean distance between mean vectors and reflecting changes in the average semantic content of image embeddings. The second is the covariance term ( $\text{Tr}(\Sigma_1 + \Sigma_2 - 2(\Sigma_1 \Sigma_2)^{1/2})$ ), which measures differences in the spread and shape of the distributions by comparing their covariance structures.

Intuitively, the FID penalises both differences in the average embedding and discrepancies in the variation structure of the features. Therefore, a lower FID indicates greater similarity in visual semantics and spatial dispersion. We computed the FID between every pair of locations and stored the results in a symmetric distance matrix. This matrix serves as a foundation for subsequent

clustering and visualisation analyses, providing insights into how closely different regions align in terms of their visual identity.

### 3 Results

#### 3.1 Evidence of Geographic Knowledge

Our analysis reveals that both FLUX 1-schnell and SD 3.5-L demonstrate substantial geographic knowledge. The clustering patterns derived from pairwise FID comparisons provide compelling evidence that these models encode meaningful spatial and visual relationships between U.S. locations.

Both models exhibit clear tendencies to group geographically proximate regions, suggesting an underlying understanding of regional similarities in built environments and landscape characteristics. As Figure 2 shows, the FLUX model demonstrates robust regional clustering for the Mountain West, where Alaska, Colorado, Idaho, Montana, Oregon, Washington, and Wyoming consistently cluster together (Cluster 1). The Desert Southwest states maintain their cohesive grouping, with Arizona, Nevada, New Mexico, and Utah forming a distinct cluster (Cluster 2) in FLUX, indicating recognition of the region's distinctive arid landscape and associated urban development patterns. The SD 3.5-L model shows refined geographic sensitivity by clustering the core Mountain West states (Alaska, Colorado, Idaho, Montana, Oregon) separately from a Plains/Mountain border region (North Dakota, South Dakota, Utah, Wyoming), while grouping the Desert Southwest (Arizona, California, Nevada, New Mexico), demonstrating nuanced recognition of landscape transitions. The New England region presents compelling evidence of geographic knowledge across both models. SD 3.5-L exhibits a particularly sophisticated regional understanding by clustering the core New England and Mid-Atlantic states together (Cluster 5), whereas FLUX maintains a tighter New England grouping (Maine, New Hampshire, Vermont) with the addition of West Virginia, potentially reflecting a shared Appalachian topography and rural character.

Evidence of geographic knowledge also emerges in the capital city clustering patterns, where both models demonstrate a deep understanding of urban hierarchy and developmental characteristics. Major metropolitan capitals cluster together, with both models grouping Atlanta, Austin, Indianapolis, and Oklahoma City, reflecting their shared characteristics as significant state capitals with substantial metropolitan development. SD 3.5-L further refines this by creating a separate cluster for western metropolitan capitals (Denver, Honolulu, Phoenix, Salt Lake City), suggesting recognition of distinct regional urban development patterns. Mid-sized capitals also show coherent clustering patterns, reflecting their shared characteristics as significant state capitals with substantial downtown cores but less metropolitan complexity than the largest cities. Lastly, the clustering of truly small capitals presents a more complex picture. When correctly generated, these locations should reflect the distinctive characteristics of American small-town state capitals: modest governmental buildings, limited high-rise development, and small downtown scales. However, as detailed in the following section, some of these locations suffer from systematic misgeneration issues.

#### 3.2 Small Capital Misgeneration

Clustering results for the smallest state capitals highlight a notable limitation of both models: instead of producing contextually accurate representations, they frequently generate imagery that resembles European cities. This tendency results in visual outputs that group together, reflecting shared architectural elements typical of Old World urban environments.

The clustering patterns provide compelling evidence of systematic misgeneration. In FLUX, three of the four misgenerated capitals (Frankfort, Montpelier, Pierre) cluster together in Cluster 10, forming a cohesive group based on their shared European visual characteristics. Dover, however, appears in Cluster 11 alongside correctly generated American capitals (Annapolis, Concord, Richmond), suggesting that FLUX's Dover misgeneration may be producing European imagery that shares certain visual characteristics with American small capitals, perhaps similar building scales and street layouts. SD 3.5-L exhibits a highly systematic pattern: Bismarck and Pierre form a distinct cluster (Cluster 5), consistently producing misgenerations that resemble German and French urban environments, respectively. Dover and Olympia appear as isolated single-member clusters (Clusters 6 and 11), with Olympia notably characterised by features of ancient Greek architecture. These misgeneration patterns reveal both toponymic confusion and deeper systematic biases. When faced with less-represented or ambiguous place names, both models default to more globally prominent references rather than correctly identifying their intended American contexts.

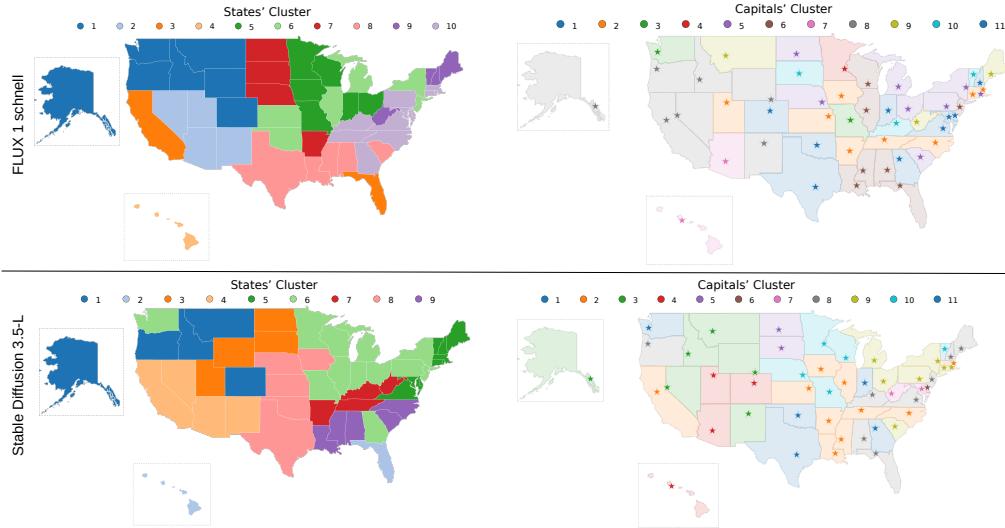
#### 3.3 Lack of Diversity

Geographic bias becomes evident when we compare how closely each location's generated image matches the models' idea of a generic 'USA'. By measuring FID distances between the images of each state and capital and those generated for the prompt 'USA,' we find consistent patterns. The models systematically represent the generic USA through dense, urban environments, favouring certain regions over others. Moreover, as Table 1 shows, both generators struggle to reconcile frontier, desert or tropical landscapes with their USA stereotype. States such as Alaska, Hawaii, Arizona and Nevada register FID scores up to eight times higher than the metropolitan tier. This metropolitan representative bias extends to capital cities, where larger capitals like Raleigh and Boston align closely with the generic "USA" representation, while smaller capitals, such as Juneau, show significant divergence. Thus, when prompted with 'USA,' the models rarely, if ever, generate images that reflect the scale or character of these smaller downtowns.

Both diffusion models possess a fairly detailed latent map of American places, yet deploy only a fraction of that knowledge when asked for a macro-regional view. The default synthesis systematically gravitates to a stylised metropolitan skyline, leaving rural, frontier, and small-town imagery outside the generative spotlight.

### 4 Discussion and Conclusion

Our study reveals a dual dynamic in the geographic knowledge of diffusion models. On one hand, both FLUX 1 and SD 3.5-L exhibit a surprisingly detailed latent understanding of U.S. geography: when prompted with states or capitals, the generated images often reflect regionally coherent features and cluster according to real-world



**Figure 2: Hierarchical clustering of FID for states and capitals for FLUX 1-schnell (top) and Stable Diffusion 3.5-L.(bottom)**

**Table 1: Top-5 and bottom-5 FID to the generic USA prompt, combining states and capital cities.**

Rank	FLUX 1-schnell			Stable Diffusion 3.5-L		
	State	FID	Capital	State	FID	Capital
<i>Most similar to "USA" (lowest FID)</i>						
1	New Jersey	331.8	Madison	408.5	Illinois	413.7
2	Michigan	404.2	Raleigh	467.5	Minnesota	796.2
3	Illinois	414.1	Jackson	484.7	New York	829.2
4	Minnesota	545.8	Nashville	501.1	Georgia	843.5
5	Texas	599.4	Little Rock	539.5	Indiana	848.3
<i>Least similar to "USA" (highest FID)</i>						
50	Hawaii	2594.5	Frankfort	2968.3	Arizona	3352.9
49	Alaska	2143.4	Saint Paul	2730.7	North Dakota	3207.3
48	Florida	1988.7	Dover	2530.8	Hawaii	3139.8
47	Arizona	1970.6	Pierre	2522.3	Nevada	3131.7
46	California	1961.9	Montpelier	2481.1	Wyoming	3088.3

spatial patterns. On the other hand, this knowledge is selectively applied. When prompted at a broader scale, the diversity of American landscapes collapses in a metropolitan-centric stereotype. Additionally, the frequent misgeneration of small capitals further reflects how data sparsity or toponymic ambiguity may amplify this bias.

These findings underscore that addressing geographic bias is no longer a marginal concern, but a requirement for the responsible application of generative technologies in real-world urban analytics and civic design. Future research will extend this work toward global generalisation and develop dynamic prompting strategies that actively elevate underrepresented geographies, paving the way for more inclusive and equitable generative place-making.

## Acknowledgments

B.L. and M.L. acknowledge the partial support of the European Union's Horizon Europe research and innovation program under grant agreement No. 101120237 (ELIAS).

## References

- [1] Ciro Beneduce, Bruno Lepri, and Massimiliano Luca. 2025. Large language models are zero-shot next location predictors. *IEEE Access* (2025).
- [2] Ciro Beneduce, Bruno Lepri, and Massimiliano Luca. 2025. Urban Safety Perception Through the Lens of Large Multimodal Models: A Persona-based Approach. *arXiv preprint arXiv:2503.00610* (2025).
- [3] Melissa Hall, Samuel J Bell, Candace Ross, Adina Williams, Michal Drozdzał, and Adriana Romero Soriano. 2024. Towards geographic inclusion in the evaluation of text-to-image models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*. 585–601.
- [4] Melissa Hall, Candace Ross, Adina Williams, Nicolas Carion, Michal Drozdzał, and Adriana Romero Soriano. 2024. DIG In: Evaluating Disparities in Image Generations with Indicators for Geographic Diversity. *arXiv:2308.06198 [cs.CV]* <https://arxiv.org/abs/2308.06198>
- [5] 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. *Advances in neural information processing systems* 30 (2017).
- [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- [7] Akshita Jha, Vinodkumar Prabhakaran, Remi Denton, Sarah Laszlo, Shachi Dave, Rida Qadri, Chandan K. Reddy, and Sunipa Dev. 2024. ViSAGe: A Global-Scale Analysis of Visual Stereotypes in Text-to-Image Generation. *arXiv:2401.06310 [cs.CV]* <https://arxiv.org/abs/2401.06310>
- [8] Black Forest Labs. 2024. FLUX. <https://github.com/black-forest-labs/flux>.
- [9] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2022. SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. *arXiv:2108.01073 [cs.CV]* <https://arxiv.org/abs/2108.01073>
- [10] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. *arXiv:2112.10741 [cs.CV]* <https://arxiv.org/abs/2112.10741>
- [11] Maxime Oquab, Timothée Daret, Thé Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaddin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2024. DINOV2: Learning Robust Visual Features without Supervision. *arXiv:2304.07193 [cs.CV]* <https://arxiv.org/abs/2304.07193>
- [12] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [13] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [14] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oran Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. 2022. Make-A-Video: Text-to-Video Generation without Text-Video Data. *arXiv:2209.14792 [cs.CV]* <https://arxiv.org/abs/2209.14792>

- [15] Yixin Wan, Arjun Subramonian, Anaelia Ovalle, Zongyu Lin, Ashima Suvarna, Christina Chance, Hritik Bansal, Rebecca Pattichis, and Kai-Wei Chang. 2024. Survey of Bias In Text-to-Image Generation: Definition, Evaluation, and Mitigation. arXiv:2404.01030 [cs.CV] <https://arxiv.org/abs/2404.01030>
- [16] Qingyi Wang, Yuebing Liang, Yunhan Zheng, Kaiyuan Xu, Jinhua Zhao, and Shenhao Wang. 2025. Generative AI for Urban Planning: Synthesizing Satellite Imagery via Diffusion Models. arXiv:2505.08833 [cs.CV] <https://arxiv.org/abs/2505.08833>
- [17] Yuxin Yang, Pengfei Zhu, Mengshi Qi, and Huadong Ma. 2024. Uncovering the human motion pattern: Pattern Memory-based Diffusion Model for Trajectory Prediction. arXiv:2401.02916 [cs.CV] <https://arxiv.org/abs/2401.02916>