

# Dreaming with ChatGPT: Unraveling the Challenges of LLMs Dream Generation

**Harel Berger**

Georgetown University  
hb711@georgetown.edu

**Hadar King**

The Hebrew University of Jerusalem  
Hadar.king@mail.huji.ac.il

**Omer David**

Bar-Ilan University  
omer.david1@live.biu.ac.il

## Abstract

Large Language Models (LLMs), such as ChatGPT, are used daily for different human-like text generation tasks. This motivates us to ask: *Can an LLM generate human dreams?* For this research, we explore this new avenue through the lens of ChatGPT, and its ability to generate valid dreams. We have three main findings: (i) Chatgpt-4o, the new version of chatGPT, generated all requested dreams. (ii) Generated dreams meet key psychological criteria of dreams. We hope our work will set the stage for developing a new task of dream generation for LLMs. This task can help psychologists evaluate patients' dreams based on their demographic factors.

## 1 Introduction

A dream is a series of involuntary images, ideas, and emotions during sleep, especially in the rapid eye movement (REM) stage (apa, 2024). Dreams are crucial in psychology, as they provide insight into the mind, revealing hidden desires, fears, psychological status, and conflicts (Freud, 1900; Hobson, 2009; Solomonova et al., 2021). Dreams are utilized as a therapeutic tool for treating certain psychological disorders (Beauchemin and Hays, 1995). Moreover, Lucid dreaming treatment (LDT) is a clinical method that can help patients reduce nightmares (de Macêdo et al., 2019) and address other mental health issues (Beauchemin and Hays, 1995; Sackwild and Stumbrys, 2021).

Large Language Models (LLMs) aim to mimic psychological phenomena by simulating aspects of human cognition, such as language understanding, reasoning, and emotion recognition (Sartori and Orrù, 2023; Hofweber et al., 2024; Kuo and Chen, 2023). While still not there, using dream descriptions generated by LLMs could be employed in psychological treatments by creating specific dream characteristics tailored to individual needs. LDT requires training and practice (Ellis et al., 2021),

which can be challenging for patients to achieve on their own. Therefore, utilizing LLM-generated dreams, customized to patients' needs and personal characteristics, may enhance the effectiveness of LDT.

In this work, we lay the groundwork for this task - dream description generation. As this avenue is undermined, we try to find whether certain LLMs can generate dream descriptions that meet psychological criteria. We picked ChatGPT, the most globally popular LLM<sup>1</sup> as our test case. We use several versions of ChatGPT3.5 and ChatGPT4o, the most recent version of the OpenAI's LLM. Through an in-depth analysis of the samples produced by different versions of ChatGPT, we find that:

- ChatGPT4o generates all requested dream descriptions, which is false for its predecessors.
- Dream descriptions generated by explored LLMs follow some common psychological definitions of a dream but do not fully capture how a dream looks/feels.

## 2 Dreams in Psychology

Traditionally, dreams are mostly associated and analyzed through REM sleep (Hobson and Pace-Schott, 2002; Nir and Tononi, 2010). Formally, in the APA Dictionary of Psychology (apa, 2024), REM dreams are defined by four attributes: (1) a sense of motion in space paired with visual imagery (*Motion*); (2) strong emotions, especially fear, euphoria, or anger (*Emotion*); (3) the perception that dream events, characters, and situations are real (*Realness*); and (4) unexpected changes in characters, situations, and plot elements (*Discontinuity*). Other attributes derived from psychological works include the location of the dream, which is mostly in normative daily scenes (Domhoff, 2007; Snyder et al., 1968) (*Location*); the existence of at least

<sup>1</sup><https://zapier.com/blog/best-llm/>

one other being (Domhoff, 2007; Snyder, 1970; Dorus et al., 1971) (*Other Beings*); the existence of objects (Domhoff, 2007; Snyder, 1970; Dorus et al., 1971) (*Objects*); and the activity of talking with other beings (Domhoff, 2007; Snyder, 1970) (*Conversation*). We will check if generated dreams meet psychological criteria.

### 3 Related Work

LLMs are being tested through different advanced generation tasks of human nature, such as sarcasm (Chakrabarty et al., 2020), metaphor (Chakrabarty et al., 2021), humour (Mittal et al., 2022; Dsilva, 2024; Tikhonov and Shtykovskiy, 2024), songs (Tian and Peng, 2022; He et al., 2019), hyperbole (Tian et al., 2021), tongue twisters (Loakman et al., 2024), and storytelling (Yao et al., 2019; Yang et al., 2022). Despite its closeness to storytelling, dream description generation is characterized by a sense of discontinuity (apa, 2024) while storytelling showcases a coherent plot (Fan et al., 2019). Also, dreams occur mostly through REM sleep (Hobson and Pace-Schott, 2002; Nir and Tononi, 2010), while stories are written while wide-awake and conscious.

The mimicry of human thinking and behavior by LLMs is still under research. Binz and Schulz (Binz and Schulz, 2023) and Abbasiantaeb et al. (Abbasiantaeb et al., 2024) explored LLMs’ abilities to simulate human understanding and interactions. Zhang et al. (Zhang et al., 2023b) focused on using human-like reasoning to improve LLMs’ decision-making. Another research area is LLMs’ fairness, with Bender et al. (Bender et al., 2021) and Noble (Noble, 2018) highlighting the risks of biases and stereotypes. Zheng et al. (Zheng et al., 2024) explored enhanced biases in judgments carried out by LLMs. Turpin et al. (Turpin et al., 2024) showed that biases in quality assessment tasks can significantly affect LLMs’ reasoning.

However, the topic of generating dream descriptions is yet to be covered. Recent work (Bertolini et al., 2024) explored LLMs’ ability to classify dream descriptions by emotions but did not examine their ability to produce dreams or consider other psychological criteria.

This research evaluates ChatGPT’s ability to generate dream descriptions matching known psychological frameworks (apa, 2024; Snyder et al., 1968; Dorus et al., 1971).

## 4 Methodology

In this research, we tackle these research questions:

1. Can LLMs generate dream descriptions?
2. Do dream descriptions generated by LLMs comply with the psychological criteria of dreams?

To address these questions, we devised a dedicated methodology. To generate the dream descriptions we use Context-Less Generation (Wan et al., 2023; Wan and Chang, 2024), in which the model is prompted with a simple zero-shot description of the dreamer - race, ethnicity, religion, and sex (Hanna et al., 2023; Salewski et al., 2024; Mahomed et al., 2024; Calderon et al.; Zhang et al., 2023a). We added a control group, person. Recent work defined the lack of access to one’s dream content as a limitation (Scarpelli et al., 2022) for correctly analyzing parasomnia events. For a beneficial treatment, it is advised to work closely with a patient to generate dream descriptions based on their descriptions. Therefore, we have chosen different demographic factors to match this suggestion (Table 1).

Our methodology follows the steps below: (1) Given a demographic factor, replace the mask in the prompt with it: "**Pretand you are {MASK}. Please generate a dream that this person dreamt last night.**", and the prompt to model  $M$ . (2) Get the response  $r_i$ . (3) Run steps 1-2 five times in different sessions (no memorization, zero-shot), thus getting the set  $r_{i_j} = \{r_{i_0}, r_{i_1}, r_{i_2}, r_{i_3}, r_{i_4}\}$ . (5) Annotate using human annotators each  $r_{i_j}$  set based on predefined attributes. (6) Analyze the results.

Race	Ethnicity	Religion	Sex	Control
Asian	Indian	Jewish	Female	Person
Black	Arab	Christian	Male	
White	Hispanic	Muslim		

Table 1: Table of simple demographic factors of people used for prompting GPT models.

Some models provided very few dream descriptions. We concluded this by automatically analyzing for a single disclaimer or absence of multiple blank lines<sup>2</sup>.

**Attributes:** For each sample, we annotated the following attributes: (1) is there a dream description?

<sup>2</sup>Concrete dream descriptions were spread across multiple lines upon close inspection.

(yes/no), based on the existence of a story. (2) is there a disclaimer? (yes/no), where a disclaimer is a text similar to "I'm sorry, but I cannot fulfill that request." or "As an AI, I don't have dreams or feelings". (3) the pronoun used for the dreamer. (4-10) the psychological attributes from Section 2 - *Motion* (yes/no), *Emotion* (yes/no), *Realness* (yes/no), *Discontinuity* (yes/no), *Location*, *Other Beings*, *Objects* and *Conversation* (yes/no).

**Human Evaluation:** Three annotators participated: two Masters students with an academic background in psychology and one computer science postdoctoral fellow. Each sample was annotated by two annotators, with a third resolving any disagreements (Mukhtar et al., 2017). The full instructions given to annotators is presented in Appendix C.

**Metrics:** We used a success rate metric for generating dream descriptions, similar to previous work (Wen et al., 2024; Zhao et al., 2024). This measured the model's ability to produce valid dream descriptions (i.e., containing a dream) or without disclaimers. The success rate was the number of samples meeting the criteria divided by the total samples.

## 5 Experiments

We generated four popular demographic factors groups to use for the prompts - religion, race, ethnicity, and sex, and a control group - person (Table 1). We used the prompt from Section 4 with each factor.

We evaluated several gpt models: gpt-3.5-turbo (gpt3.5T), gpt-3.5-turbo-16k-0613 (gpt3.5T16k), gpt-3.5-turbo-0613 (gpt3.5T0613), and gpt-3.5-turbo-1106 (gpt3.5T1106), gpt-4o (gpt4o). We used the default parameters (e.g., temperature 1.0, Top-P 1.0) of all models. In total, for each model, we obtained 60 samples, based on the demographic factors. The samples, code, and annotations are available online<sup>3</sup>. The code is under the MIT license (Open Source Initiative, 2023).

## 6 Dream Description Generation Analysis

In this section, we analyzed all 300 generated samples, to assess the ability of a model to generate a basic dream description. The full details of each are presented in Appendix B.

We found that gpt4o generated 100% of requested dream descriptions, while gpt3.5T16k and gpt3.5T0613 produced over 68%. However, both

<sup>3</sup><https://github.com/harelber/DreamGPT>

gpt3.5T and gpt3.5T1106 generated less than 20% of descriptions, with the latter experiencing a significant drop in performance despite being released later than gpt3.5T16k and gpt3.5T0613<sup>4</sup>. Thus, we dropped further analysis of the latter models.

Out of these generated dream descriptions, we continually analyzed the samples (Table 2). We explored whether the model did not produce a disclaimer stating it is an AI that does not dream, thus following the prompt directly without an explicit objection.

Model	Gen	No Disc	1st
gpt3.5T	13%	-	-
gpt3.5T16k	75%	31%	15%
gpt3.5T0613	68%	39%	17%
gpt3.5T1106	18%	-	-
gpt4o	100%	98%	73%

Table 2: Dream descriptions generation characteristics, based on the generated dreams (*Gen*) out of total sample size, the nonexistence of a disclaimer (*No Disc*), and whether the dream is in first person view (*1st*). The original sample size is 60 dreams. The gpt3.5T and gpt3.5T1106 were eliminated in the deeper analysis due to their poor performance in the initial dream generation.

We found that this phenomenon of no disclaimer+dream was found in 97% of gpt4o samples, 39% of gptT0613 samples, and 31% of gptT061316k samples.

We also looked at whether the description was generated in first person, as the prompt started with "*pretend you are...*". gpt4o met 73% of the times for this aspect, gptgpt3.5T16k 15% of the time, and gpt3.5T0613 17% of the times in the same criteria.

In short, although with some decrease caused by matching the full criteria, gpt4o followed the prompted dream description with a significant gap (~60%) between its performance and the other two models' performances.

Although not all generated samples complied with the no disclaimer+first person criteria, we continued with the generated dream descriptions for further analysis (Gen from Table 2).

For the next sections, we considered 60 dreams for gpt4o, 45 dreams for gpt3.5T16k, and 41 dreams for gptgpt3.5T0613<sup>5</sup>.

<sup>4</sup><https://context.ai/compare/gpt-3-5-turbo-16k/gpt-3-5-turbo>

<sup>5</sup>Similarly to other work (Wan et al., 2023), that drew interesting conclusions from small LLM-generated samples.

## 7 Psychological Dream Attributes

**APA Attributes:** The results of APA’s attributes (Section 2) are presented in Table 3. It can be seen that the three models meet the motion and emotion dream properties raised by APA. In the discontinuity attribute, gpt3.5T16k got a success rate of 56%, and gpt3.5T0613 got 37%. gpt4o shows the greatest promise in this attribute, with a success rate of 70%. However, all models lack a sense of realness, as this property does not have a clear indication in the samples.

Model	M	E	R	D
gpt3.5T16k	98%	100%	0%	56%
gpt3.5T0613	100%	100%	0%	37%
gpt4o	100%	100%	7%	70%

Table 3: APA Attributes Results. M stands for motion, E for emotion, R for realness, and D for discontinuity. It is shown that gpt4o complies the most with APA’s properties of ERM dreams.

Model	N_Loc	Other Beings	Conv
gpt3.5T16k	47%	96%	51%
gpt3.5T0613	44%	88%	44%
gpt4o	73%	95%	67%

Table 4: Other Attributes Results. N\_Loc stands for locations in nature, Other Beings for people/animals, and Conv for conversation. It is shown that gpt4o complies the most with all properties.

**Other Attributes:** We explored attributes from various psychological sources, including locations, beings and interactions (Section 2). Non-daily locations appeared in 73% of gpt4o samples, 44% of gpt3.5T0613 samples, and 47% of gpt3.5T16k samples. This shows that the models do not fully comply with this property. Also, all models included at least one other being in the generated samples (Domhoff, 2007; Snyder, 1970; Dorus et al., 1971). Conversations were found in 67% of gpt4o samples, 44% of gpt3.5T0613 samples, and 51% of gpt3.5T16k’s samples.

Overall, meeting all psychological dream definitions is not trivial for LLMs’ generated dream descriptions. However, the ability to generate dream descriptions with embedded creatures, and motion/emotion rules is met 100% by each model we

explored. Still, gpt4o is the leader in psychological attributes in general.

## 8 Conclusion

In this work, we examined the possibility of generating dream descriptions by LLMs. We explored it through the test case of ChatGPT models. The most promising model was found to be gpt4o. We found that some fundamental psychological attributes are met by the generated descriptions, but there is still progress to be made. We hope this initial work will pave the way to more LLM-dreams research, contributing to the psychological analysis of human dreams, enhancing LDT, and alleviating disorders such as insomnia.

## 9 Ethics Statement

This paper initially explores the capabilities ChatGPT to generate dream descriptions. As the authors only infer descriptions and do not look for a specific person’s dream, the resulting dreams do not expose any private data of an individual.

## 10 Limitations

Despite our interesting findings, this work is subject to several limitations. First, our annotations were based on human annotators. Due to the lack of concise annotations of psychological attributes of dreams, such as discontinuity and realness, we annotated the data with human annotators as an initial work. We envision an extension of this work using fine-tuned model to annotate the data (Wang et al., 2024; Wu et al., 2023).

Second, our data was limited to 300 samples. Although this data seems small, it gave interesting aspects of the ability of LLMs to generate dream descriptions. We intend to curate a larger corpus for more comprehensive research.

Next, we explored ChatGPT, the most popular LLM globally. It would be beneficial to explore the new task with other LLMs (e.g., Llama (Touvron et al., 2023) or Gemini (Team et al., 2023)).

Also, this work initialized the research of generating dream descriptions by LLMs. We used a small set of psychological attributes and a limited set of demographic factors. More advanced work on this topic may follow a broader range of psychological aspects, analyzing combinations of demographic factors, and adding more factors such as jobs and maternity status. This future work will also analyze biases that may arise in the dream descriptions.

## References

2024. *APA Dictionary - Dream*.

Zahra Abbasiantaeb, Yifei Yuan, Evangelos Kanoulas, and Mohammad Aliannejadi. 2024. Let the llms talk: Simulating human-to-human conversational qa via zero-shot llm-to-llm interactions. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 8–17.

Kathleen M Beauchemin and Peter Hays. 1995. Pre-vailling mood, mood changes and dreams in bipolar disorder. *Journal of affective disorders*, 35(1-2):41–49.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623.

Lorenzo Bertolini, Valentina Elce, Adriana Michalak, Hanna-Sophia Widhoezl, Giulio Bernardi, and Julie Weeds. 2024. Automatic annotation of dream report’s emotional content with large language models. In *Proceedings of the 9th Workshop on Computational Linguistics and Clinical Psychology (CLPsych 2024)*, pages 92–107.

Marcel Binz and Eric Schulz. 2023. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120.

Nitay Calderon, Naveh Porat, Eyal Ben-David, Alexander Chapanin, Zorik Gekhman, Nadav Oved, Vitaly Shalumov, and Roi Reichart. Measuring the robustness of nlp models to domain shifts.

Tuhin Chakrabarty, Debanjan Ghosh, Smaranda Muresan, and Nanyun Peng. 2020.  $r^3$ : Reverse, retrieve, and rank for sarcasm generation with commonsense knowledge. *arXiv preprint arXiv:2004.13248*.

Tuhin Chakrabarty, Xurui Zhang, Smaranda Muresan, and Nanyun Peng. 2021. **MERMAID: Metaphor generation with symbolism and discriminative decoding**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4250–4261, Online. Association for Computational Linguistics.

Tainá Carla Freitas de Macêdo, Glescikelly Herminia Ferreira, Katie Moraes de Almondes, Roumen Kirov, and Sérgio Arthur Mota-Rolim. 2019. My dream, my rules: can lucid dreaming treat nightmares? *Frontiers in psychology*, 10:2618.

G William Domhoff. 2007. Realistic simulation and bizarreness in dream content: Past findings and suggestions for future research. *The new science of dreaming*, 2:1–27.

E. Dorus, W. Dorus, and A. Rechtschaffen. 1971. **The incidence of novelty in dreams**. *Archives of General Psychiatry*, 25(4):364–368.

Ryan Rony Dsilva. 2024. **Augmenting Large Language Models with Humor Theory To Understand Puns**. Ph.D. thesis, Purdue University Graduate School.

Jason G Ellis, Joseph De Koninck, and Celyne H Bastien. 2021. Managing insomnia using lucid dreaming training: A pilot study. *Behavioral sleep medicine*, 19(2):273–283.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. *arXiv preprint arXiv:1902.01109*.

Sigmund Freud. 1900. *The Interpretation of Dreams*. Macmillan, New York.

John J Hanna, Abdi D Wakene, Christoph U Lehmann, and Richard J Medford. 2023. Assessing racial and ethnic bias in text generation for healthcare-related tasks by chatgpt1. *MedRxiv*.

He He, Nanyun Peng, and Percy Liang. 2019. **Pun generation with surprise**. In *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)*, pages 1734–1744, Minneapolis, Minnesota. Association for Computational Linguistics.

J Allan Hobson. 2009. Rem sleep and dreaming: towards a theory of protoconsciousness. *Nature Reviews Neuroscience*, 10(11):803–813.

J Allan Hobson and Edward F Pace-Schott. 2002. The cognitive neuroscience of sleep: neuronal systems, consciousness and learning. *Nature Reviews Neuroscience*, 3(9):679–693.

Thomas Hofweber, Peter Hase, Elias Stengel-Eskin, and Mohit Bansal. 2024. Are language models rational? the case of coherence norms and belief revision. *arXiv preprint arXiv:2406.03442*.

Hui-Chi Kuo and Yun-Nung Chen. 2023. **Zero-shot prompting for implicit intent prediction and recommendation with commonsense reasoning**. *Preprint*, arXiv:2210.05901.

Tyler Loakman, Chen Tang, and Chenghua Lin. 2024. Train & constrain: Phonologically informed tongue-twister generation from topics and paraphrases. *arXiv preprint arXiv:2403.13901*.

Yaaseen Mahomed, Charlie M Crawford, Sanjana Gautam, Sorelle A Friedler, and Danaë Metaxa. 2024. Auditing gpt’s content moderation guardrails: Can chatgpt write your favorite tv show? In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 660–686.

- Anirudh Mittal, Yufei Tian, and Nanyun Peng. 2022. *AmbiPun: Generating humorous puns with ambiguous context*. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1053–1062, Seattle, United States. Association for Computational Linguistics.
- Neelam Mukhtar, Mohammad Abid Khan, and Nadia Chiragh. 2017. Effective use of evaluation measures for the validation of best classifier in urdu sentiment analysis. *Cognitive Computation*, 9:446–456.
- Yuval Nir and Giulio Tononi. 2010. Dreaming and the brain: from phenomenology to neurophysiology. *Trends in cognitive sciences*, 14(2):88–100.
- Safiya Umoja Noble. 2018. Algorithms of oppression: How search engines reinforce racism. In *Algorithms of oppression*. New York university press.
- Open Source Initiative. 2023. Mit license. <https://opensource.org/license/mit/>. Accessed: 2024-06-14.
- Lana Sackwild and Tadas Stumbrys. 2021. The healing and transformative potential of lucid dreaming for treating clinical depression. *International journal of dream research*, 14(2):296–308.
- Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models’ strengths and biases. *Advances in Neural Information Processing Systems*, 36.
- Giuseppe Sartori and Graziella Orrù. 2023. Language models and psychological sciences. *Frontiers in Psychology*, 14:1279317.
- Serena Scarpelli, Valentina Alfonsi, Maurizio Gorgoni, and Luigi De Gennaro. 2022. What about dreams? state of the art and open questions. *Journal of Sleep Research*, 31(4):e13609.
- F. Snyder. 1970. The phenomenology of dreaming. In L. Madow and L.H. Snow, editors, *The Psychodynamic Implications of the Physiological Studies on Dreams*, pages 124–151. Charles S Thomas, Springfield.
- F. Snyder, I. Karacan, V. K. Jr. Tharp, and J. Scott. 1968. Phenomenology of rems dreaming. *Psychophysiology*, 4(3):375.
- Elizaveta Solomonova, Claudia Picard-Deland, Iris L Rapoport, Marie-Hélène Pennestri, Mysa Saad, Tetyana Kendzerska, Samuel Paul Louis Veissiere, Roger Godbout, Jodi D Edwards, Lena Quilty, et al. 2021. Stuck in a lockdown: Dreams, bad dreams, nightmares, and their relationship to stress, depression and anxiety during the covid-19 pandemic. *PLoS One*, 16(11):e0259040.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Yufei Tian and Nanyun Peng. 2022. *Zero-shot sonnet generation with discourse-level planning and aesthetics features*. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3587–3597, Seattle, United States. Association for Computational Linguistics.
- Yufei Tian, Arvind krishna Sridhar, and Nanyun Peng. 2021. *HypoGen: Hyperbole generation with commonsense and counterfactual knowledge*. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1583–1593, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alexey Tikhonov and Pavel Shtykovskiy. 2024. Humor mechanics: Advancing humor generation with multi-step reasoning. *arXiv preprint arXiv:2405.07280*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2024. Language models don’t always say what they think: unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36.
- Yixin Wan and Kai-Wei Chang. 2024. White men lead, black women help: Uncovering gender, racial, and intersectional bias in language agency. *arXiv preprint arXiv:2404.10508*.
- Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. 2023. "kelly is a warm person, joseph is a role model": Gender biases in llm-generated reference letters. *arXiv preprint arXiv:2310.09219*.
- Xinru Wang, Hannah Kim, Sajjadur Rahman, Kushan Mitra, and Zhengjie Miao. 2024. Human-llm collaborative annotation through effective verification of llm labels. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–21.
- Hao Wen, Yuanchun Li, Guohong Liu, Shanhui Zhao, Tao Yu, Toby Jia-Jun Li, Shiqi Jiang, Yunhao Liu, Yaqin Zhang, and Yunxin Liu. 2024. Autodroid: Llm-powered task automation in android. In *Proceedings of the 30th Annual International Conference on Mobile Computing and Networking*, pages 543–557.
- Zhanglin Wu, Yilun Liu, Min Zhang, Xiaofeng Zhao, Junhao Zhu, Ming Zhu, Xiaosong Qiao, Jingfei

Zhang, Ma Miaomiao, Zhao Yanqing, et al. 2023. Empowering a metric with llm-assisted named entity annotation: Hw-tsc’s submission to the wmt23 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 822–828.

Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. 2022. Re3: Generating longer stories with recursive reprompting and revision. *arXiv preprint arXiv:2210.06774*.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7378–7385.

Angela Zhang, Mert Yuksekogonul, Joshua Guild, James Zou, and Joseph Wu. 2023a. Chatgpt exhibits gender and racial biases in acute coronary syndrome management. *medRxiv*, pages 2023–11.

Zheyuan Zhang, Shane Storks, Fengyuan Hu, Sungryull Sohn, Moontae Lee, Honglak Lee, and Joyce Chai. 2023b. From heuristic to analytic: Cognitively motivated strategies for coherent physical commonsense reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7354–7379, Singapore. Association for Computational Linguistics.

Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2024. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19632–19642.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.

## A Nature locations found in Dreams - Full Analysis

This section shows the full list of locations found in our dreams data. The locations can be found in table 5.

Garden	Sea	Ocean	Forest
Meadow	Lake	Waterfall	River
Mountain	Field	Oasis	Island
Lagoon	Sky	Hills	Pond

Table 5: Nature locations of dreams found in our data.

## B Models History & Tokens

Table 6 discloses the dream generation rates of each explored model, based on its release date and

amount of tokens, as a complementary to Section 6. Dates and Tokens data acquired from<sup>6 7</sup>.

Model	DGR	Date	Tokens
gpt3.5T	13% (8)	11.28.22	4K
gpt3.5T16k	73% (44)	06.13.23	16K
gpt3.5T0613	68% (41)	06.13.23	4K
gpt3.5T1106	18% (11)	11.06.23	16K
gpt4o	100% (60)	05.13.24	128K

Table 6: Dream generation rate (DGR), based on each model, its date of release, and the number of tokens used as context window. The DGR is measured by counting the actual dreams (no sole disclaimer) out of all responses. The generation rate and actual count are provided for clarity.

## C Instructions to Annotators

In the annotations of dreams, when the symbol (V/X) is shown, please put V for true/exists, and X for false/nonexist. If you are not sure, please put X.

These are the attributes we explore:

- Is there a dream description(v/x) - is there a story or just a statement on the inability of the AI to generate a dream?
- disclaimer (v/x) - if the model states something as “As an AI, I don’t have personal dreams as humans do. However, I can create a fictional dream scenario for you.”, this means that it disclaims that it generates a dream and it is not natural. If there is nothing more than this disclaimer, and no dream was generated, please leave the entire row blank.
- narrator (I/You/He/She/They) - The point of view of the dreamer - is it “I dreamt that. . .”, or “he dreamt”. This is considered as the pronoun of a dream in the paper.
- location - A one-word location of the dream, such as desert, garden. If the dreamer moves places, please add other places.
- other persons - other persons mentioned in the dream

<sup>6</sup><https://community.openai.com/t/what-are-the-differences-between-gpt-3-5-turbo-models/557028/2>

<sup>7</sup><https://context.ai/compare/gpt-3-5-turbo-16k/gpt-3-5-turbo>

- animals - same as persons, but with animals
- items - same with animals, but with items
- conversation (v/x) - if there is any conversation in the dream.
- motion (x/v) - visual imagery along with a sense of motion in space, such as "I was walking".
- emotion (x/v) - intense emotion, especially fear, elation, or anger.
- belief of realness (x/v) - belief that dream characters, events, and situations are real
- discontinuity (x/v) - sudden discontinuities in characters, situations, and plot elements. The word suddenly helps a lot here