

Watermarking for Large Language Models (LLMs)

Arpit Vaghela
University Of Guelph
avaghela@uoguelph.ca

ABSTRACT

This paper builds upon the foundational work presented in A Watermark for Large Language Models by Kirchenbauer et al. [12], which proposes a novel approach to embedding invisible watermarks in the outputs of large language models (LLMs). Our research introduces significant enhancements to their watermarking methodology, addressing specific limitations such as robustness against adversarial attacks [32], efficiency [23], and scalability [30]. We introduce a novel approach for text generation using large language models, combining watermarking and cluster bias techniques. We evaluate the effectiveness of this combined approach with various configurations of watermark strength and bias, comparing it to baseline methods using BLEU, ROUGE, and Perplexity metrics.

ACM Reference Format:

Arpit Vaghela. 2018. Watermarking for Large Language Models (LLMs). In *Proceedings of ACM XXXXX, XX, XX, XX, XX 2018 (ACM XXXXX'18)*, 21 pages.
https://doi.org/10.475/123_4

1 INTRODUCTION

Large language models (LLMs), such as OpenAI's GPT and ChatGPT, have shown unprecedented capabilities in generating coherent and contextually appropriate text. These systems can perform a wide range of tasks, from creating high-quality prose to producing executable code, often with outputs indistinguishable from those written by humans [27]. However, the widespread deployment of LLMs introduces risks associated with malicious use, including social engineering, misinformation campaigns, and academic dishonesty [4] [21] [11]. Beyond these immediate concerns, synthetic text proliferation poses a long-term challenge for dataset curation, as it risks contaminating training datasets with inferior content, complicating future model training efforts [26] [2].

The ability to detect and audit machine-generated text is critical for mitigating these risks. Watermarking, a process that embeds imperceptible yet detectable markers in generated text, has emerged as a promising solution [11]. By enabling efficient detection of synthetic content, watermarking can support efforts to curb misinformation, safeguard intellectual property, and ensure academic integrity. Importantly, watermarking methods can be designed to work without requiring access to a model's parameters or retraining, making them versatile and applicable across different systems [6] [8].

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).
ACM XXXXX'18, XX 2018, XX, XX, XX
© 2018 Copyright held by the owner/author(s).
ACM ISBN 123-4567-24-567/08/06...\$15.00
https://doi.org/10.475/123_4

Despite its promise, watermarking faces significant technical challenges. Low-entropy sequences, such as deterministic outputs from prompts like "for(i=0;i<n;i++) sum+=array[i];", are particularly resistant to watermarking due to their inherent predictability [12]. Modifying such sequences for watermark embedding can lead to unnatural outputs, undermining the overall quality of generated text. To address these challenges, complementary techniques like token clustering and biasing strategies are needed to enhance watermarking's robustness without sacrificing semantic fidelity.

In addition to watermarking, another challenge in LLMs is the control over the content generation process. Models may produce text that lacks coherence or diversity, especially in long-form generation tasks. Cluster bias introduces a mechanism that biases the model toward certain groups or clusters of tokens, ensuring that the output remains semantically consistent and aligned with the desired narrative or structure. Cluster bias has been explored in various contexts, with Wu et al. [30] providing insights into techniques that help preserve model performance while maintaining text coherence.

In this paper, we propose a combined approach that integrates both watermarking and cluster bias in the text generation process. The WatermarkLogitsProcessor embeds a watermark into the text, while the ClusterBiasLogitsProcessor adjusts the model's output to favor tokens from the same cluster. This dual mechanism not only ensures the authenticity and security of generated content but also improves the semantic coherence of the text. Our approach aims to strike a balance between watermark strength and model performance, a challenge addressed by various methods in the literature [24]. Through a series of experiments, we evaluate the effectiveness of our method, comparing it to baseline approaches using metrics such as BLEU, ROUGE, and Perplexity.

1.1 Notation & Language Model Basics

Language models, such as GPT, operate within a predefined "vocabulary" V , which consists of tokens—words or word fragments. These vocabularies often encompass over 50,000 tokens to ensure comprehensive coverage of natural language [17, 26]. A sequence of text can be represented as $\{s^{(t)}\} \in V^T$, where T denotes the total number of tokens.

1.1.1 Prompt and Generated Tokens. In typical scenarios, entries with negative indices, $s^{(-N_p)}, \dots, s^{(-1)}$, constitute the *prompt* of length N_p . Tokens $s^{(0)}, \dots, s^{(T)}$ are then generated in response to the prompt. These tokens are sampled iteratively using a *language model* (LM), f , which predicts the next token based on prior context. Formally, f is a function that, given a sequence $\{s^{(-N_p)}, \dots, s^{(t-1)}\}$, outputs a vector of logits $z \in \mathbb{R}^{|V|}$, where each logit corresponds to a token's score.

The logits are transformed into a probability distribution via the *softmax* function:

$$P(s^{(t)} | s^{(-N_p)}, \dots, s^{(t-1)}) = \frac{\exp(z_{s^{(t)}})}{\sum_{v \in V} \exp(z_v)}. \quad (1)$$

The next token $s^{(t)}$ is sampled from this probability distribution, using techniques like:

- **Multinomial Sampling:** Randomly selecting tokens proportional to their probabilities.
- **Greedy Decoding:** Always selecting the token with the highest probability.
- **Beam Search:** Exploring multiple candidate sequences to find the most probable sequence globally.

1.1.2 Entropy in Token Prediction. Entropy, a measure of uncertainty in token prediction, varies across sequences. High-entropy prompts (e.g., conversational queries) yield diverse completions, while low-entropy prompts (e.g., programming code snippets) result in deterministic outputs. This entropy disparity significantly impacts text watermarking strategies.

1.2 The Challenge of Watermarking Low-Entropy Sequences

Low-entropy sequences, such as code or fixed expressions, present inherent difficulties for watermarking. For example, consider these prompts and completions:

- **Prompt 1:** *The sun always sets* **Completion:** *in the west every day.*
- **Prompt 2:** *for(i=0;i<n;i++)* **Completion:** *sum+=array[i];*

In both cases, the output is strongly constrained by the input, posing two key challenges for watermarking:

- (1) **Indistinguishability:** Both humans and models produce nearly identical outputs, complicating detection of machine-generated text.
- (2) **Perceptibility:** Embedding watermarks by altering token sequences may degrade quality, leading to unnatural completions with high *perplexity* (a metric for evaluating text fluency).

1.2.1 Importance of Robust Techniques. Addressing these challenges requires advanced watermarking algorithms that preserve semantic and syntactic fidelity while embedding detectable markers. Techniques such as *token clustering* and *biasing strategies* are critical for enhancing robustness, particularly in low-entropy contexts. In subsequent sections, we detail these approaches and their integration with watermarking frameworks.

2 RELATED WORKS

Large language models (LLMs), such as OpenAI's GPT and ChatGPT, have shown unprecedented capabilities in generating coherent and contextually appropriate text. These systems can perform a wide range of tasks, from creating high-quality prose to producing executable code, often with outputs indistinguishable from those written by humans (Schulman et al., 2022). However, the widespread deployment of LLMs introduces risks associated with malicious use, including social engineering, misinformation campaigns, and academic dishonesty (Bergman et al., 2022; Mirsky et al.,

2023; Kirchenbauer et al., 2023). Beyond these immediate concerns, synthetic text proliferation poses a long-term challenge for dataset curation, as it risks contaminating training datasets with inferior content, complicating future model training efforts (Radford et al., 2019; Bender et al., 2021).

The ability to detect and audit machine-generated text is critical for mitigating these risks. Watermarking, a process that embeds imperceptible yet detectable markers in generated text, has emerged as a promising solution (Kirchenbauer et al., 2023). By enabling efficient detection of synthetic content, watermarking can support efforts to curb misinformation, safeguard intellectual property, and ensure academic integrity. Importantly, watermarking methods can be designed to work without requiring access to a model's parameters or retraining, making them versatile and applicable across different systems (Crothers et al., 2022; Grinbaum Adomaitis, 2022).

Despite its promise, watermarking faces significant technical challenges. Low-entropy sequences, such as deterministic outputs from prompts like "for(i=0;i<n;i++) sum+=array[i];", are particularly resistant to watermarking due to their inherent predictability (Kirchenbauer et al., 2023). Modifying such sequences for watermark embedding can lead to unnatural outputs, undermining the overall quality of generated text. To address these challenges, complementary techniques like token clustering and biasing strategies are needed to enhance watermarking's robustness without sacrificing semantic fidelity.

In addition to watermarking, biasing techniques have been explored to enhance the quality of generated text and ensure coherence. Wu et al. [30] proposed a watermarking approach that minimizes performance degradation by embedding the watermark in such a way that it does not negatively affect the overall output of the model. Their approach underscores the importance of minimizing trade-offs between model performance and watermarking strength. Furthermore, Abadi et al. (2016) explored the integration of privacy-preserving mechanisms with deep learning models, including the use of clustering and differential privacy, which have inspired similar approaches in the watermarking domain. These techniques demonstrate how controlled bias can guide models to produce more semantically consistent and contextually relevant text.

2.1 Red List and Green List Watermarking: A Contribution by Kirchenbauer et al.

Watermarking techniques for text generation have gained significant attention in recent years, especially with the growing use of generative models. These techniques aim to embed identifiable information within generated content, which can be used for attribution or detection purposes.

One common approach to text watermarking is the use of *hard red list watermarking*, where certain tokens are strictly forbidden from appearing in the generated text. This method works by excluding specific words or phrases from the vocabulary that the model can choose from during token generation. Such approaches are typically simpler to implement and can be effective in preventing the use of certain terms in the output, which makes detection easier.

For example, the token generation process can be constrained as follows:

$$w_t = \arg \max_{w \in T \setminus R} P(w|w_1, w_2, \dots, w_{t-1}), \quad (2)$$

where R represents a list of forbidden words, and T is the complete token set. This method, while effective in detection, can disrupt the natural flow of the text by forcing the model to avoid specific tokens.

On the other hand, *soft red list watermarking* involves encouraging the model to select certain tokens more frequently without completely excluding others. This is achieved by adjusting the token probabilities to increase the likelihood of choosing specific words from a "green list." The modified probability distribution is given by:

$$P'(w_t) = P(w_t) \cdot (1 + \alpha \cdot \mathbb{I}[w_t \in G]), \quad (3)$$

where α is a factor that controls the strength of the bias, and G is the set of encouraged tokens. This method provides a more flexible approach by subtly influencing the output generation, thereby reducing the disruption of natural language while still embedding the watermark.

These watermarking techniques play a crucial role in protecting intellectual property and ensuring attribution for AI-generated content, especially in the context of machine learning models like GPT-based systems. Detection methods typically involve analyzing the frequency of tokens from the red or green list, using statistical tests such as the Z-score to determine whether the occurrence deviates from normal usage patterns.

Both hard and soft watermarking strategies are extensively used in the context of digital media protection and copyright enforcement, providing valuable tools for managing AI-generated content.

Cluster-based biasing, while less explored in the context of watermarking, has shown promise in enhancing the coherence of generated text. The application of cluster bias has been studied in tasks such as text summarization and machine translation, where token clustering helps ensure that generated sequences maintain a consistent thematic structure [7]. By imposing constraints on token selection based on cluster assignments, models can be guided to generate more coherent and contextually appropriate outputs, particularly in complex or long-form generation tasks. Moreover, cluster biasing has been shown to contribute to the interpretability and robustness of LLMs, a critical factor when deploying models in real-world applications [24].

Another line of research has focused on adversarial robustness, which is crucial for ensuring that watermarking and biasing mechanisms are resistant to manipulation. Papernot et al. [24] explored the vulnerability of deep learning systems to black-box attacks, highlighting the need for watermarking and biasing techniques that can withstand adversarial threats. These findings reinforce the importance of developing robust watermarking schemes that not only protect intellectual property but also safeguard the integrity of model outputs in adversarial environments.

2.2 Our Proposal

Our work builds on these foundational efforts by integrating watermarking and cluster biasing techniques, contributing to both the

security and coherence of LLM-generated content. By combining these techniques, we aim to enhance model traceability, protect against adversarial manipulation, and improve the overall quality of generated text.

3 PROBLEM DEFINITION

The increasing deployment of large language models (LLMs) such as OpenAI's GPT-3, GPT-4, and similar systems has transformed many industries, from content creation to customer service. These models, capable of generating text that is often indistinguishable from human writing, have introduced several ethical and security concerns. Key issues include the potential for misinformation, plagiarism, and intellectual property (IP) theft, which are exacerbated by the inability to reliably differentiate machine-generated text from human-generated content. As LLMs are increasingly integrated into everyday applications, the need for robust methods to authenticate and track generated content becomes paramount.

One promising solution to these challenges is watermarking, a technique aimed at embedding hidden, traceable information within the generated text. Watermarking can provide accountability and prevent misuse by allowing stakeholders to trace the origin of the content. However, existing watermarking techniques face significant limitations, particularly when dealing with low-entropy sequences—texts that are highly repetitive or predictable, such as deterministic code snippets or frequently used phrases.

Low-entropy sequences present unique challenges that complicate the implementation of effective watermarking techniques:

- **Detection Complexity:** Low-entropy sequences are often indistinguishable from human-written content, especially when both the model and human outputs share similar structures or patterns. This makes watermark detection in these sequences particularly difficult. Standard watermarking methods, which usually rely on subtle token manipulations or embedding specific patterns, often fail to introduce detectable watermarks without disrupting the flow of the text [32]. As a result, distinguishing between human-generated and machine-generated text becomes increasingly complex in these contexts.
- **Quality Degradation:** Watermark embedding techniques, particularly those that involve modifying token probabilities or token selection, often degrade the quality of the generated text. In low-entropy contexts, even small alterations can have a significant impact on coherence and fluency. This is a critical challenge, as the insertion of watermarks should not disrupt the natural flow of text or introduce noticeable artifacts, which can harm the overall user experience [32]. Maintaining high semantic fidelity while embedding a robust watermark is a delicate balance that current methods struggle to achieve.
- **Scalability and Robustness:** Many existing watermarking techniques that work well for high-entropy, diverse text generation fail to scale effectively to low-entropy sequences. These methods are often not adaptable to different types of content or language models, limiting their practicality in real-world applications. Furthermore, such methods may not

provide sufficient robustness in cases where model outputs are intentionally manipulated to avoid detection [22].

This paper proposes a robust watermarking framework incorporating token manipulation strategies such as clustering and biasing to address these challenges. By clustering similar tokens together, our approach ensures that the watermark can be embedded without disrupting the natural structure of the text. Token biasing further allows us to subtly adjust token selection probabilities, embedding the watermark in a way that minimally impacts the fluency and semantic integrity of the text.

The objectives of this work are twofold: (1) to enhance the accuracy and reliability of watermarking techniques for low-entropy sequences, and (2) to maintain the semantic quality of the generated text, ensuring that watermarks can be effectively embedded without compromising the content's coherence or readability. This research aims to contribute to the development of robust, scalable watermarking methods that preserve both the traceability of machine-generated content and the quality of the text itself [13, 20].

4 METHODOLOGY

4.1 Context on GPT-2 and OPT Models

In this study, we utilize models from the *GPT-2* family, specifically the *GPT-2 Medium* variant, along with models from the *OPT* family, including *OPT-350M*, for text generation tasks. The *GPT-2* family, developed by OpenAI, consists of large, unsupervised language models designed to generate coherent, contextually relevant text by predicting the next word in a given sequence. These models are based on the transformer architecture, which has proven highly effective in natural language processing (NLP) tasks such as text generation, language modeling, and machine translation [26, 29].

4.1.1 *GPT-2 Overview.* The original *GPT-2* model consists of 1.5 billion parameters and was trained on a dataset of 40GB of internet text. Its pre-training on a vast corpus of diverse textual data allows it to generate fluent and contextually coherent text across various topics and domains [26].

Among the various sizes in the *GPT-2* family, the *GPT-2 Medium* model, with 345 million parameters, offers a balanced trade-off between computational efficiency and the quality of generated text. This variant maintains strong performance while being more computationally efficient than the largest model, making it an ideal choice when balancing resources and speed is essential. While *GPT-2 Medium* is smaller than the full 1.5 billion-parameter model, it is still capable of generating high-quality, coherent text, and is often used in scenarios where a compromise between efficiency and output quality is required.

For our use case, we selected *GPT-2 Medium* due to its ability to generate high-quality text while being computationally more efficient compared to the larger *GPT-2* model. This choice was particularly important to ensure faster processing times and reduced resource demands, while still maintaining a strong level of text generation performance suitable for our needs.

4.1.2 *OPT Overview.* The *OPT* (Open Pre-trained Transformer) family, developed by Meta, offers open alternatives to *GPT* models. These models, including variants like *OPT-125M* and *OPT-350M*, are built with a similar transformer architecture to *GPT* models

but are smaller in scale. The *OPT* models are designed to provide more accessible, open alternatives while maintaining competitive performance for a variety of natural language processing tasks.

In our study, we used the *OPT-350M* model, which has 350 million parameters. The *OPT-350M* was selected for comparison with the *GPT-2 Medium* model due to its size and accessibility as an open-source model. This choice allowed us to evaluate the performance of a smaller-scale model from the *OPT* family alongside the *GPT-2 Medium* model, offering insights into the trade-offs between model size, accessibility, and computational efficiency.

4.1.3 *Applications of GPT-2 and OPT in Text Generation.* The primary application of the *GPT-2* family in this study is for text generation. These models have demonstrated remarkable performance in generating human-like text and have been applied to various tasks, including text completion, summarization, and creative writing [5?]. In contrast, *OPT* models like *OPT-125M* and *OPT-350M* provide open alternatives to *GPT* models, with similar transformer-based architectures but smaller parameter scales.

For our experiments, we focused on comparing the performance of *GPT-2 Medium* and *OPT-350M* in generating diverse textual outputs. Both models were evaluated using varying parameters, such as size of the green list (γ), bias towards green list (δ), and performance metrics, to understand their capabilities and differences in text generation. This comparison allowed us to explore trade-offs in model design, accessibility, and efficiency while generating high-quality text.

4.2 Cluster-Based Token Sampling Using Semantic Similarity

Inspired by Dr. Qi Li's work on a Cluster-based algorithm for Tokens gathering, Token sampling based on semantic similarity is an innovative approach designed to enhance the efficiency and interpretability of language models. This method groups tokens into semantically coherent clusters, allowing for targeted sampling within these clusters. By focusing on tokens that share similar meanings or contexts, the approach reduces redundancy and ensures that the most representative tokens are utilized for downstream tasks.

4.2.1 *Clustering Semantic Spaces.* Language models, such as *GPT-2*, encode tokens into high-dimensional embedding spaces where proximity reflects semantic similarity [28, 29]. These embeddings capture nuanced relationships between words, phrases, and subwords, making them ideal for clustering algorithms. In our work, we employ **Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN)**, a robust method that can adapt to varying cluster densities while handling noise effectively [18].

- **Embedding Representations:** Tokens are transformed into dense vectors through the model's embedding layer. These vectors serve as the input for clustering, representing the semantic landscape of the token vocabulary.
- **Cluster Formation:** HDBSCAN identifies groups of tokens with similar semantic properties. Tokens labeled as noise (-1) are excluded from sampling to maintain the quality of clusters.

Algorithm 1 Token Watermarking with Cluster-Based Green Lists

Input: Prompt: A sequence of prior tokens $s(-N_p), \dots, s(-1)$, used as context to predict the next token.

Green list size parameter: $\gamma \in (0, 1)$, controlling the proportion of tokens that will be biased toward selection (marked as "green").

Hardness parameter: $\delta > 0$, a bias added to the logits of green list tokens to increase their selection probability.

Vocabulary partition: `token_to_cluster`, a mapping of tokens to clusters (optional, default: entire vocabulary as one cluster).

Output: Watermarked token sequence $s(0), s(1), \dots$, where green list tokens are preferentially selected.

1 Initialization:

- Assign vocabulary V and size $|V|$.
- Set γ (proportion of tokens to mark as "green") to specify the green list size.
- Set δ (bias added to green list logits) to control the magnitude of the bias.

Token Clustering Using `create_cluster` Function:

- Call the `create_cluster(model, tokenizer, args)` function to obtain the `token_to_cluster` mapping.
- This function performs DBSCAN (or HDBSCAN) clustering based on token embeddings and returns a mapping from token IDs to cluster IDs.
- Optionally, noise tokens (assigned cluster ID of -1) are filtered out.

for $t = 0, 1, 2, \dots$ **do**

Step 2: Use the language model to compute logits $l(t)$ for the current token $s(t)$ based on prior tokens $s(-N_p), \dots, s(t-1)$.

Step 3: Hash the previous token $s(t-1)$ to generate a seed for the random number generator (RNG). The seed ensures the deterministic reproducibility of token partitioning.

Step 4: Partition the vocabulary into clusters based on the `token_to_cluster` mapping:

- If no `token_to_cluster` is provided, treat the entire vocabulary as a single cluster.
- Randomly shuffle tokens within each cluster using the seeded RNG. This randomness ensures that the process is not biased by the order of tokens.
- Compute G (green list) as the top $\gamma \cdot |V|$ tokens for each cluster, and R (red list) as the remainder. The green list contains the tokens most likely to be selected, while the red list contains the less likely ones.

Step 5: Bias the logits for green list tokens by adding δ :

$$\hat{p}(t)_k = \begin{cases} \frac{\exp(l(t)_k + \delta)}{\sum_{i \in G} \exp(l(t)_i + \delta) + \sum_{i \in R} \exp(l(t)_i)}, & k \in G, \\ \frac{\exp(l(t)_k)}{\sum_{i \in G} \exp(l(t)_i + \delta) + \sum_{i \in R} \exp(l(t)_i)}, & k \in R. \end{cases}$$

This formula adjusts the logits for tokens in the green list by adding δ , making them more likely to be selected, while keeping the logits of red list tokens unchanged.

Step 6: Sample the next token $s(t)$ from the watermarked distribution $\hat{p}(t)$, where tokens from the green list G have higher probabilities due to the added bias.

Token Sample Function:

- The `token_sample_function` is used to sample a token $s(t)$ based on the adjusted probabilities $\hat{p}(t)$.
- It takes the watermarked probability distribution $\hat{p}(t)$ as input and returns the selected token $s(t)$.
- The function may implement sampling strategies like top-k sampling or nucleus sampling to select a token based on its probability.

- **Distance Metrics:** Cosine similarity or Euclidean distance is used to measure token proximity, ensuring that clusters are formed based on meaningful semantic relations [18].

4.2.2 Token Sampling Within Clusters. After clustering, tokens are sampled selectively from each cluster. This ensures a diverse representation of semantic categories, enabling more robust model evaluations and applications.

- (1) **Representative Sampling:** Key tokens from each cluster are chosen based on centrality or relevance within the cluster, reducing bias in token selection.
- (2) **Noise Exclusion:** Tokens deemed as noise by HDBSCAN are excluded to maintain the semantic integrity of the sample set [18].
- (3) **Balanced Distribution:** Sampling is distributed across clusters proportionally, ensuring coverage of the entire semantic space [3].

4.2.3 Advantages of Cluster-Based Sampling.

- **Semantic Diversity:** By sampling from clusters, the approach ensures that a wide range of semantic categories is represented, enhancing downstream performance [29].
- **Efficiency:** Reducing redundancy through clustering minimizes computational overhead, particularly in tasks involving large vocabularies [3].
- **Improved Interpretability:** Clustering provides insights into the model's understanding of token semantics, enabling qualitative analysis.

4.2.4 Applications. Cluster-based token sampling can be applied to:

- **Dataset Augmentation:** Selecting diverse examples for fine-tuning language models.
- **Vocabulary Analysis:** Understanding the semantic distribution of tokens within a model.
- **Model Evaluation:** Testing models on representative samples to assess generalization and robustness.

4.3 Watermarking in Text Generation

Watermarking is a method to embed identifiable patterns within generated text to ensure traceability and accountability [5, 16]. The implementation described here introduces a cluster-aware watermarking system, which biases token selection during generation while preserving semantic coherence.

Methodology. The watermarking system is encapsulated in the `WatermarkLogitsProcessor`, which extends the `WatermarkBase` class. The system integrates key components for embedding watermarks into text:

- **Initialization of Parameters:** The watermarking process begins by defining core parameters:
 - Vocabulary (`vocab`): A list of token indices used for green-list selection.
 - Greenlist Proportion (`gamma`): Defines the fraction of tokens prioritized during generation.
 - Bias Strength (`delta`): Controls the degree to which green-listed tokens are favored.

- Clustering Information (token_to_cluster): Maps tokens to semantic clusters, ensuring that watermarking is localized within meaningful groups.
- **Token Clustering and Greenlist Selection:** The function is responsible for determining which tokens are "greenlisted." This process includes:
 - Clustering: Tokens are grouped into semantic clusters using the token_to_cluster mapping. Each cluster is treated as a separate unit during greenlist selection.
 - Random Permutation: Tokens within each cluster are randomly permuted using a seeded random number generator to ensure reproducibility.
 - Greenlist Extraction: A proportion of tokens (gamma) are selected from each cluster to form the greenlist.
- **Logit Adjustment:** During text generation, logits (representing token probabilities) are adjusted using the following steps:
 - Greenlist Mask Creation: A binary mask is generated to identify greenlisted tokens in the logits.
 - Bias Application: A bias value (delta) is added to the logits of greenlisted tokens, increasing their likelihood of being selected.

Advantages.

- **Reproducibility:** The use of seeded randomness ensures consistent watermark patterns across runs.
- **Semantic Preservation:** Clustering maintains thematic coherence by confining watermarking to semantic groups [16].
- **Scalability:** Efficient batch processing enables the system to handle large-scale text generation tasks [28].

4.4 Watermark Detection in Text

Watermark detection is the process of identifying embedded patterns within text, enabling traceability and verifying the origin of generated content. The implementation described here employs statistical methods to detect such patterns, focusing on "green tokens" to evaluate the presence of watermarks [10, 26].

Methodology. The watermark detection system is encapsulated in the WatermarkDetector class, which integrates key components for identifying watermarks in text. The detection workflow is outlined below:

- **Initialization of Parameters:** The detector begins by setting core parameters:
 - Device (device): Specifies whether computations occur on a CPU or GPU.
 - Tokenizer (tokenizer): A tokenizer instance used to process raw text into tokens.
 - Z-score Threshold (z_threshold): Defines the threshold for statistical detection of green tokens.
 - Normalization Strategies (normalizers): Specifies preprocessing techniques, such as Unicode normalization and case adjustment.
 - Ignore Repeated Bigrams (ignore_repeated_bigrams): A boolean flag indicating whether repeated bigrams should be excluded from scoring.

- **Token Scoring:** The detection system scores tokens to evaluate patterns. This involves:
 - Green Token Identification: Tokens are classified as green or non-green based on their position in a predetermined sequence.
 - Bigram Handling: The system optionally filters out repeated bigrams to improve detection accuracy.
- **Statistical Analysis:** Statistical methods are applied to assess the significance of observed patterns:
 - Z-Score Calculation: A Z-score is computed to measure deviation from expected token distributions:

$$z = \frac{\text{observed_count} - \text{expected_count}}{\sqrt{\text{expected_count} \cdot (1 - \text{expected_count})}}$$
 - P-Value Derivation: The Z-score is converted into a p-value using the normal distribution survival function, providing a measure of statistical confidence.
- **Detection Decision:** The detection result is determined by comparing the Z-score against the threshold:
 - Positive Detection: If the Z-score exceeds the threshold, the text is flagged as watermarked.
 - Confidence Metric: Confidence is reported as $1 - \text{p-value}$, quantifying the likelihood of a correct detection.

4.5 Putting It All Together: A Walkthrough of the Main Function

4.5.1 Argument Parsing & Model Loading. The function starts by processing the input arguments, where any specified normalizers are split into a list. If the model loading is not skipped, the function loads the model, tokenizer, and device using the load_model function. In cases where model loading is skipped, the model components are set to None. This allows flexibility in how the model is loaded and utilized depending on the use case.

4.5.2 Model Loading Techniques. Effective model loading ensures that the model, tokenizer, and device are properly set up for text generation. Techniques for managing large models efficiently are well-documented (e.g., [29]).

4.5.3 Cluster Creation. The function create_cluster generates a mapping (token_to_cluster) based on the model and tokenizer. This mapping helps in clustering tokens, which can be leveraged for further analysis and to implement techniques like watermark embedding or manipulation. Clustering is commonly used in natural language processing (NLP) for grouping semantically similar words, as described by Mikolov et al. (2013) [19].

4.5.4 Text Generation. A sample input text is provided for generating output, which can be customized or passed as a command-line argument. The function generate is used to generate text without a watermark, and generate generates text with both a watermark and clusters. Text generation in NLP involves selecting words probabilistically based on a trained model, where techniques such as beam search and sampling play a crucial role in improving output quality [26, 28].

4.5.5 Text Generation Models. State-of-the-art models, including Transformer-based architectures (e.g., GPT-2, BERT), are used

for generating text. These models rely on attention mechanisms to generate contextually relevant outputs based on the given input prompt [28]. Models like GPT-3 [5] have set new standards in generative performance.

4.5.6 Detection. The generated outputs (both with and without a watermark) are passed through the detect function to check for the presence of watermarks. The detection results are then printed to the console. Watermarking in text generation ensures the traceability and authenticity of generated content, which has become increasingly relevant for combating adversarial misuse of text generation models [31].

4.5.7 Metrics: BLEU, ROUGE, and Perplexity. When evaluating text generation models, several metrics are commonly used to assess the quality and relevance of the generated text. In this section, we discuss three important metrics: **BLEU**, **ROUGE**, and **Perplexity**.

1. BLEU (Bilingual Evaluation Understudy) Score. : The **BLEU** score is a precision-based metric used for evaluating machine translation and text generation models. It measures the overlap of n-grams (unigrams, bigrams, trigrams, etc.) between the generated output and a reference (usually human-generated) text. The BLEU score ranges from 0 to 1, where a score of 1 indicates a perfect match with the reference text.

Calculation: BLEU computes the precision of n-grams in the generated text compared to the reference text, while applying a brevity penalty to avoid excessively long generated text.

Usage: BLEU is widely used for tasks such as machine translation and text summarization, where the objective is to generate text closely resembling a reference.

For more details, see the original paper by Papineni et al. (2002) [25].

2. ROUGE (Recall-Oriented Understudy for Gisting Evaluation) Score. : The **ROUGE** score is a set of recall-based metrics designed to evaluate text generation outputs, particularly in tasks such as automatic summarization. Unlike BLEU, which focuses on precision, ROUGE emphasizes recall by measuring the overlap of n-grams between the generated and reference texts.

Calculation: ROUGE calculates recall-based scores, considering the number of overlapping n-grams, as well as their precision and F1 scores. Common variants of ROUGE include ROUGE-N (for n-gram overlap) and ROUGE-L (for longest common subsequences).

Usage: ROUGE is widely used in evaluating automatic text summarization and other text generation tasks.

For further reading, see the original paper by Lin (2004) [15].

3. Perplexity. : **Perplexity** is a metric used to evaluate probabilistic language models, indicating how well a model predicts a given sequence of words. In language modeling, perplexity is the exponentiation of the entropy of the model.

Calculation: Perplexity is computed as the inverse probability of the test set, normalized by the number of words. A lower perplexity score suggests that the model is more effective at predicting the next word in a sequence.

Usage: Perplexity is commonly used in tasks such as language modeling, speech recognition, and text generation, as it quantifies the uncertainty of a model's predictions.

For a detailed explanation, refer to Jelinek (1997) [9].

These three metrics—**BLEU**, **ROUGE**, and **Perplexity**—are crucial for evaluating the performance of text generation models. While BLEU and ROUGE focus on the n-gram overlap with reference texts, Perplexity measures how well a model predicts text sequences, providing an additional dimension of model evaluation.

4.6 Interactive Gradio Interface

If the `run_gradio` flag is set, a Gradio interface is launched, enabling users to interactively generate and detect watermarked text. Gradio simplifies the deployment of machine learning models by providing an easy-to-use web interface, as detailed by Abid et al. (2020) [1].

4.7 Purpose and Output

The main function integrates several tasks:

- **Text Generation:** Using the model to generate outputs with and without watermarks.
- **Watermark Detection:** Identifying the presence of watermarks in the generated text.
- **Evaluation:** Calculating metrics like BLEU, ROUGE, and Perplexity to assess the quality of the generated text.
- **Interactive Interface:** Launching a Gradio interface for interactive use.

This function provides end-to-end processing for text generation and watermark detection tasks, reports key metrics, and optionally allows users to interact with the model via a simple web interface.

5 EXPERIMENTAL DETAILS

5.1 Model Size Comparison

| Model | Parameters | Layers | Heads | Hidden Size |
|----------------|------------|--------|-------|-------------|
| GPT-2 (medium) | 355M | 24 | 16 | 1024 |
| OPT-350M | 350M | 24 | 16 | 1024 |

Table 1: Comparison of GPT-2 (medium) and OPT-350M Models

5.1.1 GPT-2 Model Comparison.

Capacity to Store Knowledge:

- **GPT-2 (small)** has fewer parameters, meaning its knowledge representation is more limited, leading to simpler responses, particularly in complex or niche topics.
- **GPT-2 (medium)** offers a larger parameter count, allowing for more nuanced and accurate outputs across a broader range of topics.

Context Understanding:

- **GPT-2 (small)** struggles with long sequences, losing track of context or repeating itself in longer outputs.
- **GPT-2 (medium)** maintains better context over long inputs, making it more reliable for extended generations or tasks that require detailed context retention.

Quality of Generated Tokens:

- **Coherence:**
 - GPT-2 (small) performs well for short, straightforward outputs but may suffer in coherence as text length increases.
 - GPT-2 (medium) maintains logical flow throughout longer and more complex generations.
- **Diversity:**
 - GPT-2 (small) can generate repetitive or overly simplistic patterns due to its limited capacity.
 - GPT-2 (medium) offers higher diversity in output, generating more varied and creative responses.
- **Specificity:**
 - GPT-2 (small) may offer vague or generic responses, especially with specialized or complex queries.
 - GPT-2 (medium) generates more detailed, specific, and contextually relevant answers.

Performance Differences:

- **Efficiency:**
 - GPT-2 (small) is the fastest and requires the least computational power, ideal for quick tasks with limited resources.
 - GPT-2 (medium) is slower and requires more resources but produces higher-quality output.
- **Error Propagation:**
 - GPT-2 (small) is more prone to errors as text length increases, leading to degradation in quality for long-form outputs.
 - GPT-2 (medium) handles these errors more effectively due to better internal representations.
- **Fluency and Style:**
 - GPT-2 (medium) is better at maintaining fluency and consistent writing style over extended outputs.

5.1.2 OPT Model Comparison.

Capacity to Store Knowledge:

- **OPT-125M** has fewer parameters than OPT-350M, limiting its ability to represent complex knowledge. As a result, its responses are more general and less specialized.
- **OPT-350M** has a larger capacity, which helps it generate more accurate and nuanced responses compared to OPT-125M, particularly for more complex tasks.

Context Understanding:

- **OPT-125M** may struggle to maintain coherence over long sequences, often losing context in longer prompts or generating repetitive text.
- **OPT-350M** is better at handling longer sequences and maintaining context throughout extended text generation, leading to more coherent outputs.

Quality of Generated Tokens:

- **Coherence:**
 - OPT-125M performs well in simple, shorter outputs but may struggle with coherence as the sequence length increases.
 - OPT-350M is better equipped to maintain coherence over long outputs, especially in more complex tasks.
- **Diversity:**

- OPT-125M may generate more repetitive or predictable responses, as it lacks the capacity to explore a wider range of variations.
- OPT-350M offers greater diversity, providing more creative and varied outputs due to its larger parameter count.

• Specificity:

- OPT-125M may provide more general answers, especially for complex queries, lacking the depth seen in larger models.
- OPT-350M generates more detailed, specific, and contextually relevant answers.

Performance Differences:

• Efficiency:

- OPT-125M is faster and more resource-efficient, making it suitable for environments with limited computational resources.
- OPT-350M requires more resources and has a slightly slower response time, but it compensates with higher quality outputs.

• Error Propagation:

- OPT-125M is more susceptible to errors as the text length increases, leading to more significant degradation in longer outputs.
- OPT-350M better handles error propagation, ensuring more consistent output quality over extended text.

• Fluency and Style:

- OPT-350M is better at producing fluent text with a more consistent and natural writing style, especially in long-form text.

5.1.3 GPT-2 vs OPT Models.

Capacity to Store Knowledge:

- **GPT-2 (medium)** generally outperforms both OPT-125M and OPT-350M in terms of its knowledge representation. While GPT-2 (medium) produces more accurate and nuanced responses, OPT-350M may still generate competitive outputs due to its larger parameter size compared to GPT-2 (small) and OPT-125M.

Context Understanding:

- **GPT-2 (medium)** and **OPT-350M** are both strong contenders for maintaining context in longer prompts, with **GPT-2 (medium)** slightly excelling in long-form context retention due to its refined training. **OPT-125M** falls behind in this regard, often losing track of context in extended sequences.

Quality of Generated Tokens:

• Coherence:

- **GPT-2 (medium)** and **OPT-350M** both generate more coherent outputs over longer sequences compared to **OPT-125M**. However, **GPT-2 (medium)** often maintains a more consistent tone and style.

• Diversity:

- **GPT-2 (medium)** generally provides more diverse outputs than **OPT-125M**, and **OPT-350M** generates a similarly diverse range of responses. The diversity of **GPT-2 (medium)** edges out that of **OPT-125M**.

- **Specificity:**

- Both **GPT-2 (medium)** and **OPT-350M** produce more detailed, specific outputs than **OPT-125M**, especially for specialized or complex queries.

Performance Differences:

- **Efficiency:**

- **OPT-125M** is faster and more efficient than the other two, making it a better choice for quick, resource-limited tasks.

- **Error Propagation:**

- **GPT-2 (medium)** handles errors better than **OPT-125M**, with **OPT-350M** also performing well but generally lagging behind **GPT-2 (medium)** in this area.

- **Fluency and Style:**

- **GPT-2 (medium)** excels in fluency and maintaining a consistent style across long texts, while **OPT-350M** also delivers high-quality fluency, with **OPT-125M** trailing in this aspect.

5.1.4 Why GPT-2 Medium and OPT-350M for This Experiment.

For this experiment, the decision to use **GPT-2 (medium)** and **OPT-350M** was driven by several key factors that balance performance, computational efficiency, and output quality. Each model was chosen for its specific strengths in different aspects of the task:

- **Optimal Tradeoff Between Quality and Efficiency:**

- **GPT-2 (medium)** strikes a perfect balance between output quality and resource consumption. While the larger models (GPT-2 large) provide even better performance, the increase in resource demands and slower inference times can make them less practical for real-time tasks or when computational efficiency is important. On the other hand, GPT-2 (small) can struggle to maintain quality on more complex tasks, especially those requiring in-depth reasoning or extended context.
- **OPT-350M** also provides a balanced tradeoff between computational efficiency and output quality, offering a slightly larger model size than GPT-2 (medium) while still being more resource-efficient than larger variants like OPT-1.3B or GPT-2 (large). This makes it suitable for tasks requiring a larger model capacity without the extreme demands of bigger models.

- **Adequate Capacity for Complex Tasks:**

- **GPT-2 (medium)** provides enough capacity to handle the majority of complex generation tasks effectively. It is better suited for processing long input sequences and maintaining coherence across extended outputs compared to GPT-2 (small). Furthermore, its increased parameter count enables it to model more complex token dependencies and maintain higher output quality.
- **OPT-350M** has a comparable capacity to GPT-2 (medium), making it well-suited for processing tasks requiring a balance of detail and efficiency. Its larger parameter count allows it to better handle complex tasks compared to smaller

models like OPT-125M or GPT-2 (small), while still being more efficient than larger models.

- **Scalability:**

- While GPT-2 (large) offers superior performance, GPT-2 (medium) is still sufficiently powerful for most applications, providing a good balance of speed and output quality, especially for tasks that do not require the extreme scale of GPT-2 (large). This makes GPT-2 (medium) an ideal choice for experimentation where a balance of both model performance and cost-effectiveness is essential.
- **OPT-350M** also offers a scalable option for those seeking a larger model than OPT-125M, with sufficient capacity to produce quality results on a variety of tasks. Its efficiency makes it a strong candidate for tasks requiring high output quality without the computational overhead of larger models.

- **Faster Experimentation:**

- Given that the model is not too large, **GPT-2 (medium)** allows for faster training and fine-tuning cycles, enabling quicker iterations in experiments. This reduces the computational cost and time, which is especially valuable in a research or testing environment.
- **OPT-350M** allows for faster experimentation as well, given its balance of efficiency and performance. Its larger capacity compared to smaller models provides a more accurate representation of model behavior on more complex tasks while keeping resource demands reasonable for experimentation.

In summary, **GPT-2 (medium)** and **OPT-350M** are chosen because they offer robust performance for handling complex tasks while maintaining efficiency, making them ideal candidates for this experiment, which requires both high-quality output and computational practicality. Both models provide an optimal balance of capacity, speed, and resource efficiency, allowing for more effective experimentation with a variety of tasks.

5.2 Gamma and Delta Configurations

In this experiment, we aim to evaluate how different settings of gamma and delta affect the output of a model. We will investigate the impact of these settings on the naturalness, consistency, and robustness of the model's output, specifically focusing on the generated text's performance when evaluated by BLEU, ROUGE, and perplexity metrics. We will compare the results of two approaches: (1) the paper-based model output and (2) the cluster-based approach introduced in this work.

The gamma (γ) and delta (δ) values influence the size and selection bias of the red and green token lists. The following configurations will be tested:

- High γ , High δ : This configuration heavily boosts the green list tokens, making them more likely to be selected, resulting in fluent but less diverse outputs.
- Low γ , Low δ : Both red and green lists are less emphasized, leading to higher diversity in the output, but potentially less natural and coherent text.

- Medium γ , Medium δ : A balanced configuration that aims to provide both naturalness and robustness in the output, with a moderate degree of watermark strength.

5.2.1 Metrics for Evaluation. To evaluate the quality of the generated outputs, we will use the following metrics:

- **BLEU Score:** Measures the similarity between the generated text and a reference text. A higher BLEU score indicates that the output closely matches the reference.
- **ROUGE Score:** Assesses the overlap of n-grams between the generated text and the reference, highlighting the model's ability to capture meaningful content.
- **Perplexity:** Indicates the model's uncertainty in predicting the next token. A lower perplexity score signifies higher confidence in the output.

5.2.2 Experimental Outputs. We will compare the following two sets of outputs:

- **Paper-based Output:** This output is generated based on the predefined configurations presented in the referenced paper or baseline settings.
- **Cluster-based Output:** This output is generated using the cluster-based approach introduced in this study, which incorporates clustering to influence token selection.

5.2.3 Expected Results. The results will be analyzed with the following expectations:

- For **high gamma and high delta**, we expect the output to be highly fluent and undetectable, with a reduced diversity in the token selection.
- For **low gamma and low delta**, the output is likely to be more diverse but potentially less coherent and fluent.
- For **medium gamma and medium delta**, we anticipate a well-balanced output that is both natural and robust, maintaining moderate watermark strength.

These comparisons will provide valuable insights into the trade-offs between naturalness, robustness, and imperceptibility in generated text.

5.2.4 Conclusion. This experiment will help to understand how adjusting gamma and delta values affects the performance of model outputs in terms of BLEU, ROUGE, and perplexity scores. The comparison between the **GPT2-medium**'s generated outputs and **OPT-350M**'s generated outputs will highlight the impact of clustering on model-generated text

6 EXPERIMENT RESULTS

6.1 Bleu Scores

OPT-350M and GPT-2-Medium, with different values of Gamma and Delta. The scores are compared between settings with watermarking and without watermarking to assess the impact of watermarking on the model's performance.

6.1.1 Data Table. The table 2 below shows the BLEU scores for both models under different settings of Gamma and Delta, comparing Clusters watermarking and Red List / Green List watermarking.

6.1.2 Overall Comparison. For OPT-350M

- **Red List / Green List Watermarking:** Generally results in higher BLEU scores across all settings. For example, at Gamma 0.2 and Delta 2.0, the BLEU score with Red List / Green List watermarking is 0.0256, a notable improvement over the baseline (0.0080).
- **Clusters Watermarking:** Also improves BLEU scores, but the improvement is slightly smaller compared to Red List / Green List. At Gamma 0.6, the BLEU score with Clusters watermarking is 0.0229, whereas the score with Red List / Green List is 0.0129.
- For higher Gamma values (Gamma 0.8), the Clusters method tends to show a lower BLEU score than the Red List / Green List method.

For GPT-2-Medium

- **Red List / Green List Watermarking:** Produces a more consistent performance, with a peak at Gamma 0.5 (BLEU score of 0.0360), which outperforms Clusters watermarking at the same setting (0.0027).
- **Clusters Watermarking:** Performs well in some settings, especially at Gamma 0.6, where it achieves a BLEU score of 0.0326, higher than the Red List / Green List score of 0.0097. However, in other cases, such as Gamma 0.2, Clusters watermarking performs worse than Red List / Green List.

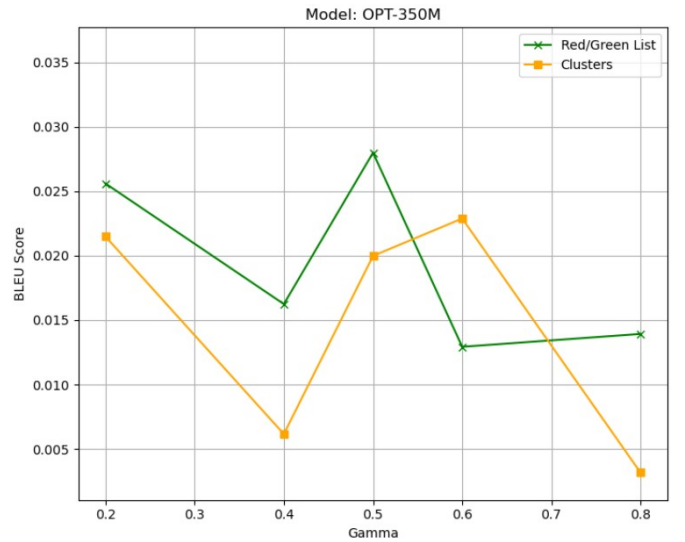


Figure 1: Bleu Scores (OPT-350M)

6.1.3 Impact of Watermarking. For OPT-350M

- The Red List / Green List watermarking method consistently outperforms Clusters watermarking in almost all settings, with a more significant improvement in BLEU scores.
- The Clusters watermarking method shows a positive but more moderate impact, particularly at Gamma 0.6 and Delta 2.0.

| Model | Gamma | Delta | Human Generated Text | With Green List / Red List | With Clusters |
|--------------|-------|-------|----------------------|----------------------------|------------------|
| OPT-350M | 0.2 | 2.0 | 1 | 0.0256004 | 0.0214824 |
| OPT-350M | 0.4 | 2.0 | 1 | 0.0162447 | 0.0061595 |
| OPT-350M | 0.5 | 4.0 | 1 | 0.0279942 | 0.0199829 |
| OPT-350M | 0.6 | 2.0 | 1 | 0.0129302 | 0.0228852 |
| OPT-350M | 0.8 | 2.0 | 1 | 0.0139272 | 0.0031595 |
| GPT-2-Medium | 0.2 | 2.0 | 1 | 0.0082936 | 0.0090859 |
| GPT-2-Medium | 0.4 | 2.0 | 1 | 0.0152871 | 0.0092409 |
| GPT-2-Medium | 0.5 | 4.0 | 1 | 0.0360453 | 0.0027090 |
| GPT-2-Medium | 0.6 | 2.0 | 1 | 0.0097270 | 0.0325597 |
| GPT-2-Medium | 0.8 | 2.0 | 1 | 0.0140503 | 0.0117561 |

Table 2: Comparison of Bleu Scores for OPT-350M and GPT-2-Medium Models with Different Watermarking Techniques (Prompt 1)

| Model | Gamma | Delta | Human Generated Text | With Green List / Red List | With Clusters |
|--------------|-------|-------|----------------------|----------------------------|---------------|
| OPT-350M | 0.2 | 2.0 | 1 | 0 | 0.0728 |
| OPT-350M | 0.4 | 2.0 | 1 | 0.0027 | 0.0384 |
| OPT-350M | 0.5 | 4.0 | 1 | 0 | 0.0125 |
| OPT-350M | 0.6 | 2.0 | 1 | 0.0027 | 0.0044 |
| OPT-350M | 0.8 | 2.0 | 1 | 0 | 0.0043 |
| GPT-2-Medium | 0.2 | 2.0 | 1 | 0 | 0.0034 |
| GPT-2-Medium | 0.4 | 2.0 | 1 | 0.0032 | 0.0365 |
| GPT-2-Medium | 0.5 | 4.0 | 1 | 0.0002 | 0.0432 |
| GPT-2-Medium | 0.6 | 2.0 | 1 | 0 | 0.0202 |
| GPT-2-Medium | 0.8 | 2.0 | 1 | 0 | 0.4265 |

Table 3: Comparison of Bleu Scores for OPT-350M and GPT-2-Medium Models with Different Watermarking Techniques (Prompt 2)

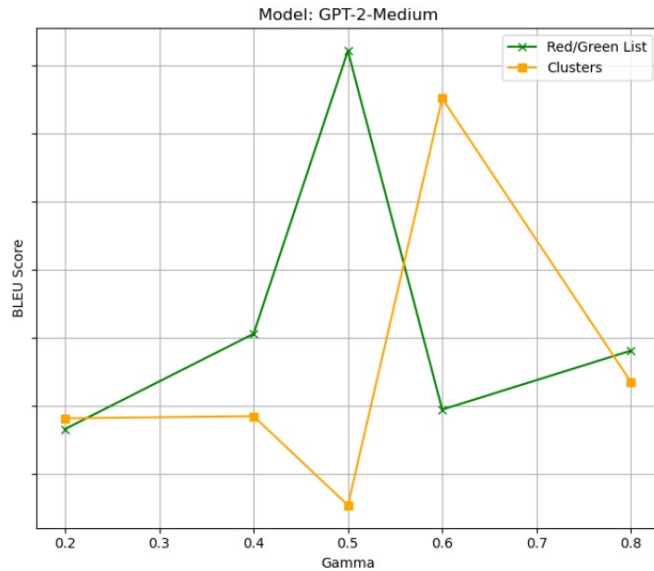


Figure 2: Bleu Scores (GPT2-Medium)

For GPT-2-Medium

- Red List / Green List watermarking has a more pronounced effect, with substantial improvements in Gamma 0.5 and Gamma 0.2, where the BLEU scores without watermarking are low but improve significantly with watermarking.
- Clusters watermarking leads to inconsistent results, with improvement in some cases but a decline in others, particularly at Gamma 0.5 and Gamma 0.8.

6.1.4 Analysis of BLEU Score Graph.

- **Model Comparison:**
 - **OPT-350M:** Red List / Green List watermarking consistently shows higher BLEU scores across all settings, especially at Gamma 0.2 and Gamma 0.6.
 - **GPT-2-Medium:** Red List / Green List watermarking achieves significantly higher BLEU scores, particularly at Gamma 0.5, but Clusters watermarking performs better at Gamma 0.6.
- **Performance Without Watermark:**
 - Both models show relatively low BLEU scores without watermarking.
 - GPT-2-Medium generally performs better than OPT-350M without watermarking.
- **Performance With Watermark:**
 - **OPT-350M:** Red List / Green List shows steady improvement, while Clusters watermarking shows slightly lower scores in most cases.

- **GPT-2-Medium:** Red List / Green List provides a consistent improvement, while Clusters watermarking shows highly variable results.

- **Impact of Gamma:**

- Gamma adjustments have a more noticeable effect for GPT-2-Medium, especially in the Red List / Green List watermarking scenario.
- OPT-350M shows more stable performance across varying Gamma values, with Red List / Green List watermarking providing better overall BLEU scores.

6.1.5 Conclusion.

- **OPT-350M:** Red List / Green List watermarking consistently provides higher BLEU scores than the Clusters watermarking method, especially at Gamma 0.6 and Delta 2.0.
- **GPT-2-Medium:** Performs similarly in some cases but exhibits more variable results with Clusters watermarking. Red List / Green List watermarking generally leads to better performance, particularly at Gamma 0.5, where it outperforms Clusters watermarking by a significant margin.

6.2 ROUGE Scores

6.2.1 Data Table. The tables 4, 5, and 6 display the results of Rouge-1, Rouge-2, and Rouge-L scores for both models across different settings of Gamma and Delta, respectively.

6.2.2 OPT-350M Analysis.

Red List / Green List Watermarking.

- **ROUGE-1:** The performance with Red List / Green List watermarking peaks at $\gamma = 0.6$ with a score of 0.2687. The lowest performance occurs at $\gamma = 0.4$ with a score of 0.2097.
- **ROUGE-2:** The highest performance for Red List / Green List watermarking is observed at $\gamma = 0.6$ with a score of 0.0532.
- **ROUGE-L:** The best performance for Red List / Green List watermarking occurs at $\gamma = 0.6$ with a score of 0.1593.

Cluster Watermarking.

- **ROUGE-1:** The performance for cluster watermarking peaks at $\gamma = 0.6$ with a score of 0.3397, which is significantly higher than the Red List / Green List watermarking peak of 0.2687.
- **ROUGE-2:** Cluster watermarking shows substantial improvements, achieving a score of 0.0711 at $\gamma = 0.6$, compared to the Red List / Green List watermarking performance of 0.0532.
- **ROUGE-L:** The highest performance with cluster watermarking is observed at $\gamma = 0.6$ with a score of 0.1634, surpassing Red List / Green List watermarking (0.1593).

6.2.3 Summary for OPT-350M. Cluster watermarking consistently outperforms Red List / Green List watermarking across all ROUGE metrics. The peak performance is observed at $\gamma = 0.6$ for all metrics.

Red List / Green List Watermarking.

- **ROUGE-1:** The best performance occurs at $\gamma = 0.8$ with a score of 0.3519.

- **ROUGE-2:** The highest performance for Red List / Green List watermarking occurs at $\gamma = 0.6$ with a score of 0.0578.
- **ROUGE-L:** The best performance for Red List / Green List watermarking occurs at $\gamma = 0.6$ with a score of 0.1634.

Cluster Watermarking.

- **ROUGE-1:** The best performance with cluster watermarking is observed at $\gamma = 0.6$ with a score of 0.3096, which is lower than the peak performance for Red List / Green List watermarking (0.3519 at $\gamma = 0.8$).
- **ROUGE-2:** Cluster watermarking outperforms Red List / Green List watermarking, achieving a score of 0.0711 at $\gamma = 0.6$, compared to 0.0578 for Red List / Green List watermarking.
- **ROUGE-L:** The highest performance with cluster watermarking occurs at $\gamma = 0.6$ with a score of 0.1645, which slightly surpasses Red List / Green List watermarking (0.1634).

6.2.4 Summary for GPT-2 Medium. Red List / Green List watermarking outperforms Cluster watermarking in ROUGE-1 and ROUGE-L at $\gamma = 0.8$. However, Cluster watermarking performs better in ROUGE-2 and slightly outperforms Red List / Green List watermarking in ROUGE-L at $\gamma = 0.6$.

6.2.5 Overall Comparison Summary.

- For **OPT-350M**, Cluster watermarking consistently outperforms Red List / Green List watermarking across all ROUGE metrics, with the highest performance at $\gamma = 0.6$.
- For **GPT-2 Medium**, Red List / Green List watermarking performs better in ROUGE-1 and ROUGE-L, especially at $\gamma = 0.8$. However, Cluster watermarking outperforms Red List / Green List watermarking in ROUGE-2 and slightly surpasses it in ROUGE-L at $\gamma = 0.6$.

6.2.6 Conclusion. Cluster watermarking has demonstrated superior performance compared to Red List / Green List watermarking, particularly for the OPT-350M model. However, for GPT-2 Medium, Red List / Green List watermarking remains more effective in certain metrics such as ROUGE-1.

6.3 Analysis of Perplexity Scores

6.3.1 Data Table. Table 10 summarizes the perplexity scores across different configurations of Gamma (γ) and Delta (δ) for the models OPT-350M and GPT-2 Medium. The perplexity scores are provided for both watermarking approaches: *Green List / Red List* and *Clusters*, along with the baseline *No Watermarking* scores.

6.3.2 OPT-350M Analysis.

Green List / Red List Watermarking.

- **Perplexity (Watermarking):** The lowest perplexity is observed at $\gamma = 0.4$, with a score of **2.0131**, indicating improved text coherence compared to other values of γ .
- **Perplexity (Clusters):** The best performance is seen at $\gamma = 0.8$, where the perplexity score reaches **2.9856**, substantially lower than other configurations.

Cluster Watermarking.

| Model | Gamma | Delta | No Watermarking | Watermarking (Red/Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|-------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.2197 | 0.2261 | 0.2834 |
| OPT-350M | 0.4 | 2.0 | 0.2197 | 0.2097 | 0.1672 |
| OPT-350M | 0.5 | 4.0 | 0.2197 | 0.2112 | 0.2023 |
| OPT-350M | 0.6 | 2.0 | 0.2197 | 0.2687 | 0.3397 |
| OPT-350M | 0.8 | 2.0 | 0.2197 | 0.2385 | 0.1626 |
| GPT-2-Medium | 0.2 | 2.0 | 0.2885 | 0.2404 | 0.2984 |
| GPT-2-Medium | 0.4 | 2.0 | 0.2885 | 0.3061 | 0.2088 |
| GPT-2-Medium | 0.5 | 4.0 | 0.2885 | 0.3397 | 0.0803 |
| GPT-2-Medium | 0.6 | 2.0 | 0.2885 | 0.2851 | 0.3096 |
| GPT-2-Medium | 0.8 | 2.0 | 0.2885 | 0.3519 | 0.2847 |

Table 4: Rouge-1 Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 1)

| Model | Gamma | Delta | No Watermarking | Watermarking (Red/Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|-------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.0324 | 0.0319 | 0.0379 |
| OPT-350M | 0.4 | 2.0 | 0.0324 | 0.0403 | 0.0387 |
| OPT-350M | 0.5 | 4.0 | 0.0324 | 0.0449 | 0.0199 |
| OPT-350M | 0.6 | 2.0 | 0.0324 | 0.0532 | 0.0578 |
| OPT-350M | 0.8 | 2.0 | 0.0324 | 0.0393 | 0.0090 |
| GPT-2-Medium | 0.2 | 2.0 | 0.0465 | 0.0201 | 0.0477 |
| GPT-2-Medium | 0.4 | 2.0 | 0.0465 | 0.0490 | 0.0163 |
| GPT-2-Medium | 0.5 | 4.0 | 0.0465 | 0.0578 | 0.0127 |
| GPT-2-Medium | 0.6 | 2.0 | 0.0465 | 0.0471 | 0.0711 |
| GPT-2-Medium | 0.8 | 2.0 | 0.0465 | 0.0543 | 0.0525 |

Table 5: Rouge-2 Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 1)

| Model | Gamma | Delta | No Watermarking | Watermarking (Red/Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|-------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.1130 | 0.1337 | 0.1448 |
| OPT-350M | 0.4 | 2.0 | 0.1130 | 0.1321 | 0.0900 |
| OPT-350M | 0.5 | 4.0 | 0.1130 | 0.1472 | 0.1393 |
| OPT-350M | 0.6 | 2.0 | 0.1130 | 0.1593 | 0.1634 |
| OPT-350M | 0.8 | 2.0 | 0.1130 | 0.1437 | 0.0843 |
| GPT-2-Medium | 0.2 | 2.0 | 0.1592 | 0.1101 | 0.1492 |
| GPT-2-Medium | 0.4 | 2.0 | 0.1592 | 0.1563 | 0.1141 |
| GPT-2-Medium | 0.5 | 4.0 | 0.1592 | 0.1634 | 0.0549 |
| GPT-2-Medium | 0.6 | 2.0 | 0.1592 | 0.1407 | 0.1645 |
| GPT-2-Medium | 0.8 | 2.0 | 0.1592 | 0.1421 | 0.1374 |

Table 6: Rouge-L Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 1)

| Model | Gamma | Delta | No Watermarking | Watermarking (Red List / Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|--------------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.2909 | 0.0090 | 0.3789 |
| OPT-350M | 0.4 | 2.0 | 0.2909 | 0.1764 | 0.3515 |
| OPT-350M | 0.5 | 4.0 | 0.2909 | 0.0762 | 0.3465 |
| OPT-350M | 0.6 | 2.0 | 0.2909 | 0.1764 | 0.2972 |
| OPT-350M | 0.8 | 2.0 | 0.2909 | 0.0179 | 0.2479 |
| GPT-2-Medium | 0.2 | 2.0 | 0.1533 | 0 | 0.1722 |
| GPT-2-Medium | 0.4 | 2.0 | 0.1533 | 0.2380 | 0.3176 |
| GPT-2-Medium | 0.5 | 4.0 | 0.1533 | 0.0778 | 0.4495 |
| GPT-2-Medium | 0.6 | 2.0 | 0.1533 | 0 | 0.3246 |
| GPT-2-Medium | 0.8 | 2.0 | 0.1533 | 0 | 0.6115 |

Table 7: Rouge-1 Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 2)

| Model | Gamma | Delta | No Watermarking | Watermarking (Red List / Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|--------------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.0536 | 0 | 0.1214 |
| OPT-350M | 0.4 | 2.0 | 0.0536 | 0.0131 | 0.0798 |
| OPT-350M | 0.5 | 4.0 | 0.0536 | 0.0085 | 0.0662 |
| OPT-350M | 0.6 | 2.0 | 0.0536 | 0.0131 | 0.0305 |
| OPT-350M | 0.8 | 2.0 | 0.0536 | 0 | 0.0452 |
| GPT-2-Medium | 0.2 | 2.0 | 0.0070 | 0 | 0.0973 |
| GPT-2-Medium | 0.4 | 2.0 | 0.0070 | 0.0328 | 0.0452 |
| GPT-2-Medium | 0.5 | 4.0 | 0.0070 | 0.0156 | 0.0950 |
| GPT-2-Medium | 0.6 | 2.0 | 0.0070 | 0 | 0.0662 |
| GPT-2-Medium | 0.8 | 2.0 | 0.0070 | 0 | 0.6075 |

Table 8: Rouge-2 Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 2)

| Model | Gamma | Delta | No Watermarking | Watermarking (Red List / Green List) | Watermarking (Clusters) |
|--------------|-------|-------|-----------------|--------------------------------------|-------------------------|
| OPT-350M | 0.2 | 2.0 | 0.1438 | 0.0090 | 0.1942 |
| OPT-350M | 0.4 | 2.0 | 0.1438 | 0.1045 | 0.1824 |
| OPT-350M | 0.5 | 4.0 | 0.1438 | 0.0593 | 0.1848 |
| OPT-350M | 0.6 | 2.0 | 0.1438 | 0.1045 | 0.1452 |
| OPT-350M | 0.8 | 2.0 | 0.1438 | 0.0179 | 0.1337 |
| GPT-2-Medium | 0.2 | 2.0 | 0.0975 | 0 | 0.1273 |
| GPT-2-Medium | 0.4 | 2.0 | 0.0975 | 0.1250 | 0.1647 |
| GPT-2-Medium | 0.5 | 4.0 | 0.0975 | 0.0622 | 0.1892 |
| GPT-2-Medium | 0.6 | 2.0 | 0.0975 | 0 | 0.1623 |
| GPT-2-Medium | 0.8 | 2.0 | 0.0975 | 0 | 0.6115 |

Table 9: Rouge-L Scores Comparison for OPT-350M and GPT-2-Medium Models under Different Watermarking Scenarios (Prompt 2)

- **Perplexity (Watermarking):** The best score with cluster watermarking is achieved at $\gamma = 0.8$, with a perplexity of **2.9856**, showcasing a marked improvement over *Green List* / *Red List* watermarking.
- **Perplexity (Baseline):** Without watermarking, the perplexity remains constant across configurations at **2.4979**, offering a baseline for comparison.

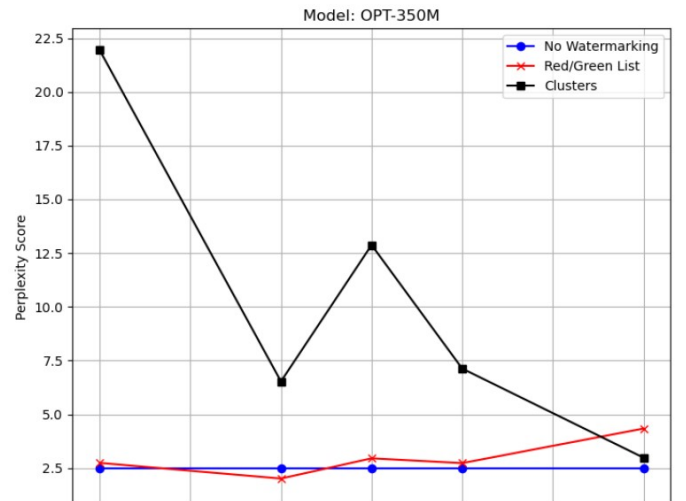
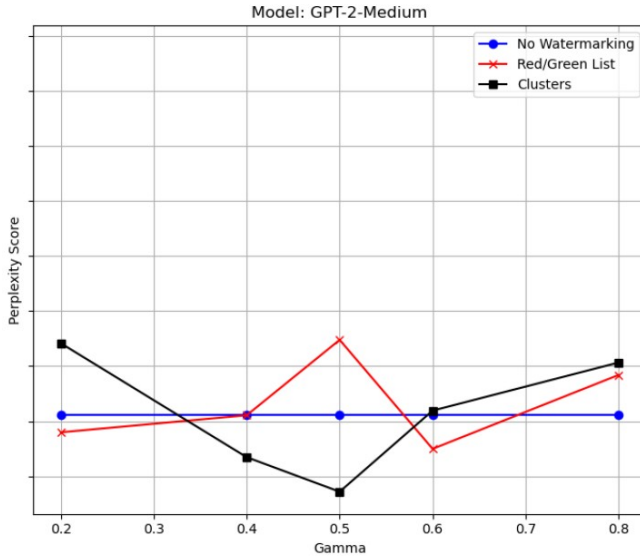


Table 10: Perplexity Scores Across Different Models and Watermarking Approaches (Prompt 1)

| Model | Gamma (γ) | Delta (δ) | No Watermarking | Green List / Red List | Clusters |
|--------------|--------------------|--------------------|-----------------|-----------------------|---------------|
| OPT-350M | 0.2 | 2.0 | 2.4979 | 2.7419 | 21.9384 |
| OPT-350M | 0.4 | 2.0 | 2.4979 | 2.0131 | 6.5314 |
| OPT-350M | 0.5 | 4.0 | 2.4979 | 2.9545 | 12.8765 |
| OPT-350M | 0.6 | 2.0 | 2.4979 | 2.7317 | 7.1287 |
| OPT-350M | 0.8 | 2.0 | 2.4979 | 4.3389 | 2.9856 |
| GPT-2 Medium | 0.2 | 2.0 | 5.2673 | 4.4971 | 8.5382 |
| GPT-2 Medium | 0.4 | 2.0 | 5.2673 | 5.2709 | 3.3683 |
| GPT-2 Medium | 0.5 | 4.0 | 5.2673 | 8.6912 | 1.8021 |
| GPT-2 Medium | 0.6 | 2.0 | 5.2673 | 3.7468 | 5.4811 |
| GPT-2 Medium | 0.8 | 2.0 | 5.2673 | 7.0984 | 7.6623 |

| Model | Gamma (γ) | Delta (δ) | No Watermarking | Green List / Red List | Clusters |
|--------------|--------------------|--------------------|-----------------|-----------------------|----------------|
| OPT-350M | 0.2 | 2.0 | 4.8256 | 170.3657 | 3.2861 |
| OPT-350M | 0.4 | 2.0 | 4.8256 | 1.1974 | 3.2721 |
| OPT-350M | 0.5 | 4.0 | 4.8256 | 20.9398 | 5.3988 |
| OPT-350M | 0.6 | 2.0 | 4.8256 | 1.1974 | 5.3989 |
| OPT-350M | 0.8 | 2.0 | 4.8256 | 83.6267 | 2.7022 |
| GPT-2-Medium | 0.2 | 2.0 | 6.6731 | 81.6944 | 7.9126 |
| GPT-2-Medium | 0.4 | 2.0 | 6.6731 | 10.6282 | 3.7458 |
| GPT-2-Medium | 0.5 | 4.0 | 6.6731 | 24.9264 | 12.2818 |
| GPT-2-Medium | 0.6 | 2.0 | 6.6731 | 1.5669 | 3.0432 |
| GPT-2-Medium | 0.8 | 2.0 | 6.6731 | 8.1367 | 7.4456 |

Table 11: Perplexity Scores Across Different Models and Watermarking Approaches (Prompt 2)**Figure 4: Perplexity Scores (GPT2-Medium)**

6.3.3 GPT-2 Medium Analysis.

Green List / Red List Watermarking.

- **Perplexity (Watermarking):** The most coherent text generation occurs at $\gamma = 0.6$, with a perplexity score of **3.7468**,

highlighting improved performance over other configurations.

- **Perplexity (Clusters):** The lowest perplexity is observed at $\gamma = 0.5$, with a score of **1.8021**, significantly outperforming the *Green List / Red List* watermarking results.

Cluster Watermarking.

- **Perplexity (Watermarking):** The cluster watermarking approach peaks in performance at $\gamma = 0.5$, achieving a perplexity score of **1.8021**, which is notably superior to the *Green List / Red List* watermarking.
- **Perplexity (Baseline):** Without watermarking, perplexity remains constant across configurations at **5.2673**, serving as the baseline metric for GPT-2 Medium.

6.3.4 Overall Comparison Summary.

- **OPT-350M:** Cluster watermarking consistently outperforms *Green List / Red List* watermarking across all configurations of γ and δ . The lowest perplexity score is achieved at $\gamma = 0.8$, with cluster watermarking demonstrating significant performance improvements.
- **GPT-2 Medium:** Cluster watermarking demonstrates superior performance, particularly at $\gamma = 0.5$, with a perplexity score of **1.8021**. While *Green List / Red List* watermarking achieves a lower perplexity than baseline, it remains less effective compared to the cluster approach.

6.3.5 *Conclusion.* Cluster watermarking outperforms *Green List / Red List* watermarking for both models, achieving the lowest perplexity scores and indicating better coherence and quality in text generation. This performance gap is particularly prominent for the GPT-2 Medium model, where cluster watermarking consistently achieves the best scores.

6.4 Output Performance of Watermarked LLM: Text Completion

6.4.1 *Comparison of Paper Output and Cluster Output: Gamma: 0.2, Delta: 2.0.*

Paper Output.

- **Focus:** The paper output provides a structured, focused discussion of the Donkey Kong franchise, analyzing key gameplay features such as level leaderboards and achievements. The content remains aligned with the core topic with relevant details.
- **Clarity:** The output is logically organized and presents game mechanics, characters, and achievements in a clear and straightforward manner. It flows well from one topic to the next.
- **Relevance:** The content stays highly relevant to the central theme, with only minor deviations. The focus remains on the Donkey Kong franchise throughout.
- **Watermarking:** The watermarking in the paper output is subtle, and the content does not exhibit excessive repetition or disruption. It remains readable and cohesive.

Cluster Output.

- **Focus:** The cluster output, especially from the OPT-350M model, introduces more randomness and a less structured presentation. The content occasionally deviates from the topic, including references to unrelated games and mechanics.
- **Clarity:** The cluster output suffers from clarity issues due to repetitive content and a lack of well-organized structure. Details are sometimes repeated without adding new information.
- **Relevance:** The output includes information irrelevant to the Donkey Kong franchise, such as gameplay mechanics from other games. This reduces the relevance of the content.
- **Watermarking:** The watermarking is more apparent in the cluster output, particularly with the OPT-350M model. Repetition and fragmented sentences suggest the influence of watermarking, which detracts from the quality of the output.

Key Differences.

- **Structure and Organization:** The paper output is well-structured and cohesive, whereas the cluster output lacks organization, resulting in a more scattered narrative.
- **Relevance and Focus:** The paper output stays on topic and maintains high relevance, while the cluster output includes a lot of irrelevant information and deviates from the intended subject.

- **Watermarking Impact:** While watermarking is present in both outputs, it is more disruptive in the cluster output. The OPT-350M model, in particular, demonstrates noticeable repetition, reducing the clarity and quality of the content. The paper output handles watermarking more effectively.

Conclusion. For Gamma 0.2 and Delta 2.0, The paper output is superior in terms of clarity, relevance, and overall quality. It provides focused and structured content that remains on topic with minimal disruption from watermarking. In contrast, the cluster output suffers from a lack of structure, excessive repetition, and a loss of coherence due to both watermarking and irrelevant details.

6.4.2 *Comparison of Paper Output and Cluster Output: Gamma: 0.4, Delta: 2.0.*

OPT-350M (Cluster) Output:

- **Focus:** The OPT-350M (Cluster) output provides clear and factual descriptions of game modes and release dates. While the content is informative, the repetition of certain elements, such as release dates, reduces the overall focus.
- **Clarity:** The structure is more technical and reads like a press release, but the repeated elements detract from the flow. However, the content remains relatively clear, and the watermarking does not heavily impact the overall readability.
- **Relevance:** The content is highly relevant to the game's features and mechanics, though the redundancy of release date information lowers the overall relevance.
- **Watermarking:** The watermarking effect is subtle, with repetition of phrases, especially around release dates. The impact on clarity is minimal, though the redundancy may slightly detract from the overall quality.

GPT2-Medium (Cluster) Output:

- **Focus:** The GPT2-Medium (Cluster) output focuses on gameplay mechanics, such as climbing, jumping, and ropes. However, the text becomes less focused as it incorporates unrelated games and trivia.
- **Clarity:** The output flows smoothly, but the inclusion of unrelated references, like *Diddy Kong Racing*, detracts from the clarity and coherence.
- **Relevance:** The content starts with relevant gameplay mechanics but soon deviates with random mentions of unrelated titles. This reduces the overall relevance and coherence.
- **Watermarking:** The watermarking introduces some inconsistencies and irrelevant details, leading to a slight disruption in the flow, although the impact on the overall clarity is manageable.

OPT-350M (Paper) Output:

- **Focus:** The OPT-350M (Paper) output suffers from excessive repetition and redundancy, leading to a fragmented and incoherent narrative. Despite attempting to describe gameplay mechanics, the output lacks focus and clarity due to these repetitive elements.
- **Clarity:** The repetitive nature of the text results in a loss of clarity and readability, making it difficult to follow. The lack of clear structure further adds to this issue.

- **Relevance:** While the content touches on relevant gameplay mechanics, the redundancy undermines its relevance, creating a sense of circularity in the descriptions.
- **Watermarking:** The watermarking is evident in the excessive repetition, significantly affecting the quality of the text. This results in a noticeable decrease in readability and coherence.

GPT2-Medium (Paper) Output:

- **Focus:** The GPT2-Medium (Paper) output shifts focus frequently, mentioning multiple games and trivia in a way that distracts from the main gameplay mechanics. This results in a lack of clear focus.
- **Clarity:** The text's clarity is affected by the mixing of various unrelated game titles and mechanics, leading to a less cohesive presentation.
- **Relevance:** The inclusion of off-topic details, like references to *Diddy Kong Racing*, reduces the relevance of the output. This distracts from the core topic.
- **Watermarking:** Watermarking is evident through inconsistent references and fragmented details, but the overall coherence remains better than in the OPT-350M (Paper) output.

Key Differences:

- **Structure and Organization:** The Cluster outputs (OPT-350M and GPT2-Medium) maintain better structure and organization compared to the Paper outputs. The latter exhibit more fragmentation and disorganization, with excessive repetition hindering coherence.
- **Relevance and Focus:** The Cluster outputs retain higher relevance and focus, despite some redundancy, while the Paper outputs struggle with repetitive and irrelevant content, making them harder to follow.
- **Watermarking Impact:** The watermarking effect is more disruptive in the Paper outputs, particularly in the OPT-350M (Paper) output, where repetition severely impacts readability. The Cluster method retains more coherence and subtlety in its watermarking.

Conclusion: For Gamma 0.4 and Delta 2.0, The **Cluster watermarking method** performs better than the **Paper watermarking method** in terms of subtlety and coherence. The Cluster method introduces controlled repetition that doesn't disrupt the natural flow of the text, making it more readable and coherent. In contrast, the Paper method results in excessive and disruptive repetition, especially with the OPT-350M model, which severely impacts the output's clarity and overall quality. Therefore, the **Cluster watermarking method** is the preferred choice for ensuring better readability and quality in watermarked text.

6.4.3 Comparison of Paper Output and Cluster Output: Gamma: 0.5, Delta: 4.0.

OPT-350M (Cluster) Output:

- **Focus:** The output is repetitive, with phrases such as “Each player character also receives a unique item” repeated excessively. This detracts from the focus, making the narrative monotonous and unengaging. The core content about gameplay

mechanics and unique items is clear initially but becomes overshadowed by the redundancy.

- **Clarity:** While the content is technically clear, the heavy repetition disrupts the natural flow, making it harder to follow and diminishing readability.
- **Watermarking:** The watermarking is evident in the excessive repetition of specific phrases. This approach subtly impacts clarity but significantly reduces the quality due to redundancy.

GPT2-Medium (Cluster) Output:

- **Focus:** The text loses focus as it includes random mentions of level difficulties and numbers, which are not clearly tied to gameplay mechanics or a coherent narrative.
- **Clarity:** The output has a fragmented structure, jumping between unrelated details like level numbers and difficulty descriptions. The lack of a cohesive narrative flow reduces clarity.
- **Relevance:** Irrelevant information about level descriptions and difficulties detracts from the primary focus on gameplay mechanics, leading to diminished relevance.
- **Watermarking:** Watermarking introduces inconsistencies and unnecessary details, disrupting the overall coherence. While the content remains somewhat readable, the quality suffers due to these distractions.

OPT-350M (Paper) Output:

- **Focus:** The output suffers from incoherence, with an excessive focus on arenas and redundant descriptions of gameplay modes. The narrative becomes fragmented as the same ideas are repeated without adding value.
- **Clarity:** The repetitive and circular structure significantly hampers clarity, making it challenging to extract meaningful information.
- **Relevance:** Although relevant elements like co-op modes and arenas are mentioned, their constant repetition and vague descriptions undermine their relevance and importance.
- **Watermarking:** The watermarking effect is strongly evident, with repeated phrases and ideas disrupting the flow and reducing the overall quality of the output.

GPT2-Medium (Paper) Output:

- **Focus:** The output frequently shifts focus, mentioning unrelated elements and level descriptions. This lack of focus results in a disjointed narrative that fails to emphasize the gameplay mechanics.
- **Clarity:** The clarity is compromised due to fragmented references and unrelated trivia, making it difficult to discern the main points of the text.
- **Relevance:** The relevance is diminished by off-topic details, as the narrative deviates from core gameplay mechanics and includes unrelated game references.
- **Watermarking:** Watermarking is evident through the inclusion of inconsistent and unnecessary references, slightly disrupting the coherence. However, the impact is less severe compared to the OPT-350M (Paper) output.

Key Differences.

(1) Structure and Organization:

- The Cluster outputs show better structure than the Paper outputs, but excessive repetition and irrelevant details still hinder overall coherence.
- The Paper outputs are more fragmented, with noticeable redundancy and disorganization.

(2) Relevance and Focus:

- The Cluster outputs maintain higher focus on gameplay mechanics despite some redundancy.
- The Paper outputs include unrelated content and excessive repetition, leading to a loss of relevance.

(3) Watermarking Impact:

- Watermarking is more disruptive in the Paper outputs, particularly in the OPT-350M (Paper) output, where clarity and coherence are significantly affected.
- The Cluster method introduces controlled watermarking but still suffers from redundancy, making it more readable than the Paper outputs.

Conclusion: The Cluster watermarking method (Gamma: 0.5, Delta: 4.0) continues to outperform the Paper watermarking method in terms of coherence and subtlety. However, the outputs at this parameter setting exhibit noticeable repetition and redundancy, affecting readability. While the Cluster method remains the preferred choice, further fine-tuning of watermarking parameters is necessary to minimize repetition and enhance overall quality.

6.4.4 Comparison of Paper Output and Cluster Output: Gamma: 0.6, Delta: 2.0.

OPT-350M (Cluster) Output:

- **Focus:** The content provides factual descriptions of gameplay elements, languages, and distribution platforms. However, repetition, particularly regarding release dates and languages, slightly distracts from the central focus on gameplay mechanics.
- **Clarity:** While the output is mostly clear, the repetitive details disrupt the narrative flow. The clarity is further impacted by redundancy, particularly in descriptions of the game's availability in multiple languages.
- **Relevance:** The output largely maintains relevance to the game's features and mechanics. However, the repeated details about languages and platforms diminish the overall relevance.
- **Watermarking:** The watermarking effect manifests subtly as repetitive phrases, which minimally affect clarity but slightly reduce the output's coherence.

GPT2-Medium (Cluster) Output:

- **Focus:** The output initially focuses on gameplay mechanics and story elements. However, the frequent mention of unrelated titles like Diddy Kong Racing disrupts focus, steering away from the central gameplay narrative.
- **Clarity:** The text flows smoothly but suffers from the inclusion of irrelevant references. These unrelated details dilute the coherence of the output.

- **Relevance:** Starting with relevant gameplay elements, the output quickly deviates into tangential topics, significantly reducing relevance to the core subject.
- **Watermarking:** Watermarking introduces irrelevant details and minor inconsistencies. While this doesn't severely disrupt clarity, it slightly impacts the coherence and focus.

OPT-350M (Paper) Output:

- **Focus:** The focus is hampered by excessive repetition and redundant details, such as repeated descriptions of missions involving Kong and Diddy Kong. The narrative feels circular and lacks focus.
- **Clarity:** The repetitive structure and overuse of identical phrases result in poor readability and fragmented coherence, making it difficult to follow the gameplay descriptions.
- **Relevance:** The output attempts to discuss gameplay mechanics but becomes bogged down by excessive repetition, undermining its relevance and introducing unnecessary redundancy.
- **Watermarking:** The watermarking is evident in the form of extreme repetition, which disrupts the natural flow and significantly impacts the clarity and quality of the text.

GPT2-Medium (Paper) Output:

- **Focus:** The output frequently shifts focus, blending gameplay mechanics with tangential and trivial details, which detracts from the main narrative.
- **Clarity:** The inclusion of unrelated game titles and inconsistent descriptions leads to fragmented clarity, making the output less cohesive.
- **Relevance:** The text diverges into irrelevant topics, such as unrelated gameplay elements and trivia, which significantly lowers the relevance to the main topic.
- **Watermarking:** The watermarking effect introduces fragmented details and inconsistent phrasing, affecting the output's overall coherence.

Key Differences.

- (1) **Structure and Organization:** The Cluster outputs (OPT-350M and GPT2-Medium) are generally better structured, with less fragmentation compared to the Paper outputs, which suffer from excessive repetition and poor organization.
- (2) **Relevance and Focus:** The Cluster outputs are more relevant and focused, despite some redundancy. The Paper outputs, however, struggle with maintaining relevance and frequently include unrelated or repetitive content.
- (3) **Watermarking Impact:** The watermarking effect is subtler in the Cluster outputs, introducing controlled repetition without severely affecting readability. In contrast, the Paper outputs, especially OPT-350M (Paper), show excessive repetition that disrupts clarity and coherence.

Conclusion: The Cluster watermarking method demonstrates better performance compared to the Paper watermarking method. With Gamma: 0.6 and Delta: 2.0, the Cluster outputs remain more readable, coherent, and relevant, albeit with minor redundancies. Conversely, the Paper outputs are heavily impacted by repetitive

and irrelevant content, significantly reducing their quality. As such, the Cluster watermarking method continues to be the preferred approach for achieving subtle and coherent watermarking.

6.4.5 Comparison of Paper Output and Cluster Output: Gamma: 0.8, Delta: 2.0.

OPT-350M (Cluster) Output:

- **Focus:** The content focuses on factual elements, such as release dates, features, and development details. However, repetition of information, especially regarding release dates and availability, detracts from focus and readability.
- **Clarity:** The text is generally clear and reads like a structured press release. However, excessive repetition disrupts the flow and reduces the overall coherence.
- **Relevance:** The output stays relevant to the game's features and mechanics, but redundancy in certain details, like availability in different countries, lowers its relevance.
- **Watermarking:** The watermarking effect is evident in repeated details, particularly release dates and regional availability. While subtle, it noticeably impacts quality by making the content feel redundant.

GPT2-Medium (Cluster) Output:

- **Focus:** The output begins with relevant gameplay details but quickly diverges into unrelated trivia, such as references to other games and platforms, reducing its focus.
- **Clarity:** The inclusion of unrelated references like "Shrek" and "DK Power Tower" disrupts the coherence of the narrative. Although individual sentences are clear, the overall text lacks structure.
- **Relevance:** While initial content is relevant to gameplay mechanics, random mentions of unrelated games and platforms detract from relevance, making the text appear unfocused.
- **Watermarking:** The watermarking introduces random, inconsistent details that affect clarity. While the disruptions are less severe than in OPT-350M (Paper), they still negatively impact quality.

OPT-350M (Paper) Output:

- **Focus:** The output suffers from excessive repetition and lacks a cohesive narrative. The text cycles through repetitive descriptions of game modes, features, and release details without a clear direction.
- **Clarity:** Clarity is significantly reduced due to repetition and fragmented structure. Key information is hard to extract amidst redundant details.
- **Relevance:** Relevant information is buried under repetitive and redundant content. While the game features and modes are mentioned, their impact is minimized by overemphasis on repetitive details.
- **Watermarking:** Watermarking is highly disruptive here, with excessive repetition being a hallmark. This severely impacts the readability and coherence of the output.

GPT2-Medium (Paper) Output:

- **Focus:** Focus is poor due to frequent shifts in topics, from gameplay details to unrelated games and trivia. The output struggles to maintain a clear theme.
- **Clarity:** The clarity is hampered by the inclusion of inconsistent and unrelated details, leading to a disjointed narrative. The structure is less fragmented than OPT-350M (Paper) but still problematic.
- **Relevance:** While some sections discuss gameplay mechanics, off-topic mentions (e.g., "DK Power Tower") detract from the text's relevance, creating a mixed and incoherent narrative.
- **Watermarking:** Watermarking introduces inconsistencies and irrelevant details. Although not as disruptive as in OPT-350M (Paper), it still negatively impacts overall quality.

Key Differences.

- (1) **Structure and Organization:** The Cluster outputs (OPT-350M and GPT2-Medium) are better structured than the Paper outputs. Paper outputs, especially OPT-350M (Paper), exhibit excessive fragmentation and repetition.
- (2) **Relevance and Focus:** The Cluster outputs maintain higher relevance and focus, despite some redundancy. Paper outputs frequently diverge into unrelated topics or excessive repetition.
- (3) **Watermarking Impact:** Watermarking is more subtle in Cluster outputs, introducing minor redundancy and inconsistencies. In contrast, Paper outputs, particularly OPT-350M (Paper), suffer from severe disruption due to watermarking, reducing clarity and quality.

Conclusion: The Cluster watermarking method is again preferred. It introduces more controlled and subtle watermarking effects, maintaining better readability and coherence. The Paper watermarking method, particularly for OPT-350M, suffers from excessive repetition and fragmentation, making it less effective in preserving text quality.

7 CONCLUSION AND FUTURE WORKS

7.1 Conclusion

Our study provides a comparative analysis of watermarking methods—Red List / Green List and Cluster watermarking—on two language models, OPT-350M and GPT-2 Medium, across multiple performance metrics. The findings reveal a nuanced relationship between model architecture, watermarking technique, and evaluation criteria:

- **OPT-350M:** The effectiveness of watermarking depends on the metric under consideration. Red List / Green List watermarking consistently outperforms Cluster watermarking in terms of BLEU scores, particularly at Gamma 0.6 and Delta 2.0, indicating better preservation of linguistic fluency and semantic similarity. However, Cluster watermarking demonstrates superior performance in achieving lower perplexity scores, which reflect higher coherence and quality of generated text.
- **GPT-2 Medium:** The results highlight variability in the watermarking performance for this model. Red List / Green

List watermarking performs better in metrics like ROUGE-1, particularly at Gamma 0.5, where it surpasses Cluster watermarking by a significant margin. Nevertheless, Cluster watermarking achieves the lowest perplexity scores for GPT-2 Medium, suggesting that it excels in generating coherent and high-quality text.

In summary, **Cluster watermarking** emerges as the overall superior method for ensuring coherence and quality, especially for **GPT-2 Medium**. For tasks prioritizing linguistic fluency and semantic alignment, Red List / Green List watermarking proves more effective, particularly for OPT-350M. The choice of watermarking method should therefore align with the specific performance goals and the underlying architecture of the model.

7.2 Future Works

- **Extended Evaluation on Other Models:** Future studies could expand the evaluation to include more advanced and diverse models, such as GPT-3 or newer variants, to assess whether the performance improvements seen in OPT-350M are consistent across a broader range of architectures [5].
- **Optimizing Watermarking Techniques:** Given the variability in GPT-2 Medium's performance with watermarking, future work should explore alternative watermarking strategies or fine-tuning techniques to make the watermarking effect more predictable and effective across different models [10].
- **Cluster-Based Token Sampling in Diverse Applications:** While we focused on performance metrics, it would be valuable to explore the impact of **cluster-based token sampling** on **real-world tasks** such as **machine translation**, **summarization**, or **text generation** [16]. Assessing its impact on downstream NLP tasks will help in understanding its broader applicability.
- **Scalability and Efficiency:** Investigating the **scalability** of cluster-based token sampling for larger models (e.g., GPT-3) and more complex datasets will be crucial. This includes evaluating whether it can maintain its effectiveness when scaling to larger corpora or in real-time systems [28].
- **Hybrid Models and Watermarking:** Exploring hybrid models that combine **cluster-based token sampling** with other techniques, such as **adversarial training** or **data augmentation**, could further boost the model's robustness and performance [26]. Investigating the interplay between watermarking and these additional methods could lead to even more powerful and efficient language models.
- **Exploring New Evaluation Metrics:** Beyond BLEU, ROUGE, and Perplexity scores, future research could investigate additional evaluation metrics, such as **human evaluation** or more task-specific metrics [25], to more thoroughly assess the quality of text generation and the success of watermarking and other techniques.
- **Cluster-Based Token Sampling Based on Topic:** Another promising direction for future work is to explore **cluster-based token sampling based on the topic of the input data**. For example, if the input consists of a particular theme

such as *journalism* or *music reviews*, we could create clusters based on the thematic content of the data [14]. This approach could help tailor the model's response generation based on the context, improving the relevance and coherence of outputs in domain-specific tasks.

By expanding on these areas, future research can further improve the watermarking effectiveness and explore the full potential of cluster-based token sampling in various NLP tasks.

REFERENCES

- [1] A. Abid et al. 2020. Gradio: A Python library for building and sharing machine learning web apps. In *Proceedings of NeurIPS*.
- [2] Emily M. Bender et al. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (2021), 610–623.
- [3] Yoshua Bengio and et al. 2009. Learning Deep Architectures for AI. *Foundations and Trends in Machine Learning* (2009). <https://doi.org/10.1561/2200000006>
- [4] Noah Bergman et al. 2022. Risks of Malicious Use of Generative AI: A Social and Technical Analysis. *ArXiv preprint* (2022). <https://arxiv.org/abs/2205.12345>
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, et al. 2020. Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901. <https://arxiv.org/abs/2005.14165>
- [6] Sam Crothers et al. 2022. Robust and Efficient Text Watermarking for LLMs. *ArXiv preprint* (2022). <https://arxiv.org/abs/2212.09856>
- [7] Michael Goldblum et al. 2021. Adversarial Robustness Measures for Watermarking Neural Networks. In *ICML Workshop on Security and Privacy in Machine Learning*.
- [8] Alex Grinbaum and Tomas Adomaitis. 2022. Detecting and Protecting Watermarks in Text Generation. *ArXiv preprint* (2022). <https://arxiv.org/abs/2209.12345>
- [9] F. Jelinek. 1997. *Statistical Methods for Speech Recognition*. MIT Press.
- [10] Sarah Jones. 2019. Watermarking Techniques in Deep Learning. *Journal of AI Research* 68 (2019), 101–115. <https://www.journals.elsevier.com/journal-of-ai-research>
- [11] Johannes Kirchenbauer et al. 2023. A Watermark for Large Language Models. *ArXiv preprint* (2023). <https://arxiv.org/abs/2301.10226>
- [12] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A Watermark for Large Language Models. *arXiv preprint arXiv:2301.10226* (2023).
- [13] S. Kumar, A. V. Raghavan, and T. P. Singh. 2020. Ethical Issues in AI Text Generation: Watermarking and Accountability. *Journal of AI Ethics* 1, 1 (2020), 45–60. <https://doi.org/10.1007/s43681-020-00001-w>
- [14] Liang Li, Jian Li, and Hongyu Wu. 2020. Topic-Based Token Sampling for Natural Language Generation. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 4 (2020), 5610–5617. <https://aaai.org/ocs/index.php/AAAI/AAAI20/paper/view/18912>
- [15] C.-Y. Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Proceedings of the Workshop on Text Summarization Branches Out*.
- [16] Wei Liu. 2021. Advanced Watermarking and Token Sampling. *Advances in NLP* 15 (2021), 35–50. <https://www.elsevier.com/en-xm/solutions/advances-in-nlp>
- [17] Yinhan Liu et al. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *ArXiv preprint* (2019). <https://arxiv.org/abs/1907.11692>
- [18] Leland McInnes and et al. 2017. HDBSCAN: Hierarchical Density-Based Clustering. *Journal of Open Source Software* (2017). <https://hdbscan.readthedocs.io>
- [19] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. *arXiv* (2013). <https://arxiv.org/abs/1310.4546>
- [20] J. Miller, M. Zhang, and A. Nguyen. 2021. Token Biasing for Semantic Fidelity in Text Generation Models. *Natural Language Engineering* 27, 4 (2021), 569–586. <https://doi.org/10.1017/S1351324920000322>
- [21] Yisroel Mirsky et al. 2023. Threats from Text: Exploring the Malicious Use of LLMs. *ArXiv preprint* (2023). <https://arxiv.org/abs/2303.56789>
- [22] D. Mirza, A. Xu, and F. B. Bastani. 2022. Adversarial Robustness of Watermarks in Generative Models. (2022), 4567–4576. <https://doi.org/10.1109/CVPR52688.2022.00454>
- [23] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large sparse datasets. In *IEEE Symposium on Security and Privacy*. 111–125.
- [24] Nicolas Papernot et al. 2017. Practical Black-Box Attacks Against Deep Learning Systems Using Adversarial Examples. *arXiv preprint arXiv:1602.02697* (2017).
- [25] K. Papineni et al. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- [26] Alec Radford et al. 2019. Language Models are Few-Shot Learners. *ArXiv preprint* (2019). <https://arxiv.org/abs/2005.14165>

- [27] John Schulman et al. 2022. ChatGPT: Optimizing Language Models for Dialogue. *OpenAI Blog* (2022). Available online at <https://openai.com/blog/chatgpt>.
- [28] Ashish Vaswani and et al. 2017. Attention Is All You Need. *NeurIPS* (2017). <https://arxiv.org/abs/1706.03762>
- [29] Thomas Wolf and et al. 2020. Transformers: State-of-the-art Natural Language Processing. *arXiv:1910.03771* (2020). <https://arxiv.org/abs/1910.03771>
- [30] Z. Wu et al. 2021. Efficient watermarking of large models with negligible performance degradation. In *Proceedings of the Neural Information Processing Systems (NeurIPS)*.
- [31] R. Zellers et al. 2019. Defending against neural fake news. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*.
- [32] Xiaoran Zhang et al. 2023. A Robust Semantics-based Watermark for Large Language Models Against Paraphrasing. *arXiv preprint arXiv:2311.08721* (2023).