

From Pixels to Tokens: Otsu Thresholding for Accurate yet Confident Decoding

Anonymous ACL submission

Abstract

This paper introduces a new decoding approach, Otsu thresholding, inspired by the well-known image processing technique, which computes a threshold between low-confidence and high-confidence tokens for each generation step. The method is tested across several tasks and instruction-based models, revealing results competitive to other state-of-the-art decoding methods, while keeping output quality. Also, it is shown that tuning specific characteristics of the model, such as the pass from a softmax function of the distribution and the limitation to K highest candidates before selecting a token, can highly affect model performance. In general, this approach is a great alternative method that offers a good balance between quality, confidence, uncertainty, and diversity.

1 Introduction

In recent years, Large Language Models (LLMs) have shown great potential in generative tasks, such as text summarization and question answering (Jurafsky and Martin, 2025). According to the latest leaderboard¹, they are able to extract high accuracy output and therefore, create high quality content. However, a common problem of LLMs is the overestimation of their confidence ("they don't always know what they don't know" (Baan et al., 2023; Hu et al., 2023)), leading to outputs with hallucinations or misleading information. Ideally, a high-performance LLM should be able to generate text that is both accurate and confident, while providing a principled estimate of uncertainty, when knowledge is insufficient.

Decoding strategies can play a key role towards the aforementioned goal, and several strategies have been proposed that account for token probabilities to control the token prediction set. The choice of the decoding method for an LLM is very important, as it impacts its performance (Shi et al., 2024).

¹<https://huggingface.co/open-llm-leaderboard>

However, most decoding methods fail to generate accurate and confident text, and they frequently result in non-representative prediction sets, with too few or, more often, too many low-probability tokens being considered. Such low-probability candidates correspond to regions of higher epistemic uncertainty in the model's predictive distribution. As a result, existing approaches overlook the need for fast, real-time decoding and for uncertainty estimation, leaving a critical gap in how uncertainty is represented in text generation (Shao et al., 2018; Hewitt et al., 2022; Holtzman et al., 2019). A solution to this problem could be the proposal of a new technique that efficiently distinguishes low-probability tokens from high-probability ones in a data-driven and uncertainty-aware manner.

The motivation of this paper is to propose an innovative and efficient decoding strategy that separates low-probability tokens from high-probability ones and mitigates uncertainty by suppressing low-confidence token candidates, and maintains strong generation quality. The strategy uses a very popular filtering method from image processing named Otsu thresholding. To the best of our knowledge, this is the first work to apply Otsu thresholding to probabilistic token filtering in LLM decoding. Thus, we check the impact of some features in Otsu, such as the usage of a softmax function on the logits of the tokens and the restriction of the number of candidates in the list. Finally, we estimate LLM performance with this approach and we compare the best combination to other sampling methods, finding competitive values with adaptive middle ground for certain and calibrated outputs.

To summarize, the main contributions of the paper are the following:

- We introduce Otsu thresholding as a novel uncertainty-aware decoding strategy that adaptively separates low-confidence from high-confidence token candidates, based on the

081 model’s predictive distribution. 130

- 082 • We analyze the proposed decoding strategy 131
- 083 over different tasks and models. 132
- 084 • We experiment with the usage of a softmax 133
- 085 function and the limitation to K candidates 134
- 086 in the approach to understand their impact on 135
- 087 the model output and we compare the best 136
- 088 approach to other state-of-the-art decoding 137
- 089 strategies to demonstrate favorable trade-offs 138
- 090 in accuracy, confidence, uncertainty, and di- 139
- 091 versity. 140

092 2 Related Work 141

093 Outstanding studies have focused on the explo- 142

094 ration of the impact of **decoding strategies on** 143

095 **text uncertainty** in Large Language Models 144

096 (Hashimoto et al., 2025). (Daheim et al., 2025) 145

097 showed how Minimum Bayes Risk (MBR) decod- 146

098 ing, which selects model generations according to 147

099 an expected risk, can be generalized into a prin- 148

100 ciple uncertainty-aware decoding method. Addi- 149

101 tionally, (He et al., 2025) proposed an uncertainty- 150

102 guided adaptive decoding framework that inte- 151

103 grates a token-level pause-then-rerank mechanism, 152

104 driven by token uncertainty (Shannon entropy). 153

105 AdaDec learns model-specific uncertainty thresh- 154

106 olds and applies a lookahead-based reranking 155

107 strategy, when uncertainty is high. Finally, the 156

108 very recent work of (Lee et al., 2025) suggested 157

109 Uncertainty-Aware Contrastive Decoding (UCD), 158

110 which introduces a cumulative energy function, 159

111 where uncertainty is quantified as the negative log- 160

112 sum-exp over logits, and decomposed into entropy 161

113 and expected logit components. 162

114 **Evaluation of factuality** in text generation is an 163

115 ongoing challenge. (Zha et al., 2023) introduced 164

116 AlignScore, a holistic metric, based on a general 165

117 function of information alignment of text and its 166

118 unified framework, which achieved substantial im- 167

119 provements over previous metrics. Additionally, 168

120 (Min et al., 2023) advocated a new evaluation met- 169

121 ric that computes factual accuracy from pieces of 170

122 generated text and was used to compare the perfor- 171

123 mance of different LLMs. Finally, (Bishop et al., 172

124 2023) proposed a new evaluation framework, Long- 173

125 DocFACTScore, for detecting human factuality, tar- 174

126 geting specifically summarized, long documents. 175

127 Previous work has comprehensively examined 176

128 **uncertainty in Natural Language Generation** 177

129 (NLG) systems (Baan et al., 2023; Hu et al., 2023) 178

130 and has explored strategies to address uncertainty 131

132 with the goal of making LLMs more trustworthy. 133

134 Firstly, (Xu et al., 2020) studied summarization de- 135

136 coders in both blackbox and whitebox ways by 136

137 focusing on the entropy of the models’ predic- 137

138 tions and revealed that features, such as the sen- 138

139 tence position and the syntactic distance between 139

140 adjacent pairs of tokens, influence uncertainty. 140

141 Moreover, (Ulmer et al., 2024) focused on token- 141

142 level uncertainty and proposed a method for non- 142

143 exchangeable conformal prediction, which was 143

144 shown to improve text generation quality. Finally, 144

145 (Fadeeva et al., 2024) introduced a token-level un- 145

146 certainty method named Claim Conditioned Prob- 146

147 ability (CCP), disentangling claim-specific uncer- 147

148 tainty from model decisions on surface forms, etc. 148

149 **Diversity of LLMs** is an ongoing challenge in 149

150 text generation, as many tokens that could be cor- 150

151 rectly selected are ignored from the sampling meth- 151

152 ods. Studies attempt to estimate it, such as (Vi- 152

153 jayakumar et al., 2016), which used diversity statis- 153

154 tics by computing the number of distinct n-grams 154

155 in a summary divided by the total number of tokens 155

156 in it. 156

157 3 Methodology 157

158 In this work, we use Otsu thresholding as an alterna- 158

159 tive decoding strategy. Otsu thresholding is based 159

160 on an automated technique from image processing 160

161 (Otsu et al., 1975), where the algorithm returns a 161

162 single intensity threshold t that separates pixels 162

163 into two classes, background (C_0) and foreground 163

164 (C_1). The goal is to find an optimal threshold that 164

165 maximizes inter-class variance: 165

$$166 \sigma_b^2(\operatorname{argmax}) = \max_{0 < t < L} \sigma_b^2(t), \quad (1) \quad 167$$

168 where L is the number of bins in the image and 168

169 $\sigma_b^2(t)$ is the inter-class variance. 169

170 Let ω_0 and ω_1 be the cumulative probabilities of 170

171 C_0 and C_1 respectively: 171

$$172 \omega_0(t) = \sum_{i=0}^{t-1} P(i), \quad (2) \quad 173$$

$$174 \omega_1(t) = \sum_{i=t}^{L-1} P(i), \quad (3) \quad 175$$

176 where $P(i|C_0)$ and $P(i|C_1)$ are the conditional 176

177 probabilities of selecting the i -th pixel in C_0 and 177

178 C_1 . 178

Also, let $\mu_0(t)$ and $\mu_1(t)$ be the mean pixel intensities of C_0 and C_1 respectively:

$$\mu_0(t) = \frac{\sum_{i=0}^{t-1} iP(i)}{\omega_0(t)}, \quad (4)$$

$$\mu_1(t) = \frac{\sum_{i=t}^{L-1} iP(i)}{\omega_1(t)}. \quad (5)$$

The algorithm performs global search, where at each pass of an intensity, it updates these variables and re-computes inter-class variance.

In summary, the steps of the algorithm are the following:

1. Compute histogram and probabilities of each intensity level.
2. Initialize $\omega_0(t)$, $\mu_0(t)$, $\omega_1(t)$, and $\mu_1(t)$.
3. For each pixel intensity, update $\omega_0(t)$, $\mu_0(t)$, $\omega_1(t)$, and $\mu_1(t)$ and then, compute $\sigma_b^2(t)$.

Consequently, in our approach, the histogram is a distribution of token logits at each decoding step. Otsu computes a threshold from the distribution and distinguishes tokens with low values (C_0) from those with high values (C_1). After that, the model samples the tokens in C_1 and uniformly chooses one as the next token. From an uncertainty perspective, Otsu thresholding identifies a data-driven boundary between these regions, enabling filtering of low-confidence token candidates that are associated with hallucinated or overconfident generations. A simple comparison between the 2 types of Otsu thresholding techniques, along with an example from the dataset, can be found in Figure 1.

We experiment with different variations of the proposed decoding strategy in an attempt to understand their impact on the results and find the best combination. Specifically, we experiment with the pass of the logit distribution from a softmax function before applying Otsu to convert logits to probabilities, and we limit the candidate set to K tokens with the highest values at each decoding step to observe whether the selected option mostly originates from the first items.

4 Experimental Setup

4.1 Datasets

The dataset that we use for our experiments is a pre-processed version² of Super-Natural Instructions

²<https://huggingface.co/datasets/Muennighoff/natural-instructions>

(Wang et al., 2022), consisting of definitions, inputs, and targets for many tasks. To test the robustness and generalization of our method, we use samples from the following 3 different tasks:

- Amazon Review: Given an Amazon product review, the goal is to generate a title from it.
- Extreme Abstract: Given the abstract of a research paper, the goal is to summarize it in no more than 30 words.
- Web Questions: Given a web question, the goal is to answer it.

The above tasks don't contain the same amount of records (Amazon Review has 13,000 records, Extreme Abstract has 10,796 records, and Web Questions has 17,850 records), so we keep the first 600 records from each task for a fair comparison between them.

The prompt that we use for these tasks is:

definition input response,

where *definition* and *input* are retrieved as features from the dataset and *response* is "Output:", "Summary:", and "Answer:" for each task respectively.

The definitions from the dataset that the above tasks have are the following:

- Amazon Review: In this task, you're given reviews from Amazon's products. Your task is to generate the Summary of the review.
- Extreme Abstract: In this task, you are given the abstract of a research paper. Your task is to generate a summary of this abstract. Your summary should not be very short, but it's better if it's not more than 30 words.
- Web Questions: A question is presented to you in this task, and your job is to write a potentially correct answer.

4.2 Large Language Models

For the scope of our research, we use two open-access, state-of-the-art, instruction-tuned LLMs, selected to cover different model architectures and design choices.

LLaMA-8B-v3.1 is an auto-regressive transformer model, trained using a combination of Supervised Fine-Tuning (SFT) and Reinforcement Learning with Human Feedback (RLHF) to align

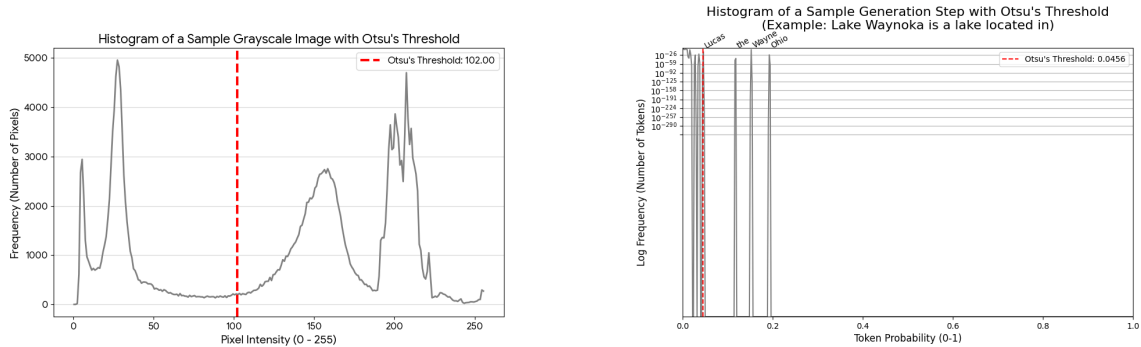


Figure 1: Example of Otsu’s method applied to an image and a token generation step.

its outputs with human preferences for helpfulness and safety (Touvron et al., 2023).

Qwen-v2.5 is a high-capacity LLM that supports context lengths of up to 128,000 tokens and a maximum generation length of 8,000 tokens. It is available in 7B and 14B parameter variants and incorporates several architectural enhancements that enable stable and high-quality generation over long contexts (Team et al., 2024).

4.3 Decoding strategies

We compare Otsu decoding against 2 decoding strategies to identify those producing outputs that are both accurate and well-calibrated in terms of quality, factuality, certainty, and diversity. Specifically, we compare:

- **Epsilon sampling:** It filters tokens whose probability is equal or greater than an epsilon cutoff, effectively removing low probability candidates from the sampling space (Hewitt et al., 2022).
- **Top-p sampling:** It chooses tokens from the smallest possible set whose cumulative probability exceeds a predefined threshold p . The probability mass is then redistributed among the selected tokens (Holtzman et al., 2019).

4.4 Hyperparameter settings

Our experiments are conducted on a Google Cloud VM instance, equipped with 208 vCPUs, 1,872 GB RAM, 100 GB memory disk, and 8 NVIDIA H100 GPUs, each with 80 GB of memory. For all the text generation strategies, we set the number of generated sequences to 5, epsilon cutoff to 9^e-4 , K to 50, and p to 0.75, wherever applicable. These values are selected to balance computational efficiency with reliable estimation of all the types of metrics.

4.5 Evaluation metrics

We evaluate LLM performance across the selected tasks using multiple types of metrics.

Quality metrics compare the generated prediction to a true reference. ROUGEL measures the longest common subsequence (LCS) between the 2 types of text (Lin, 2004). BERTScore leverages pre-trained contextual embeddings to match tokens in candidate and reference text using cosine similarity (Zhang et al., 2019) and provides an F1 score. As an additional step of **benchmarking** our work, we consider ROUGEL ideal for the review generation and abstract summarization task, while, for question answering, we use an F1 score that is computed over the individual words in the generated text against the true answer.

Factuality metrics assess the alignment between generated claims and ground-truth evidence. HHEM models extract a score from 0 to 1, indicating the level of hallucinations in text (Bao et al., 2024). AlignScore is another metric that produces values in the same range, using a unified information alignment function and splitting the 2 types of text to compute an average score from the maximum alignment scores of the pairs (Zha et al., 2023).

Certainty metrics quantify uncertainty at each generation step and are averaged per sample. Token certainty estimates the probability of the generated token, while token entropy computes the Shannon entropy from the probability distribution of the vocabulary (Shannon, 1948).

Diversity metrics capture lexical and semantic variation in the generated outputs. We estimate inter-generation diversity by comparing generated sequences for the same sample using BERTScore and keeping the remaining value from the subtraction to 1.

334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382

5 Results

5.1 Analysis of the features of the Otsu thresholding approach

As an initiative of understanding better the behavior of the proposed decoding strategy, we examine how specific design choices, such as the application of a softmax function to the logits and the limitation to K candidates, affect the performance of the model. We compare the combinations of these features across all the dataset categories, LLMs, and evaluation metrics, enabling analysis of their impact on quality, factuality, uncertainty, and diversity.

The results in Tables 1 and 2 reveal a clear trade-off between output diversity and the remaining evaluation metrics. Specifically, the use of a softmax function sharpens the output distribution, leading to more confident predictions, while consistently reducing diversity. In addition, restricting the analysis to the top 50 candidates yields performance that is comparable to, and, in some cases, slightly better than using all model predictions. This suggests that the most informative signal is concentrated within the top K candidates. Given that this restriction also offers improved computational efficiency, adopting a top 50 setting represents a practical and effective choice.

5.2 Comparison of Otsu thresholding to other decoding strategies

The best results from Otsu are then compared to other decoding strategies. We compare the different methods using all the dataset categories, LLMs, and evaluation metrics.

From Tables 3 and 4, it can be viewed that epsilon decoding generally increases quality and diversity metrics, but often at the cost of certainty. Top-p achieves the best values in uncertainty, but is not optimal in diversity. Across datasets and models, Otsu performs comparably to top-p and occasionally matches or exceeds it on key metrics without affecting diversity much, indicating that it is an effective, best-balanced, alternative decoding strategy.

5.3 Presentation of text generations

For a better understanding of the results, we present some examples of text generations from the experiments. In detail, we note the most certain summaries of the Extreme Abstract task that were generated from the different Otsu thresholding parameterizations and the other decoding strategies using

the Qwen-7B model.

Otsu Thresholding (no softmax, all candidates)

ID: task668-f7af2fc4ebf24a6888575e5e6d566a38
Summary: The CWAE model uses a new distance metric with a closed-form solution, simplifying optimization and matching or outperforming WAE and SWAE.

Otsu Thresholding (softmax, all candidates)

ID: task668-f77032d675f041fba3261703979e2353
Summary: Proposed Prox-SGD for training NNs with nonsmooth regularization and constraints, ensuring convergence to stationary points and performing well on sparse and binary NNs. To promote a NN with specific structures, we explicitly take into consideration the nonsmooth regularization (such as L1-norm) and constraints (such as interval constraint). This is formulated as a constrained nonsmooth nonconvex optimization problem, and we propose a convergent proximal-type stochastic gradient descent (Prox-SGD) algorithm.

Otsu Thresholding (no softmax, 50 candidates)

ID: task668-725120e56d404feca47a0d3a94c3ef6
Summary: Morph-Net uses morphological ops for neuron computations, offering better performance and fewer params than standard nets on various datasets.

Otsu Thresholding (softmax, 50 candidates)

ID: task668-87bc07c2e6b0448e86afa5d967e41af8
Summary: The authors propose a method to minimize the policy's Lipschitz constant during training to improve robustness and efficiency in reinforcement learning with domain randomization.

Epsilon Sampling

ID: task668-7ba32c9b1af24c91a2161def01bc5a9e
Summary: This research explores the divergence minimization perspective on GAN training, comparing original and "non-saturating" variants, and develops theoretical tools for classifying f-divergences.

Top-p Sampling

ID: task668-f7af2fc4ebf24a6888575e5e6d566a38
Summary: A new generative model, Cramer-Wold AutoEncoder (CWAE), is proposed which uses Cramer-Wold distance for improved performance compared to existing methods like WAE and SWAE. CWAE simplifies optimization without compromising on quality.

383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Dataset	Model	Softmax Usage	Number of K	ROUGEL \uparrow	BERT-F1 \uparrow	HHEM \uparrow	AlignScore \uparrow
Amazon Review	LLaMA-8B	No	All	0.0363	0.8022	0.1721	0.0628
Amazon Review	LLaMA-8B	Yes	All	0.051	0.8019	0.3544	0.1058
Amazon Review	LLaMA-8B	No	50	0.0515	0.8049	0.3497	0.0973
Amazon Review	LLaMA-8B	Yes	50	0.0493	0.8019	0.349	0.1092
Amazon Review	Qwen-7B	No	All	0.0518	0.8056	0.3614	0.0926
Amazon Review	Qwen-7B	Yes	All	0.0358	0.8025	0.1592	0.0575
Amazon Review	Qwen-7B	No	50	0.0361	0.8031	0.1723	0.0602
Amazon Review	Qwen-7B	Yes	50	0.0359	0.8029	0.1696	0.0557
Extreme Abstract	LLaMA-8B	No	All	0.0757	0.8243	0.3503	0.0451
Extreme Abstract	LLaMA-8B	Yes	All	0.0745	0.8232	0.3383	0.0461
Extreme Abstract	LLaMA-8B	No	50	0.0766	0.8229	0.3381	0.0445
Extreme Abstract	LLaMA-8B	Yes	50	0.0733	0.8233	0.3342	0.0454
Extreme Abstract	Qwen-7B	No	All	0.0647	0.8236	0.3406	0.054
Extreme Abstract	Qwen-7B	Yes	All	0.0643	0.8255	0.3542	0.055
Extreme Abstract	Qwen-7B	No	50	0.0653	0.8235	0.336	0.0556
Extreme Abstract	Qwen-7B	Yes	50	0.0651	0.824	0.3503	0.0505
Web Questions	LLaMA-8B	No	All	0.0116	0.7728	0.407	0.043
Web Questions	LLaMA-8B	Yes	All	0.0117	0.772	0.4086	0.0454
Web Questions	LLaMA-8B	No	50	0.0116	0.7736	0.3956	0.0437
Web Questions	LLaMA-8B	Yes	50	0.0119	0.7723	0.4036	0.0461
Web Questions	Qwen-7B	No	All	0.0124	0.7768	0.3711	0.0429
Web Questions	Qwen-7B	Yes	All	0.0123	0.7767	0.3852	0.0435
Web Questions	Qwen-7B	No	50	0.0119	0.7769	0.3815	0.0412
Web Questions	Qwen-7B	Yes	50	0.0119	0.7765	0.3792	0.04

Table 1: Analysis of the features of the Otsu thresholding approach for each dataset category using all the LLMs and the quality and factuality metrics.

Dataset	Model	Softmax Usage	Number of K	Token Certainty \uparrow	Token Entropy \downarrow	Diversity Level \uparrow	F1-Q/A \uparrow
Amazon Review	LLaMA-8B	No	All	0.7759	0.5013	0.4467	-
Amazon Review	LLaMA-8B	Yes	All	0.9274	0.1188	0.3469	-
Amazon Review	LLaMA-8B	No	50	0.8307	0.3434	0.3878	-
Amazon Review	LLaMA-8B	Yes	50	0.9247	0.12	0.3473	-
Amazon Review	Qwen-7B	No	All	0.8253	0.3668	0.3924	-
Amazon Review	Qwen-7B	Yes	All	0.8684	0.2189	0.4286	-
Amazon Review	Qwen-7B	No	50	0.7857	0.4496	0.4455	-
Amazon Review	Qwen-7B	Yes	50	0.8687	0.2171	0.4274	-
Extreme Abstract	LLaMA-8B	No	All	0.8403	0.3067	0.3555	-
Extreme Abstract	LLaMA-8B	Yes	All	0.9083	0.1472	0.3323	-
Extreme Abstract	LLaMA-8B	No	50	0.8371	0.2981	0.3547	-
Extreme Abstract	LLaMA-8B	Yes	50	0.9096	0.1464	0.3332	-
Extreme Abstract	Qwen-7B	No	All	0.784	0.4699	0.4112	-
Extreme Abstract	Qwen-7B	Yes	All	0.8745	0.2073	0.3862	-
Extreme Abstract	Qwen-7B	No	50	0.7889	0.4384	0.4045	-
Extreme Abstract	Qwen-7B	Yes	50	0.8758	0.2046	0.3895	-
Web Questions	LLaMA-8B	No	All	0.8342	0.3444	0.3835	0.0167
Web Questions	LLaMA-8B	Yes	All	0.9098	0.1433	0.3365	0.0171
Web Questions	LLaMA-8B	No	50	0.8342	0.329	0.3778	0.0169
Web Questions	LLaMA-8B	Yes	50	0.9107	0.1406	0.3398	0.0169
Web Questions	Qwen-7B	No	All	0.7733	0.4917	0.3666	0.0173
Web Questions	Qwen-7B	Yes	All	0.867	0.2196	0.3349	0.0176
Web Questions	Qwen-7B	No	50	0.775	0.4692	0.3608	0.0169
Web Questions	Qwen-7B	Yes	50	0.8676	0.2176	0.3341	0.0167

Table 2: Analysis of the features of the Otsu thresholding approach for each dataset category using all the LLMs and the certainty and diversity metrics.

6 Conclusions

We proposed an innovative, uncertainty-aware decoding strategy that efficiently distinguishes low-probability from high-probability tokens. Our approach accomplished results competitive to other state-of-the-art methods, while keeping quality and

diversity in text. Moreover, we showed that specific characteristics of the approach influence model output, but have a minor impact on diversity. In conclusion, the above research opens a new field of exploration in the area of probabilistic decoding using LLMs.

Dataset	Model	Method	ROUGEL \uparrow	BERT-F1 \uparrow	HHEM \uparrow	AlignScore \uparrow
Amazon Review	LLaMA-8B	Epsilon	0.0515	0.8065	0.3693	0.0942
Amazon Review	LLaMA-8B	Top-p	0.053	0.8032	0.3527	0.1096
Amazon Review	LLaMA-8B	Otsu (best)	0.0493	0.8019	0.349	0.1092
Amazon Review	Qwen-7B	Epsilon	0.0362	0.8027	0.1784	0.0607
Amazon Review	Qwen-7B	Top-p	0.0358	0.8034	0.1717	0.0601
Amazon Review	Qwen-7B	Otsu (best)	0.0359	0.8029	0.1696	0.0557
Extreme Abstract	LLaMA-8B	Epsilon	0.0753	0.8231	0.3422	0.0432
Extreme Abstract	LLaMA-8B	Top-p	0.0748	0.8226	0.3501	0.0442
Extreme Abstract	LLaMA-8B	Otsu (best)	0.0733	0.8233	0.3342	0.0454
Extreme Abstract	Qwen-7B	Epsilon	0.0642	0.8229	0.3363	0.0528
Extreme Abstract	Qwen-7B	Top-p	0.0643	0.825	0.342	0.0506
Extreme Abstract	Qwen-7B	Otsu (best)	0.0651	0.824	0.3503	0.0505
Web Questions	LLaMA-8B	Epsilon	0.0116	0.7733	0.4122	0.0431
Web Questions	LLaMA-8B	Top-p	0.0114	0.7703	0.3995	0.0459
Web Questions	LLaMA-8B	Otsu (best)	0.0119	0.7723	0.4036	0.0461
Web Questions	Qwen-7B	Epsilon	0.0123	0.7769	0.3681	0.0433
Web Questions	Qwen-7B	Top-p	0.0117	0.7762	0.3671	0.0406
Web Questions	Qwen-7B	Otsu (best)	0.0119	0.7765	0.3792	0.04

Table 3: Comparison of the best version of Otsu to the other decoding strategies for each dataset category using all the LLMs and the quality and factuality metrics.

Dataset	Model	Method	Token Certainty \uparrow	Token Entropy \downarrow	Diversity Level \uparrow	F1-Q/A \uparrow
Amazon Review	LLaMA-8B	Epsilon	0.8274	0.3468	0.3898	-
Amazon Review	LLaMA-8B	Top-p	0.9242	0.1197	0.3553	-
Amazon Review	LLaMA-8B	Otsu (best)	0.9247	0.12	0.3473	-
Amazon Review	Qwen-7B	Epsilon	0.7825	0.4683	0.445	-
Amazon Review	Qwen-7B	Top-p	0.8687	0.2186	0.427	-
Amazon Review	Qwen-7B	Otsu (best)	0.8687	0.2171	0.4274	-
Extreme Abstract	LLaMA-8B	Epsilon	0.8367	0.3019	0.3561	-
Extreme Abstract	LLaMA-8B	Top-p	0.9117	0.1319	0.3316	-
Extreme Abstract	LLaMA-8B	Otsu (best)	0.9096	0.1464	0.3332	-
Extreme Abstract	Qwen-7B	Epsilon	0.7877	0.4479	0.4059	-
Extreme Abstract	Qwen-7B	Top-p	0.8733	0.2065	0.3863	-
Extreme Abstract	Qwen-7B	Otsu (best)	0.8758	0.2046	0.3895	-
Web Questions	LLaMA-8B	Epsilon	0.8328	0.3326	0.3801	0.0171
Web Questions	LLaMA-8B	Top-p	0.9164	0.1286	0.3435	0.0162
Web Questions	LLaMA-8B	Otsu (best)	0.9107	0.1406	0.3398	0.0169
Web Questions	Qwen-7B	Epsilon	0.7756	0.4712	0.3632	0.0167
Web Questions	Qwen-7B	Top-p	0.8671	0.2155	0.3357	0.0162
Web Questions	Qwen-7B	Otsu (best)	0.8676	0.2176	0.3341	0.0167

Table 4: Comparison of the best version of Otsu to the other decoding strategies for each dataset category using all the LLMs and the certainty and diversity metrics.

444 Limitations

445 Although our work gives a great initiative for en-
446 hanced probabilistic decoding, there are still some
447 areas that could be explored, like testing on more
448 records, tasks, and models. In addition, more de-
449 scriptive, uncertainty metrics can be incorporated
450 into the experiments for further detection and elim-
451 ination. Finally, the lack of comparison of auto-
452 mated metrics to human evaluation is another lim-

453 itation, which could strength the paper claims, if
454 conducted.

References

455 Joris Baan, Nico Daheim, Evgenia Ilia, Dennis Ul-
456 mer, Haau-Sing Li, Raquel Fernández, Barbara
457 Plank, Rico Sennrich, Chrysoula Zerva, and Wilker
458 Aziz. 2023. Uncertainty in natural language gener-
459 ation: From theory to applications. *arXiv preprint*
460 *arXiv:2307.15703*.
461

462	Forrest Bao, Miaoran Li, Rogger Luo, and Ofer	Factscore: Fine-grained atomic evaluation of factual	516
463	Mendelevitch. 2024. HHEM-2.1-Open .	precision in long form text generation. <i>arXiv preprint</i>	517
464	Jennifer A Bishop, Qianqian Xie, and Sophia Anani-	<i>arXiv:2305.14251</i> .	518
465	adou. 2023. Longdocfactscore: Evaluating the fac-	Nobuyuki Otsu and 1 others. 1975. A threshold selec-	519
466	tuality of long document abstractive summarisation.	tion method from gray-level histograms. <i>Automatica</i> ,	520
467	<i>arXiv preprint arXiv:2309.12455</i> .	11(285-296):23–27.	521
468	Nico Daheim, Clara Meister, Thomas Möllenhoff,	Claude E Shannon. 1948. A mathematical theory of	522
469	and Iryna Gurevych. 2025. Uncertainty-aware de-	communication. <i>The Bell system technical journal</i> ,	523
470	coding with minimum bayes risk. <i>arXiv preprint</i>	27(3):379–423.	524
471	<i>arXiv:2503.05318</i> .	Chenze Shao, Xilin Chen, and Yang Feng. 2018. Greedy	525
472	Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem	search with probabilistic n-gram matching for neural	526
473	Shelmanov, Sergey Petrakov, Haonan Li, Hamdy	machine translation. In <i>Proceedings of the 2018 Con-</i>	527
474	Mubarak, Evgenii Tsybalov, Gleb Kuzmin, Alexan-	<i>ference on Empirical Methods in Natural Language</i>	528
475	der Panchenko, Timothy Baldwin, and 1 others. 2024.	<i>Processing</i> , pages 4778–4784.	529
476	Fact-checking the output of large language mod-	Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang,	530
477	els via token-level uncertainty quantification. <i>arXiv</i>	Yifan Wang, Yujiu Yang, and Wai Lam. 2024. A	531
478	<i>preprint arXiv:2403.04696</i> .	thorough examination of decoding methods in the era	532
479	Wataru Hashimoto, Hidetaka Kamigaito, and Taro	of llms. <i>arXiv preprint arXiv:2402.06925</i> .	533
480	Watanabe. 2025. Decoding uncertainty: The impact	Qwen Team and 1 others. 2024. Qwen2 technical report.	534
481	of decoding strategies for uncertainty estimation in	<i>arXiv preprint arXiv:2407.10671</i> , 2(3).	535
482	large language models. In <i>Findings of the Associa-</i>	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	536
483	<i>tion for Computational Linguistics: EMNLP 2025</i> ,	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	537
484	pages 14601–14613.	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	538
485	Kaifeng He, Mingwei Liu, Chong Wang, Zike Li, Yanlin	Azhar, and 1 others. 2023. Llama: Open and effi-	539
486	Wang, Xin Peng, and Zibin Zheng. 2025. Adadec:	cient foundation language models. <i>arXiv preprint</i>	540
487	Uncertainty-guided adaptive decoding for llm-based	<i>arXiv:2302.13971</i> .	541
488	code generation. <i>arXiv preprint arXiv:2506.08980</i> .	Dennis Ulmer, Chrysoula Zerva, and André FT Mar-	542
489	John Hewitt, Christopher D Manning, and Percy Liang.	tins. 2024. Non-exchangeable conformal language	543
490	2022. Truncation sampling as language model	generation with nearest neighbors. <i>arXiv preprint</i>	544
491	desmoothing. <i>arXiv preprint arXiv:2210.15191</i> .	<i>arXiv:2402.00707</i> .	545
492	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and	Ashwin K Vijayakumar, Michael Cogswell, Ram-	546
493	Yejin Choi. 2019. The curious case of neural text	prasath R Selvaraju, Qing Sun, Stefan Lee, David	547
494	degeneration. <i>arXiv preprint arXiv:1904.09751</i> .	Crandall, and Dhruv Batra. 2016. Diverse beam	548
495	Mengting Hu, Zhen Zhang, Shiwang Zhao, Minlie	search: Decoding diverse solutions from neural se-	549
496	Huang, and Bingzhe Wu. 2023. Uncertainty in natu-	quence models. <i>arXiv preprint arXiv:1610.02424</i> .	550
497	ral language processing: Sources, quantification, and	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-	551
498	applications. <i>arXiv preprint arXiv:2306.04459</i> .	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva	552
499	Daniel Jurafsky and James H. Martin. 2025. <i>Speech</i>	Naik, Arjun Ashok, Arut Selvan Dhanasekaran, An-	553
500	<i>and Language Processing: An Introduction to Natu-</i>	jana Arunkumar, David Stap, Eshaan Pathak, Gian-	554
501	<i>ral Language Processing, Computational Linguistics,</i>	annis Karamanolakis, Haizhi Lai, Ishan Purohit, Is-	555
502	<i>and Speech Recognition with Language Models</i> ,	hani Mondal, Jacob Anderson, Kirby Kuznia, Krima	556
503	3rd edition. Online manuscript released January 12,	Doshi, Kuntal Kumar Pal, and 16 others. 2022.	557
504	2025.	Super-NaturalInstructions: Generalization via declar-	558
505	Hakyung Lee, Subeen Park, Joowang Kim, Sungjun	ative instructions on 1600+ NLP tasks . In <i>Proceed-</i>	559
506	Lim, and Kyungwoo Song. 2025. Uncertainty-aware	<i>ings of the 2022 Conference on Empirical Methods</i>	560
507	contrastive decoding. In <i>Findings of the Association</i>	<i>in Natural Language Processing</i> , pages 5085–5109,	561
508	<i>for Computational Linguistics: ACL 2025</i> , pages	Abu Dhabi, United Arab Emirates. Association for	562
509	26376–26391.	Computational Linguistics.	563
510	Chin-Yew Lin. 2004. Rouge: A package for automatic	Jiacheng Xu, Shrey Desai, and Greg Durrett. 2020. Un-	564
511	evaluation of summaries. In <i>Text summarization</i>	derstanding neural abstractive summarization models	565
512	<i>branches out</i> , pages 74–81.	via uncertainty. <i>arXiv preprint arXiv:2010.07882</i> .	566
513	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike	Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu.	567
514	Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,	2023. Alignscore: Evaluating factual consistency	568
515	Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.	with a unified alignment function. <i>arXiv preprint</i>	569
		<i>arXiv:2305.16739</i> .	570

571 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q
572 Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-
573 uating text generation with bert. *arXiv preprint*
574 *arXiv:1904.09675*.