# Breaking Point Analysis of MCMC Cryptanalysis for Substitution Ciphers

Joel Mathew[1], Ashish L[2]

1(MCA, Nehru College of Engineering and Research Centre, Thrissur, Kerala, India
Email: jooeel777@gmail.com)
2(1(MCA, Nehru College of Engineering and Research Centre, Thrissur, Kerala, India
Email: ashish.mca@ncerc.ac.in)

## Abstract:

Historically, substitution ciphers have been highly vulnerable to frequency analysis [3]. Today, modern computational technique specifically Metropolis-Hastings Markov Chain Monte Carlo (MCMC) algorithms can fully automate this decryption by optimizing a fitness score based on n-gram statistics. However, these statistical methods share a fundamental limitation: they require an adequate sample size. As ciphertext length decreases, its statistical profile diverges too far from standard English, ultimately rendering the attack ineffective. This study experimentally quantifies this 'Breaking Point' the minimum ciphertext length necessary for an MCMC attack to reliably converge on the correct plaintext. Our results reveal a clear threshold: while texts exceeding 200 characters are consistently vulnerable, messages under 150 characters (such as typical SMS texts or tweets) demonstrate significant resistance to this class of automated cryptanalysis.

*Keywords ---- Cryptanalysis, Monoalphabetic Substitution Cipher, Markov Chain Monte Carlo (MCMC), Metropolis-Hastings Algorithm, Quadgram*

## I.    INTRODUCTION

The classical ciphers often rest on the secrecy of the key rather than the length of the message. However, in the context of cryptanalysis, message length is a critical factor. For a monoalphabetic substitution cipher, where each letter in the plaintext is replaced by a fixed letter in the ciphertext, the key space is $26! \approx 4 \times 10^{26}$ (approximately 88.4 bits of entropy), making brute force impossible [4].

Cryptanalysts instead exploit the non-uniform distribution of characters in natural language. In English, 'E' appears with a probability of $\approx 12.7\%$, while 'Z' appears with $\approx 0.07\%$. When a ciphertext preserves these frequency patterns, algorithms can map the high-frequency ciphertext symbols to high-frequency plaintext symbols. The Metropolis-Hastings MCMC algorithm can fully automate this decryption [1].

The "Chen & Rosenthal" style MCMC algorithm automates this by treating the decryption key as a state in a Markov Chain [1]. The algorithm randomly swaps elements of the key and accepts the new key if it produces "more English-like" text. This "English-likeness" is measured using quadgram (4-gram) statistics [2]. The unicity distance for English text is approximately 28 characters, which provides a theoretical lower bound for unique decryption [5].

### A.  The Problem

The success of any statistical attack on a substitution cipher hinge on one core assumption: the ciphertext is long enough to behave like normal English. Natural English has a highly predictable rhythm. The letter 'E' makes up roughly 12.7% of all text, while 'Q' barely breaks 0.1%. Pairs like 'TH' and 'HE' appear constantly. Together, these patterns create the statistical signature that makes frequency analysis possible.

But the signature only exists in large data samples. While dealing with short messages, the Law of Large Numbers breaks down, and random variation dominates. You would never see a novel with zero 'E's and multiple 'Z's, but that exact scenario is entirely possible in a 20-character text message.

### B. Objective

To determine the approximate threshold of text length $L_{crit}$ below which the MCMC attack fails to reliably recover the plaintext.

## II. METHODOLOGY

### A. Algorithm Design

Metropolis-Hastings Markov Chain Monte Carlo (MCMC) solver in Python 3.12, designed specifically for attacking monoalphabetic substitution ciphers. Implementation consists of four core components: fitness function, the transition rule, the acceptance criterion, and the reference data.

Fitness Function (F). We employed log-likelihood scoring based on quadgram statistics [2]. For a decrypted text T of length L, the score is computed as:

$$Score(T) = \sum_{i=0}^{L-4} \log\left(P(T[i:i+4])\right)$$

where P(gram) represents the probability of a specific four-letter sequence occurring in standard English. These probabilities were derived from a large corpus of English text [2], which provides approximately 250,000 characters of diverse. We partitioned this text into training and testing sets, ensuring no overlap between the corpus used to generate quadgram statistics and the corpus used to extract test messages. Note: A text of length L contains exactly L-3 overlapping quadgrams. Our summation runs from index 0 to L-4 inclusive,

capturing all valid four-character windows. For example, a 50-character string yields 47 quadgram evaluations.

Transition Rule. From any current key K, we generate a proposal key K' by selecting two distinct characters uniformly at random from the 26-letter alphabet and swapping their mappings

Acceptance Criterion. Acceptance of a new key is governed by the Metropolis-Hastings criterion [1]. Let $\Delta E$ = Score (Decrypted new) - Score (Decrypted old), representing the change in fitness when moving from the current key to the proposed key. The Metropolis-Hastings acceptance rule operates as follows:

If $\Delta E > 0$: Accept K' unconditionally. The new key produces more English-like text.

If $\Delta E < 0$: Accept K' with probability P = $\exp(\Delta E)$. Since $\Delta E$ is negative, this yields a probability between 0 and 1.

Reference Data. Reference data for the quadgram map was derived from standard English corpora [2]. We generated a custom English quadgram frequency map from our training partition of Sherlock Holmes text. This map contains $26^4$ = 456,976 possible quadgrams, though many have zero observed frequency.

### B. Experiment Setup

To systematically evaluate MCMC performance across practically relevant text lengths: [500, 300, 200, 150, 100, 75, 60, 50, 40, 30, 20] characters. This range spans substantial paragraphs down to brief SMS-style messages.

For each length L, we executed a rigorous protocol. First, we extracted 5 random substrings of length L from our testing partition of Sherlock Holmes text, ensuring strict separation from the training corpus used for quadgram statistics. Second,

we encrypted each substring using a randomly generated monoalphabetic substitution key uniform random permutations of the 26-letter alphabet created via Fisher-Yates shuffle. Third, we attacked each ciphertext with our MCMC solver configured for 20,000 iterations and 3 random restarts. The restart mechanism triggers when no improvement occurs for 5,000 consecutive iterations, generating a fresh random key to escape local maxima traps. Fourth, we measured accuracy as the percentage of correctly recovered characters compared to original plaintext. Note: With only 5 trials per length, variance estimates are approximate.

## III. RESULTS

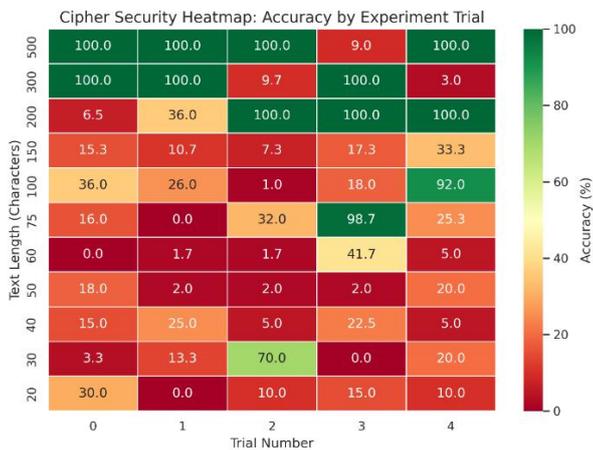"Success" is defined as >90% recovery of plaintext.



Figure 1: Heatmap visualization of decryption accuracy across all experimental trials. The vertical axis lists text lengths in descending order from 500 characters (top) to 20 characters (bottom). The horizontal axis shows trial numbers 0 through 4, representing the five independent repetitions at each length. Each cell is color-coded according to decryption accuracy: deep green indicates 100% successful recovery, bright red indicates 0%

| Len | Trail 1 | Trail 2 | Trail 3 | Trail 4 | Trail 5 | Avg Accuracy | Status |
|-----|---------|---------|---------|---------|---------|--------------|--------|
| 500 | 100% | 100% | 100% | 9% | 100% | 81.8% | Vuln |
| 300 | 100% | 100% | 10% | 100% | 3% | 62.6% | Vuln |
| 200 | 7% | 36% | 100% | 100% | 100% | 68.6% | Vuln |
| 150 | 15% | 11% | 7% | 17% | 33% | 16.6% | Resist |
| 100 | 36% | 26% | 1% | 18% | 92% | 34.6% | Unreli |
| 75 | 16% | 0% | 32% | 99% | 25% | 34.4% | Resist |
| 50 | 18% | 2% | 2% | 2% | 20% | 8.8% | Resist |

recovery (complete failure), and intermediate colors (yellow, orange) represent partial success.

Table 1

Decryption Accuracy Across Text Lengths

Note: Occasional failures at long lengths (e.g., 9% at L=500) represent local maxima. Consistent failure at L=150 indicates fundamental lack of information.

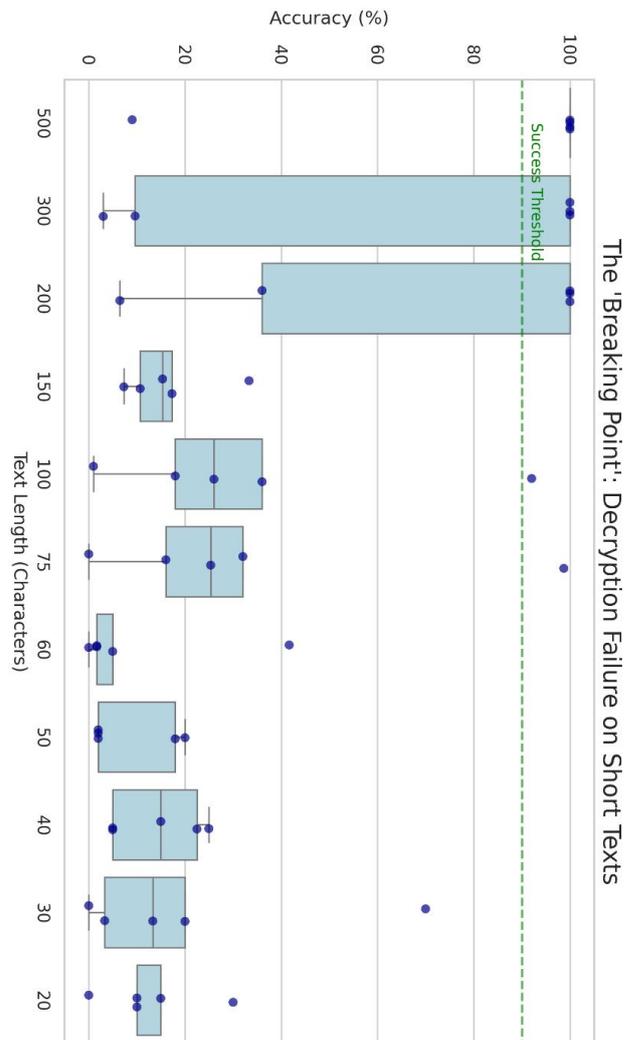### A. *Visualizing the Breaking Point*



Figure 2: A composite view of the experiment's results, combining a box-and-whisker plot with individual data points (jittered scatter plot) for each text length tested. The x-axis represents the length of the ciphertext in characters, while the y-axis measures the decryption accuracy percentage. The green dashed line at 90% serves as the success threshold; points above this line indicate a successful recovery of the original message. The plot vividly demonstrates the "Breaking Point" phenomenon. At lengths of 300 and 500 characters, the data clusters tightly near 100% accuracy, showing consistent convergence. However, around 200 characters, we observe a dramatic increase in variance illustrated by the long "whiskers" and scattered points

where the algorithm fluctuates be- tween perfect success and complete failure. As the length drops below 150, the median accuracy collapses near zero. This visual transition marks the approximate boundary where the signal (English frequency patterns) be- comes overwhelmed by noise, rendering the statistical attack ineffective. The data reveals a sharp decline in recoverability between 200 characters and 150 characters.

Three regimes appear:

• 300-500 chars: Tight clustering near 100% accuracy

• ~200 chars: High variance, mixed success

• <150 chars: Median accuracy collapses near zero

The data reveals a sharp decline between 200 and 150 characters:

• At 200 chars: 100% success in 3/5 trials

• At 150 chars: 0% success (max 33%)

• At 100 chars: One anomalous 92% success (common words), but poor average

Conclusion: The approximate breaking point is 150-175 characters. However, this is probabilistic—the 92% success at L=100 shows specific compositions may still be vulnerable.

### B. Anomalies and Limitations

The 92% success at L=100 likely came from high probability quadgrams ("THAT", "WITH", "FROM") that locked the key early. Thus, <150 chars offer high resistance but not cryptographic security; a "lucky" message might still be recovered.

## IV. CONCLUSION

Automated statistical attacks on monoalphabetic substitution ciphers face a practical hurdle: they require a minimum text length to succeed. In real-world scenarios—where message sizes are often limited by system protocols or user behavior—this boundary dictates whether automated cryptanalysis will work.

Our findings highlight a clear divide. Short, constrained texts like SMS messages, API keys, and 2FA codes prove highly resistant to automated MCMC-based attacks. However, once a message reaches paragraph length (such as an email or news snippet), it becomes extremely vulnerable. Our tests demonstrated that standard hardware can rapidly decrypt these longer texts in a matter of seconds.

Interestingly, the transition between 'secure' and 'vulnerable' is highly volatile. In the intermediate range of around 200 characters, our algorithm's success was essentially a coin toss, achieving perfect decryption 60% of the time, but failing entirely in the other 40%.

Ultimately, this research demonstrates that classical ciphers shouldn't be entirely written off as historical artifacts. Against modern automated attacks, a message's length effectively serves as its own security parameter.

## REFERENCES

[1] J. Chen and J. S. Rosenthal, "Decrypting Classical Cipher Text Using Markov Chain Monte Carlo," Statistics and Computing, vol. 22, no. 2, pp. 397–413, 2012.

[2] Practical Cryptography. (n.d.). Quadgram statistics. [Online]. Available: http://practicalcryptography.com/cryptanalysis/text-characterisation/quadgrams/

[3] C. E. Shannon, "Communication Theory of Secrecy Systems," Bell System Technical Journal, vol. 28, no. 4, pp. 656–715, 1949.

[4] A. J. Menezes, P. C. van Oorschot, and S. A. Vanstone, Handbook of Applied Cryptography, CRC Press, 1996.

[5] P. Diaconis, "The Markov Chain Monte Carlo Revolution," Bulletin of the American Mathematical Society, vol. 46, no. 2, pp. 179–205, 2009.