# Impact of Tokenization Techniques on URL Classification

M. Al-Maamari, M. Istaiti, S. Zerhoudi,
M. Dinzinger, M. Granitzer, J. Mitrovic

*October 10th 2024*

# Agenda

- Introduction

- Background

- Data

- Results

## Research Objectives

- **Study the effect of different tokenization methods on URL classification:**

- **Compare tokenization methods:**

  - Byte Pair Encoding (BPE)

  - Enhanced BPE with GPT-4 generated keywords

  - Punctuation-based splitting

  - Character-level n-grams

- **Assess effects on:**

  - Classification accuracy

  - Computational efficiency

## Importance of Webpage Classification

- **Why Webpage Classification Matters?**

    - Exponential growth of the web.

    - Challenges for search engines and web crawlers.

    - Need for improved crawling efficiency
      and targeted content indexing.



AI generated image: "Smart and cute spider-like robot crawling websites, behind it are floating web pages"

## Tokenization

- **What is Tokenization?**

    - Breaking down text into **smaller units** called **tokens**.

    - Tokens can be words, subwords, characters, or symbols

- **Problem: Lack of Whitespace in URLs**

    - URLs are continuous strings without spaces, unlike regular text.

    - Word-based tokenizers are ineffective for URLs.

## Byte Pair Encoding (BPE)

- **What is BPE?**

  - Merges frequent pairs of characters or sequences.

  - Reduces vocabulary size.

  - Captures **subword** units.

```
https://www.example.com/path/page?query=token
```

BPE

```
https://www.example.com/path/page?query=token
```

## Punctuation Split

- **How it works?**

  - Tokenize URLs at punctuation marks.

  - Splits the URL based on punctuations like:

    - Slashes `/` or `//`        - Dots `.`
    - Question marks `?`       - other marks such as `=` or `:`

```
https://www.example.com/path/page?query=token
```

Punctuation Split

```
https://www.example.com/path/page?query=token
```

## Character-level n-grams

- **How it works?**

    - Captures sequences of characters.

    - Explored n-gram ranges:

        - 1-gram
        - (1-3)-grams
        - (3-6)-grams

        ```
        https://www.example.com/path/page?query=token
        ```

        1-grams

```
['h', 't', 't', 'p', 's', ':', '/', '/', 'w', 'w', 'w', '.', 'e', 'x', 'a', 'm', 'p', 'l', 'e',
    '.', 'c', 'o', 'm', ' ... '?', 'q', 'u', 'e', 'r', 'y', '=', 't', 'o', 'k', 'e', 'n']
```

## Character-level n-grams

```
https://www.example.com/path/page?query=token
```
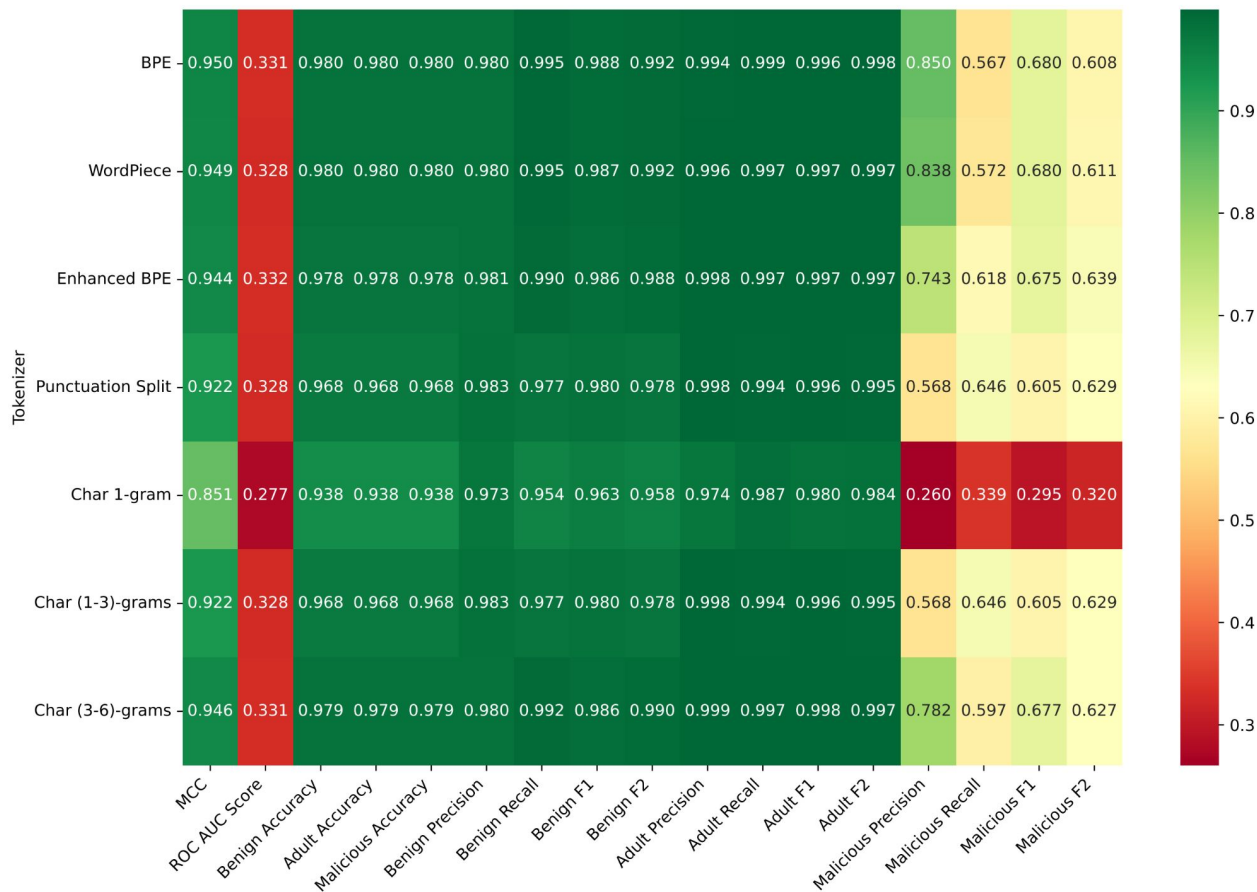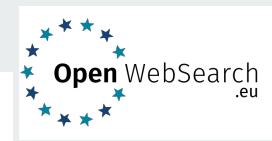
3 grams

```
['htt', 'ttp', 'tps', 'ps:', 's:/', '://', '//w', 'www', 'ww.', 'w.e', '.ex', 'exa', 'xam',
'amp', 'mpl', 'ple', 'le.', 'e.c', '.co', 'com', 'om/', 'm/p', '/pa', 'pat', 'ath', 'th/',
'h/p', '/pa', 'pag', 'age', 'ge?', 'e?q', '?qu', 'que', 'uer', 'ery', 'ry=', 'y=t', '=to',
                          'tok', 'oke', 'ken']
```
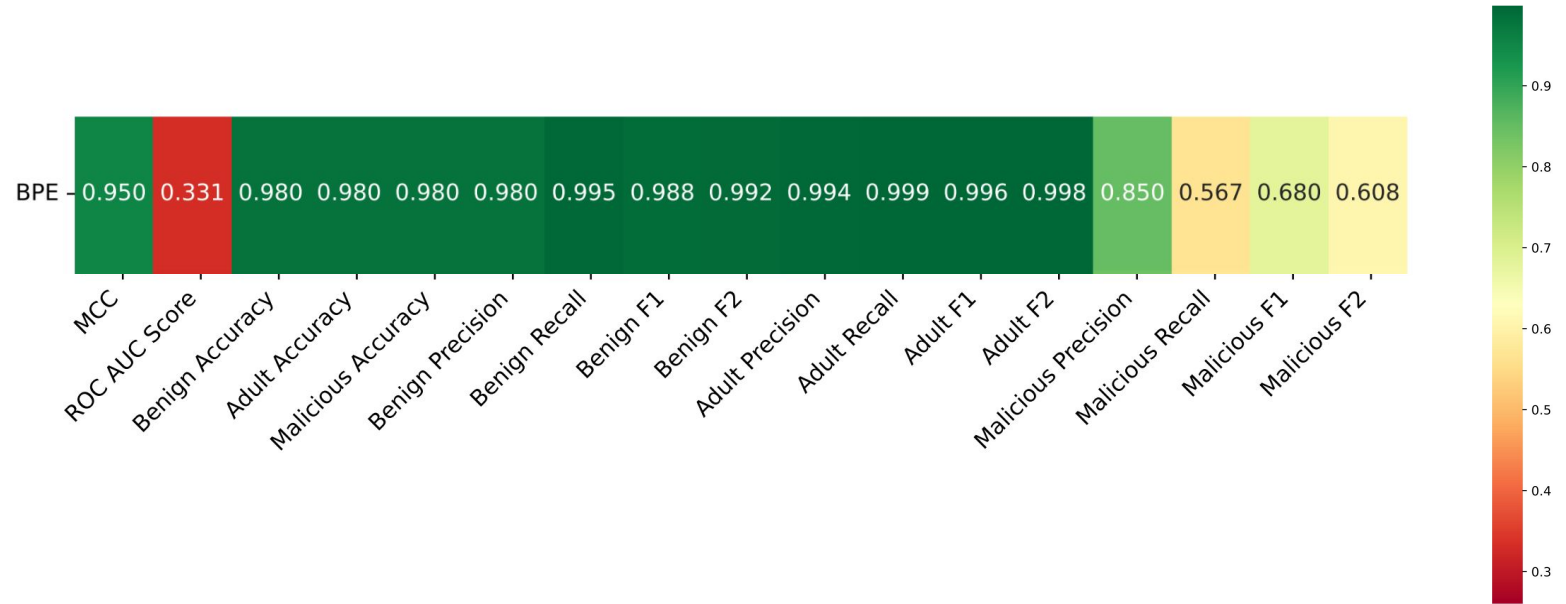
## Dataset Overview:

- **Over 1 million labeled URLs analyzed.**

- **Categories:**
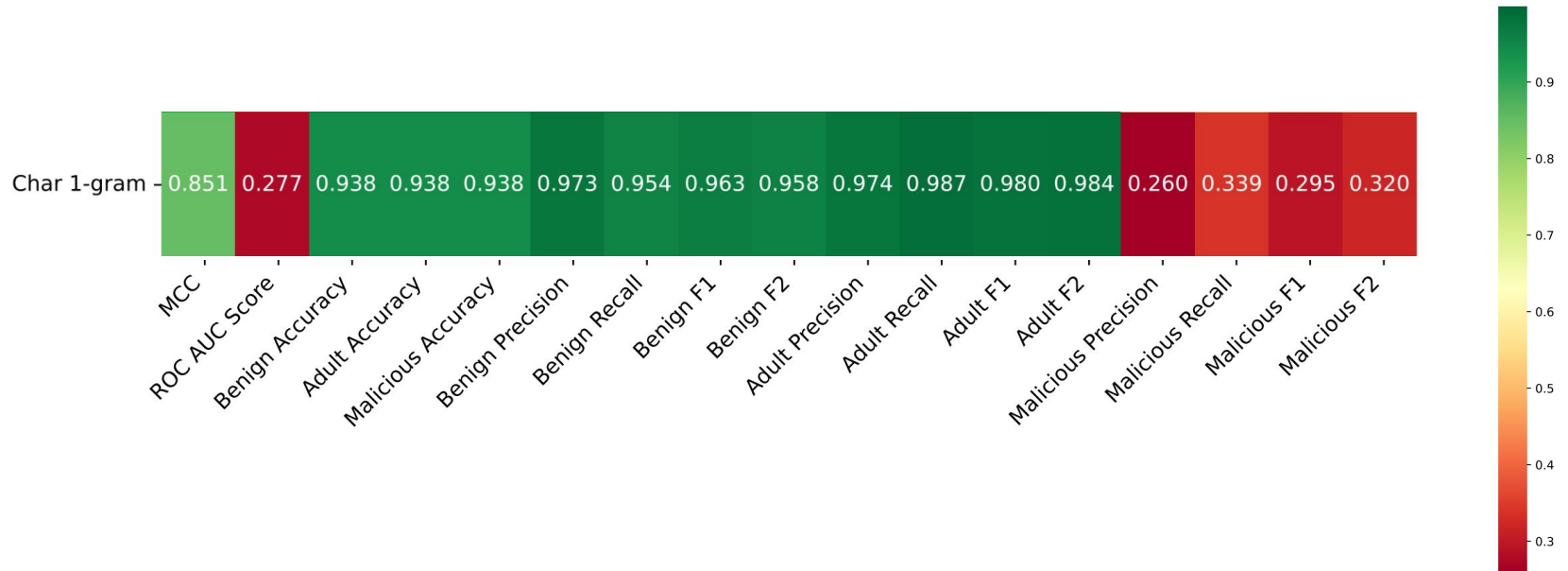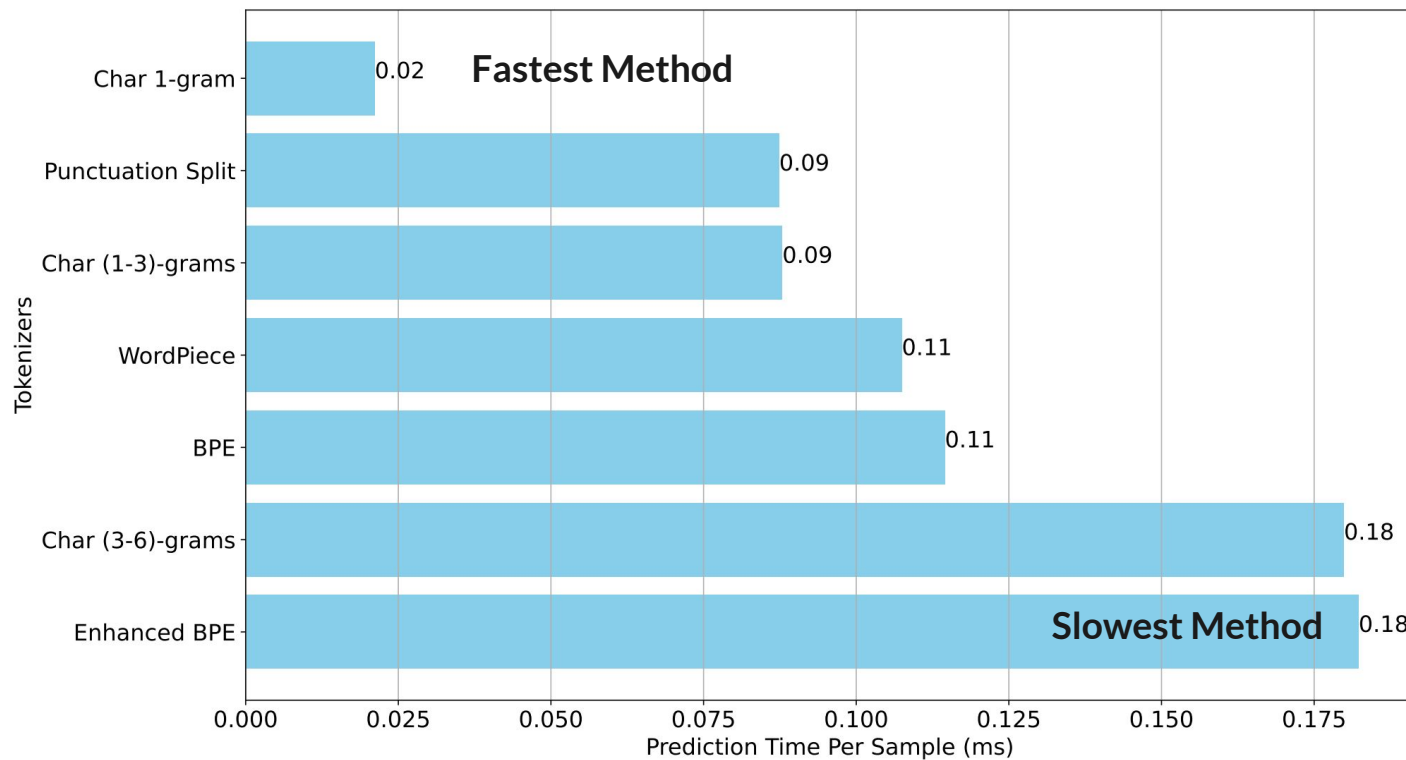
  - Malicious

  - Benign

  - Adult

## **Best Method:** BPE

# Worst Method: BPE

## Key Findings

- **Simple tokenizers are faster but less accurate:**

    - They process data quickly due to their simplicity and limited scope.

    - They lack the depth needed to recognize complex patterns, reducing accuracy.


- **Advanced tokenizers perform better overall:**

    - These methods capture meaningful subword structures, improving understanding.

    - Their complexity leads to better accuracy but increases computational time.

## Key Findings

- **All methods struggle with 'Malicious' URLs:**

    - Malicious URLs mimic benign ones, making them harder to classify accurately.

    - Tokenization alone may not detect subtle differences that indicate malicious intent.

# Thank you for listening!

**Happy to answer any question**

?