

# Text Encoding : Assignment



## What is text encoding and why is it important in computer science?

Text encoding is the method or process of converting a series of characters, i.e, letters, numbers, punctuation, and symbols into a special or unique format for transmission or storage in computers. Data is represented in computers using ASCII, UTF8, UTF32, ISCII, and Unicode encoding schemes.

Text coding is important because machine cant understand text. In order to process text data we need encoding scheme



## Explain the difference between ASCII and Unicode.

Ascii is subset of Unicode. ASCII is lemitd to digits letters from letters for English letter and has fixed length encoding (8 bit only).ASCII uses a 7-bit range to encode just 128 distinct characters

Unicode has variable lentht depending on encoding scheme and represents letters of English, Arabic, Greek, mathematical symbols, historical scripts, and emoji. Unicode represents far more characters than ASCII



## Write a Python function to encode a string using UTF-8 encoding.

```
def encode_string_to_utf8(string):  
    """Encodes a string to UTF-8 encoding.
```

Args:

string: The string to encode.



Returns:

A bytes object representing the encoded string.

```
"""
```

```
return string.encode('utf-8')
```

# Example usage:

```
string = 'Hello, world!'
```

```
encoded_string = encode_string_to_utf8(string)
```

```
print(encoded_string)
```



### **What are the potential issues that can arise when working with different text encodings?**

Unicode is a variable length encoding scheme. If we do not know the source encoding we cannot decode it. Care has to be taken when data is transmitted.

Text encoded in one format but interpreted in another can lead to misinterpretation of characters, resulting in garbled text or incorrect rendering.

Different encodings interpret bytes differently, which can lead to character corruption when a file encoded in one format is read or processed using a different encoding.



### **Explain the concept of endianness in the context of text encoding.**

Endianness is a concept that primarily relates to how bytes are ordered within larger data types, such as integers, floating-point numbers, or characters, in computer memory. In the context of text encoding, endianness becomes relevant when dealing with multibyte character encodings, where characters are represented by sequences of one or more bytes.

There are two common types of endianness:



1. **Big-Endian:** In a big-endian system, the most significant byte (the byte containing the highest-order bits) is stored at the lowest memory address, and the least significant byte is stored at the highest memory address. This means that when reading a multibyte data type, the byte with the highest value (the "big end") comes first.
2. **Little-Endian:** In a little-endian system, the least significant byte (the byte containing the lowest-order bits) is stored at the lowest memory address, and the most significant byte is stored at the highest memory address. This means that when reading a multibyte data type, the byte with the lowest value (the "little end") comes first.



### **What are the advantages and disadvantages of using UTF-16 encoding?**

UTF-16 can represent the entire Unicode character set, including characters from various scripts, symbols, emojis, and special characters. UTF-16 maintains backward compatibility with ASCII. ASCII characters are represented by a single 16-bit code unit in UTF-16, allowing for seamless integration with existing ASCII-based systems and protocols.



### **Discuss the role of byte order marks (BOM) in text encoding.**

Byte Order Mark (BOM) is a special marker used at the beginning of a text file encoded in UTF-8, UTF-16, or UTF-32 to indicate the endianness and Unicode encoding scheme used in the file. The primary role of the BOM is to signal to software reading the file which Unicode encoding variant (UTF-8, UTF-16, or UTF-32) and endianness (big-endian or little-endian) are being used, ensuring proper interpretation of the text.



### **Implement a function in Java to convert a string from UTF-8 to UTF-16 encoding.**





## **Compare and contrast UTF-8 and UTF-32 in terms of efficiency and compatibility.**

**UTF-8:** UTF-8 is a variable-width encoding scheme where each character is represented by one to four bytes. Commonly used characters, such as those in the ASCII character set, are represented by a single byte in UTF-8, making it efficient for encoding text predominantly in the ASCII range. However, characters outside the ASCII range may require multiple bytes, leading to slightly larger file sizes compared to fixed-width encodings like UTF-32, especially for non-Latin scripts or languages with many characters outside the ASCII range.

**UTF-32:** UTF-32 is a fixed-width encoding scheme where each character is represented by exactly four bytes, regardless of the character. This results in a consistent encoding size for all characters, making it less efficient in terms of storage compared to UTF-8, especially for text that primarily consists of characters in the ASCII range. However, UTF-32 is simpler to process since every character occupies a fixed amount of memory, which can lead to performance advantages in certain scenarios, particularly when random access to characters is required.



## **Explain the concept of character set and its role in text encoding.**

A character set, also known as a character repertoire or charset, is a defined collection of characters and symbols that can be used to represent textual data. Each character in a character set is assigned a unique code point, which is typically represented by a numerical value. Character sets play a fundamental role in text encoding by providing the basis for mapping characters to their corresponding binary representations.

