# A Comparative Report on ASCII and Unicode Encoding Systems

## 1. Introduction

Character encoding is a foundational concept in computer science, enabling textual data to be stored, processed, and transmitted electronically. Among the most widely used encoding standards are ASCII and Unicode, each playing a pivotal role in the evolution of modern computing systems. This report examines their structure, historical development, key differences, and technical implications.

## 2. ASCII Encoding

The American Standard Code for Information Interchange (ASCII) was introduced in 1963 by the American National Standards Institute (ANSI) [1]. It uses a 7-bit binary system to represent 128 unique characters (numbered from 0 to 127). These include:

- Uppercase and lowercase English letters
- Numeric digits (09)
- Punctuation marks
- Control characters (e.g., newline, carriage return)

Example Mapping:
Character: A
Decimal (ASCII): 65
Binary: 1000001

Initially, ASCII supported only English due to the geographical focus of early computing industries. As global demand grew, the 8th bitoriginally used for parity/error detectionwas repurposed to extend the character set to 256 characters ($2^8$). This version is known as Extended ASCII.

## 3. Unicode Encoding

To solve the fragmentation problem, the Unicode Consortium introduced the Unicode Standard in 1991 [2]. Unicode aims to provide a universal character set, assigning a unique numbercalled a code pointto every character across all writing systems.

- Unicode supports up to 1,114,112 characters (code points from 0 to 2^21  1).

- As of Unicode 15.0, more than 149,000 characters are defined [3], covering scripts such as Latin, Arabic, Chinese, Cyrillic, and emoji.

Example Code Points:

U+0041   -> A (Latin Capital Letter A)

U+03A9   ->  (Greek Capital Letter Omega)

U+0627   ->  (Arabic Letter Alef)

U+1F600  ->  (Emoji: Grinning Face)

However, Unicode is not an encoding in itselfit is a character set. To store text, Unicode code points must be encoded into bytes using one of several Unicode Transformation Formats (UTFs).

## 4. Unicode Encodings: UTF-8, UTF-16, and UTF-32

Encoding | Bit Size | Type | Characteristics

---------|----------|------|-----------------

UTF-8    | 8 to 32 bits | Variable-length | Compatible with ASCII, widely used on the web

UTF-16   | 16 or 32 bits | Variable-length | Efficient for most global scripts

UTF-32   | 32 bits | Fixed-length | Simplifies indexing, but consumes more space

UTF-8 Highlights:

- Uses 14 bytes per character.

- First 128 characters (U+0000 to U+007F) match ASCII exactly.

- Most commonly used encoding on the internet [4].

UTF-8 Byte Ranges:

ASCII (U+0000 - U+007F): 1 byte: 0xxxxxxx

Extended (U+0080 - U+07FF): 2 bytes: 110xxxxx 10xxxxxx

Extended (U+0800 - U+FFFF): 3 bytes: 1110xxxx 10xxxxxx 10xxxxxx

Extended (U+10000 - U+10FFFF): 4 bytes

UTF-16 and UTF-32 are susceptible to endiannessthe byte order in multi-byte encodings. This led to the introduction of encoding variants such as:

- UTF-16LE (Little Endian)

- UTF-16BE (Big Endian)

- UTF-32LE / UTF-32BE

A Byte Order Mark (BOM) is often used at the beginning of a text stream to indicate encoding and byte order.

## 5. Conclusion

While ASCII was sufficient for early computing, the need for multilingual, cross-platform compatibility led to the development of Unicode, which now underpins nearly all modern computing systems. With its robust structure and flexible encoding schemes, Unicode ensures that text can be stored, transmitted, and rendered accurately across diverse environments and languages.

## References

[1] ANSI X3.4-1963. "American Standard Code for Information Interchange." American National Standards Institute.

[2] The Unicode Consortium. "The Unicode Standard." https://unicode.org

[3] Unicode 15.0. "Enumerated Versions of The Unicode Standard." Unicode.org.

[4] "Usage of Unicode and UTF-8 on the Web." W3Techs - Web Technology Surveys, https://w3techs.com