

```
1: from typing import List
2: import logging as log
3: log.basicConfig(level=log.DEBUG)
4:
5:
6:
7: class Position:
8:     def __init__(self, line, character):
9:         self.line = line
10:        self.character = character
11:
12:    def __repr__(self):
13:        return f"({self.line},{self.character})"
14:
15: def is_char_beyond_multilingual_plane(char: str) → bool:
16:     return ord(char) > 0xFFFF
17:
18:
19: def utf16_unit_offset(chars: str):
20:     """Calculate the number of characters which need two utf-16 code units.
21:     Arguments:
22:         chars (str): The string to count occurrences of utf-16 code units for.
23:     """
24:     return sum(is_char_beyond_multilingual_plane(ch) for ch in chars)
25:
26:
27: def utf16_num_units(chars: str):
28:     """Calculate the length of `str` in utf-16 code units.
29:     Arguments:
30:         chars (str): The string to return the length in utf-16 code units for.
31:     """
32:     return len(chars) + utf16_unit_offset(chars)
33:
34: def position_from_utf16(lines: List[str], position: Position) → Position:
35:     """Convert the position.character from utf-16 code units to utf-32.
36:     A python application can't use the character member of `Position`
37:     directly. As per specification it is represented as a zero-based line and
38:     character offset based on a UTF-16 string representation.
39:     All characters whose code point exceeds the Basic Multilingual Plane are
40:     represented by 2 UTF-16 code units.
41:     The offset of the closing quotation mark in x="😊" is
42:     - 5 in UTF-16 representation
43:     - 4 in UTF-32 representation
44:     see: https://github.com/microsoft/language-server-protocol/issues/376
45:     Arguments:
46:         lines (list):
47:             The content of the document which the position refers to.
48:         position (Position):
49:             The line and character offset in utf-16 code units.
50:     Returns:
51:         The position with `character` being converted to utf-32 code units.
52:     """
53:     if position.line ≥ len(lines):
54:         # start of the line after last
55:         return Position(len(lines), 0) # or return position
56:
57:     line = lines[position.line]
58:
59:     _utf32_len = len(line)
60:     _utf32_index = 0
61:     _utf16_index = 0
62:     while (_utf16_index < position.character) and (_utf32_index < _utf32_len):
63:         _current_char = line[_utf32_index]
64:         is_double_width = is_char_beyond_multilingual_plane(_current_char)
```

```
65:         if (is_double_width):
66:             _utf16_index += 2
67:         else:
68:             _utf16_index += 1
69:             _utf32_index += 1
70:
71:     position = Position(
72:         line=position.line,
73:         character=_utf32_index
74:     )
75:     return position
76:
77:
78: source = """
79: 😊 😊
80: """
81:
82: lines = source.split("\n")
83: print(lines)
84: pos = Position(4,0)
85:
86: pos = position_from_utf16(lines, pos)
87: print(pos)
88:
```