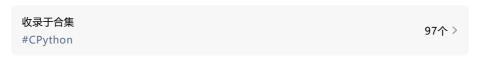
《源码探秘 CPython》22. 字符串是怎么被创建的

原创 古明地觉 古明地觉的编程教室 2022-02-06 09:30





本文我们将讨论字符串是如何被创建的,首先CPython提供了很多的**特型API**用于字符串的创建,我们一起来看一下。

PyUnicode FromString

最简单的一个特型API,根据 C 的char *来创建Python的unicode

```
1 PyObject *
2 PyUnicode_FromString(const char *u)
     // 计算C字符串的长度
4
5
     size t size = strlen(u);
     // 超过了PY_SSIZE_T_MAX,报错
6
7
     if (size > PY_SSIZE_T_MAX) {
        PyErr_SetString(PyExc_OverflowError, "input too long");
8
9
         return NULL;
10
     }
     // 调用该函数进行创建, 这个函数后面介绍
11
      return PyUnicode_DecodeUTF8Stateful(u, (Py_ssize_t)size, NULL, NULL);
12
13 }
```

所以该函数就是根据C的字符串创建Python的字符串,并且从源码中可以看到Python的字符串是由长度限制的,不能超过一个long所表示的最大范围。

PyUnicode_FromStringAndSize

也是根据 C 的**char** *来创建Python的**unicode**,但不同的是可以指定一个长度。而在 PyUnicode_FromString里面,长度写死的,就是 C 字符串的长度。

```
1 PyObject *
2 PyUnicode_FromStringAndSize(const char *u, Py_ssize_t size)
3 {
4
     if (size < 0) {
5
         PyErr_SetString(PyExc_SystemError,
                       "Negative size passed to PyUnicode_FromStringAndS
6
7 ize");
8
         return NULL;
9
10
     // 如果C的char *不为NULL, 那么还是调用这个函数创建
     if (u != NULL)
11
         return PyUnicode_DecodeUTF8Stateful(u, size, NULL, NULL);
12
13
14
      // 否则的话, 直接调用 _PyUnicode_New 申请对应大小的空间
         return (PyObject *)_PyUnicode_New(size);
15
   }
```

当然还有很多其它的创建方式,比如PyUnicode_FromUnicode,根据**wchar_t***创建Python字符串。

这里代码我们就不看了,我们使用 Cython 演示一下。

```
1 from cpython.unicode cimport PyUnicode_FromUnicode
2
3 # Py_UNICODE 理解为 C 的宽字符即可
4 cdef Py_UNICODE* s = "古明地觉的 Python小屋"
5 # 根据宽字符创建字符串,当然也可以指定字符数量
6 # 这里我们只选择前 4 个宽字符
7 cdef str name = PyUnicode_FromUnicode(s, 4)
8 print(name) # 古明地觉
```

然后我们重点看一下PyUnicode_DecodeUTF8Stateful这个函数。

PyUnicode_DecodeUTF8Stateful

先来说一下大致的流程,想象一下我们在Python里面创建一个字符串,比如: **s="你好"**,那么这个过程在 C 中就等价于将这个字符串字面量(const char* 指向)拷贝到 PyUnicode_New 所申请的堆内存中。

以上便是字符串的初始化,而字符串的初始化是以PyUnicode_DecodeUTF8Stateful函数为起点的,这个函数会在内部调用unicode_decode_utf8函数。

在unicode_decode_utf8里面会直接创建PyASCIIObject,因为解释器会假设你创建的是纯ASCII字符串,然后再调用ascii_decode进行判断,如果成功则直接返回,这在程序中属于快分支。

但如果我们创建的不是纯ASCII字符串,那么会判断这个字符串到底属于 PyUnicode_1BYTE_KIND、PyUnicode_2BYTE_KIND、PyUnicode_4BYTE_KIND 三者的哪一种,然后保存在 writer 的kind 成员中。

然后根据不同的类型调用不同的方法,比如: ucs1lib_utf8_decode、ucs2lib utf8 decode 等等。

```
1 static PyObject *
2 unicode_decode_utf8(const char *s, Py_ssize_t size,
                   _Py_error_handler error_handler, const char *errors,
3
                   Py ssize t *consumed)
4
5 {
   _PyUnicodeWriter writer;
6
7
     //直接使用 ascii_decode 进行解码
      //因为会假设是一个ASCII字符串
8
9 writer.pos = ascii_decode(s, end, writer.data);
10 s += writer.pos;
     //逐个字符进行扫描
11
12 while (s < end) {
        //然后获取 kind
13
14
         Py_UCS4 ch;
        int kind = writer.kind;
15
        //判断kind到底是哪一种
16
         //不同的kind执行不同的函数
17
18
        if (kind == PyUnicode_1BYTE_KIND) {
            if (PyUnicode_IS_ASCII(writer.buffer))
19
20
                ch = asciilib_utf8_decode(&s, end, writer.data, &writer.p
```

```
21 os);
             else
22
                 ch = ucs1lib_utf8_decode(&s, end, writer.data, &writer.po
23
24 s);
         } else if (kind == PyUnicode_2BYTE_KIND) {
25
26
              ch = ucs2lib_utf8_decode(&s, end, writer.data, &writer.pos);
27
              assert(kind == PyUnicode_4BYTE_KIND);
28
              ch = ucs4lib_utf8_decode(&s, end, writer.data, &writer.pos);
29
30
31
32
       . . . . . .
```

所以能看出创建字符串的过程就是将**char** * 指向的字节序列进行解码、然后拷贝的一个过程,因为解释器会假设**char** *指向的是一个简单的 ASCII 字节序列。

然后按照PyUnicode_1BYTE_KIND的模式去计算字符串所需的堆内存空间,而后续的字节编码检测算法发现该C级别的字节序列出现了maxchar大于255的字符,就会尝试以PyUnicode_2BYTE_KIND的模式,再次重新计算内存空间并调用malloc分配函数。

因此所以创建一个字符串, 涉及 1 到 3 次的 malloc, 效率还是比较低下的。

```
%timeit s = f"{'a' * 1000000} "

701 μs ± 132 μs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

%timeit s = f"{'a' * 1000000} 憨"

1.24 ms ± 20.9 μs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

%timeit s = f"{'a' * 1000000} \\
""

2.23 ms ± 21 μs per loop (mean ± std. dev. of 7 runs, 1000 loops each)
```

以上就是字符串的创建,事实上字符串还是很复杂的,它并没有我们想象的那么简单。

