

微信扫一扫  
关注该公众号

收录于合集

#CPython

97个 &gt;

## 效率问题

先来探究一下bytes对象的效率，我们知道Python的不可变对象在运算时，处理方式是再创建一个新的。所以三个bytes对象a、b、c在相加时，那么会先根据a + b创建临时对象，然后再根据临时对象 + c创建新的对象，最后返回指针。所以：

```
1 result = b""
2 for _ in bytes_list:
3     result += _
```

这是一种效率非常低下的做法，因为涉及大量临时对象的创建和销毁，不仅是这里字节序列，后面要分析的字符串也是同样的道理。

官方推荐的做法是，使用join，字符串和字节序列都可以对一个列表进行join，将列表里面的多个字符串或者字节序列join在一起。

举个Python中的例子，我们以字符串为例，字节序列同样如此：

```
1 def bad():
2     s = ""
3     for _ in range(1, 10):
4         s += str(_)
5     return s
6
7 def good():
8     l = []
9     for _ in range(1, 10):
10        l.append(str(_))
11    return "".join(l)
12
13 def better():
14    return "".join(str(_) for _ in range(1, 10))
15
16 def best():
17    return "".join(map(str, range(1, 10)))
```

## 字节序列缓存池

为了优化单字节bytes对象的创建效率，Python底层维护了一个缓存池，该缓存池是一个PyBytesObject\*类型的数组。

```
1 static PyBytesObject *characters[ UCHAR_MAX + 1];
```

Python内部创建单字节bytes对象时，先检查目标对象是否已在缓存池中。PyBytes\_FromStringAndSize函数是负责创建bytes对象的一个常用的Python/C API，位于Objects/bytesobject.c中：

```
1 PyObject *
2 PyBytes_FromStringAndSize(const char *str, Py_ssize_t size)
3 {
4     //PyBytesObject对象的指针
```

```

5     PyBytesObject *op;
6     if (size < 0) {
7         //显然size不可以小于0
8         PyErr_SetString(PyExc_SystemError,
9             "Negative size passed to PyBytes_FromStringAndSize");
10        return NULL;
11    }
12    //如果size为1, 表明创建的是单字节对象
13    //当然str不可以为NULL, 而且获取到的字节必须要在characters里面
14    if (size == 1 && str != NULL &&
15        (op = characters[*str & UCHAR_MAX]) != NULL)
16    {
17        #ifdef COUNT_ALLOCS
18            _Py_one_strings++;
19        #endif
20        //增加引用计数, 返回指针
21        Py_INCREF(op);
22        return (PyObject *)op;
23    }
24
25    //否则话创建新的PyBytesObject, 此时是个空
26    op = (PyBytesObject *)_PyBytes_FromSize(size, 0);
27    if (op == NULL)
28        return NULL;
29    if (str == NULL)
30        return (PyObject *) op;
31
32    //不管size是多少, 都直接拷贝即可
33    memcpy(op->ob_sval, str, size);
34    //但是size是1的话, 除了拷贝还会放到缓存池characters中
35    if (size == 1) {
36        characters[*str & UCHAR_MAX] = op;
37        Py_INCREF(op);
38    }
39    //返回其指针
40    return (PyObject *) op;
41 }

```

由此可见, 当Python程序开始运行时, 字节序列缓存池是空的。但随着**单字节bytes对象**的创建, 缓存池中的对象慢慢多了起来。

这样一来, 单字节序列首次创建后便在缓存池中缓存起来; 后续再次使用时, Python直接从缓存池中取, 避免重复创建和销毁。与前面章节介绍的小整数对象池一样, 字节序列缓存池也只能容纳为数不多的 256 个单字节序列, 但使用频率非常高。

缓存池技术作为一种以空间换时间的优化手段, 只需较小的内存为代价, 便可明显提升执行效率。

```

1  >>> a1 = b"a"
2  >>> a2 = b"a"
3  >>> a1 is a2
4  True
5  >>>
6  >>> a1 = b"ab"
7  >>> a2 = b"ab"
8  >>> a1 is a2
9  False
10 >>>

```

显然此时不需要解释了, 单字节bytes对象会缓存起来, 不是单字节则不会缓存。

## 小结

以上就是bytes对象的全部内容，我们说：

- bytes对象是一个变长、不可变对象，内部的值是通过一个C的字符数组来维护的；
- bytes也是序列型操作，它支持的操作在bytes\_as\_sequence和bytes\_as\_mapping中；
- Python内部通过维护字节序列缓存池来优化单字节bytes对象的创建和销毁操作；
- 缓存池是一种常用的以空间换时间的优化技术；

最后，Python除了bytes对象之外，还有一个bytearray对象，它的表现和bytes对象是完全一致的，但一个是可变对象、一个是不可变对象。bytearray对象底层对应的结构体是**PyByteArrayObject**，其相关操作、以及类型对象都位于bytearrayobject.c中，有兴趣可以看一下。

下一篇我们来介绍字符串。

插曲题外话，我被集中隔离了。。。。。。我在4楼办公，而5楼有一个确诊病例。最后显示我和他是密接，要在酒店隔离到二月十号，但是我没中招。所以想提醒大家，一定要注意防护哦，曾以为病例离我很远，没想到它也悄无声息地来到了我的身边。

收录于合集 #CPython 97

[< 上一篇](#)

《源码探秘 CPython》19. 字符集和字符编码

[下一篇 >](#)

《源码探秘 CPython》17. bytes 对象的行为（下）

文章已于2022-06-12修改

喜欢此内容的人还喜欢

SQL注入之-手注断了不如工具系列  
千寻安服



一文剖析MySQL主从复制异常错误代码13114  
TtrOpsStack



如何设计一个缓存函数  
Web技术学苑

