

## 《源码探秘 CPython》82. Python运行时环境的初始化，解释器在启动时都做了什么？

原创 古明地觉 古明地觉的编程教室 2022-05-07 08:30 发表于北京

收录于合集  
#CPython

97个 >



微信扫一扫  
关注该公众号



我们之前完成了 Python 的字节码、以及虚拟机的剖析工作，但这仅仅只是一部分，而其余的部分则被遮在了幕后。记得我们在分析虚拟机的时候，曾这么说过：

解释器启动时，首先会进行“运行时环境”的初始化，关于“运行时环境”的初始化是一个非常复杂的过程。并且“运行时环境”和“执行环境”是不同的，“运行时环境”是一个全局的概念，而“执行环境”是一个栈帧。关于“运行时环境”后面会单独分析，这里就假设初始化动作已经完成，我们已经站在了虚拟机的门槛外面，只需要轻轻推动第一张骨牌，整个执行过程就会像多米诺骨牌一样，一环扣一环地展开。

所以这次，我们将回到时间的起点，从 Python 的应用程序被执行开始，一步一步紧紧跟随 Python 的轨迹，完整地展示解释器在启动之初的所有动作。当我们了解所有的初始化动作之后，也就能对 Python 执行引擎执行字节码指令时的整个运行环境了如指掌了。



我们知道线程是操作系统调度的最小单元，那么 Python 的线程又是怎样的呢？

启动一个 Python 线程，底层会启动一个 C 线程，然后启动操作系统的一个原生线程（OS 线程）。所以 Python 的线程实际上是对 OS 线程的一个封装，因此 Python 的线程是货真价实的。

然后 Python 还提供了一个 PyThreadState 对象，也就是线程状态对象，维护 OS 线程执行的状态信息，相当于是 OS 线程的一个抽象描述。

虽然真正用来执行的线程及其状态肯定是由操作系统进行维护的，但是 Python 虚拟机在运行的时候总需要另外一些与线程相关的状态和信息，比如是否发生了异常等等，这些信息显然操作系统是没有办法提供的。

而 PyThreadState 对象正是为 OS 线程准备的、在虚拟机层面保存其状态信息的对象，也就是线程状态对象。在 Python 中，当前活动的 OS 线程对应的 PyThreadState 对象可以通过 PyThreadState\_GET 获得，有了线程状态对象之后，就可以设置一些额外信息了。具体内容，我们后面会说。

当然除了线程状态对象之外，还有进程状态对象，我们来看看两者在底层的定义是什么？它们位于 Include/pystate.h 中。

```
1 // ts 是 thread state 的简写
2 typedef struct _ts PyThreadState;
3 // is 是 interpreter state 的简写
4 typedef struct _is PyInterpreterState;
```

里面的 PyThreadState 表示线程状态对象，PyInterpreterState 表示进程状态对象，但它们都是 typedef 起的一个别名。前者是 struct \_ts 的别名，后者是 struct \_is 的别名，我们来看一下它们长什么样。

```

1 // Include/cpython/pystate.h
2 struct _ts {
3     //多个线程状态对象也会像链表一样串起来
4     //因为一个进程里面是可以包含多个线程的
5     //prev 指向上一个线程状态对象
6     //next 指向下一个线程状态对象
7     struct _ts *prev;
8     struct _ts *next;
9     //进程状态对象, 标识该线程属于哪一个进程
10    PyInterpreterState *interp;
11
12    //栈帧对象, 模拟线程中函数的调用堆栈
13    struct _frame *frame;
14    //递归深度
15    int recursion_depth;
16    //.....
17    //线程 id
18    uint64_t id;
19 };

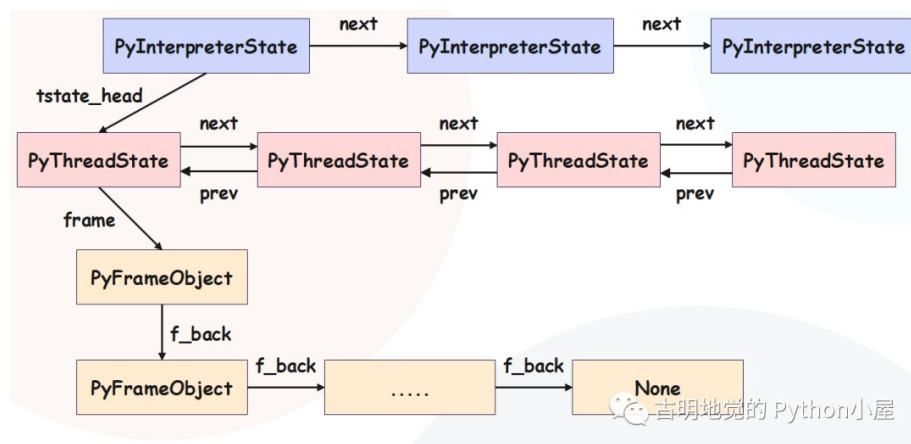
```

```

1 //Include/internal/pycore_pystate.h
2 struct _is {
3     //Python支持多进程, 多个进程也会以链表的形式进行组织
4     //当前进程的下一个进程
5     struct _is *next;
6     //进程环境中的线程状态对象的集合
7     //我们说线程状态对象会形成一个链表
8     //这里就是链表的头结点
9     struct _ts *tstate_head;
10    //进程id
11    int64_t id;
12    //....
13    PyObject *audit_hooks;
14 };

```

所以 PyInterpreterState 对象可以看成是对进程的模拟, PyThreadState 是对线程的模拟。我们之前分析虚拟机的时候说过其执行环境, 如果再将运行时环境加进去的话。



进程状态对象的 `tstate_head` 指向了线程状态对象, 对应当前活跃的 Python 线程; 每个线程状态对象的 `frame` 都指向当前正在执行的栈帧对象。

在解释器启动之后，初始化的动作是从 `Py_NewInterpreter` 函数开始的，然后这个函数调用了 `new_interpreter` 函数完成初始化。至于这两个函数长什么样一会再聊，先往后看。

我们知道当操作系统在运行一个可执行文件时，首先创建一个进程内核。同理在 Python 中亦是如此，会在 `new_interpreter` 中调用 `PyInterpreterState_New` 创建一个崭新的 `PyInterpreterState` 对象。该函数位于 `Python/pystate.c` 中。

```
1 PyInterpreterState *
2 PyInterpreterState_New(void)
3 {
4     //申请进程状态对象所需要的内存
5     PyInterpreterState *interp = PyMem_RawMalloc(sizeof(PyInterpreterSta
6 te));
7     if (interp == NULL) {
8         return NULL;
9     }
10
11     //设置属性
12     //.....
13     //.....
14     return interp;
15 }
```

关于进程状态对象我们不做过多解释，只需要知道解释器在启动时，会创建一个 `PyInterpreterState` 对象。如果开启了多进程，那么内部会继续创建，然后通过 `next` 指针将多个 `PyInterpreterState` 串成一个链表结构。

在调用 `PyInterpreterState_New` 成功创建 `PyInterpreterState` 之后，会再接再厉，调用 `PyThreadState_New` 创建一个全新的线程状态对象，相关函数定义同样位于 `Python/pystate.c` 中。

```
1 PyThreadState *
2 PyThreadState_New(PyInterpreterState *interp)
3 {
4     return new_threadstate(interp, 1);
5 }
```

我们注意到这个函数接收一个 `PyInterpreterState`，这说明了线程是依赖于进程的，因为需要进程给自己分配资源，然后这个函数又调用了 `new_threadstate`。除了传递 `PyInterpreterState` 之外，还传了一个 1，想也不用想这肯定是创建的线程数量。这里创建 1 个，也就是主线程（main thread）。

```
1 static PyThreadState *
2 new_threadstate(PyInterpreterState *interp, int init)
3 {
4     _PyRuntimeState *runtime = &_amp;PyRuntime;
5     //为线程状态对象申请内存
6     PyThreadState *tstate = (PyThreadState *)PyMem_RawMalloc(sizeof(PyTh
7 readState));
8     if (tstate == NULL) {
9         return NULL;
10    }
11    //设置从线程中获取函数调用栈的操作
12    if (_PyThreadState_GetFrame == NULL) {
13        _PyThreadState_GetFrame = threadstate_getframe;
14    }
15
16    //设置该线程所在的进程
17    tstate->interp = interp;
18
19    //下面就是设置内部的成员属性
20    //比如栈帧、递归深度、线程 id
```

```

21     tstate->frame = NULL;
22     tstate->recursion_depth = 0;
23     tstate->id = ++interp->tstate_next_unique_id;
24     //.....
25     //.....
26     //.....
27     //上一个线程状态对象
28     tstate->prev = NULL;
29     //当前线程状态对象的 next
30     //我们看到指向了线程状态对象链表的头结点
31     tstate->next = interp->tstate_head;
32     if (tstate->next)
33         //因为每个线程状态对象的prev指针都要指向上一个线程状态对象
34         //如果是头结点的话, 那么prev就是 NULL
35         //但由于新的线程状态对象在插入之后显然就变成了链表的新的头结点
36         //因此还需要将插入之前的头结点的prev指向新插入的线程状态对象
37         tstate->next->prev = tstate;
38     //将tstate设置为新的线程状态对象(链表的头结点)
39     interp->tstate_head = tstate;
40
41     //返回线程状态对象
42     return tstate;
    }

```

和 PyInterpreterState\_New 相同, PyThreadState\_New 会申请内存, 创建线程状态对象, 并且对每个成员进行初始化。其中 prev 指针和 next 指针分别指向了上一个线程状态对象和下一个线程状态对象。

而且也肯定会存在某一时刻, 存在多个 PyThreadState 对象形成一个链表, 那么什么时刻会发生这种情况呢? 显然用鼻子想也知道这是在启动多线程的时候。

此外我们看到 Python 在插入线程状态对象的时候采用的是头插法。

从源码中我们看到, 虚拟机设置了从线程中获取函数调用栈的操作, 所谓函数调用栈就是前面说的 PyFrameObject 对象链表。而且在源码中, PyThreadState 关联了 PyInterpreterState, PyInterpreterState 也关联了 PyInterpreterState。

到目前为止, 仅有的两个对象建立起了联系。对应到操作系统, 我们说进程和线程建立了联系。

而在两者建立了联系之后, 那么就很容易在 PyInterpreterState 和 PyThreadState 之间穿梭。并且在 Python 运行时环境中, 会有一个变量 (先卖个关子) 一直维护着当前活动的线程, 更准确的说是当前活动线程 (OS 线程) 对应的 PyThreadState 对象。

初始时, 该变量为 NULL。在 Python 启动之后创建了第一个 PyThreadState 之后, 会调用 PyThreadState\_Swap 函数来设置这个变量, 函数位于 Python/pystate.c 中。

```

1 PyThreadState *
2 PyThreadState_Swap(PyThreadState *newts)
3 {
4     //调用了_PyThreadState_Swap, 里面传入了两个参数
5     //第一个我们后面说, 从名字上看显然这是和GIL相关的
6     //第二个参数就是新创建的线程状态对象
7     return _PyThreadState_Swap(&_PyRuntime.gilstate, newts);
8 }

```

内部又调用了 \_PyThreadState\_Swap。

```

1 PyThreadState *
2 _PyThreadState_Swap(struct _gilstate_runtime_state *gilstate, PyThreadState
3 *newts)
4 {
5     //这里是获取当前的线程状态对象
6     //并且保证线程的安全性
7     PyThreadState *oldts = _PyRuntime.GILState_GetThreadState(gilstate);

```

```

8      //将 GIL 交给 newts, 也就是新建、即将获取执行权的线程状态对象
9      _PyRuntimeGILState_SetThreadState(gilstate, newts);
10     //....
11     return oldts;
12 }

```

所以逻辑很容易理解，有一个变量始终维护着当前活跃线程对应的线程状态对象，初始时它是个 NULL。而一旦解释器启动，并创建了第一个线程状态对象（显然对应主线程），那么就会将创建的线程状态对象交给这个变量保存。

如果调用 `_PyThreadState_Swap` 的时候，发现保存线程状态对象的变量不为 NULL，那么说明开启了多线程。变量保存的就是代码中的 `oldts`，也就是当前活动线程对应的线程状态对象，可由于它的时间片耗尽，解释器会剥夺它的执行权，然后交给 `newts`。那么 `newts` 就成为了新的当前活跃线程对应的线程状态对象，那么它也要交给变量进行保存。

而通过 `_PyThreadState_Swap` 可以看到，想要实现这一点，主要依赖两个宏。

```

// 通过&(gilstate)->tstate_current获取当前线程
#define _PyRuntimeGILState_GetThreadState(gilstate) \
    ((PyThreadState*)_Py_atomic_load_relaxed(&(gilstate)->tstate_current))

// 将newts设置为当前线程，可以理解为发生了线程的切换
#define _PyRuntimeGILState_SetThreadState(gilstate, value) \
    _Py_atomic_store_relaxed(&(gilstate)->tstate_current, \
        (uintptr_t)(value))

```

古明地觉的 Python 小屋

然后这两个宏里面出现了 `_Py_atomic_load_relaxed`, `_Py_atomic_store_relaxed` 和 `&(gilstate)->tstate_current`，这些又是什么呢？还有到底哪个变量在维护着当前活动线程对应的线程状态对象呢？其实那两个宏已经告诉你了。

```

1  //Include/internal/pycore_pystate.h
2  struct _gilstate_runtime_state {
3      //...
4      //宏里面出现的 gilstate 就是该结构体实例
5      //tstate_current指的就是当前活动的 OS 线程对应的状态对象
6      //同时也可以理解为获取到 GIL 的 Python线程
7      _Py_atomic_address tstate_current;
8      //...
9  };
10
11
12 //Include/internal/pycore_atomic.h
13 #define _Py_atomic_load_relaxed(ATOMIC_VAL) \
14     _Py_atomic_load_explicit((ATOMIC_VAL), _Py_memory_order_relaxed)
15
16 #define _Py_atomic_store_relaxed(ATOMIC_VAL, NEW_VAL) \
17     _Py_atomic_store_explicit((ATOMIC_VAL), (NEW_VAL), _Py_memory_order_
18 relaxed)
19
20 #define _Py_atomic_load_explicit(ATOMIC_VAL, ORDER) \
21     atomic_load_explicit(&((ATOMIC_VAL)->_value), ORDER)
22
23 #define _Py_atomic_store_explicit(ATOMIC_VAL, NEW_VAL, ORDER) \
24     atomic_store_explicit(&((ATOMIC_VAL)->_value), NEW_VAL, ORDER)

```

不难发现：

- `_Py_atomic_load_relaxed` 调用了 `_Py_atomic_load_explicit`, `_Py_atomic_load_explicit` 又调用了 `atomic_load_explicit`;
- `_Py_atomic_store_relaxed` 调用了 `_Py_atomic_store_explicit`, `_Py_atomic_store_explicit` 调用了 `atomic_store_explicit`;

而 `atomic_load_explicit` 和 `atomic_store_explicit` 是系统头文件 `stdatomic.h` 中定义的 api，这是在系统的 api 中修改的，所以说是线程安全的。

介绍完中间部分的内容，那么我们可以从头开始分析 Python 运行时的初始化了，我们说它是从 Py\_NewInterpreter 开始的。

```
1 // Python/pylifecycle.c
2 PyThreadState *
3 Py_NewInterpreter(void)
4 {
5     //线程状态对象
6     PyThreadState *tstate = NULL;
7     //传入线程对象, 调用 new_interpreter
8     PyStatus status = new_interpreter(&tstate);
9     //异常检测
10    if (_PyStatus_EXCEPTION(status)) {
11        Py_ExitStatusException(status);
12    }
13    //返回线程状态对象
14    return tstate;
15 }
```

另外里面出现了一个 PyStatus，表示程序执行的状态，会检测是否发生了异常。

```
1 //Include/cpython/initconfig.h
2 typedef struct {
3     enum {
4         _PyStatus_TYPE_OK=0,
5         _PyStatus_TYPE_ERROR=1,
6         _PyStatus_TYPE_EXIT=2
7     } _type;
8     const char *func;
9     const char *err_msg;
10    int exitcode;
11 } PyStatus;
```

然后我们的重点是 new\_interpreter函数，进程状态对象的创建就是在这个函数里面发生的。

```
1 //Python/pylifecycle.c
2 static PyStatus
3 new_interpreter(PyThreadState **tstate_p)
4 {
5     PyStatus status;
6
7     //运行时初始化, 如果出现异常直接返回
8     status = _PyRuntime_Initialize();
9     if (_PyStatus_EXCEPTION(status)) {
10        return status;
11    }
12    //.....
13    // 创建一个进程状态对象
14    PyInterpreterState *interp = PyInterpreterState_New();
15    //.....
16    //根据进程状态对象创建一个线程状态对象
17    //维护对应 OS 线程的状态
18    PyThreadState *tstate = PyThreadState_New(interp);
19    //将GIL的控制权交给创建的线程
20    PyThreadState *save_tstate = PyThreadState_Swap(tstate);
21    //...
22 }
```

Python在初始化运行时环境时，肯定也要对类型系统进行初始化等等，整体是一个非常庞大的过程。

到这里，我们对 new\_interpreter 算是有了一个阶段性的成功，我们创建了代表进程和线程概念

的 PyInterpreterState 和 PyThreadState 对象，并且在它们之间建立的联系。下面，new\_interpreter 将进行入另一个环节，设置系统 module。



在 new\_interpreter 中创建了 PyInterpreterState 和 PyThreadState 对象之后，就会开始设置系统的 \_\_builtins\_\_ 了。

```
1 static PyStatus
2 new_interpreter(PyThreadState **tstate_p)
3 {
4     //....
5     //申请一个PyDictObject对象, 用于 sys.modules
6     PyObject *modules = PyDict_New();
7     if (modules == NULL) {
8         return _PyStatus_ERR("can't make modules dictionary");
9     }
10
11     //然后让 interp -> modules 维护 modules
12     //由于 interp 表示的是进程状态对象, 这说明什么?
13     //显然是该进程内的多个线程共享同一个 sys.modules
14     interp->modules = modules;
15
16     //加载sys模块, 所有的module对象都在sys.modules中
17     PyObject *sysmod = _PyImport_FindBuiltin("sys", modules);
18     if (sysmod != NULL) {
19         interp->sysdict = PyModule_GetDict(sysmod);
20         if (interp->sysdict == NULL) {
21             goto handle_error;
22         }
23         Py_INCREF(interp->sysdict);
24         PyDict_SetItemString(interp->sysdict, "modules", modules);
25         if (_PySys_InitMain(runtime, interp) < 0) {
26             return _PyStatus_ERR("can't finish initializing sys");
27         }
28     }
29
30     //加载内置模块 builtins
31     //可以import builtins, 并且builtins.list等价于list
32     PyObject *bimod = _PyImport_FindBuiltin("builtins", modules);
33     if (bimod != NULL) {
34         //设置 __builtins__
35         interp->builtins = PyModule_GetDict(bimod);
36         if (interp->builtins == NULL)
37             goto handle_error;
38         Py_INCREF(interp->builtins);
39     }
40     //.....
41 }
```

整体还是比较清晰和直观的，另外我们说内置名字空间是由进程来维护的，因为进程就是用来为线程提供资源的。但是也能看出，一个进程内的多个线程共享同一个内置作用域。

显然这是非常合理的，不可能每开启一个线程，就为其创建一个 \_\_builtins\_\_。我们来从Python的角度证明这一点：

```
1 import threading
2 import builtins
3
```

```

4 def foo1():
5     builtins.list, builtins.tuple = builtins.tuple, builtins.list
6
7 def foo2():
8     print(f"猜猜下面代码会输出什么:")
9     print("list:", list([1, 2, 3, 4, 5]))
10    print("tuple:", tuple([1, 2, 3, 4, 5]))
11
12 f1 = threading.Thread(target=foo1)
13 f1.start()
14 f1.join()
15 threading.Thread(target=foo2).start()
16 """
17 猜猜下面代码会输出什么:
18 list: (1, 2, 3, 4, 5)
19 tuple: [1, 2, 3, 4, 5]
20 """

```

所有的内置对象和内置函数都在内置名字空间里面，可以通过 `import builtins` 获取、也可以直接通过 `__builtins__` 这个变量来获取。

我们在 `foo1` 中把 `list` 和 `tuple` 互换了，而这个结果显然也影响到了 `foo2` 函数。这也说明了 `__builtins__` 是属于进程级别的，它是被多个线程共享的，所以是 **interp -> modules = modules**。当然这个 `modules` 是 `sys.modules`，因为不止内置名字空间，所有的 `module` 对象都是被多个线程共享的。

而对 `__builtins__` 的初始化是在 `_PyBuiltin_Init` 函数中进行的。

```

1 //Python/builtinmodule.c
2 PyObject *
3 _PyBuiltin_Init(void)
4 {
5     PyObject *mod, *dict, *debug;
6
7     const PyConfig *config = &_amp;PyInterpreterState_GET_UNSAFE()->config;
8
9     if (PyType_Ready(&PyFilter_Type) < 0 ||
10         PyType_Ready(&PyMap_Type) < 0 ||
11         PyType_Ready(&PyZip_Type) < 0)
12         return NULL;
13
14     //创建并设置 __builtins__
15     mod = _PyModule_CreateInitialized(&builtinsmodule, PYTHON_API_VERSION);
16
17     if (mod == NULL)
18         return NULL;
19     //拿到 __builtins__ 的属性字典
20     dict = PyModule_GetDict(mod);
21     //将所有内置对象加入到 __builtins__ 中
22     //.....
23     //老铁们, 下面这些东西应该不陌生吧
24     SETBUILTIN("None", Py_None);
25     SETBUILTIN("Ellipsis", Py_Ellipsis);
26     SETBUILTIN("NotImplemented", Py_NotImplemented);
27     SETBUILTIN("False", Py_False);
28     SETBUILTIN("True", Py_True);
29     SETBUILTIN("bool", &PyBool_Type);
30     SETBUILTIN("memoryview", &PyMemoryView_Type);
31     SETBUILTIN("bytearray", &PyByteArray_Type);
32     SETBUILTIN("bytes", &PyBytes_Type);
33     SETBUILTIN("classmethod", &PyClassMethod_Type);
34     SETBUILTIN("complex", &PyComplex_Type);
35     SETBUILTIN("dict", &PyDict_Type);

```



```

36     SETBUILTIN("enumerate", &PyEnum_Type);
37     SETBUILTIN("filter", &PyFilter_Type);
38     SETBUILTIN("float", &PyFloat_Type);
39     SETBUILTIN("frozenset", &PyFrozenSet_Type);
40     SETBUILTIN("property", &PyProperty_Type);
41     SETBUILTIN("int", &PyLong_Type);
42     SETBUILTIN("list", &PyList_Type);
43     SETBUILTIN("map", &PyMap_Type);
44     SETBUILTIN("object", &PyBaseObject_Type);
45     SETBUILTIN("range", &PyRange_Type);
46     SETBUILTIN("reversed", &PyReversed_Type);
47     SETBUILTIN("set", &PySet_Type);
48     SETBUILTIN("slice", &PySlice_Type);
49     SETBUILTIN("staticmethod", &PyStaticMethod_Type);
50     SETBUILTIN("str", &PyUnicode_Type);
51     SETBUILTIN("super", &PySuper_Type);
52     SETBUILTIN("tuple", &PyTuple_Type);
53     SETBUILTIN("type", &PyType_Type);
54     SETBUILTIN("zip", &PyZip_Type);
55     debug = PyBool_FromLong(config->optimization_level == 0);
56     if (PyDict_SetItemString(dict, "__debug__", debug) < 0) {
57         Py_DECREF(debug);
58         return NULL;
59     }
60     Py_DECREF(debug);
61
62     return mod;
63 #undef ADD_TO_ALL
64 #undef SETBUILTIN
65 }

```

整个 `_PyBuiltin_Init` 函数的功能就是设置好 `__builtins__` module，而这个过程是分为两步的。

- 通过 `_PyModule_CreateInitialized` 函数创建 `PyModuleObject` 对象，我们知道这是Python模块对象的底层实现；
- 设置 module，将Python的内置对象都塞到 `__builtins__` 中；

但是我们看到设置的东西似乎少了不少，比如 `dir`、`hasattr`、`setattr` 等等，这些明显也是内置的，但是它们到哪里去了。别急，我们刚才说创建 `__builtins__` 分为两步，第一步是创建 `PyModuleObject`，而使用的函数就是 `_PyModule_CreateInitialized`，而在这个函数里面就已经完成了大部分的 `__builtins__` 设置工作。

```

1 //Object/moduleobject.c
2 PyObject *
3 _PyModule_CreateInitialized(struct PyModuleDef* module, int module_api_v
4 ersion)
5 {
6     const char* name;
7     PyModuleObject *m;
8
9     //初始化
10    if (!PyModuleDef_Init(module))
11        return NULL;
12    //拿到 module 的 name
13    //对于当前来说就是__builtins__
14    name = module->m_name;
15    //这里比较有意思，这是检测模块版本的，针对的是需要导入的py文件。
16    //我们说编译成PyCodeObject对象之后
17    //会直接从当前目录的__pycache__里面导入，那里面都是pyc文件
18    //而pyc文件的文件名是有Python解释器的版本号的
19    //这里就是比较版本是否一致，不一致则不导入pyc文件，而是会重新编译py文件
20    if (!check_api_version(name, module_api_version)) {
21        return NULL;
22    }

```

```

23     if (module->m_slots) {
24         PyErr_Format(
25             PyExc_SystemError,
26             "module %s: PyModule_Create is incompatible with m_slots", n
27         ame);
28         return NULL;
29     }
30     //创建一个PyModuleObject
31     if ((m = (PyModuleObject*)PyModule_New(name)) == NULL)
32         return NULL;
33
34     //.....
35     if (module->m_methods != NULL) {
36         //遍历methods中指定的module对象中应包含的操作集合
37         if (PyModule_AddFunctions((PyObject *) m, module->m_methods) != 0
38 ) {
39             Py_DECREF(m);
40             return NULL;
41         }
42     }
43     if (module->m_doc != NULL) {
44         //设置docstring
45         if (PyModule_SetDocString((PyObject *) m, module->m_doc) != 0) {
46             Py_DECREF(m);
47             return NULL;
48         }
49     }
50     m->md_def = module;
51     return (PyObject*)m;
52 }

```

根据上面的代码我们可以得出如下信息：

- 1) **name**: **module** 对象的名称，在这里就是 `__builtins__`；
- 2) **module\_api\_version**: Python 内部使用的version值，用于比较；
- 3) **PyModule\_New**: 用于创建一个PyModuleObject对象；
- 4) **methods**: 该module中所包含的函数的集合，在这里是**builtin\_methods**；
- 5) **PyModule\_AddFunctions**: 设置methods中的函数操作；
- 6) **PyModule\_SetDocString**: 设置docstring；



创建module对象

Python 的 module对象在底层对应的结构体是 PyModuleObject对象，我们来看看它长什么样子吧。

```

1 //Objects/moduleobject.c
2 typedef struct {
3     //头部信息
4     PyObject_HEAD
5     //属性字典，所有的属性和值都在里面
6     PyObject *md_dict;
7     //module对象包含的操作集合，里面是一些结构体
8     //每个结构体对应一个函数的相关信息
9     struct PyModuleDef *md_def;
10    //...
11    PyObject *md_name; //模块名
12 } PyModuleObject;

```

而这个对象是通过PyModule\_New创建的。

```

1 PyObject *
2 PyModule_New(const char *name)
3 {

```

```

4 //module对象的 name、PyModuleObject *
5 PyObject *nameobj, *module;
6 nameobj = PyUnicode_FromString(name);
7 if (nameobj == NULL)
8     return NULL;
9 //创建 PyModuleObject
10 module = PyModule_NewObject(nameobj);
11 Py_DECREF(nameobj);
12 return module;
13 }
14
15
16 PyObject *
17 PyModule_NewObject(PyObject *name)
18 {
19     //module对象的指针
20     PyModuleObject *m;
21     //为 module 对象申请空间
22     //模块具有属性字典, 所以是一个变长对象
23     m = PyObject_GC_New(PyModuleObject, &PyModule_Type);
24     if (m == NULL)
25         return NULL;
26     //设置相应属性, 初始化为NULL
27     m->md_def = NULL;
28     m->md_state = NULL;
29     m->md_weaklist = NULL;
30     m->md_name = NULL;
31     //属性字典
32     m->md_dict = PyDict_New();
33     //调用module_init_dict
34     if (module_init_dict(m, m->md_dict, name, NULL) != 0)
35         goto fail;
36     PyObject_GC_Track(m);
37     return (PyObject *)m;
38
39 fail:
40     Py_DECREF(m);
41     return NULL;
42 }
43
44 static int
45 module_init_dict(PyModuleObject *mod, PyObject *md_dict,
46                 PyObject *name, PyObject *doc)
47 {
48     _Py_IDENTIFIER(__name__);
49     _Py_IDENTIFIER(__doc__);
50     _Py_IDENTIFIER(__package__);
51     _Py_IDENTIFIER(__loader__);
52     _Py_IDENTIFIER(__spec__);
53
54     if (md_dict == NULL)
55         return -1;
56     if (doc == NULL)
57         doc = Py_None;
58     //模块的一些属性、__name__、__doc__ 等等
59     if (_PyDict_SetItemId(md_dict, &PyId__name__, name) != 0)
60         return -1;
61     if (_PyDict_SetItemId(md_dict, &PyId__doc__, doc) != 0)
62         return -1;
63     if (_PyDict_SetItemId(md_dict, &PyId__package__, Py_None) != 0)
64         return -1;
65     if (_PyDict_SetItemId(md_dict, &PyId__loader__, Py_None) != 0)
66         return -1;
67     if (_PyDict_SetItemId(md_dict, &PyId__spec__, Py_None) != 0)

```

```

68         return -1;
69     if (PyUnicode_CheckExact(name)) {
70         Py_INCREF(name);
71         Py_XSETREF(mod->md_name, name);
72     }
73
74     return 0;
75 }

```

这里虽然创建了一个 module 对象，但是这仅仅是一个空的 module 对象，却并没有包含相应的操作和数据。我们看到只设置了 name 和 doc 等属性。



### 设置 module 对象的属性

在 PyModule\_New 结束之后，程序继续执行 \_PyModule\_CreateInitialized 内部的代码，然后通过 PyModule\_AddFunctions 完成了对 \_\_builtins\_\_ 大部分属性的设置。这个设置的属性依赖于第二个参数 methods，在这里为 builtin\_methods。然后会遍历 builtin\_methods，并处理每一项元素，我们还是来看看长什么样子。

```

1  //Python/builtinmodule.c
2
3  static PyMethodDef builtin_methods[] = {
4      {"__build_class__", (PyCFunction)(void*)(void))builtin__build_class__,
5      METH_FASTCALL | METH_KEYWORDS, build_class_doc},
6      {"__import__", (PyCFunction)(void*)(void))builtin__import__,
7      METH_VARARGS | METH_KEYWORDS, import_doc},
8      BUILTIN_ABS_METHODDEF
9      BUILTIN_ALL_METHODDEF
10     BUILTIN_ANY_METHODDEF
11     BUILTIN_ASCII_METHODDEF
12     BUILTIN_BIN_METHODDEF
13     {"breakpoint", (PyCFunction)(void*)(void))builtin_breakpoint,
14     METH_FASTCALL | METH_KEYWORDS, breakpoint_doc},
15     BUILTIN_CALLABLE_METHODDEF
16     BUILTIN_CHR_METHODDEF
17     BUILTIN_COMPILE_METHODDEF
18     BUILTIN_DELATTR_METHODDEF
19     {"dir", builtin_dir, METH_VARARGS, dir_doc},
20     BUILTIN_DIVMOD_METHODDEF
21     BUILTIN_EVAL_METHODDEF
22     BUILTIN_EXEC_METHODDEF
23     BUILTIN_FORMAT_METHODDEF
24     {"getattr", (PyCFunction)(void*)(void))builtin_getattr, METH_FASTCALL, getattr_doc},
25     BUILTIN_GLOBALS_METHODDEF
26     BUILTIN_HASATTR_METHODDEF
27     BUILTIN_HASH_METHODDEF
28     BUILTIN_HEX_METHODDEF
29     BUILTIN_ID_METHODDEF
30     BUILTIN_INPUT_METHODDEF
31     BUILTIN_ISINSTANCE_METHODDEF
32     BUILTIN_ISSUBCLASS_METHODDEF
33     {"iter", (PyCFunction)(void*)(void))builtin_iter,
34     METH_FASTCALL, iter_doc},
35     BUILTIN_LEN_METHODDEF
36     BUILTIN_LOCALS_METHODDEF
37     {"max", (PyCFunction)(void*)(void))builtin_max,
38     METH_VARARGS | METH_KEYWORDS, max_doc},
39     {"min", (PyCFunction)(void*)(void))builtin_min,
40     METH_VARARGS | METH_KEYWORDS, min_doc},
41     {"next", (PyCFunction)(void*)(void))builtin_next,
42     METH_FASTCALL, next_doc},

```

```

45     BUILTIN_OCT_METHODDEF
46     BUILTIN_ORD_METHODDEF
47     BUILTIN_POW_METHODDEF
48     {"print",          (PyCFunction)(void(*) (void)) builtin_print,
    METH_FASTCALL | METH_KEYWORDS, print_doc},
    BUILTIN_REPR_METHODDEF
    BUILTIN_ROUND_METHODDEF
    BUILTIN_SETATTR_METHODDEF
    BUILTIN_SORTED_METHODDEF
    BUILTIN_SUM_METHODDEF
    {"vars",            builtin_vars,          METH_VARARGS, vars_doc},
    {NULL,              NULL},
};

```

怎么样，是不是看到了玄机。

总结一下就是：在 `Py_NewInterpreter` 里面调用 `new_interpreter` 函数，然后在 `new_interpreter` 这个函数里面，通过 `PyInterpreterState_New` 创建 `PyInterpreterState`，然后传递 `PyInterpreterState` 调用 `PyThreadState_New` 创建 `PyThreadState`。

接着就是执行各种初始化动作，然后在 `new_interpreter` 中调用 `_PyBuiltin_Init` 设置内置属性，在代码的最后会设置内置属性(函数、对象)。但是有很多却不在里面，比如：`dir`、`getattr`等等。

所以中间调用的 `_PyModule_CreateInitialized` 不仅仅是初始化一个 `module` 对象，还会在初始化之后将我们没有看到的一些属性设置进去。在里面先使用 `PyModule_New` 创建一个 `PyModuleObject`，此时它内部只有 `__name__` 和 `__doc__` 等属性，之后再通过 `PyModule_AddFunctions` 设置 `methods`，在这里面我们看到了 `dir`、`getattr` 等内置函数。当这些属性设置完之后，退回到 `_PyBuiltin_Init` 函数中，再设置剩余的部分属性。之后，`__builtins__` 就完成了。

另外 `builtin_methods` 是一个 `PyMethodDef` 类型的数组，里面是一个个的 `PyMethodDef` 结构体。

```

1 //Include/methodobject.h
2 struct PyMethodDef {
3     /* 内置函数的名称 */
4     const char *ml_name;
5     //实现对应逻辑的 C 函数, 但是需要转成 PyCFunction类型
6     //主要是为了更好的处理关键字参数
7     PyCFunction ml_meth;
8
9     /* 参数类型
10    #define METH_VARARGS 0x0001  扩展位置参数
11    #define METH_KEYWORDS 0x0002  扩展关键字参数
12    #define METH_NOARGS 0x0004  不需要参数
13    #define METH_O 0x0008  需要一个参数
14    #define METH_CLASS 0x0010  被classmethod装饰
15    #define METH_STATIC 0x0020  被staticmethod装饰
16    */
17    int ml_flags;
18
19    //函数的__doc__
20    const char *ml_doc;
21 };
22 typedef struct PyMethodDef PyMethodDef;

```

对于这里面每一个 `PyMethodDef`，`_PyModule_CreateInitialized` 都会基于它创建一个 `PyCFunctionObject` 对象。我们说过内置函数在底层对应 `PyCFunctionObject`。

```

1 typedef struct {
2     //头部信息
3     PyObject_HEAD
4     //PyMethodDef

```

```

5     PyMethodDef *m_ml;
6     //self参数
7     PyObject      *m_self;
8     //__module__属性
9     PyObject      *m_module;
10    //弱引用列表, 不讨论
11    PyObject      *m_weakreflist;
12    //为了效率而单独实现的调用函数
13    vectorcallfunc vectorcall;
14 } PyCFunctionObject;

```

然后 PyCFunctionObject 的创建则是通过 PyCFunction\_New 完成的。

```

1 //Objects/methodobject.c
2 PyObject *
3 PyCFunction_New(PyMethodDef *ml, PyObject *self)
4 {
5     return PyCFunction_NewEx(ml, self, NULL);
6 }
7
8
9 PyObject *
10 PyCFunction_NewEx(PyMethodDef *ml, PyObject *self, PyObject *module)
11 {
12     vectorcallfunc vectorcall;
13     //判断参数类型
14     switch (ml->ml_flags & (METH_VARARGS | METH_FASTCALL | METH_NOARGS |
15 METH_O | METH_KEYWORDS))
16     {
17         case METH_VARARGS:
18         case METH_VARARGS | METH_KEYWORDS:
19             vectorcall = NULL;
20             break;
21         case METH_FASTCALL:
22             vectorcall = cfunction_vectorcall_FASTCALL;
23             break;
24         case METH_FASTCALL | METH_KEYWORDS:
25             vectorcall = cfunction_vectorcall_FASTCALL_KEYWORDS;
26             break;
27         case METH_NOARGS:
28             vectorcall = cfunction_vectorcall_NOARGS;
29             break;
30         case METH_O:
31             vectorcall = cfunction_vectorcall_0;
32             break;
33         default:
34             PyErr_Format(PyExc_SystemError,
35                          "%s() method: bad call flags", ml->ml_name);
36             return NULL;
37     }
38
39     PyCFunctionObject *op;
40     //我们看到这里也采用了缓存池的策略
41     op = free_list;
42     if (op != NULL) {
43         free_list = (PyCFunctionObject *) (op->m_self);
44         (void)PyObject_INIT(op, &PyCFunction_Type);
45         numfree--;
46     }
47     else {
48         //否则重新申请
49         op = PyObject_GC_New(PyCFunctionObject, &PyCFunction_Type);
50         if (op == NULL)
51             return NULL;

```

```

52     }
53     //设置属性
54     op->m_weakreflist = NULL;
55     op->m_ml = ml;
56     Py_XINCRREF(self);
57     op->m_self = self;
58     Py_XINCRREF(module);
59     op->m_module = module;
60     op->vectorcall = vectorcall;
61     _PyObject_GC_TRACK(op);
62     return (PyObject *)op;
    }

```

在 `_PyBuiltin__Init` 之后，虚拟机会把 `PyModuleObject` 对象的属性字典抽取出来，赋值给 `interp -> builtins`。

```

1  //moduleobject.c
2  PyObject *
3  PyModule_GetDict(PyObject *m)
4  {
5      PyObject *d;
6      if (!PyModule_Check(m)) {
7          PyErr_BadInternalCall();
8          return NULL;
9      }
10     d = ((PyModuleObject *)m) -> md_dict;
11     assert(d != NULL);
12     return d;
13 }
14
15
16 static PyStatus
17 new_interpreter(PyThreadState **tstate_p)
18 {
19     //.....
20     PyObject *bimod = _PyImport_FindBuiltin("builtins", modules);
21     if (bimod != NULL) {
22         //通过PyModule_GetDict获取属性字典，赋值给builtins
23         interp->builtins = PyModule_GetDict(bimod);
24         if (interp->builtins == NULL)
25             goto handle_error;
26         Py_INCREF(interp->builtins);
27     }
28     else if (PyErr_Occurred()) {
29         goto handle_error;
30     }
31     //.....
32 }

```

以后Python在需要访问 `__builtins__` 时，直接访问 `interp->builtins` 就可以了，不需要再到 `interp->modules` 里面去找了。因为内置函数、属性的使用会很频繁，所以这种加速机制是有效的。



Python在创建并设置了 `__builtins__` 之后，会照猫画虎，用同样的流程来设置 `sys module`，并像

设置 interp->builtins 一样设置 interp->sysdict.

```
1 //Python/pylifecyle.c
2 static PyStatus
3 new_interpreter(PyThreadState **tstate_p)
4 {
5     //.....
6     PyObject *sysmod = _PyImport_FindBuiltin("sys", modules);
7     if (sysmod != NULL) {
8         interp->sysdict = PyModule_GetDict(sysmod);
9         if (interp->sysdict == NULL) {
10             goto handle_error;
11         }
12         Py_INCREF(interp->sysdict);
13         //设置
14         PyDict_SetItemString(interp->sysdict, "modules", modules);
15         if (_PySys_InitMain(runtime, interp) < 0) {
16             return _PyStatus_ERR("can't finish initializing sys");
17         }
18     }
19     //.....
20 }
```

Python在创建了sys module之后, 会在此module中设置一个Python搜索module时的默认路径集合。

```
1 //Python/pylifecyle.c
2 static PyStatus
3 new_interpreter(PyThreadState **tstate_p)
4 {
5     //.....
6     status = add_main_module(interp);
7     //.....
8 }
9
10
11 static PyStatus
12 add_main_module(PyInterpreterState *interp)
13 {
14     PyObject *m, *d, *loader, *ann_dict;
15     //将__main__添加进sys.modules中
16     m = PyImport_AddModule("__main__");
17     if (m == NULL)
18         return _PyStatus_ERR("can't create __main__ module");
19
20     d = PyModule_GetDict(m);
21     ann_dict = PyDict_New();
22     if ((ann_dict == NULL) ||
23         (PyDict_SetItemString(d, "__annotations__", ann_dict) < 0)) {
24         return _PyStatus_ERR("Failed to initialize __main__.__annotation
25 s__");
26     }
27     Py_DECREF(ann_dict);
28
29     if (PyDict_GetItemString(d, "__builtins__") == NULL) {
30         PyObject *bimod = PyImport_ImportModule("builtins");
31         if (bimod == NULL) {
32             return _PyStatus_ERR("Failed to retrieve builtins module");
33         }
34         if (PyDict_SetItemString(d, "__builtins__", bimod) < 0) {
35             return _PyStatus_ERR("Failed to initialize __main__.__builtin
36 s__");
37         }
38         Py_DECREF(bimod);
```



```

39     }
40     loader = PyDict_GetItemString(d, "__loader__");
41     if (loader == NULL || loader == Py_None) {
42         PyObject *loader = PyObject_GetAttrString(interp->importlib,
43                                                     "BuiltinImporter");
44         if (loader == NULL) {
45             return _PyStatus_ERR("Failed to retrieve BuiltinImporter");
46         }
47         if (PyDict_SetItemString(d, "__loader__", loader) < 0) {
48             return _PyStatus_ERR("Failed to initialize __main__.__loader__");
49         }
50     }
51     Py_DECREF(loader);
52     }
53     return _PyStatus_OK();
54 }

```

根据我们使用Python的经验，我们知道最终Python肯定会创建一个PyListObject对象，也就是sys.path，里面包含了一组PyUnicodeObject\*，指向的每一个PyUnicodeObject的内容都代表了一个搜索路径。但是这一步不是在这里完成的，至于是在哪里完成的，我们后面会说。

另外需要注意的是：在上面的逻辑中，解释器将\_\_main\_\_这个模块添加进去了，这个\_\_main\_\_估计不用我多说了。之前在PyModule\_New中，创建一个PyModuleObject对象之后，会在其属性字典中插入一个名为 "\_\_name\_\_" 的 key，value 就是 "\_\_main\_\_"。但是对于当然模块来说，这个模块也叫做 \_\_main\_\_。

```

1 name = "古明地觉"
2 import __main__
3 print(__main__.name) # 古明地觉
4
5 import sys
6 print(sys.modules["__main__"] is __main__) # True

```

我们发现这样也是可以导入的，因为这个\_\_main\_\_就是模块本身。

```

1 static PyStatus
2 add_main_module(PyInterpreterState *interp)
3 {
4     PyObject *m, *d, *loader, *ann_dict;
5     //创建__main__ module, 并将其加入到 interp->modules 中
6     m = PyImport_AddModule("__main__");
7     if (m == NULL)
8         return _PyStatus_ERR("can't create __main__ module");
9     //获取__main__的属性字典
10    d = PyModule_GetDict(m);
11
12    //获取interp->modules中的__builtins__ module
13    if (PyDict_GetItemString(d, "__builtins__") == NULL) {
14        PyObject *bimod = PyImport_ImportModule("builtins");
15        if (bimod == NULL) {
16            return _PyStatus_ERR("Failed to retrieve builtins module");
17        }
18        //将 "__builtins__", __builtins__
19        //加入到__main__ module的dict中
20        if (PyDict_SetItemString(d, "__builtins__", bimod) < 0) {
21            return _PyStatus_ERR("Failed to initialize __main__.__builtin
22 s__");
23        }
24        Py_DECREF(bimod);
25    }
26    //.....
27 }

```

因此我们算是知道了，为什么python xxx.py 执行的时候，\_\_name\_\_ 是 "\_\_main\_\_" 了，因为我

们这里设置了。而Python沿着名字空间寻找的时候，最终会在 `__main__` 的 local 空间中发现 `__name__`，且值为字符串 `"__main__"`。但如果是通过 `import` 的方式加载的，那么 `__name__` 则不是 `"__main__"`，而是模块名。

其实这个 `__main__` 我们是再熟悉不过的了，当输入 `dir()` 的时候，就会显示 `__main__` 的内容。`dir` 是可以不加参数的，如果不加参数，那么默认访问当前的 local 空间，也就是 `__main__`。

```
1 >>> __name__
2 '__main__'
3 >>>
4 >>> __builtins__.__name__
5 'builtins'
6 >>>
7 >>> import numpy as np
8 >>> np.__name__
9 'numpy'
10 >>>
```

所以说，访问模块就类似访问变量一样。`modules` 里面存放了所有的 **<模块名, 模块>**。当我们调用 `np` 的时候，是找到 `name` 为 `"numpy"` 的模块，然后这个值里面也维护了一个字典，其中也有一个 `key` 为 `__name__` 的 `entry`，`value` 为 `"numpy"`。



## 设置site-specific的module的搜索路径



Python是一个非常开放的体系，它的强大来源于丰富的第三方库，这些库由外部的 `py` 文件来提供。当使用这些第三方库的时候，只需要简单地进行 `import` 即可。一般来说，这些第三方库都放在 `Lib/site-packages` 中，如果程序想使用这些库，直接导入即可。

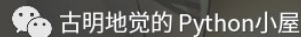
但是到目前为止，我们好像也没看到 Python 将 `site-packages` 路径设置到搜索路径里面去啊。其实在完成了 `__main__` 的创建之后，Python 才腾出手来收拾这个 `site-package`。这个关键的动作在于 Python 的一个标准库：`site.py`。

我们先来将 `Lib` 目录下的 `site.py` 删掉，然后导入一个第三方模块，看看会有什么后果。

```
Fatal Python error: init_import_size: Failed to import the site module
Python runtime state: initialized
ModuleNotFoundError: No module named 'site'

Current thread 0x0000343c (most recent call first):
<no Python frame>

Process finished with exit code 1
```



因此 Python 在初始化的过程中确实导入了 `site.py`，所以才有了如下的输出。而这个 `site.py` 也正是 Python 能正确加载位于 `site-packages` 目录下第三方包的关键所在。我们可以猜测，应该就是这个 `site.py` 将 `site-packages` 目录加入到了 `sys.path` 中，而这个动作是由 `init_import_size` 完成的。

```
1 static PyStatus
2 new_interpreter(PyThreadState **tstate_p)
3 {
4     //.....
5     if (config->site_import) {
6         status = init_import_size();
7         if (_PyStatus_EXCEPTION(status)) {
8             return status;
9         }
10    }
```

```

10     }
11     //.....
12 }
13
14
15 static PyStatus
16 init_import_size(void)
17 {
18     PyObject *m;
19     m = PyImport_ImportModule("site");
20     if (m == NULL) {
21         //这里的报错信息是不是和上图中显示的一样呢？
22         return _PyStatus_ERR("Failed to import the site module");
23     }
24     Py_DECREF(m);
25     return _PyStatus_OK();
26 }

```

在 `init_import_size` 中，只调用了 `PyImport_ImportModule` 函数，这个函数是import机制的核心所在。

比如 `PyImport_ImportModule("numpy")` 等价于 `import numpy`。

以上就是运行时环境的初始化所做的事情，但是需要注意：此时虚拟机还没有启动。对，上面的那些工作都只是前戏，是虚拟机启动之前所做的一些准备工作。而在准备工作做完之后，虚拟机就会正式启动。那么怎么启动呢？我们下一篇文章再聊。

收录于合集 [#CPython 97](#)

[< 上一篇](#)

《源码探秘 CPython》83. 激活 Python 虚拟机

[下一篇 >](#)

《源码探秘 CPython》81. import 机制是怎么实现的？

喜欢此内容的人还喜欢

从零开始学 Python 之高阶函数  
豆豆的杂货铺

(x)



从零开始学 Python 之递归函数  
豆豆的杂货铺

(x)



Django笔记二十七之数据库函数之文本函数  
Django笔记

(x)

