

CasonChan

** 学习笔记记录，精彩博文收录 **

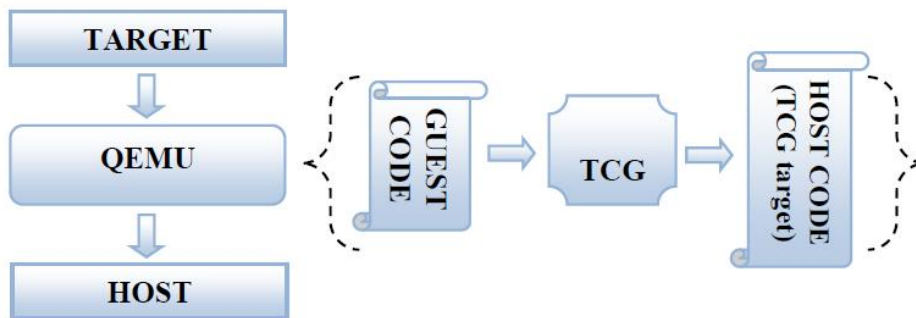
qemu源码架构

前言：本文主要概括了QEMU的代码结构，特别从代码翻译的角度分析了QEMU是如何将客户机代码翻译成TCG代码和主机代码并且最终执行的过程。并且在最后描述了QEMU和KVM之间联系的纽带。

申明：本文前面部分从qemu detailed study第七章翻译而来。

1.代码结构

如我们所知，QEMU是一个模拟器，它能够动态模拟特定架构的CPU指令，如X86，PPC，ARM等等。QEMU模拟的架构叫目标架构，运行 QEMU的系统架构叫主机架构，QEMU中有一个模块叫做微型代码生成器（TCG），它用来将目标代码翻译成主机代码。如下图所示。



我们也可以将运行在虚拟cpu上的代码叫做客户机代码，QEMU的主要功能就是不断提取客户机代码并且转化成主机指定架构的代码。整个翻译任务分为两个部分：第一个部分是将做目标代码（TB）转化成TCG中间代码，然后再将中间代码转化成主机代码。

QEMU的代码结构非常清晰但是内容非常复杂，这里先简单分析一下总体的结构

1. 开始执行：

主要比较重要的c文件有：/vl.c,/cpus.c, /exec-all.c, /exec.c, /cpu-exec.c.

QEMU的main函数定义在/vl.c中，它也是执行的起点，这个函数的功能主要是建立一个虚拟的硬件环境。它通过参数的解析，将初始化内存，需要的模拟的设备初始化，CPU参数，初始化KVM等等。接着程序就跳转到其他的执行分支文件如：/cpus.c, /exec-all.c, /exec.c, /cpu-exec.c.

2. 硬件模拟

所有的硬件设备都在/hw/ 目录下面，所有的设备都有独自的文件，包括总线，串口，网卡，鼠标等等。它们通过设备模块串在一起，在vl.c中的machine _init中初始化。这里就不讲每种设备是怎么实现的了。

3.目标机器

现在QEMU模拟的CPU架构有：Alpha, ARM, Cris, i386, M68K, PPC, Sparc, Mips, MicroBlaze, S390X and SH4.

我们在QEMU中使用./configure 可以配置运行的架构，这个脚本会自动读取本机真实机器的CPU架构，并且编译的时候就编译对应架构的代码。对于不同的QEMU做的事情都不同，所以不同架构下的代码在不同的目录下面。/target-arch/目录就对应了相应架构的代码，如/target-i386/就对应了x86系列的代码部分。虽然不同架构做法不同，但是都是为了实现将对客户机CPU架构的TBs转化成TCG的中间代码。这个就是TCG的前半部分。

4.主机

公告

昵称： [CasonChan](#)
园龄： 6年4个月
粉丝： 47
关注： 4
[+加关注](#)

导航

[博客园](#)
[首页](#)
[新随笔](#)
[联系](#)
[订阅](#) [RSS](#)
[管理](#)

< 2020年4月 >						
日	一	二	三	四	五	六
29	30	31	1	2	3	4
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30	1	2
3	4	5	6	7	8	9

统计

随笔 - 97
文章 - 0
评论 - 8
引用 - 0

搜索

常用链接

[我的随笔](#)
[我的评论](#)
[我的参与](#)
[最新评论](#)
[我的标签](#)

我的标签

[Kernel\(7\)](#)

随笔分类

[Ceph\(1\)](#)
[Docker\(1\)](#)
[Framework\(2\)](#)
[HTTP\(1\)](#)
[KVM\(4\)](#)
[Linux\(47\)](#)
[Network\(37\)](#)
[OpenFlow\(7\)](#)
[OpenStack\(9\)](#)
[OpenvSwitch\(15\)](#)
[Python\(2\)](#)
[Qemu\(3\)](#)

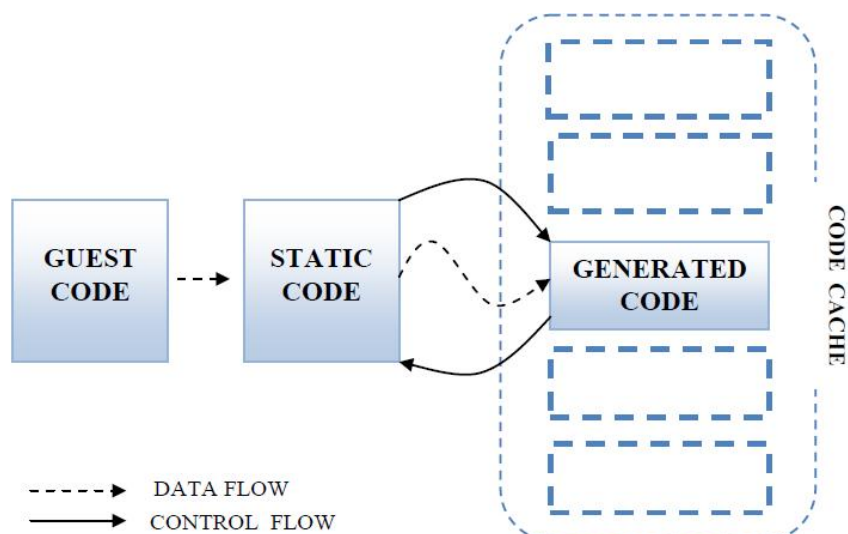
这个部分就是使用TCG代码生成主机的代码，这部分代码在/tcg/里面，在这个目录里面也对应了不同的架构，分别在不同的子目录里面，如i386就在/tcg/i386中。整个生成主机代码的过程也可以教TCG的后半部分。

5.文件总结和补充：

/vl.c:	最主要的模拟循环，虚拟机环境初始化，和CPU的执行。
/target-arch/translate.c	将客户机代码转化成不同架构的TCG操作码。
/tcg/tcg.c	主要的TCG代码。
/tcg/arch/tcg-target.c	将TCG代码转化生成主机代码
/cpu-exec.c	其中的cpu-exec()函数主要寻找下一个TB（翻译代码块），如果没找到就请求得到下一个TB，并且操作生成的代码块。

2. TCG - 动态翻译

QEMU在 0.9.1版本之前使用DynGen翻译c代码.当我们需要的时候TCG会动态的转变代码，这个想法的目的是用更多的时间去执行我们生成的代码。当新的代码从TB中生成以后， 将会被保存到一个cache中，因为很多相同的TB会被反复的进行操作，所以这样类似于内存的cache，能够提高使用效率。而 cache的刷新使用LRU算法。



编译器在执行器会从源代码中产生目标代码，像GCC这种编译器，它为了产生像函数调用目标代码会产生一些特殊的汇编目标代码，他们能够让编译器需要知道在调用函数。需要什么，以及函数调用以后需要返回什么，这些特殊的汇编代码产生过程就叫做函数的Prologue和Epilogue，这里就叫前端和后段吧。我在其他文章中也分析过汇编调用函数的过程，至于汇编里面函数调用过程中寄存器是如何变化的，在本文中就不再描述了。

函数的后端会恢复前端的状态，主要做下面2点：

1. 恢复堆栈的指针，包括栈顶和基地址。
2. 修改cs和ip，程序回到之前的前端记录点。

TCG就如编译器一样可以产生目标代码，代码会保存在缓冲区中，当进入前端和后端的时候就会将TCG生成的缓冲代码插入到目标代码中。

接下来我们就来看下如何翻译代码的：

客户机代码

SDN(8)
Windows(1)
虚拟化(20)

随笔档案

2016年12月(1)
2016年6月(4)
2016年5月(2)
2016年3月(2)
2016年1月(7)
2015年12月(15)
2015年11月(1)
2015年10月(12)
2015年9月(7)
2015年8月(21)
2015年7月(7)
2015年6月(9)
2015年5月(1)
2015年4月(1)
2015年3月(1)
2015年1月(1)
2014年11月(2)
2014年3月(1)
2014年2月(1)
2013年12月(1)

最新评论

1. Re:OpenFlow Switch学习笔记 (五)——Group Table、Meter Table 及Counters

我想问一下你关于组表的问题，我现在想用带有ovs的四台主机实现负载均衡，所以我想给ovs交换机下组表，结果总是不成功，我想请问一下，组表到底是怎么下的，我是用的手动配置的，我配置进去，但是好像丝毫不起...

--wage

2. Re:UnicodeDecodeError: 'ascii' codec can't decode byte 0xe9 in position 0: ordinal not in range(128) 解决办法

改了后就报错：
UnicodeDecodeError: 'utf8' codec can't decode byte 0xb9 in position 1: invalid start byte...

--梁鸢毓@

3. Re:整合Open vSwitch与DNSmasq为虚拟机提供DHCP功能

@ 几米憧憬在宿主机上做下SNAT就可以...

--CasonChan

4. Re:Ubuntu14.04安装配置Open vSwitch

你好，我的open vSwitch是安装在Ubuntu的物理主机上面，但是我一旦将笔记本本地有线端口加入到br0网桥中就不能访问外网了。按你说的dpclient br0后，我的电脑可以通过br0上网，...

--几米憧憬

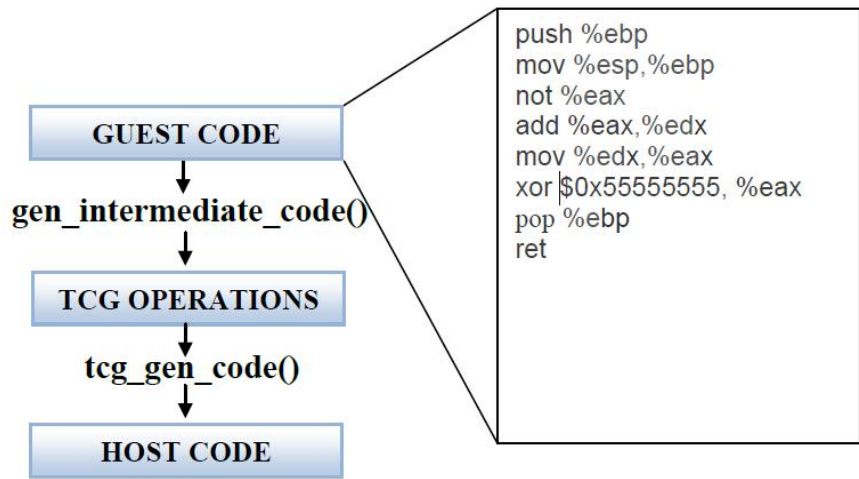
5. Re:整合Open vSwitch与DNSmasq为虚拟机提供DHCP功能

博主你好，在你这个实验中你能实现虚拟机VM1通过br0桥接通过电脑本地接口eth0上外网吗？

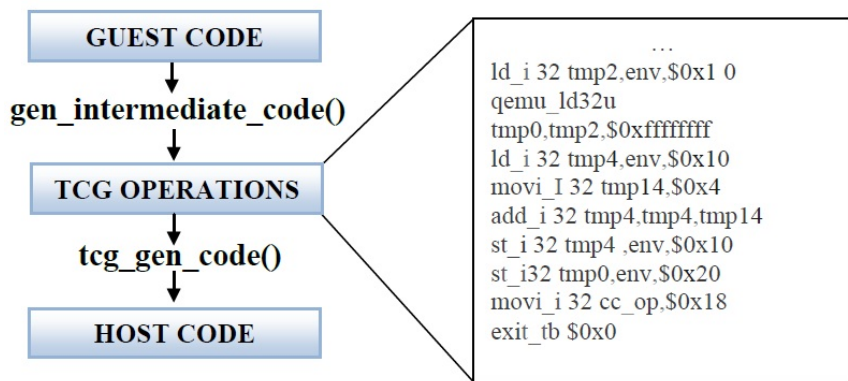
--几米憧憬

阅读排行榜

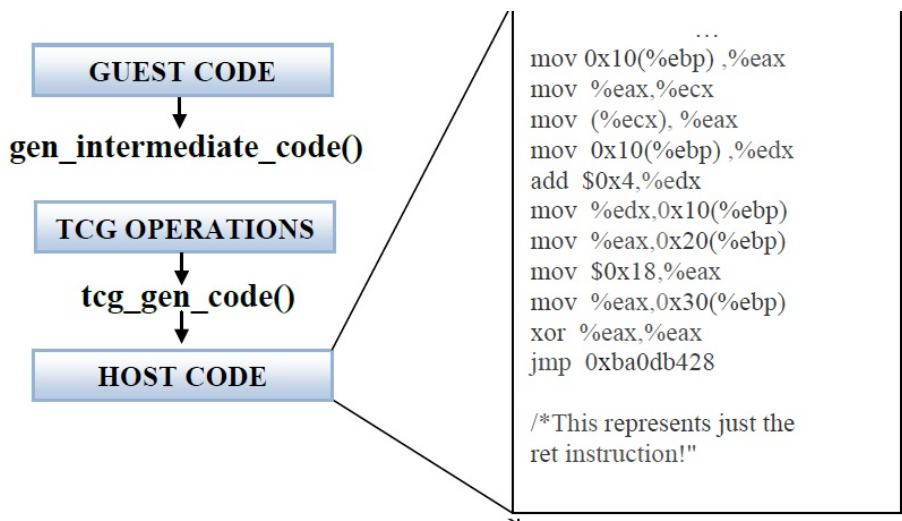
1. UnicodeDecodeError: 'ascii' codec can't decode byte 0xe9 in position 0: ordinal not in range(128) 解决办法(55787)



TCG中间代码



主机代码



3. TB链

在QEMU中，从代码cache到静态代码再回到代码cache，这个过程比较耗时，所以在QEMU中涉及了一个TB链将所有TB连在一起，可以让一个TB执行完以后直接跳到下一个TB，而不用每次都返回到静态代码部分。具体过程如下图：

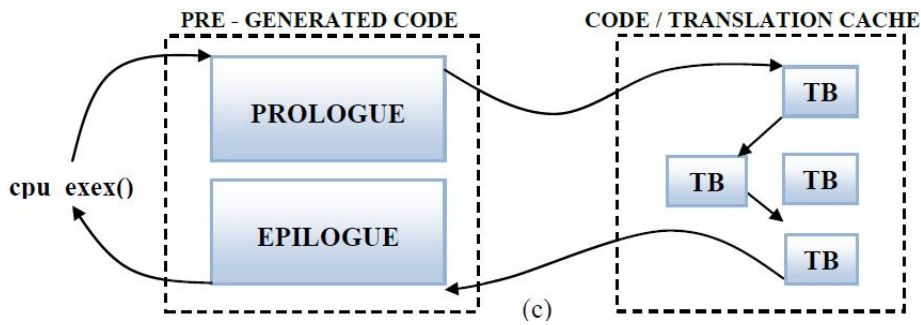
2. 利用 ipset 封禁大量 IP(16054)
3. OpenFlow Switch学习笔记(五) ——Group Table、Meter Table及Counters(8639)
4. MBR主引导扇区解析(8533)
5. OpenStack虚拟机DHCP获取不到IP地址排查(8410)

评论排行榜

1. 整合Open vSwitch与DNSmasq为虚拟机提供DHCP功能(4)
2. Ubuntu14.04安装配置Open vSwitch(2)
3. OpenFlow Switch学习笔记(五) ——Group Table、Meter Table及Counters(1)
4. UnicodeDecodeError: 'ascii' codec can't decode byte 0xe9 in position 0: ordinal not in range(128) 解决办法(1)

推荐排行榜

1. 利用 ipset 封禁大量 IP(4)
2. MBR主引导扇区解析(3)
3. 整合Open vSwitch与DNSmasq为虚拟机提供DHCP功能(3)
4. 在KCloud上轻松“玩转”Docker(2)
5. Xen虚拟机磁盘镜像模板制作(三) ——CentOS 7(2)



4. QEMU的TCG代码分析

接下来看看QEMU代码中到底怎么来执行这个TCG的，看看它是如何生成主机代码的。

main_loop(...){/vl.c} :

函数main_loop 初始化qemu_main_loop_start()然后进入无限循环cpu_exec_all()，这个是QEMU的一个主要循环，在里面会不断的判断一些条件，如虚拟机的关机断电之类的。

qemu_main_loop_start(...){/cpus.c} :

函数设置系统变量 qemu_system_ready = 1并且重启所有的线程并且等待一个条件变量。

cpu_exec_all(...){/cpus.c} :

它是cpu循环，QEMU能够启动256个cpu核，但是这些核将会分时运行，然后执行qemu_cpu_exec()。

struct CPUState{/target-xyz/cpu.h} :

它是CPU状态结构体，关于cpu的各种状态，不同架构下面还有不同。

cpu_exec(...){/cpu-exec.c}:

这个函数是主要的执行循环，这里第一次翻译之前说道德TB，TB被初始化为(TranslationBlock *tb)，然后不停的执行异常处理。其中嵌套了两个无限循环 find tb_find_fast() 和 tcg_qemu_tb_exec()。

cantb_find_fast()为客户机初始化查询下一个TB，并且生成主机代码。

tcg_qemu_tb_exec()执行生成的主机代码

struct TranslationBlock {/exec-all.h}:

结构体TranslationBlock包含下面的成员：PC, CS_BASE, Flags（表明TB），tc_ptr（指向这个TB翻译代码的指针），tb_next_offset[2], tb_jump_offset[2]（接下去的Tb），*jump_next[2], *jump_first（之前的TB）。

tb_find_fast(...){/cpu-exec.c} :

函数通过调用获得程序指针计数器，然后传到一个哈希函数从 tb_jump_cache[]（一个哈希表）得到TB的所以，所以使用tb_jump_cache可以找到下一个TB。如果没有找到下一个TB，则使用tb_find_slow。

tb_find_slow(...){/cpu-exec.c}:

这个是在快速查找失败以后试图去访问物理内存，寻找TB。

tb_gen_code(...){/exec.c}:

开始分配一个新的TB，TB的PC是刚刚从CPUState里面通过using get_page_addr_code()找到的

```
phys_pc = get_page_addr_code(env, pc);
```

```
tb = tb_alloc(pc);
```

ph当调用cpu_gen_code() 以后，接着会调用tb_link_page()，它将增加一个新的TB，并且指向它的物理页表。

cpu_gen_code(...){translate-all.c}:

函数初始化真正的代码生成，在这个函数里面有下面的函数调用：

```
gen_intermediate_code(){/target-arch/translate.c}->gen_intermediate_code_internal()  
{/target-arch/translate.c }->disas_insn(){/target-arch/translate.c}
```

disas_insn(){/target-arch/translate.c}

函数disas_insn() 真正的实现将客户机代码翻译成TCG代码，它通过一长串的switch case，将不同的指令做不同的翻译，最后调用tcg_gen_code。

tcg_gen_code(...){/tcg/tcg.c}:

这个函数将TCG的代码转化成主机代码，这个就不细细说明了，和前面类似。

#define tcg_qemu_tb_exec(...){/tcg/tcg.g}:

通过上面的步骤，当TB生成以后就通过这个函数进行执行。

```
next_tb = tcg_qemu_tb_exec(tc_ptr) :  
  
extern uint8_t code_gen_prologue[];  
  
#define tcg_qemu_tb_exec(tb_ptr) ((long REGPARAM(*)(void *))  
code_gen_prologue)(tb_ptr)
```

通过上面的步骤我们就解析了QEMU是如何将客户机代码翻译成主机代码的，了解了TCG的工作原理。接下来看看QEMU与KVM是怎么联系的。

5. QEMU中的IOCTL

在QEMU-KVM中，用户空间的QEMU是通过IOCTL与内核空间的KVM模块进行通讯的。

1. 创建KVM

在/vl.c中通过kvm_init()将会创建各种KVM的结构体变量，并且通过IOCTL与已经初始化好的KVM模块进行通讯，创建虚拟机。然后创建VCPU，等等。

2. KVM_RUN

这个IOCTL是使用最频繁的，整个KVM运行就不停在执行这个IOCTL，当KVM需要QEMU处理一些指令和IO等等的时候就会退出通过这个IOCTL退回到QEMU进行处理，不然就会一直在KVM中执行。

它的初始化过程：

vl.c中调用machine->init初始化硬件设备接着调用pc_init_pci，然后再调用pc_init1。

接着通过下面的调用初始化KVM的主循环，以及CPU循环。在CPU循环的过程中不断的执行KVM_RUN与KVM进行交互。

```
pc_init1->pc_cpus_init->pc_new_cpu->cpu_x86_init->qemu_init_vcpu->kvm_init_vcpu-  
>ap_main_loop->kvm_main_loop_cpu->kvm_cpu_exec->kvm_run
```

3.KVM_IRQ_LINE

这个IOCTL和KVM_RUN是不同步的，它也是个频率非常高的调用，它就是一般中断设备的中断注入入口。当设备有中断就通过这个IOCTL最终 调用KVM里面的kvm_set_irq将中断注入到虚拟的中断控制器。在kvm中会进一步判断属于什么中断类型，然后在合适的时机写入vmcs。当然在 KVM_RUN中会不断的同步虚拟中断控制器，来获取需要注入的中断，这些中断包括QEMU和KVM本身的，并在重新进入客户机之前注入中断。

总结：通过这篇文章能够大概的了解QEMU的代码结构，其中主要包括TCG翻译代码的过程以及QEMU和KVM的交互过程。

<http://blog.chinaunix.net/uid-26941022-id-3510672.html>

No pains, no gains

分类: [虚拟化](#), [Qemu](#)

好文要顶

关注我

收藏该文



CasonChan

关注 - 4

粉丝 - 47

+ 加关注

« 上一篇: [Generic Netlink详解](#)

» 下一篇: [关于KVM的几篇细节文档](#)

posted on 2015-10-13 17:59 [CasonChan](#) 阅读(2354) 评论(0) [编辑](#) [收藏](#)

[刷新评论](#) [刷新页面](#) [返回顶部](#)

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问](#) 网站首页。

【推荐】超50万行VC++源码：大型组态工控、电力仿真CAD与GIS源码库

【推荐】腾讯云产品限时秒杀，爆款1核2G云服务器99元/年！

相关博文:

- [qemu kvm 虚拟化](#)
- [Qemu 简述](#)
- [\[原\] KVM 虚拟化原理探究 \(2\) — QEMU启动过程](#)
- [Qemu创建KVM虚拟机内存初始化流程](#)
- [QEMU,KVM及QEMU-KVM介绍](#)
- » [更多推荐...](#)

最新 IT 新闻:

- [“5G疫情谣言”肆虐社交网络 英国1个月发生数十起基站袭击事件](#)
- [疯狂的《动森》：游戏已下架，“黑产”依旧火](#)
- [短视频借直播电商“造反”，抖音、快手挑战阿里胜算几何？](#)
- [麦当劳4月15日举行“5G”新品云发布会 智能终端要来了？](#)
- [OYO创始人曾借款购买20亿美元股票 孙正义向银行做个人担保](#)
- » [更多新闻...](#)