

ExtFUSE中lookup优化的大致过程的内核代码梳理：

Extfuse的ebpf框架分为4个部分：

- linux内核增加注册extfuse的ebpf的程序类型；
- 在FUSE的内核中增加ebfp的挂载点以及相应的钩子函数，并增加辅助函数；
- 设计相应的ebpf挂载函数：HANDLER(FUSE_xxx)；
- 在用户态文件系统建立与内核共同使用的ebpf map，并在相关的元数据操作中维护ebpf map。

1、增加extfuse ebpf程序类

linux内核的地址<https://github.com/extfuse/linux>

在bpf_types.h中定义BPF_PROG_TYPE(BPF_PROG_TYPE_EXTFUSE, extfuse) 类型

在bpf.h中增加bpf_prog_type, BPF_PROG_TYPE_EXTFUSE

在bpf_load.c中增加对程序类型extfuse的内容加载，extfuse的程序加载使用prog array的索引方式。

```
1 struct bpf_map_def SEC("maps") handlers = {
2     .type = BPF_MAP_TYPE_PROG_ARRAY,
3     .key_size = sizeof(u32),
4     .value_size = sizeof(u32),
5     .max_entries = FUSE_OPS_COUNT << 1,
6 };
```

通过 BPF_MAP_TYPE_PROG_ARRAY，可以将多个 eBPF 程序存储在一个数组中，每个程序都有一个唯一的索引。使用 BPF_MAP_TYPE_PROG_ARRAY 类型的 BPF Map，可以轻松管理和访问多个 eBPF 程序。

加载的规则如下，其中根据sec("extfuse")来搜索这个类型的ebpf程序，section字符串为"extfuse\1"就加载到索引为1的ebpf程序中，section字符串为"extfuse\2"加载到索引为2的ebpf程序中。在extfuse.c文件中说明了这种定义的方式。

```
1 #define HANDLER(F) SEC("extfuse/"__stringify(F)) int bpf_func_##F
```

int SEC("extfuse") fuse_xdp_main_handler(void *ctx) 加载的主ebpf程序

```
1 int SEC("extfuse") fuse_xdp_main_handler(void *ctx)
2 {
3     struct extfuse_req *args = (struct extfuse_req *)ctx;
4     int opcode = (int)args->in.h.opcode;
5     PRINTK("Opcode %d\n", opcode);
6     bpf_tail_call(ctx, &handlers, opcode);
7     return UPCALL;
8 }
```

HANDLER(FUSE_LOOKUP)(void *ctx) 加载的分支ebpf程序，其中#define FUSE_LOOKUP 1，其中 fuse_xdp_main_handler通过extfuse的section中的数字部分op来确定加载那个ebpf的分支程序。

2、FUSE内核中挂载点和辅助函数

```
1 static void process_init_reply(struct fuse_conn *fc, struct fuse_req *req)
   //处理FUSE文件系统的初始化回复
2 {
3     ...
4     if (arg->flags & FUSE_FS_EXTFUSE)    //FUSE_FS_EXTFUSE标志确定是否加载eBPF程
   序
5         extfuse_load_prog(fc, arg->extfuse_prog_fd); //加载ebpf
6     ...
7 }
8 int extfuse_request_send(struct fuse_conn *fc, struct fuse_req *req)    //发
   送扩展的FUSE请求
9 {
10     ...
11     if (data) {
12         struct extfuse_req ereq;
13         fuse_to_extfuse_req(req, &ereq);    //将fuse_req结构体转换为
   extfuse_req结构体
14         ret = extfuse_run_prog(data->prog, &ereq); //执行ebpf程序并执行
15         if (ret != -ENOSYS) {
16             extfuse_to_fuse_req(&ereq, req);
17             req->out.h.error = (int)ret;
18             ret = 0;
19         }
20     }
21     ...
22 }
```

建立程序需要运行的ebpf调用过程中的辅助函数

```
1 BPF_CALL_4(bpf_extfuse_read_args, void *, src, u32, type, void *, dst,
   size_t, size)
2 BPF_CALL_4(bpf_extfuse_write_args, void *, dst, u32, type, const void *,
   src, u32, size)
3 //用于在BPF程序中通过函数ID查找相应的函数原型
4 static const struct bpf_func_proto *
5 bpf_extfuse_func_proto(enum bpf_func_id func_id, const struct bpf_prog
   *prog)
6 {
7     switch (func_id) {
8     case BPF_FUNC_extfuse_read_args:
9         return &bpf_extfuse_read_args_proto;
10    case BPF_FUNC_extfuse_write_args:
11        return &bpf_extfuse_write_args_proto;
12    case BPF_FUNC_map_lookup_elem:
13        return &bpf_map_lookup_elem_proto;
14    case BPF_FUNC_map_update_elem:
15        return &bpf_map_update_elem_proto;
16    case BPF_FUNC_map_delete_elem:
17        return &bpf_map_delete_elem_proto;
18    case BPF_FUNC_tail_call:
19        return &bpf_tail_call_proto;
20    case BPF_FUNC_trace_printk:
21        return bpf_get_trace_printk_proto();
22    default:
```

```

23     return NULL;
24 }
25 }

```

向ebpf注册相应的辅助操作，用于存储与BPF验证器相关的操作函数

```

1  const struct bpf_verifier_ops extfuse_verifier_ops = {
2      .get_func_proto = bpf_extfuse_func_proto,    //用于获取BPF函数的函数原型
3      .is_valid_access = bpf_extfuse_is_valid_access, //用于判断BPF程序是否可以访问
           指定的资源
4  };

```

3、ebpf用户态程序

定义inode cache，具体用于元数据的查询、插入操作中

```

1  struct bpf_map_def SEC("maps") entry_map = {
2      .type = BPF_MAP_TYPE_HASH, // simple hash list
3      .key_size = sizeof(lookup_entry_key_t),
4      .value_size = sizeof(lookup_entry_val_t),
5      .max_entries = MAX_ENTRIES,
6      .map_flags = BPF_F_NO_PREALLOC,
7  };

```

1. `lookup_entry_key_t` 结构体：

- `uint64_t nodeid`：父节点的 ID。
- `char name[NAME_MAX]`：节点名称。

2. `lookup_entry_val_t` 结构体：

- `uint32_t stale`：标记该条目是否过时。
- `uint64_t nlookup`：引用计数。
- `uint64_t nodeid`：子节点的 ID。
- `uint64_t generation`：节点的代数。
- `uint64_t entry_valid`：条目的有效性。
- `uint32_t entry_valid_nsec`：条目有效性的纳秒部分。

定义inode attr cache，具体用于元数据信息的更新、删除、判断是否有效等操作中。

```

1  struct bpf_map_def SEC("maps") attr_map = {
2      .type = BPF_MAP_TYPE_HASH, // simple hash list
3      .key_size = sizeof(lookup_attr_key_t),
4      .value_size = sizeof(lookup_attr_val_t),
5      .max_entries = MAX_ENTRIES,
6      .map_flags = BPF_F_NO_PREALLOC,
7  };

```

1. `lookup_attr_key_t` 结构体：

- `uint64_t nodeid`：节点的 ID。

2. `lookup_attr_val_t` 结构体：

- `uint32_t stale`：标记该属性是否过时。
- `struct fuse_attr_out out`：结构体包含节点属性的详细信息。

定义 ebpf钩子函数

```

1  HANDLER(FUSE_LOOKUP)(void *ctx)
2  {
3      ...
4      lookup_entry_val_t *entry = bpf_map_lookup_elem(&entry_map, &key); //在
map中查找inode的entry条目, 找到的话, 更新输出结果, 直接返回
5      if (!entry || entry->stale) {
6          if (entry && entry->stale) //entry条目过期
7              PRINTK("LOOKUP: STALE key name: %s nodeid:
0x%llx\n", key.name, key.nodeid);
8          else //不存在entry条目
9              PRINTK("LOOKUP: No entry for node %s\n", key.name);
10         return UPCALL;
11     }
12     PRINTK("LOOKUP(0x%llx, %s): nlookup %lld\n",
13         key.nodeid, key.name, entry->nlookup);
14     /* prepare output */
15     struct fuse_entry_out out;
16     uint64_t nodeid = entry->nodeid;
17     /* negative entries have no attr */
18     if (!nodeid) {
19         create_lookup_entry(&out, entry, NULL);
20     } else {
21         lookup_attr_val_t *attr = bpf_map_lookup_elem(&attr_map, &nodeid);
//查看inodeid是否在map中
22         if (!attr || attr->stale) {
23             if (attr && attr->stale)
24                 PRINTK("LOOKUP: STALE attr for node: 0x%llx\n", nodeid);
25             else {
26                 PRINTK("LOOKUP: No attr for node 0x%llx\n", nodeid);
27                 return UPCALL;
28             }
29         }
30         PRINTK("LOOKUP nodeid 0x%llx attr ino: 0x%llx\n",
31             entry->nodeid, attr->out.attr.ino);
32         create_lookup_entry(&out, entry, &attr->out);
33     }
34     //使用bpf的辅助函数把参数写到返回请求
35     /* populate output */
36     ret = bpf_extfuse_write_args(ctx, OUT_PARAM_0, &out, sizeof(out));
37     if (ret) {
38         PRINTK("LOOKUP: Failed to write param 0: %d!\n", ret);
39         return UPCALL; //不满足条件, 继续后面的请求
40     }
41     //完成请求返回
42     return RETURN;
43 }

```

extfuse.c中还包括HANDLER(FUSE_GETATTR)、HANDLER(FUSE_READ)、HANDLER(FUSE_WRITE)、HANDLER(FUSE_SETATTR)、HANDLER(FUSE_GETXATTR)、HANDLER(FUSE_FLUSH)、HANDLER(FUSE_RENAME)、HANDLER(FUSE_RMDIR)、HANDLER(FUSE_UNLINK)这几个操作的函数框架。

其中还额外包括一个函数：

```

1  static int remove(void *ctx, int param, char *op, lookup_entry_key_t *key)

```

用于删除引用计数为0的文件元数据。

4、在用户态文件系统中更新map

在文件系统初始化时加载并初始化ebpf程序

```
1 static void stackfs_ll_init(void *userdata, struct fuse_conn_info *conn)
2     lo->ebpf_ctxt = ebpf_init("/tmp/extfuse.o");
```

插入inode到eBPF的map中，并更新缓存

```
1 int create_entry(fuse_req_t req, struct lo_inode* pinode, const char *name,
2     const char *cpath, struct fuse_entry_param *e, const char *op, int time)
3     ...
4     #ifdef ENABLE_EXTFUSE_LOOKUP
5         res = lookup_fetch(get_lo_data(req)->ebpf_ctxt, inode->pino, inode-
6             >name);
7         if (res && !errno) {
8             INFO("[%d] \t Fetched %s nlookup: %ld inode->nlookup: %ld\n",
9                 getpid(), res < 0 ? "stale" : "", res, inode->nlookup);
10            if (res < 0)
11                res = -res;
12            inode->nlookup = res + 1;
13        }
14        //插入inode到map中去
15        res = lookup_insert(get_lo_data(req)->ebpf_ctxt, inode->pino,
16            inode->name, inode->nlookup, e);
17        if (res)
18            ERROR("[%d] \t %s new node %s id: %p: %s\n",
19                getpid(), op, cpath, inode, strerror(errno));
20    #endif
21    ...
22    int lookup_insert(ebpf_context_t *ctxt, uint64_t parent_ino,
23        const char *name, uint64_t nlookup, struct fuse_entry_param *e)
24        INFO("[%d] \t Inserting node name %s (%ju) parent 0x%lx nlookup; %lu\n",
25            getpid(), name, strlen(name), parent_ino, nlookup);
26        // update entry, 调用ebpf map接口更新inode cache map, 0表示inode cache map,
27        1 表示attr cache map
28        int overwrite = 1; //XXX overwiting to update any negative entires
29        ret = ebpf_data_update(ctxt, (void *)&key, (void *)&entry, 0,
30            overwrite);
31        if (ret)
32            ERROR("[%d] \t Failed to insert %s (%ju) parent 0x%lx count %ju:
33                %s\n",
34                getpid(), name, strlen(name), parent_ino, num_entries, strerror(errno));
```

其中涉及了一些操作inode的函数，create_entry, delete_node, forget_inode, stackfs_ll_rename。