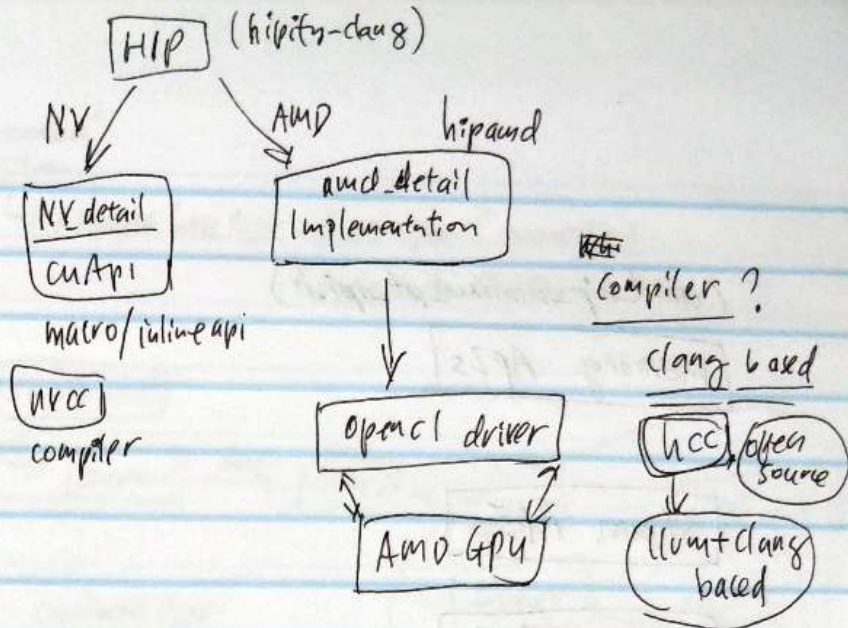## HIP

- ▷ device Context management
- ▷ memory management (Host / Device)
- ▷ Kernel Execution/launch.

- ◉ Command Queue
- ◉ (buffer) (memory)
- ● Kernel

[Allstatic]
↑
[MemObjMap]
     map    $Void^* \Rightarrow amd::Memory^*$.

amd::Monitor
     (类似 mutex)
       ↳ runtime/thread
       runtime/platform.
       runtime/os
       runtime/utils

[HIP] (hipify-clang)
   NV ↙    ↘ AMD    hipand

[NV detail / cuApi]
macro/inline api

(nvcc)
compiler

[amd_detail Implementation]
↓
[opencl driver] ⇄
[Amd GPU]

compiler?

clang based
[hcc] (open source)
↓
(llvm+clang based)

Allstatic   .ctor .dtor → 不能调用

EmbeddedObject → (new ddete) 不能调用

Stackobject → (new ddete) 不能调用

MemoryPoolObject ⇒ (new ddete) 相较调用
         → 不能调用 .dtor
         ddete

HeapObject → 显式调用 new/ddete

ReferenceCountedObject
       ↳ 显式调用 new/ddete.
       ↳ refcnt ...

( amd_hip_runtime_pt_api.h )

$\boxed{\text{Memory APIs}}$

$\boxed{\text{Stream APIs}}$

$\boxed{\text{event APIs}}$

$\boxed{\text{Launch APIs}}$

$\boxed{\text{stream}}$   ihipStream_t $\overset{\text{typedef}}{\underset{*}{\rightleftharpoons}}$ hipStream_t

$\boxed{\text{stream} = \boxed{0} \boxed{\text{Null}}}$  比 per-thread default stream 简洁了. 会被卷生阶段选择不同的

oueet 运行事.

在运行事中, 判断 == Null ⇒ getPerThreadDefaultStream

(指毛针) : 0 → 2 ⟶ hipStreamPerThread

↳ 址台新的为 hip.tls.Stream_perthread_d

每个线程 tls   $\overset{\text{class}}{\boxed{\text{stream_per_thread}}}$
对象

hipstream⟳ ⟳ ⟳ ⟳ ⟳ hipstream_t   gpu list

$\boxed{\text{开销大}}$   $\boxed{\text{含有个独立的线程}}$ (circled in pink)

hip::Stream   ⟶ ⟶ hipstream_t   ⟶ $\boxed{\text{ihip Stream_t}}$   ↑   $\boxed{\text{amd::HostQueue}}$

$\overline{\text{Stream的synchronize:}}$
Call  HostQueue的finish api   ⟶ $\boxed{\text{hip::Stream}}$
   → priority
   → HostQueue*
   → Device
   → flags

stream : $\boxed{\text{Create}}$  会又拉合创建个

每个device含有组多Queue, 添加入Device
进行管理

底层会含有一个 StreamSet (保在指针), 保存所有 device 的所有streams

$\boxed{\text{可以到建所有streams}}$

hipstream Query

and::HostQueue* → 获取 HostQueue.

command*        command*

→ Abstract base type of all openCL operations
(submit to HostQueue)

RuntimeObject

Event

Encapsulate the status of a command →

Command ← → HostQueue

next → command → next → command →

UserCommand   Marker   ....

(Api)   HostQueue

Event

(hipEvent_t)   → Event

command

command

Command → Command type

command

command → EventWaitList (dependencies)

给句约定的 eventSet
维护户的临时的event

StreamAddCallback :

setCallback (完成时候被调起来)

command ← command ← Marker ← blockMarker

                ?

last command

(维护)
Event 的 callback list (单链表)

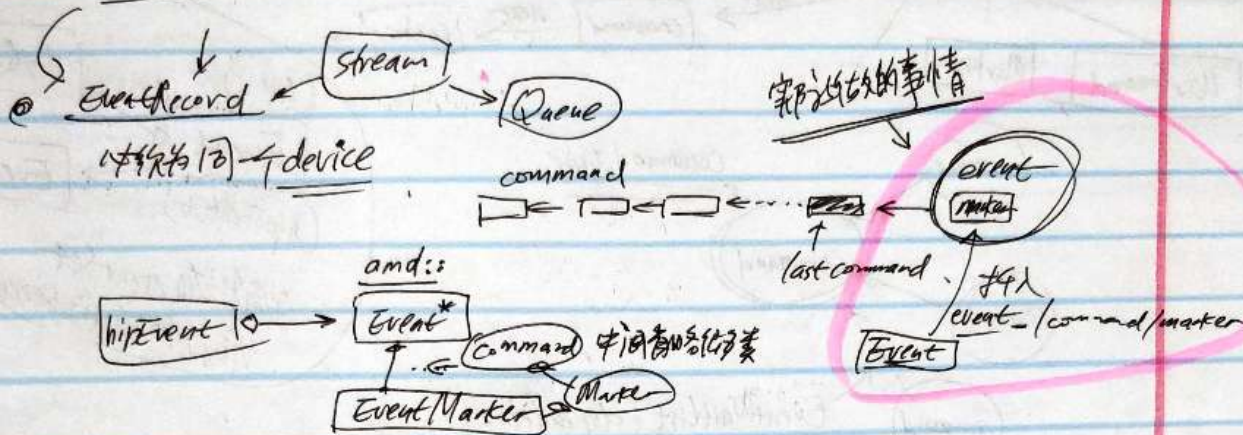atomic 类型 :

AddCallback : 首先创建一个 CallbackEntry*, 里边包含了个 callback function / userdata.

entry→next = freelist ;(head)
while (! freelist . compare-exchange-weak(entry→next, entry))
;

(怎么解决??)

ihip stream_t

hipstream_t

hip::stream*

typedef

reinterpret_cast ??

(内部)

StreamWaitEvent

stream 等待一个event 完成
所以在stream中插入一个 Marker, depend on 那个 event

到host stream queue.

Marker

Event

command

stream

**hip IPCEvent**

hip Event

@ Event Record ← stream → Queue

必须为同一个 device

command

郑这边的事情

event
marker

□ ← □ □ ← □ ← □ ← ---- □■■ ← ⊙ event marker

↑
last command

插入
event_/command/marker

Event

amd::

hipEvent ◁── → Event*
EventMarker ── Command 中间都略去了
EventMarker ── Marker

@ synchronize    record 这边其实在 synchronize 了.    调用 event_→await(completion.
                 因为在 stream 中插入了一个 marker    ↳ condition variable wait
                                                      (monitor)      (mutex)
                                                      等 status 变成
                                                                complete.

o (query)    在请求状态时,手动插入一个 marker, 告诉 queue 有人在等
             和等 complete

                 多次请求, 不再插入, 用 atomic_flag (! test_and_set)

hipEvent 中的数据域 event_ 也是在这个时候创建.
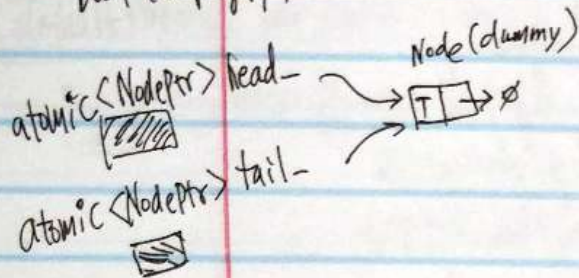
`Concurrent LinkedQueue` 基于一篇论文 `lockfree`

(non-blocking)
( un-bounded )

boost::lockfree
也有类似的实现

enqueue(T elem)

T dequeue();

bool empty();

Node


atomic<NodePtr> head_


Node(dummy)


atomic<NodePtr> tail_


$N = 5$

NodePtr 是个 tagged pointer

TaggedPointerHelper*

低5位 放入了 一个 tag

分配的指针地址是 align到 $2^5$ 的

$2^5 - 1 = mask$

所以 低5位为0，可用来 tag 其它东西。

根据 mask 可快速得到
原始指针 和 tag .



head_

tail_

next NULL

new node

```
┌─────────────────┐
│ RuntimeObject   │
└─────────────────┘
        △
┌──────────────────┐
│ CommandQueue     │ (abstract class)
└──────────────────┘
   △            △
┌───────────┐   ┌──────────────┐
│ HostQueue │   │ DeviceQueue  │
└───────────┘   └──────────────┘
```

ConcurrentLinkedQueue
   <command*> queue_

有 Thread HostThread
       在运行 run 执行 command

Hush() 循环发 thread run

(stream)
synchronize
(finish)：创建个 marker / command
       enqueue到 queue, 等待执行。

Thread：  pthread_create ⇐ Thread::entry (Thread* thread)
                    ↳ 设置 sigmask
                    ↳ 调用 thread → Main
                        ├→ (通知调用线程 我线程已启动)
                        ├→ (准备一些参数)
                        └→ 调用 run 等待线程创建完发信号 (start AP 令信号 notify)
                                                            通知。
                        ↳ run 调用 (虚函数)            start 令信一些参数
                                                       (信号准备好)
HostQueue:: Thread    ⇐
    然run 的 实现中                                    run (void*) 接受参数

含有书 virtualDevice
   ↳ 调用 queue → Loop (virtualDevice):
                        → 不断从 queue 中 dequeue个 command
         ↓                 然后调用 command → submit (virtualDevice).
      HostQueue .
                        (虚接口)      (纯例子)
                        virtualDevice.submit Kernel (*kernelcmd);
```

[hipDeviceSynchronize] ——→ 是否创建了？ 会创建marker等待
├→ (NULL Stream) queue [finish] 所有default stream结束.
└→ syncNonBlocking streams. 当前这个device) (Flag)
                            (不是所有device)

[hipDeviceReset]
├→ destroyAllStreams (current device)
└→ Purge MemObj map. (host pinned memory device memory)

[NULL]
[default stream] per device. (legacy stream)

Device class 有个 stream 上 新来时候, 代表 NULL stream.
        当 underlying HostQueue初始化为 NULL
        每个线程(?) 自己的 per-thread default stream.
        传入∅(API), renamed-function 令转成 2,
        作为 per-thread default stream.

[per-thread default stream] ⇒ 是一个flag=default stream, 不是 non-blocking的flag.

[legacy stream] 在执行 command之前, 会sync 其它 active stream (default stream only) (flag)
                                                        (non-blocking 不受影响的) (同一device)



[hip::Device] 每个device 都有一个 Memory Pool default 在 create时创建, 也有一个 memory pool 列表

 • 一个 list of HostQueue*.

 • Device Id

 • amd::Context* context; (其实共享的) Shared by all devices
                    ∘—→ (devices list)

[hip::Device] ⇐⇒ [amd::Device]        (amd::Device*)
                                        (底层类)

[hip Stream_t] ⇐⇒ [amd::HostQueue]              HostQueue, 在销毁时都是有强制同步.
(stream)

                        stream | HostQueue

hip GetLastError

- 返回的是 HS 中的 last error
  不是 gpu device 上的 last error
  (与 cuda 入符合) ??

- 拿到后 置为 success

每个 HipApi 在最开始都会 check  hip runtime 是否 init
(public)

hip_prof_api

      HIP_INIT_API_INTERNAL

         └→ HIP_INIT

            └→ trace Api call (profiling) 取当前

构建stack上的 object ctor/d...
类似于 profile run time 的事情

amd hip 用了个技巧 针对所有的 public stream. event. ...
在 runtime.h : typedef struct ihipStream_t * hipStream_t .

但是查找不到 struct ihipStream_t 的定义    这里 hipStream_t 只是一个 指针
                                 甚至可以 简单的定义为 void*

但是在 实现文件时, 用了另一个struct 来实现 hipStream_t 的功能. 比如

struct myStream {

        ...

};            t ⇒ hipStream_t

*t   ← reinterpret_cast< hipStream_t> (new myStream) .

api  传入的 hipStream_t (其实为个指针, 指向一个不存在的结构体)
        api 处理时都将其强转为 <myStream *> (t) .
        myStream* s ← reinterpret_cast
             再对s进行操作.

| Memory pool | 面向 stream 的 memory pool (全局 类 lock) | amd::Memory |
|---|---|---|

device::Memory

→ Heap busy
→ Heap free
→ Allocate Memory ⇒ (双向) virtual memory pointer ^address
→ Free Memory ⇒ 输入 amd::Memory * (+ stream*)

Heap 管理
unordered_map (amd::memory*, memoryTimestamp)

API Free Memory
• 从 busy Heap 找 amd::Memory* (map查找)，删除, 返回 memoryTimestamp
• 创建个 event / marker   associate with  memory Timestamp ??
   (到/stream)
• 将 memory 加入到 free heap 进行 track (+ timestamp)

memory Timestamp 是什么东东 ??
维护 event 和 memory object 对应的 stream

当 memory 分配时 Heap 会 track/record amd::Memory* ⟶ {stream, NULlptr}
↳ 0 那时没有 stream 在用                                      (memory Timestamp)

Free 时会更新 (event).

Kernel 1 (.. stream (memory, );
FreeMemory (memory, stream1); ← 此时前 Kernel还没完成, 但是内部会将其转入

Free Heap 同时加入 event/marker
监控 memory 在 stream 上的使用

① 一旦 event 的 status 返回 ready 代表前面的 Kernel 已结束, FreeHeap 可安全回收
                                                    在 stream 的 Kernel (同)

② 而且 同个 stream 的 Kernel 也可以直接用, 不用等待 之前 Kernel 完成
因为 Kernel 是顺序执行的.

memory pool 在重用 memory 块时
并没有进行切割, 只要大于请求的
size, 而且 stream 这时可以用, 就能复用.
(reuse)

(hipFreeAsync) ⟵ void* devPtr ( virtual memory)
　　　　　　　⟵ hipstream_t stream.

⇓

○ 根据 devPtr 找到 memory object ( amd::Memory* )
　　　└ 所有的 memory*，全部保存在一个哈希对象中 map <amd::Memory* >  void*
　　　　　　　　　　　　　　　　　　　　　　map: 虚拟地址 ⟹ amd::Memory*.
　　　　　　　　　　　　　　　　　　　　　　　　○ 排序的 按地址大小
　　　　　　　　　　　　　　　　　　　○ upper-bound 找 第1个大于这个地址 (>)
　　　　　　　　　devPtr　　　　　○ 往前进一个　　　*
　　memory ↓　　　　　　　　　　　(lower-bound) 第1个不小于 (≥)
　　　*└──(size)──┘

○ 提供的 stream 是 NULL, 里就用 NULLstream 不是 default stream per thread
　　　　　　　　　(一个设备件)　　　　, 所以用 配合创建的 stream 来用 memory pool
　　　　　　　　　　　　　　　　　　　　　也为用 perthreadstream (==2) 里提供的用.

○ 从设备对象中 FreeMemory ( loop 所有的　)　调用 memory pool 的 FreeMemory
　　( memory object 中存有 device id )　memory pool　　　成功，就 return

(hipMallocAsync) ⟵ devPtr **
　　　　　　　　　⟵ size
　　　　　　　　　⟵ stream = 0 ? 用 NULLstreams (一个设备件)
　　　　　　　　　　　　　　　　　　　　　还引于 default stream per thread

↓
(Stream) ─找到→ device ⟹ 用当前的 memory pool ⟹ allocate memory
　　　　　　　　　　　　　　( memory pool list )

(hipMemPool Create) 显然创建 memory pool, 并加入 device 中的 memory pool list.　　　　　也要加 device 中创建
(hipMemPool Destroy): 释放所有 free 的 memory, 如果还有 memory 没 Free, 因为 memory pool refcnt >1, (会 delete)
　(所以最后释放所有的 memory 用这个)　(什么时候 free 呢?)　⟵　它(们) 还会被以为 keep 在 pool 中
　　　　　　　　　　　　　　destroy

$\boxed{\text{memory pool}}$

↳ 里面可以将给不同的stream分配的内存.

hipMemallocAsync 中路径 Stream

我可以从当前device的当前memory pool 中分配.
(stream in device)
↑
default memory pool
我者 能创建个, 然后 Set Current.

② 当前也可以 显式的从特定的 memory pool (配创建的/default的) 传入 API 来分配
hipMallocFromPoolAsync ( --- memory pool, Stream)

所以 hipFreeAsync 就需要从 device in memory pool list 查找对应的序号

---

$\boxed{\text{hipCtx\_t}} \rightleftharpoons \boxed{\text{hipDevice}}$    全局含有 vector<hip::Device*> g_devices.

对外可见的类的指针(封装的)    内部表示
handle.

$\boxed{\text{hip::Device}} \leftarrow\rightarrow \boxed{\text{amd::Context}} \xleftarrow{1:n}\rightarrow \boxed{\text{amd::Device}}$

hip_context.cpp 对 devices 进行了初始化. 每个api调用时都会check是否 inited

同时还有个 host device $\rightleftharpoons$ amd::Context

---

$\boxed{\text{amd::Device}}$ 类有个静态变量 vector<Device*> 代表所有 amd devices $\Big\{$ DeviceLoad
DeviceUnload

↑ 针对
gpu类型
调用上述主动(查找加载)

$\boxed{\text{Runtime::Init}} \longrightarrow$ Flag::Init()/ option::Init() Device::Init()
↳ enumerate all devices

每个api正启动被调用时,都会查找 Runtime::Init 是否被调.

↑
$\boxed{\text{HIP\_INIT}} \rightarrow \boxed{\text{hip::init()}} \rightarrow \boxed{\text{Runtime::Init}} \rightarrow$ 初始化 hip::Devices list (+ Context)
                                                              amd::

                                        创建

$\boxed{\text{hipInit}}$ 显式调用 也达到相同目的.

● Device::Create 会创建 default memory pool.

● 每个device都 associate 一个 amd::Context 指针

**hip Ctx Create** 创建个新的 对应 device id 的 hipCtx_t。

其实就是将 g-devices [gpuid] (hip::Device*) (ret cnt++) ⇒ push到当前线程S stack 的top.

(非新创建)

每个线程有1个 stack, 维护 context.

**hip Ctx Set Current** 将 ctx 作替换当前的线程的stack的top.

为 Nullptr, 弹出stack的 top.

**hip Ctx Push Current** 将ctx push到当前线程的top of stack, 其实就是 当前线程→当前device 变成这个 c

**hip Ctx Get Device** 返回就是当前的线程对应 device id.

---

**platform Stat** 管理 CodeObject (fatbinary)

- map < hip Module_t, hip::DynCO*>
- hipStatCO // static code object
- texture reference

CodeObject

statCO    DynCO

**statCO** 管理
- modules : map <void*, $\overset{const}{}$ Fat Binary Info*> ← __hipRegisterFatBinary
- functions : map<void*, $\overset{const}{}$ Function*> ← __hipRegisterFuncs
- vars : map<$\overset{const}{void}$, Var*> ← __hipRegisterVars
- managed Vars: vector<Var*> ← __hipRegisterManagedVar

functions : host function* ⟹ Function*

↳ device function code 指针??

[code] format

.clang Offload Bundle Header {
const char magic [SIZE];
uint64_t numOfCodeObjects;
_clang Offload Bundle Info desc[1];
};

magic  numofcodeObjects  [desc]
////////  8  |  ...  |
        n   0    1    2    n-1
                              offset
                            (image)

offset [ code ] ↑ [ code ] ↑
      offset        offset

⟶ { uint64_t offset;  ⟹ 相对于 code object binary 的开头
uint64_t size;

uint64_t bundleEntryIdSize;
const char bundleEntryId[1];
}
(String)

用该 id[0] + IdSize 即可loop到
下一个 desc block

globals : {
DeviceVar
DeviceFunc  ⟺ hipFunction_t 的内部表示.
vector      含有 amd::Kernel*.
Function::<DeviceFunc*> 每个device都有一个. 从 FatBinaryInfo 中 Extract 书章.
Var
}

hip_module.cpp 怎么如何启动 kernel  **hipLaunchKernel** ( const void* [hostFunction] , ....  , shm Bytes, stream)

__global__ void KernelFunc(---) { }

hipModule_t ⟺ ihipModule_t* (amd::Program)  不然是CL
hipFunction_t* ⟺ struct ihipModuleSymbol_t* (不存在这个结构)  由这个 host Function 指针找到
(!总代指针类型)                                              对应的 hipFunction_t 指针 func

hip::DeviceFunc::asFunction(f)
⤷ hip::DeviceFunc* (在kernel()返回) [amd::Kernel*]

ihipModuleLaunchKernel (func, ...... )
图 其它每个 kernel 若要以单独 用24各起 startEvent / stopEvent 去
(profile 它的时间的)

platformState.getStatFunc
⤷ 从 statCO 中获取
⟶ functions_map
map <void*, Function>
!!! deviceFunction 在这里才从 FatBinary中提取出来 从 code object

[调用]
⤷ ihipLaunchKernelCommand  生成一个 command ( 从 DeviceFunc )

(NDRangeKernelCommand) ⤷ if (startEvent) 生成一个 marker command, 放入队列.
⤷ command → enqueue(); // 放入队列.
⤷ if (stopEvent), track bind command event 以 stopEvent.

hip::

**DeviceVar**
(OpenCL type)

hiModule_t $\Rightarrow$ cl_program $\Rightarrow$ amd::Program ( $\hookrightarrow$ gpu device 相对应)

[External]　　　[Internal] (Runtime type)

$\begin{pmatrix} name, \\ amd::Memory^* \ memobj- \\ hipDeviceptr_t \ device\_ptr-, \\ size_t \end{pmatrix}$

cl_context $\rightleftharpoons$ Amd::Context

cl_event $\rightleftharpoons$ amd::Event

cl_command_queue $\rightleftharpoons$ amd::CommandQueue

cl_kernel $\rightleftharpoons$ amd::Kernel

cl_program $\rightleftharpoons$ amd::Program

cl_device $\rightleftharpoons$ amd::Device

cl_memory $\rightleftharpoons$ amd::Memory

cl_sampler $\rightleftharpoons$ amd::Sampler

$\underline{amd::Program^* \ program} = \underline{as\_amd} \ (reinterpret\_cast \ <cl\_program>(hmod))$ .

```
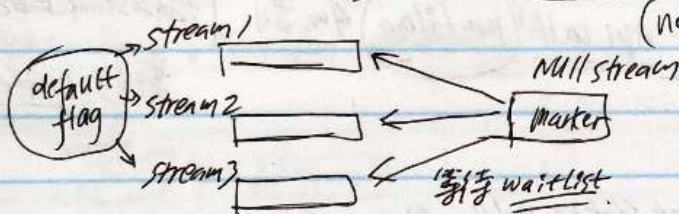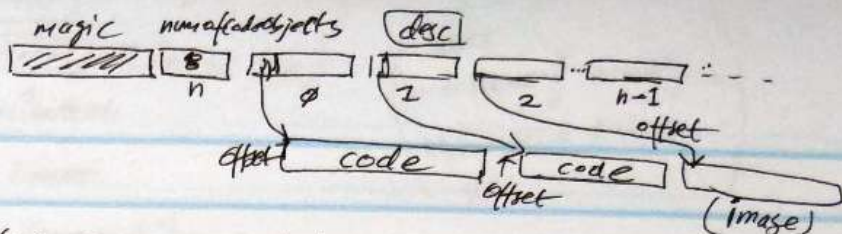template <typename cl>
typename amd::as_internal <cl>::type*      cl_program ⇒ amd::Program*
as_amd (cl* cl_obj)
{
    return amd::RuntimeObject::fromHandle (typename amd::as_internal <cl>::type>
                   (static_cast <void*>(cl_obj)) ;
}
```

( amd::Program )　map (device*, device::Program*)

( amd::Symbol )　map <device*, device::Kernel*>
　　 $\hookrightarrow$ device function
　　 (kernel).

hip::

**DeviceFun**

( name )

( amd::Kernel* )

对应的（symbol）

对应关系:

Var : - 每个 variable 相对于 device 上的 DeviceVar*

Function : - 每个 function 相对于 device 上的 DeviceFun*

构がまし

分层结构:

```
┌─────────┐     namespace
│   HIP   │      CAPI
└─────────┘
     │
     ▼
┌─────────┐     ┌─────┐
│ hip amd │◄───►│ hip │
└─────────┘     └─────┘
     │
     ▼
┌─────────────┐  ┌─────┐
│ amd open cL │◄►│ amd │
│  runtime    │  └─────┘
└─────────────┘
┌──────────────┐   ┌─────┐
│platform/device│◄►│ roc │
└──────────────┘   └─────┘
     │
     ▼
┌─────┐
│ hsa │
└─────┘
```

```
┌─────────┐
│ runtime │───────────────────┐
└─────────┘                   │
  │    │    │                 │
  ▼    ▼    ▼                 ▼
┌──────┐ ┌────┐ ┌─────────┐ ┌────────┐ ┌───────┐
│device│ │ os │ │ platform│ │ thread │ │ utils │
└──────┘ └────┘ └─────────┘ └────────┘ └───────┘
```

───────────────────────────────────────────

hip-fatbin  ┌─────────────┐  输入┌─ fname
            │FatBinaryInfo│◄──────┤
            └─────────────┘   ←── image ( ▨▨▨ memory )
                              ( 对对给新的仓个封装 )

对"所行读取地址 (建立
(program) 建立
        FatBinary DeviceInfo (program: ▨▨▨ )
                        一进制代码的一个 range (ptr+len).
        每个都建立作 program. (amd::Program)
                  └─ device program
           对抗
          (封装)


hip-clang  links device code from different translation unit together. For each device target, a
code object 's generated, dev1 ┌code┐                            embed                      ┌ELF┐
                          dav2 ├code┤ } clang-offload-bunder ┌─────────┐ ═══► global symbol. in (hip-fatbin
                               │ ⋮  │                        │fatbinary│                    `__hip_fatbin`  section
                               └────┘                        └─────────┘

Initialization code for each translation unit for host code compilation
━━━━━━━━━━━━━━
      └─► `__hipRegisterFatBinary` ═══► register the fatbinary embedded in the ELF file.
              └─► `__hipRegisterFunction` and `__hipRegisterVar`
                    ( kernel )              ( device global vars )

Termination
━━━━━━━━━━
      └─► `__hipUnregisterFatBinary`
```

extern "C"

$\underline{\text{\_hipRegisterFatBinary}}$

$\quad\hookrightarrow$ add FatBinary (data) $\Longrightarrow$ hip:: FatBinary Info**.

\_\_hipRegisterFunction $\leftarrow$ hip::FatBinaryInfo**

$\qquad\qquad\uparrow\nwarrow$ host Function (void*)
$\qquad\qquad$ device function (char*)

$\downarrow$

host function $\underset{\rightleftharpoons}{\overset{map}{}}$ hip::Functions

$\qquad\qquad\hookrightarrow$ list of DeviceFunc for all devices

$\qquad\qquad\qquad\hookrightarrow$ DeviceFunc* $\rightleftharpoons$ (hip Function_t) ☐ FatBinary code

$\qquad\qquad\qquad\hookrightarrow$ Kernel   (ModuleSymbol_t)

hsA: heterogeneous System Architecture.