

Tensorflow - 0.12

DirectSession Run

RunOptions
 namedTensorList (pair<string, Tensor>
 & output-names
 target-nodes
 Vector<Tensor> outputs

ExecutorsAndKeys

RunStateArgs

GetOrCreateExecutors

Cancellation Manager

RunStates

SendInputs

Parallel Executors

ExecutorBarrier* barrier = new ExecutorBarrier(
 num_executors, runstate.rendez, [&runstate](const Status& ret) {
 mutexlock l(runstate.mutex);
 runstate.status.Update(ret);
 runstate.executors.Done.notify();
 });

Executor::Args args;

args.step_id =

args.rendezvous = runstate.rendez;

args.cancellation_manager = &StepCancellationManager;

args.runner = [this, pool] (Executor::Args::Closure c) {

SchedClosure(pool, std::move(c));

};

for (const auto& item : executors_and_keys.items) {

Item.executor->RunAsync(args, barrier->Get());

};

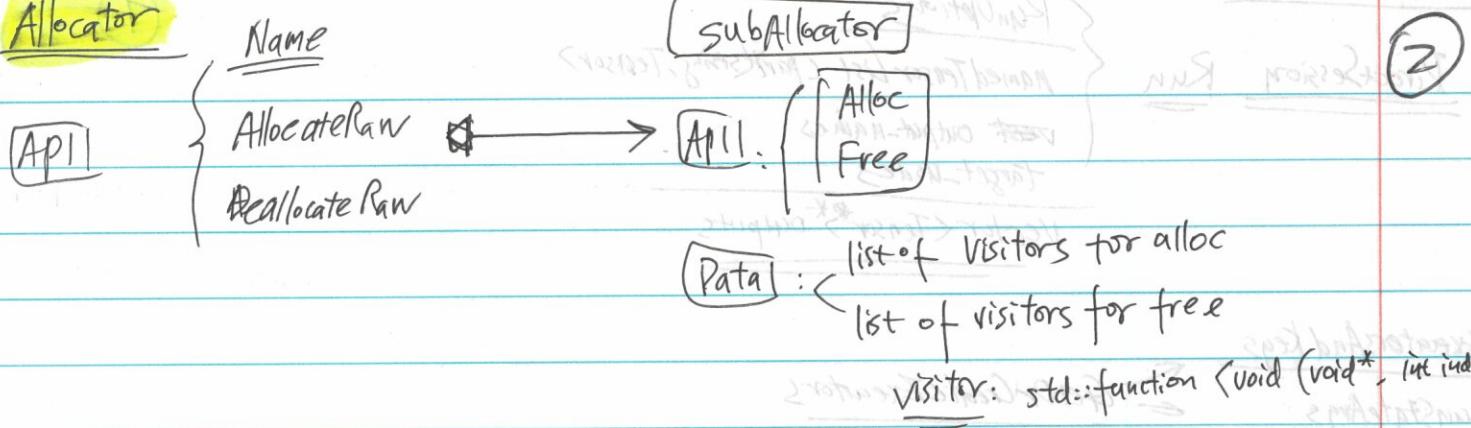
ret_millis

CBIVI

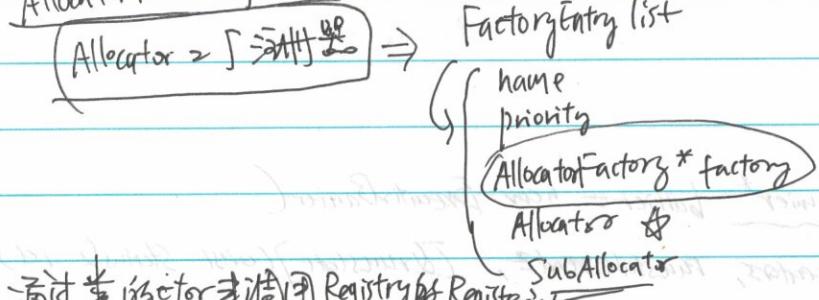
ret_millis

ret_millis

Allocator



AllocatorFactory Registry



通过类 `AllocatorFactory` 使用 `Registry` 的 `Register` 方法

`Registry` 通过提供 `singleton` 的 `factory` 在 `Registration` 构造函数中进行注册。这样在 `GetAllocator` 中第一次使用 `Allocator` 时会调用 `factory` 的 `CreateAllocator` 方法。

#define REGISTER_MEMORY_ALLOCATOR(ctr, file, line, name, priority, factor_class) |

static AllocatorFactoryRegistration allocator_factory_reg_##ctr (|
file, line, name, priority, new factory class)

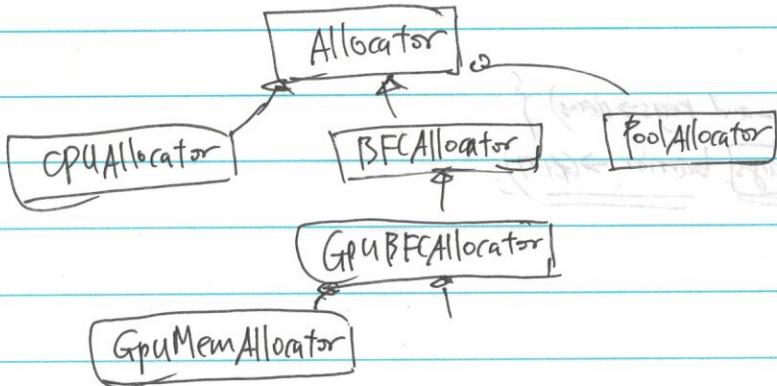
REGISTER_MEMORY_ALLOCATOR(name, priority, factory class)

entry

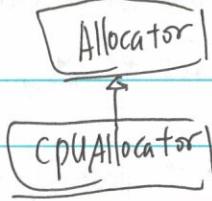
GetAllocator 在 registered list 中找一个最高优先级的 AllocatorFactory

通过 `factory` 调用 `CreateAllocator()`

插入 `entry` 到 `Allocator`, cache 逻辑



Cpu_allocatorImpl.cc



AllocateRaw : port::AlignedMalloc
DeallocateRaw : port::AlignedFree.

2/3

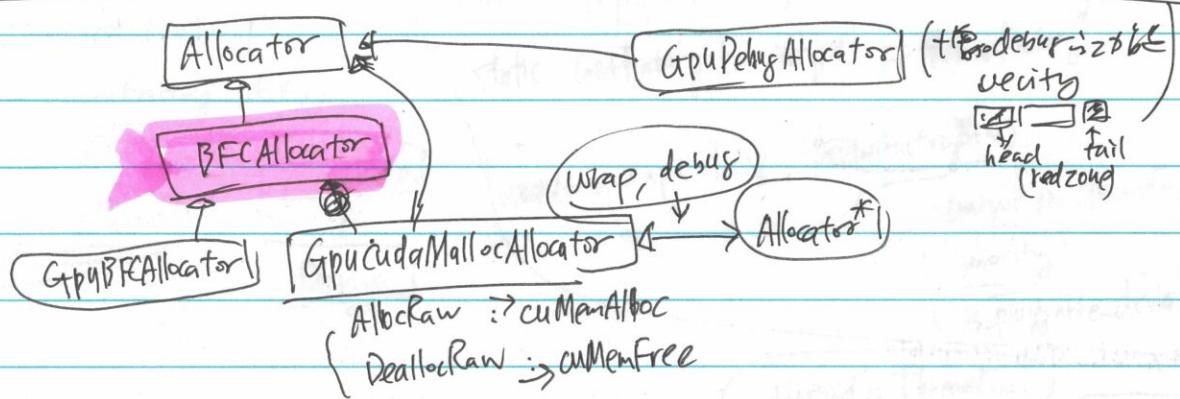
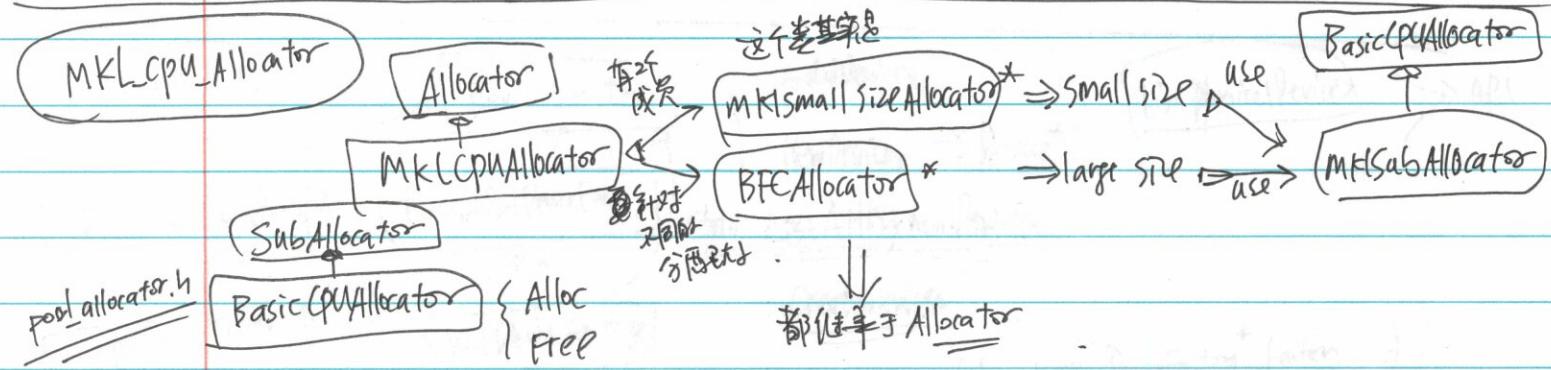
3

AllocatorFactory

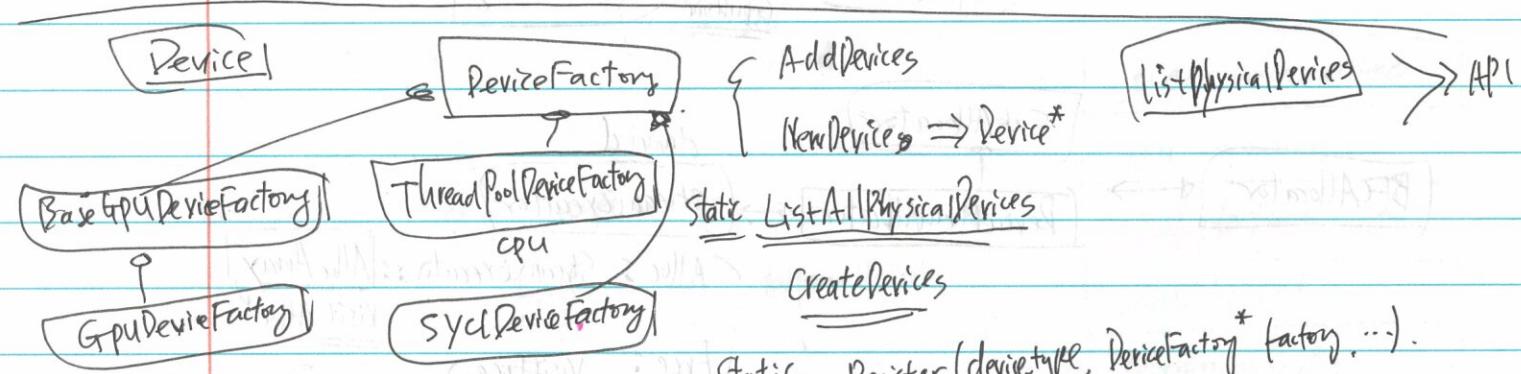
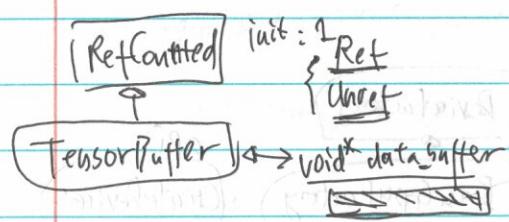
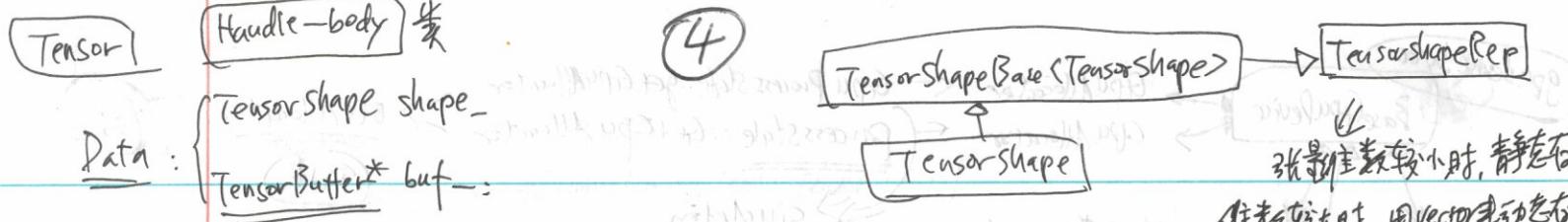
CPUAllocatorFactory

CreateAllocator : new CPUAllocator
CreateSubAllocator :
new CPUSubAllocator
(new CPUAllocator)

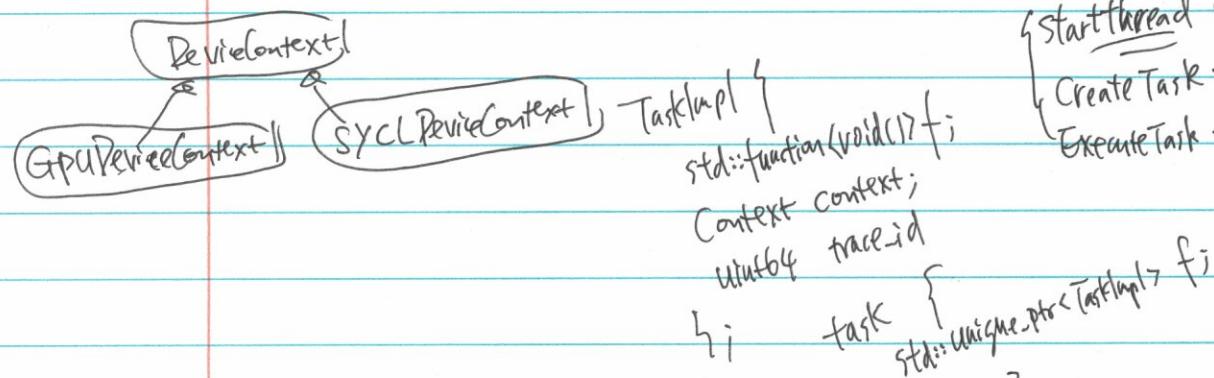
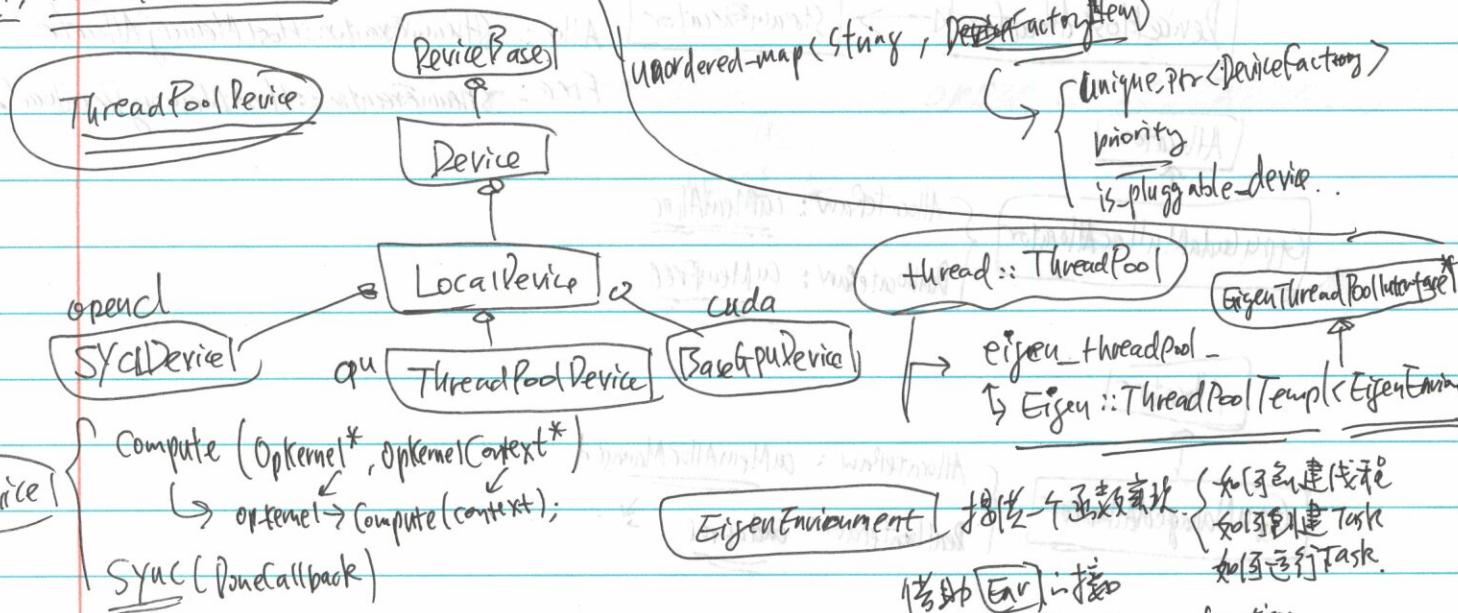
REGISTER_MEM_ALLOCATOR("DefaultCPUAllocator", 100, CPUAllocatorFactory);

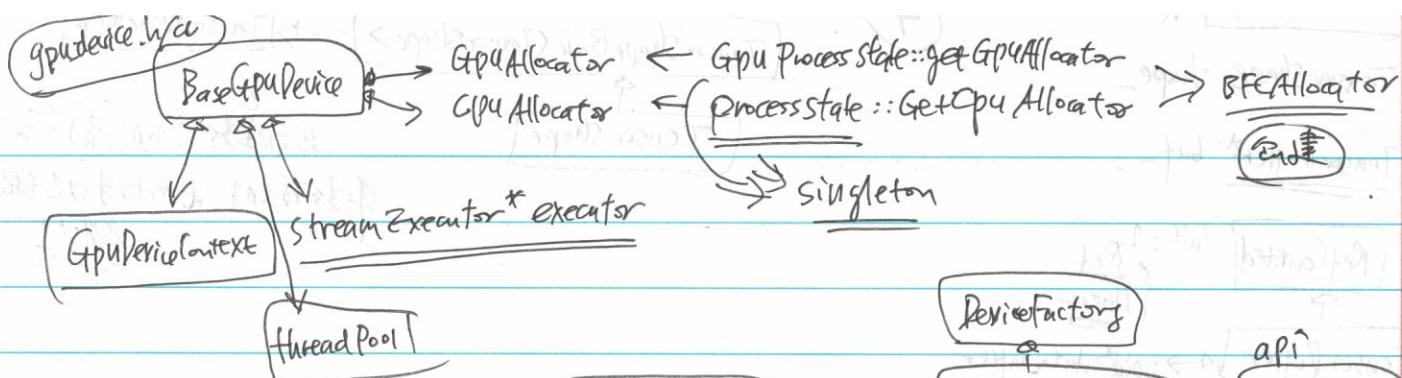


BFCAllocator

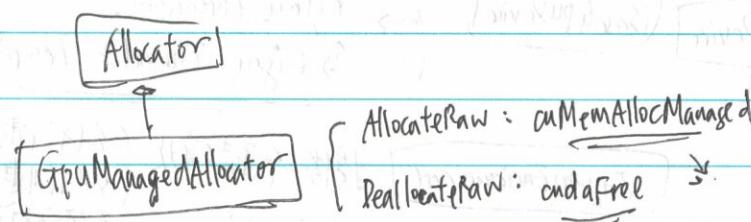
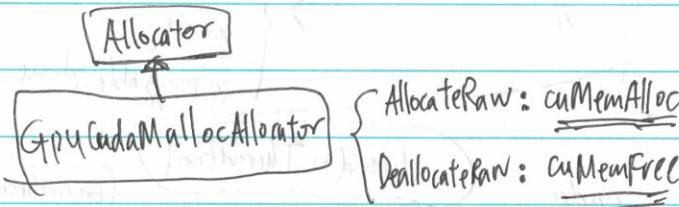
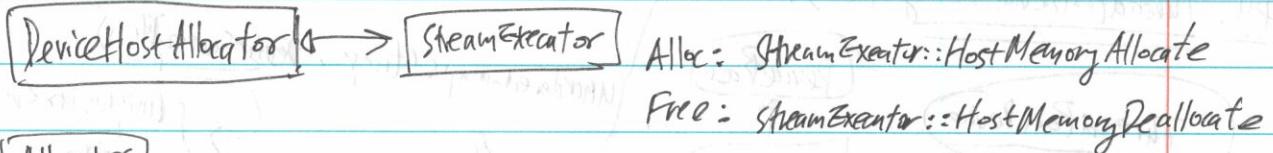
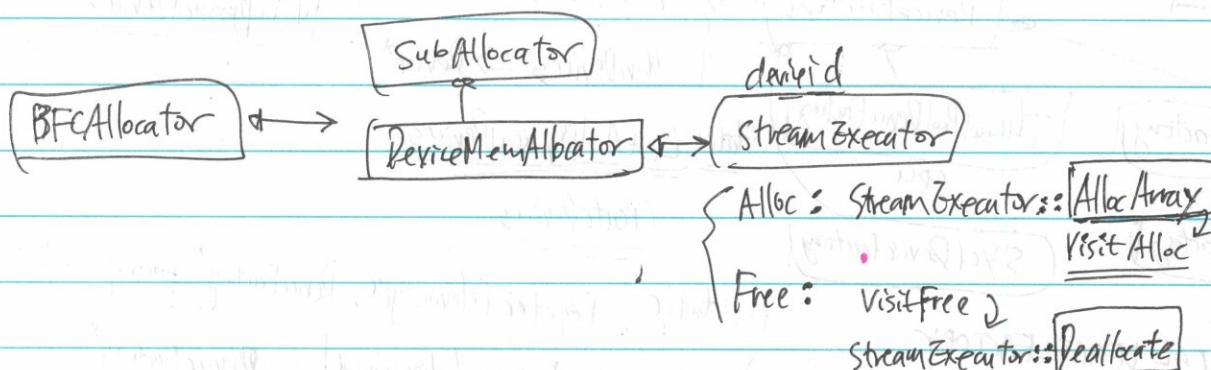
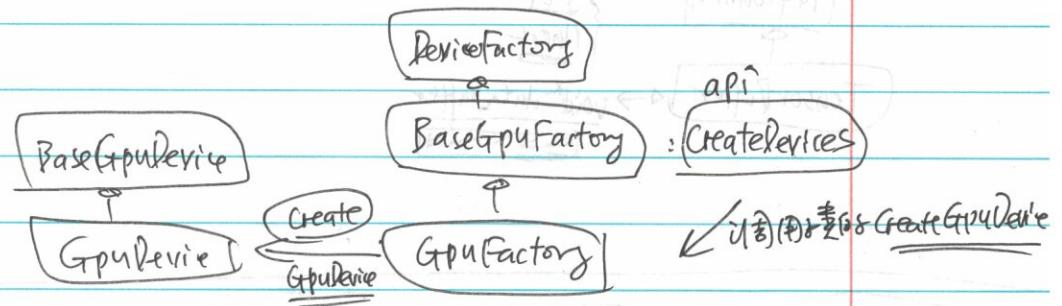


REGISTER_LOCAL_DEVICE_FACTORY
("cpu", ThreadpoolDeviceFactory, b4);





(5)



gpuUtil

template <typename T>

static se::DeviceMemory<T> AsDeviceMemory(const Tensor& t) {

T* ptr = reinterpret_cast<T*>(const_cast<void*>(DMATHelper::base(&t))).
→ t::base()

return se::DeviceMemory<T>(se::DeviceMemoryBase(ptr, t.TotalBytes()));

(6)

DeviceMemoryBase

DeviceMemory<T>

un-typed memory description

void* ::ptr.

size_t ::size.

size_t ::payload

TFToPlatformDeviceIdMap

Singleton object:

TFDeviceId → platformDeviceId

DeviceManager

tfutil API

TFToPlatformDeviceId API

core/common_runtime

core/framework

core/platform OR

core/kernel

StreamExecutor

gcc defines

little_endian :

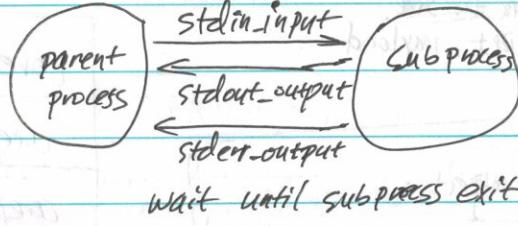
constexpr bool kLittleEndian = __BYTE_ORDER__ == __ORDER_LITTLE_ENDIAN__;

[0xdead]

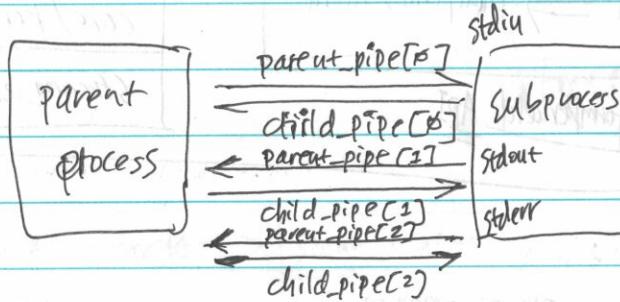
core/platform/default

(7)

SubProcess

bool Start(); pipe fork execvbool Kill(int signal);bool Wait();void SetProgram(String file, const vector<String>& args);int communicate(const string* stdin_input, string* stdcout_output, string* stderr_output), \star
exit status

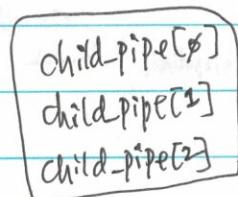
Start:



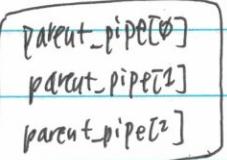
fcntl(parent_pipe[i], F_SETFL, O_NONBLOCK)

fcntl(child_pipe[i], F_SETFL, 0);

父进程管道:



子进程管道



unique_ptr<SubProcess> CreateSubProcess(const vector<String>& args) {

unique_ptr<SubProcess> proc(new SubProcess);

proc->SetProgram(args[0], args);

proc->setChannelAction(chan:STRERR, ACTION_DUPPARENT);

proc->setChannelAction(chan:STROUT, ACTION_DUPPARENT);

}

return proc;



Wait()

StackTrace

StackTraceHandler

CurrentStackTrace

If signal of type \rightarrow callstack

SIGSEGV, SIGABRT, SIGBUS, SIGILL, SIGPPE

8

[core/framework]

op.h

OpRegistryInterface

Lookup (String op_type_name, OpRegistrationData** opreg_data)

Virtual API

LookupDef (String op_type_name, const OpDef** op_def).

not op virtual

OpListRegistry

OpRegistry

Register (const OpRegistrationDataFactory& op_data_factory)

if

singleton

OpRegistry::Global() \rightarrow Register(

{ } \rightarrow (OpRegistrationData* opreg_data) \rightarrow status }

// populate *opreg_data

return status::OK();

};

\Rightarrow std::unordered_map<String, const OpRegistrationData*> registry;

ArgDef [En]

Input arg \rightarrow OpDef

ArgDef [En]

Output arg.

KernelDef

AttrDef attr[u]

input \rightarrow NodeDef
string(u) \rightarrow NodeDef
name/op/device
(String)
(map) attr



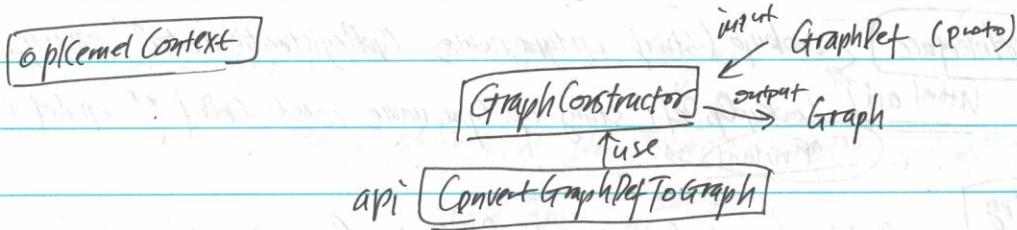
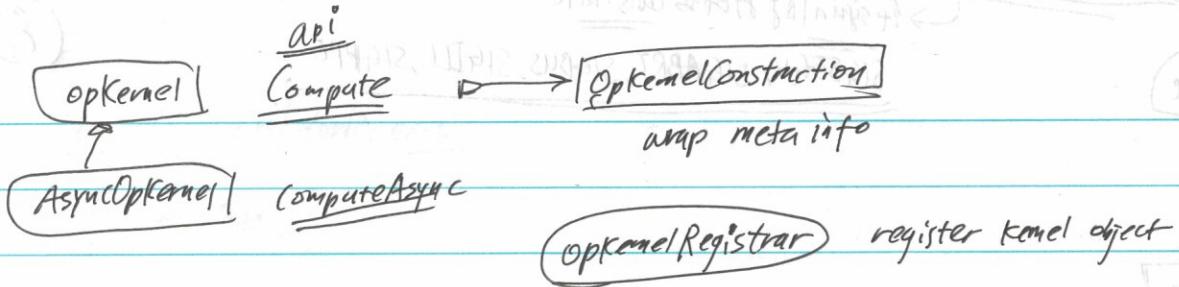
Graph {
vector<Node*>
vector<Edge*>}

Node

{
id_ : int
class_ : NodeClass
in_edges_ : EdgeSet
out_edges_ : EdgeSet.

EdgeSet { array
set }

Node* \rightarrow Edge \rightarrow dst Node*



REGISTER_OP("myopname")

- Attr("<name>: <type>")
- Attr("<name>: <type> = <default>")
- Input("<name>: <type-expr>")
- Input("<name>: Ref(<type-expr>)")
- Output("<name>: <type-expr>")
- Doc(R"(

X Y Z

...
));

static tf::InitOnStartupMarker const registerOpMarker

TF_INIT_ON_STARTUP_IF(isSystemOp() || SHOULD_REGISTER_OP(name))

<< tf::registerOp::OpDefBuilderWrapper(name)

REGISTER_OP : (name)

TF_NEW_IDFOR_INIT(RegisterOpImpl, name, false)

↳ #define TF_NEW_IDFOR_INIT_2(m, c, ...) m(c, __VA_ARGS__)
#define TF_NEW_IDFOR_INIT_1(m, c, ...) TF_NEW_IDFOR_INIT_2(m, c, __VA_ARGS__)
#define TF_NEW_IDFOR_INIT(m, ...) |
TF_NEW_IDFOR_INIT_1(m, __COUNTER__, __VA_ARGS__)

usage: M_IMPL(id, a, b) ...

M(a, b) TF_NEW_IDFOR_INIT(M_IMPL, a, b)

(10)

```

struct InitOnStartupMarker {
    InitOnStartupMarker operator<<((InitOnStartupMarker) const) {
        return *this;
    }
};

template <typename T>
InitOnStartupMarker operator<<(T && v) const {
    return std::forward<T>(v)();
}

```

! cond ? InitOnStartupMarker{} : (InitOnStartupMarker{}) << f

```
#define TF_ML_ON_STARTUP_IF(cond) \
    (std::integral_constant<bool, !(cond) >::value) \

```

? tensorflow::InitOnStartupMarker{}

: tensorflow::InitOnStartupMarker{}

```
InitOnStartupMarker OpDefBuilderMapper::operator() {
```

OpRegistersGlobal() -> Register

Tbuilder = std::move(builder_)] (OpRegistrationData* op-reg-data) -> status {

return builder_.Finalize(op-reg-data);

) ;

return {};

{ OpDef op-def (proto)

OpShape(ifference_fn, shape_inference_fn,
bool is-function_op = false).

REGISTER_OP_IMPL

```
TF_ML_ON_STARTUP_IF(cond) |
```

<< tensorflow::register_op::OpDefBuilderWrapper(name), Attr(...)

• Input(...)

• Output(...)

• Doc(...)

构造者構式

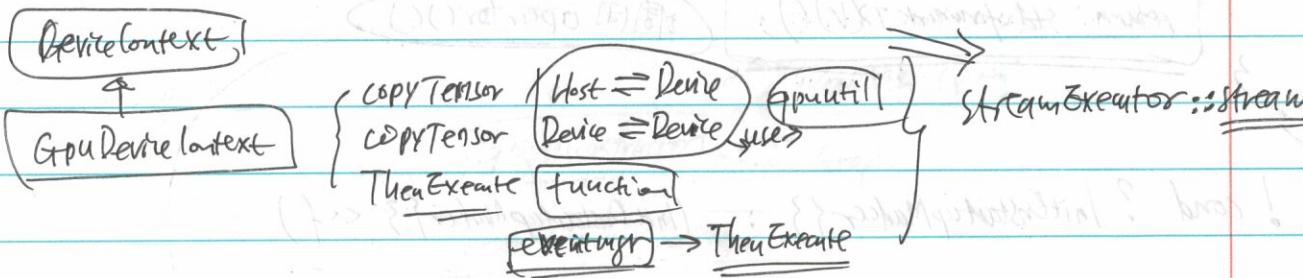
return OpDefBuilderWrapper{

core/common/runtime/device

core/common/runtime/gpu4

core/common/runtime/

(11)



DirectSession DeviceMgt

Static DeviceMgt a \rightarrow std::vector<std::unique_ptr<Device>> devices;

unordered_map<StringPiece, Device*> device_map;

CancellationManager

cuda_dun.la/ce

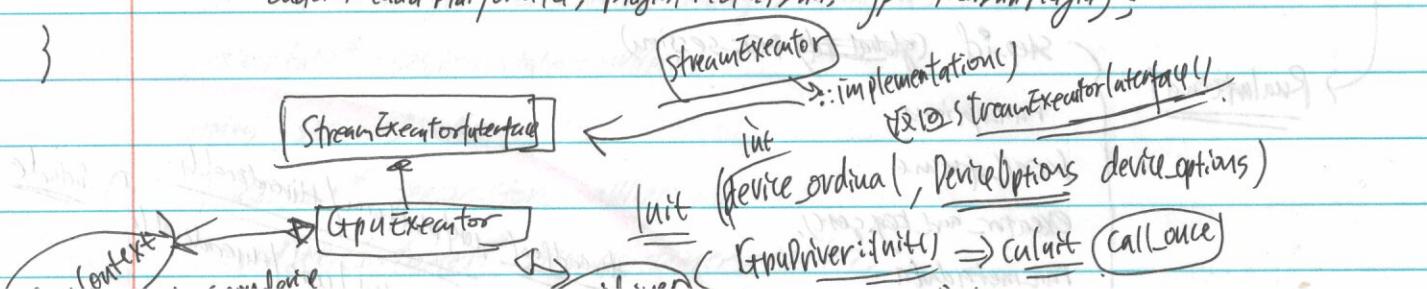
(12)

define REGISTER_MODULE_INITIALIZER(name, body) | ↴
↳ REGISTER_INITIALIZER(module, name body)
(type) = map
↳ static void google_init_##type##_##name() { body; } | ⇒ 这个函数
stream_executor::port::initializer google_initializer_##type##_##name(|) 用这个函数
google_init_##type##_##name) | ↴ 等待调用函数
(factory, 直接调用函数)
REGISTER_MODULE_INITIALIZER(register_idn, { stream_executor::initializer_cudnn(); })

```
pluginRegistry::instance() -> RegisterFactory<PluginRegistry::DnnFactory>
    cuda::KludaPlatform d, gpu::ICudnnPlugin, "cudNN",
    [ ] (internal::StreamExecutorInterface* parent) -> dnn::DnnSupport* {
        gpu::GpuExecutor* cuda_executor = dynamic_cast<gpu::GpuExecutor*>(parent);
        gpu::CudnnSupport* dnn = new gpu::CudnnSupport(cuda_executor);
        if (!dnn->nice().ok()) { ... }
    };
}
```

`PluginRegistry:::instance() -> SetDefaultFactory()`

`cuda::CudaPlatformId, PluginKind::kDnn, gpu::kCuDnnPlugin);`



DirectSession:

RunCallable (CallableHandle handle, const Vector<Tensor>& feed_tensors,

const Vector<Tensor>& fetch_tensors,

CallableHandle* collectors_of_type
CallableWithJobs

RunMetadata* run_metadata,

const ThreadPoolOptions& thread_pool_options)

{
executor_and_keys}

thread::ThreadPoolInterface* interop_threadpool
thread::ThreadPoolInterface* interop_threadpool

(Callables -)

unordered_map < int64, Callable > callable_lock (iftp)

struct Callable {

CallableHandle (int64)

std::shared_ptr<ExecutorAndKeys> executors_and_keys;

std::shared_ptr<FunctionInfo> function_info;

};

unordered_map < string, std::shared_ptr<ExecutorAndKeys> > executors, // executor lock.

unique_ptr<

vector<FunctionInfo> functions_;

protected by

(CallFrameInterface)

RanCallable(CallFrame) call_frame(this, executor_and_keys.get(), feed_tensors) fetch_tensors);

ExecutorAndKeys

input_types

output_types

RunInternal

step_id (global id per session)

run_options

&call_frame

ExecutorAndKeys get(),

run_metadata

threadpool_options

threadpool::core/platform/threadpool.h

or include

core/lib/core/threadpool.h

Executor::Args::Runner default_runner = [pool](Executor::Args::Closure c) {

pool->schedule(std::move(c));

ThreadPool

threadpool_interface + underlying_thread

= EigenThreadPool

with EigenEnvironment template arguments

TF-1.14

Executor

succeed/failture, will call DoneCallback

(14)

virtual void RunAsync(const Args& args, PoneCallback done)

→ typedef std::function<void (const Status&) >
PoneCallback

13) ~~args~~: status Run (const Args& args) {
 Status ret;

Notification n;

RunAsync (args, [ret, &n] (const Status& s) {

ret = s;

n.Notify();

});

n.WaitForNotification();

return ret;

}

[Args]: int64 step_id = 0;

Rendezvous * rendezvous = nullptr;

StepStatsCollectorInterface* stats_collector = nullptr;

CallFrameInterface * callframe = nullptr;

CancellationManager* cancellation_manager = nullptr;

SessionState* session_state = nullptr;

string session_handle;

Tensorstore* tensor_store = nullptr;

ScopedStepContainer* step_container = nullptr;

CollectiveExecutor* collective_executor = nullptr;

typedef std::function<void ()> Closure;

typedef std::function<void (Closure)> Runner;

Runner runner = nullptr;

Extractor Barrier

A class to help run multiple executors in parallel and wait until all of them are completed

```
typedef std::function<void (const Status &)> StatusCallback;
```

15

Dat21 { ReaderActions . render -
status callback done_cb
int pending_

Api { Get \Rightarrow StatsCallback : std::bind (&Executor::WhenDone, this, std::placeholders::_1).

whenDone \leftarrow 每个 executor 生成之后会调用

{ } mutex_lock (&ad_);

if (status_group_.ok() && !s.ok()) {

error-render = rend2 - ;

error-tendenz. \rightarrow Ref();

3

Status_group_.Update(s);

if (--pending_ == 0) { 11# [0-1]

`std::swap(done, done_cb_);`

States = status_group.as_summary_status(): (1) It states summary

}

3

if (error_reude2) {

~~qWer_render2~~ → StartAbort(...);

~~error->render() → href();~~

3

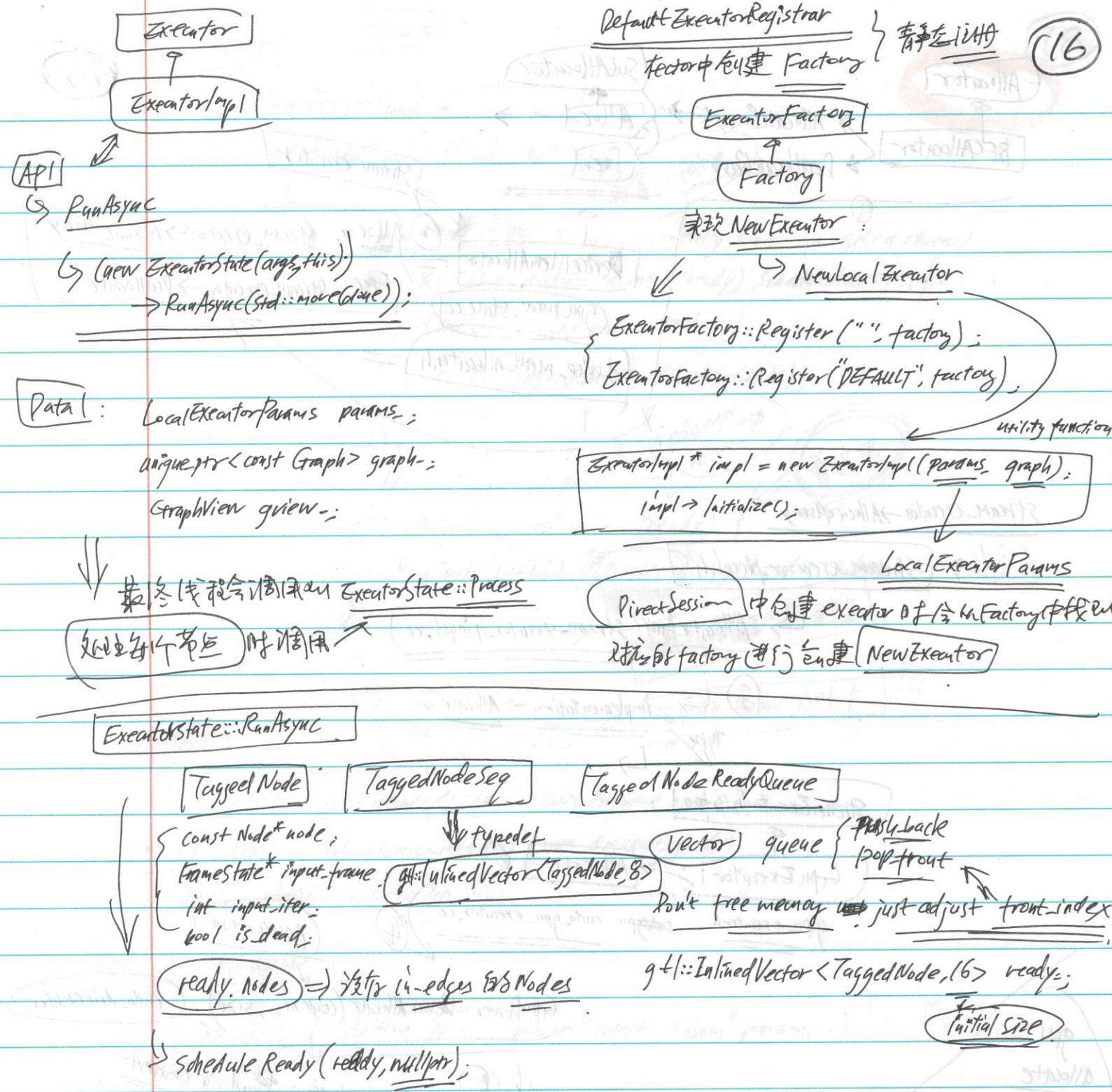
if (done) {

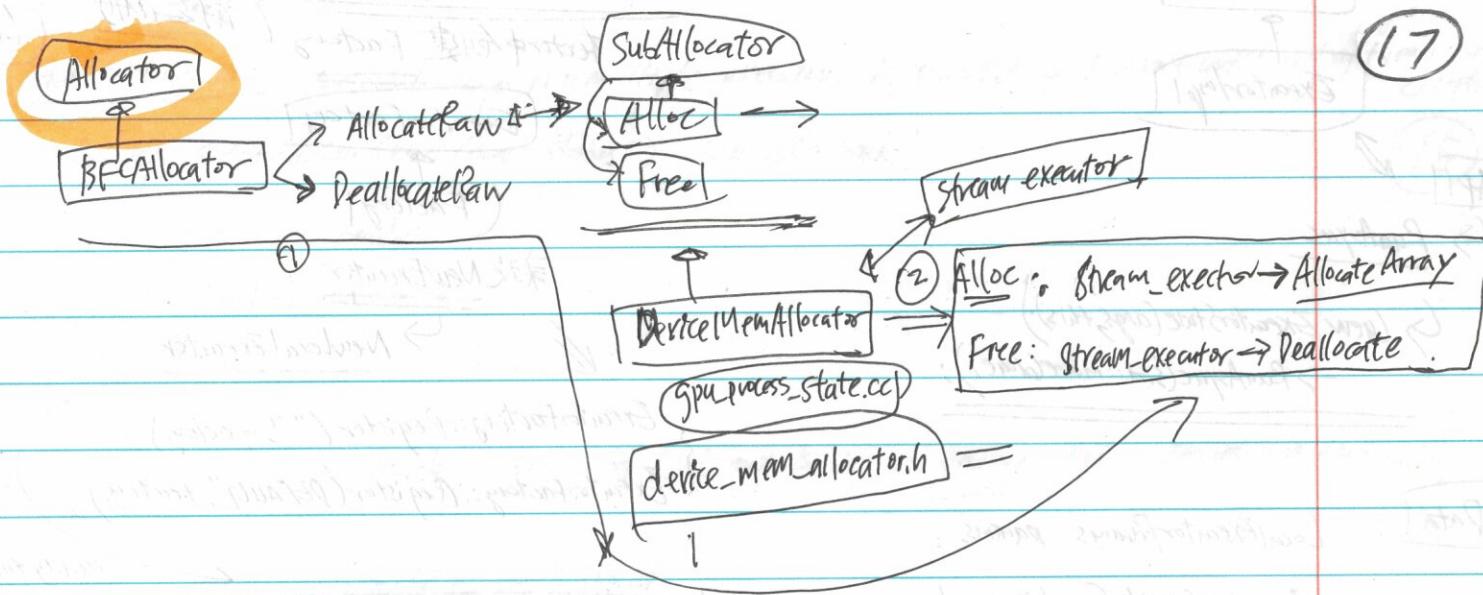
delete this;

11 delete ~~fig~~

done(status); // 调用 (用) 完成方法做回调.

}





StreamExecutor->AllocateArray

↳ StreamExecutorImpl.h

↳ AllocateArray (StreamExecutor-impl.cc)

③ ↳ Implementation → Allocate
type ↳

StreamExecutor/Interface

Gpu Executor ↳ Allocate

gpu_executor.h

cuda_gpu-executor.cc

GpuContext

GpuDriver::DeviceAllocate (context_, size);

cuda_driver.h/cc

④ ↳

⑤ ↳ cuMemAlloc

totk in Gpu device Context

Deallocate

Free

→ GpuDriver::DeviceDeallocate (context_, ptr)

⑥ ↳ cuMemFree

gpu
allocate
memory
for host/flow

Executor::RunAsync

TF2.5

78

↳ ExecutorState(T) :: RunAsync

↳ scheduleReady(&ready, nullptr);

↳ if WithAllKernelsInline: // sequentially run in single thread

RunTask([this, ready = std::move(tready), scheduledusec]()) {

for (auto& tagged_node : ready) {

process(tagged_node, scheduledusec);

}

} else: // run in thread pool

for (auto& tagged_node : tready) {

RunTask([=]() { process(tagged_node, scheduled_usec); },

3

in edges <=>
root

看这里

对一个线程处理所有节点

!!!

RunTask

Wrap.

closure function

runner_([....] { closure(c); });

(运行 thread pool) !!!

[alignas(64)] static std::atomic<int64> num_enqueue_ops(0);

[alignas(64)] static std::atomic<int64> num_dequeue_ops(0);

→ Cache line, avoid false-sharing.

process tagged_node (root)

↳ processSync (NodeItem item, OptKernelContext::Param* params...)

OptKernelContext ctx(params, item.num_outputs);

OptKernel* op_kernel = item.kernel; ← Each graph node has one kernel to associate.

Device* device = ...

(when to associate ???)

device->Compute(op_kernel, &ctx);

if NewLocalExecutor {

s = ProcessOutput(item_bctx, outputs>data(), stats);

(1) new executor
(2) call executor::Initialize()

Executor_and_keys

items [Vector]

PerPartition ExecutorAndLib

(19)

↳ Executor

const bool can_execute_synchronously = executor_and_keys->items.size() == 1 && call_timeout == 0;
graph 不会对 partition 有影响
true.

if (can_execute_synchronously) {

items.executor->Run(args)



Notification n;

[RunAsync](args, [&ret, &n])(...) {

ret = s;

n.Notify(c);

}

n.WaitForNotification();

ret

→ (通过调用 runner 的 threadpool schedule)

default_runner->[pool](closure c) {

pool->Schedule(std::move(c));

};

device.h

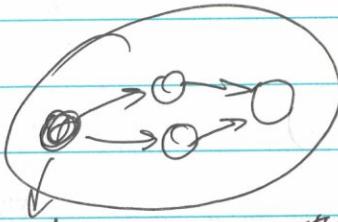
DeviceBase
of
Device

compute(OpKernel *opKernel, OpKernelContext *context) {

opKernel->compute(context);

}

core/kernels/ xx. (identity-op.h)



ready -> node step 1 process #get it.

task

gpuDevice

手写注释: 'Compute'

两个 opKernel 和 device 之间
需要 switch gpu Context
和 gpuDevice. 需要设置 cuDNN (cuDNN)

Process 之后 thread pool #get it to it.

故单线程且无线程池。

Executor

ExecutorImpl

API: RunAsync

Data : Immutable Executor State

core/common/runtime/

ImmutableExecutorState

Data

LocalExecutorParams params_

LocalExecutorParams

Local_executor_params.h

GraphView gview_;

vector<const NodeItem*> root_nodes_;

API : Initialize (const Graph & graph);

dtor

for loop gview_.num_nodes_ {

NodeItem* item = gview_.node(i);

if (item) {

params_.delete_kernel(item->kernel);

}

initialize

GraphView

Graph

Deleter
DataPointSessionMetaData
FunctionLibraryRuntime
function (status)
const shared_ptr<NodeParam>
OpKernel**>
create_kernel
function (void (OpKernel**))
delete_kernel;

DirectSessionData
CreateExecutorOp
populate

(item) NodeItem* ← Node*(n)

preprocess all nodes "and create Op kernel for each node"

item ← gview_.node(n->id())

node id ≠ index NodeItem GraphView

NodeItem :

node_id

Data |

OpKernel* kernel

boolean flags (lots of)

placement new

节点 NodeItem → NodeItem
节点 bytes NodeItemBytes
节点 offset NodeItemOffset
节点指针 NodeItem*
节点大小 NodeItemBytes()

GraphView

array of NodeItem

→ NodeItem

NodeItem* node (int32 id) {
(0 ≤ id < num_nodes_);

offset = node_offset[id];

(NodeItem*)(space_ + node_offset[id]);

(-offset)

Space

Char[]

node_offsets_

→ node offset spaces_offset

var() API

实现

NodeItemBytes()

FunctionLibrary

Simple data structure

Data → unique_ptr<FunctionLibraryDefinition> flib_def

→ unique_ptr<ProcessFunctionLibraryRuntime> proc_flr; → per-device

Container → unordered_map<Device*, std::unique_ptr<FunctionLibraryRuntime>> flr_map_.

Interface class

FunctionLibraryRuntime function.h

IOPRegistryInterface

FunctionLibraryDefinition

/core/framework

FunctionLibraryRuntimeImpl

Data

const DeviceMgr* const device_mgt;

Device* const device_;

Env* const env_;

const FunctionLibraryDefinition* const base_libdef_;

Graph Optimizer optimizer_;

Executor::Args::Runner default_runner_;

GetFunctionLibraryDefinition()

这个是(b)的 op kernel 的优化
这个是(c)的 op kernel 的优化
其中第一个是直接的
第二个是通过 optimizer

pool->schedule(task)

std::function<Status (const std::shared_ptr<const NodeProperties>&, opkernel**)> createKernels_;

(flat_hash_map<Handle, std::unique_ptr<Item>>) items_;

FunctionHandleCache function_handle_cache_;

ProcessFunctionLibraryRuntime

functionKey

Data → unordered_map<string, FunctionLibraryRuntime::Handle> table_;

→ unordered_map<FunctionLibraryRuntime::Handle, UniquePtr<FunctionData>> function_data_;

→ unordered_map<FunctionLibraryRuntime::Handle, UniquePtr<MultipleDeviceFunctionData>> mdevice_data_;

→ unordered_map<Device*, std::unique_ptr<FunctionLibraryRuntime>> flr_map_;

associate

(device->listDevices())

NewFunctionLibraryRuntime

Utility functions

New FunctionLibraryRuntime& (...).

OpKernel

• /core/framework/

IdentityOp

• /core/kernels/

IdentityOp (OpKernel construction * context)

: OpKernel(context) {}

(22)

void Compute (OpKernelContext* context) {

if (IsRefType (context->InputDType(0))) {

context->ForwardRefInputToRefOutput(0, 0);

} else {

context->SetOutput(0, context->Input(0));

}

}

bool IsExpensive() { return false; }

macro

REGISTER_KERNEL_BUILDER(Name("Identity"), Device(DEVICE_CPU), IdentityOp);

class Name

↓ ← (n) helper macro, 164478572

ctr, opname, kernel_builder_expr, is_system_kernel, ...)

ctr =
COUNTER_

Y

Y

false

IdentityOp (class name)

TF_EXTRACT_KERNEL_NAME(RegisterKernelBuilderImpl_2, kernel_builder, is_system_kernel, -VA_ARGS_)

⇒ (m, kernel_builder, ...)

TF_EXTRACT_KERNEL_NAME_IMPL(m, TF_EXTRACT_KERNEL_NAME ## kernel_builder, -VA_ARGS_)

⇒ (m, ...)

Name("Identity").Device(DEVICE_CPU)

m (-VA_ARGS_)

REGISTER_KERNEL_BUILDER_IMPL_2(TF_EXTRACT_KERNEL_NAME_Name("Identity").Device(DEVICE_CPU))

TF_NEW_ID_FOR_INIT(RegisterKernelBuilderImpl_3, opname, kernel_builder_expr, is_system_kernel, ...)

if (name_str)

name_str, tensorflow::register_kernel::Name(name_str).Device(DEVICE_CPU)

namespace register_kernel {

struct Name : public KernelRefBuilder

Name(const char* op) : KernelRefBuilder(op) {}

build_tf_ref_to_kernelrefbuilder,
Device(DEVICE_CPU).

};

REGISTER_KERNEL_BUILDER_IMPL_3 (....)

(23)

static tensorflow::InitOnStartupMarker const register_kernel_#_gt_ =

TF_MERGE_IF (KernelDef) & ← InitOnStartupMarker

[[T] (TensorFlow::KernelDef const * kernel_def) {

tensorflow::KernelFactory::OptKernelRegistrar registrar(

[kernel_def], [#_VA_ARGS_] ← 参数字符串)

[T] (TensorFlow::OptKernelConstruction* context)

→ TensorFlow::OptKernel {

return new _VA_ARGS_(context); ← 2) 方法且具体情况。
just new ← OptKernel 构造

用 OptKernelConstruction 构造

机制

InitOnStartupMarker & operator<<(fim)

(其中是函数对象
const InitOnStartupMarker)

类体

(是否是函数模板(T), 类视为
调用它的 T::operator())

);

return tensorflow::InitOnStartupMarker{};

(kernel_builder_expr.Build());

这是个 lambda 表达式。返回 InitOnStartupMarker 对象。

输入参数为 const KernelDef *

调用时参数 kernel_builder_expr.Build() 返回值

macro 展开 Name("Identity").Device(DEVICE_CPU).Build()

kernelDefBuilder->fake. 调用 Build() 函数

lambda 表达式实现

用 KernelDef 去初始化(OptKernelRegistrar 对象, 在 factory 中会调用 Register 方法)。

Opkernel Registrar

(24)

↳ for InitInternal(...)

[Opkernel factory]

[ptr OpkernelFactory]

实现原理中 Opkernel 叫法是 map 中
factory = factory.

指向 Opkernel factory.

[API] : Opkernel* Create (OpkernelConstruction* context) override

GlobalKernelRegistry() ← new KernelRegistry (static)

global_registry → registry, emplace (key, KernelRegistration (*kernel_def, kernel_class_name, facto

Kernel Registry { mutex
unordered_multimap<string, KernelRegistration> registry;

[Data] { KernelDef
KernelClassName
OpkernelFactory

temp通过 IdentityOp 建立连接，并没有 Opkernel Factory 与其对应，临时创建一个

lambda表达式。Opkernel Registrar 提供两种构造方法。
1. 从 factory * 构造
2. 从指针构造

→ 故障后可以同样构造 ← ptr Opkernel Factory

Registrar 是类

在实现的 Create 函数 override 就是转调用

这个函数

(DirectSession)

func_info → proc_fn → GetFLR(clevename)
functionLibrary runtime

FunctionLibraryRuntimeImpl (lib) function w/cc

input output

↳ CreateKernel

[NodeProperties] → Opkernel *

wrap

[LocalExecutorParams]

↳ create_kernel [function pointer]

FunctionLibraryDefinition? (map)

find fail ↳ Create Non Cached Kernel

↳ CreateOpkernel
(in op-kernel.cc)

① according to nodeProperties's node_def
find KernelRegistration(item) from KernelRegistry
② KernelRegistration's factory's Create
call
new -> Opkernel

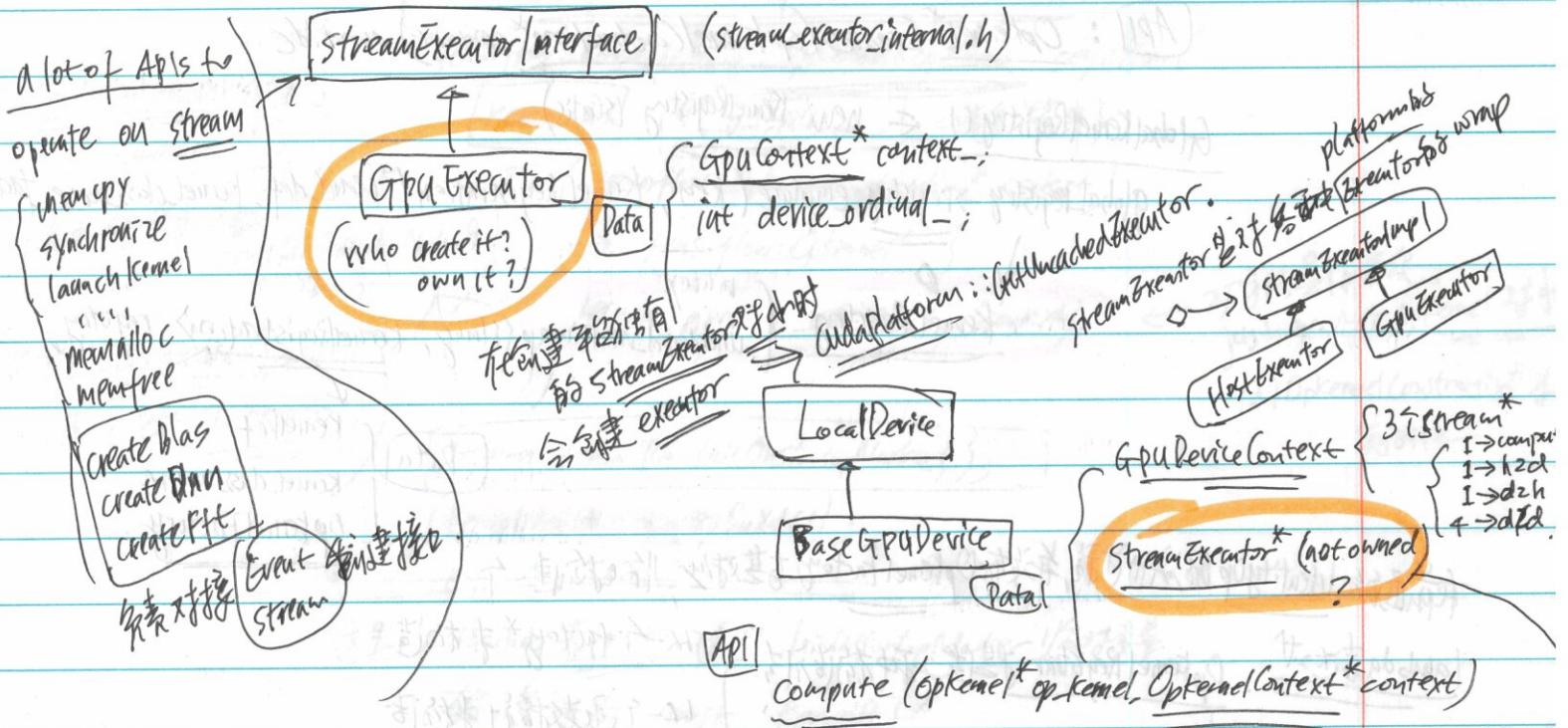
OpKernelConstruction 是一个聚合类，聚合大量的参数，(用于导出为 opDef)

(25)

StreamExecutor

/ StreamExecutor / ...

StreamExecutor > ExecutorCache & refer to
BaseGPUDevice & refer to.



在 BaseGPUDeviceFactory::CreateDevices

Each GPU

BaseGPUDevice

.init

BaseGPUDevice::Init

MultipPlatformManager<cup>

mutex
unordered_map<string, Platform*>

DeviceUtil::ExecutorForTfDeviceId

GpuMachineManager()

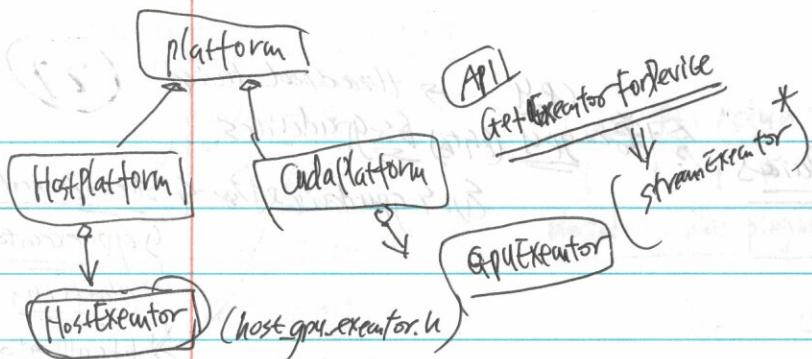
stream_executor::Platform* (gpuinfo)

se::MultipPlatformManager::platformWithName(GpuPlatform)

((StreamExecutor / multipPlatformManager) "CUDA")

stream_executor
Platform (platform)

enum PlatformKind
Cuda
Rocm
Opencl
Host
None



REGISTER_MODULE_INITIALIZER(host_platform,
StreamExecutor::host::initializeHostPlatform());
(init body)

new HostPlatform();

↳ MultiplatformManager::RegisterPlatform();

StreamInterface

HostStream (对线程池-种抽象)
(有线程池的抽象)

API:
BlockUntilDone: 等待为任务都做完。

Notification done;
EnqueueTask([&done](){ done.Notify();});
done.WaitForNotification();

exit task, let task自己通知。 Good

platformId : void*. 和其他platform有不同的是platformId?

namespace {
int plugin_id_value; // 可以是不同的插件有不同的值

3. const platform::Id VAR_NAME = &plugin_id_value;

HostExecutor

Memcpy
memset
MemZero

向后是生成一个lambdaTask

放入由 HostStream 处理 (通过调用)

fft

blas

dau 不支持

ExecutorCache

container
data

mutex

map<int, Entry> cache,

device ordinal

mutex

vector<

pair<StreamExecutorConfig,
unique_ptr<StreamExecutor>> table;

支持 int key

per-entry mutex allow concurrent initialization of
different entries, since initializing an executor may be
expensive

typedef std::unique_ptr<StreamExecutor> ExecutorFactory();
std::function<ExecutorFactory> factory;

StreamExecutor

func : {
platform*
StreamExecutorInterface*
device ordinal

每个 platform 会创建一个 StreamExecutor

new StreamExecutor(*this, new HostExecutor_

device ordinal)

不支持

共享

ExecutorForDevice

↳ 通过 DeviceId 得到 platformId

↳ 通过 gpu platformId 得到 DeviceId

↳ Get ExecutorForDevice.

↳ MultiPlatform 会根据名字 "cuda"

找 platform

DirectSessionFactory

⇒ NewSession

- ① DeviceFactory::AddDevices
- ② new DirectSession

CPU → threadpooldevice

copy to gpu (ptr) for gpudevices.

GPU & gpudeviceptr ← StreamExecutor

(27)

GpuExecutor

Platform

bfcallocator

GpuDevice

Data

Streamgroup {
 - compute stream
 - host2device stream
 - device2host stream
 - device2device stream.

GpuDeviceInfo {
 - Stream*
 - DeviceContext*
 - GpuId
 - (ann)

threadPool*

reference counted

 - gpuAllocator*
 - GpuAllocator*
 - StreamCreator*
 - GpuDeviceContext*

processstate::GetGpuAllocator

GpuProcessState (singleton) own (由于它是singleton, 所以 delete allocator)

GetGpuAllocator

gpuAllocator*
 - not owned

StreamCreator* ← init of gpus, PC, Cache, Town.

(Session to 3 device (GpuDevice) to 1 GPU (Compute))

refcounted

refcounted

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

↑

↓

StreamExecutor::cuda

(28)

cuda blas

cuda fft

Cuda event

cuda dnn

Cuda Stream

Cuda timer

blas/fft/dnn, unplugin isn't register by pluginRegistry

↓ ↓
pluginkind::kblas pluginkind::kfft pluginkind::kdnn

REGISTER_MODULE_INITIALIZER(register_xxx, { streamExecutor::initialize_xxx,

RegisterFactory[]

new gpu::CUDABLAS(gpu)

new gpu::CUDAFFT(gpu)

new gpu::CUDADNN(gpu),

lambda
expression

3)

调用

每个 plugin 都有一个 plugin id 其实是 void*

定义在 cc 之中的全局变量，有唯一性。
(偏译单元)

StreamExecutor::AsDnn()

get Dnn API
call CreateDnn()

[dnn] ← void* / pluginId

Implementation → CreateDnn()

StreamExecutorInterface

GpuExecutor | CreateDnn

On demand call pluginfactory to create it.

pluginRegistry* registry = pluginRegistry::instance();

pluginRegistry::PluginFactory* fn = Registry::GetFactory<PluginRegistry::PluginFactory>

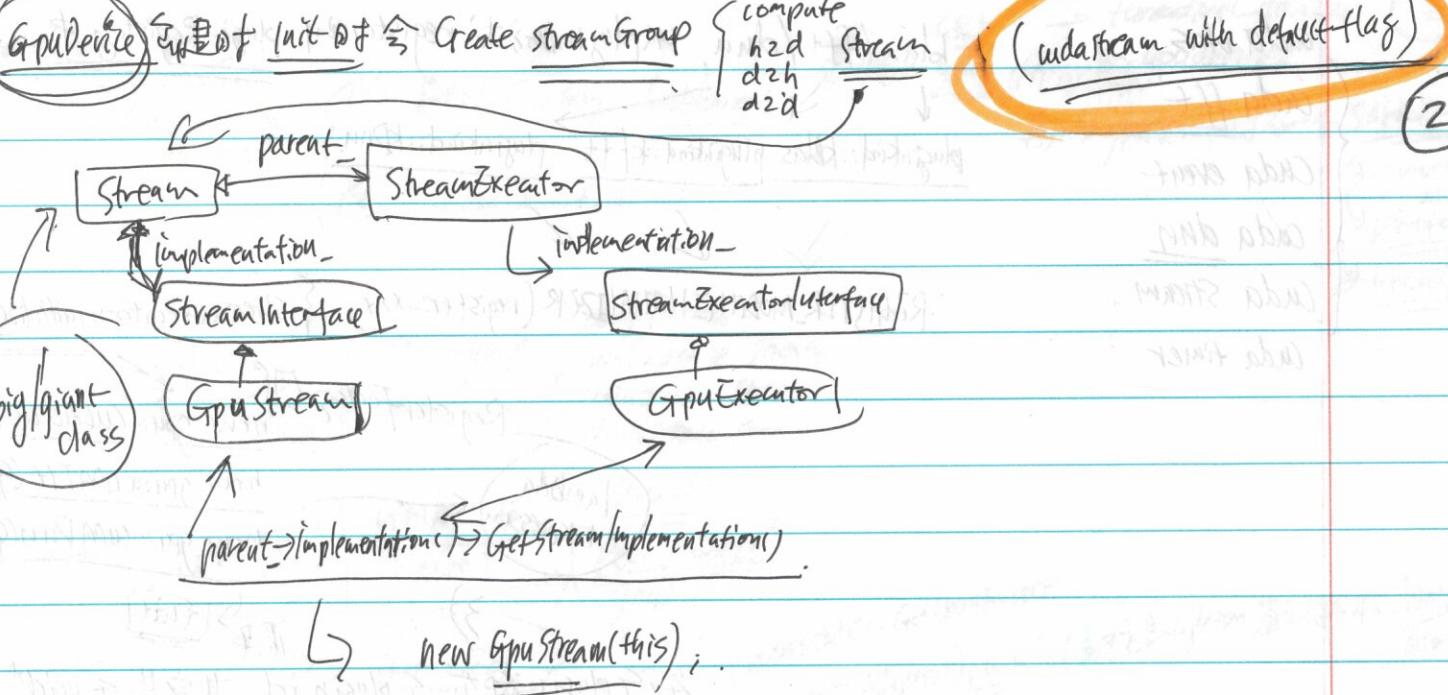
(cuda::KludaPlatformId,

pluginConfig_dnn);

↓ void* / pluginId

fn(this)

new CUDNNSupport



Stream::init() {
 parent → AllocateStream(this) ⇒ gpuExecutor → AllocateStream(..)
 {
 ↳ As GpuStream(stream) → init();
 ↓ Implementation
 GpuStream object pointer.
 }
 ↳ init()
 if (P) GpuDriver::CreateStream

onCreateStream(stream, 0)

→ It's non-blocking stream!!! default

non-blocking stream → default stream or (legacy stream)

Sync

Sync

default stream → legacy stream

per-thread default stream

→ non-blocking stream,
default flag Stream, Sync

legacy stream Sync