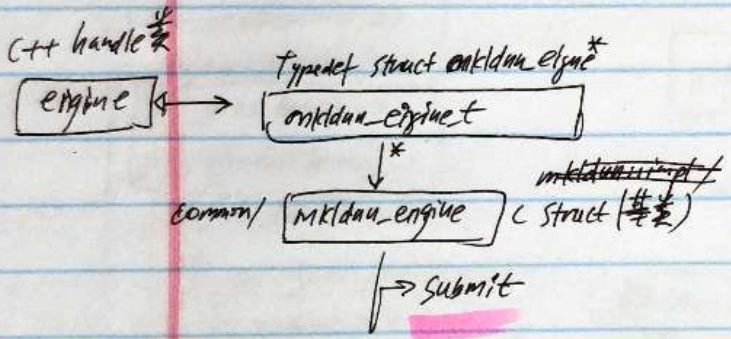
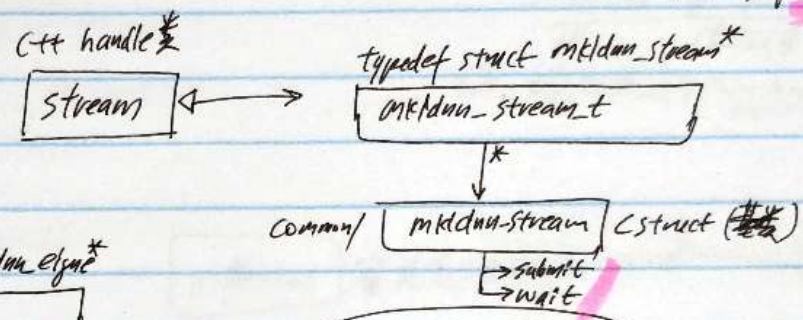
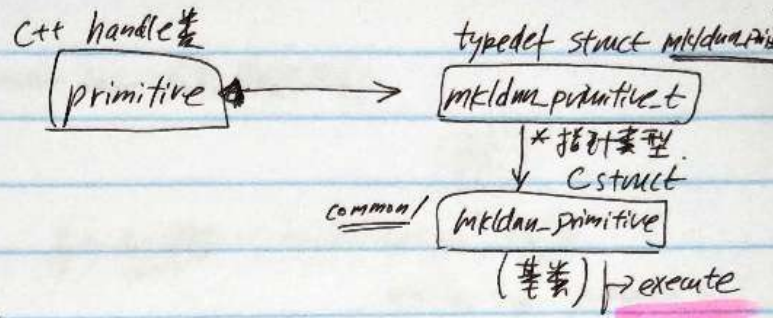
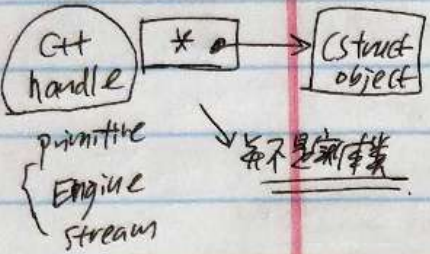
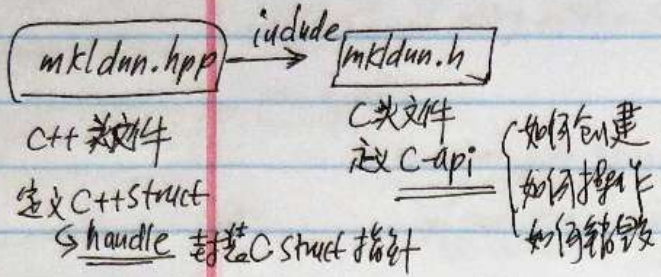


# onePNN-0.1

实现文件  
实现这些api

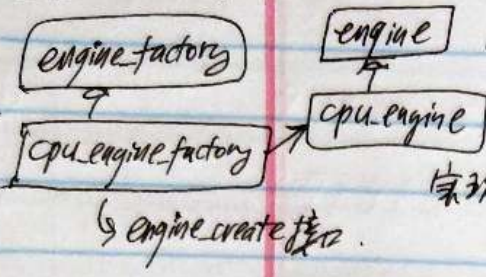


- 封装  
利与所不面  
类包含
- mkldnn\_primitive\_base\_desc\_t
  - mkldnn\_memory\_primitive\_desc\_t
  - mkldnn\_reorder\_primitive\_desc\_t
  - mkldnn\_convolution\_primitive\_desc\_t
  - mkldnn\_pooling\_primitive\_desc\_t
  - mkldnn\_relu\_primitive\_desc\_t
  - mkldnn\_lrn\_primitive\_desc\_t
  - mkldnn\_inner\_product\_primitive\_desc\_t



cpu/cpu\_engine

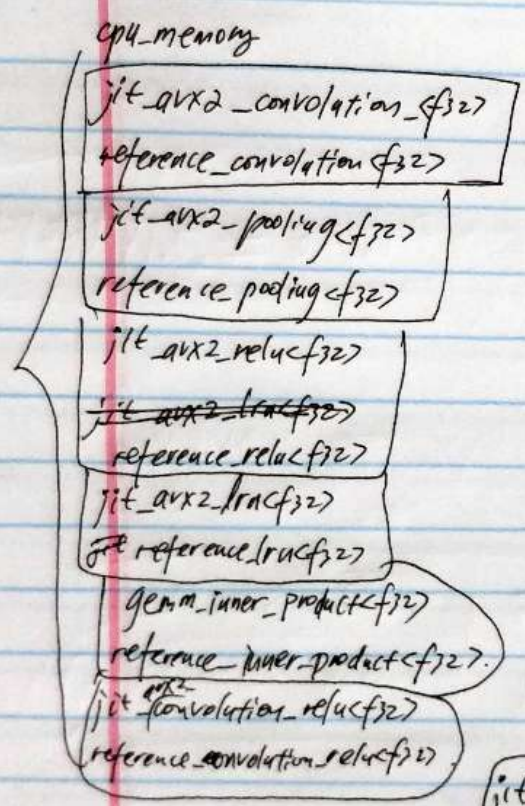
mkldnn/impl/cpu



// using mkldnn::impl::engine = mkldnn\_engine.

每次 submit 接口 提交 primitive\* batch  
(primitive -> execute)

每个 engine 支持不同的 primitives. get\_primitive\_init 返回 数组 指针



多态

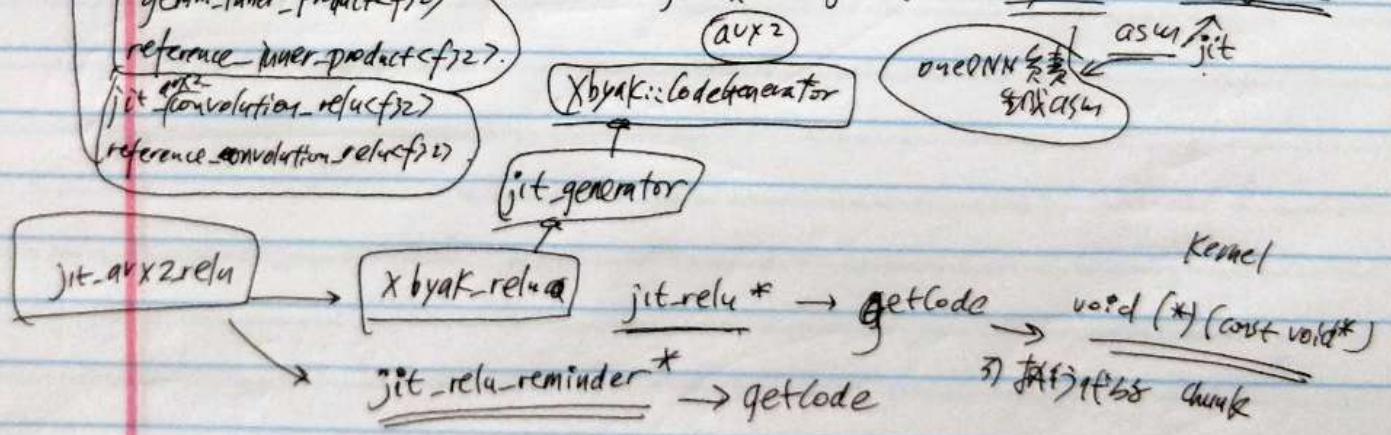
每个 description 一个 implementation\*, 接口

mkldnn::impl::primitive\_impl

包含 接口 创建这个 primitive 子类 接口

reference 实现是用 omp parallel directive 来改的

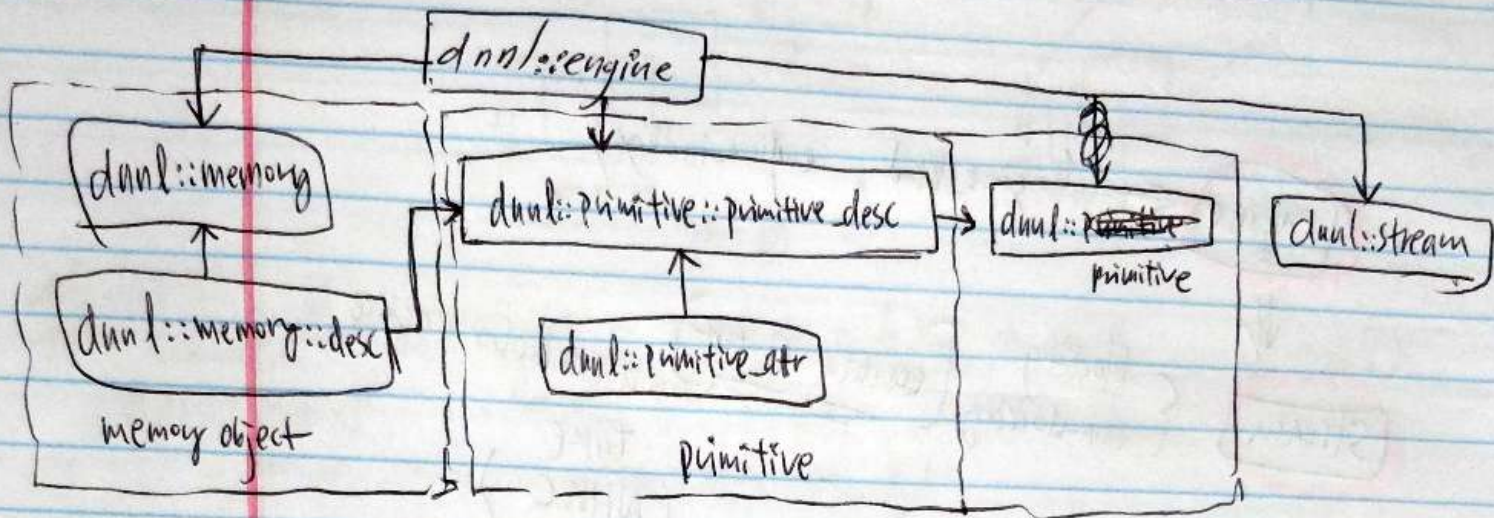
jit 实现是 手写 汇编 是 xbyak 生成 binary code





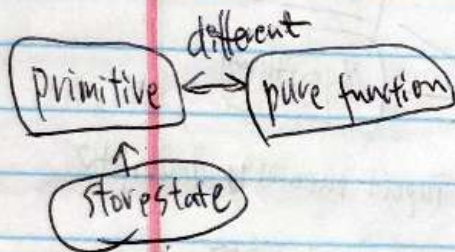
# OneDNN

programming model: primitives, engines, streams, and memory objects



## primitive

encapsulates a particular computation:   
 { forward convolution   
 backward LSTM computations   
 a data transformation.



**engine**: An abstraction of a computational device: a cpu, a specific gpu.

**streams**: encapsulates execution context tied to a particular engine.

**memory objects**: encapsulate handles to memory allocated on a specific engine.   
 tensor dimensions, data type, and memory format.

dnnl\_memory\_desc\_t  $\iff$  dnnl::memory::desc



engine ← (engine\_kind, engine\_index)

Stream

{ blocking  
non-blocking } execution context → a particular engine.

Memory object

memory descriptor

type  
(NHWC...)  
format

engine

engine

algorithm

primitive-descriptor

input memory descriptor

output memory descriptor

move  
argument

primitive

→ execute (Stream, { input memory object  
output memory object })

Stream.wait()



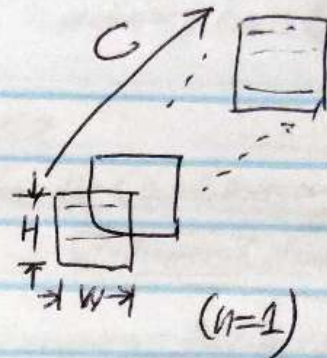
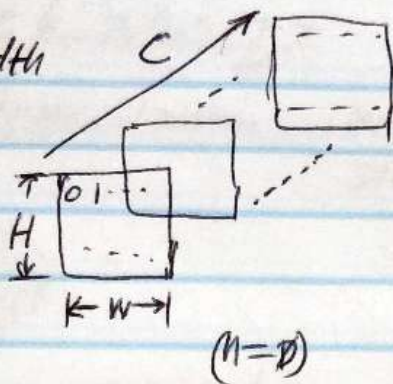
Data format (layout)

(Caffe2 default) (TF support)

$$N=2, C=16, H=5, W=4$$

$$\underline{N} \times \underline{C} \times \underline{H} \times \underline{W}$$

batch channel High Width



$$\underline{\text{Value}(n, c, h, w)} = n \times \underline{CHW} + c \times \underline{HW} + h \times \underline{W} + w \times \underline{1} \quad (= \text{offset})$$

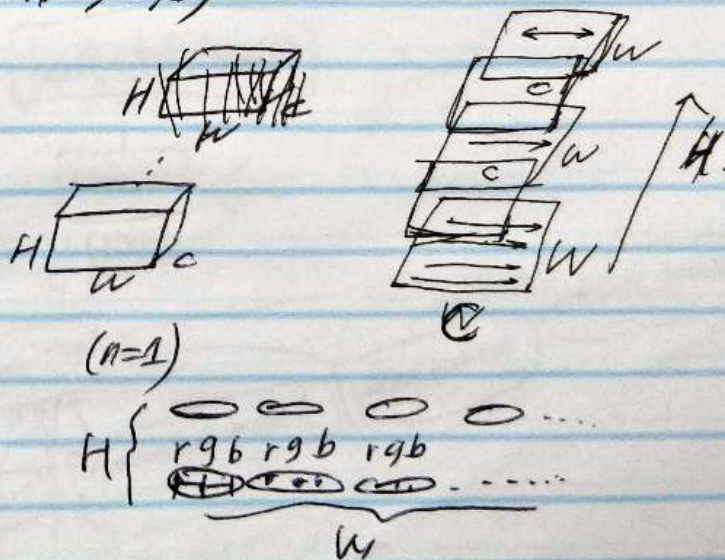
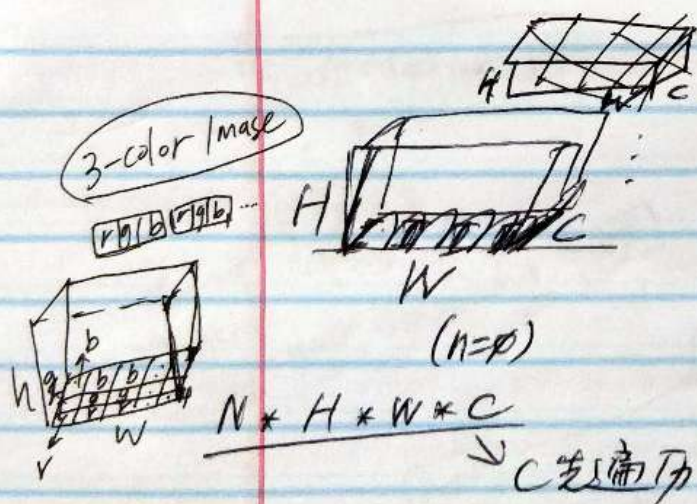
(logical offset) index strides = (CHW, HW, W, 1)

inner-most dimension  
outer-most dimension

$N \times H \times W \times C$  (tensorflow default)

$$\text{offset\_nhwc}(n, c, h, w) = n \times \underline{HWC} + h \times \underline{WC} + w \times \underline{C} + \underline{C+1}$$

strides = (HWC, WC, C, 1)





handle

template 类

template <typename T, typename traits = handle\_traits(T)>

primary traits 类

包装一个指针 & handle\_trait<T>

对于不同的 T, 可得到 specialization

template <>

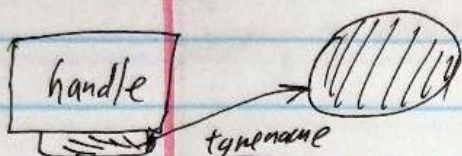
struct handle<dnnl\_memory\_t> {

static destructor(dnnl\_memory\_t,

dnnl\_memory\_destroy(p);

};

traits 类, 提供 destructor

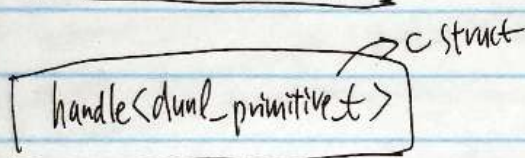


shared\_ptr<std::remove\_pointer(T)::type> data;

handle<dnnl\_memory\_t> hdl(一个具体 T 类型, weak?);

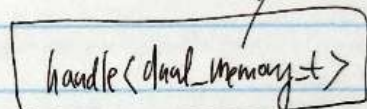
typedef struct dnnl\_memory\* dnnl\_memory\_t;

属于 shared\_ptr

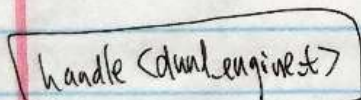


primitive

execute(const stream & astream, const std::unordered\_map<int, memory> & args) const;

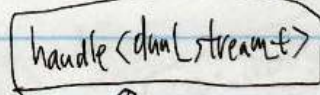


C++ class

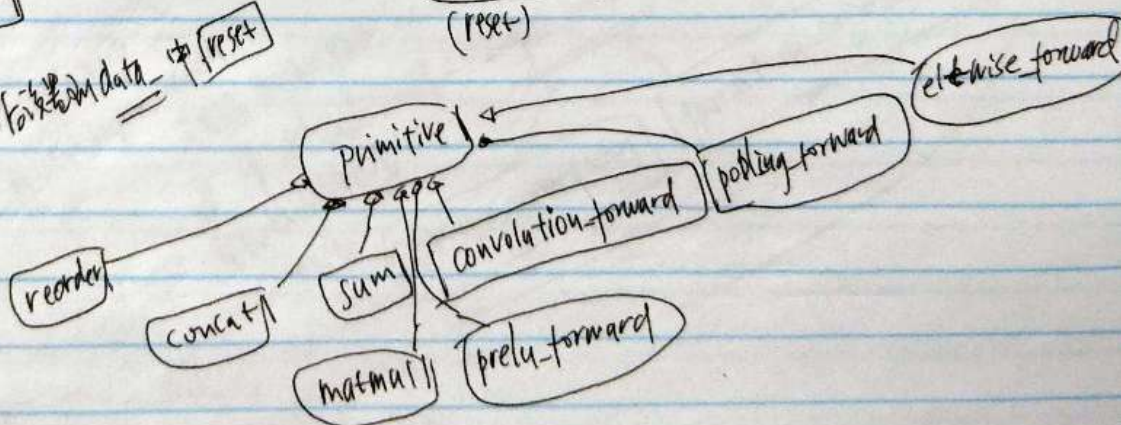


engine (reset)

技术  
调用 C-API  
封装 C-struct, 8 个 data 中 reset

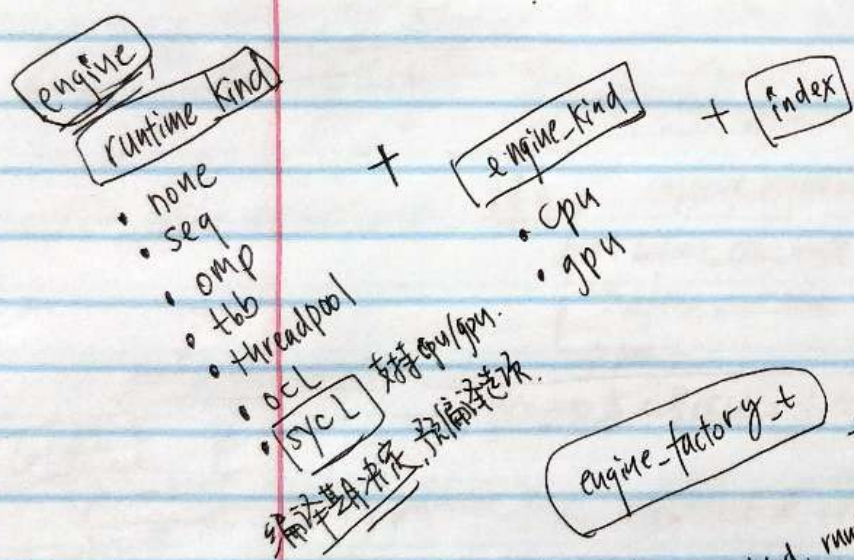
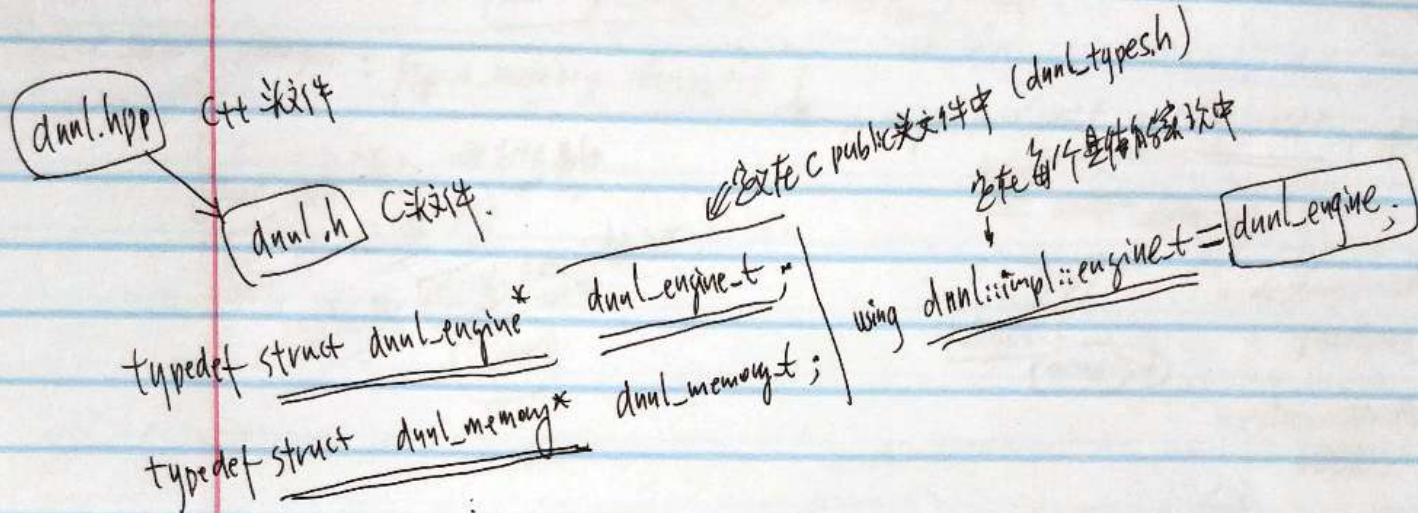
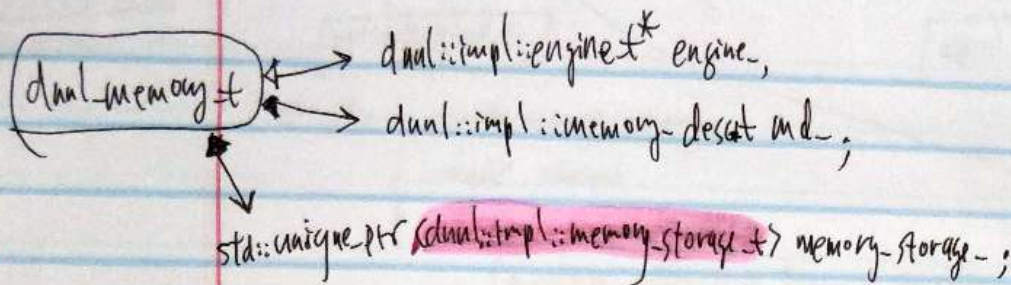


stream (reset)



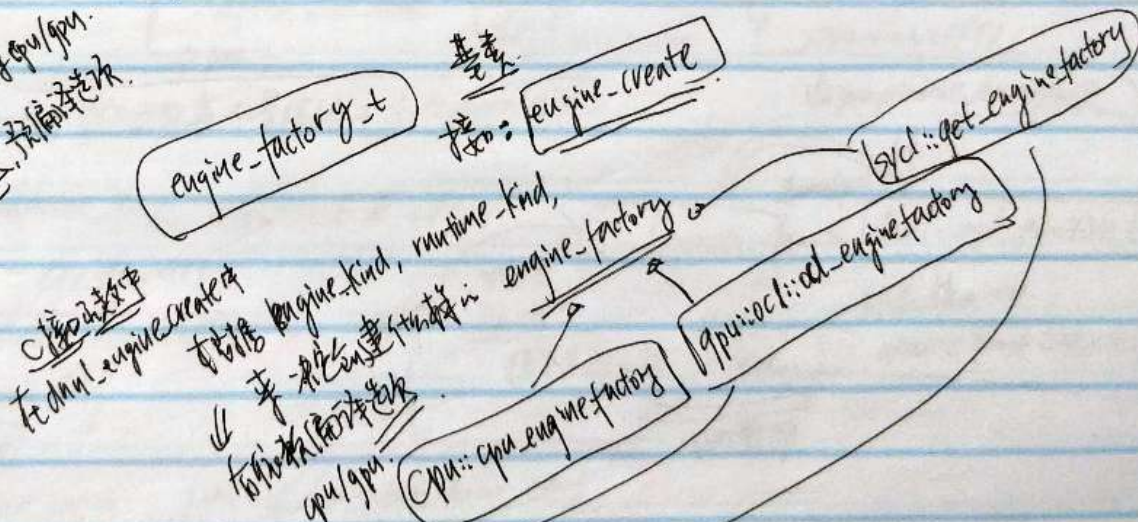


struct dnnl\_memory\_t \* = dnnl\_memory\_t;



**engine:** base 类

- 提供 memory allocation
- 提供 stream create



代码结构:

src / common: 通用基类..

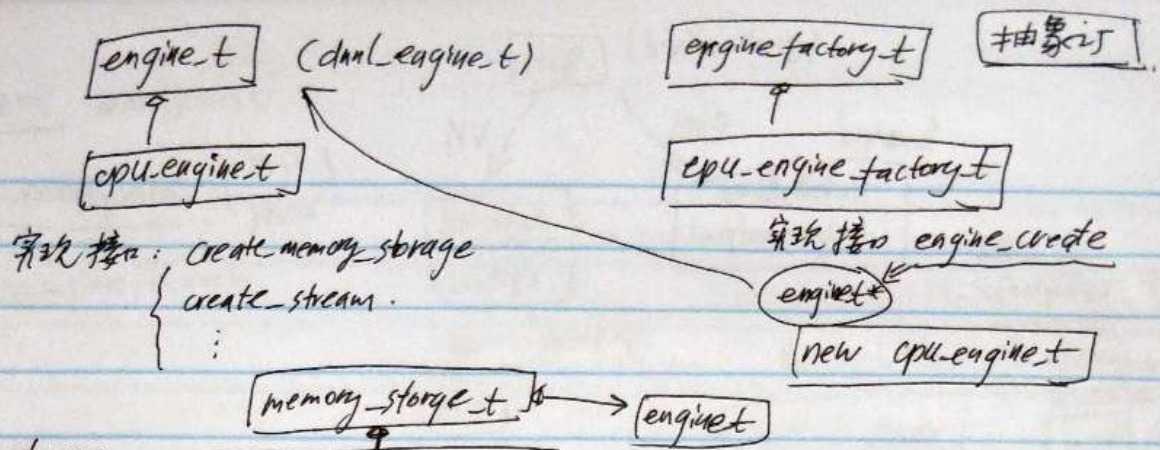
cpu: cpu-specific subclass

gpu: gpu-specific subclass

sycl: .....



cpu  
dnnl/impl/cpu  
namespace



create\_memory\_storage: cpu\_memory\_storage\_t

unique\_ptr<void, void (\*)(void\*)> data;

void\* data[...];

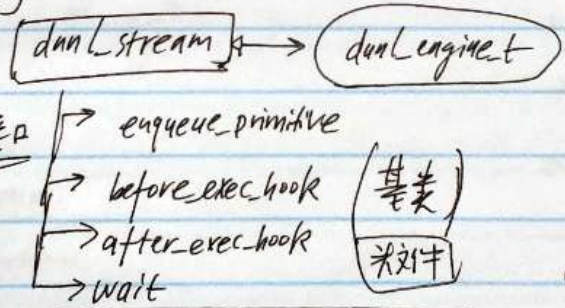
malloc, memalign

destroy(void\* ptr) { free(ptr); }

release(void\* ptr) { none; }

either 请求分配 => free 释放  
or 传入 handle (ptr) 传入 handle release data\_idelctor do nothing

generic 通用 primitive / stream



using dnnl::impl::stream\_t = dnnl::stream

实现文件

实现 C-API

dnnl::stream::wait(stream\_t\* stream)

stream->wait(1);

由 primitive\_desc 创建

dnnl::primitive\_execute

G-API

primitive\_execute: (primitive\_iface, exec\_ctx\_t & ctx) {

auto stream = ctx.stream();

stream->before\_exec\_hook();

if (get\_verbose()) {

profile

stream->enqueue\_primitive(primitive\_iface, ctx);

print profile runtime.

} else {

stream->enqueue\_primitive(primitive\_iface, ctx);

}

wrap

shared\_ptr

dnnl::primitive

dnnl::impl::primitive\_t

相同 execute 但不是虚函数

提供接口 = 0. execute(const exec\_ctx\_t & ctx)

转调用

=> dnnl::primitive::execute(exec\_ctx\_t & ctx)