

sketch2sky

What I Cannot Create, I Do Not Understand —Richard Feynman And I



≡ Primary Menu

Tensorflow XLA Service 详解 II

🔗 1190 👤 [Jiang XIAO](#)

📅 2019年9月26日 at pm1:55 (last edited 📅 2020年6月20日 at am9:34)

本文主要介绍在XLA service阶段针对HloInstruction做的一些显存优化, 对于训练框架来说, 显存优化的工作至关重要, 主要是由于现阶段GPU+CUDA远没有CPU+Linux组合强大, 后者有完善的建立在虚拟内存基础上的内存管理机制, 内存的高效使用由linux kernel来负责, 即便物理内存不足, 还可以使用swap, 内存压缩等技术确保内存的高效供应, 而在GPU+CUDA里, 这方面的工作很大程度让渡给了程序员自己来搞定, GPU程序接触到的就是物理显存, 如果程序的显存申请超过显存容量, 整个程序就会直接coredump, 此外, 显存本身就集成在GPU板卡上, 无法像内存一样扩展, 而GPU本身造价昂贵, 最后, 在深度学习训练中, 大力出奇迹的现状下, 显存的消耗明显超过的摩尔定律, 这也加剧了显存供求关系的矛盾, 正式由于训练框架做了大量的优化, 才能让模型跑起来.

XLA Service的显存优化设计思想与tensorflow整体一样遵循“静态图”的设计: 先整体优化, 再落地实施. 其中, xla/service/buffer_assignment.cc 是整个显存优化的核心, 在1.14版本中, xla/service/支持两种后端: cpu和gpu, 纷纷针对两种backend有进一步的优化算法, 本文主要针对GPU的优化逻辑进行分析

核心文件

内存优化公共:

xla/service/buffer_assignment 内存优化核心文件

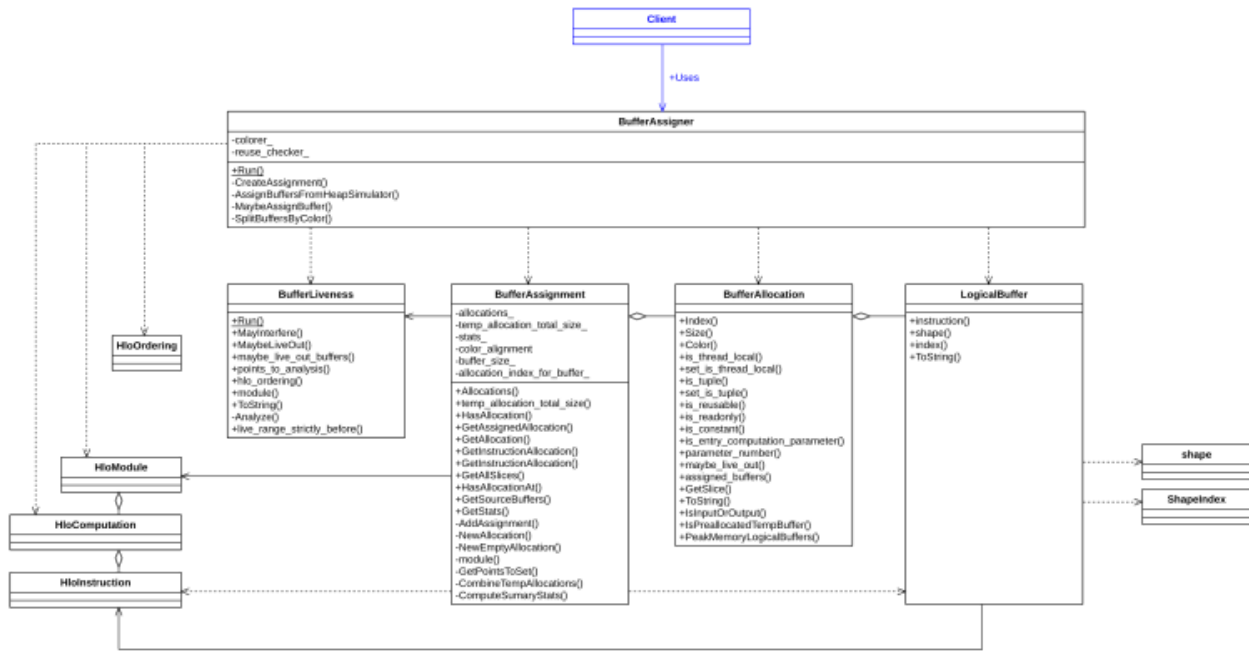
xla/service/buffer_liveness.cc 内存片生命周期分析

GPU相关:

xla/service/gpu/buffer_allocations.cc BufferAllocation的组合

xla/service/gpu/gpu_hlo_schedule.cc Hlo的处理顺序, 与显存的优化策略息息相关, 简单地说, 按照BFS并行执行的HloInstruction消耗的显存肯定大于所有的HloInstruction都顺序执行.

核心抽象



HloSchedule

XxxSchedule是TF的代码风格, 类似的有更底层用于Thunk调度的ThunkSchedule, 以及Service提供的HloSchedule. HloSchedule中最重要的就是封装了HloOrdering.

HloSchedule

XLAService内存优化的本质是处理LogicalBuffer和BufferAllocation之间的关系, 原则是使用尽可能少的BufferAllocation去承载尽可能多的LogicalBuffer, 而如何使用的更少, 就涉及到了对Hlo图的分析, 就涉及到了Ordering的问题, 使用不同策略生成Ordering, 直接影响两个LogicalBuffer之间的约束关系, 最简单的, 在图遍历中, 使用DFS和BFS的2种方式遍历会导致图上节点的内存依赖关系大有不同.

HloOrdering是描述HloInstruction加载序列的基类, 派生类有

PredecessorHloOrdering, DependencyHloOrdering 和 SequentialHloOrdering, 其中, DependencyHloOrdering基于依赖关系, 所以可以并行, 性能更高, 但耗更多的内存, 而SequentialHloOrdering完全串行, 性能相对低, 但可以节约更多内存, 而 PredecessorHloOrdering 是个虚类, 需要子类进一步填充predecessors_, 这也是GPU后端使用的方式.不同的Ordering会影响内存的依赖关系, 进一步影响Launch到GPU后Kernel的执行序列.

```

// An HLO ordering based on data dependencies in the HLO graph. In this partial
// order, instruction A executes before instruction B only if there is a path
// from A to B in the HLO graph. For example, given the following graph:
/*
      param
     /   \
  negate  exp
   \     /
    add
*/
// DependencyHloOrdering gives the following executes-before relations:
// param executes before negate, exp, and add
// negate executes before add
// exp executes before add
// add executes before nothing

```

```

// negate and exp are not ordered because the dependencies allow either to
// execute before the other (or in parallel). DependencyHloOrdering ordering
// allows maximum parallelism and enables any execution order which satisfies
// data dependencies. This requires pessimistic assumptions about buffer live
// ranges and can result in more memory used than more constrained orderings.
class DependencyHloOrdering : public PredecessorHloOrdering {

// An HLO ordering based on a total order of instructions in each computation.
// The computation total order is a sequencing of all of its instructions in
// the computation (eg, {inst0, inst1, inst2,...}) as in single-threaded
// execution. For example, given the following HLO graph:
/*
      param
     /   \
  negate exp
   \     /
    add

*/
// and the following sequence:
//
// {param, negate, exp, add}
//
// SequentialHloOrdering gives the following executes-before relations:
//   param executes before negate, exp, and add
//   negate executes before exp and add
//   exp executes before add
//   add executes before nothing
// This is more constrained than DependencyHloOrdering in this example because
// negate and exp are ordered (negate before exp). This enables param to share
// the same buffer as exp (param buffer is dead after exp). Generally, this
// ordering enables more buffer sharing (reduced memory usage) because buffer
// interference is reduced relative to DependencyHloOrdering.
class SequentialHloOrdering : public HloOrdering {

```

GpuHloSchedule

Build()	构造GpuHloSchedule单例, 同时会根据Stream的数量采用BFS或Sequence的LaunchOrder,, 这个LaunchOrder会用于构造GpuHloOrdering, 以及thunk_launch_order_
ThunkLaunchOrder()	拿到上述thunk_launch_order_
ConsumeHloOrdering()	拿到上述GpuHloOrdering

GpuHloOrdering

```

class GpuHloOrdering : public PredecessorHloOrdering : public HloOrdering

```

在使用Multi-stream的情况下, 一方面stream no的分配算法会遍历图给每个HloInstruction分配不同的stream no, 另一方面, GpuHloOrdering也会据此选择BFS算法来生成Ordering存入thunk_launch_order_以及HloOrdering.predecessors_, 后续的显存优化会据此优化显存, 确保Multi-Stream在实际执行的过程中不会产生显存问题.

BufferAssigner

BufferAssigner 用于构造BufferAssignment对象, 可以理解为BufferAssignment的wrapper

BufferAssignment

内存管理入口类. 所有关于内存优化的内容均在其Run()中实现, 在nvptx_compiler.cc:RunBackend()调用

BufferLiveness

描述一块内存的生命周期, “computes liveness of the output buffers of HLOs and their interference”, 为显存优化的决策提供支持. BufferLiveness 的构造过程经历3层分析, 从内而外, 即从前到后:

LogicalBufferAnalysis: 构造LogicalBuffer, 每个Instruction都可以有多个LogicalBuffer, 通过Accept机制遍历所有的Instruction

TuplePointsToAnalysis: 构造PointToSet, 需要之前的LogccalBuffer, 最终也是通过Accept机制遍历所有的Instruction二者其实的都是Hlo的DFS遍历handler

```
class TuplePointsToAnalysis : public DfsHloVisitorWithDefault
class LogicalBufferAnalysis : public DfsHloVisitorWithDefault
```

```
1. NVPTXCompiler::RunBackend()
2.   BufferAssigner::Run()
3.     assigner::CreateAssignment()
4.       liveness = BufferLiveness::Run(module, std::move(hlo_ordering)) //class BufferL
5.       liveness = new BufferLiveness()
6.       liveness->Analyze()
7.         points_to_analysis_ = TuplePointsToAnalysis::Run()
8.         logical_buffer_analysis = LogicalBufferAnalysis::Run()
9.         analysis = new LogicalBufferAnalysis()
10.        analysis.Analyze()
11.        return analysis
12.        analysis = new TuplePointsToAnalysis(logical_buffer_analysis)
13.        analysis.Analyze()
14.        return analysis
15.        maybe_live_out_buffers_ = points_to_analysis_->GetPointsToSet(root).CreateF
16.        return liveness
17.        assignment(new BufferAssignment(module, std::move(liveness))
```

-2- 规划内存的总入口, 注意入参

-15- liveout的部分只发生在root instruction(TODO)

LogicalBuffer

“虚拟内存”, LogicalBuffer 对内存需求的抽象, 内存优化的实质就是调整LogicalBuffer和BufferAllocation的关系, 将来会被HloValue替换掉. 同时, 整个内存优化过程都使用color的概念, 相同的color表示可以融合, 这点在BufferAllocation和BufferSet中都有用到, 显然, 处于一个BufferAllocation之下的LogicalBuffer都有相同的color, 要不会不会在一起, 相应的BufferAllocation实例也使用想用的color进行标记

什么样的两个LogicalBuffer会被分配给同一个BufferAllocation?

1. kWhile, kCall, kConditional 的LogicalBuffer

2. `cannot_merge_buffer_sets()`: !(两个buffer的id不同(不是一个buffer) && 会发生干涉)

```

1. void LogicalBufferAnalysis::NewLogicalBuffer(HloInstruction* instruction, const Shape
2. CHECK_EQ(logical_buffers_.size(), next_buffer_id_);
3. logical_buffers_.emplace_back( absl::make_unique<LogicalBuffer>(instruction, index,
4. //LogicalBuffer
5. : BufferValue(instruction, index, id),
6. id_(id)
7. const Shape& shape = ShapeUtil::GetSubshape(instruction->shape(), index);
8. index_(index) {}
9. output_buffers_[std::make_pair(instruction, index)] = logical_buffers_.back().get()
10. ++next_buffer_id_;
11. }

```

- 1- LogicalBuffer只存Shape(确切的是SubShape) id 和inst指针
- 3- 所有的LogicalBuffers对象都会丢到, LogicalBufferAnalysis::logical_buffers中
- 7- 所以, 一个LogicalBuffer只对应一个Subshape的内存, 不是整个Instruction

BufferAllocation

BufferAllocation 物理内存的抽象, 一个 BufferAllocation 实例最终对应一块连续物理内存, 通常, 一块物理内存会被多个不干涉的LogicalBuffer公用. 一个BufferAllocation实例可以分为:

```

Status BufferAssignment::ComputeSummaryStats() {
  for (auto& allocation : Allocations()) {
    if (allocation.is_entry_computation_parameter()) {
      stats_.parameter_allocation_count++;
      stats_.parameter_allocation_bytes += allocation.size();
    }
    if (allocation.is_constant()) {
      stats_.constant_allocation_count++;
      stats_.constant_allocation_bytes += allocation.size();
    }
    if (allocation.maybe_live_out()) {
      stats_.maybe_live_out_allocation_count++;
      stats_.maybe_live_out_allocation_bytes += allocation.size();
    }
    if (allocation.IsPreallocatedTempBuffer()) {
      stats_.preallocated_temp_allocation_count++;
      stats_.preallocated_temp_allocation_bytes += allocation.size();
    }
    stats_.total_allocation_count++;
    stats_.total_allocation_bytes += allocation.size();
  }
}

```

BufferAllocation::Slice

BufferAllocation里的一片, 关联在该BufferAllocation的每一个LogicalBuffer都会有对应的Slice

GetSlice 根据 buffer, 重新构造一个 Slice 对象

```

1. BufferAllocation::Slice BufferAllocation::GetSlice(
2. const LogicalBuffer& buffer) const {
3. const OffsetSize os = FindOrDie(assigned_buffers_, &buffer);
4. return Slice(this, os.offset, os.size);
5. }

```

ShapeIndex

一个 ShapeIndex 对应的是一个 LogicalBufferList, 真正使用的时候是用的 ShapeIndex 来索引到buffer

```
1.   StatusOr<BufferAllocation::Slice> BufferAssignment::GetUniqueSlice
2.   for (const LogicalBuffer* buffer : GetPointsToSet(instruction).element(index))
3.       //elements()
4.   return tree_.element(index).buffers
```

-4- index是ShapeIndex, 描述一个Shape的索引, 不是Intuction. buffers 是 BufferList, 此处可以看到, 直接用index 来索引得到的BufferList,

```
1.   IrEmitterUnnested::BuildInitializerThunk()
2.   return {absl::make_unique<MemzeroThunk>(a, nullptr)};
```

-2- 使用 ShapeIndex 索引 直接->Allocation::Slice->Thunk,

BufferAllocations

GPU后端对一组BufferAllocation的抽象

```
1.   std::vector<se::DeviceMemoryBase> buffers_
```

-1- 物理内存对象, se::DeviceMemoryBase是StreamExecutor对一块线性内存资源的抽象, 提供了最基础 offset+size的计算方法

```
1.   StatusOr<std::unique_ptr<BufferAllocations>> BufferAllocations::Builder::Build()
2.   for BufferAllocation::Index i in num_buffers:
3.       BufferAllocation allocation = buffer_assignment->GetAllocation():
4.       if address already registered:
5.           buffer_allocations->SetBuffer(i, *address);
6.       if (allocation.maybe_live_out() || allocation.IsPreallocatedTempBuffer()):
7.           se::DeviceMemoryBase buffer_address;
8.           buffer = memory_allocator->Allocate(device_ordinal, buffer_size)
9.           buffer_address = buffer.Release();
10.          buffer_allocations->SetBuffer(i, buffer_address);
```

-2- 每个 buffer 都有一个 BufferAllocation

-8- 真正分配新的buffer

优化过程

```
2   NVPTXCompiler::RunBackend()
3   hlo_schedule = GpuHloSchedule::Build(*module, *stream_assignment, pointer_size_)
4   std::unique_ptr<GpuHloSchedule> schedule(new GpuHloSchedule)
5   //this analysis figures out which temp buffers are required to run the computation
6   BufferAssigner::Run(hlo_schedule->ConsumeHloOrdering(), BufferSizeBytesFunction(), B
7   //BufferSizeBytesFunction()
8   shape_size = ShapeSizeBytesFunction()
9   ShapeUtil::ByteSizeOf(shape, pointer_size);
10  ByteSizeOfTupleIndexTable(shape, pointer_size);
11  return pointer_size * shape.tuple_shapes_size();
12  //Run()
```

```

13 BufferAssigner assigner
14 assigner.CreateAssignment(HloModule, hlo_ordering, buffer_size)
15   liveness = BufferLiveness::Run(module, std::move(hlo_ordering)
16     liveness = new BufferLiveness()
17     liveness->Analyze()
18     points_to_analysis_ = TuplePointsToAnalysis::Run()
19     logical_buffer_analysis = LogicalBufferAnalysis::Run(module)
20     std::unique_ptr<LogicalBufferAnalysis> analysis(new LogicalBufferAnalysis
21     analysis->Analyze();
22     NewLogicalBuffer()
23     analysis(new TuplePointsToAnalysis(module, logical_buffer_analysis.Consum
24     analysis->Analyze()
25     for computation in module_->MakeNoFusionComputations():
26       computation->Accept(this)
27       for instruction in computation->instructions():
28         if instruction is not kFusion:
29           continue
30         GatherFusionInstructions(instruction, fusion_instructions)
31       for instruction in fusion_instructions:
32         instruction->fused_expression_root()->Accept(this)
33     return analysis
34   for computation in module->computation():
35     if computation->IsFusionComputation():
36       continue
37     for instruction in computation->instructions():
38       for alias_buffer in points_to_analysis_->GetPointsToSet().CreateFlatten
39       if aliased_buffer->instruction() != instruction:
40         alias_buffers_.insert(aliased_buffer)
41     if computation == module_->entry_computation():
42       maybe_live_out_buffers_ = points_to_analysis_->getPointsToSet().Creat
43   return liveness
44 assignment(new BufferAssignment(module, std::move(liveness)
45
46 flat_hash_set<const LogicalBuffer*> colocated_buffers;
47 BuildColocatedBufferSets(&colocated_buffer_sets)
48 points_to_analysis = buffer_liveness.points_to_analysis()
49 module->input_output_alias_config().ForEachAlias([]{AddSetToColocatedBufferSets
50 for computation in module->MakeComputationPostOrder():
51   for instruction in computation->MakeInstructionPostOrder():
52     if instruction->opcode() == HloOpCode::kWhile:
53       AddBufferToColocatedSet()
54       AddSetToColocatedBufferSets()
55     else if == HloOpCode::kCall:
56       //...
57     else if == HloOpCode::kConditional:
58       //...
59   // Given a list of colocated buffer sets (each colocated buffer set represents
60   // the logical buffers that would be assigned to the same physical buffer),
61   // try to merge the sets if the buffers can be shared. Returns the merged set.
62   new_colocated_buffer_sets = MergeColocatedBufferSets()
63   // Given the interference map of a graph (the list of interfering node indices
64   // for each node), perform graph coloring such that interfering nodes are
65   // assigned to different colors. Returns the assigned color of the nodes, where
66   // the colors are represented as integer values [0, color_count).
67   std::vector<int64> ColorInterferenceGraph()
68   // Assign a color to each colocation set in colocated_buffer_sets, such that
69   // the sets that can be merged are assigned with the same color.
70   auto assigned_colors = ColorInterferenceGraph(interference_map);
71   for node in nodes: //nodes: 0 ~ node_count-1
72     for neighbor in interference_map[node]:

```



```

73         color = assigned_colors[neighbor]
74         if color != kColorUnassigned:
75             available_colors[color] = false
76         ...
77     std::vector<CollocatedBufferSet> new_collocated_buffer_sets(num_sets)
78     for i in colocated_buffer_sets.size():
79         buffer_sets = colocated_buffer_sets[i]
80         new_collocated_buffer_sets[assigned_colors[i]].insert(buffer_set.begin(), b
81     return new_collocated_buffer_sets;
82     swap(collocated_buffer_sets, new_collocated_buffer_sets)
83     colorer_(assignment->liveness())
84     GatherComputationsByAllocationType(module, &thread_local_computations, &global_co
85     // For each buffer set in 'collocated_buffer_sets', assigns all buffers in the
86     // same set to the same buffer allocation in 'assignment'.
87     AssignCollocatedBufferSets(collocated_buffer_sets, &collocated_buffers, &collocated_
88     for colocated_buffer_set in colocated_buffer_sets:
89         BufferAllocation* allocation = nullptr
90         for LogicalBuffer buffer in colocated_buffer_set:
91             instruction->buffer->instruction()
92             computation = instruction->parent()
93         for LogicalBuffer buffer in colocated_buffer_set:
94             buffer_size = assignment->buffer_size_(*buffer);
95             if allocation == nullptr:
96                 allocation = assignment->NewAllocation(*buffer, buffer_size)
97                 colocated_allocations->insert(allocation->index())
98             else:
99                 assignment->AddAssignment(allocation, *buffer, 0, buffer_size)
100                 allocation->AddAssignment(buffer, offset, size);
101                 assigned_buffers_.emplace(&buffer, offset_size);
102                 colocated_buffers->insert(buffer)
103     GatherComputationsByAllocationType(module, &thread_local_computations, &global_co
104     // First assign buffers for global computations. Temporary buffers for
105     // sequential computations are collected in 'buffers_to_assign_sequentially'.
106     flat_hash_map<const HloComputation*, flat_hash_set<const LogicalBuffer*>> buffers_
107     for computation in global_computations:
108         // Assigns buffers to the instructions in the given computation. "assignment"
109         // is modified to reflect the new buffer assignments. If is_thread_local is
110         // true, then all assigned buffers have the is_thread_local flag set to
111         // true.
112     AssignBuffersForComputation(collocated_buffers, buffers_to_assign_sequentially,
113     std::vector<const LogicalBuffer*> sorted_buffers;
114     for (auto* instruction : computation->instructions()):
115         for buffer in assignment->points_to_analysis().GetBuffersDefinedByInstruct
116             sorted_buffers.push_back()
117     for instruction : computation->MakeInstructionPostOrder():
118         post_order_position.emplace(instruction, position);
119     // Every sequential computation must get an entry in the
120     // buffers_to_assign_sequentially map, even if we end up with an empty set
121     // of buffers. This ensures we can correctly determine whether to run
122     // whole-module heap simulation.
123     buffers_to_assign_sequentially->emplace(computation, flat_hash_set<const Log
124     std::vector<BufferAllocation::Index> new_allocation_indices;
125     // A sorted multimap from size to indices of colocated allocations.
126     std::multimap<int64, BufferAllocation::Index> colocated_allocation_size_to_in
127     for index in colocated_allocations:
128         for buffer_offset_size in assignment->GetAllocations(index).assign_buffers
129             if IsTuple():
130                 consider_reuse = false
131                 break;
132     if considering_reusing:

```



```

133         sorted_colocated_indices.push(index)
134     while !sorted_colocated_indices.empty():
135         auto index = sorted_colocated_indices.top()
136         sorted_colocated_indices.pop()
137         // Returns the size of the allocation. Necessarily this must be at least as
138         // large as any LogicalBuffer assigned to this allocation.
139         colocated_allocation_size_to_indices.emplace(assigne...size(), index)
140     for buffer in sorted_buffers:
141         if colocated_buffers.contains(buffer)
142             continue
143         instruction = buffer->instruction()
144         // Function which returns the buffer size for a given logical buffer (shape
145         buffer_size = assignment->buffer_size_(*buffer)
146         if instruction->opcode() == HloOpcode::kConstant:
147             BufferAllocation* allocation = assignment->NewAllocation(*buffer)
148             ...
149             if ... ==
150             if ... ==
151             if ... ==
152         if buffer->...:
153             for operand in instruction->operands():
154                 for operand_slice in assignment->GetAllSlices():
155                     BufferAllocation allocation = assign.GetMutableAllocation(operand_sl
156                     if(MaybeAssignBuffer(allocation, buffer, assignment)):
157                         //bool BufferAssigner::MaybeAssignBuffer()
158                         ...
159                         assigned_operand = true
160                     if assigned_operand
161                         break
162         if reuse_colocated_checker_ ... && !assignment->HasAllocation(*buffer):
163             it = colocated_allocation_size_to_indices.lower_bound(buffer_size)
164             while(it != colocated_allocation_size_....)
165                 allocation = assignment->GetMutableAllocation(it->second)
166                 if (MaybeAssignBuffer(allocation)):
167                     colocated_allocation_size_to_indices.erase(it)
168                     break
169                 ++it;
170         if(!assignment->HasAllocation(*buffer)):
171             // Find the smallest buffer which can be reused iterating from end of
172             // new_allocation_indices (smallest) to beginning (largest).
173             for allocation_index = new_allocation_indices.size() - 1; allocation_inde
174                 BufferAllocation* allocation = assignment->GetMutableAllocation(new_al
175                 if MaybeAssignBuffer(allocation, *buffer, assignment):
176                     break
177         if !assignment->HasAllocation(buffer) && has_sequential_order && !liveness
178             // There is a sequential instruction ordering, so we delay assignment of
179             // temp buffers until after the loop. We do this right before we decide t
180             // create a new allocation, to ensure we've exhausted all the buffer
181             // re-use cases above.
182             //
183             // Entry parameters and thread local buffers were already handled earlier
184             // in this loop iteration. See BufferAllocation::IsPreallocatedTempBuffer
185             // for the definition of temp buffers.
186             (*buffers_to_assign_sequentially)[computation].insert(buffer);
187             continue
188         if !assignment->HasAllocation()
189             allocation = assignment->NewAllocation(buffer, buffer_size)
190             new_allocation_indices.push_back(allocation->index())
191     // Assigns 'buffers_to_assign_sequentially' using heap simulation, assuming
192     // the HLO instructions will be executed in the sequential order given by

```

```

193 // assignment->liveness().hlo_ordering().SequentialOrder. If
194 // 'run_whole_module_heap_simulation' is true, the heap simulation will be run
195 // assuming all global computations are sequentially ordered.
196 // Assign buffers with sequential ordering, if any. If all global computations
197 // are sequential, we can run heap simulation on the whole module, which
198 // reduces memory usage.
199 AssignBuffersWithSequentialOrdering()
200   HloSchedule schedule(&assignment->module());
201   hlo_ordering = assignment->liveness().hlo_ordering();
202   auto get_heap_algorithm = ...
203   if (run_whole_module_heap_simulation):
204     for (const auto& pair : buffers_to_assign_sequentially):
205       HloInstructionSequence* instruction_sequence = hlo_ordering.SequentialOrder
206       schedule.set_sequence()
207       all_buffers_to_assign.insert()
208       color_map = SplitBuffersByColor()
209     for (auto& single_colored_set : color_map) :
210       result = HeapSimulator::Run()
211       AssignBuffersFromHeapSimulator()
212   else:
213     for (const auto& pair : buffers_to_assign_sequentially):
214       instruction_sequence = hlo_ordering.SequentialOrder(*computation)
215       color_map = SplitBuffersByColor(buffers_to_assign)
216     for (auto& single_colored_set : color_map):
217       buffer_value_set = ToBufferValueFlatSet(single_colored_set.second)
218       options.buffers_to_assign = &buffer_value_set;
219       result = HeapSimulator::Run(get_heap_algorithm(alignment),
220                                   assignment->points_to_analysis(),
221                                   assignment->buffer_size_)
222       HeapSimulator heap(std::move(algorithm), size_fn, options, &schedule)
223       const HloComputation* entry_computation = module.entry_computation();
224       HloInstructionSequence& instruction_sequence = schedule.sequence(entry_
225       heap.RunComputation(entry_computation, instruction_sequence, points_to_
226       for instruction in instruction_sequence.instructions():
227         points_to = points_to_analysis.GetPointsToSet(instruction);
228         for user in instruction.users():
229           ...
230       for instruction in instruction_sequence.instructions():
231         buffers_defined_by_instruction = points_analysis.GetBufferDefinedBy
232         for buffer in buffers_defined_by_instruction:
233           ...
234         for operand_buffer in used_buffers[instruction]:
235           ...
236         for buffer in buffers_defined_by_instruction:
237           if !shared: -->
238             alloc_size_by_instruction += size_fn(*buffer);
239             Alloc(buffer, instruction);
240       heap->finish()
241       Result result = algorithm->Finish();
242       AssignBuffersFromHeapSimulator(result)
243       BufferAllocation* allocation = assignment->NewEmptyAllocation(result.l
244       for (const auto& buffer_chunk : result.chunk_map):
245         assignment->AddAssignment(allocation, buffer, chunk.offset, chunk.si
246         allocation->AddAssignment(bufferm, offset, size)
247         allocation->peak_buffers_ = ComputePeakMemoryLogicalBuffers(*allocatio
248   for computation in thread_local_computation:
249     if computation.IsFusionComputation():
250       continue
251     AssignBuffersForComputation(/*is_thread_local=*/true)
252   for buffer in assignment_.liveness().maybe_live_out_buffers()

```

```

253         if assignment->HasAllocation():
254             alloc = assignment->GetMutableAssignedAllocation()
255             alloc -> set_maybe_live_out(true)
256         assignment->CombineTempAllocations();
257         std::partition(allocations_.begin(), allocations_.end(),
258             [](const BufferAllocation& allocation) {
259                 return !allocation.IsPreallocatedTempBuffer();
260             });
261         allocations_.erase(first_temp_it, )
262         for combined:combined_allocation_map:
263             allocations_.push_back()
264         for index in allocations_.size():
265             allocation->set_index(index):
266                 for buffer_offset_size in allocation->assigned_buffers_:
267                     LogicalBuffer buffer = buffer_offset_size.first
268                     allocation_index_for_buffer_[buffer] = index
269         ir_emitter_context();
270         IrEmitterUnnested ir_emitter
271         entry_computation->Accept(&ir_emitter)
272         ptx = CompileToPtx()
273         const std::vector<uint8> cubin = CompilePtxOrGetCachedResult()
274         ir_emitter.ConsumeThunkSequence(), std::move(stream_assignment), hlo_schedule->Thunk
275         auto* gpu_executable = new GpuExecutable()
276         gpu_executable->set_ir_module_string(ir_module_string_before_opt);

```

-2- 这个backend其实不是llvm的backend 而是HLO的backend?

-6- //buffer_assignment.cc , 都是对元数据的管理, 没有真的分配, 包含了所有的内存管理代码, 显存优化的入口, 这里ConsumerHloOrdering()返回的是hlo_ordering

-8- 就是下文的buffer_size

-11- 对于一个tuple描述的shape, size的计算

-15- 遍历 module//buffer_liveness.cc -> DFS analysis, 构造一个liveness, 需要针对HloModule的 TuplePointsToAnalysis::Run()->LogicalBufferAnalysis::Run()等一系列分析,

-19- 遍历 HloModule, 为每个HloInstruction构造相应的Logicalbuffer, 存入LogicalBufferAnalysis, 并返回其指针, `class LogicalBufferAnalysis : public DfsHloVisitorWithDefault`, 说明也是用来遍历的handler

-21- 遍历节点, 执行各种handler, 建立logicalbuffer和instruction及其subshape的关系, 并存储在logicalbuffer中

-22- 构造LogicalBuffer对象

-26- 所有非fusion的buffer构造完毕

-27- 这个循环收集所有的fusion_inst

-30- gather 该computation下的所有的fusion instruction

-32- 给fusion的inst分配logical buffer

-44- 构造BufferAssignment

-52- 当前只有这3种op可以colocated

-53- 将这个instruction对应的LogicalBuffer放到BufferSet

-54- 将BufferSet放到BufferSets

-60- colocate意味着最终分配的是同一块物理内存

-62- 使用邻接矩阵的方式表示一张图, 用int64标识一个节点. 遍历每一个node, 首先遍历其每一个neighbor, 这个neighbor不是相邻节点, 而是通过 cannot_merge_buffer_set 设置的, 比如id()不同&&... 他们用过的color就标记为不可用, 取一个neighbor都没用过的color, 最小未用, 给这个node, 最终的效果是, 每一个node和他的neighbors节点的颜色都不同, 但可能和其他节点相同, 相邻的不能共用buffer, 猜测应该是不能共用输入及输出, 前面的输出还是可以作为后面的输入的. 最终, 就是在这里解决内存冲突: 不相邻的标记为相同颜色, 说明后续可能可以合并. MergeColocatedBufferSets 遍历所有的BufferSet, 将能合并的set合并了, 合并之前, 每个BufferSet都和Instruction/node关联, 合并之后, BufferSets中的一个Set里的buffer就不是和instruction强相关了, 可能是几个合在

一起的, 不过, 不存在把一个BufferSet拆成几个的情况, 两个bufferSet, 要么能合并, 要么不能合并, 不会把一个BufferSet内的LogicalBuffer拆分.

-70- interference_map表示vector的vector, 表征的是colocated_buffer_sets中两个set的关系, 而从这里看, 每个set里存储的一组logicalbuffer其实都是一个node使用, 这一点其实也可以从colocated_buffer_sets的构造过程可以看出: 3种op每种op都会先添加一组上logical_buffer到colocated_buffer_set, 再把set放到sets中

-72- 将 interference_map中表征的有冲突的node用不同的color表示

-77- num_sets是最大的color值, 即能将nodes划分为互斥最小集合的数量

-80- 将不冲突的, 即color相同的node都放一起

-87- AssignColocatedBufferSets, 代码里的AssignXXX都是构造/分配BufferAllocation的意思, 比如这

个 AssignColocatedBufferSets, 就是根据BufferSets分配BufferAllocation, 遍历colocated_buffer_sets 里的每一个Set里的每一个LogicalBuffer, 将一个Set的LogicalBuffer都放在一个BufferAllocation中, 一个Allocation里的LogicalBuffer不一定一样大, 将二者关联的时候会传入buffer_size的, 此时, offset都还是0

-92- Module->computation->instruction的嵌套关系

-94- 根据assigner构造时传入的函数对象计算, □ 这里是遍历所有的buffersize, 将相同的合并到一个allocation中

-95- 为 buffer 分配allocation

-96- ****每个allocation内的buffer或buffer_size是固定的?****

-99- 将这个buffer 赋给已有BufferAllocation对象, 这里 offset是0, size就是buffer_size, 是第一次存储这两个值

-106- 第一个顺序执行的点

-107- 至此, 就完成了初始的内存分配代码, 接下来 AssignBufferForComputation(), 就是将 Computation都关联到Allocation, 首先是 global_computations, 遍历所有的Instruction->points_to_analysis-

>GetBuffersDefinedByInstruction(), 将所有的Buffer按照从大到小排序, 之所以从大到小排序, 应该

和 BufferAllocation的Size 要比里面的所有 Buffer的size 都大的原因, 这样就能充分利用 allocation, 毕竟越少越好, 前面的用大块 buffer 搞了大的, 后面就可以乘凉, NewAllocation 本身就会将自己链入Assignment, 总之, 这里会通过一系列判断, 将大部分buffer存入尽可能少的Allocation中. global_computations 进行Assign的过程中, 有一些buffer被加入了sequential的范围, 这部分buffer在AssignBuffersWithSequtialOrdering, 同样, 这里Assign也是找Allocation的意思, 这里会用到堆模拟器, 会计算很多offset, size之类的, 将buffer加入到Heap中, 'NewEmptyAllocation(result.heap_size,' 就能看出, 此外, 由于根据hlo_ordering, stream_no> 1时不会所有的都是顺序的, 所以, 会存在根据hlo_ordering最终有多个chunk的情况, 更具体的, 一个computation对应多个LogicalBuffer, 并进行一个simulate, 得到一组chunk, 和一个allocation.

-112- 有关于顺序执行的代码

-116- 将传入的computation下的每一个instruction的每一个LogicalBuffer都入队

-118- 保存传入的computation下每一个instruction及其位置

-130- 根据注释, // Output tuple table may be allocated at run-time, so make sure we don't // overwrite them.

-140- 处理在sorted_buffer中, 但不在 colocated_buffers 中的, 前文的 colocated_buffers 只有3种HloOpCode可以用

-155- 获取每一个已有的Allocation

-161- 说明该buffer已经找到了 allocation

-162- 说明之前没有找到能用的

-163- 找到最接近的大于的key

-172- 从小到大(即从后向前) 找合适的allocation存储buffer

-186- 不是liveout的buffer, 即temp的buffer, 先存着, 没有相应的Allocation, 这里存的是comutation和buffer们的映射关系

-189- 实在不行了, 分配个新的Allocation

-193- 计算一个buffer的 live out 或 interfere 等

-199- buffer 顺序执行, 重点关注

-202- 这里用的LazyBestFitHeap, ChoostBestHeapAlgorithm

-208- 默认构造一个buffer只有0, 通过default_colorer(), 还有其他逻辑导致可能没有0吗?-213- out/trace.log:2019-08-02 10:44:04.437956: I tensorflow/compiler/xla/service/buffer_assignment.cc: 1262] Running per-

computation heap simulation

-215- 这里的是LogicalBuffer->BufferValue里的color, 不是前文用来合并BufferSet的color, 是通过buffer_liveness.h 中的 DefaultColorer赋值的. 这部分buffers都有color, 但没有被colocate合并过的, 也没有被合并单奥已有allocation中的temp buffer

-216- single_colored_set里的buffer是可以合并的.

-237- 分析时候可以共享内存

-243- 用整个优化过的堆大小来定义新的allocation, 一个allocation里是可以存放很多不同offset和size的buffer的. allocation有自己的size, 全文来看, 同样的color分配一个allocation. 那直接取最大的不就好了。 所以最终, Assignment里存储的tmp buffer是最小数量的Allocation和他的size, 用color区分 + 之前3种已经处理好的Op, + 其他

-245- 会存入allocation

-256- 把多个allocation里的buffer存到同一个combined_allocation里

-257- 划分的原则

-270- XLA HLO Visitor机制

-271- 使用上述代码构造的buffer来构造Thunk

Related:

[Tensorflow XLA Service Buffer优化详解](#)

[Tensorflow XLA Service 详解 II](#)

[Tensorflow XLA Service 详解 I](#)

[Tensorflow XLA Client | HloModuleProto 详解](#)

[Tensorflow XlaOpKernel | tf2xla 机制详解](#)

[Tensorflow JIT 技术详解](#)

[Tensorflow JIT/XLA UML](#)

[Tensorflow OpKernel机制详解](#)

[Tensorflow Op机制详解](#)

[Tensorflow Optimization机制详解](#)

[Tensorflow 图计算引擎概述](#)

 技术  Tensorflow, XLA, 技术

[Tensorflow XLA Service 详解 I](#)

[Tensorflow-Horovod安装部署checklist](#)

One comment on “Tensorflow XLA Service 详解 II”

Pingback: [Tensorflow OptimizationPassRegistry机制详解 - sketch2sky](#)

Leave a Reply