[View code on Github](#)

\#

# Zero-DP Memory Optimization

This is an implementation of Zero-DP introduced in the paper [ZeRO: Memory Optimization Towards Training A Trillion Parameter Models](#),

It keeps shards of the optimizer state, gradients and parameters into multiple devices/nodes. It reduces the memory consumption to $\frac{(2+2+K)\Psi}{N_d}$ of the original model, where $\Psi$ is the number of parameters, $N_d$ is the number of shards, and $K$ is number of optimizer bytes per parameter. $2 + 2$ are the parameter and gradient memory assuming 16-bit precision; i.e. 2 bytes per parameter and gradient. $K = 12$ for Adam optimizer because it maintains a copy of parameters, and two moments per parameter in fp32.

The communication volume of Zero-DP is $\mathcal{O}(3\Psi)$. For comparison data-parallel training has a communication volume of $\mathcal{O}(2\Psi)$.

```python
32  import functools
33  from typing import List, Optional, Tuple
34
35  import torch
36  import torch.distributed as dist
37  from torch import nn
```

Although this is named `Zero3` , we have only implemented the Zero-DP part of it and not the Zero-R memory optimizations which target residual memory consumption. Out implementation supports training only a subset of parameters.

This implementation is inspired by [Fairscale FSDP](#).

[Here's a script to fine-tune](#) GPT NeoX using Zero-DP memory optimization.

## Zero3 Layer

Each layer of the model (or a combination of a few consecutive layers) should be wrapped in this module.

```python
40  class Zero3Layer(nn.Module):
```

Each shard keeps parameters in `chunk` list. The `chunk[0]` is for trainable parameters and `chunk[1]` is for fixed parameters.

```python
49      chunk: List[nn.Parameter]
```

This is the sizes of the chunks in `chunk` list.

```python
51      chunk_size: List[int]
```

The first chunk is for trainable parameters.

```python
53      TRAINING_PARAMS_IDX = 0
```

This is the list of parameters split into lists as trainable and fixed parameters.

```python
56      param_refs: List[List[nn.Parameter]]
```

CUDA stream to featch parameters

```python
59      fetch_stream: Optional[torch.cuda.Stream]
```

| | | |
|---|---|---|
| CUDA stream to backup/accumulate gradients | 61 | `backup_stream: Optional[torch.cuda.Stream]` |
| List of layers right before this layer | 63 | `prev_layer: List['Zero3Layer']` |
| List of layers right after this layer | 65 | `next_layer: List['Zero3Layer']` |
| The position of the current layer; used this for debugging logs | 67 | `layer_idx: int` |
| Whether parameters have been fetched | 70 | `is_fetched: bool` |
| Device of the layer | 73 | `device: torch.device` |
| Data type of the layer | 75 | `dtype: torch.dtype` |
| The module to be wrapped | 77 | `module: nn.Module` |
| Number of nodes/devices the data is sharded across | 79 | `world_size: int` |

- `module` The module to be wrapped.
- `rank` The rank of the current node.
- `world_size` The number of nodes/devices the data is sharded across.
- `device` The device of the layer.
- `dtype` The data type of the layer.

```
81    def __init__(self, module: nn.Module, rank: int, world_size: int, device: torch.device, dtype: torch.dtype):
```

```
89              super().__init__()
```

Initialize the properties

```
92              self.device = device
93              self.dtype = dtype
94              self.module = module
95              self.prev_layer = []
96              self.next_layer = []
97              self.is_fetched = False
98              self.world_size = world_size
99              self.layer_idx = -1
100             self.fetch_stream = None
101             self.backup_stream = None
102
103             with torch.no_grad():
```

Collect all the parameters of the layer

```
105                 all_param_refs = [p for p in self.parameters()]
```

Store the shape of the parameters because we need it later to reconstruct them

```
108                 for p in all_param_refs:
109                     p._orig_shape = p.shape
```

All parameters should have the same type

```
112                 for p in all_param_refs:
113                     assert p.dtype == dtype, "All parameters should have same dtype"
```

Separate parameters as trainable and fixed

```
116                 self.param_refs = [[p for p in all_param_refs if p.requires_grad],
117                                    [p for p in all_param_refs if not p.requires_grad]]
118                 del all_param_refs
```

The `rank = 0` node will calculate the size each device/node should store, and distribute the parameters accordingly.

```
122                 if rank == 0:
```

| Description | Line | Code |
| --- | --- | --- |
| Merge and pad trainable (`merged_params[0]`) and fixed (`merged_params[1]`) parameters | 124 | `merged_params = [self._merge_and_pad_params(ps) for ps in self.param_refs]` |
| Calculate the chunk sizes of trainable and fixed params | 126 | `self.chunk_size = [(len(p) // world_size if p is not None else 0) for p in merged_params]` |
| Broadcast the sizes | 128 | `dist.broadcast(torch.tensor(self.chunk_size, device=device), src=0)` |
| | 129 | `else:` |
| Create an empty tensor to receive the sizes | 131 | `chunk_size = torch.tensor([0, 0], device=device)` |
| Receive the sizes | 133 | `dist.broadcast(chunk_size, src=0)` |
| | 134 | `self.chunk_size = chunk_size.tolist()` |
| Create parameters for trainable (`self.chunk[0]`) and fixed (`self.chunk[1]`) parameters to be stored in current device/node | 138 | `self.chunk = [nn.Parameter(self._empty((s,)), requires_grad=i == self.TRAINING_PARAMS_IDX)` |
| | 139 | `for i, s in enumerate(self.chunk_size)]` |
| An empty tensor to receive the trainable and fixed parameters combined | 142 | `chunk = self._empty((sum(self.chunk_size),))` |
| | 143 | |
| | 144 | `if rank == 0:` |
| Concatenate both trainable and fixed params | 146 | `all_params = torch.cat([p.view(world_size, -1) for p in merged_params], dim=-1).view(-1)` |
| | 147 | `del merged_params` |
| Scatter them to all the nodes/devices | 150 | `dist.scatter(chunk, list(all_params.split(sum(self.chunk_size))))` |
| | 151 | `del all_params` |
| | 152 | `else:` |

**Receive the parameters**

```
154                    dist.scatter(chunk)
```

**Collect the chunk data**

```
157              chunk = chunk.split(self.chunk_size)
158              for i, c in enumerate(chunk):
159                  self.chunk[i].data[:] = c
160              del chunk
```

**Cleanup the normal parameters**

```
163              self._cleanup_params()
```

**Add a backward hook. This gets called when the gradients relative to the module are computed.**

```
166          self._backward_hook_ref = self.register_full_backward_hook(self._backward_h
ook)  # type: ignore
```

**Merge all the parameters and pad it so that it's divisible by** `world_size` **.**

```
168      def _merge_and_pad_params(self, params: List[nn.Parameter]) -> torch.Tensor:
```

**Total number of parameters**

```
173          size = sum(p.shape.numel() for p in params)
```

**If it is not divisible by** `world_size` **, pad it**

```
176          if size % self.world_size != 0:
177              padding_fixed = self.world_size - (size % self.world_size)
```

**Otherwise, no need to pad**

```
179          else:
180              padding_fixed = 0
```

**Create an empty padding tensor**

```
182          padding = self._empty((padding_fixed,))
```

**Concatenate all the parameters and pad it**

```
184          return torch.cat([p.view(-1) for p in params] + [padding], dim=0)
```

**Get trainable chunk/shard of the parameters.**

```
186      def get_trainable_chunk(self) -> List[nn.Parameter]:
```

This is what we pass on to the optimizer on the current node.

Return and empty list if there are no trainable parameters

```python
193        if len(self.chunk[self.TRAINING_PARAMS_IDX]) == 0:
194            return []
```

Return the trainable chunk as a list

```python
197        return [self.chunk[self.TRAINING_PARAMS_IDX]]
```

**Create an empty tensor of the given shape.**

```python
199    def _empty(self, shape: Tuple[int, ...]) -> torch.Tensor:
```

```python
203        return torch.empty(shape, device=self.device, dtype=self.dtype)
```

**Cleanup the parameter data**

This will release all the memory used by the layer parameters.

```python
205    @torch.no_grad()
206    def _cleanup_params(self):
```

Set the flag to indicate that the parameters are not fetched

```python
214        self.is_fetched = False
```

Iterate through all parameters

```python
217        for ps in self.param_refs:
218            for p in ps:
```

Wait for operations on the parameters to complete before any new operations

```python
220                p.data.record_stream(torch.cuda.current_stream())
```

Check to make sure the parameter is not sharing storage with anything else

```python
222                assert p.data.storage_offset() == 0, "The tensor is not the sole occupant of the storage."
```

```python
226                p.data.storage().resize_(0)  # This is what actually clears the memory
```

Resize the storage to $0$. This will release the memory used by the parameter.

**Setting `p.data` will not release the memory, since the autograd graph keeps a reference to it.**

Make sure the parameter has no gradient data

```
228                 assert p.grad is None, 'Gradients should be None'
```

## Fetch the parameters from all shards

This will fetch all the parameter data from all the nodes and rebuild the parameters on each node.

```
230         @torch.no_grad()
231         def fetch_params(self):
```

Skip is already fetched

```
239             if self.is_fetched:
240                 return
```

Set the flag

```
243             self.is_fetched = True
```

Skip if there's nothing to fetch or share.

```
246             if sum(self.chunk_size) == 0:
247                 return
```

Use `fetch_stream` to fetch the parameters from all the shards

```
250             with torch.cuda.stream(self.fetch_stream):
```

Create an empty tensor to receive the parameters

```
252                 buffer = self._empty((self.world_size * sum(self.chunk_size),))
```

Split the continuous buffer into the number of nodes. These splits are views of `buffer'.

```
254                 buffers = list(buffer.split(sum(self.chunk_size)))
```

Concatenate both trainable and fixed chunks

```
257                 chunk = torch.cat(self.chunk, dim=0)
```

Gather the parameters from all the nodes/devices

```
260        dist.all_gather(buffers, chunk)
```

Split the gathered parameters into the trainable and fixed chunks

```
263        params = buffer.view(-1, sum(self.chunk_size)).split(self.chunk_size, dim=
1)
```

Wait for the gather operation to complete and then clear the references to the buffers

```
265        buffer.record_stream(self.fetch_stream)
266        for b in buffers:
267            b.record_stream(self.fetch_stream)
268        buffer.record_stream(self.fetch_stream)
269        del buffer
270        del buffers
```

Reshape the trainable and fixed parameters to continuous tensors

```
273        params = [p.reshape(-1) for p in params]
```

Collect the individual parameter tensors

```
276        for cont, ps in zip(params, self.param_refs):
```

If there are no parameters, skip

```
278            if not ps:
279                continue
```

Offset of the continuous tensor

```
282            offset = 0
```

Iterate through model parameters and assign the values from the continuous tensor

```
284            for p in ps:
```

Original parameter shape

```
286                shape = p._orig_shape  # type: ignore[attr-defined]
```

Change the storage size of the parameter. This was set to 0 when we cleaned up the parameters.

```
288                p.data.storage().resize_(shape.numel())
```

Assign the values from the continuous tensor

```
290        p.data[:] = cont[offset: offset + shape.numel()].reshape(shape)
```

Wait for the operations to complete before other operations can be performed

```
292        p.data.record_stream(self.fetch_stream)
```

Update the offset

```
294        offset += shape.numel()
```

Wait for the operation to complete before other operations can be performed

```
297        cont.record_stream(self.fetch_stream)
```

```
300        del params
```

## Forward pass

```
302    def forward(self, *args, **kwargs):
```

Fetch all the parameters of the current node. This gets called by the previous layer so this call is just to make sure parameters are fetched.

```
309        self.fetch_params()
```

Wait for parameter fetching to complete.

```
312        torch.cuda.current_stream().wait_stream(self.fetch_stream)
```

Start fetching parameters of the proceeding layers, so that they will fetch them which the current layer does its computations.

```
316        for layer in self.next_layer:
317            layer.fetch_params()
```

Add backward hooks to the parameters of the current layer if autograd is enabled.

```
320        if torch.is_grad_enabled():
321            self._add_backward_hooks()
```

Compute the outputs of the current layer

```
324        res = self.module(*args, **kwargs)
```

Cleanup the parameters of the layer.

*Skip cleaning up if autograd is enabled and this is the last layer in the network, because we will need to fetch the parameters again for the backward pass.*

```python
330        if not torch.is_grad_enabled() or self.next_layer:
331            self._cleanup_params()
332
333        return res
```

**Add backward hooks to the parameters of the current layer.**

```python
335    def _add_backward_hooks(self):
```

Number of backward hooks added

```python
341        self._backward_hook_handles = 0
```

Loop through trainable parameters of the current layer

```python
344        for p in self.param_refs[self.TRAINING_PARAMS_IDX]:
```

Make sure a hook hasn't already been added

```python
346            assert not hasattr(p, "_hook_handle"), 'Parameter has already been hooked'
```

Use `expand_as` to create an autograd step which we can intercept

```python
348            p_tmp = p.expand_as(p)
```

Get a handle to add the backward hook. This blog discusses about `grad_acc`.

```python
351            grad_acc = p_tmp.grad_fn.next_functions[0][0]
```

Add the backward hook

```python
353            handle = grad_acc.register_hook(
354                functools.partial(self._post_backward_hook, p))
```

Keep a reference to the handle

```python
356            p._hook_handle = handle
```

Increment the number of hooks added

```python
358            self._backward_hook_handles += 1
```

### Handle a backward event

This gets called by parameter backward hooks and the module backward hook.

```
360    def _backward_event(self):
```

Decrement the hooks counter

```
368        self._backward_hook_handles -= 1
```

If all the hooks (including the module hook) have been called, then we can back up gradients and clean up the parameters.

```
372        if self._backward_hook_handles == -1:
373            self._backup_grads()
374            self._cleanup_params()
```

Start fetch parameters of the previous layer, because autograd will next process the gradients of it.

```
377        for layer in self.prev_layer:
378            layer.fetch_params()
```

### Parameter backward hook

```
380    def _post_backward_hook(self, p: nn.Parameter, *args):
```

Remove the handle from the parameter

```
385        p._hook_handle.remove()  # type: ignore[attr-defined]
386        delattr(p, "_hook_handle")
```

Handle a backward event

```
389        self._backward_event()
```

### Module backward hook

```
391    def _backward_hook(self, *args, **kwargs):
```

Handle a backward event

```
396        self._backward_event()
```

The previous layer will start computing gradients. We need to make sure it has finished fetching params.

```
399        torch.cuda.current_stream().wait_stream(self.fetch_stream)
```

```
                                                    402            return None


Backup the gradients of the current layer          404        @torch.no_grad()
                                                    405        def _backup_grads(self):


Skip if there are no trainable parameters          410            if self.chunk_size[self.TRAINING_PARAMS_IDX] == 0:
                                                    411                return


Use the backup stream to backup the gradients      414            with torch.cuda.stream(self.backup_stream):


Buffer to store the gradients                      416                buffer = self._empty((self.world_size * self.chunk_size[self.TRAINING_PARAM
                                                                   S_IDX],))


Split the continuous buffer into number of nodes.  418                buffers = list(buffer.split(self.chunk_size[self.TRAINING_PARAMS_IDX]))
These splits are views of `buffer'.


Offset of the continuous buffer                     421                offset = 0


Iterate through trainable parameters                423                for p in self.param_refs[self.TRAINING_PARAMS_IDX]:


Collect gradients                                   425                    shape = p._orig_shape  # type: ignore[attr-defined]
                                                    426                    buffer[offset: offset + shape.numel()] = p.grad.view(-1)


Update the offset                                   428                    offset += shape.numel()


Clean the gradients                                 430                    p.grad = None


Empty tensor to accumulate the gradients of the     433                grad = self._empty((self.chunk_size[self.TRAINING_PARAMS_IDX],))
current shard
```

Accumulate the gradients of each shard. It scatters the buffers across the nodes, and each node accumulates (reduces) the tensors it receives.

```
436                    dist.reduce_scatter(grad, buffers)
```

Wait for the operation to complete and then clear the references to the buffers

```
439                for b in buffers:
440                    b.record_stream(self.fetch_stream)
441                buffer.record_stream(self.fetch_stream)
442                del buffer
443                del buffers
```

Set the chunk gradients. This is what the optimizer sees.

```
446                self.chunk[self.TRAINING_PARAMS_IDX].grad = grad
447                del grad
```

# Sequential module for `Zero3Layer` layers

```
450  class Zero3Sequential(nn.Module):
```

- `modules` List of `Zero3Layer` layers

```
454      def __init__(self, modules: List[Zero3Layer]):
```

```
458          super().__init__()
```

CUDA stream to fetch parameters

```
461          self.fetch_stream = torch.cuda.Stream()
```

CUDA stream to back up (accumulate) gradients

```
463          self.backup_stream = torch.cuda.Stream()
```

Set the streams and preceding and proceeding layers for each `Zero3Layer` layer

```
466          for i in range(len(modules)):
```

| | | |
|---|---|---|
| Set layer index | 468 | `modules[i].layer_idx = i` |
| Set streams | 470 | `modules[i].fetch_stream = self.fetch_stream` |
| | 471 | `modules[i].backup_stream = self.backup_stream` |
| Set proceeding layers | 473 | `if i + 1 < len(modules):` |
| | 474 | `    modules[i].next_layer.append(modules[i + 1])` |
| Set preceding layers | 476 | `if i - 1 >= 0:` |
| | 477 | `    modules[i].prev_layer.append(modules[i - 1])` |
| Store list of modules | 480 | `self.module_list = nn.ModuleList(modules)` |
| | 482 | `def get_trainable_chunk(self):` |
| Return the list of trainable chunks from each layer | 484 | `return sum([m.get_trainable_chunk() for m in self.module_list], [])` |
| | 486 | `def forward(self, x: torch.Tensor):` |
| Make sure gradient back up is complete | 488 | `torch.cuda.current_stream().wait_stream(self.backup_stream)` |
| Forward pass | 491 | `for m in self.module_list:` |
| | 492 | `    x = m(x)` |
| | 495 | `return x` |