

🦥 Unsloth is growing! Come join us :)



unsloth

Join our Discord

Documentation

Up to \$500K USD salary + bonus equity, health care benefits + other benefits, USA relocation etc! Complete some puzzles and earn points!

- We encourage you to use AI for coding! No experience or PhD / Masters needed - just get enough points for consideration!
- There are negative points for incorrect submissions. Read each criteria! Read [Submission](#) steps.

Role	Compensation	Role Description	Points Needed
Founding Engineer	\$400K to \$500K & equity	Help push Unsloth forward - bug fixes, core features, UI, kernels, nearly anything!	47
ML Engineer	\$250K to \$300K & equity	Help with FSDP2, Float8, Float4, kernels, Unsloth core and more!	32
ML Intern	up to \$150K py	Implementing specific features in Unsloth core. Can be remote.	18

1. [Convert nf4 to Triton](#) [Difficulty: Hard] [Max points: 14]
2. [Make QLoRA work with FSDP2](#) [Difficulty: Medium to Hard] [Max points: 12]
3. [Make torch.compile work without graph breaks for QLoRA](#) [Difficulty: Easy to Medium] [Max points: 9]
4. [Help solve 🦥 Unsloth issues!](#) [Difficulty: Varies] [Max points: 12]
5. [Memory Efficient Backprop](#) [Difficulty: Medium to Hard] [Max points: 10]
6. [Submission steps](#)

🦥 Who are we?

- 1.58bit DeepSeek R1 GGUFs [Tweet](#) and [HF Model Page](#)
- GRPO Llama 3.1 8B on a free Colab [Tweet](#)
- Gemma bug fixes [Tweet](#) and bug fixes for Llama 3, Phi 3, Qwen 2.5 [Details](#) Llama-fying Phi-4 [Details](#)
- Gradient accumulation bug fixes [Tweet](#) 4bit Dynamic Quantization [Details](#)
- Unsloth Gradient Checkpointing async offloads activations [Details](#)
- 30K Github Stars [Github](#) & 7 million monthly downloads on [Hugging Face](#)
- PyTorch conference [video](#) AI Engineer World's Fair [video](#) GPU / CUDA MODE [talk](#)

Clarifications:

1. We'll compensate you if we interview you but don't hire you
2. \$100-\$1000 bounties for Task 4
3. Submissions must be Apache-2 licensed
4. Task 4 involves solving Github issues for OSS Unsloth
5. No time limit: rolling basis
6. US based preferred

Hello! I attempted all the questions except the GitHub issues!

Comments:

✓ Used trainer and hugging face code as a base, so it's compatible with huggingface code.

1) **Convert NF4 to Triton:**

- Changed `rtol=0.01`, `atol=0.01`, and applied these settings only for `torch.float16` in options for T4 GPU.
- Implemented two versions of the kernel:
  - One with `absmax2` as blocks, which includes `tl.debug_barrier`. This version could not be compiled using `torch.compile`.
  - Another with `absmax` as blocks, which does not include `tl.debug_barrier`. This version works with `torch.compile` but is slightly slower.
- Both versions are provided below.

2) **Make QLoRA work with FSDP2:**

- Executed the code on Kaggle and displayed the implementation here.

3) **Completed.**

4) **GitHub issues skipped.**

5) **Memory-Efficient Backpropagation:**

- Improved efficiency by integrating `torch.compile`.
- For **LLaMA-1B**, the training loss aligns with the same loss function used.
- **GRPO** also matches and functions correctly. In GRPO, I included the clipping step as part of the loss and only returned the loss. In a real implementation, we should also return KL divergence, completion length, rewards, and additional factors if needed.

A lot of monkey patching, so 🚀🚀🚀🚀🚀🚀🚀🚀🚀🚀 lol

---

```
1 # Code to install Unsloth, Triton, Torch etc
2 %%capture
3 !pip install --no-deps bitsandbytes accelerate xformers==0.0.29 peft trl triton
4 !pip install --no-deps cut_cross_entropy unsloth_zoo
5 !pip install sentencepiece protobuf datasets huggingface_hub hf_transfer
6 !pip install --no-deps unsloth

1 # Helpful functions used through the entire notebook
2 import torch
3 import torch.nn as nn
4 from transformers import set_seed
5 import time
6 import inspect
7 import os
8 major_version, minor_version = torch.cuda.get_device_capability()
9 HAS_BFLOAT16 = (major_version >= 8)
10 from inspect import currentframe as _C, getframeinfo
11 _F = lambda c: getframeinfo(c).lineno # Gets line number
12 WARN = lambda x: print(f"\033[31m{x}\033[0m") # Red colored warnings
13
14 # https://stackoverflow.com/questions/18425225/getting-the-name-of-a-variable-as-a-string
15 def NAME(var):
16     callers_local_vars = inspect.currentframe().f_back.f_locals.items()
17     names = [var_name for var_name, var_val in callers_local_vars if var_val is var]
18     return names[0] if len(names) != 0 else ""
19
20 def assert_same(x, y, line, dtype):
21     assert(x.dtype == dtype)
22     try: torch.testing.assert_close(x, y, check_stride = True, rtol=0.01, atol=0.01)
23     except Exception as error:
24         raise RuntimeError(
25             f"Failed allclose at line [{line}]: {NAME(x)}, {NAME(y)}\n{str(error)}"
26         )
27
28 os.environ["HF_HUB_ENABLE_HF_TRANSFER"] = "1"
```

---

---

## ✓ A) Convert nf4 to Triton. [Difficulty: Hard] [Max points: 14]


1. Goal: Convert a nf4 quantized tensor into fp16 or bf16 into a *single* Triton kernel The double dequant of the `absmax` and weight forming must be done in 1 Triton kernel. Must work on Tesla T4.
2. Must be faster than Unsloth's `fast_dequantize` by 1.15x or more, and not use large intermediate memory buffers.
3. Must not use `torch.compile`, but can use `trace.enabled` to help on writing Triton kernels.
4. Good material: [Unsloth fast dequantize function](#), also [bitsandbytes dequantize blockwise](#)
5. Use `test_dequantize_function` to test your implementation.
6. No CUDA allowed. Custom CUDA inside of the Triton is allowed.
7. Watch Tim's videos on Youtube: [8-bit Optimizers](#)

```
1 from bitsandbytes.nn import Linear4bit
2 from transformers.activations import ACT2FN
3 from unsloth.kernels.utils import fast_dequantize
4 from peft.utils.integrations import dequantize_module_weight as peft_dequantize
5 def unsloth_dequantize(weight):
6     return fast_dequantize(weight.weight, weight.weight.quant_state)
7
8 def bnb_Linear4bit(hd, m, dtype = torch.float16):
9     return Linear4bit(
10         hd, m, bias = None,
11         compute_dtype = dtype,
12         compress_statistics = True,
13         quant_type = "nf4",
14     )
15
16 # [NEW] as at 18th Feb 2025
17 def assert_correct_bnb(weight, dtype):
18     assert(weight.weight.dtype == torch.uint8)
19     assert(weight.weight.quant_state.dtype == dtype)
20     assert(weight.weight.quant_state.absmax.dtype == torch.uint8)
21     assert(weight.weight.quant_state.code.dtype == torch.float32)
22     assert(weight.weight.quant_state.offset.dtype == torch.float32)
23     assert(weight.weight.quant_state.blocksize == 64)
24     assert(weight.weight.quant_state.state2.absmax.dtype == torch.float32)
```

```


25     assert(weight.weight.quant_state.state2.code.dtype == torch.float32)
26     assert(weight.weight.quant_state.state2.blocksize == 256)
27
28 class MLP(nn.Module):
29     def __init__(self, hd = 4096, m = 14336, dtype = torch.float16):
30         super().__init__()
31         self.gate_proj = bnb_Linear4bit(hd, m, dtype = dtype).to("cuda")
32         self.up_proj = bnb_Linear4bit(hd, m, dtype = dtype).to("cuda")
33         self.down_proj = bnb_Linear4bit(m, hd, dtype = dtype).to("cuda")
34         # [NEW] as at 18th Feb 2025
35         self.gate_proj.weight.quant_state.dtype = dtype
36         self.up_proj.weight.quant_state.dtype = dtype
37         self.down_proj.weight.quant_state.dtype = dtype
38         self.act_fn = ACT2FN["silu"]
39     def forward(self, x):
40         return self.down_proj(self.act_fn(self.gate_proj(x)) * self.up_proj(x))
41
42 def mlp_forward(X, mlp, fx):
43     up = X @ fx(mlp.up_proj).t()
44     gate = X @ fx(mlp.gate_proj).t()
45     h = mlp.act_fn(gate) * up
46     down = h @ fx(mlp.down_proj).t()
47     return down
48
49 def mlp_dequantize(X, mlp, fx):
50     a = fx(mlp.up_proj).t(); torch.cuda.synchronize()
51     b = fx(mlp.gate_proj).t(); torch.cuda.synchronize()
52     c = fx(mlp.down_proj).t(); torch.cuda.synchronize()
53     return a, b, c
54
55 def test_dequantize(dequantize_fx):
56     elapsed = 0
57     options = [
58         (2, 3333, 2048, 8192, 3407, torch.float16),
59         (5, 777, 1024, 4096, 3409, torch.float16),
60         (3, 2048, 4096, 14336, 3408, torch.float16),
61     ]
62     for (bsz, qlen, hd, m, seed, dt) in options:
63         set_seed(seed)
64         torch.set_default_dtype(torch.float32)
65         mlp = MLP(hd = hd, m = m, dtype = dt)
66         X = torch.randn((bsz, qlen, hd), device = "cuda", dtype = dt)
67         torch.cuda.synchronize()
68
69         # Warmup
70         for _ in range(2):
71             assert_same(mlp_forward(X, mlp, dequantize_fx), mlp(X), _F(_C()), dt)
72             # [NEW] as at 18th Feb 2025
73             assert_correct_bnb(mlp.up_proj, dt)
74             assert_correct_bnb(mlp.gate_proj, dt)
75             assert_correct_bnb(mlp.down_proj, dt)
76             a, b, c = mlp_dequantize(X, mlp, dequantize_fx)
77             A, B, C = mlp_dequantize(X, mlp, unsloth_dequantize)
78             assert_same(a, A, _F(_C()), dt)
79             assert_same(b, B, _F(_C()), dt)
80             assert_same(c, C, _F(_C()), dt)
81
82         # Benchmarking
83         torch.cuda.synchronize()
84         start = time.time()
85         for _ in range(1000): mlp_dequantize(X, mlp, dequantize_fx)
86         elapsed += time.time() - start
87     return elapsed


```

 <ipython-input-2-65cf205adda9>:3: UserWarning: WARNING: Unsloth should be imported before transformers to ensure all opt

Please restructure your imports with 'import unsloth' at the top of your file.

from unsloth.kernels.utils import fast\_dequantize

 Unsloth: Will patch your computer to enable 2x faster free finetuning.

 Unsloth Zoo will now patch everything to make training faster!

Write your Triton kernel below, and test it:

```

1 from triton import jit
2 import triton
3 import triton.language as tl
4 import pdb
5
6 @triton.jit
7 def _your_dequantize_nf4_kernel(#pointers
8                               absmax2_ptr, out2_ptr, absmax1_ptr, out1_ptr, weight_ptr, code2_ptr, code1_ptr, offset1

```

```

9             #Dimensions
10            blocksize1, blocksize2, num_elems_1, num_elems_2, shift1, shift2,
11            #meta param
12            BLOCK_SIZE: tl.constexpr, BLOCK_SIZE_2:tl.constexpr
13        ):
14
15    #init
16    pid = tl.program_id(axis=0)
17    block_start = pid * BLOCK_SIZE
18    offsets = block_start + tl.arange(0, BLOCK_SIZE)
19    mask_for_abs2 = offsets < num_elems_2
20    indexs = offsets >> shift2
21
22    #de_doublequant
23    absmx2_max = tl.load(absmax2_ptr + indexs, mask = mask_for_abs2, other = 128)
24    absmx1_values = tl.load(absmax1_ptr + offsets, mask = mask_for_abs2, other = 128)
25    x = tl.load(code2_ptr + tl.cast(absmax1_values, tl.int32), mask = mask_for_abs2, other = 128)
26    val_offset = tl.fma(x, absmx2_max, tl.load(offset1))
27
28    tl.store(out2_ptr + offsets, val_offset, mask = mask_for_abs2, eviction_policy='evict_last')
29    tl.debug_barrier()
30
31
32    new_offsets = pid * BLOCK_SIZE_2 + tl.arange(0, BLOCK_SIZE_2)
33
34    #indexs = new_offsets >> shift1
35    indexs = tl.inline_asm_elementwise(
36        asm="""
37        shr.b32 $0, $4, $8;
38        shr.b32 $1, $5, $9;
39        shr.b32 $2, $6, $10;
40        shr.b32 $3, $7, $11;
41        """,
42        constraints=(
43            "=r,=r,=r,=r,"
44            "r,r,r,r,r,r,r,r"),
45        args=[new_offsets, shift1],
46        dtype=(tl.int32),
47        is_pure=True,
48        pack=4,
49    )
50
51    mask_for_abs1 = new_offsets < num_elems_1
52    gathered_max = tl.load(out2_ptr + indexs, mask = mask_for_abs1, other = 128)
53    weights_values = tl.load(weight_ptr + new_offsets, mask = mask_for_abs1, other=128)
54
55    #high_bits = weights_values >> 4
56    high_bits, lowtem = tl.inline_asm_elementwise(
57        asm="""
58        {
59            // --- Unpack the 32-bit input into 4 bytes ---
60            .reg .b8 a<4>;
61            // The input packed value is in $0.
62            mov.b32 {a0, a1, a2, a3}, $8;
63
64            // --- Convert each 8-bit value to a 32-bit integer ---
65            cvt.u32.u8 $4, a0;
66            cvt.u32.u8 $5, a1;
67            cvt.u32.u8 $6, a2;
68            cvt.u32.u8 $7, a3;
69        }
70        // --- Compute high nibble for each (byte >> 4) ---
71        shr.b32 $0, $4, 4;
72        shr.b32 $1, $5, 4;
73        shr.b32 $2, $6, 4;
74        shr.b32 $3, $7, 4;
75
76        // --- Compute low nibble for each (byte & 0x0F) ---
77        and.b32 $4, $4, 0x0F;
78        and.b32 $5, $5, 0x0F;
79        and.b32 $6, $6, 0x0F;
80        and.b32 $7, $7, 0x0F;
81
82        """,
83        constraints=(
84            "=r,=r,=r,=r,=r,=r,=r,=r,"
85            "r"),
86        args=[weights_values],
87        dtype=(tl.int32, tl.int32),
88        is_pure=True,
89        pack=4,
90    )

```

```

91
92     x = tl.load(code1_ptr + high_bits, mask = mask_for_abs1, other=4)
93     hi_val = x * (gathered_max)
94
95     x = tl.load(code1_ptr + lowtem, mask = mask_for_abs1, other=4)
96     lo_val = x * (gathered_max)
97
98     #storing
99     out_hi2 = new_offsets * 2
100    out_lo2 = new_offsets * 2 + 1
101    tl.store(out1_ptr + out_hi2, hi_val, mask = mask_for_abs1, eviction_policy='evict_last')
102    tl.store(out1_ptr + out_lo2, lo_val, mask = mask_for_abs1, eviction_policy='evict_last')
103
104    #@torch.compile(fullgraph=True)
105    def _your_dequantize_nf4(weight, quant_state):
106        absmx2 = quant_state.state2.absmax
107        blocksize2 = quant_state.state2.blocksize
108        absmx1 = quant_state.absmax
109        out2 = torch.empty(absmx1.shape, dtype=quant_state.dtype, device=absmx1.device)
110        offset1 = quant_state.offset
111        out1 = torch.empty(quant_state.shape, dtype=quant_state.dtype, device=weight.device)
112        blocksize1 = quant_state.blocksize//2
113        num_elems_1 = out1.numel()
114        num_elems_2 = out2.numel()
115        code2 = quant_state.state2.code
116        code1 = quant_state.code
117        shift1 = blocksize1.bit_length() - 1
118        shift2 = blocksize2.bit_length() - 1
119
120        grid = lambda meta: (triton.cdiv(num_elems_2, meta['BLOCK_SIZE']),)
121
122        kernal = _your_dequantize_nf4_kernel[grid](
123            #wanna_be_pointers
124            absmx2, out2, absmx1, out1, weight, code2, code1, offset1,
125            #Dimensions
126            blocksize1, blocksize2, num_elems_1, num_elems_2, shift1, shift2,
127            #meta_param
128            BLOCK_SIZE = 64, BLOCK_SIZE_2 = 64*32
129        )
130
131        return out1
132
133    def your_dequantize_nf4(weight):
134        return _your_dequantize_nf4(weight.weight.data, weight.weight.quant_state)


1 #-----other version -----
2 # from triton import jit
3 # import triton
4 # import triton.language as tl
5 # import pdb
6
7 # @triton.jit
8 # def _your_dequantize_nf4_kernel(#pointers
9 #                                absmx2_ptr, absmx1_ptr, out1_ptr, weight_ptr, code2_ptr, code1_ptr, offset1,
10 #                                #Dimensions
11 #                                num_elems_1, num_elems_2, shift1, shift2,
12 #                                #meta param
13 #                                BLOCK_SIZE: tl.constexpr,
14 #                                ):
15
16 #     #pdb.set_trace()
17 #     pid = tl.program_id(axis=0)
18
19 #     #tl.static_print("pid", pid)
20
21 #     new_offsets = pid * BLOCK_SIZE + tl.arange(0, BLOCK_SIZE)
22
23 #     #tl.static_print("new_offsets", new_offsets)
24
25 #     mask_for_abs1 = new_offsets < num_elems_1
26
27 #     #tl.static_print("mask_for_abs1", mask_for_abs1)
28
29 #     g = new_offsets >> shift1
30
31 #     #tl.static_print("g", g)
32
33 #     mask_for_g = g < num_elems_2
34
35 #     #tl.static_print("mask_for_g", mask_for_g)
36
37 #     # Recompute the intermediate value on the fly:

```

```

38 #     absmax2_val = tl.load(absmax2_ptr + (g >> shift2), mask=mask_for_g, other=128)
39 #     #tl.static_print("absmax2_val", absmax2_val)
40
41 #     absmax1_val = tl.load(absmax1_ptr + g, mask=mask_for_g, other=128)
42 #     #tl.static_print("absmax1_val", absmax1_val)
43 #     x_val = tl.load(code2_ptr + tl.cast(absmax1_val, tl.int32), mask=mask_for_g, other=128)
44 #     #tl.static_print("x_val", x_val)
45 #     gathered_max = tl.fma(x_val, absmax2_val, tl.load(offset1).cast(tl.float32))
46
47 #     #tl.static_print("gathered_max", gathered_max)
48
49
50 #     # Now process weights based on the recomputed gathered_max
51 #     weights_vals = tl.load(weight_ptr + new_offsets, mask=mask_for_abs1, other=128)
52 #     #tl.static_print("weights_vals", weights_vals)
53 #     high_bits = weights_vals >> 4
54 #     #tl.static_print("high_bits", high_bits)
55
56 #     x_hi = tl.load(code1_ptr + high_bits, mask=mask_for_abs1, other=4)
57 #     #tl.static_print("x_hi", x_hi)
58 #     hi_val = x_hi * gathered_max
59 #     #tl.static_print("hi_val", hi_val)
60
61 #     low_part = weights_vals & 0x0F
62 #     #tl.static_print("low_part", low_part)
63 #     x_lo = tl.load(code1_ptr + low_part, mask=mask_for_abs1, other=4)
64 #     #tl.static_print("x_lo", x_lo)
65 #     lo_val = x_lo * gathered_max
66 #     #tl.static_print("lo_val", lo_val)
67
68 #     # Compute output indices and store results
69 #     out_hi2 = new_offsets * 2
70 #     out_lo2 = new_offsets*2 + 1
71 #     tl.store(out1_ptr + out_hi2, hi_val, mask=mask_for_abs1)
72 #     tl.store(out1_ptr + out_lo2, lo_val, mask=mask_for_abs1)
73
74
75 # def _your_dequantize_nf4(weight, quant_state):
76 #     absmax2 = quant_state.state2.absmax
77 #     blocksize2 = quant_state.state2.blocksize
78 #     absmax1 = quant_state.absmax
79 #     #out2 = torch.empty(absmax1.shape, dtype=quant_state.dtype, device=absmax1.device)
80 #     offset1 = quant_state.offset
81 #     out1 = torch.empty(quant_state.shape, dtype=quant_state.dtype, device=weight.device)
82 #     blocksize1 = quant_state.blocksize//2
83 #     num_elems_1 = out1.numel()
84 #     num_elems_2 = absmax1.numel()
85 #     code2 = quant_state.state2.code
86 #     code1 = quant_state.code
87
88
89 #     shift1 = blocksize1.bit_length() - 1
90 #     shift2 = blocksize2.bit_length() - 1
91
92 #     grid = lambda meta: (triton.cdiv(weight.numel(), meta['BLOCK_SIZE']),)
93
94 #     _your_dequantize_nf4_kernel[grid](
95 #         #pointers
96 #         absmax2 , absmax1 , out1 , weight , code2 , code1 , offset1,
97 #         #Dimensions
98 #         num_elems_1, num_elems_2, shift1, shift2,
99 #         #meta param
100 #         BLOCK_SIZE = 2048,
101 #     )
102
103 #     return out1
104
105 # def your_dequantize_nf4(weight):
106 #     return _your_dequantize_nf4(weight.weight.data, weight.weight.quant_state)

1 ## TEST IT BELOW:
2 test_dequantize(your_dequantize_nf4)
3
4 ## CALCULATE SPEEDUP (hopefully 1.15x faster or more)
5 test_dequantize(unsloth_dequantize) / test_dequantize(your_dequantize_nf4)

```

🔄 1.1295745450117198

Marking Criteria for A) Max points = 14

```

if attempted_A:
    A_score = 0
    if single_triton_kernel: A_score += 3
    speedup = old_time / new_time
    if speedup <= 1.00: A_score -= 3
    if speedup >= 1.05: A_score += 1
    if speedup >= 1.10: A_score += 2
    if speedup >= 1.15: A_score += 2
    if kernel_works_in_torch_compile: A_score += 1
    else: A_score -= 1
    if custom_asm_works: A_score += 3
    if uses_cache_eviction: A_score += 1
    if tested_in_f16_and_bf16: A_score += 1
    else: A_score -= 1
    final_score += A_score
else:
    final_score += 0

```

---



---



---

## ✓ B) Make QLoRA work with FSDP2 [Difficulty: Medium to Hard] [Max points: 10]

1. Goal: Write a single Python script to finetune Llama 3.1 8B on 2x or more GPUs with FSDP2.
2. You must showcase this working in a free **Kaggle notebook with 2 x Tesla T4 GPUs**.
3. Pipeline parallelism is also fine, but must utilize [zero bubble scheduling](#) somehow.
4. Can use a pre-quantized 4bit BnB safetensor file from [Unsloth's HF page](#) or a full 16bit one, but must do QLoRA.
5. Can use `accelerate` but must be FSDP2 or related - you can investigate <https://github.com/huggingface/accelerate/pull/3394>, Torch Titan, other repos etc.
6. Must be fully `transformers` compatible - so we must use `TrainingArguments` and `Trainer`, or TRL related classes.
7. The loss must be equivalent to single GPU training.
8. You must enable all features in FSDP2 - ie showcase offloading, checkpointing, mixed precision training etc.
9. You can use `nf4` from `torch` A0, but best from `bitsandbytes`.
10. Finally showcase everything working in a free Kaggle 2x Tesla T4 notebook.

Ran on Kaggle :

```

1 # HELPFUL functions to undo Unsloth patches:
2 import sys
3
4 def remove_patched_module(package_name):
5     modules_to_delete = [
6         name for name in sys.modules
7         if name == package_name or name.startswith(package_name + ".")
8     ]
9     for name in modules_to_delete: del sys.modules[name]
10
11 remove_patched_module("trl")
12 remove_patched_module("transformers")
13 remove_patched_module("peft")
14 remove_patched_module("bitsandbytes")

```

```

1 def dequantize_bnb_weight_temp_patched(weight: "torch.nn.Parameter", dtype: "torch.dtype", state=None):
2     """
3     Helper function to dequantize 4bit or 8bit bnb weights.
4
5     If the weight is not a bnb quantized weight, it will be returned as is.
6     """
7     if not isinstance(weight, torch.nn.Parameter):
8         raise TypeError(f"Input weight should be of type nn.Parameter, got {type(weight)} instead")
9
10    cls_name = weight.__class__.__name__
11    if cls_name not in ("Params4bit", "Int8Params"):
12        return weight
13
14    if cls_name == "Params4bit":
15        output_tensor = your_dequantize_nf4(weight.data, weight.quant_state)
16        logger.warning_once(

```

```

17         f"The model is going to be dequantized in {output_tensor.dtype} – if you want to upcast it to another dtype,
18     )
19     return output_tensor.to(dtype)
20
21     if state.SCB is None:
22         state.SCB = weight.SCB
23
24     if hasattr(bnb.functional, "int8_vectorwise_dequant"):
25         # Use bitsandbytes API if available (requires v0.45.0+)
26         dequantized = bnb.functional.int8_vectorwise_dequant(weight.data, state.SCB)
27     else:
28         # Multiply by (scale/127) to dequantize.
29         dequantized = weight.data * state.SCB.view(-1, 1) * 7.874015718698502e-3
30
31     return dequantized.to(dtype)

```

The script below runs fine in Kaggle 2x Telsa T4s(check Kaggle notebook attached):

```

1
2 from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig, TrainingArguments
3 import transformers.integrations.bitsandbytes
4 transformers.integrations.bitsandbytes.dequantize_bnb_weight = dequantize_bnb_weight_temp_patched
5
6 import os
7 import torch
8 os.environ["HF_HUB_ENABLE_HF_TRANSFER"] = "1"
9 os.environ["CUDA_VISIBLE_DEVICES"] = "0,1"
10 os.environ["PYTORCH_CUDA_ALLOC_CONF"] = "expandable_segments:True,\\"
11     "roundup_power2_divisions:[32:256,64:128,256:64,>:32]"
12 from peft import get_peft_model, LoraConfig, TaskType
13 from trl import SFTTrainer, SFTConfig
14 from datasets import load_dataset
15 import os
16 max_seq_length = 2048
17 torch.set_default_dtype(torch.float16)
18 model_name = "unsloth/meta-Llama-3.1-8B-Instruct-bnb-4bit"
19 dtype = torch.float16
20 bnb_config = BitsAndBytesConfig(
21     load_in_4bit = True,
22     bnb_4bit_use_double_quant = True,
23     bnb_4bit_quant_type = "nf4",
24     bnb_4bit_compute_dtype = dtype,
25 )
26 model = AutoModelForCausalLM.from_pretrained(
27     model_name,
28     device_map = "auto",
29     attn_implementation = "sdpa",
30     quantization_config = bnb_config,
31 )
32 tokenizer = AutoTokenizer.from_pretrained(model_name)
33 tokenizer.padding_side = "right"
34
35 lora_config = LoraConfig(
36     r = 64,
37     lora_alpha = 128,
38     target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
39         "gate_proj", "up_proj", "down_proj"],
40     lora_dropout = 0,
41     bias = "none",
42     task_type = TaskType.CAUSAL_LM,
43 )
44
45 # Get LoRA and setup model
46 model = get_peft_model(model, lora_config)
47 with torch.no_grad():
48     for name, param in model.named_parameters():
49         if ".lora_A." in name or ".lora_B." in name: param.requires_grad_(True)
50         else: param.requires_grad_(False)
51 model.gradient_checkpointing_enable()
52 model.enable_input_require_grads()
53
54 # Get dataset
55 from datasets import load_dataset
56 from trl import SFTTrainer, SFTConfig
57 url = "https://huggingface.co/datasets/laion/OIG/resolve/main/unified_chip2.jsonl"
58 dataset = load_dataset("json", data_files = {"train": url}, split = "train[:10%]")
59
60 fsdp_config = {
61     "compute_environment": "LOCAL_MACHINE",
62     "debug": False,
63     "distributed_type": "FSDP",

```



```

64 "downcast_bf16": "no",
65 "enable_cpu_affinity": False,
66 "fsdp_config": {
67     "fsdp_activation_checkpointing": True,
68     "fsdp_auto_wrap_policy": "TRANSFORMER_BASED_WRAP",
69     "fsdp_backward_prefetch": "BACKWARD_PRE",
70     "fsdp_cpu_ram_efficient_loading": True,
71     "fsdp_forward_prefetch": True,
72     "fsdp_offload_params": True,
73     "fsdp_sharding_strategy": "FULL_SHARD",
74     "fsdp_state_dict_type": "SHARDED_STATE_DICT",
75     "fsdp_sync_module_states": True,
76     "fsdp_use_orig_params": True
77 },
78 "machine_rank": 0,
79 "main_process_ip": "",
80 "main_process_port": 8000,
81 "main_training_function": "main",
82 "mixed_precision": "fp16",
83 "num_machines": 2,
84 "num_processes": 2,
85 "rdzv_backend": "static",
86 "same_network": True,
87 "tpu_env": [],
88 "tpu_use_cluster": False,
89 "tpu_use_sudo": False,
90 "use_cpu": False
91 }
92
93 training_args = SFTConfig(
94     per_device_train_batch_size = 2,
95     gradient_accumulation_steps = 4,
96     warmup_steps = 1,
97     max_steps = 10,
98     logging_steps = 1,
99     output_dir = "outputs",
100     seed = 3407,
101     max_seq_length = max_seq_length,
102     fp16 = model.get_input_embeddings().weight.dtype == torch.float16,
103     bf16 = model.get_input_embeddings().weight.dtype == torch.bfloat16,
104     report_to = "none", # For W&B
105     dataset_num_proc = 4,
106
107     fsdp_config=fsdp_config
108 )
109
110 trainer = SFTTrainer(
111     model = model,
112     train_dataset = dataset,
113     processing_class = tokenizer,
114     args =training_args,
115 )
116 trainer.train()
117

```

Reminder your code must have the same loss curve over 60 steps or so.

```

1 #del model
2 import gc
3 gc.collect()
4 torch.cuda.empty_cache()

```

## Marking Criteria for B) Max points = 10

```

if attempted_B:
    B_score = 0
    if FSDP2_works_with_QLoRA:
        if torch_compile_works: B_score += 5
        else: B_score += 3
        if uses_part_A_and_single_kernel_and_faster: B_score += 3
        elif uses_torchAO:
            if torchAO_slower_than_BnB: B_score -= 3
    elif TP_or_PP_with_QLoRA:
        if zero_bubble: B_score += 3
        else: B_score += 2
    elif FSDP1_works_with_QLoRA:

```

```

        B_score += 1
    if kaggle_notebook_2_tesla_t4_example:
        B_score += 2
    else:
        B_score = 0
    final_score += B_score
else:
    final_score -= 2

```

---

## ✓ C) Make `torch.compile` work without graph breaks for QLoRA [Difficulty: Easy to Medium] [Max points: 9]

1. Goal: Write a single Python script like task B), except the goal is to `torch.compile` all modules if possible.
2. There must NOT be graph breaks, and excessive re-compilations should not be seen.
3. You should have say max 30 compilations. Over 60 is definitely wrong.
4. The loss must match with the non compiled module.
5. Utilize patching as much as possible.
6. Think about which areas might need disabling for compilation. Think about regional compilation. How do we compile sections efficiently?
7. Log memory / VRAM usage, and monitor speedups as well.
8. Must work for QLoRA.

We provided a script below, and showcased how to detect if graph breaks are seen. We also torch compiled the MLP for Llama:

```

1 from warnings import warn
2 from typing import Callable, Optional, Tuple
3 import bitsandbytes.functional as F
4 from functools import reduce # Required in Python 3
5 import operator
6
7 def prod(iterable):
8     return reduce(operator.mul, iterable, 1)
9
10 torch_compile_options = torch_compile_options = {
11     "epilogue_fusion" : True,
12     "max_autotune" : True,
13     "shape_padding" : True,
14     "trace.enabled" : True,
15     "triton.cudagraphs" : False,
16 }
17 @torch.compile(fullgraph=False, dynamic=True, options = torch_compile_options)
18 class MatMul4Bit(torch.autograd.Function):
19     # forward is the same, but we added the fallback for pre-turing GPUs
20     # backward is mostly the same, but adds one extra clause (see "elif state.CxB is not None")
21
22     @staticmethod
23     def forward(ctx, A, B, out=None, bias=None, quant_state: Optional[F.QuantState] = None):
24         # default of pytorch behavior if inputs are empty
25         ctx.is_empty = False
26         if prod(A.shape) == 0:
27             ctx.is_empty = True
28             ctx.A = A
29             ctx.B = B
30             ctx.bias = bias
31             B_shape = quant_state.shape
32             if A.shape[-1] == B_shape[0]:
33                 return torch.empty(A.shape[:-1] + B_shape[1:], dtype=A.dtype, device=A.device)
34             else:
35                 return torch.empty(A.shape[:-1] + B_shape[1:], dtype=A.dtype, device=A.device)
36
37     # 1. Dequantize
38     # 2. MatmulN
39     output = torch.nn.functional.linear(A, F.dequantize_4bit(B, quant_state).to(A.dtype).t(), bias)
40
41     # 3. Save state
42     ctx.state = quant_state
43     ctx.dtype_A, ctx.dtype_B, ctx.dtype_bias = A.dtype, B.dtype, None if bias is None else bias.dtype
44
45     if any(ctx.needs_input_grad[2:]):

```

```

46         ctx.tensors = (None, B)
47     else:
48         ctx.tensors = (None, None)
49
50     return output
51
52     @staticmethod
53     def backward(ctx, grad_output):
54         if ctx.is_empty():
55             bias_grad = None if ctx.bias is None else torch.zeros_like(ctx.bias)
56             return torch.zeros_like(ctx.A), torch.zeros_like(ctx.B), None, bias_grad, None
57
58         req_gradA, _, _, req_gradBias, _ = ctx.needs_input_grad
59         _, B = ctx.tensors
60
61         grad_A, grad_B, grad_bias = None, None, None
62
63         if req_gradBias:
64             # compute grad_bias first before changing grad_output dtype
65             grad_bias = grad_output.sum(0, dtype=ctx.dtype_bias)
66
67         # not supported by PyTorch. TODO: create work-around
68         # if req_gradB: grad_B = torch.matmul(grad_output.t(), A)
69         if req_gradA:
70             grad_A = torch.matmul(grad_output, F.dequantize_4bit(B, ctx.state).to(grad_output.dtype).t())
71
72         return grad_A, grad_B, None, grad_bias, None
73
74 import bitsandbytes as bnb
75 from bitsandbytes.nn import Linear4bit
76
77 bnb.MatMul4Bit = MatMul4Bit
78
79 import torch
80 from torch.nn.attention.flex_attention import (
81     _DEFAULT_SPARSE_BLOCK_SIZE,
82     create_block_mask,
83     create_mask,
84     flex_attention,
85 )
86 from transformers.models.llama.modeling_llama import repeat_kv
87
88 torch_compile_options = torch_compile_options = {
89     "epilogue_fusion" : True,
90     "max_autotune" : True,
91     "shape_padding" : True,
92     "trace.enabled" : True,
93     "triton.cudagraphs" : False,
94 }
95
96 @torch.compile(fullgraph=False, dynamic=True, options = torch_compile_options)
97 def compiled_llama_mlp(self, x):
98     down_proj = self.down_proj(self.act_fn(self.gate_proj(x)) * self.up_proj(x))
99     return down_proj
100
101 @torch.compile(fullgraph=False, dynamic=True, options = torch_compile_options)
102 def flex_attention_forward(
103     module: nn.Module,
104     query: torch.Tensor,
105     key: torch.Tensor,
106     value: torch.Tensor,
107     attention_mask: Optional[torch.Tensor],
108     scaling: float,
109     dropout: float = 0.0,
110     **kwargs,
111 ):
112     key_states = repeat_kv(key, module.num_key_value_groups)
113     value_states = repeat_kv(value, module.num_key_value_groups)
114
115     # 1. Define score_mod (This replaces the explicit weight calculation)
116     def score_mod(score, b, h, q_idx, kv_idx):
117         return score * scaling
118
119     # 2. Define mask_mod (This handles masking, including causal masking)
120     def mask_mod(b, h, q_idx, kv_idx):
121         if attention_mask is None:
122             return True # No Masking
123         causal_mask = attention_mask[b, h, q_idx, kv_idx] # Accessing the mask for the specific batch, head, query and
124         return causal_mask.bool() # Convert to boolean
125
126     # 3. Call flex_attention
127     attn_output = flex_attention(query, key_states, value_states, score_mod=score_mod, mask_mod=mask_mod)

```

```

128
129     attn_output = attn_output.transpose(1, 2).contiguous()
130
131     # 4. Apply Dropout (Explicitly)
132     attn_output = F.dropout(attn_output, p=dropout, training=module.training) # Apply dropout here
133
134     return attn_output, None
135
136 # @torch.compile(fullgraph=False, dynamic=True, options = torch_compile_options)
137 # class LlamaRMSNorm(nn.Module):
138 #     def __init__(self, hidden_size, eps=1e-6):
139 #         """
140 #         LlamaRMSNorm is equivalent to T5LayerNorm
141 #         """
142 #         super().__init__()
143 #         self.weight = nn.Parameter(torch.ones(hidden_size))
144 #         self.variance_epsilon = eps
145
146 #     def forward(self, hidden_states):
147 #         input_dtype = hidden_states.dtype
148 #         hidden_states = hidden_states.to(torch.float32)
149 #         variance = hidden_states.pow(2).mean(-1, keepdim=True)
150 #         hidden_states = hidden_states * torch.rsqrt(variance + self.variance_epsilon)
151 #         return self.weight * hidden_states.to(input_dtype)
152
153 #     def extra_repr(self):
154 #         return f"{{tuple(self.weight.shape)}}, eps={{self.variance_epsilon}}"
155
156
157 # def fixed_cross_entropy(source, target, num_items_in_batch: int = None, ignore_index: int = -100, **kwargs):
158 #     reduction = "sum" if num_items_in_batch is not None else "mean"
159 #     loss = nn.functional.cross_entropy(source, target, ignore_index=ignore_index, reduction=reduction)
160 #     if reduction == "sum":
161 #         loss = loss / num_items_in_batch
162 #     return loss
163
164
165 # @torch.compile(fullgraph=False, dynamic=True, options = torch_compile_options)
166 # def ForCausalLMLoss(
167 #     logits, labels, vocab_size: int, num_items_in_batch: int = None, ignore_index: int = -100, **kwargs
168 # ):
169 #     # Upcast to float if we need to compute the loss to avoid potential precision issues
170 #     logits = logits.float()
171 #     labels = labels.to(logits.device)
172 #     # Shift so that tokens < n predict n
173 #     labels = nn.functional.pad(labels, (0, 1), value=ignore_index)
174 #     shift_labels = labels[... , 1:].contiguous()
175
176 #     # Flatten the tokens
177 #     logits = logits.view(-1, vocab_size)
178 #     shift_labels = shift_labels.view(-1)
179 #     # Enable model parallelism
180 #     shift_labels = shift_labels.to(logits.device)
181 #     loss = fixed_cross_entropy(logits, shift_labels, num_items_in_batch, ignore_index, **kwargs)
182 #     return loss
183
184
185 import transformers.models.llama.modeling_llama
186 #transformers.models.llama.modeling_llama.LlamaMLP.forward = compiled_llama_mlp
187 transformers.models.llama.modeling_llama.eager_attention_forward = flex_attention_forward
188 transformers.models.llama.modeling_llama.LlamaRMSNorm = torch.compile(transformers.models.llama.modeling_llama.LlamaRMS
189 transformers.loss.loss_utils.ForCausalLMLoss = torch.compile(transformers.loss.loss_utils.ForCausalLMLoss, fullgraph=False)

```

```

1 from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig
2 from peft import get_peft_model, LoraConfig, TaskType
3 os.environ["HF_HUB_ENABLE_HF_TRANSFER"] = "1"
4 os.environ["CUDA_VISIBLE_DEVICES"] = "0,1"
5 os.environ["PYTORCH_CUDA_ALLOC_CONF"] = \
6     "expandable_segments:True,"\
7     "roundup_power2_divisions:[32:256,64:128,256:64,>:32]"
8
9 max_seq_length = 1024
10 torch.set_default_dtype(torch.float16)
11 model_name = "unsloth/Llama-3.2-1B-Instruct-bnb-4bit"
12 dtype = torch.float16
13 bnb_config = BitsAndBytesConfig(
14     load_in_4bit = True,
15     bnb_4bit_use_double_quant = True,
16     bnb_4bit_quant_type = "nf4",
17     bnb_4bit_compute_dtype = dtype,
18 )
19 model = AutoModelForCausalLM.from_pretrained(

```

```

20     model_name,
21     device_map = "auto",
22     attn_implementation = "sdpa",
23     quantization_config = bnb_config,
24 )
25 tokenizer = AutoTokenizer.from_pretrained(model_name)
26 tokenizer.padding_side = "right"
27
28 lora_config = LoraConfig(
29     r = 32,
30     lora_alpha = 64,
31     target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
32                       "gate_proj", "up_proj", "down_proj"],
33     lora_dropout = 0,
34     bias = "none",
35     task_type = TaskType.CAUSAL_LM,
36 )
37
38 # Get LoRA and setup model
39 model = get_peft_model(model, lora_config)
40 with torch.no_grad():
41     for name, param in model.named_parameters():
42         if ".lora_A." in name or ".lora_B." in name: param.requires_grad_(True)
43         else: param.requires_grad_(False)
44
45 # Currently GC will cause torch.compile to be disabled, so disable it
46 # model.gradient_checkpointing_enable()
47 model.enable_input_require_grads()
48
49 # Get dataset
50 from datasets import load_dataset
51 from trl import SFTTrainer, SFTConfig
52 url = "https://huggingface.co/datasets/laion/OIG/resolve/main/unified_chip2.jsonl"
53 dataset = load_dataset("json", data_files = {"train" : url}, split = "train[:10%]")

```

```

🔄 config.json: 100%                               1.52k/1.52k [00:00<00:00, 149kB/s]
Unused kwargs: ['_load_in_4bit', '_load_in_8bit', 'quant_method']. These kwargs are not used in <class 'transformers.uti
Unused kwargs: ['_load_in_4bit', '_load_in_8bit', 'quant_method']. These kwargs are not used in <class 'transformers.uti
model.safetensors: 100%                             1.03G/1.03G [00:08<00:00, 215MB/s]
generation_config.json: 100%                         234/234 [00:00<00:00, 25.2kB/s]
tokenizer_config.json: 100%                          54.7k/54.7k [00:00<00:00, 5.71MB/s]
tokenizer.json: 100%                                17.2M/17.2M [00:00<00:00, 39.5MB/s]
special_tokens_map.json: 100%                       454/454 [00:00<00:00, 27.6kB/s]
unified_chip2.jsonl: 100%                           95.6M/95.6M [00:00<00:00, 213MB/s]
Generating train split:    210289/0 [00:00<00:00, 342782.99 examples/s]

```

We provide full logging for torch.compile like below:

```

1 # Must show all graph breaks are not seen with torch.compile
2 import os
3 os.environ["TORCHDYNAMO_VERBOSE"] = "1"
4 os.environ["TORCHINDUCTOR_FORCE_DISABLE_CACHES"] = "1"
5 os.environ["TORCHINDUCTOR_COMPILE_THREADS"] = "1"
6
7 import logging
8 torch._inductor.config.debug = True
9 torch._logging.set_logs(
10     dynamo = logging.WARN,
11     inductor = logging.WARN,
12     graph_breaks = True,
13     recompiles = True,
14     recompiles_verbose = True,
15     compiled_autograd_verbose = True,
16     # aot_joint_graph = True, # Enable for more logs
17     # aot_graphs = True,
18 )
19 torch._dynamo.config.verbose = True
20 torch._dynamo.config.suppress_errors = False

```

When we execute the code below, we can see graph breaks - remove them.

```


1 trainer = SFTTrainer(
2     model = model,
3     train_dataset = dataset,

```

```

4     processing_class = tokenizer,
5     args = SFTConfig(
6         per_device_train_batch_size = 1,
7         gradient_accumulation_steps = 2,
8         warmup_steps = 1,
9         max_steps = 10,
10        logging_steps = 1,
11        output_dir = "outputs",
12        seed = 3407,
13        max_seq_length = max_seq_length,
14        fp16 = model.get_input_embeddings().weight.dtype == torch.float16,
15        bf16 = model.get_input_embeddings().weight.dtype == torch.bfloat16,
16        report_to = "none", # For W&B
17        dataset_num_proc = 4,
18    ),
19 )
20 trainer.train()

```

 Converting train dataset to ChatML (num\_proc=4): 100% 21029/21029 [00:00<00:00, 31315.24 examples/s]  
 Applying chat template to train dataset (num\_proc=4): 100% 21029/21029 [00:04<00:00, 12486.34 examples/s]  
 Tokenizing train dataset (num\_proc=4): 100% 21029/21029 [00:13<00:00, 708.60 examples/s]  
 Truncating train dataset (num\_proc=4): 100% 21029/21029 [00:05<00:00, 3604.97 examples/s]  
 W0310 01:33:43.652000 6856 torch/\_inductor/debug.py:434] [0/0] model\_\_0\_forward\_1 debug trace: /content/torch\_compile\_de  
 W0310 01:33:45.528000 6856 torch/\_inductor/debug.py:434] [0/0] model\_\_0\_backward\_2 debug trace: /content/torch\_compile\_d  
[10/10 00:04, Epoch 0/1]

#### Step Training Loss

1	1.519300
2	2.393000
3	2.499000
4	3.531300
5	2.135300
6	2.973600
7	2.242700
8	1.621300
9	2.216400
10	2.676600

```

TrainOutput(global_step=10, training_loss=2.3808391451835633, metrics={'train_runtime': 14.8098,
'train_samples_per_second': 1.35, 'train_steps_per_second': 0.675, 'total_flos': 10592155496448.0, 'train_loss':
2.3808391451835633})

```

```

1 del model
2 import gc
3 gc.collect()
4 torch.cuda.empty_cache()

```

Log all your steps for debugging in a Colab (maybe this one). Edward's blog <http://blog.ezyang.com/>, Horace's blogs <https://www.thonking.ai/>, Slaying OOMs by Jane & Mark: <https://www.youtube.com/watch?v=UvRI4ansfCg> could be useful.

## Marking Criteria for C) Max points = 9

```

if attempted_C:
    C_score = 0
    if uses_flex_attention:
        if dynamic_sequence_length_works: C_score += 3
        else: C_score += 1
    if no_torch_compile_BnB: C_score -= 2
    elif use_part_A: C_score += 1
    elif torch_compile_BnB: C_score += 1

    if attention_compiled:
        if excessive_recompilation: C_score -= 3
        else: C_score += 2
    if mlp_compiled:
        if excessive_recompilation: C_score -= 3
        C_score += 1

```

```

if not loss_compiled: C_score -= 1
if not layernorms_compiled: C_score -= 3

if max_autotune_triton_matmul:
    if excessive_recompilation: C_score -= 2
    else: C_score += 2

final_score += C_score
else:
    final_score -= 1

```

---



---

## D) Help solve 🐢 Unsloth issues! [Difficulty: Varies] [Max points: 12]

Head over to <https://github.com/unslothai/unsloth>, and find some issues which are still left standing / not resolved. The tag **currently fixing** might be useful.

Each successfully accepted and solved issue will also have \$100 to \$1000 of bounties.

It's best to attempt these features:

- **Tool Calling** [Points = 1] Provide a tool calling Colab notebook and make it work inside of Unsloth. Bounty: \$1000
- **GGUF Vision support** [Points = 1] Allow exporting vision finetunes to GGUF directly. Llava and Qwen VL must work. Bounty: \$500
- **Refactor Attention** [Points = 2] Refactor and merge xformers, SDPA, flash-attn, flex-attention into a simpler interface. Must work seamlessly inside of Unsloth. Bounty: \$350
- **[DONE] Windows support** [Points = 2] Allow `pip install unsloth` to work in Windows - Triton, Xformers, bitsandbytes should all function. You might need to edit `pyproject.toml`. Confirm it works. Bounty: \$300
- **Support Sequence Classification** [Points = 1] Create patching functions to patch over `AutoModelForSequenceClassification`, and allow finetuner to use `AutoModelForSequenceClassification`. Bounty: \$200
- **VLMs Data Collator** [Points = 1] Make text & image mixing work efficiently -so some inputs can be text only. Must work on Qwen, Llama, Pixtral. Bounty: \$100
- **VLMs image resizing** [Points = 1] Allow finetuner to specify maximum image size, or get it from the `config.json` file. Resize all images to specific size to reduce VRAM. Bounty: \$100
- **Support Flex Attention** [Points = 2] Allow dynamic sequence lengths without excessive recompilation. Make this work on SWAs and normal causal masks. Also packed sequence masks. Bounty: \$100
- **VLMs train only on completions** [Points = 1] Edit `train_on_responses_only` to allow it to work on VLMs. Bounty: \$100

## Marking Criteria for D) Max points = 12

```

if attempted_D:
    D_score = 0
    for subtask in subtasks:
        if successfully_completed_subtask:
            D_score += score_for_subtask
    final_score += D_score

```

---



---

## ✓ E) Memory Efficient Backprop [Difficulty: Medium to Hard] [Max points: 10]

In LLMs, the last layer is a projection matrix to calculate the probabilities of the next token, ie  $\sigma(XW)$ . However, if the vocabulary size is very large, say 128K, then the materialization of the logits causes VRAM spikes.

For example, if the `bsz = 4`, `qlen = 4096`, `hd = 4096`, `vocab = 128K`, then the memory usage for the logits in `bfloat16` would be 4GB. In the worst case, we might even need to upcast logits to `float32`, so 8GB is needed.

In Unsloth, we utilize [Apple's Cut Cross Entropy Loss](#) to reduce VRAM usage, by allowing a Triton kernel to create the logits on the fly to calculate the cross entropy loss. But this does not generalize well to other functions.

Our goal is to generalize this ultimately, but directly creating logits on the fly will be hard. Instead, let's take a slightly less complex approach. Let's first review some stuff. We first notice that during the normal case after forming the intermediate logits for 2 batches, we then do a

gather function to aggregate the intermediate results into a single column:

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \times W = \begin{bmatrix} x_1 W \\ x_2 W \end{bmatrix}$$

$$f\left(\begin{bmatrix} x_1 W \\ x_2 W \end{bmatrix}\right) = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

So, if we can somehow skip the materialization of the intermediate logits, and just output the output of  $f$ , we can save a lot of VRAM!

Notice during backpropagation we can use the chain rule:

$$\frac{dL}{dX} = \frac{dL}{dy} \frac{dy}{dX}; \frac{dL}{dW} = \frac{dL}{dy} \frac{dy}{dW}$$

$$\frac{dL}{dy} = \text{Downstream from backprop}$$

$$\frac{dy}{dX} = W^T$$

$$\frac{dy}{dW} = X^T$$

$$\frac{dL}{dX} = \frac{dL}{dy} W^T$$

$$\frac{dL}{dW} = X^T \frac{dL}{dy}$$

If we simply compute the intermediate tensors on the fly via batches, say we do batch 1, then batch 2, we can reduce VRAM usage from 4GB to 2GB!

$$\frac{dL}{dX} = \begin{bmatrix} \frac{dL_1}{dy_1} W^T \\ \frac{dL_2}{dy_2} W^T \end{bmatrix}$$

$$\frac{dL}{dW} = \left( X_1^T \frac{dL_1}{dy_1} + X_2^T \frac{dL_2}{dy_2} \right)$$

1. Your goal is to write a `torch.autograd.Function` with a forward and backward pass showcasing this memory efficient implementation.
2. You must NOT hard code the derivatives - move the transformation function from the logits / intermediate tensors to a smaller tensor as a separate function which can allow `autograd` to pass through it.
3. As a hint, look at `torch.checkpoint` at <https://github.com/pytorch/pytorch/blob/main/torch/utils/checkpoint.py>. Also, don't forget about the upstream gradients! We need to multiply them to the current gradients!
4. Make the Cross Entropy Loss work. You must show other functions working as well.

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import matplotlib.pyplot as plt
5
6
7 class MemoryEfficientLinear(torch.autograd.Function):
8     @staticmethod
9     def forward(ctx, X, linear, labels, forward_function, chunk_size=None):
10
11         bsz, seq_len, hd = X.shape
12         # Flatten the batch and sequence dims.
13         X_flat = X.reshape(-1, hd) # shape: (bsz * seq_len, hd)
14         labels_flat = labels.reshape(-1) # shape: (bsz * seq_len,)
15         total_tokens = X_flat.shape[0]
16
17         if chunk_size is None:
18             chunk_size = total_tokens // 2 if total_tokens > 1 else 1
19         ctx.chunk_size = chunk_size
20
21         # Determine valid tokens (assume -100 indicates positions to ignore).
22         valid_mask = (labels_flat != -100)
23         total_valid_tokens = valid_mask.sum().to(dtype=X.dtype)
24         ctx.total_valid_tokens = total_valid_tokens
25
26         total_loss = 0.0
27         grad_chunks = []
28
29         # Split the flattened inputs into chunks.
30         X_chunks = torch.split(X_flat, chunk_size, dim=0)
31         label_chunks = torch.split(labels_flat, chunk_size, dim=0)
32
33         for chunk, lab in zip(X_chunks, label_chunks):
34             # Count valid tokens in this chunk.
35             valid_tokens_chunk = (lab != -100).sum().to(dtype=X.dtype)

```



```

36         # Calculate the scaling factor for this chunk.
37         scaling_factor = valid_tokens_chunk / total_valid_tokens if total_valid_tokens > 0 else 0.0
38
39         # Here we define a lambda that wraps the transformation_function.
40         compute_chunk_loss = lambda chunk_in: forward_function(chunk_in, linear, lab) * scaling_factor
41
42         # Compute both the loss value and gradient for the chunk.
43         chunk_grad, chunk_loss = torch.func.grad_and_value(
44             compute_chunk_loss, argnums=0, has_aux=False
45         )(chunk)
46         total_loss += chunk_loss
47         grad_chunks.append(chunk_grad)
48
49         # Save the precomputed gradients to be used in backward.
50         precomputed_grad = torch.cat(grad_chunks, dim=0)
51         ctx.precomputed_grad = precomputed_grad
52         ctx.input_shape = X.shape # to reshape the gradient back to original dimensions
53         return total_loss
54
55     @staticmethod
56     def backward(ctx, dY):
57         # Retrieve the precomputed gradient and multiply by the upstream gradient.
58         precomputed_grad = ctx.precomputed_grad
59         grad_X_flat = precomputed_grad * dY
60         grad_X = grad_X_flat.view(ctx.input_shape)
61
62         return grad_X, None, None, None, None
63
64
65 # (1) Original Transformation Function using CrossEntropyLoss.
66 def transformation_function(flat_batch, linear, flat_labels):
67     """
68     Computes standard cross-entropy loss.
69
70     Parameters:
71         flat_batch: Tensor of shape (N, hd)
72         linear: Linear layer mapping from hd to vocab
73         flat_labels: Tensor of shape (N,), containing target class indices.
74
75     Returns:
76         loss: Scalar tensor representing the cross-entropy loss.
77     """
78     x = linear(flat_batch).float() # shape: (N, vocab)
79
80
81     ce_loss = nn.CrossEntropyLoss(ignore_index=-100, reduction="mean")
82     loss = ce_loss(x, flat_labels)
83     return loss
84
85 # (2) Label-Smoothing Transformation Function.
86 def label_smoothing_transformation_function(flat_batch, linear, flat_labels, smoothing=0.1):
87     """
88     Computes a label-smoothed cross-entropy loss.
89
90     Parameters:
91         flat_batch: Tensor of shape (N, hd)
92         linear: Linear layer mapping from hd to vocab
93         flat_labels: Tensor of shape (N,), containing target class indices.
94         smoothing: The label smoothing factor (default is 0.1)
95
96     Returns:
97         loss: Scalar tensor representing the label-smoothed loss.
98     """
99     x = linear(flat_batch).float() # shape: (N, vocab)
100     log_probs = F.log_softmax(x, dim=1)
101
102     num_classes = x.shape[1]
103     with torch.no_grad():
104         true_dist = torch.full_like(x, smoothing / (num_classes - 1))
105         true_dist.scatter_(1, flat_labels.unsqueeze(1), 1.0 - smoothing)
106
107     loss = torch.mean(torch.sum(-true_dist * log_probs, dim=1))
108     return loss
109
110 # --- Experiment Settings ---
111 def run_experiment(chunk_size, transformation_fn, memory_efficient=True, device="cuda"):
112     # Use smaller dimensions for testing purposes.
113     bsz = 4
114     qlen = 1000 # sequence length
115     hd = 4096 # hidden dimension
116     vocab = 128*1024 # vocabulary size
117     total_tokens = bsz * qlen
118

```

```

119 # Reset and record initial GPU memory stats if using CUDA.
120 if device.type == "cuda":
121     torch.cuda.empty_cache() # Clear cache before measurement.
122     initial_memory = torch.cuda.memory_allocated(device)
123     torch.cuda.reset_peak_memory_stats(device)
124
125 # Create new input tensor and corresponding labels.
126 X = torch.randn(bsz, qlen, hd, device=device, requires_grad=True)
127 labels = torch.randint(0, vocab, (bsz, qlen), device=device)
128
129 # Create a linear layer.
130 linear = nn.Linear(hd, vocab, bias=False).to(device)
131
132 # Run forward/backward pass.
133 if memory_efficient:
134     loss = MemoryEfficientLinear.apply(X, linear, labels, transformation_fn, chunk_size)
135 else:
136     loss = transformation_fn(X.view(-1, hd), linear, labels.view(-1))
137 loss.backward()
138
139 if device.type == "cuda":
140     # Peak and final VRAM usage in MB.
141     peak_memory = torch.cuda.max_memory_allocated(device) / (1024**2)
142     final_memory = torch.cuda.memory_allocated(device) / (1024**2)
143     # Retrieve additional memory statistics.
144     stats = torch.cuda.memory_stats(device)
145     mem_stats = {
146         "allocated_peak": stats["allocated_bytes.all.peak"] / (1024**2),
147         "reserved_peak": stats["reserved_bytes.all.peak"] / (1024**2),
148         "active_peak": stats["active_bytes.all.peak"] / (1024**2),
149         "allocated_current": stats["allocated_bytes.all.current"] / (1024**2),
150         "reserved_current": stats["reserved_bytes.all.current"] / (1024**2),
151         "active_current": stats["active_bytes.all.current"] / (1024**2),
152     }
153 else:
154     peak_memory = None
155     final_memory = None
156     mem_stats = None
157 return loss.item(), peak_memory, final_memory, mem_stats, total_tokens
158
159 if __name__ == "__main__":
160     # Determine device.
161     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
162     if device.type != "cuda":
163         print("CUDA is not available; VRAM measurement requires a GPU. Running on CPU without VRAM stats.")
164
165     # Dictionary mapping function names to transformation functions.
166     transformation_functions = {
167         "Cross Entropy": transformation_function,
168         "Label Smoothing": label_smoothing_transformation_function
169     }
170
171     # Run experiments for each transformation function.
172     results = {}
173     chunk_sizes = [1, 2, 4, 8, 16, 32, 64, 128, 256, 512]
174     for name, func in transformation_functions.items():
175         print(f"--- Testing {name} Transformation Function ---")
176         mem_usage = []
177         final_mem_usage = []
178         loss_vals = []
179         extra_stats = []
180         effective_cs = []
181         for cs in chunk_sizes:
182             loss_val, peak, final_mem, stats, total_tokens = run_experiment(cs, func, memory_efficient=True, device=device)
183             eff_cs = cs if cs is not None else total_tokens // 2
184             effective_cs.append(eff_cs)
185             loss_vals.append(loss_val)
186             if peak is not None:
187                 mem_usage.append(peak)
188                 final_mem_usage.append(final_mem)
189                 extra_stats.append(stats)
190             print(f"Memory-Efficient | Chunk size: {eff_cs:>4} | Loss: {loss_val:>8.4f} | Peak VRAM: {peak:6.2f} MB")
191             print(f"Extra Stats: Allocated Peak: {stats['allocated_peak']:.2f} MB, Reserved Peak: {stats['reserved_peak']:.2f} MB, Active Peak: {stats['active_peak']:.2f} MB")
192         else:
193             print(f"Memory-Efficient | Chunk size: {eff_cs:>4} | Loss: {loss_val:>8.4f}")
194         results[name] = (effective_cs, mem_usage, final_mem_usage)
195
196     # Run the baseline (full, non-chunked) experiment.
197     base_loss, base_peak, base_final, base_stats, _ = run_experiment(chunk_size=None, transformation_fn=func, memory_efficient=True, device=device)
198     if base_peak is not None:
199         print(f"Baseline (full) | Loss: {base_loss:>8.4f} | Peak VRAM: {base_peak:6.2f} MB | Final VRAM: {base_final:6.2f} MB")
200         print(f"Extra Stats: Allocated Peak: {base_stats['allocated_peak']:.2f} MB, Reserved Peak: {base_stats['reserved_peak']:.2f} MB, Active Peak: {base_stats['active_peak']:.2f} MB")

```

```

201     else:
202         print(f"Baseline (full)           | Loss: {base_loss:>8.4f}")
203         results[name + "_baseline"] = (base_peak, base_final)
204         print("\n")
205
206 # Plot the VRAM usage vs. chunk size for each transformation function (if running on CUDA).
207 if device.type == "cuda":
208     plt.figure(figsize=(12, 6))
209     for name in transformation_functions:
210         cs, usage, final_usage = results[name]
211         base_peak, base_final = results[name + "_baseline"]
212         plt.plot(cs, usage, marker="o", label=f"{name} (Peak, Memory-Efficient)")
213         plt.plot(cs, final_usage, marker="s", label=f"{name} (Final, Memory-Efficient)")
214         plt.axhline(y=base_peak, linestyle="--", label=f"{name} (Peak, Baseline)")
215         plt.axhline(y=base_final, linestyle=":", label=f"{name} (Final, Baseline)")
216     plt.xlabel("Chunk Size")
217     plt.ylabel("VRAM Usage (MB)")
218     plt.title("VRAM Usage vs. Chunk Size for Different Transformation Functions")
219     plt.legend()
220     plt.grid(True)
221     plt.tight_layout()
222     plt.show()
223 else:
224     print("Skipping plot because CUDA is not available.")
225

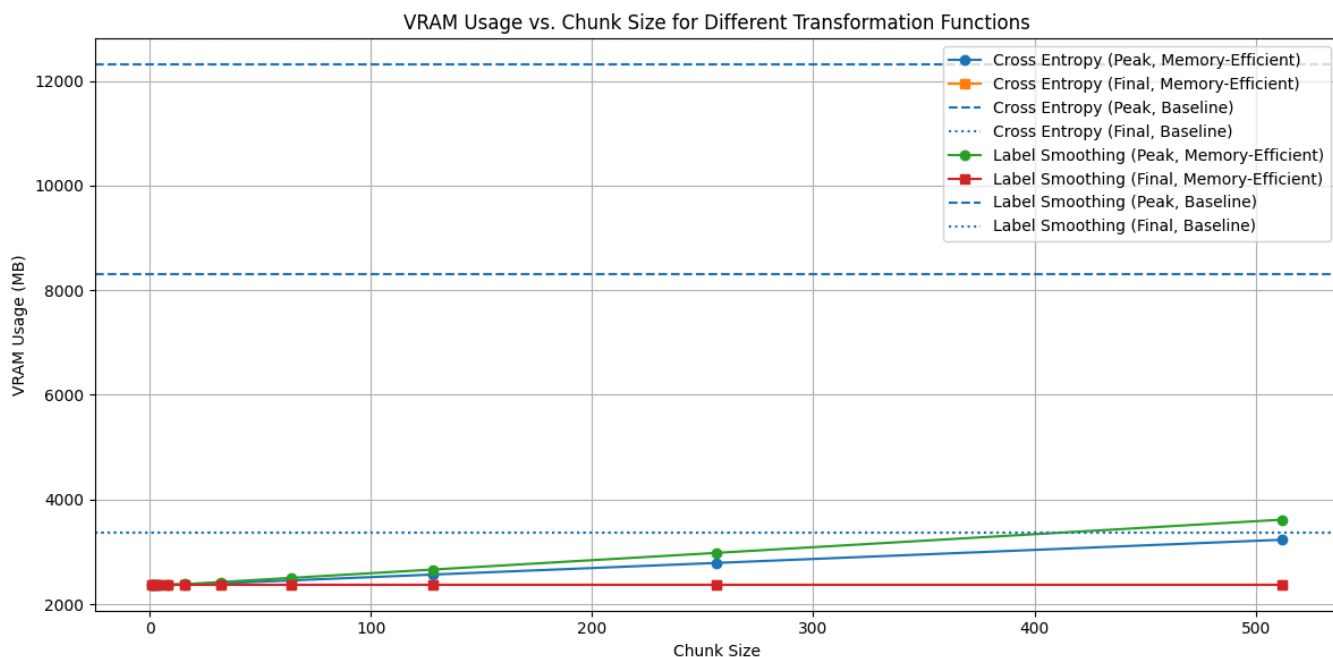
```

```

--- Testing Cross Entropy Transformation Function ---
Memory-Efficient | Chunk size: 1 | Loss: 11.9514 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 2 | Loss: 11.9572 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 4 | Loss: 11.9495 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 8 | Loss: 11.9614 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 16 | Loss: 11.9503 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 32 | Loss: 11.9393 | Peak VRAM: 2397.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2397.62 MB, Reserved Peak: 2612.00 MB, Active Peak: 2397.62 MB
Memory-Efficient | Chunk size: 64 | Loss: 11.9559 | Peak VRAM: 2453.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2453.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2453.12 MB
Memory-Efficient | Chunk size: 128 | Loss: 11.9533 | Peak VRAM: 2564.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2564.62 MB, Reserved Peak: 2680.00 MB, Active Peak: 2564.62 MB
Memory-Efficient | Chunk size: 256 | Loss: 11.9484 | Peak VRAM: 2786.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2786.62 MB, Reserved Peak: 2868.00 MB, Active Peak: 2786.62 MB
Memory-Efficient | Chunk size: 512 | Loss: 11.9662 | Peak VRAM: 3230.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 3230.62 MB, Reserved Peak: 3380.00 MB, Active Peak: 3230.62 MB
Baseline (full) | Loss: 11.9544 | Peak VRAM: 8310.62 MB | Final VRAM: 3365.87 MB
Extra Stats: Allocated Peak: 8310.62 MB, Reserved Peak: 9612.00 MB, Active Peak: 8310.62 MB

--- Testing Label Smoothing Transformation Function ---
Memory-Efficient | Chunk size: 1 | Loss: 11.9629 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 2 | Loss: 11.9707 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 4 | Loss: 11.9458 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 8 | Loss: 11.9567 | Peak VRAM: 2373.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2373.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2373.12 MB
Memory-Efficient | Chunk size: 16 | Loss: 11.9512 | Peak VRAM: 2381.75 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2381.75 MB, Reserved Peak: 2612.00 MB, Active Peak: 2381.75 MB
Memory-Efficient | Chunk size: 32 | Loss: 11.9522 | Peak VRAM: 2421.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2421.62 MB, Reserved Peak: 2612.00 MB, Active Peak: 2421.62 MB
Memory-Efficient | Chunk size: 64 | Loss: 11.9374 | Peak VRAM: 2501.12 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2501.12 MB, Reserved Peak: 2612.00 MB, Active Peak: 2501.12 MB
Memory-Efficient | Chunk size: 128 | Loss: 11.9609 | Peak VRAM: 2660.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2660.62 MB, Reserved Peak: 2808.00 MB, Active Peak: 2660.62 MB
Memory-Efficient | Chunk size: 256 | Loss: 11.9589 | Peak VRAM: 2978.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 2978.62 MB, Reserved Peak: 3124.00 MB, Active Peak: 2978.62 MB
Memory-Efficient | Chunk size: 512 | Loss: 11.9546 | Peak VRAM: 3614.62 MB | Final VRAM: 2373.12 MB
Extra Stats: Allocated Peak: 3614.62 MB, Reserved Peak: 3892.00 MB, Active Peak: 3614.62 MB
Baseline (full) | Loss: 11.9542 | Peak VRAM: 12310.63 MB | Final VRAM: 3365.87 MB
Extra Stats: Allocated Peak: 12310.63 MB, Reserved Peak: 13612.00 MB, Active Peak: 12310.63 MB

```



```

1 import os
2 os.kill(os.getpid(), 9)

```

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import matplotlib.pyplot as plt
5
6
7 class MemoryEfficientLinear(torch.autograd.Function):
8     @staticmethod
9     def forward(ctx, X, linear, labels, forward_function, chunk_size=None):
10
11         bsz, seq_len, hd = X.shape
12         # Flatten the batch and sequence dims.
13         X_flat = X.reshape(-1, hd)          # shape: (bsz * seq_len, hd)
14         labels_flat = labels.reshape(-1)     # shape: (bsz * seq_len,)
15         total_tokens = X_flat.shape[0]
16
17         if chunk_size is None:
18             chunk_size = total_tokens // 2 if total_tokens > 1 else 1
19         ctx.chunk_size = chunk_size
20
21         # Determine valid tokens (assume -100 indicates positions to ignore).
22         valid_mask = (labels_flat != -100)
23         total_valid_tokens = valid_mask.sum().to(dtype=X.dtype)
24         ctx.total_valid_tokens = total_valid_tokens
25
26         total_loss = 0.0
27         grad_chunks = []
28
29         # Split the flattened inputs into chunks.
30         X_chunks = torch.split(X_flat, chunk_size, dim=0)
31         label_chunks = torch.split(labels_flat, chunk_size, dim=0)
32
33         for chunk, lab in zip(X_chunks, label_chunks):
34             # Count valid tokens in this chunk.
35             valid_tokens_chunk = (lab != -100).sum().to(dtype=X.dtype)
36             # Calculate the scaling factor for this chunk.
37             scaling_factor = valid_tokens_chunk / total_valid_tokens if total_valid_tokens > 0 else 0.0
38
39             # Here we define a lambda that wraps the transformation_function.
40             compute_chunk_loss = lambda chunk_in: forward_function(chunk_in, linear, lab) * scaling_factor
41
42             # Compute both the loss value and gradient for the chunk.
43             chunk_grad, chunk_loss = torch.func.grad_and_value(
44                 compute_chunk_loss, argnums=0, has_aux=False
45             )(chunk)
46             total_loss += chunk_loss
47             grad_chunks.append(chunk_grad)
48
49             # Save the precomputed gradients to be used in backward.
50             precomputed_grad = torch.cat(grad_chunks, dim=0)
51             ctx.precomputed_grad = precomputed_grad
52             ctx.input_shape = X.shape # to reshape the gradient back to original dimensions
53         return total_loss
54
55     @staticmethod
56     def backward(ctx, dY):
57         # Retrieve the precomputed gradient and multiply by the upstream gradient.
58         precomputed_grad = ctx.precomputed_grad
59         grad_X_flat = precomputed_grad * dY
60         grad_X = grad_X_flat.view(ctx.input_shape)
61
62         return grad_X, None, None, None, None
63 # --- Transformation Functions ---
64
65 # (1) Original Transformation Function using CrossEntropyLoss.
66 def transformation_function(flat_batch, linear, flat_labels):
67     """
68     Computes standard cross-entropy loss.
69
70     Parameters:
71         flat_batch: Tensor of shape (N, hd)
72         linear: Linear layer mapping from hd to vocab
73         flat_labels: Tensor of shape (N,), containing target class indices.
74
75     Returns:
76         loss: Scalar tensor representing the cross-entropy loss.
77     """
78     x = linear(flat_batch).float() # shape: (N, vocab)
79
80     ce_loss = nn.CrossEntropyLoss(ignore_index=-100, reduction="mean")
81     loss = ce_loss(x, flat_labels)

```

```

83     return loss
84
85 # (2) Label-Smoothing Transformation Function.
86 def label_smoothing_transformation_function(flat_batch, linear, flat_labels, smoothing=0.1):
87     """
88     Computes a label-smoothed cross-entropy loss.
89
90     Parameters:
91         flat_batch: Tensor of shape (N, hd)
92         linear: Linear layer mapping from hd to vocab
93         flat_labels: Tensor of shape (N,), containing target class indices.
94         smoothing: The label smoothing factor (default is 0.1)
95
96     Returns:
97         loss: Scalar tensor representing the label-smoothed loss.
98     """
99     x = linear(flat_batch).float() # shape: (N, vocab)
100     log_probs = F.log_softmax(x, dim=1)
101
102     num_classes = x.shape[1]
103     with torch.no_grad():
104         true_dist = torch.full_like(x, smoothing / (num_classes - 1))
105         true_dist.scatter_(1, flat_labels.unsqueeze(1), 1.0 - smoothing)
106
107     loss = torch.mean(torch.sum(-true_dist * log_probs, dim=1))
108     return loss

```

```

1 from transformers.modeling_outputs import (
2     CausalLMOutputWithPast,
3 )
4 import torch
5 import torch.nn as nn
6 import torch.nn.functional as F
7 from typing import Callable, List, Optional, Tuple, Union
8 def forward(
9     self,
10    input_ids: torch.LongTensor = None,
11    attention_mask: Optional[torch.Tensor] = None,
12    position_ids: Optional[torch.LongTensor] = None,
13    past_key_values = None,
14    inputs_embeds: Optional[torch.FloatTensor] = None,
15    labels: Optional[torch.LongTensor] = None,
16    use_cache: Optional[bool] = None,
17    output_attentions: Optional[bool] = None,
18    output_hidden_states: Optional[bool] = None,
19    return_dict: Optional[bool] = None,
20    cache_position: Optional[torch.LongTensor] = None,
21    logits_to_keep: Union[int, torch.Tensor] = 0,
22    **kwargs#): Unpack[KwargsForCausalLM],
23    ): #-> Union[Tuple, CausalLMOutputWithPast]:
24
25    output_attentions = output_attentions if output_attentions is not None else self.config.output_attentions
26    output_hidden_states = (
27        output_hidden_states if output_hidden_states is not None else self.config.output_hidden_states
28    )
29    return_dict = return_dict if return_dict is not None else self.config.use_return_dict
30    #print(10000)
31    # decoder outputs consists of (dec_features, layer_state, dec_hidden, dec_attn)
32    outputs = self.model(
33        input_ids=input_ids,
34        attention_mask=attention_mask,
35        position_ids=position_ids,
36        past_key_values=past_key_values,
37        inputs_embeds=inputs_embeds,
38        use_cache=use_cache,
39        output_attentions=output_attentions,
40        output_hidden_states=output_hidden_states,
41        return_dict=return_dict,
42        cache_position=cache_position,
43        **kwargs,
44    )
45    num_chunks = 16
46    hidden_states = outputs[0]
47    # Only compute necessary logits, and do not upcast them to float if we are not computing the loss
48    slice_indices = slice(-logits_to_keep, None) if isinstance(logits_to_keep, int) else logits_to_keep
49    #my_param = kwargs.get("num_items_in_batch", -1)
50    labels = nn.functional.pad(labels, (0, 1), value=-100)
51    shift_labels = labels[..., 1:].contiguous()
52    loss = MemoryEfficientLinear.apply(hidden_states[:, slice_indices, :], self.lm_head, shift_labels, transformation_fur
53
54    # logits = self.lm_head(hidden_states[:, slice_indices, :])
55    # ce_loss = nn.CrossEntropyLoss(reduction="mean")

```

```

56 # loss = ce_loss(logits.view(-1, logits.shape[-1]), shift_labels.view(-1))
57 return (loss, )
58 # loss = None
59 # if labels is not None:
60 #     loss = self.loss_function(logits=logits, labels=labels, vocab_size=self.config.vocab_size, **kwargs)
61
62 # if not return_dict:
63 #     output = (logits,) + outputs[1:]
64 #     return (loss,) + output if loss is not None else output
65
66 # return CausalLMOutputWithPast(
67 #     loss=loss,
68 #     logits=logits,
69 #     past_key_values=outputs.past_key_values,
70 #     hidden_states=outputs.hidden_states,
71 #     attentions=outputs.attentions,
72 # )
73
74
75 import transformers.models.llama.modeling_llama
76 transformers.models.llama.modeling_llama.LlamaForCausalLM.forward = forward
77 from trl import SFTTrainer, SFTConfig
78 def compute_loss(self, model, inputs, return_outputs=False, num_items_in_batch=None):
79     """
80     Compute training loss and additionally compute token accuracies
81     """
82     (loss, outputs) = Trainer.compute_loss(
83         self, model = model, inputs = inputs, return_outputs=True, num_items_in_batch=num_items_in_batch
84     )
85
86     # # Compute token accuracy if we have labels and if the model is not using Liger (no logits)
87     # if "labels" in inputs and not self.args.use_liger:
88     #     shift_logits = outputs.logits[..., :-1, :].contiguous()
89     #     shift_labels = inputs["labels"][..., 1:].contiguous()
90
91     #     # Get predictions
92     #     predictions = shift_logits.argmax(dim=-1)
93
94     #     # Create mask for non-padding tokens (assuming ignore_index is -100)
95     #     mask = shift_labels != -100
96
97     #     # Calculate accuracy only on non-padding tokens
98     #     correct_predictions = (predictions == shift_labels) & mask
99     #     total_tokens = mask.sum()
100     #     correct_tokens = correct_predictions.sum()
101
102     #     # Gather the correct_tokens and total_tokens across all processes
103     #     correct_tokens = self.accelerator.gather_for_metrics(correct_tokens)
104     #     total_tokens = self.accelerator.gather_for_metrics(total_tokens)
105
106     #     # Compute the mean token accuracy and log it
107     #     accuracy = (correct_tokens.sum() / total_tokens.sum()).item() if total_tokens.sum() > 0 else 0.0
108     #     self._metrics["mean_token_accuracy"].append(accuracy)
109
110     return (loss, outputs) if return_outputs else loss
111
112
113 SFTTrainer.compute_loss = compute_loss
114 from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig, Trainer
115 from peft import get_peft_model, LoraConfig, TaskType
116 import os
117 from trl import SFTTrainer, SFTConfig
118 os.environ["HF_HUB_ENABLE_HF_TRANSFER"] = "1"
119 os.environ["CUDA_VISIBLE_DEVICES"] = "0"
120 os.environ["PYTORCH_CUDA_ALLOC_CONF"] = \
121     "expandable_segments:True,\\"
122     "roundup_power2_divisions:[32:256,64:128,256:64,>:32]"
123
124 max_seq_length = 1024
125 #torch.set_default_dtype(torch.float16)
126 model_name = "unsloth/Llama-3.2-1B-Instruct"
127 dtype = torch.float16
128 # bnb_config = BitsAndBytesConfig(
129 #     load_in_4bit = True,
130 #     bnb_4bit_use_double_quant = True,
131 #     bnb_4bit_quant_type = "nf4",
132 #     bnb_4bit_compute_dtype = dtype,
133 # )
134 model = AutoModelForCausalLM.from_pretrained(
135     model_name,
136     device_map = "auto",
137     attn_implementation = "sdpa",


```

```

138     #quantization_config = bnb_config,
139 )
140 tokenizer = AutoTokenizer.from_pretrained(model_name)
141 tokenizer.padding_side = "right"
142 tokenizer.pad_token = tokenizer.eos_token
143
144 lora_config = LoraConfig(
145     r = 32,
146     lora_alpha = 64,
147     target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
148                     "gate_proj", "up_proj", "down_proj"],
149     lora_dropout = 0,
150     bias = "none",
151     task_type = TaskType.CAUSAL_LM,
152 )
153
154 # Get LoRA and setup model
155 model = get_peft_model(model, lora_config)
156 with torch.no_grad():
157     for name, param in model.named_parameters():
158         if ".lora_A." in name or ".lora_B." in name : param.requires_grad_(True)
159         else: param.requires_grad_(False)
160
161 # Currently GC will cause torch.compile to be disabled, so disable it
162 # model.gradient_checkpointing_enable()
163 model.lm_head.weight.requires_grad = True
164 model.enable_input_require_grads()
165
166 # Get dataset
167 from datasets import load_dataset
168 from trl import SFTTrainer, SFTConfig
169 url = "https://huggingface.co/datasets/laion/OIG/resolve/main/unified_chip2.jsonl"
170 dataset = load_dataset("json", data_files = {"train" : url}, split = "train[:10%]")
171
172
173 trainer = SFTTrainer(
174     model = model,
175     train_dataset = dataset,
176     processing_class = tokenizer,
177     args = SFTConfig(
178         per_device_train_batch_size = 1,
179         gradient_accumulation_steps = 2,
180         warmup_steps = 1,
181         max_steps = 60,
182         logging_steps = 1,
183         output_dir = "outputs",
184         seed = 3407,
185         fp16 = model.get_input_embeddings().weight.dtype == torch.float16,
186         bf16 = model.get_input_embeddings().weight.dtype == torch.bfloat16,
187         max_seq_length = max_seq_length,
188         report_to = "none", # For W&B
189         dataset_num_proc = 4,
190     ),
191 )
192
193 trainer.train()
194

```



 /usr/local/lib/python3.11/dist-packages/huggingface\_hub/utils/\_auth.py:94: UserWarning:  
The secret 'HF\_TOKEN' does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>),  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.  
warnings.warn(  

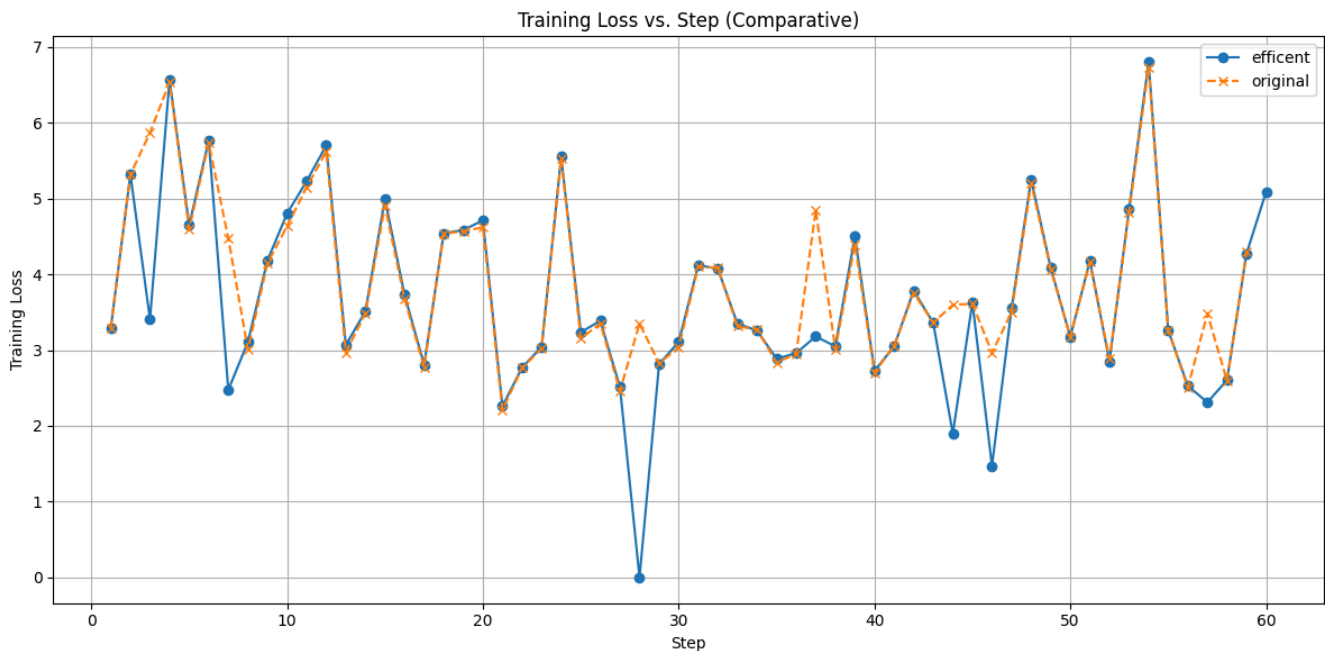
[60/60 00:57, Epoch 0/1]

#### Step Training Loss

1	3.296300
2	5.318100
3	3.403900
4	6.565700
5	4.644800
6	5.757900
7	2.461400
8	3.091100
9	4.175800
10	4.782600
11	5.220400
12	5.690000
13	3.060300
14	3.508500
15	4.989200
16	3.723100
17	2.802300
18	4.543800
19	4.581000
20	4.711800
21	2.256900
22	2.770200
23	3.034600
24	5.553500
25	3.238300
26	3.381400
27	2.510700
28	0.000000
29	2.814000
30	3.110800
31	4.122300
32	4.073800
33	3.350100
34	3.256200
35	2.885400
36	2.962600
37	3.179400
38	3.047500
39	4.500500
40	2.725700
41	3.044000
42	3.788100
43	3.357400
44	1.895200
45	3.629000

46	1.469600
47	3.551900
48	5.253300
49	4.088900
50	3.170800
51	4.178600
52	2.851900
53	4.861800
54	6.804000
55	3.261000
56	2.530000
57	2.310100
58	2.605000
59	4.274100
60	5.092100

```
TrainOutput(global_step=60, training_loss=3.685314436753591, metrics={'train_runtime': 58.4604,
'train_samples_per_second': 2.053, 'train_steps_per_second': 1.026, 'total_flos': 69180575662080.0, 'train_loss':
3.685314436753591})
```



```
1 del model
2 import gc
3 gc.collect()
4 torch.cuda.empty_cache()
```

GRPO\_memory\_efficient\_linear\_works below:

```
1 def transformation_function_grpo(new_hidden_states, old_hidden_states, linear, input_ids, completion_mask, advantages, ep
2
3     nlogits = linear(new_hidden_states).float()
4     ologits = linear(old_hidden_states).float()
5
6     selected_logits_n = torch.gather(nlogits, dim=-1, index=input_ids.unsqueeze(-1)).squeeze(-1)
7     selected_logits_o = torch.gather(ologits, dim=-1, index=input_ids.unsqueeze(-1)).squeeze(-1)
8
9     per_token_logps_n = selected_logits_n - torch.logsumexp(selected_logits_n, dim=-1)
10    per_token_logps_o = selected_logits_o - torch.logsumexp(selected_logits_o, dim=-1)
11    # this needs a ref_model (we should replace this with ref model for real impl)
12    per_token_kl = torch.exp(per_token_logps_o - per_token_logps_n) - (per_token_logps_o - per_token_logps_n) - 1
13    coef_1 = torch.exp(per_token_logps_n - per_token_logps_o.detach())
14    #print("coef", coef_1)
15    coef_2 = torch.clamp(coef_1, 1 - epsilon, 1 + epsilon)
16    per_token_loss1 = coef_1 * advantages.unsqueeze(1)
```

```

17     per_token_loss2 = coef_2 * advantages.unsqueeze(1)
18     per_token_loss = -torch.min(per_token_loss1, per_token_loss2)
19     if beta != 0.0:
20         per_token_loss = per_token_loss1 + beta * per_token_kl
21     loss = (per_token_loss1 * completion_mask).sum() / completion_mask.sum()
22     #print(loss)
23     return loss
24
25
26 import torch
27
28 class GRPO_memory_efficient_linear(torch.autograd.Function):
29     @staticmethod
30     def forward(ctx, new_hidden_states, old_hidden_states, linear, input_ids, mask, advantages, beta = 0, epsilon = 0.2,
31                 device = new_hidden_states.device):
32
33         # Get original dimensions.
34         B, T, H = new_hidden_states.shape
35         # print(new_hidden_states.shape)
36         # print(old_hidden_states.shape)
37         # print(input_ids.shape)
38         # print(mask.shape)
39         # print(advantages.shape)
40         # Exclude the last time step from new/old hidden states and associated tensors.
41         new_hs_flat = new_hidden_states[:, :-1, :].reshape(-1, H)          # shape: [B*(T-1), H]
42         old_hs_flat = old_hidden_states[:, :-1, :].reshape(-1, H)          # shape: [B*(T-1), H]
43         ids_flat = input_ids.reshape(-1)                                   # shape: [B*(T)]
44         mask_flat = mask.reshape(-1)                                       # shape: [B*(T)]
45         adv_flat = advantages                                              # shape: [B]
46
47         # Allocate a gradient tensor for the flattened new_hidden_states.
48         grad_flat = torch.empty_like(new_hs_flat)
49         accumulated_loss = torch.zeros(1, device=device)
50
51         # Chunk the flattened tensors along the first (combined) dimension.
52         grad_chunks = torch.chunk(grad_flat, n_chunks, dim=0)
53         new_chunks = torch.chunk(new_hs_flat, n_chunks, dim=0)
54         old_chunks = torch.chunk(old_hs_flat, n_chunks, dim=0)
55         ids_chunks = torch.chunk(ids_flat, n_chunks, dim=0)
56         mask_chunks = torch.chunk(mask_flat, n_chunks, dim=0)
57         adv_chunks = torch.chunk(adv_flat, n_chunks, dim=0)
58
59         # Process each chunk using torch.func.grad_and_value.
60         for grad_chunk, new_chunk, old_chunk, ids_chunk, mask_chunk, adv_chunk in zip(
61             grad_chunks, new_chunks, old_chunks, ids_chunks, mask_chunks, adv_chunks):
62             # Ensure new_chunk is detached and requires gradient.
63             new_chunk = new_chunk.detach().requires_grad_()
64
65             # Compute gradients and loss in one shot.
66             (chunk_grad, chunk_loss) = torch.func.grad_and_value(
67                 transformation_function_grpo,
68                 argnums=0,
69                 has_aux=False,
70             )(new_chunk, old_chunk, linear, ids_chunk, mask_chunk, adv_chunk, epsilon, beta)
71
72             grad_chunk.copy_(chunk_grad)
73             accumulated_loss += chunk_loss # accumulate the unscaled loss
74
75         # Average gradients and loss over all chunks.
76         grad_flat.div_(n_chunks)
77         accumulated_loss.div_(n_chunks)
78
79         # Reassemble the gradient for new_hidden_states:
80         # Create a zero tensor for full new_hidden_states gradient.
81         grad_new_hidden = torch.zeros_like(new_hidden_states)
82         # Insert the computed gradient for tokens [:-1].
83         grad_new_hidden[:, :-1, :] = grad_flat.view(B, T - 1, H)
84
85         # Save the gradient for backward.
86         ctx.save_for_backward(grad_new_hidden)
87         return accumulated_loss
88
89     @staticmethod
90     def backward(ctx, grad_output, dcompletion_length=None, dmean_kl=None):
91         (grad_input,) = ctx.saved_tensors
92         return grad_input, None, None, None, None, None, None, None
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

```

```

4
5 from typing import Callable, List, Optional, Tuple, Union
6
7 def forward_grpo(
8     self,
9     input_ids: torch.LongTensor = None,
10    attention_mask: Optional[torch.Tensor] = None,
11    position_ids: Optional[torch.LongTensor] = None,
12    past_key_values = None,
13    inputs_embeds: Optional[torch.FloatTensor] = None,
14    labels: Optional[torch.LongTensor] = None,
15    use_cache: Optional[bool] = None,
16    output_attentions: Optional[bool] = None,
17    output_hidden_states: Optional[bool] = None,
18    return_dict: Optional[bool] = None,
19    cache_position: Optional[torch.LongTensor] = None,
20    logits_to_keep: Union[int, torch.Tensor] = 0,
21    **kwargs,
22 ) -> Union[Tuple, CausalLMOutputWithPast]:
23
24    output_attentions = output_attentions if output_attentions is not None else self.config.output_attentions
25    output_hidden_states = (
26        output_hidden_states if output_hidden_states is not None else self.config.output_hidden_states
27    )
28    return_dict = return_dict if return_dict is not None else self.config.use_return_dict
29
30    # decoder outputs consists of (dec_features, layer_state, dec_hidden, dec_attn)
31    outputs = self.model(
32        input_ids=input_ids,
33        attention_mask=attention_mask,
34        position_ids=position_ids,
35        past_key_values=past_key_values,
36        inputs_embeds=inputs_embeds,
37        use_cache=use_cache,
38        output_attentions=output_attentions,
39        output_hidden_states=output_hidden_states,
40        return_dict=return_dict,
41        cache_position=cache_position,
42        **kwargs,
43    )
44
45    hidden_states = outputs[0]
46    # Only compute necessary logits, and do not upcast them to float if we are not computing the loss
47    slice_indices = slice(-logits_to_keep, None) if isinstance(logits_to_keep, int) else logits_to_keep
48
49    if hidden_states.shape[0] * hidden_states.shape[1] <= 512:
50        logits = self.lm_head(hidden_states[:, slice_indices, :])
51    else:
52        logits = hidden_states[:, slice_indices, :]
53
54    loss = None
55
56    return CausalLMOutputWithPast(
57        loss=loss,
58        logits=logits,
59        past_key_values=outputs.past_key_values,
60        hidden_states=outputs.hidden_states,
61        attentions=outputs.attentions,
62    )
63 import transformers.models.llama.modeling_llama
64 transformers.models.llama.modeling_llama.LlamaForCausalLM.forward = forward_grpo
65
66
67 def _get_per_token_logps(self, model, input_ids, attention_mask, logits_to_keep):
68     # We add 1 to `logits_to_keep` because the last logits of the sequence is later excluded
69     # logits = model(input_ids=input_ids, attention_mask=attention_mask, logits_to_keep=logits_to_keep + 1).logits
70     # logits = hidden_states[:, :-1, :] # (B, L-1, V), exclude the last logit: it corresponds to the next token pred
71
72     # input_ids = input_ids[:, -logits_to_keep:]
73     # For transformers<=4.48, logits_to_keep argument isn't supported, so here we drop logits ourselves.
74     # See https://github.com/huggingface/trl/issues/2770
75     # logits = logits[:, -logits_to_keep:]
76     return None
77
78 from trl import GRPOConfig, GRPOTrainer
79 GRPOTrainer._get_per_token_logps = _get_per_token_logps
80
81 def compute_loss(self, model, inputs, return_outputs=False, num_items_in_batch=None):
82     if return_outputs:
83         raise ValueError("The GRPOTrainer does not support returning outputs")
84     # Compute the per-token log probabilities for the model
85

```

```

86 prompt_ids, prompt_mask = inputs["prompt_ids"], inputs["prompt_mask"]
87 completion_ids, completion_mask = inputs["completion_ids"], inputs["completion_mask"]
88 input_ids = torch.cat([prompt_ids, completion_ids], dim=1)
89 attention_mask = torch.cat([prompt_mask, completion_mask], dim=1)
90 logits_to_keep = completion_ids.size(1) # we only need to compute the logits for the completion tokens
91
92 #per_token_logits = self._get_per_token_logits(model, input_ids, attention_mask, logits_to_keep)
93
94 # Compute the KL divergence between the model and the reference model
95 # if self.beta != 0.0:
96 #     ref_per_token_logits = inputs["ref_per_token_logits"]
97 #     per_token_kl = (
98 #         torch.exp(ref_per_token_logits - per_token_logits) - (ref_per_token_logits - per_token_logits) - 1
99 #     )
100
101 # Compute the loss
102 advantages = inputs["advantages"]
103 # When using num_iterations == 1, old_per_token_logits == per_token_logits, so we can skip it's computation (see
104 # _generate_and_score_completions) and use per_token_logits.detach() instead.
105 # old_per_token_logits = inputs["old_per_token_logits"] if self.num_iterations > 1 else per_token_logits.detach()
106 # coef_1 = torch.exp(per_token_logits - old_per_token_logits)
107 # coef_2 = torch.clamp(coef_1, 1 - self.epsilon, 1 + self.epsilon)
108 # per_token_loss1 = coef_1 * advantages.unsqueeze(1)
109 # per_token_loss2 = coef_2 * advantages.unsqueeze(1)
110 # per_token_loss = -torch.min(per_token_loss1, per_token_loss2)
111 # if self.beta != 0.0:
112 #     per_token_loss = per_token_loss + self.beta * per_token_kl
113 # loss = (per_token_loss * completion_mask).sum() / completion_mask.sum()
114 new_hidden_states = model(input_ids = input_ids, logits_to_keep = logits_to_keep + 1).logits
115 with torch.no_grad():
116     old_hidden_states = model(input_ids = input_ids, logits_to_keep = logits_to_keep + 1).logits
117
118 n_chunks = 2
119
120 loss = GRPO_memory_efficient_linear.apply(
121     new_hidden_states, old_hidden_states, model.lm_head,
122     completion_ids, completion_mask, advantages, self.beta,
123     0.2,
124     n_chunks,
125 )
126 # No metrics for this simple impl
127 # Log the metrics
128 # mode = "eval" if self.control.should_evaluate else "train"
129
130 # if self.beta != 0.0:
131 #     mean_kl = (per_token_kl * completion_mask).sum() / completion_mask.sum()
132 #     self._metrics[mode]["kl"].append(self.accelerator.gather_for_metrics(mean_kl).mean().item())
133
134 # is_clipped = (per_token_loss1 < per_token_loss2).float()
135 # clip_ratio = (is_clipped * completion_mask).sum() / completion_mask.sum()
136 # self._metrics[mode]["clip_ratio"].append(self.accelerator.gather_for_metrics(clip_ratio).mean().item())
137 return loss
138
139 GRPOTrainer.compute_loss = compute_loss

```

```

1 import re
2 from datasets import load_dataset, Dataset
3
4 # Load and prep dataset
5 SYSTEM_PROMPT = """
6 Respond in the following format:
7 <reasoning>
8 ...
9 </reasoning>
10 <answer>
11 ...
12 </answer>
13 """
14
15 XML_COT_FORMAT = """\
16 <reasoning>
17 {reasoning}
18 </reasoning>
19 <answer>
20 {answer}
21 </answer>
22 """
23
24 def extract_xml_answer(text: str) -> str:
25     answer = text.split("<answer>")[-1]
26     answer = answer.split("</answer>")[0]
27     return answer.strip()

```

```

28
29 def extract_hash_answer(text: str) -> str | None:
30     if "####" not in text:
31         return None
32     return text.split("####")[1].strip()
33
34 # uncomment middle messages for 1-shot prompting
35 def get_gsm8k_questions(split = "train") -> Dataset:
36     data = load_dataset('openai/gsm8k', 'main')[split] # type: ignore
37     data = data.map(lambda x: { # type: ignore
38         'prompt': [
39             {'role': 'system', 'content': SYSTEM_PROMPT},
40             {'role': 'user', 'content': x['question']}
41         ],
42         'answer': extract_hash_answer(x['answer'])
43     }) # type: ignore
44     return data # type: ignore
45
46 dataset = get_gsm8k_questions()
47
48 # Reward functions
49 def correctness_reward_func(prompts, completions, answer, **kwargs) -> list[float]:
50     responses = [completion[0]['content'] for completion in completions]
51     q = prompts[0][-1]['content']
52     extracted_responses = [extract_xml_answer(r) for r in responses]
53     print('-'*20, f"Question:\n{q}", f"\nAnswer:\n{answer[0]}", f"\nResponse:\n{responses[0]}", f"\nExtracted:\n{extracted_responses}")
54     return [2.0 if r == a else 0.0 for r, a in zip(extracted_responses, answer)]
55
56 def int_reward_func(completions, **kwargs) -> list[float]:
57     responses = [completion[0]['content'] for completion in completions]
58     extracted_responses = [extract_xml_answer(r) for r in responses]
59     return [0.5 if r.isdigit() else 0.0 for r in extracted_responses]
60
61 def strict_format_reward_func(completions, **kwargs) -> list[float]:
62     """Reward function that checks if the completion has a specific format."""
63     pattern = r"^(?<reasoning>\n.*?</reasoning>\n<answer>\n.*?</answer>\n$)
64     responses = [completion[0]["content"] for completion in completions]
65     matches = [re.match(pattern, r) for r in responses]
66     return [0.5 if match else 0.0 for match in matches]
67
68 def soft_format_reward_func(completions, **kwargs) -> list[float]:
69     """Reward function that checks if the completion has a specific format."""
70     pattern = r"^(?<reasoning>.*?</reasoning>\n.*?<answer>.*?</answer>)"
71     responses = [completion[0]["content"] for completion in completions]
72     matches = [re.match(pattern, r) for r in responses]
73     return [0.5 if match else 0.0 for match in matches]
74
75 def count_xml(text) -> float:
76     count = 0.0
77     if text.count("<reasoning>\n") == 1:
78         count += 0.125
79     if text.count("\n</reasoning>\n") == 1:
80         count += 0.125
81     if text.count("\n<answer>\n") == 1:
82         count += 0.125
83     count -= len(text.split("\n</answer>\n")[-1])*0.001
84     if text.count("\n</answer>") == 1:
85         count += 0.125
86     count -= (len(text.split("\n</answer>")[-1]) - 1)*0.001
87     return count
88
89 def xmlcount_reward_func(completions, **kwargs) -> list[float]:
90     contents = [completion[0]["content"] for completion in completions]
91     return [count_xml(c) for c in contents]

```



README.md: 100%	7.94k/7.94k [00:00<00:00, 594kB/s]
train-00000-of-00001.parquet: 100%	2.31M/2.31M [00:00<00:00, 14.2MB/s]
test-00000-of-00001.parquet: 100%	419k/419k [00:00<00:00, 29.5MB/s]
Generating train split: 100%	7473/7473 [00:00<00:00, 69116.15 examples/s]
Generating test split: 100%	1319/1319 [00:00<00:00, 51867.95 examples/s]
Map: 100%	7473/7473 [00:00<00:00, 11381.94 examples/s]

```

1 from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig, Trainer
2 from peft import get_peft_model, LoraConfig, TaskType
3 import os
4
5 os.environ["HF_HUB_ENABLE_HF_TRANSFER"] = "1"
6 os.environ["CUDA_VISIBLE_DEVICES"] = "0"

```

```

7 os.environ["PYTORCH_CUDA_ALLOC_CONF"] = \
8     "expandable_segments:True,"\
9     "roundup_power2_divisions:[32:256,64:128,256:64,>:32]"
10
11 max_seq_length = 1024
12 #torch.set_default_dtype(torch.float16)
13 model_name = "unsloth/Llama-3.2-1B-Instruct"
14 dtype = torch.float16
15 # bnb_config = BitsAndBytesConfig(
16 #     load_in_4bit = True,
17 #     bnb_4bit_use_double_quant = True,
18 #     bnb_4bit_quant_type = "nf4",
19 #     bnb_4bit_compute_dtype = dtype,
20 # )
21 model = AutoModelForCausalLM.from_pretrained(
22     model_name,
23     device_map = "auto",
24     attn_implementation = "sdpa",
25     #quantization_config = bnb_config,
26 )
27 tokenizer = AutoTokenizer.from_pretrained(model_name)
28 tokenizer.padding_side = "right"
29 tokenizer.pad_token = tokenizer.eos_token
30
31 lora_config = LoraConfig(
32     r = 32,
33     lora_alpha = 64,
34     target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
35                      "gate_proj", "up_proj", "down_proj"],
36     lora_dropout = 0,
37     bias = "none",
38     task_type = TaskType.CAUSAL_LM,
39 )
40
41 # Get LoRA and setup model
42 model = get_peft_model(model, lora_config)
43 with torch.no_grad():
44     for name, param in model.named_parameters():
45         if ".lora_A." in name or ".lora_B." in name : param.requires_grad_(True)
46         else: param.requires_grad_(False)
47
48 # Currently GC will cause torch.compile to be disabled, so disable it
49 # model.gradient_checkpointing_enable()
50 model.lm_head.weight.requires_grad = True
51 model.enable_input_require_grads()
52
53
54 training_args = GRPOConfig(
55     use_vllm = False,
56     learning_rate = 5e-6,
57     adam_beta1 = 0.9,
58     adam_beta2 = 0.99,
59     weight_decay = 0.1,
60     warmup_ratio = 0.1,
61     lr_scheduler_type = "cosine",
62     optim = "paged_adamw_8bit",
63     logging_steps = 1,
64     fp16 = model.get_input_embeddings().weight.dtype == torch.float16,
65     bf16 = model.get_input_embeddings().weight.dtype == torch.bfloat16,
66     per_device_train_batch_size = 4,
67     gradient_accumulation_steps = 1,
68     num_generations = 4,
69     max_prompt_length = 256,
70     max_completion_length = 200,
71     max_steps = 40,
72     save_steps = 20,
73     max_grad_norm = 0.1,
74     report_to = "none",
75     output_dir = "outputs",
76 )
77
78 trainer = GRPOTrainer(
79     model = model,
80     processing_class = tokenizer,
81     reward_funcs = [
82         xmlcount_reward_func,
83         soft_format_reward_func,
84         strict_format_reward_func,
85         int_reward_func,
86         correctness_reward_func,
87     ],
88     args = training_args,

```