```
1 import os
 2 DEVELOPMENT_MODE = False
 3 # Detect if we're running in Google Colab
 4 try:
5
       import google.colab
      IN COLAB = True
      print("Running as a Colab notebook")
8 except:
      IN COLAB = False
9
10
11 # Install if in Colab
12 #if IN COLAB:
13 #
      %pip install transformer_lens
14 #
       %pip install circuitsvis
15 #
      %pip install nnsight
16
      # Install a faster Node version
17 #
       !curl -fsSL https://deb.nodesource.com/setup_16.x | sudo -E bash -; sudo apt-get install
19 \# Hot reload in development mode \& not running on the CD
20 if not IN_COLAB:
      from IPython import get_ipython
21
22
      ip = get_ipython()
23
      if not ip.extension_manager.loaded:
24
           ip.extension_manager.load('autoreload')
25
           %autoreload 2
27 IN_GITHUB = os.getenv("GITHUB_ACTIONS") == "true"
Running as a Colab notebook
 1 ! pip install jaxtyping
 2 ! pip install git+https://github.com/TransformerLensOrg/TransformerLens
 3 ! pip install nnsight
 4 ! pip install datasets
 5 ! pip install BitsAndBytes
 1 import torch
 2 import torch.nn as nn
 3 from typing import Dict, List, Optional, Union, NamedTuple
 4 from transformers import PreTrainedModel
 5 import einops
 6 from dataclasses import dataclass
 7 from typing import Literal
 8 import logging
10 # Named tuple for storing both logits and loss, matching TransformerLens's Output class
11 class Output(NamedTuple):
       logits: torch.Tensor # [batch, pos, d_vocab]
13
      loss: Optional[torch.Tensor] = None # Either scalar or [batch, pos-1]
15 @dataclass
16 class ActivationCache:
      """A wrapper around a dictionary of activations with helper functions for analysis.
17
18
      This provides similar functionality to TransformerLens's ActivationCache but works with
19
20
      Hugging Face models. Key operations include:
21
      - Storing activations from model forward passes
22
      - Decomposing residual streams
23
      - Computing head/neuron contributions
24
       - Applying LayerNorm scaling
25
26
      cache_dict: Dict[str, torch.Tensor]
27
      model: PreTrainedModel
      has_batch_dim: bool = True
28
29
30
            _getitem__(self, key: str) -> torch.Tensor:
           """"Get activation by key name."""
31
32
           return self.cache_dict[key]
33
34
      def keys(self):
35
           """Get all activation keys."""
36
           return self.cache_dict.keys()
37
38
      def decompose_resid(self, layer: Optional[int] = None, mlp_input: bool = False,
39
                          apply_ln: bool = False, pos_slice: Optional[Union[slice, int]] = None,
40
                          mode: Literal['all', 'mlp', 'attn'] = 'all') -> torch.Tensor:
41
           """Decompose residual stream into components from each layer.
42
43
           Args:
               layer: Which layer to decompose up to (None means all layers)
45
               mlp_input: Whether to include attention output for current layer
               apply_ln: Whether to apply LayerNorm scaling
               pos_slice: Which positions to keep
               mode: Which components to include ('all', 'mlp', or 'attn')
50
           Returns:
51
               Tensor of shape [components, batch, pos, d_model] containing the decomposed
52
               residual stream
53
54
           components = []
55
56
           # Add embeddings if needed
           if mode in ['all', 'attn']:
57
               components.append(self['embed'])
58
               if 'pos_embed' in self.cache_dict:
59
```

Resources X

You are subscribed to Colab Pro. Learn more Available: 48 compute units Usage rate: approximately 1.66 per hour You have 1 active session.

Manage sessions

Python 3 Google Compute Engine backend (GPU) Showing resources from 15:06 to 15:56





Disk 39.3 / 235.7 GB

```
61
 62
            # Process each layer
            n_layers = layer if layer is not None else self.model.config.num_hidden_layers
 63
 64
            for l in range(n_layers):
 65
                # Add attention components
 66
                if mode in ['all'. 'attn']:
                    if f'blocks.{l}.attn.hook_result' in self.cache_dict:
 67
 68
                        components.append(self[f'blocks.{l}.attn.hook_result'])
 69
 70
                # Add MLP components
 71
                if mode in ['all', 'mlp']:
                    if f'blocks.{l}.mlp.hook_post' in self.cache_dict:
 72
 73
                        components.append(self[f'blocks.{l}.mlp.hook_post'])
 74
 75
            # Stack components
 76
            components = torch.stack(components)
 77
 78
            # Apply position slicing if needed
 79
            if pos_slice is not None:
 80
                components = components[..., pos_slice, :]
            # Apply layer norm if requested
 83
            if apply_ln:
                components = self.apply_ln_to_stack(components, layer=layer)
 85
 86
            return components
 87
 88
        def apply_ln_to_stack(self, residual_stack: torch.Tensor, layer: Optional[int] = None,
 89
                             mlp_input: bool = False) -> torch.Tensor:
            """Apply appropriate LayerNorm scaling to a stack of residual stream components."""
 90
 91
            # Get LayerNorm weights/bias for target layer
 92
            if laver is None:
 93
                layer = self.model.config.num_hidden_layers
 94
 95
            ln_weight = self.model.transformer.layers[layer].ln1.weight
 96
            ln_bias = self.model.transformer.layers[layer].ln1.bias
 97
 98
            # Calculate LaverNorm statistics
            mean = residual_stack.mean(dim=-1, keepdim=True)
 99
100
            var = ((residual_stack - mean) ** 2).mean(dim=-1, keepdim=True)
101
            normed = (residual_stack - mean) / (var + 1e-5).sqrt()
102
103
            # Apply weights and bias
104
            return \ ln\_weight * normed + ln\_bias
105
106
        def stack_activation(self, activation_name: str, layer: int = -1) -> torch.Tensor:
107
            """Stack activations with a given name from all layers up to layer.""
108
            components = []
109
            n_layers = layer if layer >= 0 else self.model.config.num_hidden_layers
110
111
            for l in range(n_layers):
                key = f"blocks.{l}.{activation_name}"
112
113
                if key in self.cache_dict:
                    components.append(self.cache_dict[key])
114
115
            return torch.stack(components)
117
118 def create_caching_hooks(model, hook_points):
119
        """Create forward hooks to cache activations at specified points in a Qwen model."""
120
        cache = \{\}
        hooks = []
121
122
123
        def get activation(name):
            def hook(module, input, output):
124
125
                cache[name] = output
126
            return hook
127
128
        # Parse hook points and attach hooks
129
        for hook point in hook points:
            # Split into components (e.g. "model.layers.0.input_layernorm")
130
            components = hook_point.split('.')
131
132
133
            # Navigate model architecture to find module
134
            current module = model
135
            for comp in components:
136
                if comp == "model":
137
                    current_module = current_module.model # Access the inner model
138
                elif comp.isdigit():
139
                    current_module = current_module[int(comp)]
140
                else:
141
                    current_module = getattr(current_module, comp)
142
143
            # Register hook
144
            hook = current_module.register_forward_hook(get_activation(hook_point))
145
            hooks.append(hook)
146
147
        return hooks, ActivationCache(cache, model)
148
149 def run_with_cache(model, inputs, names_filter=None):
150
        """Run model with activation caching, properly handling quantized models."""
        # Create list of hook points if not specified
151
152
        if names_filter is None:
153
           names_filter = [
154
                 "model.embed_tokens",
155
                f"model.layers.{model.config.num_hidden_layers-1}.input_layernorm"
156
157
        # Create and attach hooks
158
```

components.append(self['pos_embed'])

```
159
       hooks, cache = create_caching_hooks(model, names_filter)
160
161
            # Ensure inputs are properly formatted
162
163
            if not isinstance(inputs, dict):
164
                inputs = {
165
                    "input_ids": inputs,
                    "attention_mask": torch.ones_like(inputs),
166
167
                    "return_dict": True
168
                }
169
170
            # Run model without autocast
171
            with torch.no_grad():
172
                outputs = model(**inputs)
173
174
            # Create Output object
175
            if hasattr(outputs, 'loss'):
176
                return Output(outputs.logits, outputs.loss), cache
177
            else:
178
                return Output(outputs.logits), cache
179
180
        finally:
181
            # Always remove hooks
            for hook in hooks:
182
                hook.remove()
  1 from nnsight import LanguageModel
  2 import torch
  3 from transformers import {\tt BitsAndBytesConfig}
  5 \# First, let's set up the 4-bit quantization configuration
  6 # This matches the settings used by Unsloth for their quantized models
  7 quantization_config = BitsAndBytesConfig(
        load_in_4bit=True,
        bnb_4bit_quant_type="nf4", # Normal Float 4 format
 10
        bnb_4bit_compute_dtype=torch.float16,
 11
        bnb_4bit_use_double_quant=True # Nested quantization for additional memory savings
 12 )
 13
 14 # Disable gradients since we're doing inference only
 15 torch.set grad enabled(False)
 16
 17 # Load the base instruction model (4-bit quantized version)
 18 base model = LanguageModel(
        "unsloth/Qwen2.5-14B-Instruct-bnb-4bit",
 19
        device_map="cuda:0",
 20
 21
        trust_remote_code=True, # Required for Qwen models
        {\tt quantization\_config=quantization\_config}
 22
 23 )
 24
 25 \# Load the COT model (4-bit quantized version)
 26 cot_model = LanguageModel(
        "unsloth/DeepSeek-R1-Distill-Qwen-14B-unsloth-bnb-4bit",\\
 27
 28
 29
        device_map="cuda:1",
 30
        trust_remote_code=True,
 31
        \verb| quantization_config= quantization_config| \\
 32 )
 1 # %%
  2 import os
  3 from IPython import get_ipython
  5 ipython = get_ipython()
  6 # Code to automatically update the HookedTransformer code as its edited without restarting the
  7 if ipython is not None:
        ipython.magic("load_ext autoreload")
  9
        ipython.magic("autoreload 2")
 10
 11 import plotly.io as pio
 12 pio.renderers.default = "jupyterlab"
 13
 14 # Import stuff
 15 import einops
 16 import json
 17 import argparse
 19 from datasets import load_dataset
 20 from pathlib import Path
 21 import plotly.express as px
 22 from torch.distributions.categorical import Categorical
 23 from tqdm import tqdm
 24 import torch
 25 import numpy as np
 26 #from transformer_lens import HookedTransformer
 27 from jaxtyping import Float
 28 #from transformer_lens.hook_points import HookPoint
 29
 30 from functools import partial
 31
 32 from IPython.display import HTML
 33
 34 from transformer_lens.utils import to_numpy
 35 import pandas as pd
 36
 37 from html import escape
 38 import colorsys
```

```
39
 40
 41 import wandb
 42
 43 import plotly.graph_objects as go
 44
 45 update_layout_set = {
        "xaxis_range", "yaxis_range", "hovermode", "xaxis_title", "yaxis_title", "colorbar", "colortitle_x", "bargap", "bargroupgap", "xaxis_tickformat", "yaxis_tickformat", "title_y", "i" "xaxis_gridwidth", "xaxis_gridcolor", "yaxis_showgrid", "yaxis_gridwidth"
 46
 47
 48
 49 }
 50
 51 def imshow(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
 52
        if isinstance(tensor, list):
 53
            tensor = torch.stack(tensor)
 54
        kwargs_post = {k: v for k, v in kwargs.items() if k in update_layout_set}
 55
        kwargs_pre = {k: v for k, v in kwargs.items() if k not in update_layout_set}
 56
        if "facet_labels" in kwargs_pre:
 57
            facet_labels = kwargs_pre.pop("facet_labels")
             facet_labels = None
        if "color_continuous_scale" not in kwargs_pre:
            kwargs_pre["color_continuous_scale"] = "RdBu"
         fig = px.imshow(to_numpy(tensor), color_continuous_midpoint=0.0,labels={"x":xaxis, "y":ya>
 62
        if facet_labels:
 64
             for i, label in enumerate(facet_labels):
 65
                 fig.layout.annotations[i]['text'] = label
 66
 67
        fig.show(renderer)
 68
 69 def line(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
        px.line(y=to_numpy(tensor), labels={"x":xaxis, "y":yaxis}, **kwargs).show(renderer)
 70
 71
 72 def scatter(x, y, xaxis="", yaxis="", caxis="", renderer=None, return_fig=False, **kwargs):
 73
        x = to numpv(x)
 74
        y = to_numpy(y)
 75
        fig = px.scatter(y=y, x=x, labels={"x":xaxis, "y":yaxis, "color":caxis}, **kwargs)
 76
        if return_fig:
            return fia
 77
 78
        fig.show(renderer)
 79
 80 def lines(lines_list, x=None, mode='lines', labels=None, xaxis='', yaxis='', title = '', log_Y
 81
        # Helper function to plot multiple lines
 82
        if type(lines_list)==torch.Tensor:
 83
             lines_list = [lines_list[i] for i in range(lines_list.shape[0])]
 84
        if x is None:
 85
            x=np.arange(len(lines_list[0]))
 86
        fig = go.Figure(layout={'title':title})
 87
        fig.update_xaxes(title=xaxis)
 88
        fig.update_yaxes(title=yaxis)
        for c, line in enumerate(lines_list):
 89
 90
            if type(line)==torch.Tensor:
                 line = to_numpy(line)
 91
 92
            if labels is not None:
 93
                 label = labels[c]
 94
             else:
             fig.add_trace(go.Scatter(x=x, y=line, mode=mode, name=label, hovertext=hover, **kwarg:
        if log_y:
            fig.update_layout(yaxis_type="log")
 99
        fig.show()
100
101 def bar(tensor, renderer=None, xaxis="", yaxis="", **kwargs):
102
        px.bar(
103
            y=to numpy(tensor),
             labels={"x": xaxis, "y": yaxis},
104
            template="simple_white"
105
106
            **kwargs).show(renderer)
107
108 \ def \ create\_html(strings, \ values, \ saturation=0.5, \ allow\_different\_length=False):
        # escape strings to deal with tabs, newlines. etc.
109
110
        escaped_strings = [escape(s, quote=True) for s in strings]
        processed_strings = [
    s.replace("\n", "<br/>").replace("\t", "&emsp;").replace(" ", "&nbsp;")
111
112
113
             for s in escaped\_strings
114
115
116
        if isinstance(values, torch.Tensor) and len(values.shape)>1:
117
             values = values.flatten().tolist()
118
119
        if not allow_different_length:
120
            assert len(processed_strings) == len(values)
121
122
        # scale values
123
        max_value = max(max(values), -min(values))+1e-3
124
        scaled_values = [v / max_value * saturation for v in values]
125
126
        # create html
127
128
        for i, s in enumerate(processed_strings):
129
            if i<len(scaled_values):</pre>
130
                 v = scaled_values[i]
131
            else:
132
                v = 0
133
            if v < 0:
134
                hue = 0 # hue for red in HSV
135
            else:
               hue = 0.66 # hue for blue in HSV
136
137
             rgb_color = colorsys.hsv_to_rgb(
```

```
138
                hue, v, 1
            ) \# hsv color with hue 0.66 (blue), saturation as v, value 1
139
            hex_color = "#%02x%02x%02x" % (
140
                int(rgb\_color[0] * 255),
141
142
                int(rgb_color[1] * 255),
143
                int(rgb\_color[2] * 255),
144
145
            html += f'<span style="background-color: {hex_color}; border: 1px solid lightgray; for
146
147
        display(HTML(html))
148
149 # crosscoder stuff
150
151 def arg_parse_update_cfg(default_cfg):
152
153
        Helper function to take in a dictionary of arguments, convert these to command line argume
154
155
        If in Ipython, just returns with no changes
156
157
        if get_ipython() is not None:
158
            # Is in IPython
            print("In IPython - skipped argparse")
159
160
            return default_cfg
161
        cfg = dict(default_cfg)
162
        parser = argparse.ArgumentParser()
163
        for key, value in default_cfg.items():
164
            if type(value) == bool:
165
                # argparse for Booleans is broken rip. Now you put in a flag to change the default
166
                if value:
167
                    parser.add_argument(f"--{key}", action="store_false")
                else:
168
                    parser.add_argument(f"--{key}", action="store_true")
169
170
171
            else:
                parser.add_argument(f"--{key}", type=type(value), default=value)
172
173
        args = parser.parse_args()
174
        parsed_args = vars(args)
175
        cfg.update(parsed_args)
176
        print("Updated config")
177
        print(json.dumps(cfg, indent=2))
178
        return cfg
179
180 def load_pile_lmsys_mixed_tokens():
181
            print("Loading data from disk")
182
183
            all_tokens = torch.load("/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-2.pt")
184
185
            print("Data is not cached. Loading data from HF")
186
            data = load_dataset(
187
                "ckkissane/pile-lmsys-mix-1m-tokenized-gemma-2",
188
                split="train",
189
                cache_dir="/workspace/cache/"
190
191
            data.save_to_disk("/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-2.hf")
192
            data.set_format(type="torch", columns=["input_ids"])
193
            all_tokens = data["input_ids"]
            torch.save(all_tokens, "/workspace/data/pile-lmsys-mix-1m-tokenized-gemma-2.pt")
194
195
            print(f"Saved tokens to disk")
        return all_tokens
  1 class Buffer:
  3
        Buffer implementation that properly handles activation collection from nnsight models.
  5
             _init__(self, cfg, model_A, model_B, all_tokens):
            self.cfg = cfg
            self.model_A = model_A
            self.model_B = model_B
  9
            self.all_tokens = all_tokens
 10
 11
            # Buffer setup
 12
            self.buffer_size = cfg["batch_size"] * cfg["buffer_mult"]
            self.buffer_batches = self.buffer_size // (cfg["seq_len"] - 1)
 13
            self.buffer_size = self.buffer_batches * (cfg["seq_len"] - 1)
 14
 15
            # Initialize buffer with appropriate dtype
 16
 17
            self.buffer = torch.zeros(
                (self.buffer_size, 2, model_A.config.hidden_size),
 18
                dtype=torch.bfloat16,
 19
                device=cfg["device"],
 20
                requires_grad=False
 21
 22
 23
 24
            self.token_pointer = 0
 25
            self.first = True
 26
            self.normalize = True
 27
 28
            # Pre-compute normalization factors
            \verb"print("Calculating normalization factors...")"
 29
 30
            norm_A = self._estimate_norm_scaling_factor(model_A, "Model A")
 31
            norm_B = self._estimate_norm_scaling_factor(model_B, "Model B")
 32
 33
            self.normalisation_factor = torch.tensor(
 34
                [norm_A, norm_B],
 35
                device=cfg["device"],
 36
                dtype=torch.float32
 37
 38
```

```
39
            self.refresh()
 40
 41
        def collect activations(self, model, tokens):
              ""Collect activations from a model using nnsight's context manager."""
 42
 43
            with model.trace() as trace:
 44
                # Get the target layer
 45
                laver = model
                for comp in self.cfg["hook_point"].split('.'):
 46
 47
                    if comp == "model":
 48
                        continue
                    elif comp == "layers":
 49
 50
                        continue
 51
                    elif comp.isdigit():
 52
                        layer = layer.layers[int(comp)]
 53
 54
                        layer = getattr(layer, comp)
 55
 56
                # Run forward pass
 57
                model(input_ids=tokens, attention_mask=torch.ones_like(tokens))
                # Get activations - now they're materialized
                return layer.output[0].to(dtype=torch.bfloat16)
 61
        @torch.no_grad()
 62
        def _estimate_norm_scaling_factor(self, model, desc, n_batches=100):
            """Calculate norm scaling factor with proper tensor handling."""
 64
 65
            print(f"\nEstimating norm scaling factor for {desc}")
 66
 67
            batch_size = min(8, self.cfg["model_batch_size"])
 68
 69
            try:
 70
                model.eval()
 71
 72
                for i in range(n batches):
 73
                    # Get a batch of tokens
                    tokens = self.all\_tokens[i*batch\_size: (i+1)*batch\_size].to(model.devices) \\
 74
 75
 76
                        # Collect activations using nnsight's context manager
 77
 78
                        acts = self._collect_activations(model, tokens)
 79
                        # Calculate norm - activations are now materialized
 80
                        norm = acts.norm(dim=-1).mean().item()
 81
 82
                        norms.append(norm)
 83
 84
                    except Exception as e:
 85
                        print(f"Batch {i} failed: {str(e)}")
 86
                        continue
 87
 88
                    if i % 10 == 0:
                        torch.cuda.empty_cache()
 89
 90
 91
            except Exception as e:
 92
                print(f"Error during norm estimation: {str(e)}")
 93
 94
                print(f"Warning: Using default scaling factor for {desc}")
 95
                return np.sqrt(model.config.hidden_size)
            mean_norm = np.mean(norms)
 99
            return np.sqrt(model.config.hidden_size) / mean_norm
100
101
        @torch.no_grad()
102
        def refresh(self):
            """Refresh buffer with new activations."""
103
            self.pointer = 0
104
105
            print("Refreshing the buffer!")
106
            num_batches = self.buffer_batches if self.first else self.buffer_batches // 2
107
108
            self.first = False
109
            batch_size = min(8, self.cfg["model_batch_size"])
110
111
112
            for i in range(0, num_batches, batch_size):
113
                tokens = self.all tokens[
                    self.token_pointer : min(
114
115
                        self.token_pointer + batch_size,
116
                        num_batches
117
                ].to(self.model_A.device)
118
119
120
121
                    # Collect activations from both models
122
                    acts_A = self._collect_activations(self.model_A, tokens)
123
                    acts_B = self._collect_activations(self.model_B, tokens)
124
125
                    # Stack and process activations
                    acts = torch.stack([acts_A, acts_B], dim=0)
126
127
                    acts = acts[:, :, 1:, :]  # Drop BOS token
128
129
                    # Reshape for buffer storage
130
                    acts = einops.rearrange(
131
132
                        "n_layers batch seq_len d_model -> (batch seq_len) n_layers d_model"
133
134
135
                    self.buffer[self.pointer : self.pointer + acts.shape[0]] = acts
                    self.pointer += acts.shape[0]
136
                    self.token_pointer += batch_size
137
```

```
139
                except Exception as e:
                    print(f"Error in batch {i}: {str(e)}")
140
141
                    continue
142
                if i % 10 == 0:
143
144
                    torch.cuda.empty_cache()
145
            # Shuffle buffer
146
147
            self.pointer = 0
148
            self.buffer = self.buffer[torch.randperm(self.buffer.shape[0]).to(self.cfg["device"])]
149
150
        @torch.no_grad()
        def next(self):
    """Get next batch of activations."""
151
152
153
            out = self.buffer[self.pointer : self.pointer + self.cfg["batch_size"]]
154
            self.pointer += self.cfg["batch_size"]
155
156
            if self.pointer > self.buffer.shape[0] // 2 - self.cfg["batch_size"]:
                self.refresh()
157
158
            if self.normalize:
159
160
                out = out * self.normalisation_factor[None, :, None]
161
  1 import wandb
  2 from google.colab import userdata
  3 import os
  4 os.environ["WANDB_API_KEY"] = userdata.get('WANDB_API_KEY')
  5 wandb.login()
  6 from huggingface_hub.hf_api import HfFolder
  8 HfFolder.save_token(userdata.get('HF_TOKEN'))
 wandb: Currently logged in as: jacktpayne51 (jacktpayne51-macquarie-university)
  1 from tqdm.auto import tqdm # Properly import tqdm
  2 import torch
  3 from datasets import load_dataset
  4 from pathlib import Path
  5 import os
  7 def load_or_create_qwen_tokens():
  8
        Creates and caches a dataset of tokens specifically for Qwen models using the
        uncopyrighted version of the Pile and LMSys chat data.
 10
 11
        # First, ensure the cache directory exists
 12
 13
        cache dir = Path("/workspace/data")
 14
        cache dir.mkdir(parents=True, exist ok=True)
        cache_path = cache_dir / "pile-lmsys-mix-1m-tokenized-qwen.pt"
 15
 16
 17
            \verb|print("Looking for previously cached Qwen-tokenized data...")|\\
 18
 19
            # Add weights_only=True for security
 20
            all_tokens = torch.load(cache_path, weights_only=True)
 21
            print(f"Found cached data! Loading {len(all_tokens)} tokenized sequences...")
 22
 23
        {\tt except (FileNotFoundError, RuntimeError):}
 24
            print("No cached data found. Creating new dataset...")
 25
 26
            print("Loading uncopyrighted Pile data...")
 27
            pile = load_dataset(
 28
                "monology/pile-uncopyrighted",
 29
                split="train",
 30
                streaming=True
 31
 32
            print("Loading LMSys chat data...")
 33
 34
            lmsys = load_dataset(
 35
                "lmsys/lmsys-chat-1m",
                split="train",
                streaming=True
 37
 39
 40
            target_samples = 10_000 # 500k from each source for balance
 42
 43
            print(f"Collecting balanced dataset of {target_samples * 2} samples...")
 44
            pile_count = 0
 45
            lmsys_count = 0
 46
 47
            # Collect Pile samples with progress tracking
 48
            pile iter = iter(pile)
            print("Collecting Pile samples...")
 49
 50
            while pile_count < target_samples:</pre>
 51
                try:
 52
                    sample = next(pile iter)
 53
                    texts.append(sample['text'])
 54
                    pile count += 1
 55
                    if pile_count % 1000 == 0:
                        \verb|print(f"Collected {pile\_count}|/{target\_samples}| Pile samples")|\\
 56
 57
                except StopIteration:
 58
                    break
 59
```

```
# Collect LMSys samples with progress tracking
            lmsys iter = iter(lmsys)
 61
 62
            print("\nCollecting LMSys samples...")
 63
            while lmsys_count < target_samples:
 64
                try:
                     sample = next(lmsys_iter)
 65
 66
                     chat_text = format_chat_for_qwen(sample["conversation"])
 67
                     texts.append(chat_text)
 68
                     lmsys_count += 1
 69
                     if lmsys_count % 1000 == 0:
 70
                         print(f"Collected \{lmsys\_count\}/\{target\_samples\} \ LMSys \ samples")
 71
                except StopIteration:
 72
                     break
 73
 74
            print(f"\nTokenizing \{len(texts)\}\ sequences with Qwen tokenizer...")
 75
            # Process in smaller batches to manage memory
 76
            batch_size = 1000
 77
            all_tokens = []
 78
            total_batches = len(texts) // batch_size + (1 if len(texts) % batch_size else 0)
 79
            for i in range(0, len(texts), batch_size):
                batch = texts[i:i + batch_size]
                print(f"Processing batch {i//batch_size + 1}/{total_batches}")
                tokenized = base_model.tokenizer(
                     padding=True,
                     truncation=True,
 86
                     max_length=1024,
                     return_tensors="pt"
 88
 89
                all_tokens.append(tokenized.input_ids)
 90
 91
 92
            # Concatenate all batches
            all_tokens = torch.cat(all_tokens, dim=0)
 93
 94
 95
            print("Caching tokenized data...")
 96
            torch.save(all_tokens, cache_path)
 97
        print(f"Successfully prepared {len(all_tokens)} sequences")
 98
 99
        return all_tokens
100
101 def format_chat_for_qwen(conversation):
102
103
        Formats LMSys chat data to match Qwen's expected chat format.
104
        This ensures the model will process chat data similarly to how it was trained.
105
106
107
            conversation: A list of message dictionaries with 'role' and 'content' keys
108
109
        str: A formatted chat string in Qwen's expected format
110
111
112
        formatted_chat = ""
113
        for message in conversation:
114
            role = message["role"]
115
            content = message["content"]
117
            # Format based on Qwen's chat template
            if role == "user":
118
119
                formatted\_chat += f'' < |im\_start| > user \setminus \{content\} < |im\_end| > \n''
            elif role == "assistant":
120
                formatted_chat += f"<|im_start|>assistant\n{content}<|im_end|>\n"
121
122
123
        return formatted chat
  1 cfg = {
            # Core training parameters
            "batch_size": 1024, # Reduced from 4096
"buffer_mult": 32, # Reduced from 128
            "model_batch_size": 8, # Significantly reduced
  6
            "lr": 5e-5,
            "l1_coeff": 2,
            "beta1": 0.9,
  8
            "beta2": 0.999,
 10
            # Model architecture settings
 11
            "dict size": 16384,
 12
            "hook_point": "model.layers.24.input_layernorm",
 13
            "d_in": base_model.config.hidden_size,
 14
 15
 16
            # Hardware settings
            "device": "cuda:0",
 17
            "seq_len": 1024,
 18
 19
            "enc_dtype": "fp16",
 20
 21
            # Training schedule
            "num_tokens": 400_000_000,
"save_every": 30000,
 22
 23
            "log_every": 100,
 24
 25
 26
            # Additional settings
            "seed": 42,
 27
            "wandb_project": "crosscoder-cot-analysis",
 28
 29
            "wandb_entity": "jacktpayne51",
            "dec_init_norm": 0.08
 30
 31
 32 def setup_training(base_model, cot_model):
```

```
Prepares the crosscoder training environment with Qwen-specific settings.
35
       print("Beginning training setup...")
37
       # Load our Qwen-specific tokenized data
39
       all_tokens = load_or_create_qwen_tokens()
40
41
       cfg = {
           # Core training parameters
43
           "batch_size": 1024,  # Reduced from 4096
           "buffer_mult": 32,
44
                                # Reduced from 128
           "model_batch_size": 8, # Significantly reduced
45
46
           "lr": 5e-5,
47
           "l1_coeff": 2,
           "beta1": 0.9,
48
49
           "beta2": 0.999,
50
51
           # Model architecture settings
           "dict_size": 16384,
52
           "hook_point": "model.layers.24.input_layernorm",
53
           "d_in": base_model.config.hidden_size,
54
55
56
           # Hardware settings
           "device": "cuda:0",
"seq_len": 1024,
57
58
           "enc_dtype": "fp16",
59
60
61
           # Training schedule
62
           "num_tokens": 400_000_000,
           "save_every": 30000,
63
64
           "log_every": 100,
65
66
           # Additional settings
67
           "seed": 42,
68
           "wandb_project": "crosscoder-cot-analysis",
69
           "wandb_entity": "jacktpayne51",
70
           "dec_init_norm": 0.08
71
72
73
       \verb"print("Creating activation buffer...")"
74
75
           cfg=cfg,
76
           model_A=base_model,
77
           model_B=cot_model,
78
           all_tokens=all_tokens
79
80
      print("Training setup complete!")
81
       return cfg, buffer
2 from torch import nn
3 import pprint
 4 import torch.nn.functional as F
 5 from typing import Optional, Union
 6 from huggingface_hub import hf_hub_download
8 from typing import NamedTuple
10 DTYPES = {"fp32": torch.float32, "fp16": torch.float16, "bf16": torch.bfloat16}
11 SAVE_DIR = Path("/workspace/crosscoder-model-diff-replication/checkpoints")
12
13 class LossOutput(NamedTuple):
14
       # loss: torch.Tensor
15
       12 loss: torch.Tensor
16
       l1_loss: torch.Tensor
       10 loss: torch.Tensor
17
       explained_variance: torch.Tensor
18
19
      explained variance A: torch. Tensor
20
      {\tt explained\_variance\_B: torch.Tensor}
21
22 class CrossCoder(nn.Module):
23
       def __init__(self, cfg):
24
           super().__init__()
25
           self.cfg = cfg
26
           d_hidden = self.cfg["dict_size"]
27
           d_in = self.cfg["d_in"]
28
           self.dtype = DTYPES[self.cfg["enc_dtype"]]
29
           torch.manual_seed(self.cfg["seed"])
30
           \# hardcoding n_models to 2
31
           self.W_enc = nn.Parameter(
32
               torch.empty(2, d_in, d_hidden, dtype=self.dtype)
33
34
           self.W_dec = nn.Parameter(
35
               torch.nn.init.normal_(
                   torch.empty(
37
                       d_hidden, 2, d_in, dtype=self.dtype
41
           self.W_dec = nn.Parameter(
               torch.nn.init.normal_(
43
                   torch.empty(
44
                       d_hidden, 2, d_in, dtype=self.dtype
45
46
48
           # Make norm of W_dec 0.1 for each column, separate per layer
```

```
49
                           self.W_dec.data = (
                                   self.W_dec.data / self.W_dec.data.norm(dim=-1, keepdim=True) * self.cfg["dec_init_
  50
  51
  52
                          # Initialise W_enc to be the transpose of W_dec
  53
                           self.W_enc.data = einops.rearrange(
  54
                                    self.W dec.data.clone(),
  55
                                    "d_hidden n_models d_model \rightarrow n_models d_model d_hidden",
  56
  57
                           {\tt self.b\_enc = nn.Parameter(torch.zeros(d\_hidden, dtype=self.dtype))}
  58
                          self.b_dec = nn.Parameter(
  59
                                   torch.zeros((2, d_in), dtype=self.dtype)
  60
  61
                          self.d_hidden = d_hidden
  62
  63
                           self.to(self.cfg["device"])
  64
                           self.save_dir = None
  65
                           self.save_version = 0
  66
  67
                  def encode(self, x, apply_relu=True):
                           # x: [batch, n_models, d_model]
                          x_{enc} = einops.einsum(
  70
  71
                                   self.W_enc,
  72
                                    "batch n_models d_model, n_models d_model d_hidden -> batch d_hidden",
  73
  74
                          if apply relu:
  75
                                   acts = F.relu(x_enc + self.b_enc)
  76
                           else:
  77
                                   acts = x_enc + self.b_enc
  78
                          return acts
  79
  80
                 def decode(self, acts):
  81
                          # acts: [batch, d_hidden]
  82
                          acts dec = einops.einsum(
  83
                                   acts.
                                   self.W dec.
  84
  85
                                    "batch d_hidden, d_hidden n_models d_model -> batch n_models d_model",
  86
  87
                          return acts_dec + self.b_dec
  88
  89
                 def forward(self, x):
  90
                          # x: [batch, n_models, d_model]
  91
                          acts = self.encode(x)
  92
                          return self.decode(acts)
  93
  94
                  def get_losses(self, x):
  95
                          # x: [batch, n_models, d_model]
  96
                          x = x.to(self.dtype)
  97
                          acts = self.encode(x)
  98
                          # acts: [batch, d_hidden]
                           x_reconstruct = self.decode(acts)
  99
100
                          diff = x_reconstruct.float() - x.float()
                           squared_diff = diff.pow(2)
101
102
                           l2_per_batch = einops.reduce(squared_diff, 'batch n_models d_model -> batch', 'sum')
103
                           l2_loss = l2_per_batch.mean()
104
                           total_variance = einops.reduce((x - x.mean(0)).pow(2), 'batch n_models d_model -> batc
105
106
                          explained_variance = 1 - l2_per_batch / total_variance
107
108
                          per\_token\_l2\_loss\_A = (x\_reconstruct[:, 0, :] - x[:, 0, :]).pow(2).sum(dim=-1).squeezer(in the construct[:, 0, :]).squeezer(in t
                           total_variance_A = (x[:, 0, :] - x[:, 0, :].mean(0)).pow(2).sum(-1).squeeze()
109
110
                          explained_variance_A = 1 - per_token_l2_loss_A / total_variance_A
111
                           per\_token\_l2\_loss\_B = (x\_reconstruct[:, 1, :] - x[:, 1, :]).pow(2).sum(dim=-1).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :]).pow(2).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :]).pow(2).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :]).pow(2).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :]).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :] - x[:, 1, :]).squeez{(instruct)} = (x_reconstruct[:, 1, :] - x[:, 1, :] - x[
112
                           total_variance_B = (x[:, 1, :] - x[:, 1, :].mean(0)).pow(2).sum(-1).squeeze()
113
                           explained_variance_B = 1 - per_token_l2_loss_B / total_variance_B
114
115
                          decoder norms = self.W dec.norm(dim=-1)
116
117
                          # decoder_norms: [d_hidden, n_models]
                           total_decoder_norm = einops.reduce(decoder_norms, 'd_hidden n_models -> d_hidden', 'su
118
                           l1 loss = (acts * total decoder norm[None, :]).sum(-1).mean(0)
119
120
121
                           l0 loss = (acts>0).float().sum(-1).mean()
122
                           return LossOutput(l2_loss=l2_loss, l1_loss=l1_loss, l0_loss=l0_loss, explained_variand
123
124
125
                  def create_save_dir(self):
126
                          base_dir = Path("/workspace/crosscoder-model-diff-replication/checkpoints")
127
                           version list = [
                                   \verb|int(file.name.split("\_")[1]|)|\\
128
129
                                    for file in list(SAVE_DIR.iterdir())
130
                                   if "version" in str(file)
131
132
                          if len(version_list):
133
                                   version = 1 + max(version_list)
134
135
                           self.save_dir = base_dir / f"version_{version}"
136
137
                           self.save_dir.mkdir(parents=True)
138
139
                 def save(self):
                          if self.save_dir is None:
140
141
                                   self.create_save_dir()
142
                           weight_path = self.save_dir / f"{self.save_version}.pt"
143
                          cfg_path = self.save_dir / f"{self.save_version}_cfg.json"
144
145
                           torch.save(self.state_dict(), weight_path)
                          with open(cfg_path, "w") as f:
146
                                   json.dump(self.cfg, f)
147
```

```
148
            print(f"Saved as version {self.save_version} in {self.save_dir}")
149
150
            self.save_version += 1
151
152
        @classmethod
153
        def load_from_hf(
154
            cls.
            repo_id: str = "ckkissane/crosscoder-gemma-2-2b-model-diff",
path: str = "blocks.14.hook_resid_pre",
155
156
157
            device: Optional[Union[str, torch.device]] = None
158
        ) -> "CrossCoder":
159
160
            Load CrossCoder weights and config from HuggingFace.
161
162
163
                repo_id: HuggingFace repository ID
164
                path: Path within the repo to the weights/config
165
                model: The transformer model instance needed for initialization
166
                device: Device to load the model to (defaults to cfg device if not specified)
167
168
            Initialized CrossCoder instance
169
170
171
172
            # Download config and weights
173
            config_path = hf_hub_download(
174
                repo_id=repo_id,
175
                filename=f"{path}/cfg.json"
176
            weights_path = hf_hub_download(
177
178
                repo_id=repo_id,
                filename=f"{path}/cc_weights.pt"
179
180
181
182
            # Load config
            with open(config_path, 'r') as f:
183
184
                cfg = json.load(f)
185
186
            # Override device if specified
187
            if device is not None:
188
                cfg["device"] = str(device)
189
190
            # Initialize CrossCoder with config
191
            instance = cls(cfg)
192
193
            # Load weights
194
            state_dict = torch.load(weights_path, map_location=cfg["device"])
195
            instance.load_state_dict(state_dict)
196
197
            return instance
198
199
        @classmethod
200
        def load(cls, version_dir, checkpoint_version):
201
            save_dir = Path("/workspace/crosscoder-model-diff-replication/checkpoints") / str(vers)
202
            cfg_path = save_dir / f"{str(checkpoint_version)}_cfg.json"
203
            weight_path = save_dir / f"{str(checkpoint_version)}.pt
204
205
            cfg = json.load(open(cfg_path, "r"))
            pprint.pprint(cfg)
207
            self = cls(cfg=cfg)
208
            self.load_state_dict(torch.load(weight_path))
            return self
  2 import tqdm
  4 from torch.nn.utils import clip_grad_norm_
  5 class Trainer:
        def __init__(self, cfg, model_A, model_B, all_tokens):
            self.cfg = cfg
            self.model_A = model_A
            self.model_B = model_B
 10
            self.crosscoder = CrossCoder(cfg)
            self.buffer = Buffer(cfg, model_A, model_B, all_tokens)
 11
 12
            self.total_steps = cfg["num_tokens"] // cfg["batch_size"]
 13
            self.optimizer = torch.optim.Adam(
 14
                self.crosscoder.parameters(),
 15
                lr=cfq["lr"],
 16
                betas=(cfg["beta1"], cfg["beta2"]),
 17
 18
            self.scheduler = torch.optim.lr scheduler.LambdaLR(
 19
                self.optimizer, self.lr_lambda
 20
 21
 22
            self.step_counter = 0
 23
            wandb.init(project=cfg["wandb_project"], entity=cfg["wandb_entity"])
 24
 25
        def lr_lambda(self, step):
 26
 27
            if step < 0.8 * self.total_steps:</pre>
 28
                return 1.0
 29
            else:
 30
                return 1.0 - (step - 0.8 * self.total_steps) / (0.2 * self.total_steps)
 31
 32
        def get_l1_coeff(self):
 33
            # Linearly increases from 0 to cfg["l1_coeff"] over the first 0.05 * self.total_steps
            if self.step_counter < 0.05 * self.total_steps:</pre>
                return self.cfg["l1_coeff"] * self.step_counter / (0.05 * self.total_steps)
 35
```

```
36
            else:
                return self.cfg["l1_coeff"]
37
38
39
       def step(self):
            acts = self.buffer.next()
40
            losses = self.crosscoder.get_losses(acts)
41
            loss = losses.l2_loss + self.get_l1_coeff() * losses.l1_loss
42
43
            loss.backward()
            \verb|clip_grad_norm_(self.crosscoder.parameters(), max_norm=1.0)|\\
44
45
            self.optimizer.step()
46
            self.scheduler.step()
47
            self.optimizer.zero_grad()
48
49
            loss_dict = {
50
                "loss": loss.item(),
                "l2_loss": losses.l2_loss.item(),
51
52
                "l1_loss": losses.l1_loss.item(),
53
                "l0_loss": losses.l0_loss.item(),
54
                "l1_coeff": self.get_l1_coeff(),
                "lr": self.scheduler.get_last_lr()[0],
                "explained_variance": losses.explained_variance.mean().item(),
"explained_variance_A": losses.explained_variance_A.mean().item(),
                "explained_variance_B": losses.explained_variance_B.mean().item(),
59
            self.step_counter += 1
61
            return loss_dict
62
63
       def log(self, loss_dict):
64
            wandb.log(loss_dict, step=self.step_counter)
65
            print(loss_dict)
66
       def save(self):
67
68
            self.crosscoder.save()
69
70
       def train(self):
71
            self.step_counter = 0
72
73
                for i in tqdm.trange(self.total_steps):
                     loss_dict = self.step()
if i % self.cfg["log_every"] == 0:
74
75
76
                         self.log(loss_dict)
                     if (i + 1) % self.cfg["save_every"] == 0:
77
                         self.save()
78
            finally:
79
                self.save()
 1 def train_crosscoder(base_model, cot_model):
       # Set up our configuration and buffer
cfg, buffer = setup_training(base_model, cot_model)
 3
       # Initialize our crosscoder
 6
       crosscoder = CrossCoder(cfg)
       # Create trainer instance
 8
       trainer = Trainer(
 9
10
            cfg=cfg,
11
            model_A=base_model,
12
            model_B=cot_model,
13
            all_tokens=buffer
14
15
16
       # Start training
17
       \verb|print("Beginning crosscoder training...")| \\
       trainer.train()
19
20
        return crosscoder
21
22 # Main execution
23 if __name__ == "_
                      __main__":
       # Train the crosscoder
24
25
       crosscoder = train_crosscoder(base_model, cot_model)
26
27
       # Save the final model
28
       crosscoder.save()
29
30
```

Beginning training setup...

Looking for previously cached Qwen-tokenized data...

Found cached data! Loading 20000 tokenized sequences...

Successfully prepared 20000 sequences

Creating activation buffer...

Calculating normalization factors...

Estimating norm scaling factor for Model ${\bf A}$

Batch 0 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 0 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 1 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 2 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 3 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 4 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 5 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 6 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 7 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 8 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 9 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 10 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 11 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 12 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 13 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 14 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 15 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 16 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 17 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 18 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 19 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 20 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 21 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 22 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 23 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 24 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 25 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 26 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 27 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 28 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 29 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 30 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 31 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 32 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 33 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 34 failed: 'Qwen2ForCausallM' object has no attribute 'layers' Batch 35 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 36 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 37 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 38 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 39 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 40 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 41 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 42 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 43 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 44 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 45 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 46 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 47 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 48 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 49 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 50 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 51 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 52 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 53 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 54 failed: 'Owen2ForCausalLM' object has no attribute 'layers' Batch 55 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 56 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 57 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 58 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 59 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 60 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 61 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 62 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 63 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 64 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 65 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'
Batch 66 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 67 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 68 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 69 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 70 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 71 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 72 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 73 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 74 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 75 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 76 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 77 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 78 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 79 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 80 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 81 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 82 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 83 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' Batch 84 failed: 'Owen2ForCausalLM' object has no attribute 'layers'

Batch 85 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 86 failed: 'Qwen2ForCausalLM' object has no attribute 'layers'

Batch 87 failed: 'Qwen2ForCausalLM' object has no attribute 'layers' 'Owen? For Caucall M' object has no attribute Ratch 00 failed: