

✓ RNA_code_train_Data_

```
1 import pandas as pd
2 import os, gc
3 import numpy as np
4 from sklearn.model_selection import KFold
5
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9 from torch.utils.data import Dataset, DataLoader
10
11 import torch
12 from fastai.vision.all import *
13 def flatten(o):
14     "Concatenate all collections and items as a generator"
15     for item in o:
16         if isinstance(o, dict): yield o[item]; continue
17         elif isinstance(item, str): yield item; continue
18         try: yield from flatten(item)
19         except TypeError: yield item
20
21 from torch.cuda.amp import GradScaler, autocast
22 @delegates(GradScaler)
23 class MixedPrecision(Callback):
24     "Mixed precision training using Pytorch's `autocast` and `GradScaler`"
25     order = 10
26     def __init__(self, **kwargs): self.kwarg = kwargs
27     def before_fit(self):
28         self.autocast,self.learn.scaler,self.scales = autocast(),GradScaler(**self.kwarg),L()
29     def before_batch(self): self.autocast.__enter__()
30     def after_pred(self):
31         if next(flatten(self.pred)).dtype==torch.float16: self.learn.pred = to_float(self.pred)
32     def after_loss(self): self.autocast.__exit__(None, None, None)
33     def before_backward(self): self.learn.loss_grad = self.scaler.scale(self.loss_grad)
34     def before_step(self):
35         "Use `self` as a fake optimizer. `self.skipped` will be set to True `after_step` if gradients overflow"
36         self.skipped=True
37         self.scaler.step(self)
38         if self.skipped: raise CancelStepException()
39         self.scales.append(self.scaler.get_scale())
40     def after_step(self): self.learn.scaler.update()
41
42     @property
43     def param_groups(self):
44         "Pretend to be an optimizer for `GradScaler`"
45         return self.opt.param_groups
46     def step(self, *args, **kwargs):
47         "Fake optimizer step to detect whether this batch was skipped from `GradScaler`"
48         self.skipped=False
49     def after_fit(self): self.autocast,self.learn.scaler,self.scales = None,None,None
50
51 import fastai
52 fastai.callback.fp16.MixedPrecision = MixedPrecision
53
54 def seed_everything(seed):
55     random.seed(seed)
56     os.environ['PYTHONHASHSEED'] = str(seed)
57     np.random.seed(seed)
58     torch.manual_seed(seed)
59     torch.cuda.manual_seed(seed)
60     torch.backends.cudnn.deterministic = True
61     torch.backends.cudnn.benchmark = True
62
63 #fname = 'example0'
64 #PATH = '/kaggle/input/stanford-ribonanza-rna-folding-converted/'
65 OUT = './'
66 bs = 256
```

```

67 num_workers = 2
68 SEED = 2023
69 nfold = 4
70 device = 'cuda' if torch.cuda.is_available() else 'cpu'
71
72 class RNA_Dataset(Dataset):
73     def __init__(self, df, mode='train', seed=2023, fold=0, nfold=4,
74                 mask_only=False, **kwargs):
75         self.seq_map = {'A':0, 'C':1, 'G':2, 'U':3}
76         self.Lmax = 206
77         df['L'] = df.sequence.apply(len)
78         df_2A3 = df.loc[df.experiment_type=='2A3_MaP']
79         df_DMS = df.loc[df.experiment_type=='DMS_MaP']
80
81         split = list(KFold(n_splits=nfold, random_state=seed,
82                             shuffle=True).split(df_2A3))[fold][0 if mode=='train' else 1]
83         df_2A3 = df_2A3.iloc[split].reset_index(drop=True)
84         df_DMS = df_DMS.iloc[split].reset_index(drop=True)
85
86         m = (df_2A3['SN_filter'].values > 0) & (df_DMS['SN_filter'].values > 0)
87         df_2A3 = df_2A3.loc[m].reset_index(drop=True)
88         df_DMS = df_DMS.loc[m].reset_index(drop=True)
89
90         self.seq = df_2A3['sequence'].values
91         self.L = df_2A3['L'].values
92
93         self.react_2A3 = df_2A3[[c for c in df_2A3.columns if \
94                                 'reactivity_0' in c]].values
95         self.react_DMS = df_DMS[[c for c in df_DMS.columns if \
96                                 'reactivity_0' in c]].values
97         self.react_err_2A3 = df_2A3[[c for c in df_2A3.columns if \
98                                     'reactivity_error_0' in c]].values
99         self.react_err_DMS = df_DMS[[c for c in df_DMS.columns if \
100                                     'reactivity_error_0' in c]].values
101         self.sn_2A3 = df_2A3['signal_to_noise'].values
102         self.sn_DMS = df_DMS['signal_to_noise'].values
103         self.mask_only = mask_only
104
105     def __len__(self):
106         return len(self.seq)
107
108     def __getitem__(self, idx):
109         seq = self.seq[idx]
110         if self.mask_only:
111             mask = torch.zeros(self.Lmax, dtype=torch.bool)
112             mask[:len(seq)] = True
113             return {'mask':mask}, {'mask':mask}
114         seq = [self.seq_map[s] for s in seq]
115         seq = np.array(seq)
116         mask = torch.zeros(self.Lmax, dtype=torch.bool)
117         mask[:len(seq)] = True
118         seq = np.pad(seq, (0, self.Lmax-len(seq)))
119
120         react = torch.from_numpy(np.stack([self.react_2A3[idx],
121                                             self.react_DMS[idx]], -1))
122         react_err = torch.from_numpy(np.stack([self.react_err_2A3[idx],
123                                                 self.react_err_DMS[idx]], -1))
124         sn = torch.FloatTensor([self.sn_2A3[idx], self.sn_DMS[idx]])
125
126         return {'seq':torch.from_numpy(seq), 'mask':mask}, \
127               {'react':react, 'react_err':react_err,
128                'sn':sn, 'mask':mask}
129
130 class LenMatchBatchSampler(torch.utils.data.BatchSampler):
131     def __iter__(self):
132         buckets = [[]] * 100
133         yielded = 0
134
135         for idx in self.sampler:
136             s = self.sampler.data_source[idx]
137             if isinstance(s, tuple): L = s[0]["mask"].sum()

```

```

138         else: L = s["mask"].sum()
139         L = max(1, L // 16)
140         if len(buckets[L]) == 0: buckets[L] = []
141         buckets[L].append(idx)
142
143         if len(buckets[L]) == self.batch_size:
144             batch = list(buckets[L])
145             yield batch
146             yielded += 1
147             buckets[L] = []
148
149     batch = []
150     leftover = [idx for bucket in buckets for idx in bucket]
151
152     for idx in leftover:
153         batch.append(idx)
154         if len(batch) == self.batch_size:
155             yielded += 1
156             yield batch
157             batch = []
158
159     if len(batch) > 0 and not self.drop_last:
160         yielded += 1
161         yield batch
162
163 def dict_to(x, device='cuda'):
164     return {k:x[k].to(device) for k in x}
165
166 def to_device(x, device='cuda'):
167     return tuple(dict_to(e, device) for e in x)
168
169 class DeviceDataLoader:
170     def __init__(self, dataloader, device='cuda'):
171         self.dataloader = dataloader
172         self.device = device
173
174     def __len__(self):
175         return len(self.dataloader)
176
177     def __iter__(self):
178         for batch in self.dataloader:
179             yield tuple(dict_to(x, self.device) for x in batch)
180
181 class SinusoidalPosEmb(nn.Module):
182     def __init__(self, dim=16, M=10000):
183         super().__init__()
184         self.dim = dim
185         self.M = M
186
187     def forward(self, x):
188         device = x.device
189         half_dim = self.dim // 2
190         emb = math.log(self.M) / half_dim
191         emb = torch.exp(torch.arange(half_dim, device=device) * (-emb))
192         emb = x[..., None] * emb[None, ...]
193         emb = torch.cat((emb.sin(), emb.cos()), dim=-1)
194         return emb
195
196 class RNA_Model(nn.Module):
197     def __init__(self, dim=192, depth=12, head_size=32, **kwargs):
198         super().__init__()
199         self.emb = nn.Embedding(4, dim)
200         self.pos_enc = SinusoidalPosEmb(dim)
201         self.transformer = nn.TransformerEncoder(
202             nn.TransformerEncoderLayer(d_model=dim, nhead=dim//head_size, dim_feedforward=4*dim,
203                 dropout=0.1, activation=nn.GELU(), batch_first=True, norm_first=True), depth)
204         self.proj_out = nn.Linear(dim, 2)
205
206     def forward(self, x0):
207         mask = x0['mask']
208         Lmax = mask.sum(-1).max()
209         mask = mask[:, :Lmax]

```

```

209         mask = mask[:, :Lmax]
210         x = x0['seq'][:, :Lmax]
211
212         pos = torch.arange(Lmax, device=x.device).unsqueeze(0)
213         pos = self.pos_enc(pos)
214         x = self.emb(x)
215         x = x + pos
216
217         x = self.transformer(x, src_key_padding_mask=~mask)
218         x = self.proj_out(x)
219
220         return x
221
222 def loss(pred, target):
223     p = pred[target['mask'][:, :pred.shape[1]]]
224     y = target['react'][target['mask']].clip(0, 1)
225     loss = F.l1_loss(p, y, reduction='none')
226     loss = loss[~torch.isnan(loss)].mean()
227
228     return loss
229
230 class MAE(Metric):
231     def __init__(self):
232         self.reset()
233
234     def reset(self):
235         self.x, self.y = [], []
236
237     def accumulate(self, learn):
238         x = learn.pred[learn.y['mask'][:, :learn.pred.shape[1]]]
239         y = learn.y['react'][learn.y['mask']].clip(0, 1)
240         self.x.append(x)
241         self.y.append(y)
242
243     @property
244     def value(self):
245         x, y = torch.cat(self.x, 0), torch.cat(self.y, 0)
246         loss = F.l1_loss(x, y, reduction='none')
247         loss = loss[~torch.isnan(loss)].mean()
248         return loss
249
250 #df = pd.read_csv(os.path.join(PATH, r'C:\Users\VENKATESH\Downloads\Data_train.csv'))
251
252
253
254 seed_everything(SEED)
255 os.makedirs(OUT, exist_ok=True)
256 PATH="/content/drive/MyDrive/projects_RNA/train_data.parquet"
257 df=pd.read_parquet(PATH)
258
259 for fold in [0]: # running multiple folds at kaggle may cause OOM
260     ds_train = RNA_Dataset(df, mode='train', fold=fold, nfolds=nfolds)
261     ds_train_len = RNA_Dataset(df, mode='train', fold=fold,
262                                nfolds=nfolds, mask_only=True)
263     sampler_train = torch.utils.data.RandomSampler(ds_train_len)
264     len_sampler_train = LenMatchBatchSampler(sampler_train, batch_size=bs,
265                                              drop_last=True)
266     dl_train = DeviceDataLoader(torch.utils.data.DataLoader(ds_train,
267                                                             batch_sampler=len_sampler_train, num_workers=num_workers,
268                                                             persistent_workers=True), device)
269
270     ds_val = RNA_Dataset(df, mode='eval', fold=fold, nfolds=nfolds)
271     ds_val_len = RNA_Dataset(df, mode='eval', fold=fold, nfolds=nfolds,
272                              mask_only=True)
273     sampler_val = torch.utils.data.SequentialSampler(ds_val_len)
274     len_sampler_val = LenMatchBatchSampler(sampler_val, batch_size=bs,
275                                           drop_last=False)
276     dl_val= DeviceDataLoader(torch.utils.data.DataLoader(ds_val,
277                                                         batch_sampler=len_sampler_val, num_workers=num_workers), device)
278     gc.collect()
279
280     data = DataLoaders(dl_train, dl_val)

```

```

281 model = RNA_Model()
282 model = model.to(device)
283 learn = Learner(data, model, loss_func=loss, cbs=[GradientClip(3.0)],
284                 metrics=[MAE()]).to_fp16()
285 #fp16 doesn't help at P100 but gives x1.6-1.8 speedup at modern hardware
286
287 learn.fit_one_cycle(32, lr_max=5e-4, wd=0.05, pct_start=0.02)
288 torch.save(learn.model.state_dict(), os.path.join(OUT, f'{fname}_{fold}.pth'))
289 gc.collect()
290
291 import gc
292 import os
293 import time
294 import pandas as pd
295 import numpy as np
296 import json
297 import torch
298 from fastai.data.load import DataLoader
299
300 from datasets import DatasetEightInfer, DatasetTenInfer
301 from models import ModelThirtyNine, ModelThirtyTwo
302 from seed_all import seed_everything
303
304 SUBMISSION_NUMBER = 27 # the setup is shown in this repository for 27 and 23 only
305 MODEL_EPOCH_NUMBER = 27 # 27 for submission number 27, and 44 for submission number 23
306 # (how many epochs the model was trained, starting from zero)
307
308 BATCH = 128
309 COL_A = 'reactivity_2A3_MaP'
310 COL_D = 'reactivity_DMS_MaP'
311
312
313 def batch_to_csv(output, ids, main_path_for_parquets):
314     # received a batch of outputs (B, 459, 2) and ids (B, 4) as numpy arrays
315     name_of_csv = ids[0][0]
316     dfs = []
317     for i in range(output.shape[0]):
318         start_id = ids[i][0]
319         end_id = ids[i][1]
320         start_index = ids[i][2]
321         num_reactivities = ids[i][3]
322         # Extract relevant reactivities from output[i]
323         reactivities_a = output[i, start_index: start_index + num_reactivities, 0]
324         reactivities_d = output[i, start_index: start_index + num_reactivities, 1]
325         # Create a DataFrame for the current datapoint
326         datapoint_df = pd.DataFrame({
327             'id': np.arange(start_id, end_id + 1),
328             COL_D: reactivities_d,
329             COL_A: reactivities_a
330         })
331         dfs.append(datapoint_df)
332     small_df = pd.concat(dfs, ignore_index=True)
333     # the df will be written into .parquet
334     path = os.path.join(main_path_for_parquets, f'{name_of_csv}.parquet')
335     small_df.to_parquet(path, index=False, engine='pyarrow')
336     return
337
338
339 # before running, folder ../submissions/{SUBMISSION_NUMBER}/all needs to already exist
340 # for submission number 23, it runs for a very long time (eight plus hours) because bpps are not saved
341 # and need to be calculated in dataset
342 if __name__ == '__main__':
343     seed_everything()
344     with open('SETTINGS.json') as f:
345         data = json.load(f)
346         path_to_test_data = data["TEST_DATA"]
347         model_dir = data["MODEL_DIR"]
348         submission_dir = data["SUBMISSION_DIR"]
349         model_string = f'{SUBMISSION_NUMBER}/models/model_{MODEL_EPOCH_NUMBER}.pth'
350         path_to_model = os.path.join(model_dir, model_string)
351         main_path_string = f'{SUBMISSION_NUMBER}/all/'
352         main_path_for_parquets = os.path.join(submission_dir, main_path_string)

```

```

352 main_path_for_parquets = os.path.join(submission_dir, main_path_string)
353
354 df = pd.read_parquet(path_to_test_data, engine='pyarrow')
355 # device
356 device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
357 print(device)
358
359 if SUBMISSION_NUMBER == 27:
360     dataset_skeleton = DatasetEightInfer
361     model_skeleton = ModelThirtyNine
362     num_workers = 0
363 elif SUBMISSION_NUMBER == 23:
364     dataset_skeleton = DatasetTenInfer
365     model_skeleton = ModelThirtyTwo
366     num_workers = 40
367
368 # dataset and dataloader
369 dataset = dataset_skeleton(df=df)
370 loader = DataLoader(dataset=dataset, batch_size=BATCH, pin_memory=False, shuffle=False, device=device,
371                     num_workers=num_workers) # num_workers is set to 40 for bpps (submission number 23)
372
373 # model
374 model = model_skeleton()
375 # load the state dict
376 model.load_state_dict(torch.load(path_to_model))
377 model.eval()
378 model.to(device)
379
380 # Start timer
381 start_time = time.time()
382 with torch.no_grad():
383     i = 0
384     for data, ids in loader:
385         i += 1
386         out = model(data)
387         batch_to_csv(out.detach().cpu().numpy(), ids.detach().cpu().numpy(), main_path_for_parquets)
388         if i % 50 == 0:
389             print(f"step {i}")
390 # End timer
391 end_time = time.time()
392 # Calculate elapsed time
393 elapsed_time = end_time - start_time
394 print("Elapsed time: ", elapsed_time)
395
396
397
398
399
400
401
402
403

```

```

➡ /usr/local/lib/python3.10/dist-packages/torch/nn/modules/transformer.py:286: UserWarning: enable_nested_tensor
  warnings.warn(f"enable_nested_tensor is True, but self.use_nested_tensor is False because {why_not_sparsity_
/usr/local/lib/python3.10/dist-packages/torch/amp/autocast_mode.py:250: UserWarning: User provided device_type
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/torch/cuda/amp/grad_scaler.py:126: UserWarning: torch.cuda.amp.GradSca
  warnings.warn(
0.00% [0/32 00:00<?]

epoch train_loss valid_loss mae time
0.00% [0/531 00:00<?]

```

Double-click (or enter) to edit

1 # install packages

```
1 pip install datasets
```



Collecting datasets

Downloading datasets-2.18.0-py3-none-any.whl (510 kB)

510.5/510.5 kB 6.0 MB/s eta 0:00:00

Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from datasets) (3.13.1)

Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from datasets) (1.25.2)

Requirement already satisfied: pyarrow>=12.0.0 in /usr/local/lib/python3.10/dist-packages (from datasets) (14.0.2)

Requirement already satisfied: pyarrow-hotfix in /usr/local/lib/python3.10/dist-packages (from datasets) (0.6)

Collecting dill<0.3.9,>=0.3.0 (from datasets)

Downloading dill-0.3.8-py3-none-any.whl (116 kB)

116.3/116.3 kB 13.1 MB/s eta 0:00:00

Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from datasets) (1.5.3)

Requirement already satisfied: requests>=2.19.0 in /usr/local/lib/python3.10/dist-packages (from datasets) (2.31.0)

Requirement already satisfied: tqdm>=4.62.1 in /usr/local/lib/python3.10/dist-packages (from datasets) (4.66.2)

Requirement already satisfied: xxhash in /usr/local/lib/python3.10/dist-packages (from datasets) (3.4.1)

Collecting multiprocessing (from datasets)

Downloading multiprocessing-0.70.16-py310-none-any.whl (134 kB)

134.8/134.8 kB 13.5 MB/s eta 0:00:00

Requirement already satisfied: fsspec[http]<=2024.2.0,>=2023.1.0 in /usr/local/lib/python3.10/dist-packages (from datasets) (2024.2.0)

Requirement already satisfied: aiohttp in /usr/local/lib/python3.10/dist-packages (from datasets) (3.9.3)

Requirement already satisfied: huggingface-hub>=0.19.4 in /usr/local/lib/python3.10/dist-packages (from datasets) (0.23.0)

Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from datasets) (23.2)

Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from datasets) (6.0.1)

Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.3.1)

Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (23.2.0)

Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.4.1)

Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (6.0.5)

Requirement already satisfied: yarl<2.0,>=1.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.9.7)

Requirement already satisfied: async-timeout<5.0,>=4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (4.0.3)

Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (4.9.0)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->datasets) (3.3.2)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->datasets) (3.6)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->datasets) (2.0.7)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->datasets) (2024.2.2)

Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas->datasets) (2.8.2)

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil->datasets) (1.16.0)

Installing collected packages: dill, multiprocessing, datasets

Successfully installed datasets-2.18.0 dill-0.3.8 multiprocessing-0.70.16

```
1 !pip install rotary_embedding_torch
```



Collecting rotary_embedding_torch

Downloading rotary_embedding_torch-0.5.3-py3-none-any.whl (5.3 kB)

Collecting beartype (from rotary_embedding_torch)

Downloading beartype-0.17.2-py3-none-any.whl (872 kB)

872.4/872.4 kB 18.8 MB/s eta 0:00:00

Collecting einops>=0.7 (from rotary_embedding_torch)

Downloading einops-0.7.0-py3-none-any.whl (44 kB)

44.6/44.6 kB 6.3 MB/s eta 0:00:00

Requirement already satisfied: torch>=2.0 in /usr/local/lib/python3.10/dist-packages (from rotary_embedding_torch) (2.3.0)

Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (3.13.1)

Requirement already satisfied: typing-extensions>=4.8.0 in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (4.9.0)

Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (1.12.0)

Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (3.2.1)

Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (3.1.3)

Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch) (2024.2.0)

Collecting nvidia-cuda-nvrtc-cu12==12.1.105 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cuda_nvrtc_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (23.7 MB)

23.7/23.7 MB 49.7 MB/s eta 0:00:00

Collecting nvidia-cuda-runtime-cu12==12.1.105 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cuda_runtime_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (823 kB)

823.6/823.6 kB 60.9 MB/s eta 0:00:00

Collecting nvidia-cuda-cupti-cu12==12.1.105 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cuda_cupti_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (14.1 MB)

14.1/14.1 MB 63.4 MB/s eta 0:00:00

Collecting nvidia-cudnn-cu12==8.9.2.26 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cudnn_cu12-8.9.2.26-py3-none-manylinux1_x86_64.whl (731.7 MB)

731.7/731.7 MB 1.6 MB/s eta 0:00:00

Collecting nvidia-cublas-cu12==12.1.3.1 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cublas_cu12-12.1.3.1-py3-none-manylinux1_x86_64.whl (410.6 MB)

410.6/410.6 MB 2.9 MB/s eta 0:00:00

Collecting nvidia-cufft-cu12==11.0.2.54 (from torch>=2.0->rotary_embedding_torch)

Downloading nvidia_cufft_cu12-11.0.2.54-py3-none-manylinux1_x86_64.whl (121.6 MB)

121.6/121.6 MB 8.5 MB/s eta 0:00:00

```

Collecting nvidia-curand-cu12==10.3.2.106 (from torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_curand_cu12-10.3.2.106-py3-none-manylinux1_x86_64.whl (56.5 MB)
    56.5/56.5 MB 11.7 MB/s eta 0:00:00
Collecting nvidia-cusolver-cu12==11.4.5.107 (from torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_cusolver_cu12-11.4.5.107-py3-none-manylinux1_x86_64.whl (124.2 MB)
    124.2/124.2 MB 5.5 MB/s eta 0:00:00
Collecting nvidia-cuspars-cu12==12.1.0.106 (from torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_cuspars_cu12-12.1.0.106-py3-none-manylinux1_x86_64.whl (196.0 MB)
    196.0/196.0 MB 2.4 MB/s eta 0:00:00
Collecting nvidia-nccl-cu12==2.19.3 (from torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_nccl_cu12-2.19.3-py3-none-manylinux1_x86_64.whl (166.0 MB)
    166.0/166.0 MB 6.3 MB/s eta 0:00:00
Collecting nvidia-nvtx-cu12==12.1.105 (from torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_nvtx_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (99 kB)
    99.1/99.1 kB 11.3 MB/s eta 0:00:00
Requirement already satisfied: triton==2.2.0 in /usr/local/lib/python3.10/dist-packages (from torch>=2.0->rotary_embedding_torch)
Collecting nvidia-nvjitlink-cu12 (from nvidia-cusolver-cu12==11.4.5.107->torch>=2.0->rotary_embedding_torch)
  Downloading nvidia_nvjitlink_cu12-12.4.99-py3-none-manylinux2014_x86_64.whl (21.1 MB)
    21.1/21.1 MB 69.1 MB/s eta 0:00:00
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch>=2.0->rotary_embedding_torch)
Requirement already satisfied: mpmath>=0.19 in /usr/local/lib/python3.10/dist-packages (from sympy->torch>=2.0->rotary_embedding_torch)
Installing collected packages: nvidia-nvjitlink-cu12, nvidia-nvtx-cu12, nvidia-nccl-cu12, nvidia-curand-cu12,
Successfully installed beartype-0.17.2 einops-0.7.0 nvidia-cublas-cu12-12.1.3.1 nvidia-cuda-cupti-cu12-12.1.106

```

```

1 import torch
2 import torch.nn as nn
3 import math
4 from rotary_embedding_torch import RotaryEmbedding
5
6 LEN = 457
7 LEN_EOS = 459
8 LEN_FOR_GENERALIZATION = 722
9
10 #####
11 # the code for building transformer (building blocks) is from
12 # https://towardsdatascience.com/build-your-own-transformer-from-scratch-using-pytorch-84c850470dcb
13
14 # the way how sinusoidal embedding is calculated is from https://www.kaggle.com/code/iafoss/rna-starter-0-186
15 class PosEnc(nn.Module):
16     """
17     sinusoidal embeddings
18     """
19     def __init__(self, dim=192, M=10000, num_tokens=LEN_EOS):
20         super().__init__()
21         positions = torch.arange(num_tokens).unsqueeze(0)
22         half_dim = dim // 2
23         emb = math.log(M) / half_dim
24         emb = torch.exp(torch.arange(half_dim) * (-emb))
25         emb = positions[..., None] * emb[None, ...]
26         emb = torch.cat((emb.sin(), emb.cos()), dim=-1)
27         self.pos = emb
28
29     def forward(self, x):
30         device = x.device
31         pos = self.pos.to(device)
32         res = x + pos
33         return res
34
35
36 # https://towardsdatascience.com/build-your-own-transformer-from-scratch-using-pytorch-84c850470dcb
37 class MultiHeadAttention(nn.Module):
38     def __init__(self, d_model, num_heads):
39         super(MultiHeadAttention, self).__init__()
40         assert d_model % num_heads == 0, "d_model must be divisible by num_heads"
41
42         self.d_model = d_model
43         self.num_heads = num_heads
44         self.d_k = d_model // num_heads
45
46         self.W_q = nn.Linear(d_model, d_model)
47         self.W_k = nn.Linear(d_model, d_model)

```



```

48     self.W_v = nn.Linear(d_model, d_model)
49     self.W_o = nn.Linear(d_model, d_model)
50
51     def scaled_dot_product_attention(self, Q, K, V, mask=None):
52         attn_scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d_k)
53         if mask is not None:
54             _MASKING_VALUE = -1e+30 if attn_scores.dtype == torch.float32 else -1e+4
55             attn_scores = attn_scores.masked_fill(mask == 0, _MASKING_VALUE)
56         attn_probs = torch.softmax(attn_scores, dim=-1)
57         output = torch.matmul(attn_probs, V)
58         return output
59
60     def split_heads(self, x):
61         batch_size, seq_length, d_model = x.size()
62         return x.view(batch_size, seq_length, self.num_heads, self.d_k).transpose(1, 2)
63
64     def combine_heads(self, x):
65         batch_size, _, seq_length, d_k = x.size()
66         return x.transpose(1, 2).contiguous().view(batch_size, seq_length, self.d_model)
67
68     def forward(self, Q, K, V, mask=None):
69         Q = self.split_heads(self.W_q(Q))
70         K = self.split_heads(self.W_k(K))
71         V = self.split_heads(self.W_v(V))
72
73         attn_output = self.scaled_dot_product_attention(Q, K, V, mask)
74         output = self.W_o(self.combine_heads(attn_output))
75         return output
76
77
78 class AttentionRotary(nn.Module):
79     def __init__(self, d_model, num_heads, rotary_emb):
80         super(AttentionRotary, self).__init__()
81         assert d_model % num_heads == 0, "d_model must be divisible by num_heads"
82
83         self.d_model = d_model
84         self.num_heads = num_heads
85         self.d_k = d_model // num_heads
86         self.rotary_emb = rotary_emb
87
88         self.W_q = nn.Linear(d_model, d_model)
89         self.W_k = nn.Linear(d_model, d_model)
90         self.W_v = nn.Linear(d_model, d_model)
91         self.W_o = nn.Linear(d_model, d_model)
92
93     def scaled_dot_product_attention(self, Q, K, V, mask=None):
94         attn_scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d_k)
95         if mask is not None:
96             _MASKING_VALUE = -1e+30 if attn_scores.dtype == torch.float32 else -1e+4
97             attn_scores = attn_scores.masked_fill(mask == 0, _MASKING_VALUE)
98         attn_probs = torch.softmax(attn_scores, dim=-1)
99         output = torch.matmul(attn_probs, V)
100         return output
101
102     def split_heads(self, x):
103         batch_size, seq_length, d_model = x.size()
104         return x.view(batch_size, seq_length, self.num_heads, self.d_k).transpose(1, 2)
105
106     def combine_heads(self, x):
107         batch_size, _, seq_length, d_k = x.size()
108         return x.transpose(1, 2).contiguous().view(batch_size, seq_length, self.d_model)
109
110     def forward(self, Q, K, V, mask=None):
111         Q = self.split_heads(self.W_q(Q))
112         Q = self.rotary_emb.rotate_queries_or_keys(Q)
113         K = self.split_heads(self.W_k(K))
114         K = self.rotary_emb.rotate_queries_or_keys(K)
115         V = self.split_heads(self.W_v(V))
116
117         attn_output = self.scaled_dot_product_attention(Q, K, V, mask)
118         output = self.W_o(self.combine_heads(attn_output))

```

```

119         return output
120
121
122 class CustomAttentionBPP(nn.Module):
123     def __init__(self, d_model, num_heads=1):
124         super(CustomAttentionBPP, self).__init__()
125         assert d_model % num_heads == 0, "d_model must be divisible by num_heads"
126
127         self.d_model = d_model
128         self.num_heads = num_heads
129         self.d_k = d_model // num_heads
130
131         self.W_v = nn.Linear(d_model, d_model)
132         self.W_o = nn.Linear(d_model, d_model)
133
134     def scaled_dot_product_attention(self, bpp, V, mask=None):
135         attn_scores = bpp.unsqueeze(1)
136         _MASKING_VALUE = -1e+30 if attn_scores.dtype == torch.float32 else -1e+4
137         attn_scores = attn_scores.masked_fill(attn_scores == 0, _MASKING_VALUE)
138         if mask is not None:
139             attn_scores = attn_scores.masked_fill(mask == 0, _MASKING_VALUE)
140         attn_probs = torch.softmax(attn_scores, dim=-1)
141         output = torch.matmul(attn_probs, V)
142         return output
143
144     def split_heads(self, x):
145         batch_size, seq_length, d_model = x.size()
146         return x.view(batch_size, seq_length, self.num_heads, self.d_k).transpose(1, 2)
147
148     def combine_heads(self, x):
149         batch_size, _, seq_length, d_k = x.size()
150         return x.transpose(1, 2).contiguous().view(batch_size, seq_length, self.d_model)
151
152     def forward(self, bpp, V, mask=None):
153         V = self.split_heads(self.W_v(V))
154
155         attn_output = self.scaled_dot_product_attention(bpp, V, mask)
156         output = self.W_o(self.combine_heads(attn_output))
157         return output
158
159
160 # https://towardsdatascience.com/build-your-own-transformer-from-scratch-using-pytorch-84c850470dcb
161 # gelu is used instead of relu
162 class PositionWiseFeedForward(nn.Module):
163     def __init__(self, d_model, d_ff):
164         super(PositionWiseFeedForward, self).__init__()
165         self.fc1 = nn.Linear(d_model, d_ff)
166         self.fc2 = nn.Linear(d_ff, d_model)
167         self.gelu = nn.GELU()
168
169     def forward(self, x):
170         return self.fc2(self.gelu(self.fc1(x)))
171
172
173 # https://towardsdatascience.com/build-your-own-transformer-from-scratch-using-pytorch-84c850470dcb
174 # with minor modifications
175 class EncoderLayer(nn.Module):
176     def __init__(self, d_model, num_heads, d_ff, dropout):
177         super(EncoderLayer, self).__init__()
178         self.self_attn = MultiHeadAttention(d_model, num_heads)
179         self.feed_forward = PositionWiseFeedForward(d_model, d_ff)
180         self.norm1 = nn.LayerNorm(d_model)
181         self.norm2 = nn.LayerNorm(d_model)
182         self.dropout = nn.Dropout(dropout)
183
184     def forward(self, x, mask):
185         attn_output = self.self_attn(x, x, x, mask)
186         x = self.norm1(x + self.dropout(attn_output))
187         ff_output = self.feed_forward(x)
188         x = self.norm2(x + self.dropout(ff_output))
189         return x

```

```

190
191
192 class EncoderLayerRotary(nn.Module):
193     def __init__(self, d_model, num_heads, d_ff, dropout, rotary_emb):
194         super(EncoderLayerRotary, self).__init__()
195         self.self_attn = AttentionRotary(d_model, num_heads, rotary_emb)
196         self.feed_forward = PositionWiseFeedForward(d_model, d_ff)
197         self.norm1 = nn.LayerNorm(d_model)
198         self.norm2 = nn.LayerNorm(d_model)
199         self.dropout = nn.Dropout(dropout)
200
201     def forward(self, x, mask):
202         attn_output = self.self_attn(x, x, x, mask)
203         x = self.norm1(x + self.dropout(attn_output))
204         ff_output = self.feed_forward(x)
205         x = self.norm2(x + self.dropout(ff_output))
206         return x
207
208
209 # https://towardsdatascience.com/build-your-own-transformer-from-scratch-using-pytorch-84c850470dcb
210 class DecoderLayer(nn.Module):
211     def __init__(self, d_model, num_heads, d_ff, dropout):
212         super(DecoderLayer, self).__init__()
213         self.self_attn = MultiHeadAttention(d_model, num_heads)
214         self.cross_attn = MultiHeadAttention(d_model, num_heads)
215         self.feed_forward = PositionWiseFeedForward(d_model, d_ff)
216         self.norm1 = nn.LayerNorm(d_model)
217         self.norm2 = nn.LayerNorm(d_model)
218         self.norm3 = nn.LayerNorm(d_model)
219         self.dropout = nn.Dropout(dropout)
220
221     def forward(self, x, enc_output, src_mask, tgt_mask):
222         attn_output = self.self_attn(x, x, x, tgt_mask)
223         x = self.norm1(x + self.dropout(attn_output))
224         attn_output = self.cross_attn(x, enc_output, enc_output, src_mask)
225         x = self.norm2(x + self.dropout(attn_output))
226         ff_output = self.feed_forward(x)
227         x = self.norm3(x + self.dropout(ff_output))
228         return x
229
230
231 class DecoderLayerRotary(nn.Module):
232     def __init__(self, d_model, num_heads, d_ff, dropout, rotary_emb):
233         super(DecoderLayerRotary, self).__init__()
234         self.self_attn = AttentionRotary(d_model, num_heads, rotary_emb)
235         self.cross_attn = AttentionRotary(d_model, num_heads, rotary_emb)
236         self.feed_forward = PositionWiseFeedForward(d_model, d_ff)
237         self.norm1 = nn.LayerNorm(d_model)
238         self.norm2 = nn.LayerNorm(d_model)
239         self.norm3 = nn.LayerNorm(d_model)
240         self.dropout = nn.Dropout(dropout)
241
242     def forward(self, x, enc_output, src_mask, tgt_mask):
243         attn_output = self.self_attn(x, x, x, tgt_mask)
244         x = self.norm1(x + self.dropout(attn_output))
245         attn_output = self.cross_attn(x, enc_output, enc_output, src_mask)
246         x = self.norm2(x + self.dropout(attn_output))
247         ff_output = self.feed_forward(x)
248         x = self.norm3(x + self.dropout(ff_output))
249         return x
250
251
252 # similar to DecoderLayer, but as cross_attn, it uses CustomAttentionBPP
253 class DecoderLayerTwo(nn.Module):
254     def __init__(self, d_model, num_heads, d_ff, dropout):
255         super(DecoderLayerTwo, self).__init__()
256         self.self_attn = MultiHeadAttention(d_model, num_heads)
257         self.cross_attn = CustomAttentionBPP(d_model)
258         self.feed_forward = PositionWiseFeedForward(d_model, d_ff)
259         self.norm1 = nn.LayerNorm(d_model)
260         self.norm2 = nn.LayerNorm(d_model)

```

```

261         self.norm3 = nn.LayerNorm(d_model)
262         self.dropout = nn.Dropout(dropout)
263
264     def forward(self, x, bpp, mask):
265         attn_output = self.self_attn(x, x, x, mask)
266         x = self.norm1(x + self.dropout(attn_output))
267         attn_output = self.cross_attn(bpp=bpp, V=x, mask=mask)
268         x = self.norm2(x + self.dropout(attn_output))
269         ff_output = self.feed_forward(x)
270         x = self.norm3(x + self.dropout(ff_output))
271         return x
272
273 #####
274 # models:
275
276
277 # first it is decoder layer to use bpp (with sinusoidal pos embeds), then uses rotary embeddings
278 # tgt, info1: seq_inds; info2: bpp; src or info3: struct_inds
279 class ModelThirtyTwo(nn.Module):
280     def __init__(self, tgt_vocab=7, src_vocab=6, d_model=192, num_heads=6, num_layers=8,
281                 d_ff=(192*4), dropout=0.1, num_tokens=LEN_EOS):
282         super(ModelThirtyTwo, self).__init__()
283         self.tgt_embedding = nn.Embedding(tgt_vocab, d_model)
284         self.src_embedding = nn.Embedding(src_vocab, d_model)
285         self.positional_enc = PosEnc(dim=d_model, num_tokens=num_tokens)
286         self.rotary = RotaryEmbedding(dim=d_model//num_heads)
287         self.decoder_one = DecoderLayerTwo(d_model, num_heads, d_ff, dropout)
288         self.decoder = DecoderLayerRotary(d_model, num_heads, d_ff, dropout, self.rotary)
289         self.encoder_layers = nn.ModuleList([EncoderLayerRotary(d_model, num_heads, d_ff, dropout, self.rotary)
290         self.fc = nn.Linear(d_model, 2)
291
292     def forward(self, data):
293         tgt = data['info1']
294         bpp = data['info2']
295         src = data['info3']
296         mask = data['mask']
297
298         mask = mask.unsqueeze(1).unsqueeze(2)
299         src = self.src_embedding(src)
300         tgt = self.positional_enc(self.tgt_embedding(tgt))
301
302         output = self.decoder_one(x=tgt, bpp=bpp, mask=mask)
303
304         output = self.decoder(x=output, enc_output=src, src_mask=mask, tgt_mask=mask)
305         for enc_layer in self.encoder_layers:
306             output = enc_layer(output, mask)
307
308         output = self.fc(output)
309         return output
310
311
312 class ModelThirtyNine(nn.Module):
313     def __init__(self, tgt_vocab=7, src_vocab=6, d_model=384, num_heads=6, num_layers=8, d_ff=384, dropout=0):
314         super(ModelThirtyNine, self).__init__()
315         self.tgt_embedding = nn.Embedding(tgt_vocab, d_model)
316         self.src_embedding = nn.Embedding(src_vocab, d_model)
317         self.positional_enc = RotaryEmbedding(dim=d_model//num_heads)
318         self.decoder = DecoderLayerRotary(d_model, num_heads, d_ff, dropout, self.positional_enc)
319         self.encoder_layers = nn.ModuleList([EncoderLayerRotary(d_model, num_heads, d_ff, dropout, self.positional_enc)
320         self.fc = nn.Linear(d_model, 2)
321
322     def forward(self, data):
323         tgt = data['info1']
324         src = data['info2']
325         mask = data['mask']
326
327         mask = mask.unsqueeze(1).unsqueeze(2)
328         tgt = self.tgt_embedding(tgt)
329         src = self.src_embedding(src)
330
331         output = self.decoder(x=tgt, enc_output=src, src_mask=mask, tgt_mask=mask)

```

```
332     for enc_layer in self.encoder_layers:
333         output = enc_layer(output, mask)
334
335     output = self.fc(output)
336     return output
```