

```

1 from typing import Tuple
2
3 from torch import LongTensor, cat
4 from transformers import PreTrainedTokenizer
5 from transformers.tokenization_utils_base import TruncationStrategy, BatchEncoding
6
7
8 class PreProcessor:
9     framework: str = "pt"
10
11     def __init__(
12         self,
13         tokenizer: PreTrainedTokenizer,
14         max_input_len: int,
15     ):
16         self.tokenizer: PreTrainedTokenizer = tokenizer
17         self.max_input_len: int = max_input_len
18
19     def get_token_tensor(
20         self,
21         text: str,
22         truncation: TruncationStrategy = TruncationStrategy.DO_NOT_TRUNCATE,
23     ) -> LongTensor:
24         """Complexity: O(n) where n is the number of characters in the text"""
25         if len(text) == 0:
26             return LongTensor([])
27         # tokenizing a string is O(n) where n is the length of the string
28         tokenized_text = self.tokenizer(
29             text,
30             return_tensors=self.framework,
31             padding=False,
32             add_special_tokens=False,
33             truncation=truncation,
34             max_length=self.max_input_len,
35         )
36         if isinstance(tokenized_text, dict) or isinstance(tokenized_text, BatchEncoding):
37             token_tensor: LongTensor = tokenized_text["input_ids"]
38             # O(1) because we are accessing a single element
39             # in a dictionary and saving the reference to it.
40         elif isinstance(tokenized_text, LongTensor):
41             token_tensor: LongTensor = tokenized_text
42             # O(1) because we are saving the reference to the tensor
43         else:
44             raise TypeError(
45                 "The tokenizer output is not one of:"
46                 "dict, BatchEncoding, LongTensor"
47             )
48         token_tensor = token_tensor.squeeze()
49         # O(n) where n is the number of tokens in the text
50         # because we are copying n elements from one tensor to the other
51         # the number of tokens in the text
52         # is always less than the number of characters in the text
53
54         return token_tensor
55
56     def __call__(
57         self,
58         prompt: str,
59         prefix: str = "",
60         truncation: TruncationStrategy = TruncationStrategy.DO_NOT_TRUNCATE,
61         postfix: str = "",
62     ) -> Tuple[LongTensor, int, int]:
63         """A helper method for __call__ that tokenize the prompt
64         all the arguments are sent directly from the __call__ method
65         Returns:
66             the tokenized prompt as a list of ints,
67             the length of the prompt,
68             the length of the prefix,
69             the length of the postfix
70
71         Complexity: O(a + b + c) where:
72             'a' is the number of characters in the prefix
73             'b' is the number of characters in the prompt
74             'c' is the number of characters in the postfix"""
75
76         prefix_tokens: LongTensor = self.get_token_tensor(prefix, truncation)
77         # O(A) where A is the number of characters in the prefix.
78         postfix_tokens: LongTensor = self.get_token_tensor(postfix, truncation)
79         # O(B) where B is the number of characters in the postfix.
80         prompt_tokens: LongTensor = self.get_token_tensor(prompt, truncation)
81         # O(C) where C is the number of characters in the prompt.
82         token_tensor: LongTensor = cat((prefix_tokens, prompt_tokens, postfix_tokens))
83         # O( + b + c) where 'a' is the number of tokens in the prefix.
84         # 'b' is the number of tokens in the prompt.
85         # 'c' is the number of tokens in the postfix.
86         # we know that the number of tokens is less than the number of characters

```

```

87         # We know that the number of tokens is less than the number of characters
88         return token_tensor, len(prefix_tokens), len(prompt_tokens), len(postfix_tokens)
89
90 from typing import Optional
91
92 from torch import tensor
93 from transformers import PreTrainedTokenizer
94
95 from src.grouped_sampling import TokenIDS
96
97
98 class PostProcessor:
99     def __init__(
100         self,
101         tokenizer: PreTrainedTokenizer,
102     ):
103         self.tokenizer: PreTrainedTokenizer = tokenizer
104
105     def __call__(
106         self,
107         token_ids: TokenIDS,
108         num_new_tokens: Optional[int],
109         prompt_len: int,
110         return_text: bool,
111         return_tensors: bool,
112         return_full_text: bool,
113         clean_up_tokenization_spaces: bool,
114         prefix_len: int = 0,
115         postfix_len: int = 0,
116     ):
117         """A helper method for __call__
118         that converts the token ids to dictionary
119         token_ids - the token ids from the _forward method
120         prompt_len - the length of the tokenized prompt
121         the rest of the arguments are the arguments
122         from the __call__ method
123         look up the documentation of the __call__ method for more info
124         Complexity: O(n) where n is the number of tokens in token_ids
125         """
126         # define n as the length of the token_ids
127         full_prompt_len = prefix_len + prompt_len + postfix_len
128         if num_new_tokens is None:
129             shorten_token_list = token_ids
130             # O(1)
131         else:
132             final_num_tokens = full_prompt_len + num_new_tokens
133             # O(1)
134             if len(token_ids) > final_num_tokens:
135                 shorten_token_list = token_ids[:final_num_tokens]
136                 # O(final_num_tokens)
137                 # because we are copying final_num_tokens
138                 # elements from one list to the other
139             else:
140                 shorten_token_list = token_ids
141                 # O(1)
142
143         generated_tokens = shorten_token_list[full_prompt_len:]
144         # O(num_new_tokens) because we are copying
145         # num_new_tokens elements from one list to the other
146         if return_full_text:
147             prompt_tokens = shorten_token_list[prefix_len:prefix_len + prompt_len]
148             # Without prefix and postfix
149             # O(prompt_len) because we are copying prompt_len
150             # elements from one list to the other
151             final_token_list = prompt_tokens + generated_tokens
152             # O(prompt_len + num_new_tokens)
153             # because we are copying elements from one list to the other
154         else:
155             final_token_list = generated_tokens
156             # O(1) because we are saving the reference to the list
157         final_ans = {}
158         if return_tensors:
159             final_ans["generated_token_ids"] = tensor(final_token_list)
160             # O(prompt_len + num_new_tokens)
161             # because we are copying at maximum
162             # prompt_len + num_new_tokens elements from a list to a tensor
163         if return_text:
164             final_ans["generated_text"] = self.tokenizer.decode(
165                 final_token_list,
166                 skip_special_tokens=True,
167                 clean_up_tokenization_spaces=clean_up_tokenization_spaces
168             )
169             # decoding is O(m)
170             # where m is the number of tokens in the text
171             # so the complexity of this line is O(n)
172             # because final_token_list could be at most n tokens long
173         return final_ans

```

```

173         return final_ans
174
175 from abc import ABC, abstractmethod
176 from typing import Set, Optional
177
178 from torch import Tensor
179
180 from .token_ids import TokenIDS
181
182
183 class RepetitionPenaltyStrategy(ABC):
184     theta: float
185
186     @staticmethod
187     def get_already_generated_tokens(
188         tokens: TokenIDS,
189         generation_start: int
190     ) -> Set[int]:
191         return set(tokens[generation_start:])
192
193     @abstractmethod
194     def __call__(
195         self,
196         logits: Tensor,
197         tokens: TokenIDS,
198         generation_start: int
199     ) -> Tensor:
200         raise NotImplementedError
201
202     def __str__(self) -> str:
203         return f"{self.__class__.__name__}(theta={self.theta})"
204
205     def __repr__(self) -> str:
206         return str(self)
207
208
209 class LogitScalingRepetitionPenalty(
210     RepetitionPenaltyStrategy
211 ):
212     def __init__(self, theta: float):
213         if theta <= 1:
214             raise ValueError("Theta must be > 1")
215         self.theta = theta
216
217     def __call__(
218         self,
219         logits: Tensor,
220         tokens: TokenIDS,
221         generation_start: int
222     ) -> Tensor:
223         """applies repetition penalty,
224         using the repetition_penalty_theta parameter defined in the class
225         the formula from the paper is:
226         softmax(original_logits, T) =
227             exp(original_logits_i / (T * h(i)))
228             / sum(exp(original_logits_i / (T * h(i))))
229         which is equivalent to:
230         pi = exp((xi/h(i)) / T / sum(exp(xj/(T * h(j))))
231         and to:
232         penalized_logits = original_logits / h(i)
233         pi = softmax(original_logits / h(i), T)
234         where:
235             h(i) = 0 if i in generated_tokens else 1
236             T is the temperature parameter of softmax
237         Args:
238             logits: the logits matrix of shape (group_size, vocab_size)
239             tokens: the sequence of tokens from the prompt and the generated
240             generation_start: The index of the first
241         Returns:
242             original_logits / h(i)
243             complexity: O(group_size * n)
244             where n is the number of tokens generated by the algorithm
245         """
246         generated_tokens = self.get_already_generated_tokens(
247             tokens,
248             generation_start
249         )
250         for token_id in generated_tokens:
251             # len(generated_tokens) < max(n, vocab_size)
252             logits[:, token_id] /= self.theta
253             # O(group_size) because we are taking a slice of size group_size
254         return logits
255
256
257 class NoRepetitionPenalty(
258     RepetitionPenaltyStrategy
259 ):

```

```

260     theta = 1.0
261
262     def __call__(
263         self,
264         logits: Tensor,
265         tokens: TokenIDS,
266         generation_start: int
267     ) -> Tensor:
268         return logits
269
270
271 DEFAULT_REPETITION_PENALTY = LogitScalingRepetitionPenalty(1.2)
272
273
274 def repetition_penalty_factory(
275     theta: Optional[float]
276 ) -> RepetitionPenaltyStrategy:
277     if theta is None:
278         return DEFAULT_REPETITION_PENALTY
279     if theta == 1:
280         return NoRepetitionPenalty()
281     elif theta > 1:
282         return LogitScalingRepetitionPenalty(theta)
283     else:
284         raise ValueError("theta must be greater than 1")
285
286
287 from enum import Enum
288
289
290 class GenerationType(Enum):
291     """The type of generation to use"""
292     GREEDY = "greedy"
293     TOP_K = "top_k"
294     TOP_P = "top_p"
295     TREE = "tree"
296     RANDOM = "random"
297
298     def requires_softmax(self) -> bool:
299         """Whether the generation type requires a softmax"""
300         return self in (self.TOP_K, self.TOP_P, self.TREE, self.RANDOM)
301
302
303 from typing import Dict
304 from warnings import warn
305
306 from torch import cuda, LongTensor, ones, long, Tensor, cat, no_grad
307 from transformers import AutoModelForCausalLM
308 from torch.nn import Softmax
309
310 from .token_ids import TokenIDS
311 from .repetition_penalty import RepetitionPenaltyStrategy
312
313
314 class GroupedGenerationUtils:
315     descriptive_attrs = (
316         "group_size",
317         "end_of_sentence_stop",
318         "temp",
319         "repetition_penalty_strategy",
320     )
321
322     def __init__(
323         self,
324         model_name: str,
325         group_size: int,
326         max_input_len: int,
327         end_of_sentence_id: int,
328         padding_id: int,
329         vocab_size: int,
330         repetition_penalty_strategy: RepetitionPenaltyStrategy,
331         end_of_sentence_stop: bool = True,
332         temp: float = 1.0,
333         use_softmax: bool = True,
334         **kwargs
335     ):
336         """initializes the model wrapper
337         Args:
338             model_name: the name of the model to be used
339             repetition_penalty_strategy:
340                 the strategy to be used for the repetition penalty
341             group_size: the number of next tokens to be generated
342             max_input_len: the maximum length of the input to the model
343             padding_id: the id of the padding token
344             vocab_size: the size of the vocabulary
345             end_of_sentence_stop:
346                 whether to stop when the end of sentence token is generated

```

```

346         whether to stop when the end of sentence token is generated
347     end_of_sentence_id: int
348         the id of the end of sequence token
349     use_softmax: bool
350         true if the model should use softmax,
351         false if it should return the logits
352     **kwargs: the arguments to be passed to the model
353     Complexity: O(1)
354     """
355     self.use_softmax: bool = use_softmax
356     self.end_of_sentence_id: int = end_of_sentence_id
357     self.repetition_penalty_strategy: RepetitionPenaltyStrategy = repetition_penalty_strategy
358     self.group_size: int = group_size
359     self.max_input_len: int = max_input_len
360     self.padding_id: int = padding_id
361     self.end_of_sentence_stop: bool = end_of_sentence_stop
362     self.model = AutoModelForCausalLM.from_pretrained(
363         model_name, **kwargs
364     )
365     self.temp: float = temp
366     self.vocab_size: int = vocab_size
367     if cuda.is_available():
368         self.model = self.model.cuda()
369
370 @property
371 def padding_tokens(self) -> LongTensor:
372     cpu_tokens = ones(self.group_size, dtype=long) * self.padding_id
373     if cuda.is_available():
374         return cpu_tokens.cuda()
375     return cpu_tokens
376
377 def prepare_model_kwargs(
378     self, tokens: TokenIDS
379 ) -> Dict[str, LongTensor]:
380     """preparing the arguments for the model call
381     Args:
382         tokens: the tokens to be sent to the model
383     Returns:
384         a dictionary of the arguments for the model call
385     Complexity: O(group_size + n) where n is the number of tokens
386     """
387     if not isinstance(tokens, Tensor):
388         tokens = LongTensor(tokens) # O(n)
389     padded_tokens: LongTensor = cat(
390         (tokens, self.padding_tokens), dim=0
391     ).unsqueeze(0)
392     # the length of padded_tokens is n + group_size - 1
393     # so creating it is O(n + group_size)
394     attention_len = padded_tokens.shape[1] # n + group_size - 1
395     if attention_len > self.max_input_len:
396         padded_tokens = padded_tokens[:, -self.max_input_len:]
397         # O(self.max_input_len) which is constant so O(1)
398         attention_len = self.max_input_len
399     attention_mask: LongTensor = ones(
400         [1, attention_len], dtype=long
401     )
402     # O(attention_len) so O(n + group_size)
403     if cuda.is_available():
404         padded_tokens = padded_tokens.cuda() # O(n + group_size)
405         attention_mask = attention_mask.cuda() # O(n + group_size)
406     else:
407         warn("CUDA is not available, using CPU")
408     return {
409         "input_ids": padded_tokens,
410         "attention_mask": attention_mask,
411     }
412
413 def get_logits_matrix(self, tokens: TokenIDS) -> Tensor:
414     """Given a sequence of tokens,
415     returns the logits matrix of shape (group_size, vocab_size)
416     where logits[i] is the logits vector of the i-th next token
417     complexity: O(n^2 + group_size^2) where n is the length of the tokens
418     notice that the logits are not divided by the temperature in this function."""
419     # define n as the number of tokens in tokens
420     model_kwargs = self.prepare_model_kwargs(tokens) # O(n + group_size)
421     with no_grad():
422         # The time complexity of causal language model's __call__ function
423         # is O(n^2) where n is the length of the inputs
424         outputs = self.model(
425             **model_kwargs
426         )
427         # the length of all the inputs is n + group_size - 1
428         # so the complexity of this line is O((n + group_size - 1)^2)
429         # which is O(n^2 + group_size^2 + group_size * n)
430         # we now that if a > b and a, b > 1 then a^2 > ab
431         # so the complexity is O(n^2 + group_size^2)
432     unscaled_relevant_logits: Tensor

```

```

433     unscaled_relevant_logits = outputs.logits[0, -self.group_size:, :self.vocab_size]
434     # The shape of unscaled_relevant_logits is (group_size, vocab_size)
435     # So the complexity of this line should be
436     # O(group_size) because we are coping group_size * vocab_size
437     # elements from one tensor to the another
438     return unscaled_relevant_logits
439
440 def get_prob_mat(
441     self, tokens: TokenIDS, generation_start: int
442 ) -> Tensor:
443     """Returns the probability matrix
444     as a list of lists of floats
445     Time complexity: O(n^2 + group_size^2)
446     where n is the number of tokens"""
447     unscaled_relevant_logits = self.get_logits_matrix(tokens)
448     # O(n^2 + group_size^2)
449     # unscaled_relevant_logits is a tensor of shape (group_size, vocab_size)
450     if not self.end_of_sentence_stop:
451         unscaled_relevant_logits[:, self.end_of_sentence_id] = -float('inf')
452         # setting a vector of size vocab_size
453         # so the complexity is O(group_size)
454         # setting the logits to -inf so the probability will be 0
455     penalized_logits = self.repetition_penalty_strategy(
456         unscaled_relevant_logits, tokens, generation_start
457     )
458     # O(group_size * n)
459     # where n is the number of tokens generated by the algorithm
460     prob_tensor = self.logits_to_probs(penalized_logits) # O(group_size)
461     # We are doing a softmax operator
462     # of group_size different vectors of size vocab_size
463     # The complexity of the softmax for each vector is
464     # O(1) because the size of the vector size is constant
465     # the complexity of this line is O(group_size)
466     # because we are doing group_size softmax operations
467     return prob_tensor
468
469 def logits_to_probs(self, penalized_logits: Tensor) -> Tensor:
470     """Gets the logits matrix and returns the probability matrix
471     Time complexity: O(group_size)"""
472     if self.use_softmax:
473         # if the generation type
474         # is not greedy then we need to apply softmax
475         # to the penalized logits
476         if self.temp != 1.0:
477             penalized_logits /= self.temp
478             # O(group_size * vocab_size)
479             # because we are dividing a matrix of size (group_size, vocab_size)
480         return Softmax(dim=1)(penalized_logits)
481         # O(group_size * vocab_size) so the complexity is O(group_size)
482     return penalized_logits
483
484 def __str__(self):
485     return f"GroupedGenerationUtils({self.as_dict()})"
486
487 def __repr__(self):
488     return str(self)
489
490 def as_dict(self):
491     return {attr_name: getattr(self, attr_name)
492             for attr_name in self.descriptive_attrs}
493
494 from __future__ import annotations
495
496 from abc import ABC, abstractmethod
497 from collections.abc import Callable
498 from typing import Optional, List, Union, Dict, Any
499
500 from torch import LongTensor
501 from transformers import (
502     AutoTokenizer,
503     AutoConfig,
504     PreTrainedTokenizer,
505 )
506 from transformers.tokenization_utils_base import TruncationStrategy
507
508 from .generation_type import GenerationType
509 from .generation_utils import GroupedGenerationUtils
510 from .postprocessor import PostProcessor
511 from .preprocessor import PreProcessor
512 from .repetition_penalty import RepetitionPenaltyStrategy, DEFAULT_REPETITION_PENALTY
513 from .completion_dict import CompletionDict
514 from .token_ids import TokenIDS
515
516 MAX_MODEL_INPUT_SIZE = 32768
517
518
519 def remove_nones(d: Dict[str, Any]) -> Dict[str, Any]:

```

```

520     """Returns a copy of a dictionary with all the not None values"""
521     return {key: d[key] for key in d.keys() if d[key] is not None}
522
523
524 def get_padding_id(tokenizer: PreTrainedTokenizer):
525     padding_id = tokenizer.pad_token_id
526     if not isinstance(padding_id, int):
527         padding_id = tokenizer.unk_token_id
528     if not isinstance(padding_id, int):
529         padding_id = tokenizer.mask_token_id
530     if not isinstance(padding_id, int):
531         raise RuntimeError(f"padding_id is {padding_id} and its type is {type(padding_id)}")
532     return padding_id
533
534
535 class GroupedGenerationPipeLine(Callable, ABC):
536     """An abstract base class for
537     A callable object that given a func_prompt
538     and length of wanted answer,
539     generates text
540     the text generator has a model,
541     and some parameters
542     (Defined in the subclasses)"""
543
544     framework: str = "pt"
545     descriptive_attrs = (
546         "model_name",
547         "generation_type",
548         "answer_length_multiplier",
549         "wrapped_model",
550     )
551
552     def __init__(
553         self,
554         model_name: str,
555         group_size: int,
556         temp: Optional[float] = None,
557         end_of_sentence_stop: Optional[bool] = None,
558         repetition_penalty_strategy: RepetitionPenaltyStrategy = DEFAULT_REPETITION_PENALTY,
559         answer_length_multiplier: float = 16,
560     ):
561         """Model name: the name of the model
562         used for loading from hugging face hub
563         group size: int
564         the number of tokens to be predicted at each model call
565         temp: float
566         temperature parameter for the softmax function
567         answer_length_multiplier: int
568             if the answer length is not given,
569             the maximum answer length is set to:
570             the length of the prompt * answer_length_multiplier
571         repetition_penalty_strategy: RepetitionPenaltyStrategy
572             The strategy for the repetition penalty
573         """
574         self.model_name: str = model_name
575         tokenizer = AutoTokenizer.from_pretrained(model_name)
576         end_of_sentence_id = tokenizer.eos_token_id
577         end_of_sentence_stop = end_of_sentence_stop and end_of_sentence_id is not None
578         max_input_len = tokenizer.model_max_length
579         max_len_is_huge = max_input_len > MAX_MODEL_INPUT_SIZE
580         if max_len_is_huge or max_input_len is None:
581             config = AutoConfig.from_pretrained(model_name)
582             max_input_len = config.max_position_embeddings
583             max_len_is_still_huge = max_input_len > MAX_MODEL_INPUT_SIZE
584             if max_len_is_still_huge or max_input_len is None:
585                 raise ValueError(
586                     "The maximum length of the model is too big"
587                 )
588         self.pre_processing_strategy: PreProcessor = PreProcessor(
589             tokenizer=tokenizer,
590             max_input_len=max_input_len,
591         )
592         self.post_processing_strategy: PostProcessor = PostProcessor(
593             tokenizer=tokenizer,
594         )
595         wrapped_model_kwargs: Dict[str, Any] = {
596             "model_name": model_name,
597             "group_size": group_size,
598             "max_input_len": max_input_len,
599             "end_of_sentence_id": end_of_sentence_id,
600             "end_of_sentence_stop": end_of_sentence_stop,
601             "repetition_penalty_strategy": repetition_penalty_strategy,
602             "padding_id": get_padding_id(tokenizer),
603             "temp": temp,
604             "use_softmax": self.generation_type.requires_softmax(),
605             "vocab_size": tokenizer.vocab_size,
606         },

```

```

606         ,
607         self.wrapped_model: GroupedGenerationUtils = GroupedGenerationUtils(**remove_nones(wrapped_model_kwargs))
608         self.answer_length_multiplier: float = answer_length_multiplier
609
610     @property
611     @abstractmethod
612     def generation_type(self) -> GenerationType:
613         """A method that chooses the generation type
614         Returns:
615             a GenerationType object"""
616         raise NotImplementedError
617
618     @abstractmethod
619     def _forward(
620         self,
621         tokenized_prompt: LongTensor,
622         num_new_tokens: Optional[int] = None,
623         num_return_sequences: int = 1,
624     ) -> List[TokenIDs]:
625         """A helper method for __call__ that generates the new tokens
626         Has a unique implementation for each subclass
627         Args:
628             tokenized_prompt: List[int]
629                 the tokenized prompt from the preprocess method
630             num_new_tokens: int - the number of new tokens to generate
631                 from the __call__ method
632         Returns:
633             the prompt + generated text as a list/tuple of ints"""
634         pass
635
636     def __call__(
637         self,
638         prompt_s: Union[str, List[str]],
639         max_new_tokens: Optional[int] = None,
640         return_tensors: bool = False,
641         return_text: bool = True,
642         return_full_text: bool = True,
643         clean_up_tokenization_spaces: bool = False,
644         prefix: str = "",
645         num_return_sequences: int = 1,
646         truncation: TruncationStrategy = TruncationStrategy.DO_NOT_TRUNCATE,
647         postfix: str = "",
648     ) -> CompletionDict | List[CompletionDict] | List[List[CompletionDict]]:
649         """The function that outside code should call to generate text
650         Args:
651             prompt_s: str or list of str - the prompt(s) to start the generation from
652                 (the text given by the user)
653                 if many prompts are given as a list,
654                 the function process each one independently and returns them as a list.
655                 (the same as calling [_call__(prompt, *args, **kwargs)
656                 for prompt in prompts]))
657             max_new_tokens: Optional[int] > 0
658                 - the number of tokens to generate
659                 if None, the function will generate tokens
660                 until one of them is the end of sentence token
661             return_tensors: bool - whether to return the generated token ids
662             return_text: bool - whether to return the generated string
663             return_full_text: bool - whether to return the full text
664                 (prompt + generated text)
665                 (if false, it will return only the generated text)
666             clean_up_tokenization_spaces: bool
667                 - whether to clean up tokenization spaces
668                 This parameter is forwarded to the decode function of the AutoTokenizer class
669             prefix (str, defaults to an empty string):
670                 Prefix added to prompt.
671             num_return_sequences (int, defaults to 1):
672                 The number of independently generated answers to return for each prompt.
673                 For GroupedSamplingPipeLine:
674                     each answer will be generated with different seed.
675                 For GroupedTreePipeLine:
676                     the num_return_sequences with the highest scores will be returned.
677             truncation: TruncationStrategy
678                 - whether to truncate the prompt
679             postfix: str
680                 - a postfix to add to the prompt
681         Returns:
682             Each result comes as a dictionary with the following keys:
683             - "generated_text"
684               (str, present when return_text=True)
685               -- The generated text.
686             - "generated_token_ids" (
687               torch.tensor, present when return_tensors=True)
688               -- The token
689               ids of the generated text.
690             """
691
692         if max_new_tokens is None and \

```



```

693         not self.wrapped_model.end_of_sentence_stop:
694         raise ValueError(
695             "max_new_tokens must be given if end_of_sentence_stop is False"
696         )
697     if isinstance(prompt_s, list):
698         return [self.__call__(
699             prompt_s=prompt,
700             max_new_tokens=max_new_tokens,
701             return_text=return_text,
702             return_tensors=return_tensors,
703             return_full_text=return_full_text,
704             clean_up_tokenization_spaces=clean_up_tokenization_spaces,
705             truncation=truncation,
706             postfix=postfix,
707         ) for prompt in prompt_s]
708     tokens: LongTensor
709     prefix_len: int
710     postfix_len: int
711     prompt_len: int
712
713     tokens, prefix_len, prompt_len, postfix_len = self.pre_processing_strategy(
714         prompt=prompt_s,
715         prefix=prefix,
716         truncation=truncation,
717         postfix=postfix
718     )
719     # O(len(prompt) + len(prefix) + len(postfix))
720
721     if max_new_tokens is None:
722         max_new_tokens = int(prompt_len * self.answer_length_multiplier)
723         # O(1)
724     tokenized_answers: List[TokenIDS]
725     tokenized_answers = self._forward(
726         tokens,
727         max_new_tokens,
728         num_return_sequences
729     )
730
731     if num_return_sequences > 1:
732         # O(sum(len(tokenized_answer) for tokenized_answer in tokenized_answers))
733         return [self.post_processing_strategy(
734             token_ids=tokenized_answer,
735             num_new_tokens=max_new_tokens,
736             prompt_len=prompt_len,
737             return_text=return_text,
738             return_tensors=return_tensors,
739             return_full_text=return_full_text,
740             clean_up_tokenization_spaces=clean_up_tokenization_spaces,
741             prefix_len=prefix_len,
742             postfix_len=postfix_len,
743         ) for tokenized_answer in tokenized_answers]
744     else:
745         # O(len(tokenized_answers[0]))
746         return self.post_processing_strategy(
747             token_ids=tokenized_answers[0],
748             num_new_tokens=max_new_tokens,
749             prompt_len=prompt_len,
750             return_text=return_text,
751             return_tensors=return_tensors,
752             return_full_text=return_full_text,
753             clean_up_tokenization_spaces=clean_up_tokenization_spaces,
754             prefix_len=prefix_len,
755             postfix_len=postfix_len,
756         )
757
758     def __repr__(self):
759         attrs_description = ", ".join(
760             f"{attr}={getattr(self, attr)}" for attr in self.descriptive_attrs
761         )
762         return f"{self.__class__.__name__}: " + attrs_description
763
764     def __str__(self):
765         return repr(self)
766
767     def as_dict(self) -> Dict[str, Any]:
768         """Returns a dictionary representation
769         of the generator
770         such that it can be saved and loaded
771         using the from_dict method"""
772         return {
773             key: getattr(self, key)
774             for key in self.descriptive_attrs
775         }
776
777     @classmethod
778     def from_dict(cls, my_dict: Dict[str, Any]):
779         """Creates an GroupedGenerationPipeLine from a dictionary

```

```

780         The dictionary should have the same format
781         as the dictionary returned by the as_dict method"""
782         if "generation_type" in my_dict.keys():
783             my_dict.pop("generation_type")
784         wrapped_model: GroupedGenerationUtils = my_dict.pop("wrapped_model")
785         wrapped_model_dict = wrapped_model.as_dict()
786         my_dict.update(wrapped_model_dict)
787         return cls(**my_dict)
788
789 import heapq
790 from collections.abc import Iterator
791 from random import seed
792 from typing import Callable, List, Dict, Optional, Any
793
794 from torch import Tensor, zeros, argmax, multinomial, manual_seed
795
796 from .generation_type import GenerationType
797 from .base_pipeline import GroupedGenerationPipeLine
798
799
800 class ChangingSeed(Iterator):
801     """Context manager for changing the seed of the random module.
802     How to use:
803     with ChangingSeed(first_seed, number_of_different_seeds) as changing_seed:
804         for _ in changing_seed:
805             # do something with random module"""
806     def __init__(self, default_seed: int, max_num_calls: int):
807         self.default_seed: int = default_seed
808         self.curr_seed: int = self.default_seed
809         self.max_num_calls: int = max_num_calls
810         self.curr_num_calls: int = 0
811
812     def __enter__(self):
813         self.curr_num_calls = 0
814         self.curr_seed = self.default_seed
815         return self
816
817     def __exit__(self, *args):
818         self.curr_seed = self.default_seed
819         seed(self.default_seed)
820         manual_seed(self.default_seed)
821
822     def __iter__(self):
823         self.curr_seed = self.default_seed
824         return self
825
826     def __next__(self):
827         self.curr_seed += 1
828         seed(self.curr_seed)
829         manual_seed(self.curr_seed)
830         self.curr_num_calls += 1
831         if self.curr_num_calls > self.max_num_calls:
832             raise StopIteration
833
834
835 class TokenProb:
836     """Class for storing the probability of a token and the token itself.
837     Used to store the probabilities of the next tokens in the sampling generator.
838     Is useful because it supports the < and > operators, which are used in the
839     heapq module
840     The < and > are the opposite of each other because the heapq module is only supporting minimum heaps
841     and I need a maximum heap"""
842     __slots__ = ['token_id', 'prob']
843
844     def __init__(self, token_id: int, prob: Tensor):
845         self.token_id: int = token_id
846         self.prob: Tensor = prob
847
848     def __lt__(self, other: "TokenProb"):
849         """Overrides the < operator
850         Comparison is done by the probability"""
851         return self.prob > other.prob
852
853     def __gt__(self, other: "TokenProb"):
854         """Overrides the > operator
855         Comparison is done by the probability"""
856         return self.prob < other.prob
857
858
859 class GroupedSamplingPipeLine(GroupedGenerationPipeLine):
860     """A GroupedGenerationPipeLine that generates text
861     using random sampling
862     with top-k or top-p filtering."""
863     default_seed: int = 0
864     seed(default_seed)
865     manual_seed(default_seed)

```

```

866 unique_attrs = "top_k", "top_p"
867
868 def __init__(self, top_k: Optional[int] = None,
869             top_p: Optional[float] = None, *args, **kwargs):
870     self.top_p: Optional[float] = top_p
871     self.top_k: Optional[int] = top_k
872     super().__init__(*args, **kwargs)
873
874 def __setattr__(self, key, value):
875     super().__setattr__(key, value)
876     if key == "default_seed":
877         seed(value)
878         manual_seed(value)
879
880 @property
881 def generation_type(self) -> GenerationType:
882     if self.top_k is None and self.top_p is None:
883         return GenerationType.RANDOM
884     if self.top_k == 1 or self.top_p == 0.0:
885         return GenerationType.GREEDY
886     if self.top_k is not None:
887         return GenerationType.TOP_K
888     if self.top_p is not None:
889         if self.top_p < 1.0:
890             return GenerationType.TOP_P
891         return GenerationType.RANDOM
892     raise RuntimeError("Uncovered case in generation_type property")
893
894 @property
895 def sampling_func(self) -> Callable[[Tensor], int]:
896     gen_type_to_filter_method: Dict[GenerationType, Callable[[Tensor], int]] = {
897         GenerationType.TOP_K: self.top_k_sampling,
898         GenerationType.TOP_P: self.top_p_sampling,
899         GenerationType.GREEDY: GroupedSamplingPipeLine.highest_prob_token,
900         GenerationType.RANDOM: GroupedSamplingPipeLine.unfiltered_sampling,
901     }
902     return gen_type_to_filter_method[self.generation_type]
903
904 @staticmethod
905 def unfiltered_sampling(prob_vec: Tensor) -> int:
906     """A sampling function that doesn't filter any tokens.
907     returns a random token id sampled from the probability vector"""
908     return multinomial(prob_vec, 1).item()
909
910 @staticmethod
911 def highest_prob_token(prob_vec: Tensor) -> int:
912     """Gets a probability vector of shape (vocab_size,)
913     returns the token id with the highest probability"""
914     return argmax(prob_vec).item()
915
916 def top_p_sampling(self, prob_vec: Tensor) -> int:
917     """Gets a probability vector of shape (vocab_size,)
918     computes a probability vector with the top p tokens
919     such that their sum in the original vector is <= self.top_p.
920     and samples from that vector.
921     If token with the highest probability
922     have a probability higher than top_p, it will be sampled"""
923     prob_sum: float = 0.0
924     converted_probs = [
925         TokenProb(i, prob) for i, prob in enumerate(prob_vec)
926     ]
927     heapq.heapify(converted_probs)
928     new_probs = zeros(prob_vec.shape, dtype=float)
929     while prob_sum < self.top_p and len(converted_probs) > 0:
930         curr_token_prob: TokenProb = heapq.heappop(converted_probs)
931         token_id = curr_token_prob.token_id
932         if curr_token_prob.prob <= 0.0:
933             break
934         if curr_token_prob.prob > 1:
935             raise ValueError(
936                 f"Probability of token {token_id} "
937                 f"in the vector {prob_vec} "
938                 f"is {curr_token_prob.prob} "
939                 f"which is higher than 1")
940         prob_sum += curr_token_prob.prob
941         new_probs[token_id] = curr_token_prob.prob
942     if prob_sum == 0.0:
943         return converted_probs[0].token_id
944     return GroupedSamplingPipeLine.unfiltered_sampling(new_probs)
945
946 def top_k_sampling(self, prob_vec: Tensor) -> int:
947     """Gets a token id: probability mapping
948     returns the TOP_K tokens
949     with the highest probability.
950     this is the bottleneck of the sampling generator."""
951     top_k_keys: List[int] = heapq.nlargest(
952         self.top_k,

```

```

953         range(prob_vec.shape[0]),
954         key=lambda x: prob_vec[x]
955     )
956     prob_sum = sum(prob_vec[token_id] for token_id in top_k_keys)
957     new_probs = zeros(prob_vec.shape, dtype=float)
958     for token_id in top_k_keys:
959         new_probs[token_id] = prob_vec[token_id] / prob_sum
960     return GroupedSamplingPipeLine.unfiltered_sampling(new_probs)
961
962 def generate_group(self, prob_mat: Tensor) -> List[int]:
963     """Generates a group of tokens
964     using the choice_function.
965     Complexity: O(group_size)"""
966     prob_mat.cpu()
967     # coping a tensor of size (group_size, vocab_size)
968     # so the complexity is O(group_size)
969     # (vocab_size is constant)
970     new_group: List[int] = [
971         self.sampling_func(prob_vec)
972         for prob_vec in prob_mat
973     ]
974     # the complexity of the loop is O(group_size)
975     # because self.sampling_func gets a tensor
976     # of constant size (vocab_size,)
977     # and therefore must be O(1) in complexity
978     # and the loop has group_size iterations.
979     del prob_mat
980     for i, token_id in enumerate(new_group):
981         if token_id == self.wrapped_model.end_of_sentence_id:
982             return new_group[:i + 1]
983         # return the group until the end of sentence token included
984         # the complexity of this line is O(group_size)
985         # because it is coping a list with maximum size of group_size
986     return new_group
987
988 def _forward(
989     self,
990     tokenized_prompt: Tensor,
991     num_new_tokens: Optional[int] = None,
992     num_return_sequences: int = 1,
993 ) -> List[List[int]]:
994     """Complexity:
995     O(num_return_sequences * (
996         ((n ^ 3) / group_size) +
997         ((n * 1 ^ 2) / group_size) +
998         group_size +
999         n)
1000     )
1001     where l is the number of tokens in the prompt
1002     and n is the number of new tokens to generate"""
1003     # let's define l = len(tokenized_prompt), n = num_new_tokens
1004     answers: List[List[int]] = []
1005     curr_token_list: List[int] = tokenized_prompt.tolist()
1006     # coping a tensor of size lso O(1)
1007     if num_new_tokens is None:
1008         raise RuntimeError("num_new_tokens is None")
1009     for _ in ChangingSeed(
1010         default_seed=self.default_seed,
1011         max_num_calls=num_return_sequences):
1012         # num_return_sequences iterations
1013         for _ in range(num_new_tokens // self.wrapped_model.group_size):
1014             # and each iteration is
1015             # O(n ^ 2 + 1 ^ 2 + group_size ^ 2 + group_size)
1016             # so the complexity of the loop is
1017             # O((n ^ 3) / group_size + (n * 1 ^ 2) / group_size + group_size + n)
1018             prob_mat: Tensor = self.wrapped_model.get_prob_mat(
1019                 curr_token_list, len(tokenized_prompt)
1020             )
1021             # complexity: O(group_size ^ 2 + len(curr_token_list) ^ 2)
1022             # len(curr_token_list) <= n + 1
1023             # so the complexity is
1024             # O(group_size ^ 2 + (n + 1) ^ 2) is equals to
1025             # O(n ^ 2 + n1 + 1 ^ 2 + group_size ^ 2)
1026             # but n1 <= max(n^2, 1^2) so the complexity
1027             # is O(n ^ 2 + 1 ^ 2 + group_size ^ 2)
1028             new_tokens = self.generate_group(prob_mat)
1029             # complexity: O(group_size)
1030             # len(curr_token_list) <= n + 1
1031             # so the complexity is O(group_size * (n + 1 + group_size))
1032             # len(new_tokens) = group_size
1033             if self.wrapped_model.end_of_sentence_id in new_tokens:
1034                 # the check is O(group_size)
1035                 end_of_sentence_index = new_tokens.index(
1036                     self.wrapped_model.end_of_sentence_id)
1037                 # O(group_size) because len(new_tokens) <= group_size
1038                 new_tokens = new_tokens[:end_of_sentence_index]
1039                 # O(group_size) because end_of_sentence_index < group_size

```

```

1039         # O(group_size) because end_of_sentence_index < group_size
1040         curr_token_list.extend(new_tokens)
1041         # O(group_size) because len(new_tokens) <= group_size
1042         answers.append(curr_token_list)
1043         # O(1)
1044     return answers
1045
1046 def __repr__(self):
1047     super_representation = super().__repr__()
1048     unique_representation = '/n'.join(
1049         f"{unique_attr_name}={getattr(self, unique_attr_name)}"
1050         for unique_attr_name in self.unique_attrs)
1051     return super_representation + unique_representation
1052
1053 def as_dict(self) -> Dict[str, Any]:
1054     super_dict = super(GroupedSamplingPipeLine, self).as_dict()
1055     super_dict.update(
1056         {unique_attr: self.__getattr__(unique_attr)
1057          for unique_attr in self.unique_attrs}
1058     )
1059     return super_dict

```