

الأمر نستخدم ثم بنا، الخاصة الرسومات لبطاقة بنجاح تشغيل برنامج بتثبيت نقوم عندم
الناتالية. الامعلومات تظهر أن يمكن تفاصيله، لعرض nvidia-smi

```
(base) lzw@lzw-MS-7E01:~$ nvidia-smi
```

```
17 2023 04:15:43
```

```

+-----+
| NVIDIA-SMI 535.86.10      :      535.86.10      CUDA: 12.2      |
+-----+-----+-----+-----+
| GPU                          Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC | |
|                          : /          |              |    GPU-Util    Compute M. |
|                          |              |              |              MIG M. |
+=====+=====+=====+=====+
|   0   NVIDIA GeForce RTX 4070                | 00000000:01:00.0      |              | N/A |
| 0%    34C    P8                          9W / 215W |    666MiB / 12282MiB |    15%      |    |
|                          |              |              |              N/A |
+-----+-----+-----+-----+

```

```

+-----+
| : |
| GPU  GI  CI      PID      G  /usr/lib/xorg/Xorg      GPU |
|      ID  ID                |                          |
+=====+=====+=====+=====+
|   0   N/A  N/A    1926    G  /usr/lib/xorg/Xorg      381MiB |
|   0   N/A  N/A    2065    G  /usr/bin/gnome-shell    120MiB |
|   0   N/A  N/A    3482    G  gnome-control-center    2MiB |
|   0   N/A  N/A    3803    G  ...irefox/2987/usr/lib/firefox/firefox 149MiB |
+-----+

```

22.04. ملحظات بعناية، هنا الرابط إلى الرجوع يرجى المرحلة. هذه إلى الوصول الصعب من الواقع، في
3.

تعلم

طورتها التي الكبيرة اللغة نماذج من مجموعة هي
مهام أداء على وقادرة الحجم حيث من فعالة لتكون مصممة النماذج هذه. باسم سابقا الممعرفة
استخدامها، ولكي فيية عن المزيد تعلم في ترغب كنت إذا جي. بشكل الطبعية اللغمة معالجة
اتباعها: يمكنك التي الخطوات بعض إليك

7. الإضافية الموارد

و `torch.cuda.is_available()` مثل منصات على الاصطناعي بالذكاء والمهتمين المطورين مجتمعات إلى انضم `torch.cuda.is_available()`.

فهذا، لتعميق الإنترنت على المتاحات التعليلية والبرامج التعليلية الدورات تابع

الخاصة. مشاريعك في فعال بشرك واستخدامها `torch.cuda.is_available()` تعلم من ستتمكن الخطوات، هذه باتباع

التالي، الخطأ سنواجه الأمر، تشغيل ومحاولة النماذج، تنزيل بعد

`torch.cuda.is_available()`: ذاكرة نفاذ. `torch.cuda.is_available()` 86.00 تخصيص محاولة تمت. `torch.cuda.is_available()` 64.81 بالفعّل؛ جيغابايت 9.70 تخصيص تم جيغابايت؛ 11.69 إجمالي سرعة 0: `torch.cuda.is_available()` بكتير أكبر المحجوزة الذاكرة كانت إذا. `torch.cuda.is_available()` بواسطة إجماليًا محجوزة جيغابايت 9.70 متاحة؛ التجزئة. لتجنب `max_split_size_mb` ضرب حاول المخصصة، الذاكرة من

13 حوالي يبلغ `torch.cuda.is_available()` نموذج وحجم فقط، جيغابايت 12 هي لدينا الرسومات بطاقة ذاكرة لأن نظرًا بنا. الخاصة الرسومات بطاقة باستخدام تشغيله نستطيع لافإننا جيغابايت،

`torch.cuda.is_available()`، `torch.cuda.is_available()` الآخر، المشروع استخدام نحاول `torch.cuda.is_available()`.

التالي: الخطأ نواجه

`torch.cuda.is_available()`: على العتور تم ولكن الجهاز، نفس على الموترات جميع تكون أن توقع `torch.cuda.is_available()` الذاكرة في `torch.cuda.is_available()` الوسيط من التحقق عنده `torch.cuda.is_available()` والأقل، على جهازين `torch.cuda.is_available()`

هذا. عن `torch.cuda.is_available()` نسأل ونحن

التالي. الكود إضافة إلى نحتاج جدًا. جميلًا إصلاحًا `torch.cuda.is_available()` لنا يوفر

```
input_ids = input_ids.to(model.device)
```

`model.device`، النموذج عليه يعمل الذي الجهاز إلى `input_ids` تحوّل تم

تشغيله. يمكننا أخيرًا،

```
(llama) lzw@lzw-MS-7E01:~/Projects/open_llama_3b$ python run.py
```

```
:
```

```
: .
```

```
:
```

```
: .
```



```

pad_id = self.tokenizer.pad_id
tokens = torch.full((bsz, total_len), pad_id, dtype=torch.long, device="cuda")
for k, t in enumerate(prompt_tokens):
    tokens[k, : len(t)] = torch.tensor(t, dtype=torch.long, device="cuda")
if logprobs:
    token_logprobs = torch.zeros_like(tokens, dtype=torch.float)

prev_pos = 0
eos_reached = torch.tensor([False] * bsz, device="cuda")
input_text_mask = tokens != pad_id
for cur_pos in range(min_prompt_len, total_len):
    logits = self.model.forward(tokens[:, prev_pos:cur_pos], prev_pos)
    if logprobs:
        token_logprobs[:, prev_pos + 1 : cur_pos + 1] = -F.cross_entropy(
            input=logits.transpose(1, 2),
            target=tokens[:, prev_pos + 1 : cur_pos + 1],
            reduction="none",
            ignore_index=pad_id,
        )
    if temperature > 0:
        probs = torch.softmax(logits[:, -1] / temperature, dim=-1)
        next_token = sample_top_p(probs, top_p)
    else:
        next_token = torch.argmax(logits[:, -1], dim=-1)

    next_token = next_token.reshape(-1)
    #
    next_token = torch.where(
        input_text_mask[:, cur_pos], tokens[:, cur_pos], next_token
    )
    tokens[:, cur_pos] = next_token
    eos_reached |= (~input_text_mask[:, cur_pos]) & (
        next_token == self.tokenizer.eos_id
    )
    prev_pos = cur_pos
    if all(eos_reached):
        break

```

```

if logprobs:
    token_logprobs = token_logprobs.tolist()
out_tokens, out_logprobs = [], []
for i, toks in enumerate(tokens.tolist()):
    #
    start = 0 if echo else len(prompt_tokens[i])
    toks = toks[start : len(prompt_tokens[i]) + max_gen_len]
    probs = None
    if logprobs:
        probs = token_logprobs[i][start : len(prompt_tokens[i]) + max_gen_len]
    #
    if self.tokenizer.eos_id in toks:
        eos_idx = toks.index(self.tokenizer.eos_id)
        toks = toks[:eos_idx]
        probs = probs[:eos_idx] if logprobs else None
    out_tokens.append(toks)
    out_logprobs.append(probs)
return (out_tokens, out_logprobs if logprobs else None)

```

المقدم: للكود شرحاً إليك بالطلب،

```

def generate(
    self,
    prompt_tokens: List[List[int]],
    max_gen_len: int,
    temperature: float = 0.6,
    top_p: float = 0.9,
    logprobs: bool = False,
    echo: bool = False,
) -> Tuple[List[List[int]], Optional[List[List[float]]]:

```

معاملات: عدة وتقبل generate تُسمى الوظيفة هذه

- prompt_tokens الرموز من تسلسلات على تحتوي قائمة هي
- max_gen_len إنشأه. يتم الذي للنص الأقصى ال طول هو
- temperature الإنشاء. عملية أثناء العشوائية في للتحكم تُستخدم مع لمة هي
- top_p العيونات. أخذ عملية أثناء التنوع في للتحكم تُستخدم و 1 0 بين احتمالية عتبة هو

الاحتمالية. اللوغاريتمات لتخزين token_logprobs موثر لإنشاء يتم True، تساوي logprobs كانت إذا

```
prev_pos = 0
eos_reached = torch.tensor([False] * bsz, device="cuda")
input_text_mask = tokens != pad_id
```

إلى: أعلاه الكود ترجمته تم

```
prev_pos = 0
eos_reached = torch.tensor([False] * bsz, device="cuda")
input_text_mask = tokens != pad_id
```

ترجمتها يتم لا والتي، `□□□□□□□□` بلغة ودوال متغيرات أسماء على يحتوي لأنه ترجمته يتم لم الكود ملاحظة: عادة.

نهاية رموز إلى الوصول تم إذا `□` ما `eos_reached` السابق `□□`، `□□` الموضوع `prev_pos` المتغيرات تهئية يتم بالحوش `□□`. ملؤه التي لم التي `□□` الموضوع `input_text_mask` والتسلسل `□□`.

الكلي. الطول إلى وصولاً المطالبة لطول الأدنى الحد من بدءاً الرموز تولد التالية الحلقة

```
for cur_pos in range(min_prompt_len, total_len):
    logits = self.model.forward(tokens[:, prev_pos:cur_pos], prev_pos)
    if logprobs:
        token_logprobs[:, prev_pos + 1 : cur_pos + 1] = -F.cross_entropy(
            input=logits.transpose(1, 2),
            target=tokens[:, prev_pos + 1 : cur_pos + 1],
            reduction="none",
            ignore_index=pad_id,
        )
    if temperature > 0:
        probs = torch.softmax(logits[:, -1] / temperature, dim=-1)
        next_token = sample_top_p(probs, top_p)
    else:
        next_token = torch.argmax(logits[:, -1], dim=-1)

next_token = next_token.reshape(-1)
next_token = torch.where(
    input_text_mask[:, cur_pos], tokens[:, cur_pos], next_token
)
```

```

tokens[:, cur_pos] = next_token
eos_reached |= (~input_text_mask[:, cur_pos]) & (
    next_token == self.tokenizer.eos_id
)
prev_pos = cur_pos
if all(eos_reached):
    break

```

إلى: أعلاه الكود تترجمه تمت

```

next_token = next_token.reshape(-1)
next_token = torch.where(
    input_text_mask[:, cur_pos], tokens[:, cur_pos], next_token
)
tokens[:, cur_pos] = next_token
eos_reached |= (~input_text_mask[:, cur_pos]) & (
    next_token == self.tokenizer.eos_id
)
prev_pos = cur_pos
if all(eos_reached):
    break

```

يتم لا والتي `tokens` البرمجة لغة في محددة ودوال أسماء على يحتوي لأنه تترجمته يتم لم الكود ملاحظة: عادةً، تترجمتها

فإنه صحيحًا، `logprobs` كان إذا التالي. للرمز `logits` احتمالي توزيع بإنشاء النموذج يقوم الحلقه، هذه داخل يستخدم، `temperature` على اعتماده. المقتطعة. الانتروبيا باستخدام الاحتمالية اللوغاريتمات يحسب التسلسل. نهائية رموز ويفحص `tokens` موترت بتحديث يقوم ثم التالي. الرمز لتحديد `tokens` أو `tokens`.

```

if logprobs:
    token_logprobs = token_logprobs.tolist()

```

قائمة إلى اللوغاريتمات احتمالات موترت تحوي لي يتم، `True` تساوي `logprobs` كانت إذا

```

out_tokens, out_logprobs = [], []
for i, toks in enumerate(tokens.tolist()):
    ...

```

بها. المربطة اللوغاريتمية والاحتمالات الناتج الرموز تسلسل توليدي يتم إدخال، عينة لكل الحلقه، هذه في

```
return (out_tokens, out_logprobs if logprobs else None)
```

ذلك. طلب تم إذا له المقابلة اللوغاريتمية والاحتمالات المُنشأة الرموز تسلسلات الدالة شعبي أخيرًا،