openai / tiktoken

<> Code    Issues 35    Pull requests 21    Actions    Projects    Security    Insights

⚠ You only have a single verified email address. We recommend verifying at least one more email address to ensure you can recover your account if you lose access to your primary email.    Email settings    ✕

**tiktoken** / tiktoken_ext / openai_public.py  ⧉

hauntsaninja  Sync codebase  ✕                                    9d01e56 · 2 months ago    🕐 History

```
1    from tiktoken.load import data_gym_to_mergeable_bpe_ranks, load_tiktoken_bpe
```

**tiktoken** / tiktoken_ext / openai_public.py                                      ↑ Top

Code    Blame    127 lines (112 loc) · 4.53 KB                    Raw ⧉ ⬇    ✏ ⌄    <>

```python
 6    FIM_SUFFIX = "<|fim_suffix|>"
 7    ENDOFPROMPT = "<|endofprompt|>"
 8
 9
10    def gpt2():
11        mergeable_ranks = data_gym_to_mergeable_bpe_ranks(
12            vocab_bpe_file="https://openaipublic.blob.core.windows.net/gpt-2/encodings/main/vocab.bpe",
13            encoder_json_file="https://openaipublic.blob.core.windows.net/gpt-2/encodings/main/encoder.json",
14            vocab_bpe_hash="1ce1664773c50f3e0cc8842619a93edc4624525b728b188a9e0be33b7726adc5",
15            encoder_json_hash="196139668be63f3b5d6574427317ae82f612a97c5d1cdaf36ed2256dbf636783",
16        )
17        return {
18            "name": "gpt2",
19            "explicit_n_vocab": 50257,
20            # The pattern in the original GPT-2 release is:
21            # r"""'s|'t|'re|'ve|'m|'ll|'d| ?[\p{L}]+| ?[\p{N}]+| ?[^\s\p{L}\p{N}]+|\s+(?!\S)|\s+"""
22            # This is equivalent, but executes faster:
23            "pat_str": r"""(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+|\s+(?!\S)|\s+""",
24            "mergeable_ranks": mergeable_ranks,
25            "special_tokens": {ENDOFTEXT: 50256},
26        }
27
28
29    def r50k_base():
30        mergeable_ranks = load_tiktoken_bpe(
31            "https://openaipublic.blob.core.windows.net/encodings/r50k_base.tiktoken",
32            expected_hash="306cd27f03c1a714eca7108e03d66b7dc042abe8c258b44c199a7ed9838dd930",
33        )
34        return {
35            "name": "r50k_base",
36            "explicit_n_vocab": 50257,
37            "pat_str": r"""(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+|\s+(?!\S)|\s+""",
38            "mergeable_ranks": mergeable_ranks,
39            "special_tokens": {ENDOFTEXT: 50256},
40        }
41
42
43    def p50k_base():
44        mergeable_ranks = load_tiktoken_bpe(
45            "https://openaipublic.blob.core.windows.net/encodings/p50k_base.tiktoken",
46            expected_hash="94b5ca7dff4d00767bc256fdd1b27e5b17361d7b8a5f968547f9f23eb70d2069",
47        )
48        return {
49            "name": "p50k_base",
50            "explicit_n_vocab": 50281,
51            "pat_str": r"""(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+|\s+(?!\S)|\s+""",
52            "mergeable_ranks": mergeable_ranks,
53            "special_tokens": {ENDOFTEXT: 50256},
54        }
55
56
57    def p50k_edit():
58        mergeable_ranks = load_tiktoken_bpe(
59            "https://openaipublic.blob.core.windows.net/encodings/p50k_base.tiktoken",
```

**Symbols**    ✕

Find definitions and references for functions and other symbols in this file by clicking a symbol below or in the code.

⧨ Filter symbols                                    r

const  ENDOFTEXT

const  FIM_PREFIX

const  FIM_MIDDLE

const  FIM_SUFFIX

const  ENDOFPROMPT

func  gpt2

func  r50k_base

func  p50k_base

func  p50k_edit

func  cl100k_base

func  o200k_base

const  ENCODING_CONSTRUCT...

```python
60            expected_hash="94b5ca7dff4d00767bc256fdd1b27e5b17361d7b8a5f968547f9f23eb70d2069",
61        )
62        special_tokens = {ENDOFTEXT: 50256, FIM_PREFIX: 50281, FIM_MIDDLE: 50282, FIM_SUFFIX: 50283}
63        return {
64            "name": "p50k_edit",
65            "pat_str": r"""'(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+|\s+(?!\S)|\s+""",
66            "mergeable_ranks": mergeable_ranks,
67            "special_tokens": special_tokens,
68        }
69
70
71  ⌄ def cl100k_base():
72        mergeable_ranks = load_tiktoken_bpe(
73            "https://openaipublic.blob.core.windows.net/encodings/cl100k_base.tiktoken",
74            expected_hash="223921b76ee99bde995b7ff738513eef100fb51d18c93597a113bcffe865b2a7",
75        )
76        special_tokens = {
77            ENDOFTEXT: 100257,
78            FIM_PREFIX: 100258,
79            FIM_MIDDLE: 100259,
80            FIM_SUFFIX: 100260,
81            ENDOFPROMPT: 100276,
82        }
83        return {
84            "name": "cl100k_base",
85            "pat_str": r"""'(?i:[sdmt]|ll|ve|re)|[^\r\n\p{L}\p{N}]?+\p{L}+|\p{N}{1,3}| ?[^\s\p{L}\p{N}]++[\r\n]*|\s*[\r\n]|\s+(?!\S)|\s+""",
86            "mergeable_ranks": mergeable_ranks,
87            "special_tokens": special_tokens,
88        }
89
90
91  ⌄ def o200k_base():
92        mergeable_ranks = load_tiktoken_bpe(
93            "https://openaipublic.blob.core.windows.net/encodings/o200k_base.tiktoken",
94            expected_hash="446a9538cb6c348e3516120d7c08b09f57c36495e2acfffe59a5bf8b0cfb1a2d",
95        )
96        special_tokens = {
97            ENDOFTEXT: 199999,
98            ENDOFPROMPT: 200018,
99        }
100       # This regex could be made more efficient
101       pat_str = "|".join(
102           [
103               r"""[^\r\n\p{L}\p{N}]?[\p{Lu}\p{Lt}\p{Lm}\p{Lo}\p{M}]*[\p{Ll}\p{Lm}\p{Lo}\p{M}]+(?i:'s|'t|'re|'ve|'m|'ll|'d)?""",
104               r"""[^\r\n\p{L}\p{N}]?[\p{Lu}\p{Lt}\p{Lm}\p{Lo}\p{M}]+[\p{Ll}\p{Lm}\p{Lo}\p{M}]*(?i:'s|'t|'re|'ve|'m|'ll|'d)?""",
105               r"""\p{N}{1,3}""",
106               r""" ?[^\s\p{L}\p{N}]+[\r\n/]*""",
107               r"""\s*[\r\n]+""",
108               r"""\s+(?!\S)""",
109               r"""\s+""",
110           ]
111       )
112       return {
113           "name": "o200k_base",
114           "pat_str": pat_str,
115           "mergeable_ranks": mergeable_ranks,
116           "special_tokens": special_tokens,
117       }
118
119
120 ⌄ ENCODING_CONSTRUCTORS = {
121       "gpt2": gpt2,
122       "r50k_base": r50k_base,
123       "p50k_base": p50k_base,
124       "p50k_edit": p50k_edit,
125       "cl100k_base": cl100k_base,
126       "o200k_base": o200k_base,
127   }
```