

README

本模型使用时需输入小分子的SMILES结构式，获取其所预测的RNA适配体输出

由于核糖开关适配体数据量有限，故将DNA原件全部转为RNA进行探索预测

模型采用类Transformer架构，将小分子与aptamer序列视为翻译任务

详细代码见 **RNA_smiles (1).ipynb** 文件

调用模型的代码如下，可进一步封装为api

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from tokenizers import Tokenizer
5
6 # -----
7 # 1. 重新定义模型（确保嵌入层维度=256）
8 # -----
9 class Mol2Aptamer(nn.Module):
10     def __init__(self, smiles_vocab_size, rna_vocab_size, d_model=256,
11                  nhead=8, num_encoder_layers=2, num_decoder_layers=3):
12         super().__init__()
13         self.d_model = d_model # 显式保存模型维度，便于验证
14
15         # Encoder: 2层 TransformerEncoder (d_model=256)
16         self.encoder = nn.TransformerEncoder(
17             encoder_layer=nn.TransformerEncoderLayer(
18                 d_model=d_model,           # 注意力层期望的维度=256
19                 nhead=nhead,
20                 dim_feedforward=2048,
21                 batch_first=True
22             ),
23             num_layers=num_encoder_layers
24         )
25
26         # Decoder: 3层 TransformerDecoder (d_model=256)
27         self.decoder = nn.TransformerDecoder(
28             decoder_layer=nn.TransformerDecoderLayer(
29                 d_model=d_model,
30                 nhead=nhead,
31                 dim_feedforward=2048,
32                 batch_first=True
33             ),
34             num_layers=num_decoder_layers
```

```
34     )
35
36     # 嵌入层：明确设置输出维度=d_model=256
37     self.smiles_embedding = nn.Embedding(smiles_vocab_size, d_model)
38 ) # 输出 (seq_len, 256)
39     self.rna_embedding = nn.Embedding(rna_vocab_size, d_model)
40     self.pos_embedding = nn.Embedding(512, d_model) # 位置编码维度
41 也=256
42
43     # 输出层
44     self.fc_out = nn.Linear(d_model, rna_vocab_size)
45
46     def forward(self, smiles_ids, rna_inp):
47         batch_size, seq_len_smi = smiles_ids.shape
48         batch_size, seq_len_rna = rna_inp.shape
49         device = smiles_ids.device
50
51         # -----
52         # 关键：验证嵌入层输出维度（确保=256）
53         # -----
54         smiles_emb = self.smiles_embedding(smiles_ids)
55         assert smiles_emb.shape[-1] == self.d_model, \
56             f"嵌入层输出维度错误：期望 {self.d_model}，实际 {smiles_emb.s
57 hape[-1]}"
58
59         # 叠加位置编码（维度需与嵌入层一致）
60         smiles_pos = torch.arange(seq_len_smi, device=device).unsqueeze
61         (0).repeat(batch_size, 1)
62         smiles_emb += self.pos_embedding(smiles_pos)
63         assert smiles_emb.shape[-1] == self.d_model, \
64             f"位置编码后维度错误：期望 {self.d_model}，实际 {smiles_emb.s
65 hape[-1]}"
66
67         # Encoder前向（输入维度=256，匹配注意力层期望）
68         memory = self.encoder(smiles_emb)
69
70         # RNA嵌入与位置编码（同样验证维度）
71         rna_emb = self.rna_embedding(rna_inp)
72         assert rna_emb.shape[-1] == self.d_model, \
73             f"RNA嵌入层输出维度错误：期望 {self.d_model}，实际 {rna_emb.s
74 hape[-1]}"
75         rna_pos = torch.arange(seq_len_rna, device=device).unsqueeze(0)
76         .repeat(batch_size, 1)
77         rna_emb += self.pos_embedding(rna_pos)
```

```
71
72     # Decoder前向
73     tgt_mask = nn.Transformer.generate_square_subsequent_mask(seq_length_rna, device=device)
74     decoder_out = self.decoder(tgt=rna_emb, memory=memory, tgt_mask=tgt_mask)
75
76     # 输出层
77     logits = self.fc_out(decoder_out)
78     return logits
79
80 # -----
81 # 2. 加载分词器 + 初始化模型（显式确认嵌入层维度）
82 #
83 smiles_tokenizer = Tokenizer.from_file("/root/autodl-tmp/rna/smiles_tokenizer.json")
84 rna_tokenizer = Tokenizer.from_file("/root/autodl-tmp/rna/rna_tokenizer.json")
85
86 # 词汇表大小
87 smiles_vocab_size = len(smiles_tokenizer.get_vocab())
88 rna_vocab_size = len(rna_tokenizer.get_vocab())
89
90 # 设备与模型初始化 (d_model=256, 嵌入层输出维度=256)
91 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
92 model = Mol2Aptamer(
93     smiles_vocab_size=smiles_vocab_size,
94     rna_vocab_size=rna_vocab_size,
95     d_model=256,           # 强制嵌入层输出维度=256
96     nhead=8,
97     num_encoder_layers=2,
98     num_decoder_layers=3
99 ).to(device)
100
101 #
102 # 3. 修复权重加载：确保嵌入层权重被正确加载
103 #
104 # 加载原始权重
105 state_dict = torch.load("/root/autodl-tmp/model_epoch_59.pth", map_location=device)
106
107 # 查看所有权重键，确认嵌入层键是否存在（如 "smiles_embedding.weight"）
108 print("原始权重中的嵌入层相关键: ")
109 for key in state_dict.keys():
```

```
110     if "embedding" in key:
111         print(f"- {key}: 形状 {state_dict[key].shape}")
112 print("-" * 50)
113
114 # 过滤权重: 保留所有与当前模型匹配的键 (包括嵌入层、位置编码、输出层)
115 filtered_state_dict = {}
116 model_keys = set(model.state_dict().keys()) # 当前模型需要的键
117 for key, value in state_dict.items():
118     if key in model_keys:
119         # 额外验证嵌入层权重维度是否正确 (如 smiles_embedding.weight 形状
120         # 应为 (vocab_size, 256))
121         if "embedding.weight" in key:
122             assert value.shape[-1] == model.d_model, \
123                 f"{key} 维度错误: 期望 {model.d_model}, 实际 {value.shap
124             e[-1]}"
125             filtered_state_dict[key] = value
126     else:
127         print(f"跳过冗余键: {key}")
128
129 # 加载过滤后的权重
130 model.load_state_dict(filtered_state_dict, strict=False) # strict=False
131 # 允许模型有未加载的键 (如无)
132 model.eval()
133 print("权重加载成功!")
134
135 # -----
136 # 4. 修复生成函数: 确保SMILES编码后维度正确
137 # -----
138
139 def generate_aptamers(
140     model, smiles, smiles_tokenizer, rna_tokenizer,
141     max_len=80, num_return=5,
142     strategy="topk", top_k=10, top_p=0.9, temperature=0.8,
143     device=None
144 ):
145     if device is None:
146         device = next(model.parameters()).device
147
148     model.eval()
149     model.to(device)
150
151     # 验证特殊token
152     required_tokens = ["<pad>", "<bos>", "<eos>", "<unk>"]
153     for token in required_tokens:
154         token_id = rna_tokenizer.token_to_id(token)
```

```
151     if token_id is None:
152         raise ValueError(f"RNA tokenizer missing required token: {token}")
153
154     bos_id = rna_tokenizer.token_to_id("<bos>")
155     eos_id = rna_tokenizer.token_to_id("<eos>")
156     pad_id = rna_tokenizer.token_to_id("<pad>")
157     unk_id = rna_tokenizer.token_to_id("<unk>")
158     smiles_pad_id = smiles_tokenizer.token_to_id("<pad>")
159     if smiles_pad_id is None:
160         raise ValueError("SMILES tokenizer missing <pad> token")
161
162     # -----
163     # 修复SMILES编码: 确保输入到嵌入层的张量格式正确
164     # -----
165     try:
166         smi_encoded = smiles_tokenizer.encode(smiles)
167         max_smi_len = 128 # 与训练时一致
168         # 补全/截断SMILES到max_smi_len
169         smi_ids = smi_encoded.ids[:max_smi_len]
170         smi_ids += [smiles_pad_id] * (max_smi_len - len(smi_ids))
171         # 转换为张量: (batch_size=1, seq_len=128)
172         smi_ids = torch.tensor(smi_ids, dtype=torch.long).unsqueeze(0).
to(device)
173         assert smi_ids.shape == (1, max_smi_len), \
174             f"SMILES张量形状错误: 期望 (1, {max_smi_len}), 实际 {smi_ids.
shape}"
175     except Exception as e:
176         raise ValueError(f"Failed to encode SMILES: {str(e)}")
177
178     results = []
179     with torch.no_grad():
180         # -----
181         # 验证Encoder输入维度(嵌入层输出应为256)
182         # -----
183         smiles_emb = model.smiles_embedding(smi_ids)
184         assert smiles_emb.shape == (1, max_smi_len, model.d_model), \
185             f"SMILES嵌入后形状错误: 期望 (1, {max_smi_len}, {model.d_mod
el}), 实际 {smiles_emb.shape}"
186
187         # 预计算Encoder输出
188         memory = model.encoder(smiles_emb)
189
190         for _ in range(num_return):
```

```
191     generated = [bos_id]
192     has_unk = False
193
194     for _ in range(max_len - 1):
195         # RNA输入张量: (1, current_len)
196         rna_inp = torch.tensor([generated], dtype=torch.long).t
197         o(device)
198
199         # 模型前向传播
200         logits = model(smi_ids, rna_inp)[:, -1, :] # (1, rna_v
201         ocab_size)
202
203         # 温度调节
204         logits = logits / temperature
205
206         # 采样策略
207         if strategy == "greedy":
208             next_id = torch.argmax(logits, dim=-1).item()
209         elif strategy == "topk":
210             topk_probs, topk_ids = torch.topk(logits, k=top_k)
211             topk_probs = F.softmax(topk_probs, dim=-1)
212             idx = torch.multinomial(topk_probs, 1).item()
213             next_id = topk_ids[0, idx].item()
214         elif strategy == "topp":
215             sorted_logits, sorted_ids = torch.sort(logits, desc
216             ending=True)
217             sorted_probs = F.softmax(sorted_logits, dim=-1)
218             cum_probs = torch.cumsum(sorted_probs, dim=-1)
219             cutoff = max(1, torch.sum(cum_probs <= top_p).item(
220             ))
221             filtered_probs = sorted_probs[:, :cutoff]
222             filtered_ids = sorted_ids[:, :cutoff]
223             idx = torch.multinomial(filtered_probs, 1).item()
224             next_id = filtered_ids[0, idx].item()
225         else:
226             raise ValueError(f"Unknown strategy: {strategy}")
227
228         # 终止条件
229         if next_id == unk_id:
230             has_unk = True
231             if next_id == eos_id:
232                 break
233             generated.append(next_id)
```

```
231         # 后处理
232         filtered_ids = [id for id in generated if id not in [bos_id
233 , eos_id, pad_id, unk_id]]
234         try:
235             seq = rna_tokenizer.decode(filtered_ids, skip_special_t
236 okens=True)
237             if seq and not has_unk and seq not in results:
238                 results.append(seq)
239         except Exception as e:
240             print(f"Warning: Failed to decode sequence: {str(e)}")
241
242     # 补全候选数量
243     while len(results) < num_return:
244         results.append(results[-1] if results else "")
245     return results[:num_return]
246
247     # -----
248     # 5. 测试生成（成功运行）
249     # -----
250
251     if __name__ == "__main__":
252         # 输入SMILES（苯酚）
253         smiles = "C1=CC=C(C=C1)O"
254
255         # 生成前先验证模型嵌入层维度
256         print(f"模型嵌入层输出维度: {model.d_model}")
257         print(f"SMILES嵌入层权重形状: {model.smiles_embedding.weight.shape}")
258     )
259         print(f"RNA嵌入层权重形状: {model.rna_embedding.weight.shape}")
260
261     # 生成Aptamer
262     candidates = generate_aptamers(
263         model=model,
264         smiles=smiles,
265         smiles_tokenizer=smiles_tokenizer,
266         rna_tokenizer=rna_tokenizer,
267         max_len=80,
268         num_return=5,
269         strategy="topk",
270         top_k=10,
271         temperature=0.8,
272         device=device
273     )
274
275     # 打印结果
```

```
272     print("\n候选Aptamer序列: ")
273     for i, seq in enumerate(candidates, 1):
274         print(f"{i}. {seq}")
```

输出:

```
1 原始权重中的嵌入层相关键:
2 - smiles_embedding.weight: 形状 torch.Size([188, 256])
3 - rna_embedding.weight: 形状 torch.Size([100, 256])
4 - pos_embedding.weight: 形状 torch.Size([512, 256])
5 -----
6 跳过冗余键: decoder_layer.self_attn.in_proj_weight
7 跳过冗余键: decoder_layer.self_attn.in_proj_bias
8 跳过冗余键: decoder_layer.self_attn.out_proj.weight
9 跳过冗余键: decoder_layer.self_attn.out_proj.bias
10 跳过冗余键: decoder_layer.multihead_attn.in_proj_weight
11 跳过冗余键: decoder_layer.multihead_attn.in_proj_bias
12 跳过冗余键: decoder_layer.multihead_attn.out_proj.weight
13 跳过冗余键: decoder_layer.multihead_attn.out_proj.bias
14 跳过冗余键: decoder_layer.linear1.weight
15 跳过冗余键: decoder_layer.linear1.bias
16 跳过冗余键: decoder_layer.linear2.weight
17 跳过冗余键: decoder_layer.linear2.bias
18 跳过冗余键: decoder_layer.norm1.weight
19 跳过冗余键: decoder_layer.norm1.bias
20 跳过冗余键: decoder_layer.norm2.weight
21 跳过冗余键: decoder_layer.norm2.bias
22 跳过冗余键: decoder_layer.norm3.weight
23 跳过冗余键: decoder_layer.norm3.bias
24 权重加载成功!
25 模型嵌入层输出维度: 256
26 SMILES嵌入层权重形状: torch.Size([188, 256])
27 RNA嵌入层权重形状: torch.Size([100, 256])
28
29 候选Aptamer序列:
30 1. GGGU AAU AC GCA GA CGU GAGG GAU GCA CU CGG AU GCGU AGG GG GUU GAU C
   A
31 2. GAA GCA GCA CA GA GGU CAU CUU GAU CU CGG CUU GAU AGG GU CGU CC GUAA
   CU CC C
32 3. GCC GG AA CU A CUU CA CGU AC GACU GU CACA CCU GAGG GGU CGU AA CAA G
   U GGU AU GCGU
33 4. GCC GG GGU GG AA AC GAU CU CAU CC GG GUU GU
34 5. GAA GCA GCA CA GA GGU CA GAU GCA CU CGG ACC CC AUU CU CC UU CC AUCC
```

35 | CU CAU CC GUCC ACC CU AU

可进一步更改生成参数调控输出:

```
1 # 自定义生成参数（适合调整多样性、长度等）
2 aptamers = generate_aptamers(
3     model=model,
4     smiles="CC(=O)OC1=CC=CC=C1C(=O)O",
5     smiles_tokenizer=smiles_tokenizer,
6     rna_tokenizer=rna_tokenizer,
7     max_len=100,                      # 最长生成100个token
8     num_return=3,                     # 返回3个候选序列
9     strategy="topp",                  # 使用top-p采样（核采样）
10    top_p=0.95,                      # 累积概率阈值0.95
11    temperature=0.8,                 # 降低多样性（值越小越确定）
12 )
13
14 print("生成的Aptamer序列: ")
15 for i, seq in enumerate(aptamers):
16     print(f"候选{i+1}: {seq}")
```

之后调用RNAfold包计算适配体序列稳定性并排序获得最终输出，后续可接入3D可视化工具，并进行分子动力学验证

```
1 #过滤函数（RNAfold计算ΔG）
2 import RNA
3
4 def filter_by_rnafold(sequences, min_length=20, max_length=80, max_homo
polymer=6, max_candidates=5):
5     """
6         sequences: list of str
7             返回过滤+排序后的序列
8     """
9     results = []
10    for seq in sequences:
11        # 长度限制
12        if len(seq) < min_length or len(seq) > max_length:
13            continue
14        # 去掉长同聚核苷酸 (AAAAAAA)
15        if any(base*max_homopolymer in seq for base in "ACGU"):
16            continue
17        # 用 RNAfold 预测 ΔG
18        structure, mfe = RNA.fold(seq)
```

```

19     results.append((seq, mfe))
20
21     # 按 ΔG 从低到高排序 (越低越稳定)
22     results = sorted(results, key=lambda x: x[1])
23     return results[:max_candidates]
24
25 def generate_and_filter(
26     model, smiles, smiles_tokenizer, rna_tokenizer,
27     num_generate=50, return_top=5,
28     strategy="topk", top_k=10, top_p=0.9, temperature=0.8
29 ):
30     # Step1: 生成候选
31     candidates = generate_aptamers(
32         model, smiles, smiles_tokenizer, rna_tokenizer,
33         max_len=80, num_return=num_generate,
34         strategy=strategy, top_k=top_k, top_p=top_p, temperature=temperature,
35         device=device
36     )
37
38     # Step2: 过滤 & 打分
39     filtered = filter_by_rnafold(candidates, max_candidates=return_top)
40
41     return filtered
42
43 smiles = "C1=CC=C(C=C1)O" # phenol
44
45 top_candidates = generate_and_filter(
46     model, smiles, smiles_tokenizer, rna_tokenizer,
47     num_generate=100, return_top=5,
48     strategy="topp", top_p=0.9, temperature=0.7
49 )
50
51 print("Top 适配体候选 (按ΔG排序) : ")
52 for seq, mfe in top_candidates:
53     print(f"{seq} ΔG={mfe:.2f}")

```

输出

```

1 Top 适配体候选 (按ΔG排序) :
2   G CGA GAGG AGU GGU GG GGU CA GAU GCA CU CGG ACC CC AUU CU CC C   ΔG=-4.
3   10
4   G CA AUGG CC ACC CC GG GGU GG GCGC GAA AGU GGU   ΔG=-3.10

```

4 G CU CU CGG GA CGA CC CA CGU CC GGGU GG CUU GAU AGG GG GGU GGU CC AUCC
 CU CC ΔG=-2.20

5 G CC GGU ACA CA GG AGG CU GGU GCGC GGU GAA GU GCC GA GU CGU AA C ΔG=-
 2.10

6 G CGU AC GACU CA GG GCC A GAGG GAU CGG GU GGU CGU GGU CCU GAU GCA AUCU
 CU CC C ΔG=-1.80