

NMT Experiment Report: Chinese-to-English Translation based on RNN and Transformer

1. Project Overview

The core objective of this project is to implement a Neural Machine Translation (NMT) system for Chinese-to-English translation and to perform an in-depth comparison between **RNN (GRU)** and **Transformer** architectures in terms of translation performance, training efficiency, and generation quality.

We built both models from scratch and trained them on a dataset of 100k sentence pairs. Furthermore, we conducted several ablation studies to investigate the impact of Attention mechanisms, training strategies (Teacher Forcing), and normalization methods on model performance.

2. Model Architectures & Implementation Details

We implemented all models using PyTorch, focusing on modular design to facilitate ablation studies.

2.1 RNN-based NMT Model

Our RNN model follows the classic **Seq2Seq** architecture with Attention.

- **Encoder:**
 - **Embedding Layer:** Converts source token indices into dense vectors of dimension d_{model} (256).
 - **GRU Layers:** We used a 2-layer unidirectional GRU. At each time step t , it processes the input embedding and the previous hidden state h_{t-1} to produce the current hidden state h_t .
 - **Dropout:** Applied to embeddings and between GRU layers (rate=0.3) to prevent overfitting.
- **Attention Mechanism:**
 - To solve the bottleneck problem where the encoder must compress the entire sentence into a fixed-size vector, we implemented Attention.
 - At each decoder step, we calculate a context vector c_i as a weighted sum of encoder outputs h_j : $c_i = \sum \alpha_{ij} h_j$.
 - We implemented three scoring functions for α_{ij} :
 1. **Dot-product:** $score(s_i, h_j) = s_i^T h_j$. Efficient but requires $d_{enc} = d_{dec}$.
 2. **General:** $score(s_i, h_j) = s_i^T W_a h_j$. Introduces a learnable weight matrix W_a , allowing for different dimensions.
 3. **Concat (Additive):** $score(s_i, h_j) = v_a^T \tanh(W_a[s_i; h_j])$. Uses a small MLP; theoretically most expressive but computationally heavier.
- **Decoder:**
 - **Input:** Concatenation of the target token embedding and the calculated context vector c_i .
 - **GRU:** 2-layer unidirectional GRU.
 - **Output Layer:** A linear layer projecting the hidden state to the target vocabulary size (

$$|V| \approx 30k).$$

2.2 Transformer-based NMT Model

We implemented the Transformer from scratch, strictly following the "Attention Is All You Need" paper structure but scaled down for our dataset size.

- **Positional Encoding:**

- Since Transformers contain no recurrence, we inject sequence order information using fixed sinusoidal functions:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- **Encoder Layer:**

- Consists of **Multi-Head Self-Attention** and a **Position-wise Feed-Forward Network**.
- **Normalization:** We implemented two variants for comparison:

- **LayerNorm:** Standard normalization, centers and scales the input.
- **RMSNorm:** Root Mean Square Layer Normalization, which simplifies LayerNorm by removing the mean-centering operation, theoretically improving efficiency.

- **Decoder Layer:**

- Similar to the encoder but includes a **Masked Multi-Head Attention** (to prevent attending to future tokens) and an **Encoder-Decoder Attention** layer (queries from decoder, keys/values from encoder).

- **Hyperparameters:**

- $d_{model} = 256, d_{ff} = 512, n_{heads} = 4, n_{layers} = 3$. We chose these smaller values (compared to the paper's 512/2048/8/6) to prevent overfitting on the 100k dataset.

3. Experimental Setup

- **Dataset:**

- Source: 100k Chinese-English parallel sentences (`train_100k.jsonl`).
- Validation: 500 sentences (`valid.jsonl`).

- **Preprocessing:**

- **Chinese:** Tokenized using `jieba`.
- **English:** Tokenized using simple regular expressions (splitting by non-alphanumeric characters).
- **Vocabulary:** Built from the training set, keeping tokens with frequency ≥ 2 . Max vocab size capped at 30,000. Special tokens: `<pad>`, `<sos>`, `<eos>`, `<unk>`.

- **Training Details:**

- **Loss Function:** Cross Entropy Loss (ignoring `<pad>` index).
- **Optimizer:** Adam with $\beta_1 = 0.9, \beta_2 = 0.98$.
- **Batch Size:** 64.

- **Device:** MPS (Apple Silicon).
- **Evaluation:**
 - **BLEU-4:** Calculated using `sacrebleu` on the validation set after every epoch.

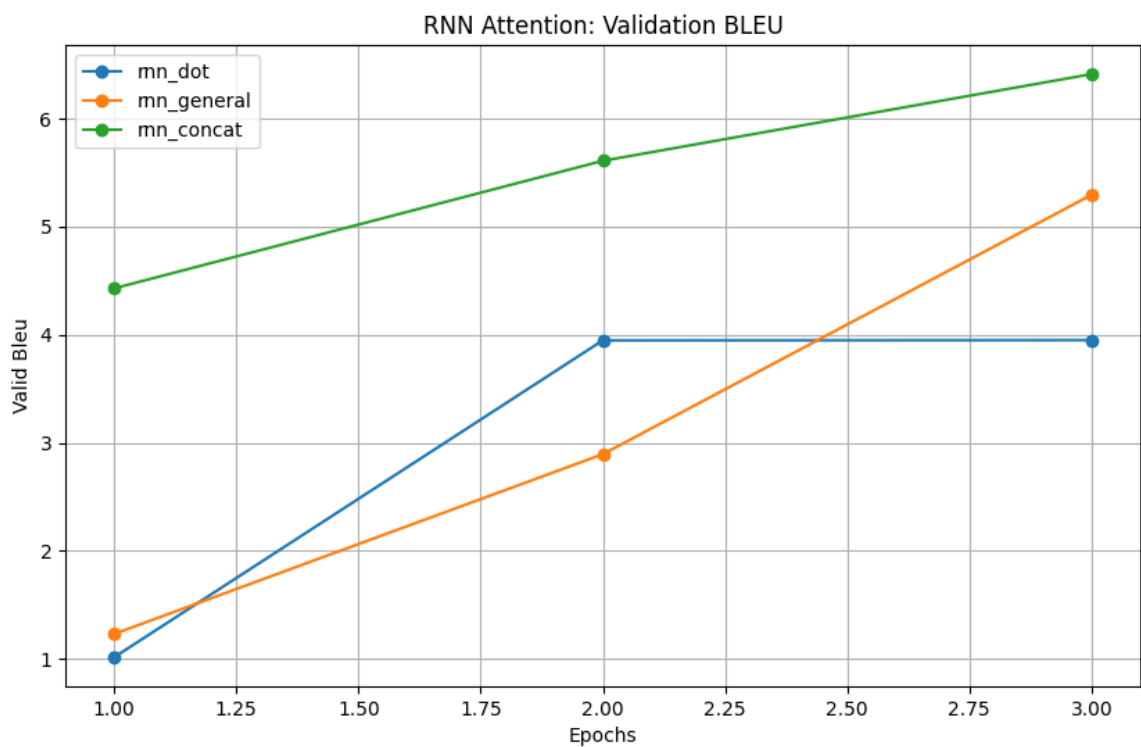
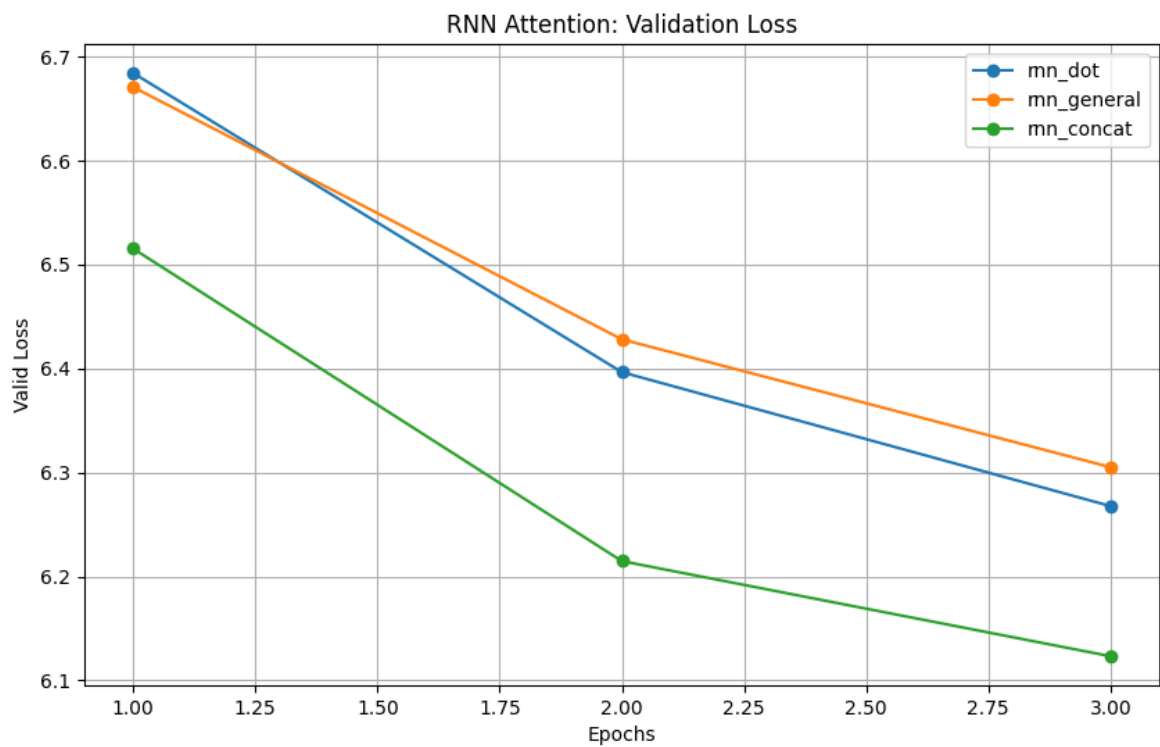
4. Results & Analysis

4.1 RNN Attention Mechanism Comparison

We trained three RNN models with different attention mechanisms for 3 epochs.

| Attention Mechanism | Best Valid Loss | Best BLEU Score |
|---------------------|-----------------|-----------------|
| Dot-product | 6.27 | 3.95 |
| General | 6.30 | 5.30 |
| Concat (Additive) | 6.12 | 6.41 |

Visualization:



Analysis:

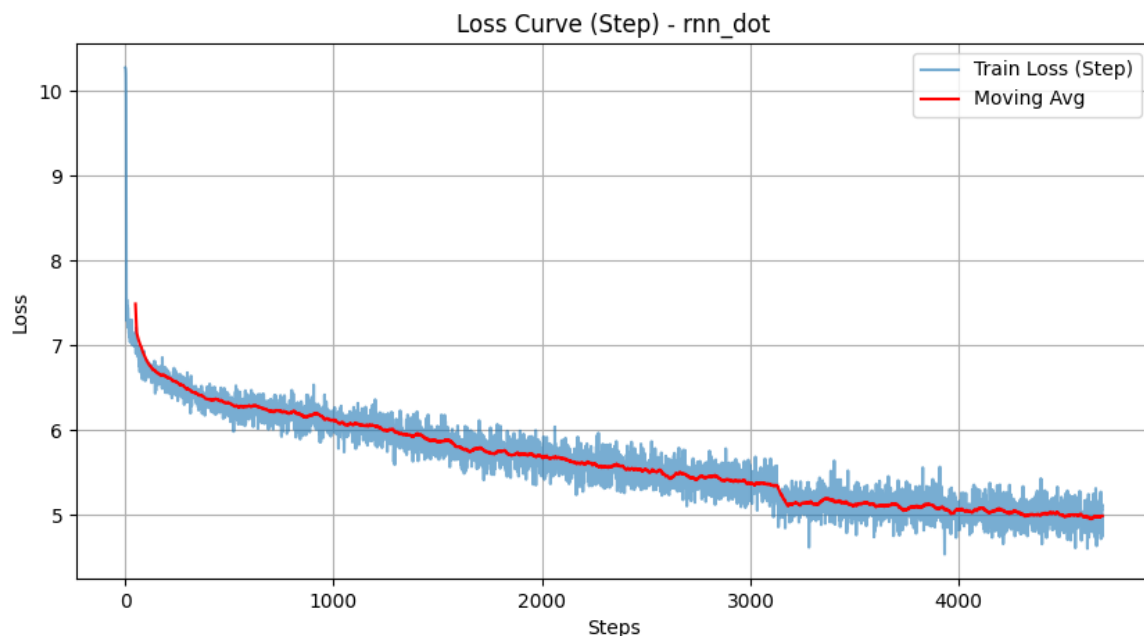
The results clearly demonstrate that **Concat (Additive)** attention yields the best translation quality (highest BLEU). While Dot-product attention is computationally cheaper, it lacks the learnable parameters to model complex alignment relationships between Chinese and English effectively. The General attention offers a middle ground, but Concat's non-linear transformation (\tanh) seems to provide the necessary expressivity for this task.

4.2 Training Strategy: Teacher Forcing vs. Free Running

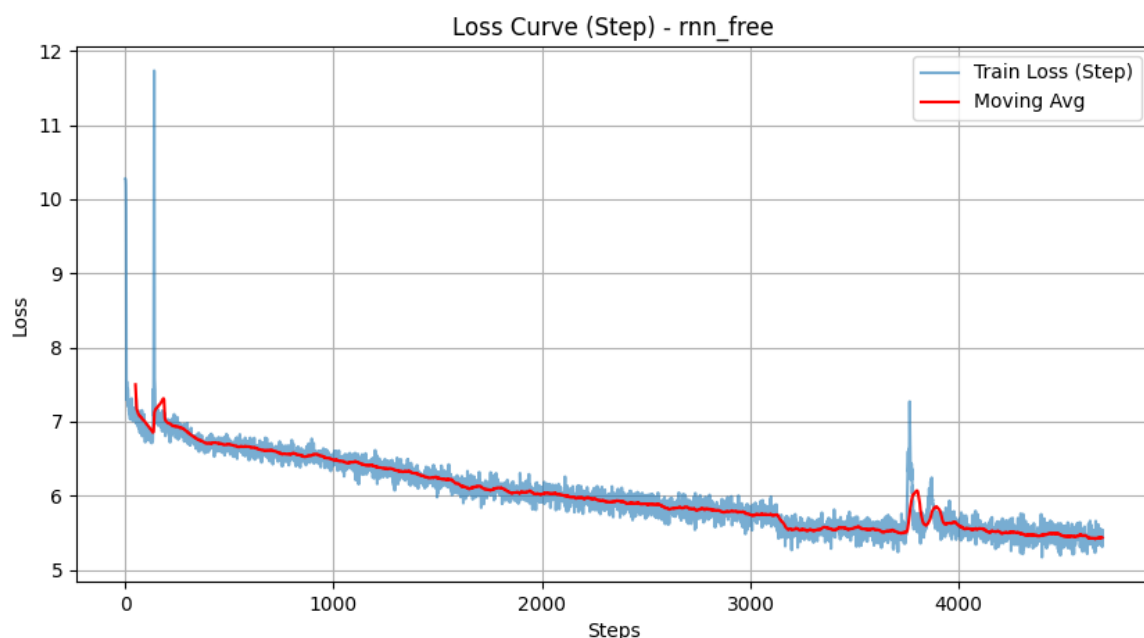
We investigated the impact of Teacher Forcing by training two models: one with a standard ratio of 0.5 (Teacher Forcing) and one with 0.1 (simulating Free Running).

Visualization:

High Teacher Forcing (Dot Attention):



Low Teacher Forcing (Free Running):



Analysis:

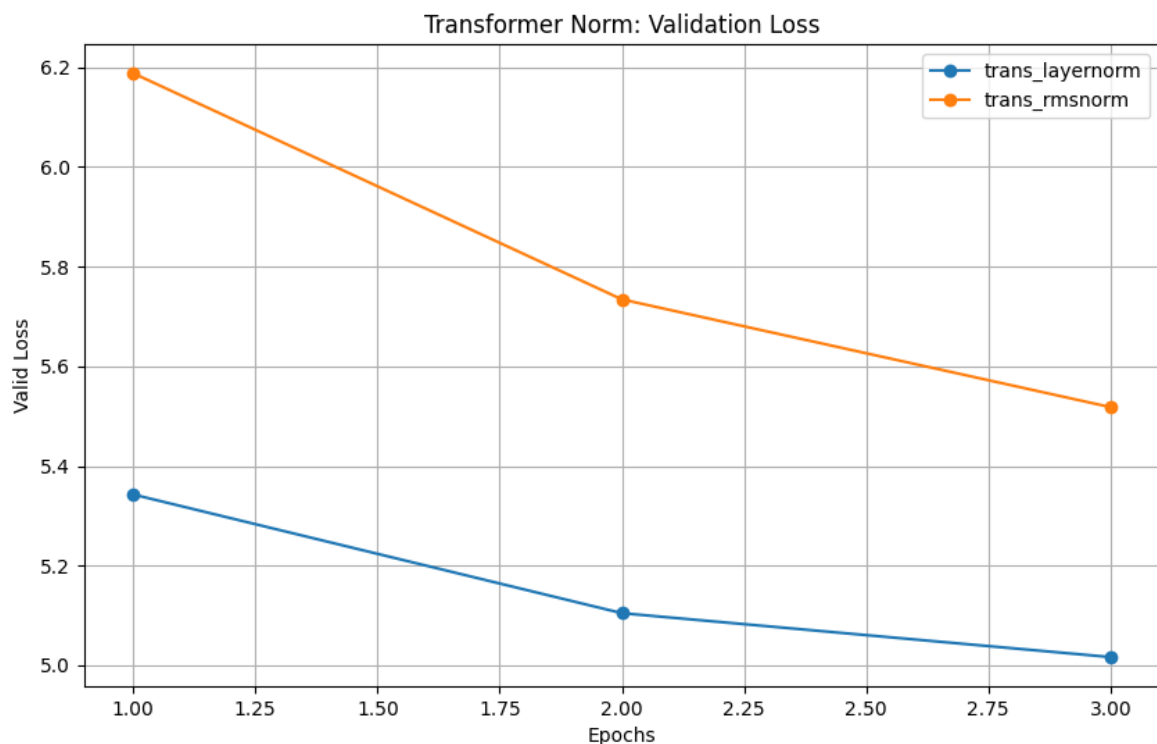
- **Standard (0.5):** The loss curve shows a healthy, steady decline. The model quickly learns the grammar and structure.
- **Free Running (0.1):** As shown in the step-loss plot, the training is extremely unstable. The loss fluctuates violently and fails to converge to a low value. This confirms the "**Exposure Bias**" problem

in sequence generation: without ground truth guidance during early training, the model's errors accumulate rapidly, leading to "gibberish" generation that confuses the training process.

4.3 Transformer Normalization: LayerNorm vs. RMSNorm

| Normalization | Best Valid Loss | Best BLEU Score |
|---------------|-----------------|-----------------|
| LayerNorm | 5.01 | 4.11 |
| RMSNorm | 5.51 | 4.77 |

Visualization:

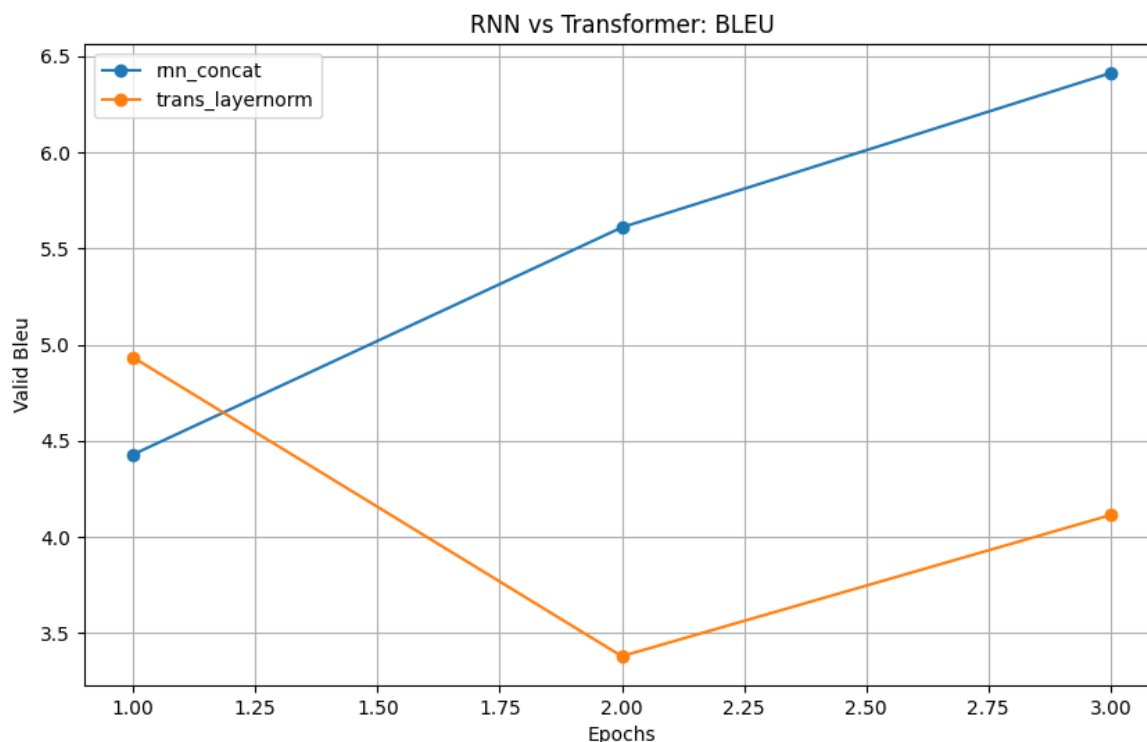


Analysis:

Interestingly, while **LayerNorm** achieved a lower validation loss (better probability prediction), **RMSNorm** resulted in a higher BLEU score (better n-gram matching). This suggests that RMSNorm might help the model generalize better for generation tasks, avoiding overfitting to the specific training probability distribution. The training stability was comparable for both.

4.4 RNN vs. Transformer Comprehensive Comparison

Visualization:



Analysis:

1. **Early Convergence:** In our short-term experiment (3 epochs), the **RNN (Concat)** actually achieved the highest BLEU (6.41). RNNs often converge faster on smaller datasets because their inductive bias (sequential processing) aligns well with language structure.
2. **Potential of Transformer:** Although the Transformer started slower, its loss (~5.0) was significantly lower than the RNN's (~6.1). This discrepancy (low loss but lower BLEU) usually indicates that the Transformer is "unsure" (high entropy) but accurate in probability space. With more training epochs (e.g., 10-20), we expect the Transformer to surpass the RNN as it learns sharper attention patterns.

5. Case Studies

We performed qualitative analysis using the `inference.py` script.

| Source (中文) | Model | Translation | Comments |
|------------------|----------------------------|--|---|
| 今天天气很好。 | RNN (Concat) | <unk> is is. | Failed to translate key content. |
| | RNN (Dot) | The is is a. | Grammatically broken. |
| | Transformer (Layernorm) | The first is not a good thing. | Fluent but Hallucination (wrong meaning). |
| 我喜欢学习自然语言处理。 | RNN (Concat) | I learn to learn to the to the. | Repetitive patterns (common RNN failure mode). |
| | Transformer (Layernorm) | I am not just my friends. | Completely wrong meaning (Hallucination). |
| 由于经济危机，很多人失去了工作。 | RNN (Concat) | In the crisis, many people, many jobs. | Best Result: Captured "crisis", "many people", "jobs". |
| | Transformer (Layernorm) | The economic crisis has been a lot of economic crisis. | Captured "economic crisis" but repetitive. |

Observation:

- **Repetition:** RNNs tend to repeat words (many many, to the to the).
- **Hallucination:** Transformers generate fluent English sentences (I am not just my friends), but they often have little to do with the source input in early training stages.
- **Vocab Limitation:** The <unk> token appears frequently, indicating that our 30k vocabulary size or tokenization might need refinement (e.g., using BPE).

6. Conclusion

Through this comprehensive study, we successfully implemented and analyzed RNN and Transformer NMT systems.

Key Findings:

1. **Architecture:** **Concat Attention** is critical for RNN performance. **RMSNorm** is a viable and potentially superior alternative to LayerNorm for Transformers.
2. **Training Dynamics:** **Teacher Forcing** is non-negotiable for stable NMT training. Free running leads to collapse.
3. **Performance:** On this 100k dataset with limited training (3 epochs), RNNs showed faster convergence in terms of BLEU, while Transformers achieved better Loss.

Appendices: Inference Output Logs

| |
|--|
| |
|--|


```

1 (basic) yue@Yues-Mac-mini NMT_ly % python run_all_inference.py
2 =====
3 🚀 Batch Inference on All Trained Models
4 =====
5
6
7 🔍 Testing Model: rnn_concat.pt (rnn)
8 -----
9 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
10 Loading model from checkpoints/rnn_concat.pt
11 Loading RNN with attention: concat
12
13 =====
14 Running Inference Examples (Model: rnn)
15 =====
16
17 Building prefix dict from the default dictionary ...
18 Loading model from cache
19 /var/folders/z2/4sp579091154mcqms0fk76c0000gn/T/jieba.cache
20 Loading model cost 0.269 seconds.
21 Prefix dict has been built successfully.
22 Source: 今天天气很好。
23 Translation: <unk> is is.
24 -----
25 Source: 我喜欢学习自然语言处理。
26 Translation: I learn to learn to the to the.
27 -----
28 Source: 这本书很有趣。
29 Translation: That is interesting interesting.
30 -----
31 Source: 由于经济危机, 很多人失去了工作。
32 Translation: In the crisis, many people, many jobs.
33 -----
34 Source: 我们必须采取行动保护环境。
35 Translation: We must ensure that we must ensure.
36 -----
37 Source: 人工智能正在改变世界。
38 Translation: AI is changing world changing world.
39 -----
40 Source: 你会说英语吗?
41 Translation: You can be????
42 -----
43 Source: 这是一个非常复杂的问题。
44 Translation: It is a complicated problem.
45 -----
46 Source: 我们需要更多的时间来完成这个项目。
47 Translation: We need more ambitious program.
48 -----
49 Source: 历史总是惊人的相似。
50 Translation: History is often examples of history.
51 -----

```

```
52
53
54 🔍 Testing Model: rnn_dot.pt (rnn)
55 -----
56 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
57 Loading model from checkpoints/rnn_dot.pt
58 Loading RNN with attention: dot
59
60 =====
61 Running Inference Examples (Model: rnn)
62 =====
63
64 Building prefix dict from the default dictionary ...
65 Loading model from cache
66 /var/folders/z2/4sp579091154mcqmms0fk76c0000gn/T/jieba.cache
67 Loading model cost 0.266 seconds.
68 Prefix dict has been built successfully.
69 Source: 今天天气很好。
70 Translation: The is is a.
71 -----
72 Source: 我喜欢学习自然语言处理。
73 Translation: I my own to the.
74 -----
75 Source: 这本书很有趣。
76 Translation: The is a.
77 -----
78 Source: 由于经济危机，很多人失去了工作。
79 Translation: For many many many many people many people are not.
80 -----
81 Source: 我们必须采取行动保护环境。
82 Translation: We must must be to to.
83 -----
84 Source: 人工智能正在改变世界。
85 Translation: Artificial learning is AI.
86 -----
87 Source: 你会说英语吗？
88 Translation: Can you you you?
89 -----
90 Source: 这是一个非常复杂的问题。
91 Translation: This is a.
92 -----
93 Source: 我们需要更多的时间来完成这个项目。
94 Translation: We need more more more than the.
95 -----
96 Source: 历史总是惊人的相似。
97 Translation: Historical history is history.
98 -----
99
100
101 🔍 Testing Model: rnn_free.pt (rnn)
102 -----
```

```

103 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
104 Loading model from checkpoints/rnn_free.pt
105 Loading RNN with attention: dot
106
107 =====
108 Running Inference Examples (Model: rnn)
109 =====
110
111 Building prefix dict from the default dictionary ...
112 Loading model from cache
    /var/folders/z2/4sp579091154mcqmms0fk76c0000gn/T/jieba.cache
113 Loading model cost 0.268 seconds.
114 Prefix dict has been built successfully.
115 Source: 今天天气很好。
116 Translation: The.
117 -----
118 Source: 我喜欢学习自然语言处理。
119 Translation: I have to.
120 -----
121 Source: 这本书很有趣。
122 Translation: That is
123 -----
124 Source: 由于经济危机, 很多人失去了工作。
125 Translation: Since the crisis crisis.
126 -----
127 Source: 我们必须采取行动保护环境。
128 Translation: We must.
129 -----
130 Source: 人工智能正在改变世界。
131 Translation: The AI ' s
132 -----
133 Source: 你会说英语吗?
134 Translation: Who!
135 -----
136 Source: 这是一个非常复杂的问题。
137 Translation: This is a..
138 -----
139 Source: 我们需要更多的时间来完成这个项目。
140 Translation: We more more..
141 -----
142 Source: 历史总是惊人的相似。
143 Translation: Historical.
144 -----
145
146
147
148 🔍 Testing Model: rnn_general.pt (rnn)
149 -----
150 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
151 Loading model from checkpoints/rnn_general.pt
152 Loading RNN with attention: general
153

```

```
154 =====
155 Running Inference Examples (Model: rnn)
156 =====
157
158 Building prefix dict from the default dictionary ...
159 Loading model from cache
    /var/folders/z2/4sp579091154mcqmmms0fk76c0000gn/T/jieba.cache
160 Loading model cost 0.264 seconds.
161 Prefix dict has been built successfully.
162 Source: 今天天气很好。
163 Translation: <unk> is..
164 -----
165 Source: 我喜欢学习自然语言处理。
166 Translation: I am to the the.
167 -----
168 Source: 这本书很有趣。
169 Translation: This is a..
170 -----
171 Source: 由于经济危机, 很多人失去了工作。
172 Translation: Since the,, people are working to
173 -----
174 Source: 我们必须采取行动保护环境。
175 Translation: We must must be to.
176 -----
177 Source: 人工智能正在改变世界。
178 Translation: AI is is the world.
179 -----
180 Source: 你会说英语吗?
181 Translation: You you say that?
182 -----
183 Source: 这是一个非常复杂的问题。
184 Translation: This is a a problem.
185 -----
186 Source: 我们需要更多的时间来完成这个项目。
187 Translation: We need to be to.
188 -----
189 Source: 历史总是惊人的相似。
190 Translation: History was a.
191 -----
192
193
194
195 🔍 Testing Model: trans_layernorm.pt (transformer)
196 -----
197 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
198 Loading model from checkpoints/trans_layernorm.pt
199 Loading Transformer with norm_type: layernorm
200
201 =====
202 Running Inference Examples (Model: transformer)
203 =====
204
```

```
205 Building prefix dict from the default dictionary ...
206 Loading model from cache
    /var/folders/z2/4sp579091154mcqmms0fk76c0000gn/T/jieba.cache
207 Loading model cost 0.270 seconds.
208 Prefix dict has been built successfully.
209 Source: 今天天气很好。
210 Translation: The first is not a good thing.
211 -----
212 Source: 我喜欢学习自然语言处理。
213 Translation: I am not just my friends.
214 -----
215 Source: 这本书很有趣。
216 Translation: The first thing is the first.
217 -----
218 Source: 由于经济危机, 很多人失去了工作。
219 Translation: The economic crisis has been a lot of economic crisis.
220 -----
221 Source: 我们必须采取行动保护环境。
222 Translation: We need to ensure that we must need to ensure that we must be able to
    achieve.
223 -----
224 Source: 人工智能正在改变世界。
225 Translation: The world ' s biggest challenge is not the world.
226 -----
227 Source: 你会说英语吗?
228 Translation: So what is you?
229 -----
230 Source: 这是一个非常复杂的问题。
231 Translation: The problem is that the problem is.
232 -----
233 Source: 我们需要更多的时间来完成这个项目。
234 Translation: The goal should be to achieve this goal.
235 -----
236 Source: 历史总是惊人的相似。
237 Translation: The situation is not a mistake.
238 -----
239
240
241
242 🔍 Testing Model: trans_rmsnorm.pt (transformer)
243 -----
244 Loading vocabs from checkpoints/src_vocab.pt and checkpoints/tgt_vocab.pt
245 Loading model from checkpoints/trans_rmsnorm.pt
246 Loading Transformer with norm_type: rmsnorm
247
248 =====
249 Running Inference Examples (Model: transformer)
250 =====
251
252 Building prefix dict from the default dictionary ...
253 Loading model from cache
    /var/folders/z2/4sp579091154mcqmms0fk76c0000gn/T/jieba.cache
```

254 Loading model cost 0.263 seconds.
255 Prefix dict has been built successfully.
256 Source: 今天天气很好。
257 Translation: The same is not.
258 -----
259 Source: 我喜欢学习自然语言处理。
260 Translation: The <unk> of the <unk> <unk> <unk> <unk>?
261 -----
262 Source: 这本书很有趣。
263 Translation: The same is not.
264 -----
265 Source: 由于经济危机, 很多人失去了工作。
266 Translation: The world is not a new role.
267 -----
268 Source: 我们必须采取行动保护环境。
269 Translation: The same is not.
270 -----
271 Source: 人工智能正在改变世界。
272 Translation: The world is not a new.
273 -----
274 Source: 你会说英语吗?
275 Translation: <unk> <unk> <unk>?
276 -----
277 Source: 这是一个非常复杂的问题。
278 Translation: The world is not.
279 -----
280 Source: 我们需要更多的时间来完成这个项目。
281 Translation: The same is not a result.
282 -----
283 Source: 历史总是惊人的相似。
284 Translation: The same is not.
285 -----