

1. Model analysis

In this Project 2.a assignment, the model I used consists of an embedding layer that converts each tokenized word into input vectors, which are then fed into an RNN layer. For the subsequent RNN layer, a single layer is sufficient for training, as the task involves simple arithmetic operations with three numbers (addition, subtraction, multiplication, and division). Finally, a fully connected layer compresses the output into a single class. For the loss function, I chose mean squared error to address the regression problem in this assignment.

The optimizer used is Adam.

2. Dataset analysis

I designed two datasets, each containing 100,000 samples. The first dataset consists of 90% addition and subtraction operations with three numbers, and 10% includes addition, subtraction, multiplication, and division operations.

The second dataset contains only addition and subtraction operations, with 50% consisting of three two-digit numbers and 50% of three one-digit numbers.

I believe that the quantity of data significantly impacts training results. For instance, the predictions for the addition and subtraction operations in the first dataset are likely to be better than those for the mixed operations. This is because the model is less familiar with the multiplication and division operations, leading to inherently poorer predictions compared to the more common addition and subtraction operations.

Left: Test results for 10,000 samples of multiplication and division with three numbers.

Right: Test results for 10,000 samples of addition and subtraction with three numbers.

```
(86*7786=)277
(86/86=)86
(45+54/20=)151
(46-42-41=)78
(7-56-40=)93
(51+99/81=)945
(32-55-54=)78
(25*19/45=)110
(11/36-44=)43
(58-33/57=)366
(58/91-36=)11785
(40-40-66=)460
(11/100-36=)21
(27/94*91=)12
(40+59-80=)111
(39+89-46=)82
(//72+23=)24
(100-33*26=)833
(63*68*21=)11785
(86/58/38=)8
(86/98/86=)10
(13+85-70=)27
(85/62/82=)174
(45/26-94=)93
(20+26-5=)45
(32*64+44=)1178
(8-10/11=)21
(42+4+21=)129
(83*33-25=)11785
(60-26/41=)1287
(68-27-86=)56
(29/50-40=)11744
(21/45+35=)229
(15*78/39=)48
(22/32-48=)51
(35+14*350
(79/184=)88
(51+8*56=)414
(35-27-15=)17
(30/83+55=)11785
(78/57/9=)38
(78/74*44=)36
(75+91/56=)75
(18-50-57=)87
(66-95*90=)1702
(37*32+100=)11785
(2-7+54=)24
(79/80+4=)112
(17-48+67=)1168
(30+39/29=)20
(17-58-10=)23
(12/56-78=)78
(71/39-47=)46
(100/10-3=)10
Test Loss:1864254000.8509466
```

```
(786-31-29=)755
(89-67+95=)86
(75+67-51=)96
(20-87+92=)20
(95-80+64=)75
(53+9+9=)1111
(79+19+34=)1132
(22-69-46=)85
(74-88-72=)349
(65-83+85=)70
(51-35+57=)73
(47-88-2=)1154
(37-37+19=)23
(19-80-93=)147
(22+1+51=)70
(17+90+33=)1337
(49+12-52=)11
(23-24+77=)78
(80-8-70=)20
(50+10-70=)131
(25+17-56=)48
(82+69-9=)140
(83-46-57=)34
(17+72-19=)70
(13-82-52=)119
(75-21-6+9=)11
(68-97-20=)49
(35-47+59=)43
(59+10-32=)38
(8-7-66=)59
(75+42-54=)57
(63+100+38=)202
(53+67+9=)131
(2+96+35=)130
(70+4+58=)131
(23-44-3=)22
(36+50-50=)32
(50-70+17=)98
(32-55-13=)40
(10-78-83=)45
(80+35-51=)63
(90+61-75=)76
(1+40-12=)27
(7-76-9=)82
(17-76-62=)43
(51+99+52=)254
(62+14-32=)42
(5+95+11=)109
(40+11+88=)1377
(80-53-88=)43
(84-88-14=)59
(11+82-62=)39
(70-30+23=)65
Test Loss:15936389267886959
```

Secondly, the range of the data results can impact training outcomes. For instance, the predictions for the one-digit operations in the second dataset are likely to be better than those for the two-digit operations. This is because datasets with a larger range of results require a greater volume of data to train the model on a more diverse set of combinations. As a result, the predictions for two-digit operations tend to be poorer than those for one-digit operations.

Left: Test results for 10,000 samples of three two-digit numbers.

Right: Test results for 10,000 samples of three one-digit numbers.

```
1.271100099e-0260
[85+68-58]=8
[82-21+581]=1142
[74-16-41]=118
[71-22-39]=100
[64-62-59]=38
[97+41-13]=165
[60-84+78]=54
[82+58+27]=167
[79+46-58]=67
[54+18-84]=173
[12+65-53]=24
[76+84+17]=1177
[48-27+98]=1115
[54-32+84]=106
[95+33-52]=76
[48-41+44]=121
[26-10-21]=5
[75-93-70]=-87
[33+61-57]=-47
[42+45+17]=104
[41+18-28]=52
[96+77-13]=136
[39+48+40]=127
[79+12-28]=22
[71-81-13]=-22
[38+10+87]=135
[48-66-51]=-69
[54+87-48]=94
[69-69-21]=-31
[76-62-78]=-53
[33+42-50]=-24
[85-49-25]=-12
[31+51+67]=1589
[48+73-59]=62
[75-67-86]=-78
[10+14-77]=-52
[27-69+40]=2
[17+60+96]=172
[78-48+68]=78
[36+2-23]=-20
[1+49+63]=127
[82-23+98]=157
[83+7-42]=-6
[53-46-24]=-18
[77+46-69]=54
[17+52+46]=111
[24-53+92]=63
[98+46-27]=115
[80+7+23]=66
[77+69+21]=98
[69+85-42]=121
[68-71-84]=-87
[41-82+46]=5
Test Loss: 0.3379486189672927
```

```
[6-3+4]=-3
[8+9+8]=24
[2+5-2]=-5
[2+9+0]=110
[1-0+0]=11
[6-9-2]=-5
[3+4+4]=11
[0-7+4]=-11
[2-0+8]=110
[5+0-8]=-3
[7-0-8]=-15
[3+4-2]=-5
[8+7-8]=7
[1+6-7]=-20
[5-0+2]=7
[7+2-5]=4
[4-9+7]=2
[7-1+9]=115
[3-4+1]=-2
[2+7+2]=11
[9+1+2]=18
[6-5-6]=-5
[1+8-9]=-20
[3+3+5]=15
[6+8+8]=22
[9+9-5]=13
[5+3+7]=9
[8+2+4]=6
[5-2+3]=6
[6-7-5]=-4
[7-0-1]=6
[7+3-0]=10
[7-9-9]=-11
[0+9-7]=2
[8+9-8]=9
[4+4+4]=12
[3+9-9]=5
[6+3+0]=9
[0+0-4]=-4
[7-8-2]=-3
[0-8-0]=-17
[2+9+1]=12
[3+5+5]=13
[1+6-9]=-4
[8+7-7]=8
[0-1+4]=5
[5+2-1]=6
[0+4+4]=13
[2+7+4]=113
[7-3+2]=13
[7+7-5]=9
[2+0-6]=-4
[8+8-0]=16
Test Loss: 0.06415668916736326
```

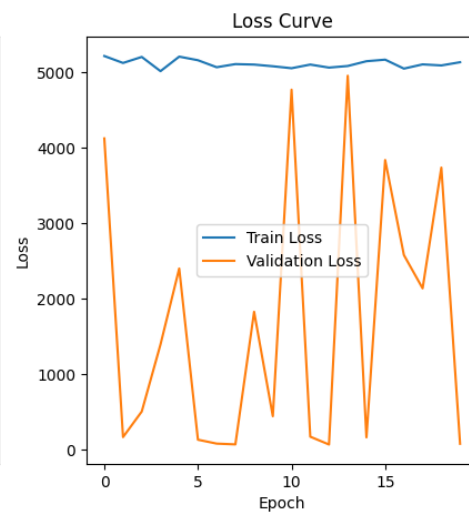
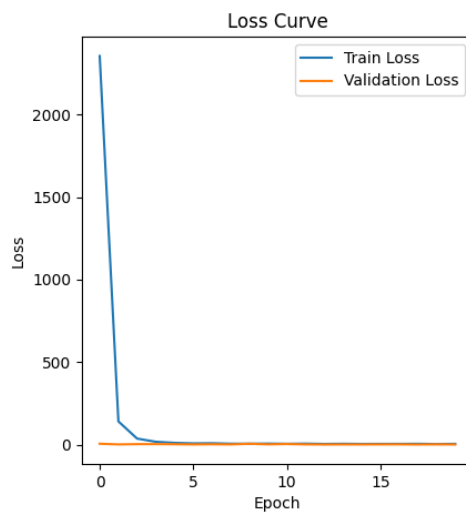
3. Discussion

In model training, a larger learning rate makes it more difficult for the gradient to converge to the minimum value and can lead to oscillations. However, the training process will be faster. Conversely, a smaller learning rate facilitates convergence of the gradient to the minimum value, but results in a slower training process.

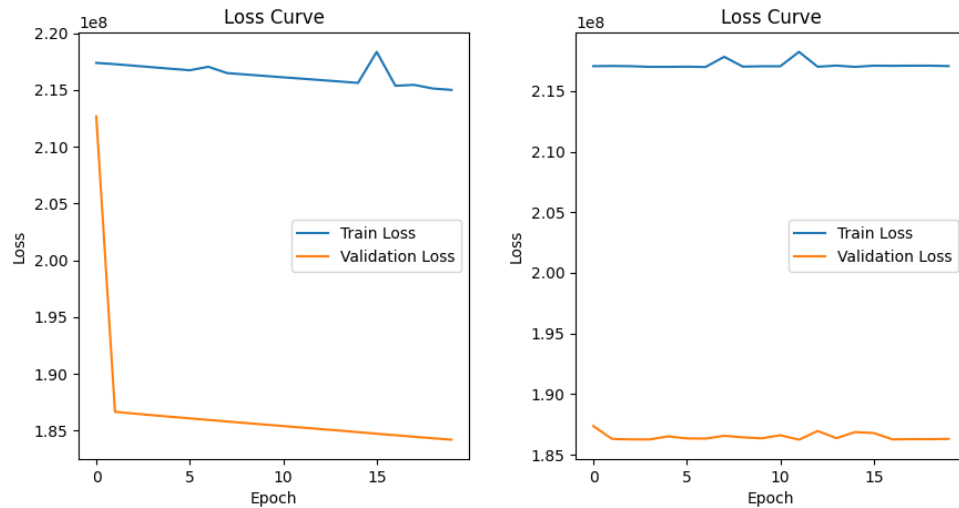
Left image: Learning rate = 0.001

Right image: Learning rate = 1

Dataset: 50% consists of operations with three two-digit numbers, and 50% consists of operations with three one-digit numbers.



Dataset: 90% consists of addition and subtraction operations with three numbers, and 10% consists of addition, subtraction, multiplication, and division operations with three numbers.

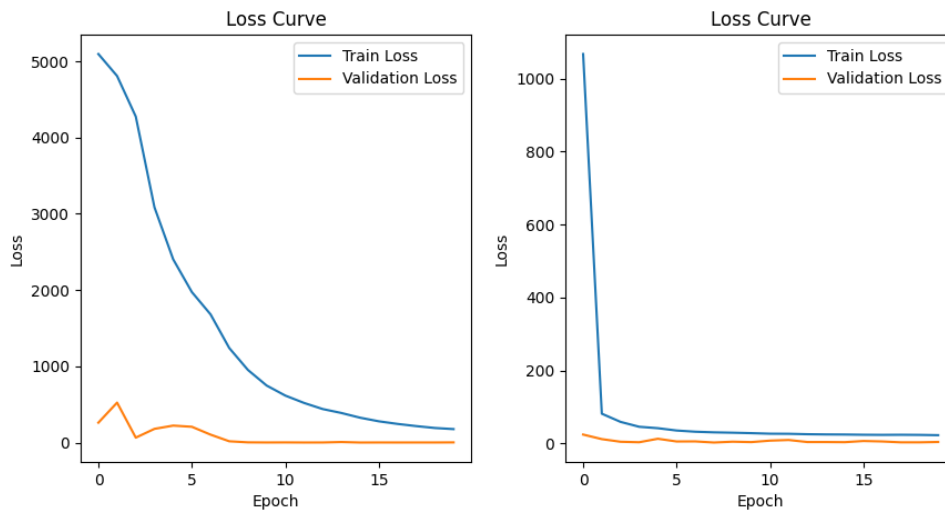


A larger batch size typically results in a smaller range of learning rates for convergence, and the loss tends to be higher. In contrast, a smaller batch size can successfully converge over a wider range of base learning rates, indicating that a smaller batch size allows for a larger range of learning rates for model convergence, making it easier to achieve convergence.

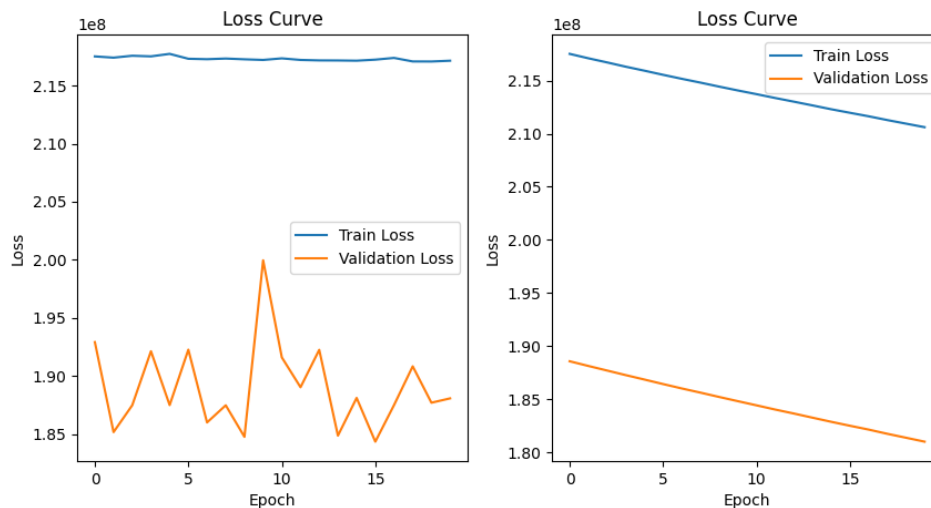
Left image: Batch size = 1024

Right image: Batch size = 16

Dataset: 50% consists of operations with three two-digit numbers, and 50% consists of operations with three one-digit numbers



Dataset: 90% consists of addition and subtraction operations with three numbers, and 10% consists of addition, subtraction, multiplication, and division operations with three numbers.



The model I used is a single-layer RNN. RNNs are particularly well-suited for tasks with sequential characteristics, and since this task only requires predicting the addition and subtraction results of three two-digit numbers, a single-layer RNN is capable of effectively handling this task.

4. I compared the model loss differences among LSTM, RNN, and GRU under the same conditions of batch size, learning rate, and dataset. Since LSTM and GRU have gates that control the filtering of historical data, they are better at preserving long-term memory compared to RNN, which lacks such gates.

Additionally, LSTM has a long-term memory cell that GRU does not have, allowing LSTM to achieve better prediction results than GRU when dealing with long sequences of data.

Below are the training loss and test results.

Left image: LSTM_train_loss

Center image: GRU_train_loss

Right image: RNN_train_loss

```
gpu
Epoch 1/20, Train Loss: 2485.33316190909485
Validation Loss: 852.1814228511625
Epoch 2/20, Train Loss: 542.4871326411166
Validation Loss: 276.9143177660293
Epoch 3/20, Train Loss: 206.63062674707308
Validation Loss: 128.71765714848178
Epoch 4/20, Train Loss: 166.47776262488232
Validation Loss: 76.51806679664005
Epoch 5/20, Train Loss: 64.46848702972585
Validation Loss: 48.55317265689373
Epoch 6/20, Train Loss: 39.64319886062904
Validation Loss: 30.760595946219585
Epoch 7/20, Train Loss: 25.28712718330535
Validation Loss: 21.5284721031785
Epoch 8/20, Train Loss: 17.611710683148
Validation Loss: 15.26335711877556
Epoch 9/20, Train Loss: 11.80874355983887
Validation Loss: 10.820367180625471
Epoch 10/20, Train Loss: 8.796291812614004
Validation Loss: 8.046138029139042
Epoch 11/20, Train Loss: 6.487427123862803
Validation Loss: 6.930682366747626
Epoch 12/20, Train Loss: 5.424467479474735
Validation Loss: 5.6337299156056
Epoch 13/20, Train Loss: 4.086877065795389
Validation Loss: 4.2761583707011795
Epoch 14/20, Train Loss: 3.7131003701173296
Validation Loss: 3.9614970952774717
Epoch 15/20, Train Loss: 3.2816884952637277
Validation Loss: 3.144448970151475
Epoch 16/20, Train Loss: 2.648863200980821
Validation Loss: 2.3803962614279213
Epoch 17/20, Train Loss: 2.212427727642723
Validation Loss: 2.563663986325264
Epoch 18/20, Train Loss: 1.8172687644308263
Validation Loss: 1.8881881576317165
Epoch 19/20, Train Loss: 1.5385176700911796
Validation Loss: 1.6261639858399535
Epoch 20/20, Train Loss: 1.1268559418280015
Validation Loss: 1.2997114099562168
```

```
gpu
Epoch 1/20, Train Loss: 2334.8241068288046
Validation Loss: 824.3733024597168
Epoch 2/20, Train Loss: 514.6542138186368
Validation Loss: 260.72211140336394
Epoch 3/20, Train Loss: 181.77358598939398
Validation Loss: 111.4207495611206
Epoch 4/20, Train Loss: 81.67846391931017
Validation Loss: 57.614952969551084
Epoch 5/20, Train Loss: 40.0040966607376
Validation Loss: 33.070898489521
Epoch 6/20, Train Loss: 28.35607987154326
Validation Loss: 20.4023570015277
Epoch 7/20, Train Loss: 17.140618570318278
Validation Loss: 13.347754579782485
Epoch 8/20, Train Loss: 10.881443574362613
Validation Loss: 8.375298416146532
Epoch 9/20, Train Loss: 8.222422123870888
Validation Loss: 6.375298416146532
Epoch 10/20, Train Loss: 5.730487226000563
Validation Loss: 5.207543715453625
Epoch 11/20, Train Loss: 4.033718787856238
Validation Loss: 3.763177741384983
Epoch 12/20, Train Loss: 3.0565973786224107
Validation Loss: 1.947823389426697
Epoch 13/20, Train Loss: 2.6073757320274717
Validation Loss: 1.13493147928248
Epoch 14/20, Train Loss: 2.07869171513434
Validation Loss: 2.055666895747187
Epoch 15/20, Train Loss: 1.748400664602696
Validation Loss: 1.709600009037188
Epoch 16/20, Train Loss: 1.347457330878074
Validation Loss: 1.94319588896124
Epoch 17/20, Train Loss: 1.21227482050163
Validation Loss: 1.185354831304882
Epoch 18/20, Train Loss: 1.08394304666233
Validation Loss: 1.421112375161172
Epoch 19/20, Train Loss: 1.004225480404151
Validation Loss: 1.259777541458066
Epoch 20/20, Train Loss: 1.1268335155966924
Validation Loss: 1.167155818749428
```

```
*****
gpu
Epoch 1/20, Train Loss: 3120.1541202099744
Validation Loss: 192.1651489297811
Epoch 2/20, Train Loss: 510.47384238243103
Validation Loss: 220.6257582021232
Epoch 3/20, Train Loss: 158.76157944310796
Validation Loss: 86.5145834455633
Epoch 4/20, Train Loss: 71.5563600077058
Validation Loss: 47.46667387199402
Epoch 5/20, Train Loss: 38.39054435762492
Validation Loss: 28.78952346314187
Epoch 6/20, Train Loss: 23.34774777700278
Validation Loss: 16.081708495789237
Epoch 7/20, Train Loss: 14.661707812673557
Validation Loss: 11.499975323677063
Epoch 8/20, Train Loss: 10.991881454194134
Validation Loss: 8.28107011651982
Epoch 9/20, Train Loss: 8.982188263455102
Validation Loss: 6.69745836640167
Epoch 10/20, Train Loss: 6.4773452404421
Validation Loss: 7.928260564804077
Epoch 11/20, Train Loss: 4.726984254184907
Validation Loss: 4.728048950897899
Epoch 12/20, Train Loss: 6.111019644886255
Validation Loss: 5.54247151374818
Epoch 13/20, Train Loss: 5.5552073966777325
Validation Loss: 7.157616048069
Epoch 14/20, Train Loss: 4.6661612838831814
Validation Loss: 2.91338180353736
Epoch 15/20, Train Loss: 3.4509913248839597
Validation Loss: 3.63029677093525
Epoch 16/20, Train Loss: 2.78441186651748
Validation Loss: 3.14219875287029
Epoch 17/20, Train Loss: 3.3094917033428104
Validation Loss: 2.2212479077279568
Epoch 18/20, Train Loss: 3.33627678788723
Validation Loss: 2.925554132461518
Epoch 19/20, Train Loss: 3.2946212794615463
Validation Loss: 3.98143602612646
Epoch 20/20, Train Loss: 3.374931660505463
Validation Loss: 6.02193198800087
```

We can observe that the loss results for LSTM and GRU are quite similar, both around 1. This is due to the relatively short length of the data. In contrast, RNN performs worse and experiences loss fluctuations during training.

Left image: LSTM_test_loss

Center image: GRU_test_loss

Right image: RNN_test_loss

```
[-1.8,-1.12
(9.0,-3.12
(0.0,-9.12
(8.6,-1.15
(2.2,9.12
(9.4,-9.12
(6.9,-1.18
(0.3,-8.15
(5.6,-9.18
(5.7,-16.18
(7.1,-8.13
(9.5,-8.122
(0.3,-9.16
(0.8,-2.110
(1.8,-4.113
(4.4,-4.14
(9.4,-4.18
(4.3,-1.18
(9.7,-1.115
(6.8,-1.17
(5.0,-4.14
(9.7,-5.17
(5.5,-1.19
(9.1,-5.113
(6.8,-5.13
(4.7,-2.10
(7.2,-1.18
(9.6,-0.13
(7.9,-0.17
(0.4,-5.11
(7.0,-5.14
(7.2,-6.111
(0.3,-9.16
(9.7,-1.115
(3.4,-4.14
(1.7,8.10
(0.5,-1.17
(3.4,-3.15
(7.3,-0.14
(2.3,-2.15
(9.2,-0.111
(1.1,-0.10
(1.1,-0.12
(7.7,-8.110
(9.3,-5.117
(4.8,-4.12
(3.4,-1.11
(2.2,-2.12
(1.7,-8.11
(6.0,-3.17
(7.5,-16.17
(0.0,-4.14
(8.3,-3.12
Test Loss: 0.902663154734538
```

```
[-1.8,-1.12
(9.0,-3.12
(0.0,-9.12
(8.6,-1.15
(2.2,9.12
(9.4,-9.12
(6.9,-1.18
(0.3,-8.15
(5.6,-9.18
(5.7,-16.18
(7.1,-8.13
(9.5,-8.122
(0.3,-9.16
(0.8,-2.110
(1.8,-4.113
(4.4,-4.14
(9.4,-4.18
(4.3,-1.18
(9.7,-1.115
(6.8,-1.17
(5.0,-4.14
(9.7,-5.17
(5.5,-1.19
(9.1,-5.113
(6.8,-5.13
(4.7,-2.10
(7.2,-1.18
(9.6,-0.13
(7.9,-0.17
(0.4,-5.11
(7.0,-5.14
(7.2,-6.111
(0.3,-9.16
(9.7,-1.115
(3.4,-4.14
(1.7,8.10
(0.5,-1.17
(3.4,-3.15
(7.3,-0.14
(2.3,-2.15
(9.2,-0.111
(1.1,-0.10
(1.1,-0.12
(7.7,-8.110
(9.3,-5.117
(4.8,-4.12
(3.4,-1.11
(2.2,-2.12
(1.7,-8.11
(6.0,-3.17
(7.5,-16.17
(0.0,-4.14
(8.3,-3.12
Test Loss: 0.9625463007714389
```

```
[-1.8,-1.12
(9.0,-3.12
(0.0,-9.12
(8.6,-1.15
(2.2,9.12
(9.4,-9.12
(6.9,-1.18
(0.3,-8.15
(5.6,-9.18
(5.7,-16.18
(7.1,-8.13
(9.5,-8.122
(0.3,-9.16
(0.8,-2.110
(1.8,-4.113
(4.4,-4.14
(9.4,-4.18
(4.3,-1.18
(9.7,-1.115
(6.8,-1.17
(5.0,-4.14
(9.7,-5.17
(5.5,-1.19
(9.1,-5.113
(6.8,-5.13
(4.7,-2.10
(7.2,-1.18
(9.6,-0.13
(7.9,-0.17
(0.4,-5.11
(7.0,-5.14
(7.2,-6.111
(0.3,-9.16
(9.7,-1.115
(3.4,-4.14
(1.7,8.10
(0.5,-1.17
(3.4,-3.15
(7.3,-0.14
(2.3,-2.15
(9.2,-0.111
(1.1,-0.10
(1.1,-0.12
(7.7,-8.110
(9.3,-5.116
(4.8,-4.12
(3.4,-1.11
(2.2,-2.12
(1.7,-8.11
(6.0,-3.17
(7.5,-16.17
(0.0,-4.14
(8.3,-3.12
Test Loss: 1.3851681714883288
```

During testing, we can see that the results for LSTM and GRU are quite similar,

while RNN performs worse.