**1.a)** Pearson Correlation: 0.8645



Training and Validation Loss (Batch-wise)



Batchwise Losses
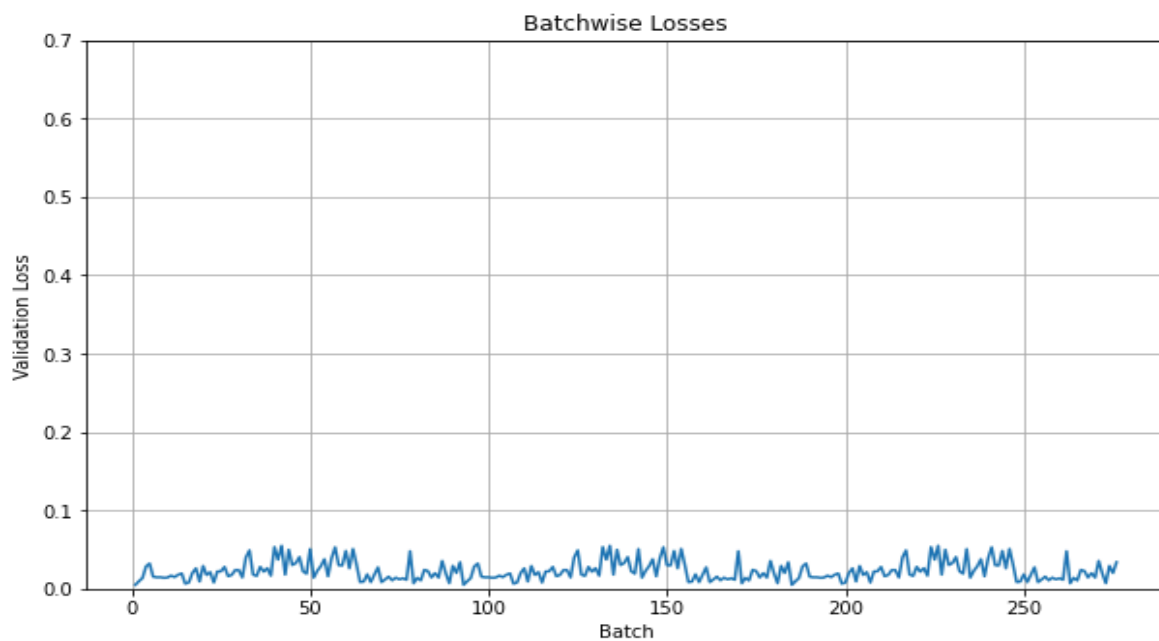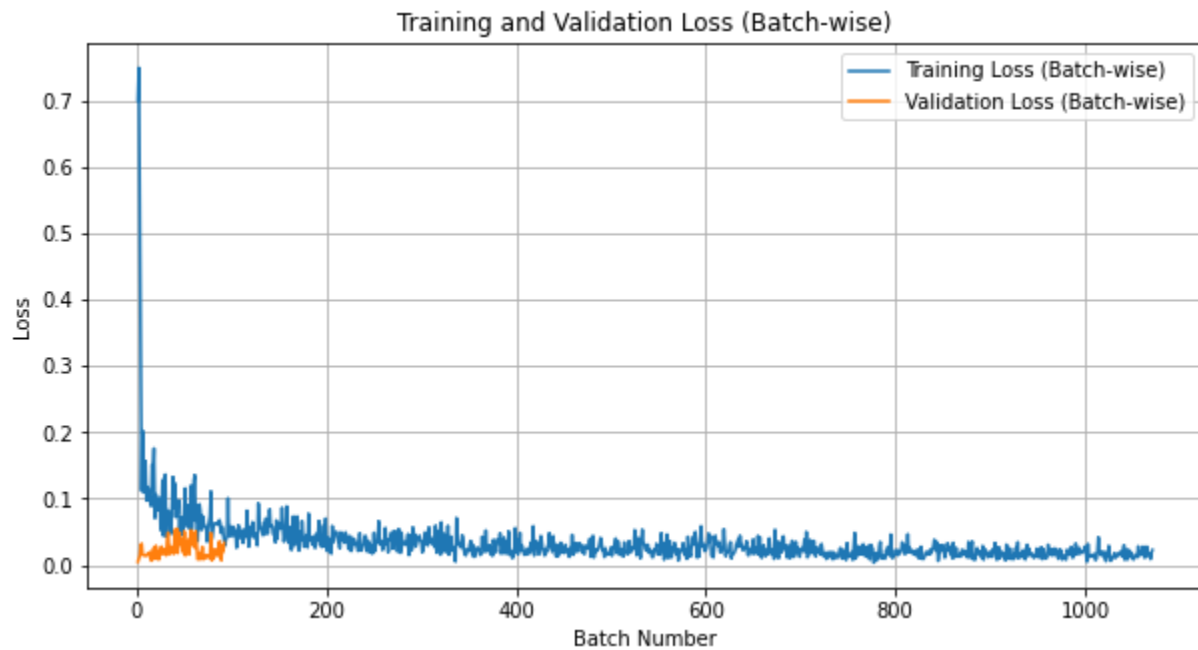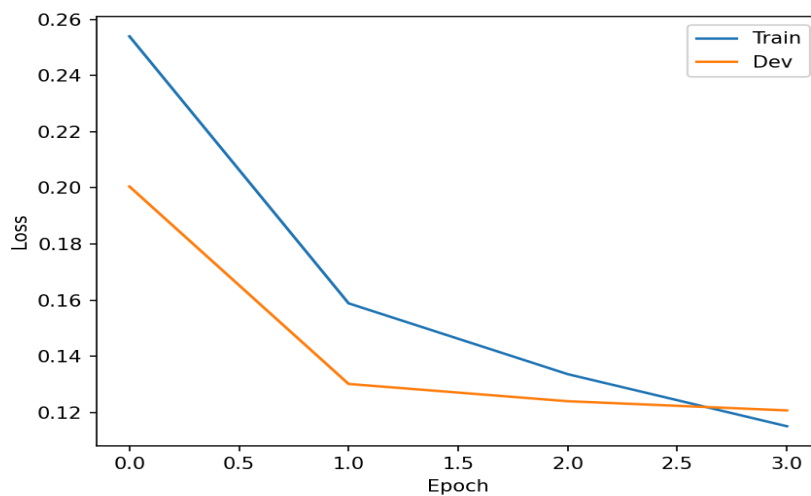
**Average Training Loss: 0.0183, Average Validation Loss: 0.0228**
**Pearson Correlation: 0.8645**

The very slightly higher average validation loss of 0.02 indicates that the model's performance on unseen data is not as strong as on the training data, highlighting the need to address minor overfitting through regularization techniques or model simplification.

**1.b**) Pearson Correlation: 0.8491
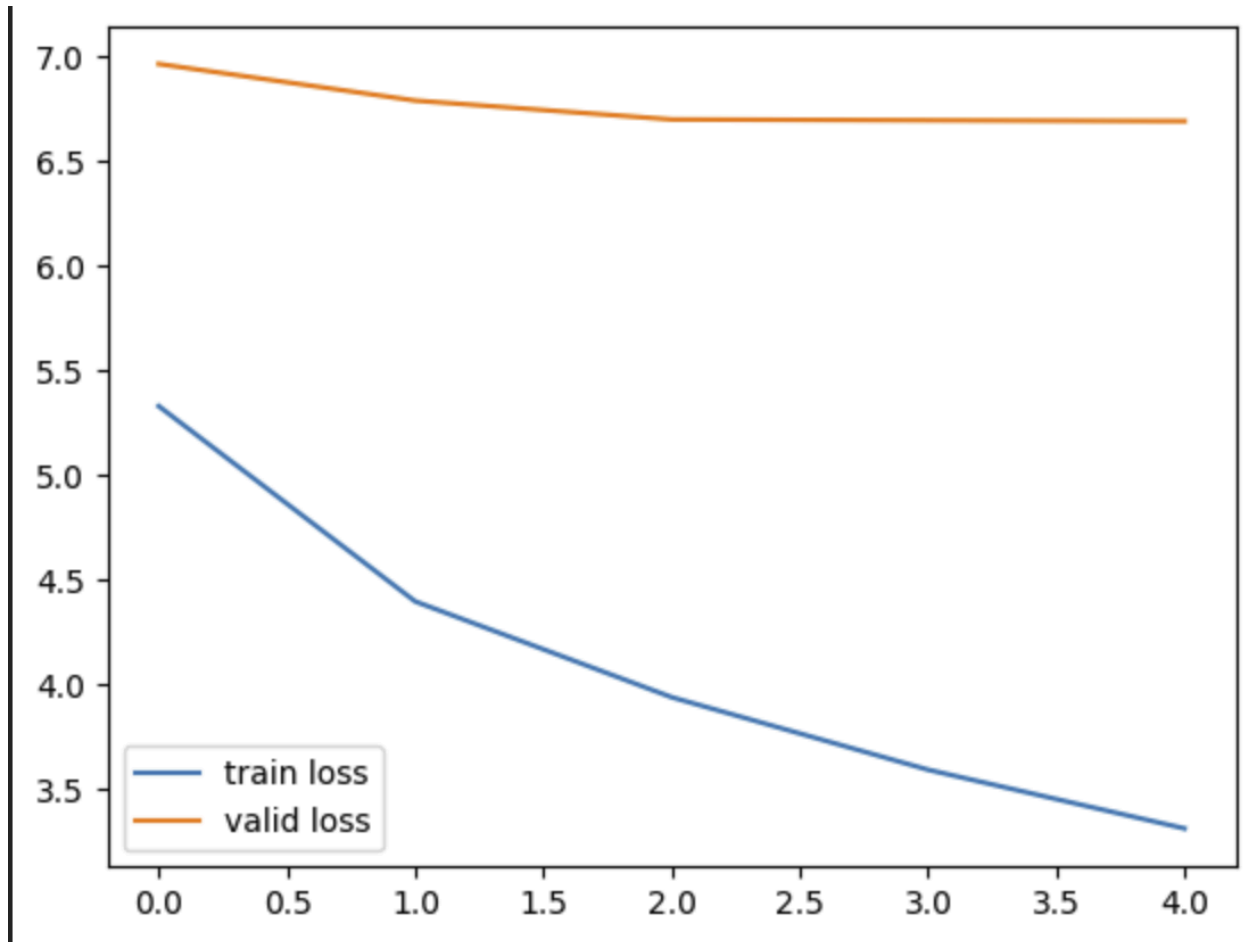
**1.c)**



Pearson Correlation: 0.8793

Avg Training and Validation Losses:  0.1653841813416994, 0.1433219905637348

Setup 1B achieved the lowest Pearson Correlation (0.8491) due to the limited capture of semantic nuances by cosine similarity and embeddings of the pre trained Sentence-BERT model alone. In contrast, Setup 1A utilized training on BERT's semantic understanding and specific configurations for better similarity predictions, resulting in a higher correlation (0.8645). Setup 1C's fine-tuning of Sentence-BERT with a task-specific loss function led to the highest correlation (0.879) by capturing task nuances and leveraging Sentence-BERT's strong semantic representations.

2.a) Consistent decrease in losses

**Bleu Score**

Val: {'bleu': 0.02, 'precisions': [0.02], 'brevity_penalty': 1.0, 'length_ratio': 10.0, 'translation_length': 100, 'reference_length': 10}

Test: {'bleu': 0.0, 'precisions': [0.0], 'brevity_penalty': 1.0, 'length_ratio': 31.0, 'translation_length': 93, 'reference_length': 3}

**Bert Score**

Val: {'precision': [0.6361724138259888], 'recall': [0.6481119394302368], 'f1': [0.6420866847038269], 'hashcode': 'distilbert-base-uncased_L5_no-idf_version=0.3.12(hug_trans=4.37.2)'}

Test: {'precision': [0.6335083842277527], 'recall': [0.6447191834449768], 'f1': [0.639064610004425], 'hashcode': 'distilbert-base-uncased_L5_no-idf_version=0.3.12(hug_trans=4.37.2)'}

**Meteor Score**

Val: {'meteor': 0.052631578947368425}

Test: {'meteor': 0.0}

**2.b)** meteor: test 0.5772148087831891
Val: 0.5458072768029277


Bert val: Average Precision: 0.9218533658772473
Average Recall: 0.9203263581376188
Average F1-score: 0.9210045642707687

Test: Average Precision: 0.9273441431521574
Average Recall: 0.9249510829271417
Average F1-score: 0.9260744284534105

Bleu: tets: 0.6173815483698936, 0.4767904215770271, 0.3812786762509227,
0.31045530684044015

Val: 0.5896815800080613, 0.4406826405673354, 0.34410406607004834,
0.27436246022629457

2. c) Consistent decrease in losses

### eval/loss



### train/loss



Test Bleu:
0.5262450627132879
0.36582935681843914
0.26621435971067786
0.198371742684942

Val Bleu:
0.49367290122530916

0.32905018030188277
0.23163742460062606
0.16812812937367191

Bert Score:

Val:Average Precision: 0.8538484167481084
Average Recall: 0.8531928047324832
Average F1-score: 0.853293219625043

{'precision': [0.7615596055984497, 0.9409469366073608, 0.9765104651451111, 0.8342380523681641, 0.8931359648704529, 0.847620964050293, 0.8608109283440727, 0.8261321783065796, 0.8516960144042969, 0.9173384308815002, 0.8862889409065247, 0.778535008430481, 0.8676823377609253, 0.8386300802230835, 0.9089933037757874, 0.9204580187797546, 0.8180539011955261, 0.8242896795272827, 0.8404986585836792, 0.9005566835403442, 0.9201635122299194, 0.8513404130935669, 0.8843384981153396, 0.8449311256408691, 0.8072896080372145, 0.8719667196273804, 0.8648838996887207, 0.9142832159996833, 0.8509263396263123, 0.889707647614664, 0.8613384962081909, 0.8456195592880249, 0.9224295616149902, 0.8686690330505371, 0.8532232046127319, 0.8777071237564087, 0.8942518302...

{'precision': [0.9565405844564209, 0.9432793855667114, 0.9297959804534912, 0.8565741777420044, 0.9748955965042114, 0.8853832483291626, 0.8974635004997253, 0.8996868239212036, 0.8640740565338135, 0.7279601097106934, 0.8768697381019592, 0.7714994781559494019, 0.8951216936111145, 0.7954033017158508, 0.8706485629081726, 0.8036237955093384, 0.7843533158302307, 0.8634216189338446, 0.9420335292816162, 0.8879314661026001, 0.8969099521636963, 0.9107899665833313, 0.8871927261353539, 0.9014237523887028, 0.8765199184417725, 0.908748886512756...

Test: Average Precision: 0.86248078213886
Average Recall: 0.8638476273464178
Average F1-score: 0.8629324492036045
Meteor:
Test 0.48825410834957206

Val: 0.44629732228246627


Model 2a is the worst. This model is not a good choice for machine translation as the complexity of this model is not adequate for translating german to english as both of these are really complex languages whereas our model is quite simple.
Model 2b performs very good as t5 is trained on huge datasets. T5 is known to outperform the state-of-the-art for a variety of tasks
Model 2c performs better than 2a but worse than 2b as we have fine tuned t5 for a different translation task compared to 2b.