

# Introduction

In this repository sparsity in the feed-forward architecture was tested. The method was implemented in an encoder based transformer model for NLP and a simple neural network for Mnist problem. Experiments were made to compare sparse architecture to the baseline feed-forward on speed and accuracy. The sparse feed forward component was a variation of that presented in the article: Jaszczur, S., et al. (2021). Sparse is Enough in Scaling Transformers.

## Encoder based transformer

Transformer model parameters:  $d\_model=384$ ,  $n\_heads=6$ ,  $n\_layers=6$ ,  $context=256$ ,  $tr\_batch\_size=64$ ,  $lr=3e-4$ ,  $dropout=0.2$ . During experiments with transformers in sparse models 1/64 of all parameters were used in feed forward layers per token inference (equals to 24 ReLu connections selection out of 1536 in  $FF[d\_model \times d\_model^4 - ReLu - d\_model^4 \times d\_model]$ ).

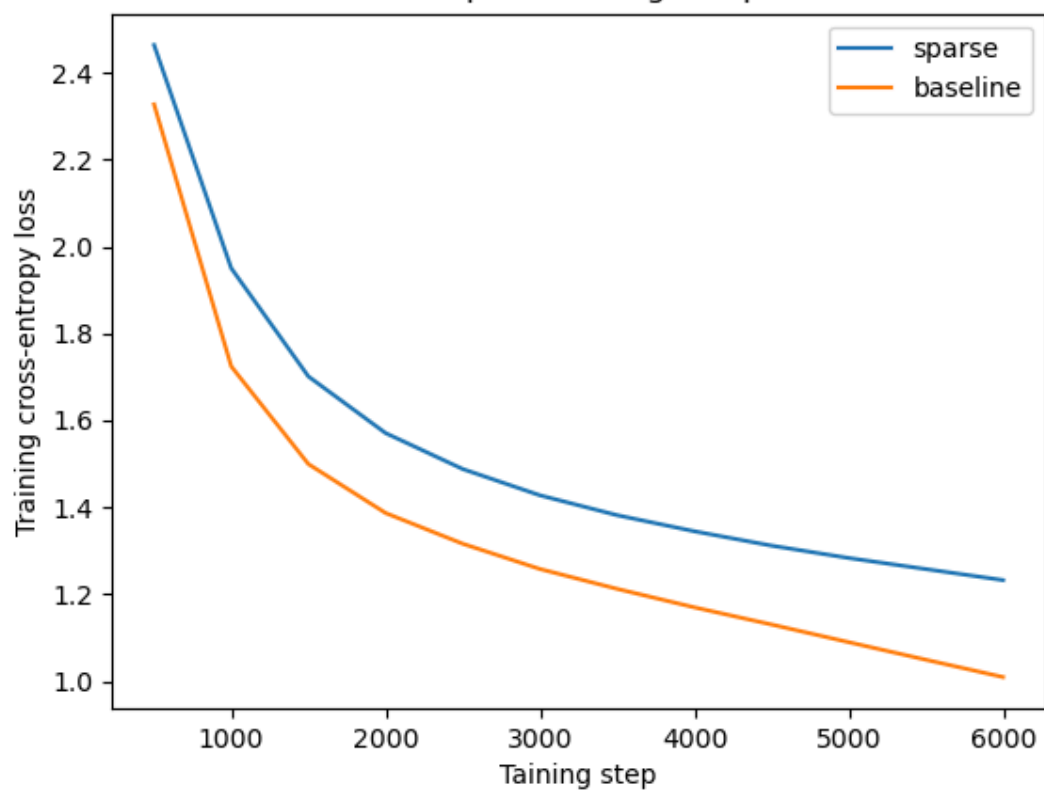
Problem: Tiny Shakespeare text generation

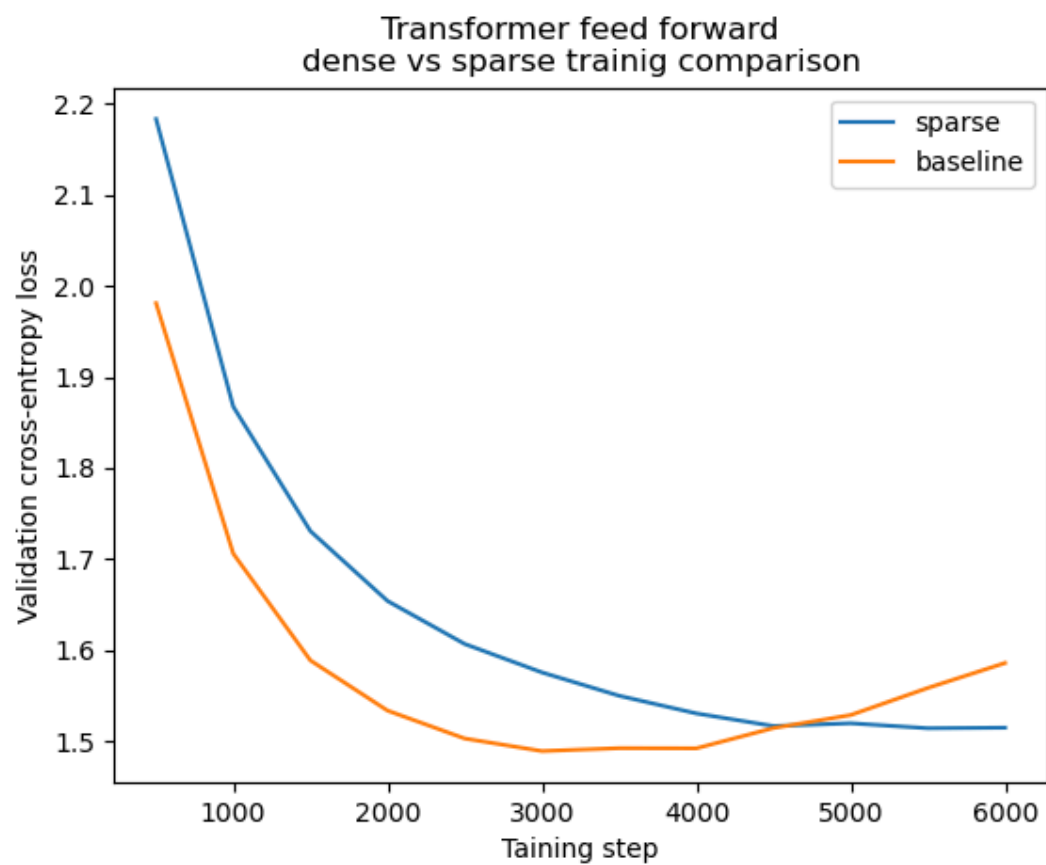
## Training

Baseline model training time (6K steps): 21.35 min

Sparse model training time (6K steps): 25.28 min

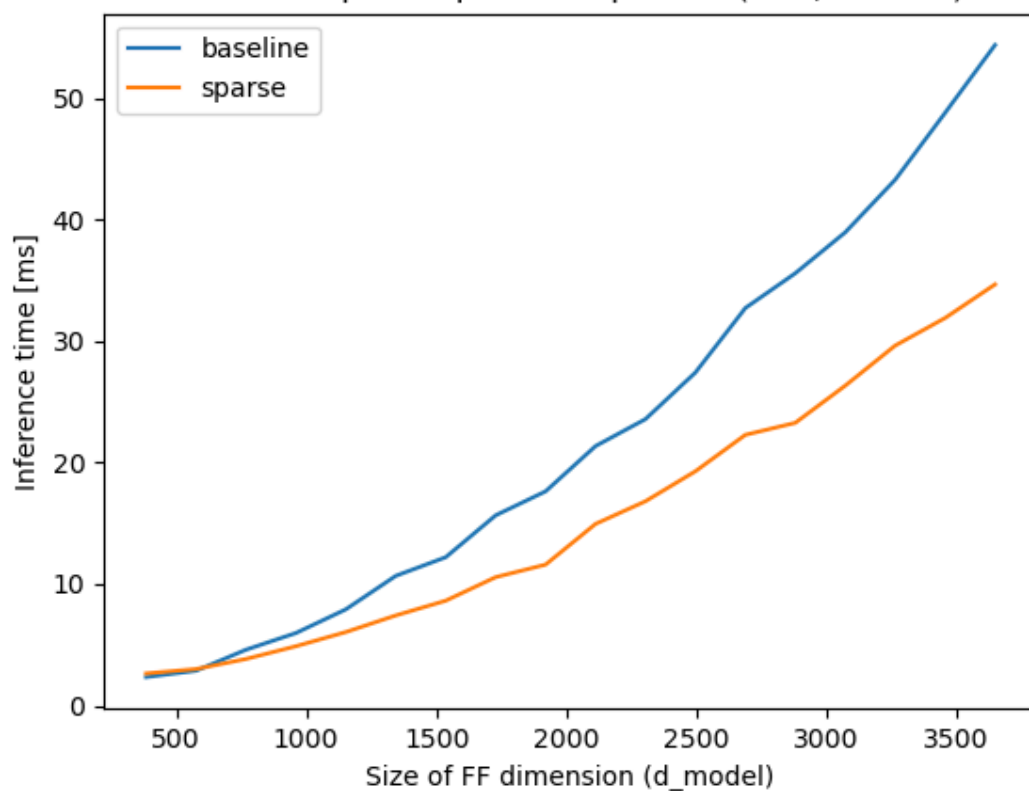
Transformer feed forward  
dense vs sparse trainig comparison



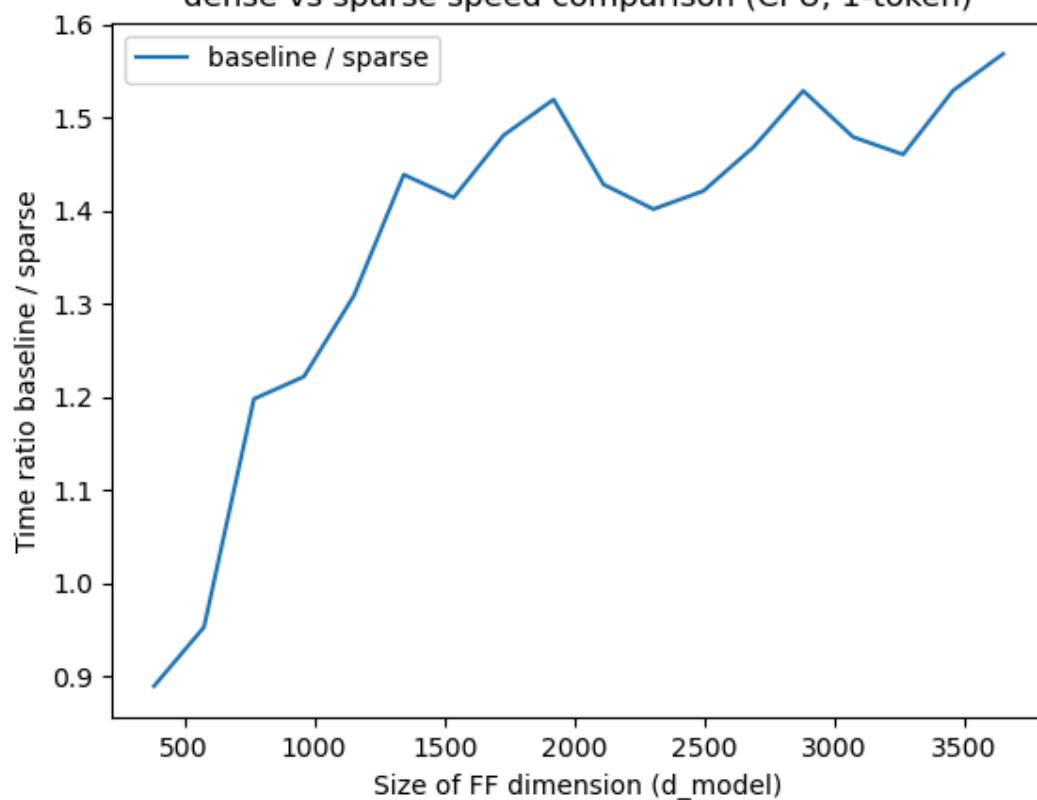


Inference

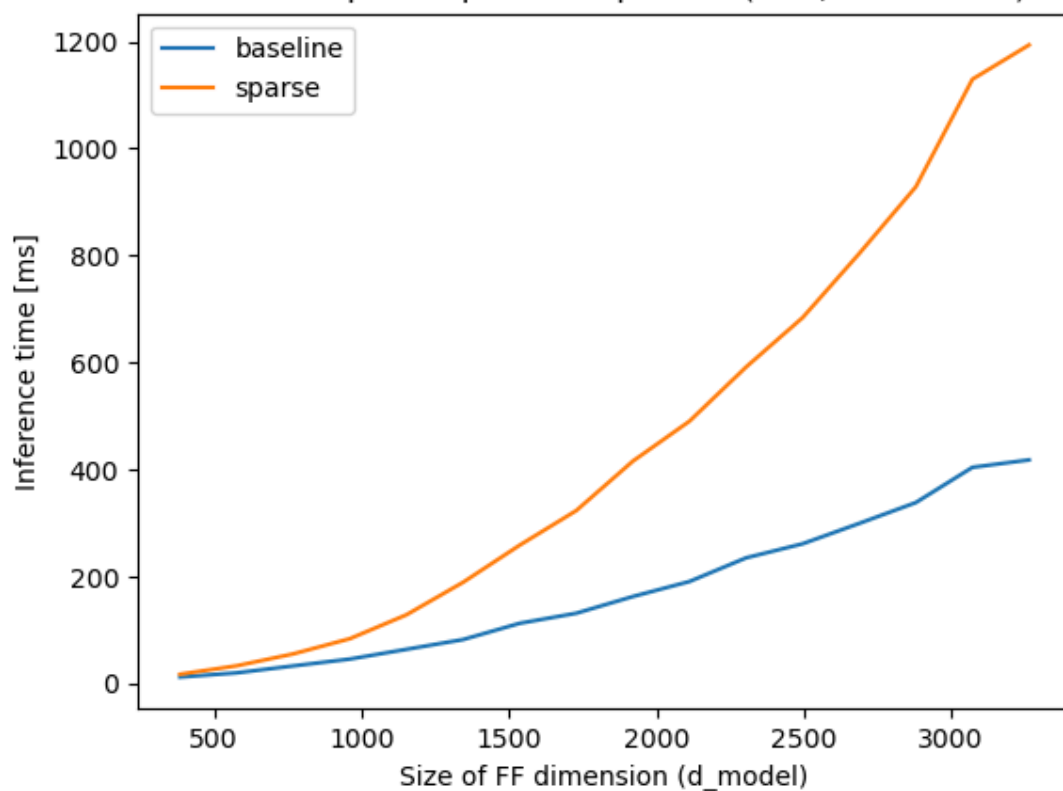
Transformer feed forward  
dense vs sparse speed comparison (CPU, 1-token)

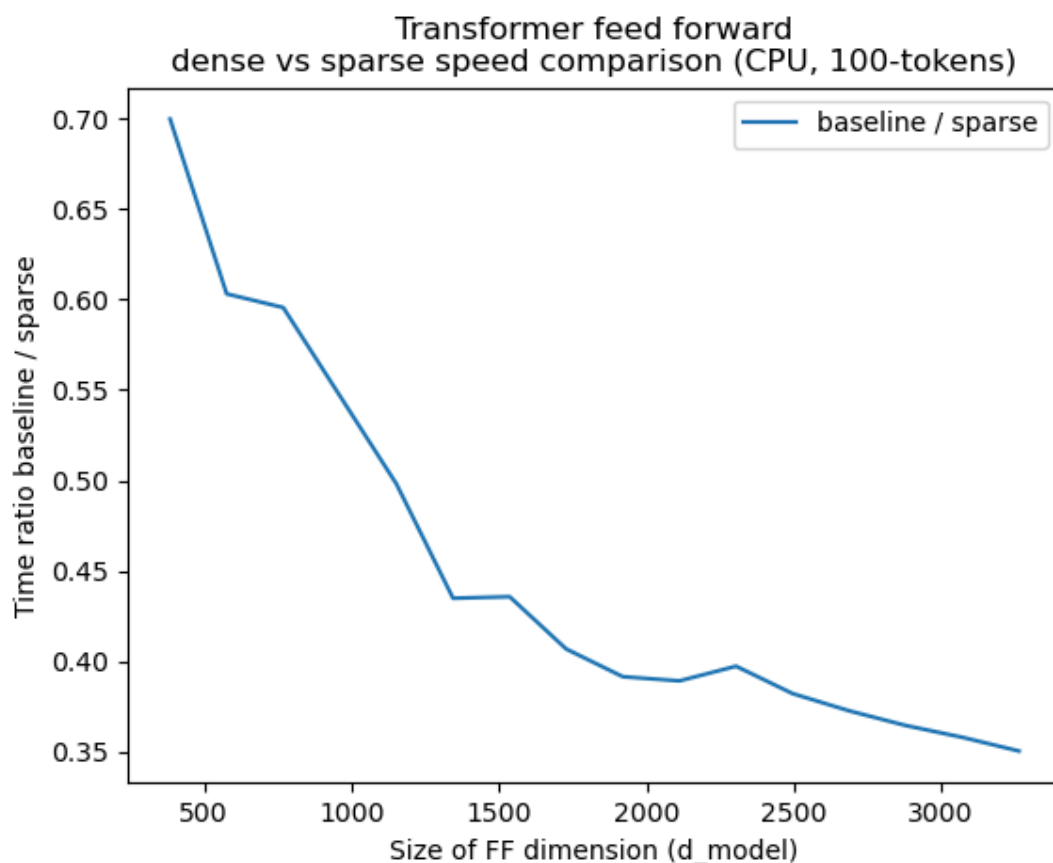


Transformer feed forward  
dense vs sparse speed comparison (CPU, 1-token)



Transformer feed forward  
dense vs sparse speed comparison (CPU, 100-tokens)





## Summary

Sparse transformer model is 10% bigger than the baseline (11M vs 10M params). Training of sparse models is 15.5% slower and it takes significantly more steps to train it properly than baseline.

During inference when loading only 1/64 FF parameters per token, the sparse model is 40% faster when processing a batch of 1 token, and much slower when it needs to process many tokens or batches concurrently. Solutions to that could be to process models input tokens with all parameters loaded, implement caching of previously computed tokens and infer next tokens loading 1/64 FF parameters. It would speed-up generation of new tokens by 40% and does not slow down computation of input tokens.

## Experiments:

### Simple NN

Architecture of the model:

linear(784, d\_model)

feed-forward or sparse FF (d\_model, d\_model\*4)

linear(d\_model, 10)

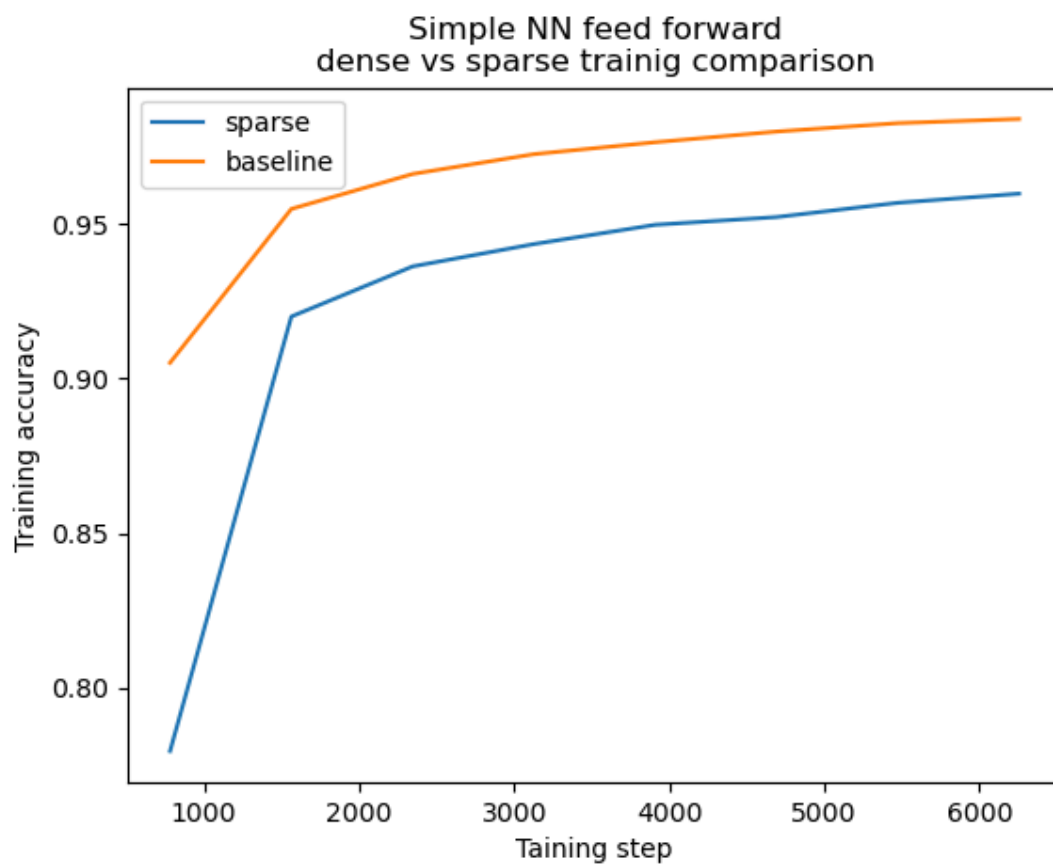
Problem:

Mnist images classification

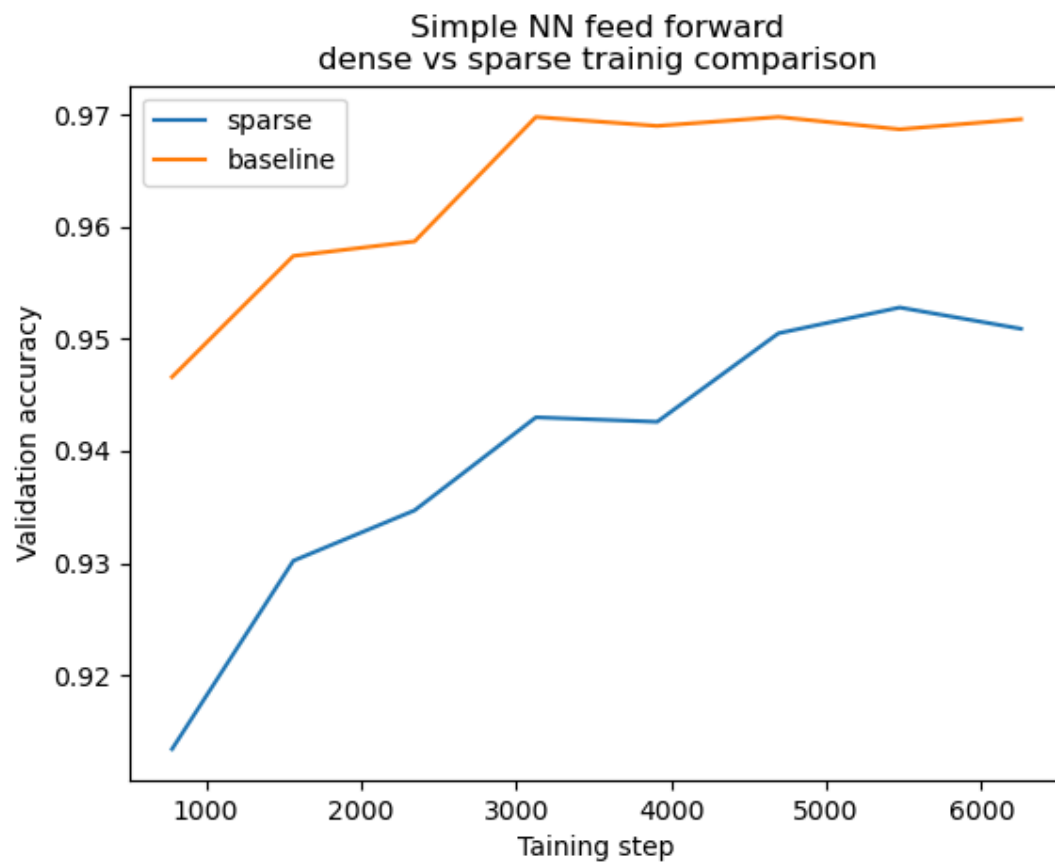
During training (accuracy validation) d\_model of 100 was used

Sparse model uses and trains 1/50 FF layer parameters per sample/image.

Training

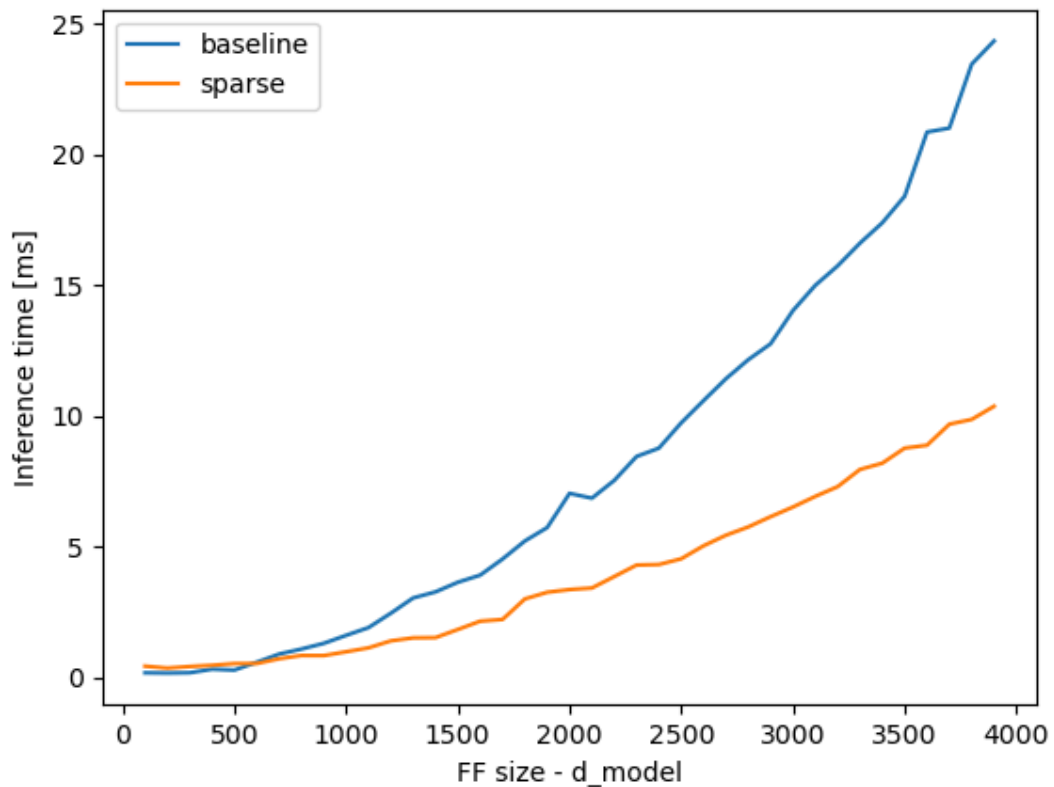




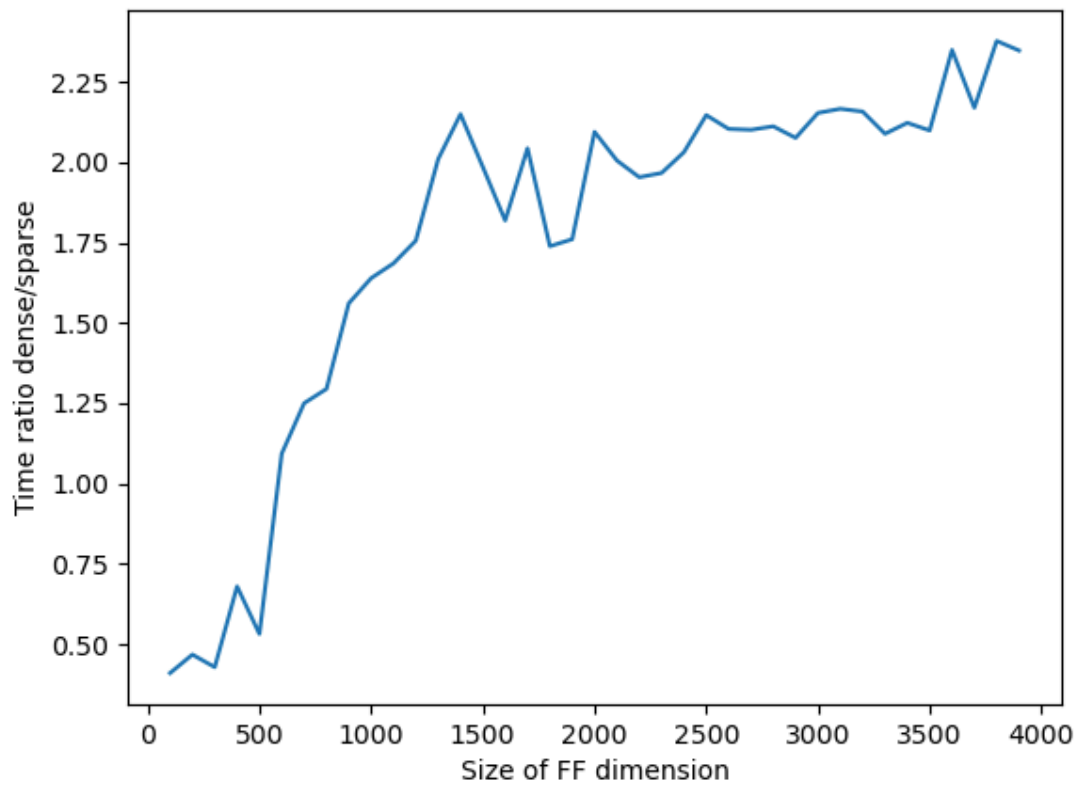


Inference

Feed forward - dense vs sparse speed comparison (CPU, single-batch)



Feed forward - dense vs sparse speed comparison (CPU, single-batch)



## Summary

Using 1/50 FF params speedup of 100% was achieved.

## Conclusion

Sparsity in feed forward layers is a promising method. It could find usage in specific situations, where longer training and slightly less accuracy can be traded for a faster inference. Same sparse FF can be 2x faster than baseline FF during inference.

For a beneficial usage of this method in a transformer architecture, additional implementation of caching processed tokens and distinctive processes for input and generated token processing needs to be made. Sparse FF in the transformer can speed up its generated token computing by 40%.

Quality and speed of o model is highly dependable on sparsity percentage. 2% sparsity provides presented speedups but when tested 4% a speedup was not observed. Sparsity of 10% has a similar quality to the baseline, but is significantly slower.