



# Using BERT for NLU in production





# What is BERT

## Bidirectional Encoder Representations from Transformers

- Novel language representation model
- Fine-tuned BERT models achieved state-of-the-art results on multiple NLP tasks: SQUAD 2.0 question answering, GLUE multi task learning, Google natural questions, etc.:

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

*GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>)*



# Pre-training BERT

- **Unsupervised**, performed on raw text data
- **Very compute-intensive**, takes days on TPUv2, months on GPU

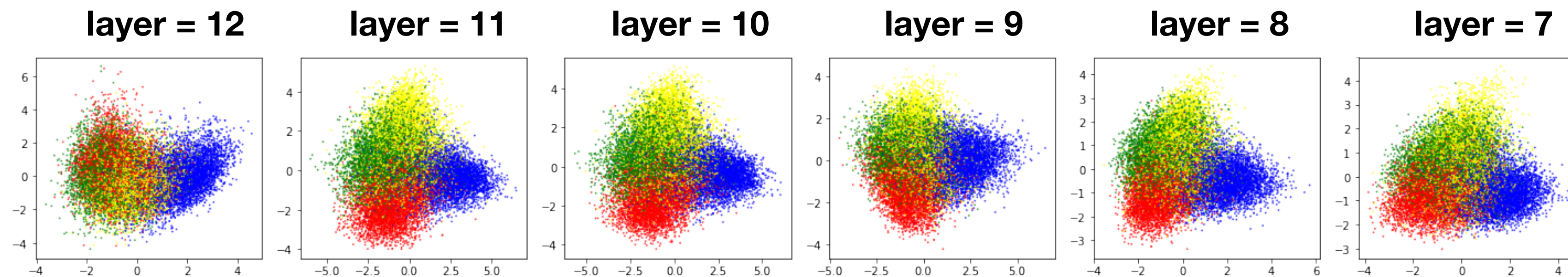


- Pre-trained models are available in open-source
- Guide to pre-training with cloud TPU is available on Medium



# Text representation

- Even without fine-tuning, a pre-trained BERT model can be used to generate text representations (embeddings)
- These representations can be used for token-level or sentence-level tasks
- Encoder layer output shape == [encoder\_dim, seq\_len]
- Sentence embeddings are obtained by pooling encoder layer output



*Mean-pooled representations for samples from UCI-news-aggregator dataset (4 classes)*



# Visualising embeddings

Even without fine-tuning, BERT can be used for efficient text representation



# Fine-tuning with BERT

## ★ End-to-end

The whole model is fine-tuned on a new task

- Superior results on large datasets
- Need **k** BERT models to solve **k** different tasks

## ★ Frozen

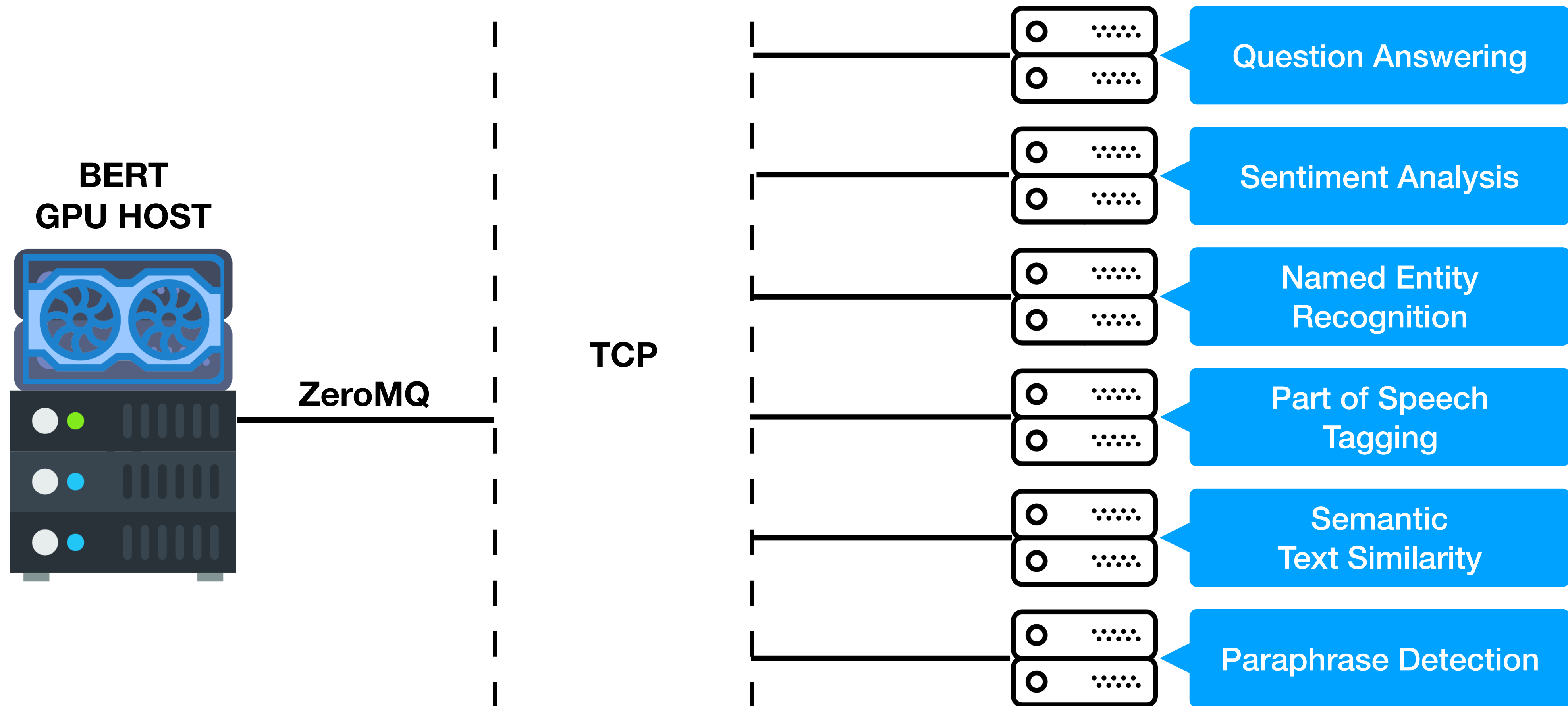
BERT model is fixed (frozen), encoder output used as features

- Good results with smaller datasets
- Need a single BERT + **k** lightweight models to solve **k** different tasks



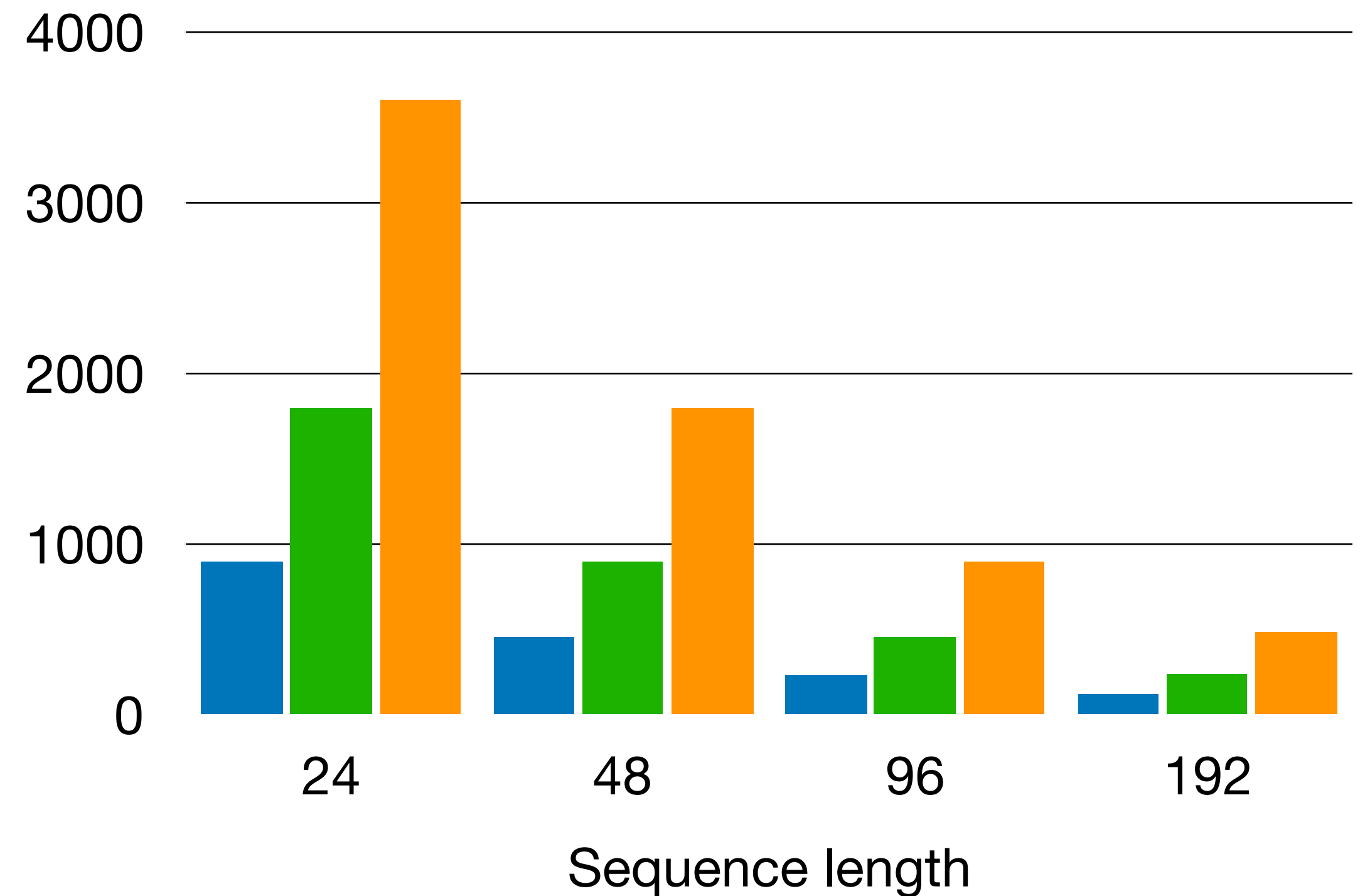
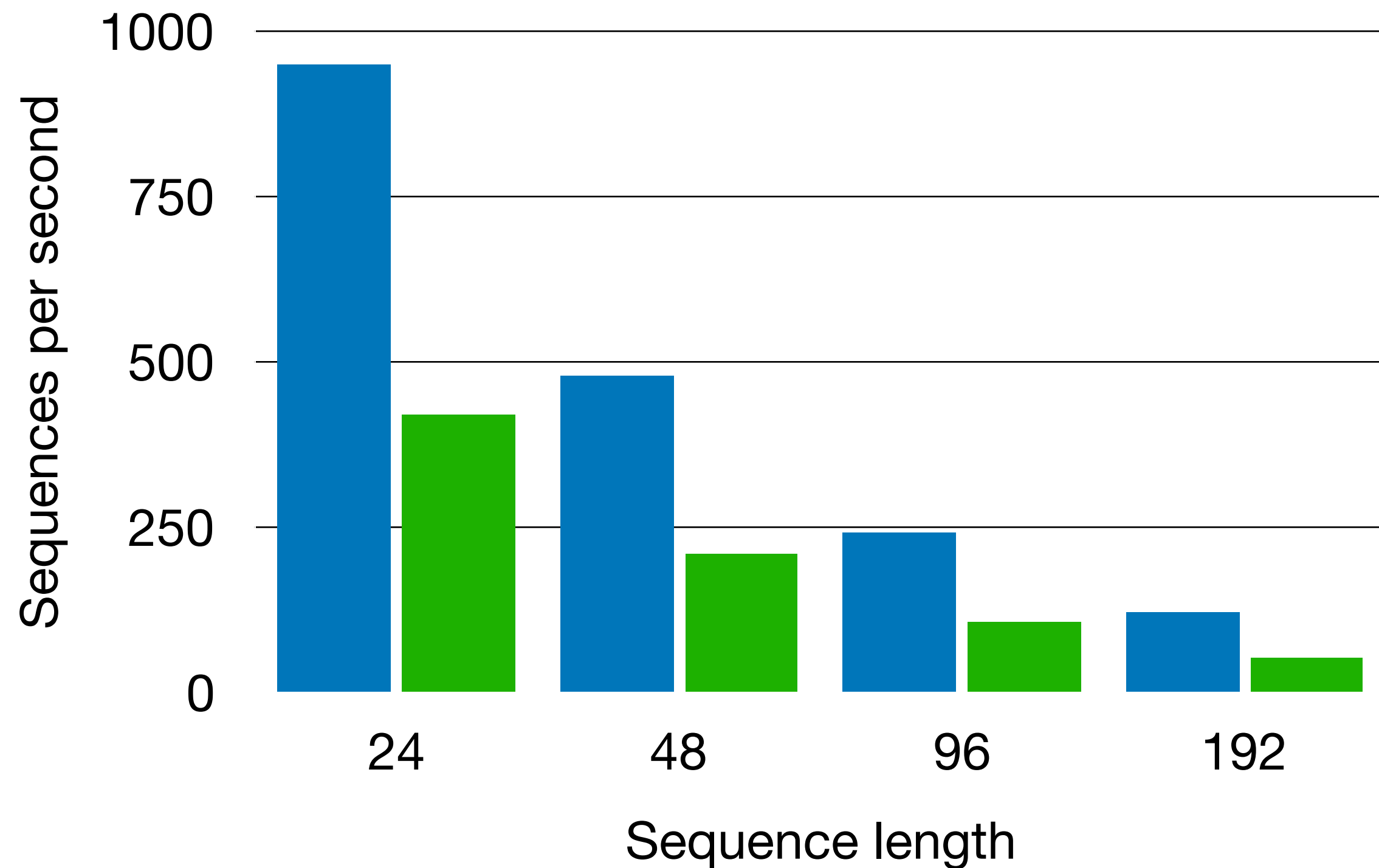
# Serving architecture

A single model can serve multiple clients solving a wide array of NLP tasks





# BERT-base performance on Tesla V100



■ batch\_size=256  
1 worker      ■ batch\_size=1  
5 workers

■ 1 GPU      ■ 2 GPU      ■ 4 GPU







# Service reliability

- Client-server infrastructure introduces additional failure modes
- A standalone solution is required for mission-critical systems
- Local model only requires TensorFlow
- Inference is performed in-process with `tf.Estimator`
- Model graph and weights are serialized and distributed in **Protobuf**





# Info

- [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#)
- [Pre-training BERT from scratch with cloud TPU](#)
- [bert\\_experimental](#)
- [bert-as-a-service](#)
- [Embedding Projector](#)
- [Serving Google-BERT in Production using Tensorflow and ZeroMQ/](#)

**Denis Antyukhov**  
NLP Solutions, Sberbank  
[dantyukhov@icloud.com](mailto:dantyukhov@icloud.com)  
Telegram: @aphex34

