# SELECTIVE PARAMETER UPDATING - MEETING 2
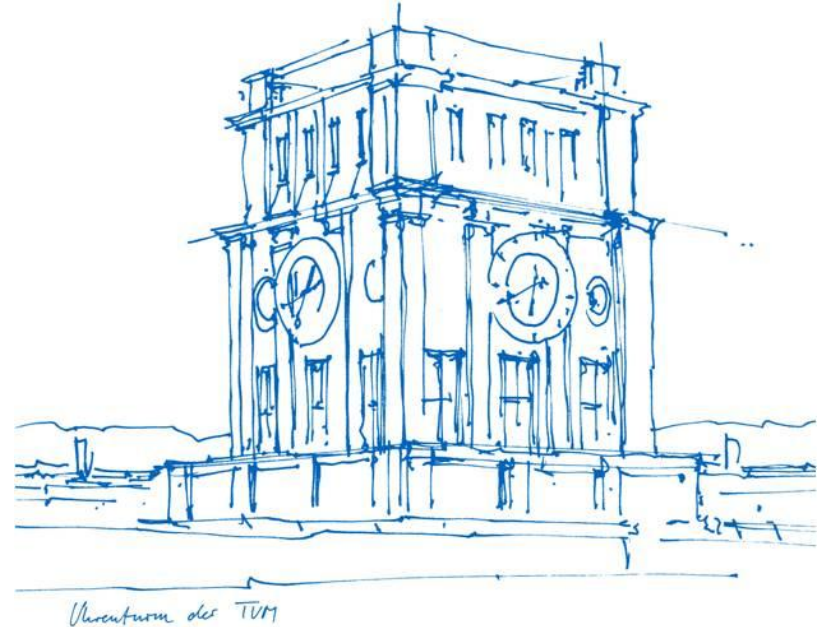
Coşku Barış Coşlu

Mato Gudelj

Baykam Say

17 May 2023

# Contents

# Plan from last meeting…

- Train a full baseline, compare with known numbers to validate training code
- Extend training code with proper logging (TensorBoard)
- Extend training code with an extensible config system
- Brainstorm potential approaches
  - By next meeting shortlist a few and implement at least one
- …keep reading literature

# Plan from last meeting…

- Train a full baseline, compare with known numbers to validate training code
- Extend training code with proper logging (TensorBoard)
- Extend training code with an extensible config system
- Brainstorm potential approaches
  - By next meeting shortlist a few and implement at least one
- …keep reading literature

# Extensible config system

- YAML based
- Encodes all information needed to replicate a training run
- Easy to extend with new datasets
- Adding a new strategy is 1 LOC
    - Currently freezing strategies
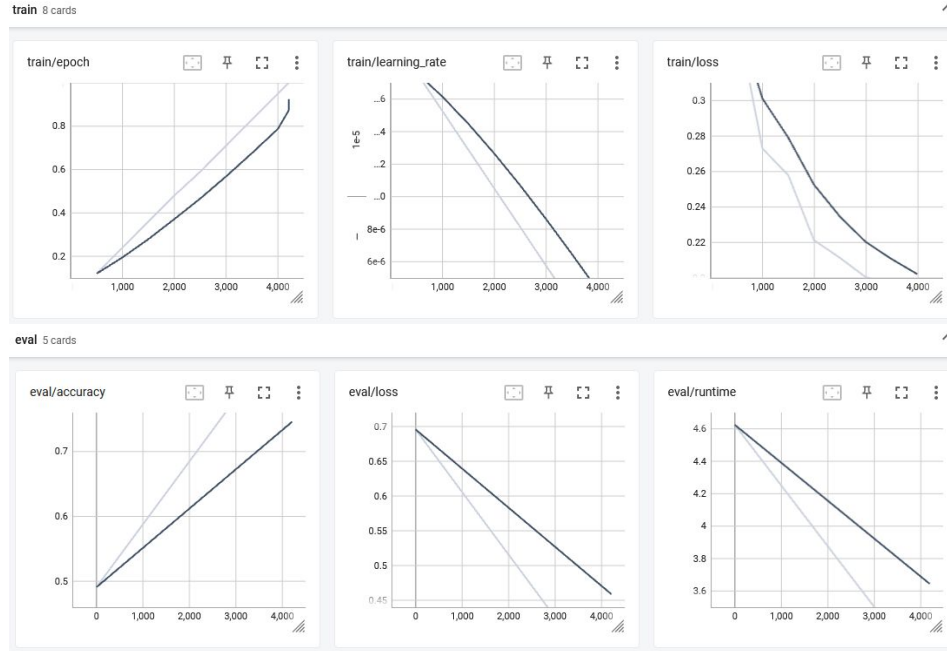    - In the future: additive PEFT, etc.

```yaml
code >  ! example.yaml
 1  %YAML 1.2
 2  ---
 3  dataset:
 4      name: "squad"
 5      n_train: 5000
 6      n_val: 500
 7
 8  model:
 9      base_model: "distilbert-base-cased"
10
11  freeze:
12      strategy: "all_but_last_n"
13      args:
14          n: 1
15
16  train:
17      weight_decay: 0.01
18      num_train_epochs: 5
19      learning_rate: 0.00002
20      per_device_train_batch_size: 16
21      per_device_eval_batch_size: 16
22      output_dir: "results"
23      evaluation_strategy: "epoch"
24
25  evaluate:
26      metric_function: "none"
27  ...
28
```

# Baseline ~~SQuAD~~, **SST-2**

- We had trouble replicating SQuAD results
    - Training takes too much time and memory
    - Original DistilBERT paper does not provide configuration
- This caused a bit of a delay
- Went another route: added support for SST-2
    - Verified pipeline works
    - Replicated results from the paper "Freeze And Reconfigure"
        - Paper result: 91.1%
        - Baseline result: 90.7%

```yaml
code >  ! sst2.yaml
1    %YAML 1.2
2    ---
3    dataset:
4        name: "sst2"
5
6    model:
7        base_model: "distilbert-base-cased"
8
9    freeze:
10       strategy: "none"
11
12   train:
13       weight_decay: 0.01
14       num_train_epochs: 5
15       learning_rate: 0.00002
16       per_device_train_batch_size: 16
17       per_device_eval_batch_size: 16
18       output_dir: "results"
19       evaluation_strategy: "epoch"
20
21   evaluate:
22       metric_function: "accuracy"
23   ...
24
```

# TensorBoard

# Brainstorming potential approaches

- Shortlisted 3 potential approaches
    - Applying a different Kronecker product approximation (KoPA [1]) for PEFT as in KronA [2].
    - Iterate on LST [3]
        - Try different input strategies
        - Try fusing last couple of layers
    - Iterate on (IA)$^3$ [4]
        - Data-dependent weights of K/V/FF vectors, inspired by Involution [5]
- LST interesting due to low VRAM footprint, (IA)$^3$ due to simplicity and high performance
- Due to baseline delays didn't manage to try an approach yet

# Open Issues & Solutions

- Issues
  - T5 model family support needed to compare to LST
    - OOM locally even with the smallest model (Full FT)
  - Not feasible to replicate $(IA)^3$ paper results
    - Very large models (3B)
- Solutions
  - T5/LST
    - Might not be an issue for LST, as it uses significantly less VRAM than Full FT.
    - If problem persists, move to Colab for these experiments
  - $(IA)^3$
    - Find other literature with reliable $(IA)^3$ numbers on smaller models

# Plan for the Next Two Weeks

- Implement LST and/or (IA)$^3$ and start doing experiments
- Try to get T5 running and ideally replicate LST
- Add performance-centric metrics to Tensorboard (memory footprint, forward pass latency, etc.)
- Brainstorm more approaches
- …keep reading literature

# References

[1]    Chencheng Cai, Rong Chen, Han Xiao: "KoPA: Automated Kronecker Product Approximation", 2019; arXiv:1912.02392

[2]    Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Partovi Nia, James J. Clark, Mehdi Rezagholizadeh: "KronA: Parameter Efficient Tuning with Kronecker Adapter", 2022; arXiv:2212.10650

[3]    Yi-Lin Sung, Jaemin Cho, Mohit Bansal: "LST: Ladder Side-Tuning for Parameter and Memory Efficient Transfer Learning", 2022; arXiv:2206.06522

[4]    Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, Colin Raffel: "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning", 2022; arXiv:2205.05638

[5]    Duo Li, Jie Hu, Changhu Wang, Xiangtai Li, Qi She, Lei Zhu, Tong Zhang, Qifeng Chen: "Involution: Inverting the Inherence of Convolution for Visual Recognition", 2021; arXiv:2103.06255