# 1 Question answering task on the SQUADv2 dataset

|  | SQUADv2 (Exact Match) | SQUADv2 (F1) |
|---|---|---|
| Adam | $48.41 \pm 0.57$ | $49.99 \pm 0.54$ |
| M-FAC | $49.80 \pm 0.43$ | $52.18 \pm 0.20$ |

Table 1: Comparing M-FAC optimizer (without weight decay) against HuggingFace's Adam baseline on the **bert-tiny** model.

|  | SQUADv2 (Exact Match) | SQUADv2 (F1) |
|---|---|---|
| Adam | $54.80 \pm 0.47$ | $58.13 \pm 0.31$ |
| M-FAC | $58.02 \pm 0.39$ | $61.35 \pm 0.24$ |

Table 2: Comparing M-FAC optimizer (without weight decay) against HuggingFace's Adam baseline on the **bert-mini** model.

# 2 Text classification on a subset of GLUE tasks

|  | SST-2 (Acc.) | MRPC (F1) | MRPC (Acc.) | STS-B (Pearson) | STS-B (Spearman) |
|---|---|---|---|---|---|
| Adam | $80.11 \pm 0.65$ | $81.68 \pm 0.33$ | $69.90 \pm 0.32$ | $64.39 \pm 5.02$ | $66.52 \pm 5.67$ |
| M-FAC | $81.86 \pm 0.76$ | $82.77 \pm 0.22$ | $72.94 \pm 0.37$ | $80.15 \pm 0.52$ | $80.62 \pm 0.43$ |

|  | QQP (F1) | QQP (Acc.) | MNLI-m (Acc.) | MNLI-mm (Acc.) | QNLI (Acc.) |
|---|---|---|---|---|---|
| Adam | $81.09 \pm 0.15$ | $77.58 \pm 0.08$ | $65.36 \pm 0.13$ | $66.78 \pm 0.15$ | $77.85 \pm 0.15$ |
| M-FAC | $84.29 \pm 0.08$ | $79.71 \pm 0.13$ | $68.28 \pm 3.29$ | $68.98 \pm 3.05$ | $81.17 \pm 0.43$ |

Table 3: Comparing M-FAC optimizer (without weight decay) against HuggingFace's Adam baselines on the **bert-tiny** model.

|  | SST-2 (Acc.) | MRPC (F1) | MRPC (Acc.) | STS-B (Pearson) | STS-B (Spearman) |
|---|---|---|---|---|---|
| Adam | $85.46 \pm 0.58$ | $84.57 \pm 0.36$ | $76.57 \pm 0.80$ | $82.09 \pm 0.54$ | $82.64 \pm 0.71$ |
| M-FAC | $84.20 \pm 0.58$ | $85.06 \pm 1.63$ | $78.87 \pm 2.33$ | $84.66 \pm 0.30$ | $84.65 \pm 0.30$ |

|  | QQP (F1) | QQP (Acc.) | MNLI-m (Acc.) | MNLI-mm (Acc.) | QNLI (Acc.) |
|---|---|---|---|---|---|
| Adam | $86.45 \pm 0.12$ | $82.43 \pm 0.10$ | $73.30 \pm 0.20$ | $74.85 \pm 0.09$ | $83.85 \pm 0.10$ |
| M-FAC | $86.75 \pm 0.20$ | $82.67 \pm 0.23$ | $74.59 \pm 0.41$ | $75.95 \pm 0.14$ | $83.70 \pm 0.13$ |

Table 4: Comparing M-FAC optimizer (without weight decay) against HuggingFace's Adam baselines on the **bert-mini** model.

# 3 Text classification on a subset of GLUE tasks (evaluation on the official test sets)

|  | SST-2 (Acc.) | MRPC (F1) | MRPC (Acc.) | STS-B (Pearson) | STS-B (Spearman) |
|---|---|---|---|---|---|
| AdamW | 83.2 | 81.1 | 71.1 | 74.3 | 73.6 |
| M-FAC | 83.4* | 81.9* | 72.7* | 75.3* | 73.2* |

|  | QQP (F1) | QQP (Acc.) | MNLI-m (Acc.) | MNLI-mm (Acc.) | QNLI (Acc.) |
|---|---|---|---|---|---|
| AdamW | 62.2 | 83.4 | 70.2 | 70.3 | 81.5 |
| M-FAC | 62.8 | 83.9 | 71.0 | 70.5 | 81.7 |

Table 5: Comparing M-FAC optimizer (without weight decay) against authors' (https://github.com/google-research/bert) **tuned bert-tiny** competitive baselines on a subset of GLUE benchmark test sets. * Modest tuning of learning rate and dampening because of an extremely low number of samples (*i.e.* gradients) in the dataset.