

## 1 Question answering task on the SQUADv2 dataset

	SQUADv2 (Exact Match)	SQUADv2 (F1)	SQUADv2 (Train loss)
Adam	48.41 $\pm$ 0.57	49.99 $\pm$ 0.54	2.73 $\pm$ 0.01
M-FAC	49.80 $\pm$ 0.43	52.18 $\pm$ 0.20	2.44 $\pm$ 0.02

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

	SQUADv2 (Exact Match)	SQUADv2 (F1)	SQUADv2 (Train loss)
Adam	54.80 $\pm$ 0.47	58.13 $\pm$ 0.31	1.86 $\pm$ 0.02
M-FAC	58.02 $\pm$ 0.39	61.35 $\pm$ 0.24	1.75 $\pm$ 0.01

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.)	SST-2 (Train loss)	MRPC (F1)	MRPC (Acc.)	MRPC (Train loss)
Adam	80.11 $\pm$ 0.65	0.41 $\pm$ 0.01	81.68 $\pm$ 0.33	69.90 $\pm$ 0.32	0.61 $\pm$ 0.01
M-FAC	81.86 $\pm$ 0.76	0.32 $\pm$ 0.01	82.77 $\pm$ 0.22	72.94 $\pm$ 0.37	0.55 $\pm$ 0.01

	STS-B (Pearson)	STS-B (Spearman)	STS-B (Train loss)	QNLI (Acc.)	QNLI (Train loss)
Adam	64.39 $\pm$ 5.02	66.52 $\pm$ 5.67	4.04 $\pm$ 0.45	77.85 $\pm$ 0.15	0.50 $\pm$ 0.01
M-FAC	80.15 $\pm$ 0.52	80.62 $\pm$ 0.43	1.09 $\pm$ 0.02	81.17 $\pm$ 0.43	0.48 $\pm$ 0.01

	QQP (F1)	QQP (Acc.)	QQP (Train loss)
Adam	77.58 $\pm$ 0.08	81.09 $\pm$ 0.15	0.42 $\pm$ 0.01
M-FAC	79.71 $\pm$ 0.13	84.29 $\pm$ 0.08	0.40 $\pm$ 0.01

	MNLI-m (Acc.)	MNLI-mm (Acc.)	MNLI (Train loss)
Adam	65.36 $\pm$ 0.13	66.78 $\pm$ 0.15	0.85 $\pm$ 0.01
M-FAC	68.28 $\pm$ 3.29	68.98 $\pm$ 3.05	0.81 $\pm$ 0.05

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

	SST-2 (Acc.)	SST-2 (Train loss)	MRPC (F1)	MRPC (Acc.)	MRPC (Train loss)
Adam	85.46 $\pm$ 0.58	0.31 $\pm$ 0.01	84.57 $\pm$ 0.36	76.57 $\pm$ 0.80	0.54 $\pm$ 0.01
M-FAC	84.20 $\pm$ 0.58	0.29 $\pm$ 0.01	85.06 $\pm$ 1.63	78.87 $\pm$ 2.33	0.46 $\pm$ 0.01

  

	STS-B (Pearson)	STS-B (Spearman)	STS-B (Train loss)	QNLI (Acc.)	QNLI (Train loss)
Adam	82.09 $\pm$ 0.54	82.64 $\pm$ 0.71	1.58 $\pm$ 0.10	83.85 $\pm$ 0.10	0.41 $\pm$ 0.01
M-FAC	84.66 $\pm$ 0.30	84.65 $\pm$ 0.30	0.85 $\pm$ 0.03	83.70 $\pm$ 0.13	0.42 $\pm$ 0.01

  

	QQP (F1)	QQP (Acc.)	QQP (Train loss)
Adam	82.43 $\pm$ 0.10	86.45 $\pm$ 0.12	0.34 $\pm$ 0.01
M-FAC	82.67 $\pm$ 0.23	86.75 $\pm$ 0.20	0.35 $\pm$ 0.01

  

	MNLI-m (Acc.)	MNLI-mm (Acc.)	MNLI (Train loss)
Adam	73.30 $\pm$ 0.20	74.85 $\pm$ 0.09	0.70 $\pm$ 0.01
M-FAC	74.59 $\pm$ 0.41	75.95 $\pm$ 0.14	0.68 $\pm$ 0.01

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