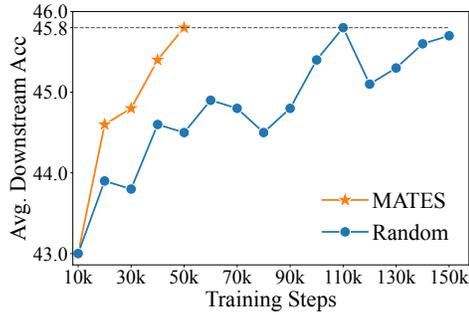


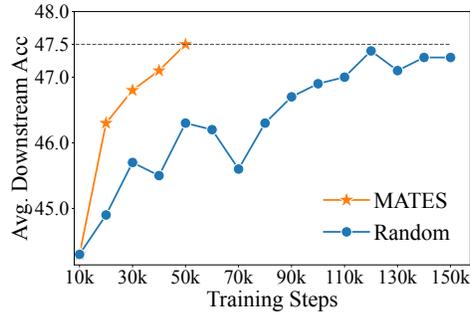
Table 1: Zero-/Two-shot evaluation of pretraining Pythia-1B with different data selection methods.

Methods (#FLOPs *1e19)	SciQ	ARC-E	ARC-C	LogiQA	OBQA
Random (17.67)	65.8 _(1.5) /74.6 _(1.4)	43.7 _(1.0) /45.1 _(1.0)	25.6 _(1.3) /25.5 _(1.3)	27.5 _(1.8) /24.6 _(1.7)	31.8 _(2.1) /30.0 _(2.1)
DSIR (17.67)	65.8 _(1.5) /75.0 _(1.4)	42.6 _(1.0) /44.4 _(1.0)	24.7 _(1.3) /25.8 _(1.3)	28.7 _(1.8) / 25.8 _(1.7)	29.2 _(2.0) / 30.6 _(2.1)
SemDeDup (19.13)	66.8 _(1.5) /75.9 _(1.4)	45.5 _(1.0) /45.1 _(1.0)	25.3 _(1.3) /25.2 _(1.3)	27.6 _(1.8) /24.1 _(1.7)	30.6 _(2.1) / 30.6 _(2.1)
DsDm (22.04)	68.2 _(1.5) / 76.6 _(1.3)	45.0 _(1.0) /45.5 _(1.0)	26.5 _(1.3) /26.5 _(1.3)	26.6 _(1.7) /25.2 _(1.7)	29.4 _(2.0) /30.0 _(2.1)
QuRating (37.67)	67.1 _(1.5) /76.2 _(1.3)	45.5 _(1.0) / 46.7 _(1.0)	25.6 _(1.3) /26.0 _(1.3)	26.9 _(1.7) /24.4 _(1.7)	29.8 _(2.0) /30.2 _(2.1)
MATES (Ours) (19.97)	67.3 _(1.5) /75.7 _(1.4)	44.9 _(1.0) /46.1 _(1.0)	25.9 _(1.3) / 26.8 _(1.3)	28.7 _(1.8) /25.2 _(1.7)	32.2 _(2.1) / 30.6 _(2.1)

Methods (#FLOPs *1e19)	BoolQ	HellaSwag	PIQA	WinoGrande	Average
Random (17.67)	60.2 _(0.9) /56.8 _(0.9)	43.8 _(0.5) /42.9 _(0.5)	68.9 _(1.1) /68.3 _(1.1)	50.7 _(1.4) /52.1 _(1.4)	46.4 _(1.4) /46.7 _(1.3)
DSIR (17.67)	59.7 _(0.9) /55.6 _(0.9)	44.2 _(0.5) /43.6 _(0.5)	68.3 _(1.1) /67.4 _(1.1)	53.2 _(1.4) /52.8 _(1.4)	46.3 _(1.4) /46.8 _(1.3)
SemDeDup (19.13)	60.2 _(0.9) /57.3 _(0.9)	45.3 _(0.5) /44.4 _(0.5)	69.7 _(1.1) /68.4 _(1.1)	52.5 _(1.4) / 53.5 _(1.4)	47.1 _(1.4) /47.2 _(1.3)
DsDm (22.04)	59.0 _(0.9) /56.1 _(0.9)	44.8 _(0.5) /44.6 _(0.5)	68.9 _(1.1) /68.0 _(1.1)	51.9 _(1.4) /52.7 _(1.4)	46.7 _(1.3) /47.2 _(1.3)
QuRating (37.67)	60.3 _(0.9) /58.0 _(0.9)	45.2 _(0.5) /44.7 _(0.5)	70.2 _(1.1) /68.3 _(1.1)	51.6 _(1.4) /51.8 _(1.4)	46.9 _(1.3) /47.4 _(1.3)
MATES (Ours) (19.97)	60.9 _(0.9) / 58.1 _(0.9)	45.3 _(0.5) / 44.8 _(0.5)	69.5 _(1.1) / 68.7 _(1.1)	52.4 _(1.4) /51.0 _(1.4)	47.5 _(1.4) / 47.5 _(1.3)

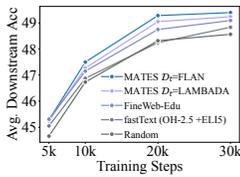


(a) Pythia-410M.

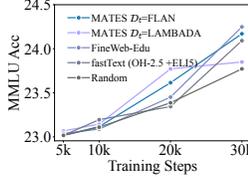


(b) Pythia-1B.

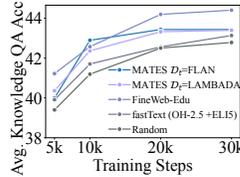
Figure 1: Full data training performances of Pythia-410M and 1B.



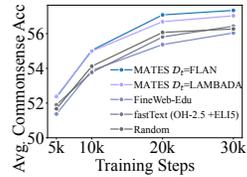
(a) Average acc.



(b) MMLU acc.

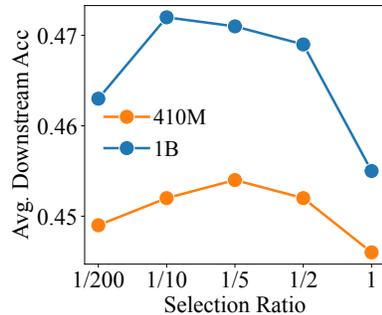


(c) Knowledge QA acc.

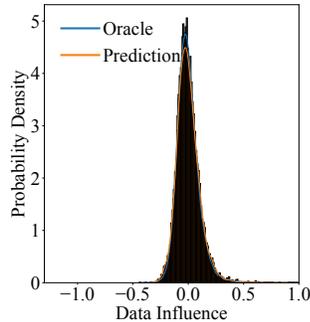


(d) Commonsense acc.

Figure 2: Zero-shot evaluation of pretraining Pythia-1B on FineWeb. We perform different data selection methods starting from 5k steps.



(a) Zero-shot performances of MATES with different data selection ratios at 40k steps. Main setup: 1/5; Random: 1.



(b) Distribution of oracle and predicted data influence. Experiments from Pythia-410M at 40k steps.

Figure 3: (a): Ablation of the selection ratio; (b): Data influence distribution.