

Categories		Methods	Sub-Methods	Related Works
3 Algorithm	3.1 Compression Techniques	3.1.1 Sparsification	3.1.1.1 Weight Pruning	Llm-pruner, Sheared Llama, SliceGPT, Wanda, RIA, FastPruning, oBERT, BESA, SparseGPT
			3.1.1.2 Sparse Attention	SparseFlash, DynamicPruning, SnapKV, PyramidKV, Quest, InfilM, H2O, FastGen, FlexGen, Cissorhands, StreamingLLM, Lm-infinite, Moa
			3.1.2.1 Weight-only	GPTQ, LUT, AWQ, OWQ, SpQR, SqueezeLLM, QuIP, FineQuant, QuantEase, LLM-MQ
			3.1.2.2 Weight-Activation	ZeroQuant series, FlexGen, RPTQ, OliVe, QLLM, Atom, QServe, FlatQuant, SpinQuant. etc.
		3.1.2 Quantization	3.1.2.3 KV Cache	MKLV, WKVQuant, SKVQ, KIVI, RotateKV
			3.1.3 Low-rank Decomposition	LoRD, TensorGPT, LoSparse, LPLR, ASVD, SVD-LLM
	3.2 Meta-Arch. Design	3.2.1 RNN-based		RWKV series
		3.2.2 Mamba		Mamba series, Jamba, Zamba
		3.2.3 TTT-based		TTT, TTTN, TTT++, TTAC, sTTT