

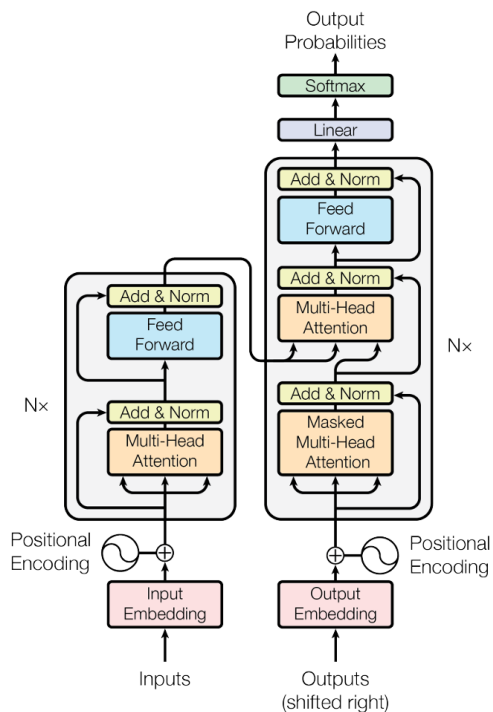
LittleBit: Ultra Low-Bit Quantization via Latent Factorization

Banseok Lee*, Dongkyu Kim*, Youngcheon You, Youngmin Kim†

Samsung Research

*: Equal Contribution, †: Corresponding author
{bs93.lee, dongkyu.k, y01000.you, ym1012.kim}@Samsung.com

Motivation: LLMs on Mobile Devices?

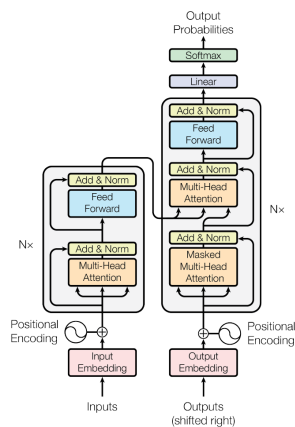


70B LLMs [**138GB**] VRAM



Mobile [**12GB**] VRAM

Motivation: LLMs on Mobile Devices?



1bit?

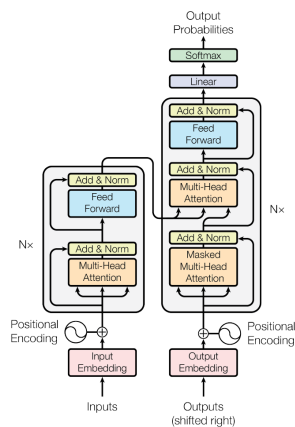


Mobile [12GB] VRAM

LittleBit: Ultra Low-Bit Quantization via Latent Factorization

Samsung Research

Motivation: LLMs on Mobile Devices?



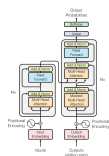
70B 1Bit [14GB] VRAM

Mobile [12GB] VRAM

LittleBit: Ultra Low-Bit Quantization via Latent Factorization

Samsung Research

Motivation: LLMs on Mobile Devices?



70B 0.1Bit [2GB] VRAM



Mobile [12GB] VRAM

LittleBit: Ultra Low-Bit Quantization via Latent Factorization

Samsung Research

Our Approach

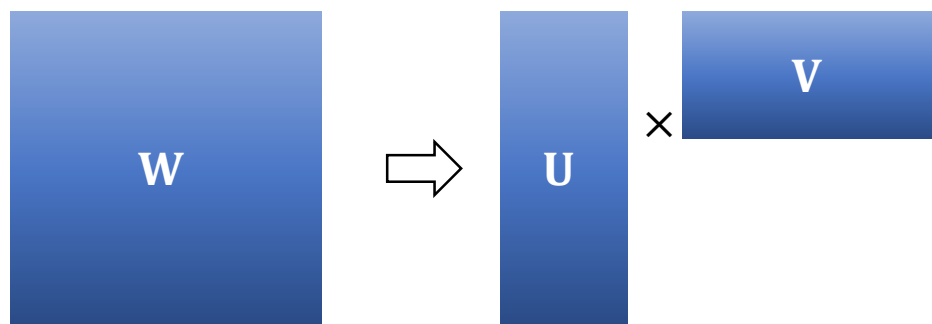


1. Compress
(Latent Factorization)



2. Recover
(Quantization-Aware Training (QAT))

Factorization & Initialization (Dual-SVID)

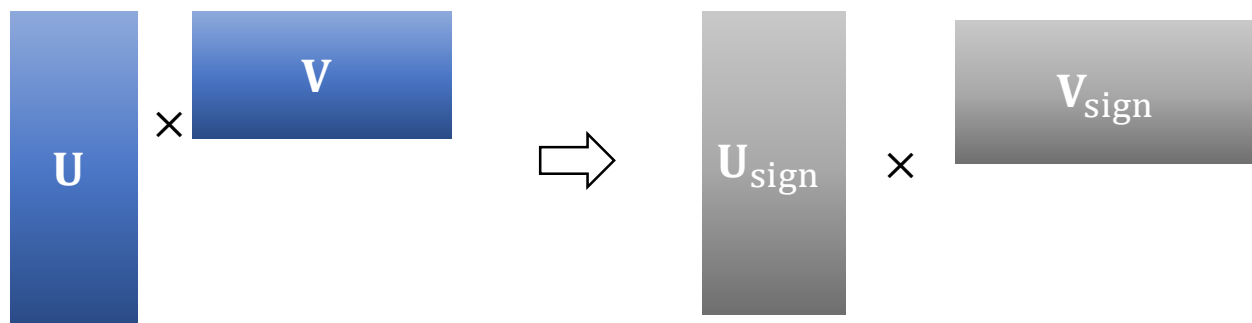


1. Low-rank Factorization (SVD)

 FP16

 Binary(1bit)

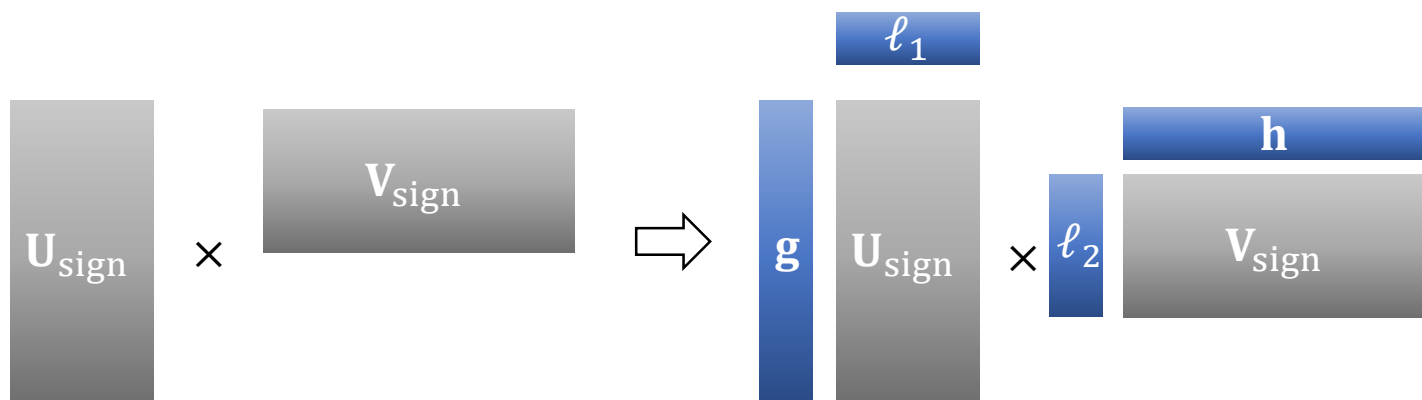
Factorization & Initialization (Dual-SVID)



2. Binarization

■ FP16
■ Binary(1bit)

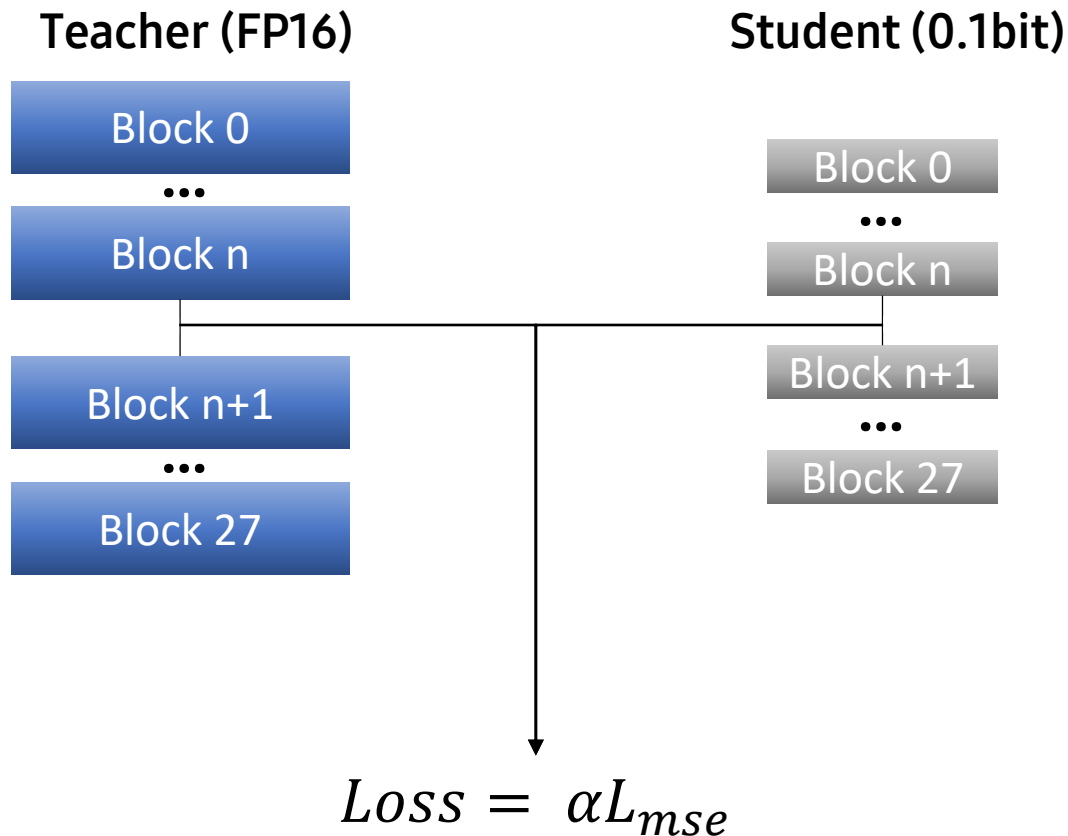
Factorization & Initialization (Dual-SVID)



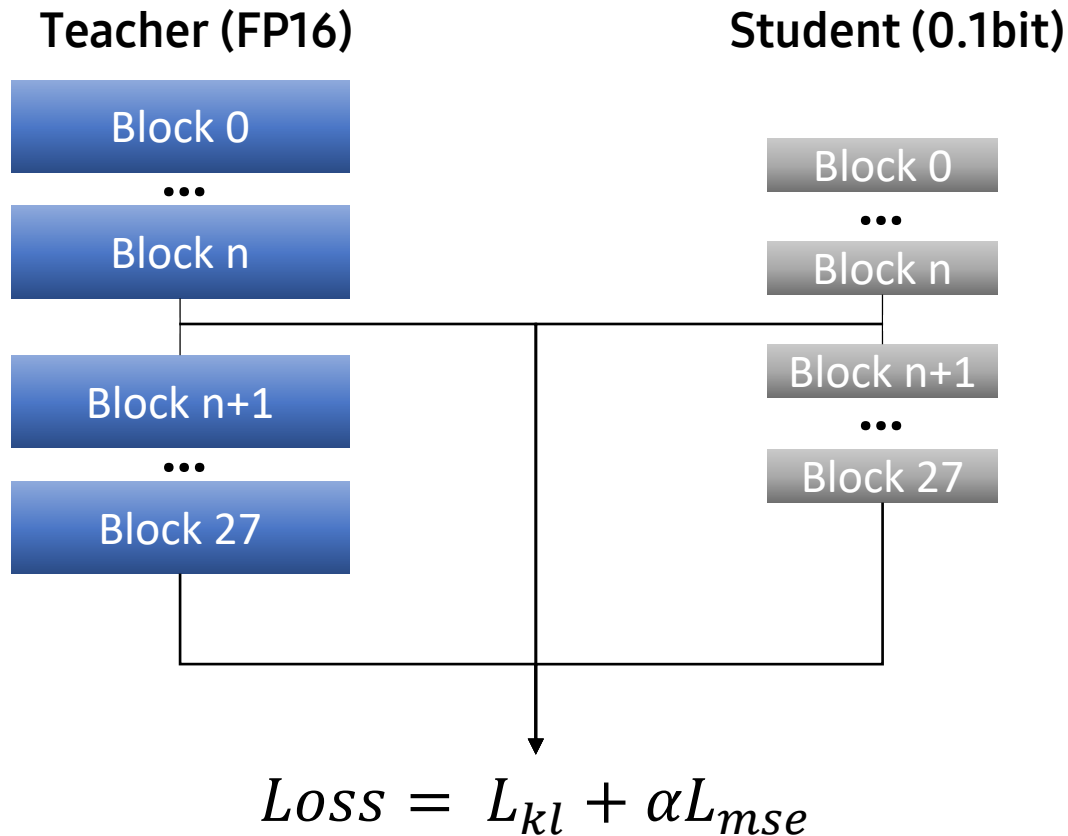
3. Multi-scale Compensation

FP16
Binary(1bit)

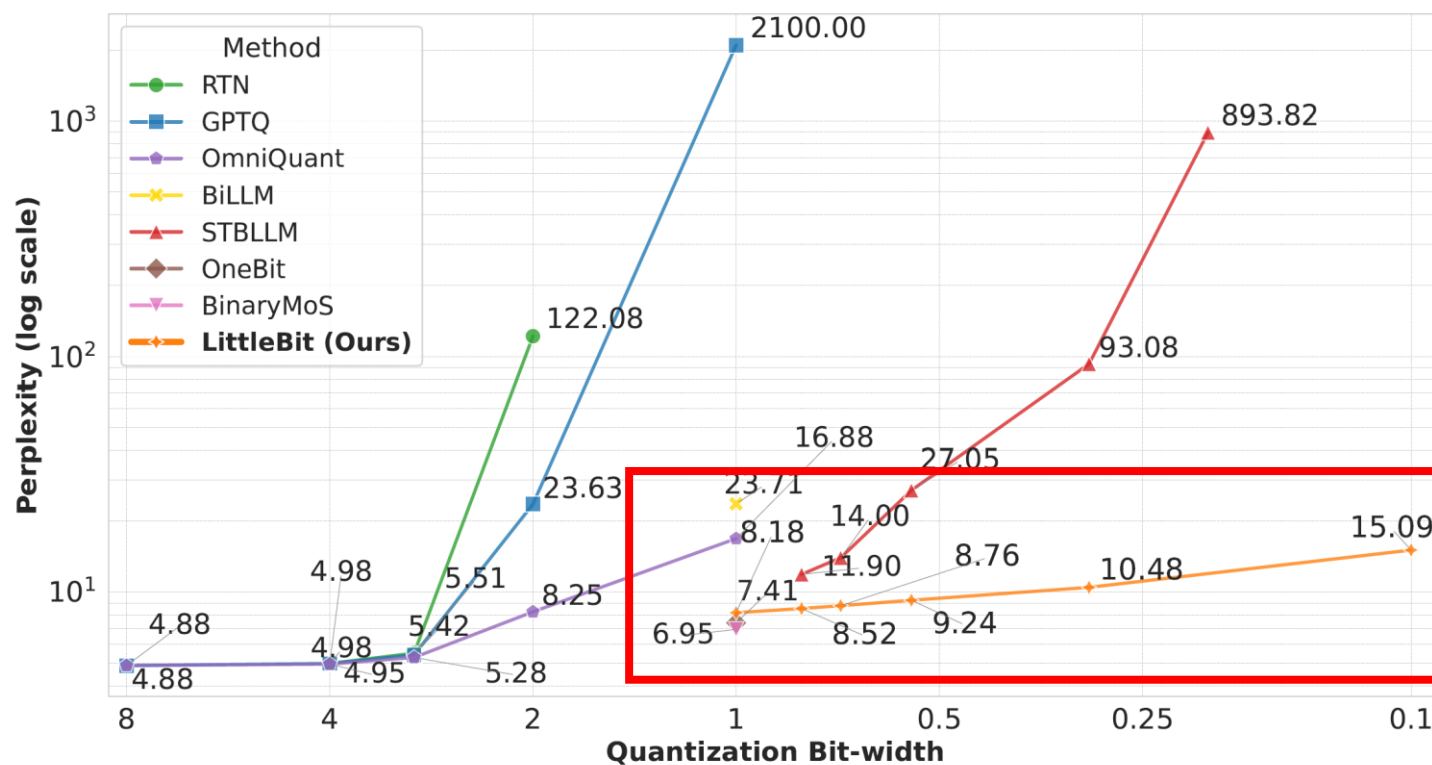
Recovery via QAT (Knowledge Distillation)



Recovery via QAT (Knowledge Distillation)



Result1: SOTA Perplexity



LittleBit: Ultra Low-Bit Quantization via Latent Factorization

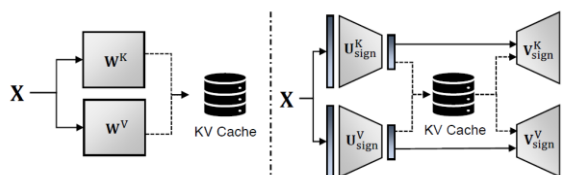
Samsung Research

Result2: Memory Reduction

GPU Memory ↓

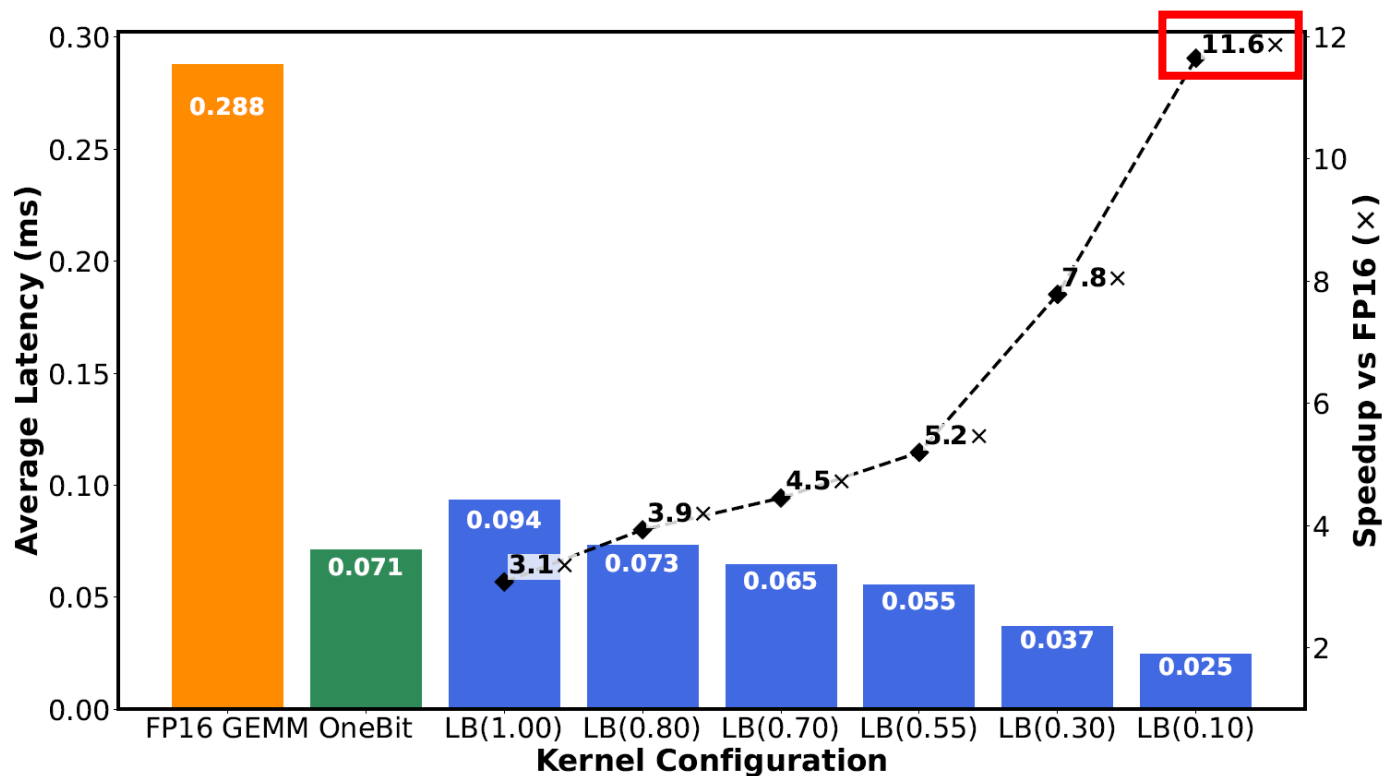
Model	FP16	OneBit [13]	LittleBit (Ours)	
BPW	16	1	0.3	0.1
Llama2-7B	13.49 GB	1.36 GB (9.92×	0.79 GB (17.08×	0.63 GB (21.41×
Llama2-13B	26.06 GB	2.28 GB (11.43×	1.15 GB (22.66×	0.84 GB (31.02×
Llama2-70B	138.04 GB	9.72 GB (14.20×	3.70 GB (37.31×	1.98 GB (69.72×

KV Cache ↓



BPW	KV Latent Rank (r)	KV Cache Reduction Factor
0.80	1,624	$\sim 2.5\times$
0.55	1,112	$\sim 3.7\times$
0.30	600	$\sim 6.8\times$
0.10	192	$\sim 21.3\times$

Result3: Up to 11.6x Kernel Speed-Up



measured on RTX 4090

Conclusion



LittleBit: SOTA Sub-1-bit Quantization

LittleBit: Ultra Low-Bit Quantization via Latent Factorization

Samsung Research