

DUO: No Compromise to Accuracy Degradation

Jinda Jia, Cong Xie*, Fanjiang Ye, Hao Feng, Hanlin Lu, Daoce Wang, Haibin Lin, Zhi Zhang, Xin Liu*



Motivation: Gradient Quantization Results in Accuracy Degradation

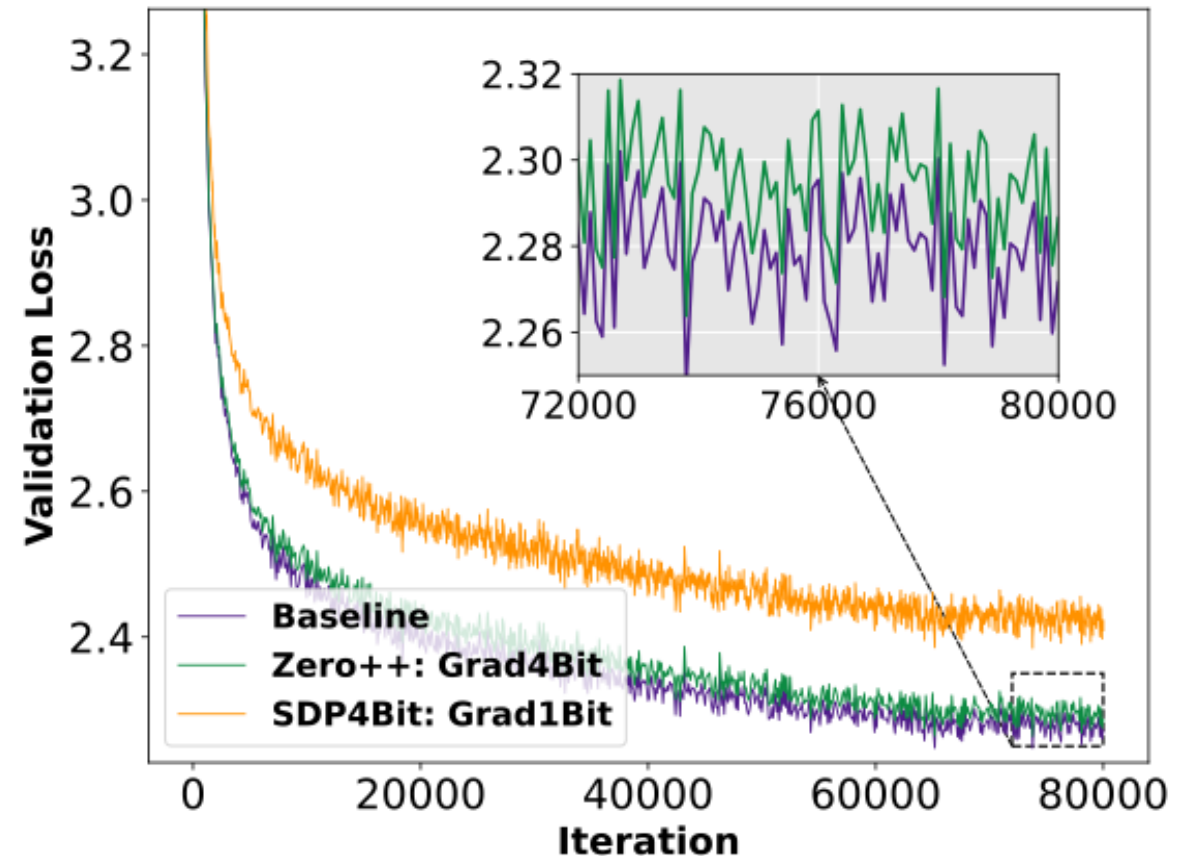
ZeRO++ with 4bit Gradient Quantization:

Intra/Inter Grad: 32 bits \rightarrow 4 bits

SDP4Bit with 1bit Gradient Quantization:

Intra Node Grad: 32 bits \rightarrow 8 bits

Inter Node Grad: 32 bits \rightarrow 1 bits



Loss of ZeRO++, SDP4Bit

**DUO = Fast-Slow (algorithm)
+
Comp-Comm Overlapping
+
Gradient Offload**

Design 1 Algorithm: Fast-Slow Reduction

“Slow” Gradient Reduction:

- 1) High Precision FP32
- 2) Async with forward-backward
- 3) On the CPU

“Fast” Gradient Reduction:

- 1) Low Precision INT1/INT4/INT8...
- 2) Sync before Optimizer Step
- 3) On the GPU

Algorithm 1 Distributed training with **DUO**

Require: parameter weight (main copy): w , weight for forward-backward: \tilde{w} , weight for communication: w' , weight difference: d , gradient: g , worker: p , p th shard of weight or gradient: $[p]$, gradient produced by worker p : g^p , gradients with compression \tilde{g} , averaged gradient: \bar{g} , current iteration: t

```

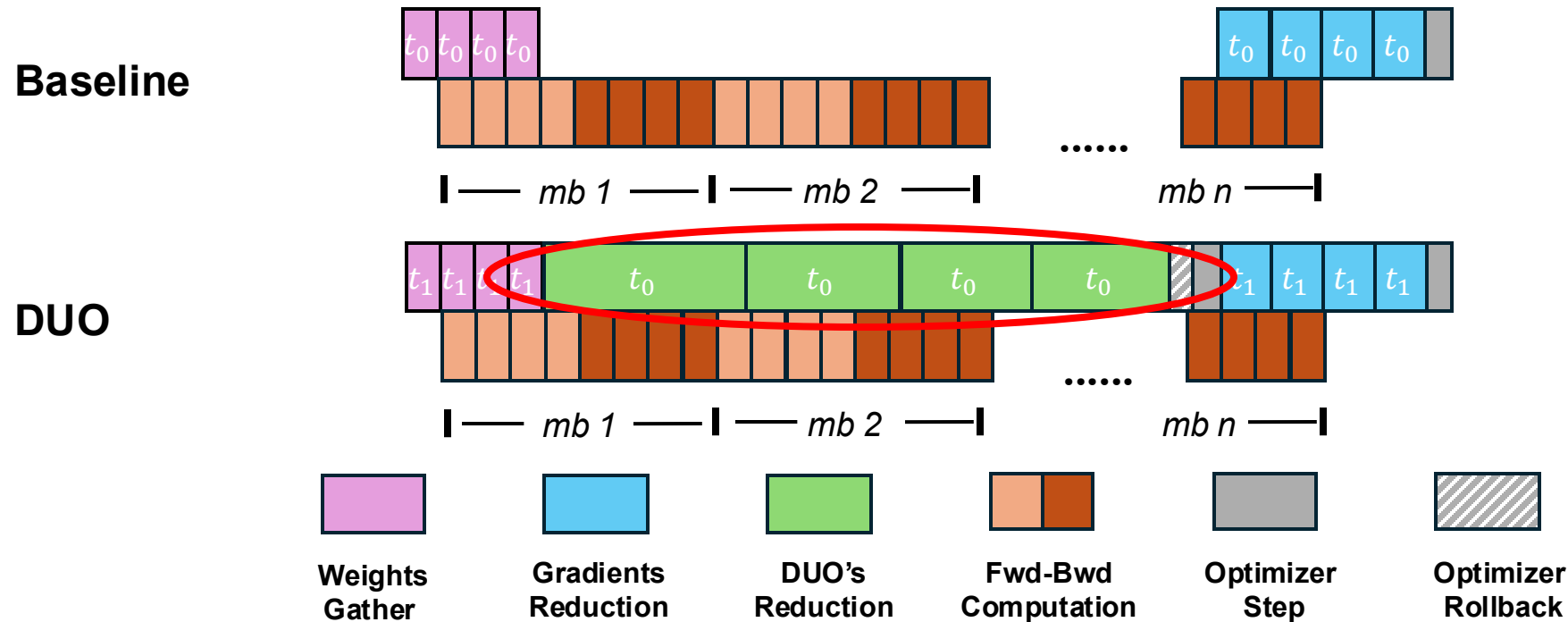
1: function ForwardPass
2:    $d_t[p] = w'_t[p] - \tilde{w}_{t-1}[p]$ 
3:    $\tilde{d}_t[p] \leftarrow \text{QuantizeWeightsDiff}(d_t[p])$ 
4:    $\tilde{d}_t \leftarrow \text{AllGather}(\tilde{d}_t[p])$ 
5:   Start AsyncReduceScatter( $g_{t-1}^p$ )
6:    $\tilde{w}_t \leftarrow \tilde{w}_{t-1} + \tilde{d}_t$ 
7:    $\text{output}^p \leftarrow \text{ForwardPass}(\tilde{w}_t, \text{input}^p)$ 
8: function BackwardPass
9:    $g_t^p \leftarrow \text{Gradient}(\tilde{w}_t, \text{output}_t^p)$ 
10:  SaveReplica( $g_t^p$ )
11:   $\tilde{g}_t^p \leftarrow \text{CompressGradient}(g_t^p)$ 
12:   $\tilde{g}_t[p] \leftarrow \text{ReduceScatter}(\tilde{g}_t^p)$ 
13:   $\bar{g}_{t-1}[p] \leftarrow \text{Wait AsyncReduceScatter Finish}$ 
14:  Recover  $w_{t-1}[p]$ 
15:   $w_t[p] \leftarrow \text{Optimizer}(\bar{g}_{t-1}[p], w_{t-1}[p])$ 
16:   $w'_{t+1}[p] \leftarrow \text{Optimizer}(\tilde{g}_t[p], w_t[p])$ 

```

Design 2 System: Overlapping Communication with Computation

Existing LLM training overlaps only with the first and last micro-batch.

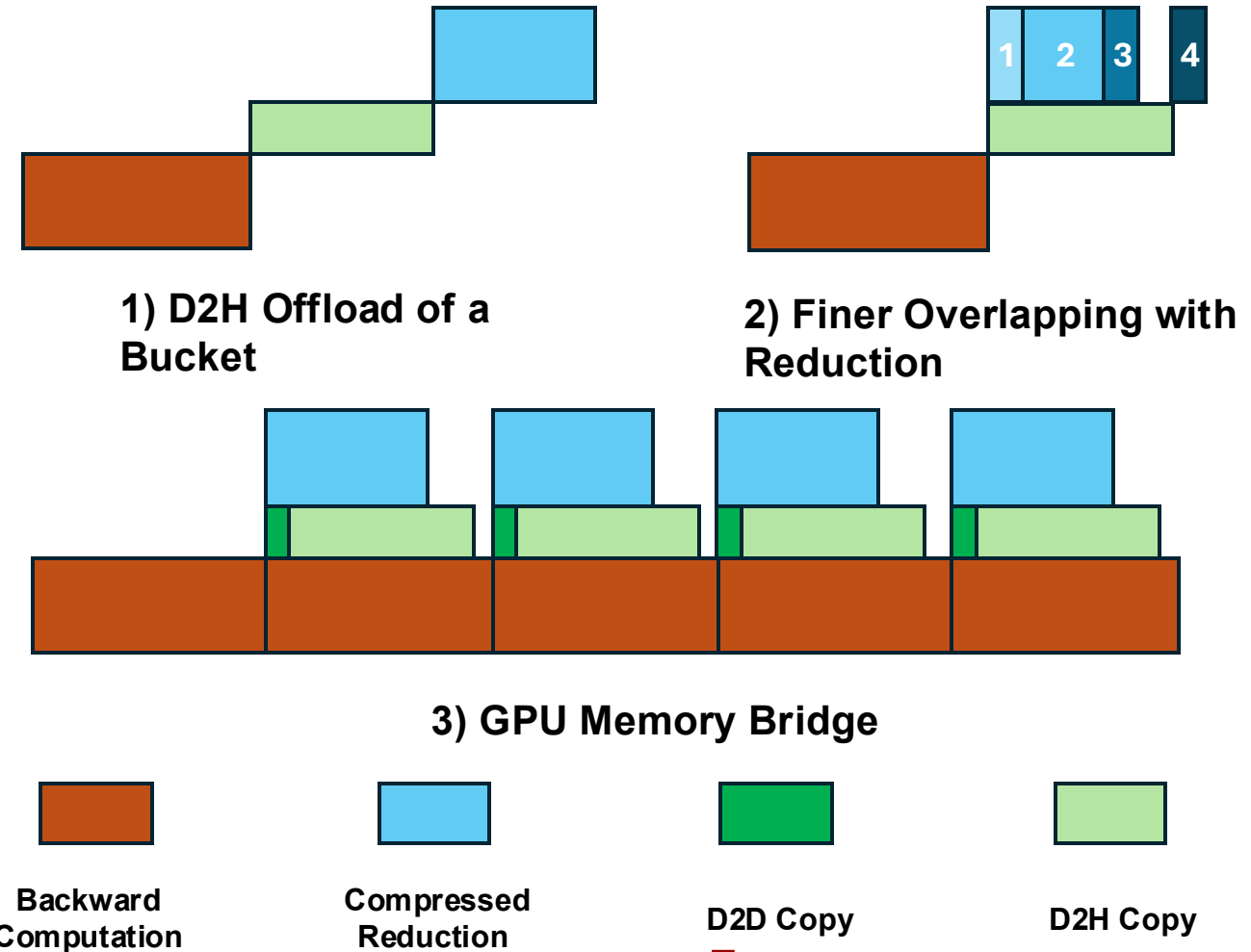
Ours overlaps “slow” reduction with other idle phases.



Design 3 System: Gradient Offload Optimization

DUO requires copying gradients to CPU before starting the “fast” reduction (D2H copy bottleneck).

- 1) Bucket-wise offloading
- 2) Fine-grained D2H overlap
- 3) GPU memory bridge



Experiments

➤ End-to-end Training Loss with DUO

Grad Bits	Strategy	125M	350M	1.3B	6.7B
32	Full Precision	2.2716	2.0582	1.8854	1.7527
4	SDP4Bit DUO	2.2757	2.0629	1.8944	1.7570
		2.2727	2.0592	1.8907	1.7535
1	SDP4Bit DUO	2.4204	2.1843	2.0489	1.8434
		2.2712	2.0686	1.8943	1.7572
0	DUO	2.2761	2.0628	1.8941	1.7550

Final loss is close to the baseline. Even fully remove “fast” reduction (0-bit)

Experiments

➤ E2E Training Throughput with Different Models

Model	Method	H20 (8×8 H20 Nodes)		A100 (8×4 A100 Nodes)	
		Throughput (TFLOPS)	Mem (MB)	Throughput (TFLOPS)	Mem (MB)
1.3B	DUO	70.82 ± 0.63	15880	117.53 ± 0.49	21270
		SDP4Bit	71.02 ± 0.63	15832	117.61 ± 2.35
2.7B	DUO	78.71 ± 0.92	25600	125.31 ± 0.37	35298
		SDP4Bit	79.15 ± 0.80	126.42 ± 0.82	34422
6.7B	DUO	84.09 ± 2.04	16162	122.67 ± 3.76	26480
		SDP4Bit	85.61 ± 1.57	121.05 ± 1.89	25684
13B	DUO	93.02 ± 1.08	27596	132.88 ± 0.95	27332
		SDP4Bit	94.68 ± 0.77	136.73 ± 1.35	27142
18B	DUO	100.85 ± 1.29	36540	116.37 ± 1.28	27378
		SDP4Bit	102.93 ± 1.04	122.91 ± 1.16	27088

DUO achieves similar throughput to SDP4Bit, showing that our high-precision communication adds minimal overhead.

Looking forward to see you on Dec 3



DUO: No Compromise to Accuracy Degradation

Exhibit Hall C,D,E

Wed 3 Dec 4:30 p.m. PST — 7:30 p.m. PST