Supplemental material

Table 1: Compression without sparse encoding. Simplified weights without sparse encoding (CSR and Bloomier encoding) can be compressed for transmission. This table presents compression results for all models considered using only Huffman and arithmetic coding on the pruned and clustered weights.

Model	Pruning Method	Layer	Compression I Huffman	Cactor (Size KB) Arithmetic
LeNet-300-100	Magnitude	FC-0 FC-1	$\begin{array}{ c c c c c }\hline 28.6\times & (32.2) \\ 28.5\times & (4.1) \\ \hline \end{array}$	$44.8 \times (20.5)$ $40.0 \times (2.9)$
	DNS	FC-0 FC-1	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$97.7 \times (9.4)$ $91.0 \times (1.3)$
LeNet5	Magnitude	CNN-1 FC-0	$ \begin{vmatrix} 26.7 \times & (1.4) \\ 27.6 \times & (78.7) \end{vmatrix} $	$30.6 \times (1.2)$ $36.3 \times (59.8)$
	DNS	CNN-1 FC-0	$\begin{array}{c cc} 29.6 \times & (3.3) \\ 31.4 \times & (49.9) \end{array}$	$61.3 \times (1.6)$ $186 \times (8.4)$
VGG-16	Magnitude	FC-0 FC-1	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$79.2 \times (4950)$ $62.4 \times (1027)$

Table 2: **Network compression with CSR encoding.** In the original Deep Compression paper, CSR encoded weights were compressed with Huffman coding. Below are results from applying both Huffman and arithmetic coding to CSR encoded weights for all models considered. This was done to show the relative benefits of different compression techniques independent of the CSR encoding scheme.

Model	Pruning Method	Layer	Compression Factor (Size K CSR Huffman A		Size KB) Arithmetic
LeNet-300-100	Magnitude	FC-0 FC-1	$ \begin{array}{ccc} 40.2 \times & (22.9) \\ 46.8 \times & (2.5) \end{array} $	·	$73.6 \times (12.5)$ $53.2 \times (2.2)$
	DNS	FC-0 FC-1	$ 112 \times (8.2) $ $ 99.2 \times (1.2) $	'	$156 \times (5.9)$ $138 \times (0.85)$
LeNet5	Magnitude	CNN-1 FC-0	$\begin{array}{ c c c c c }\hline 40.4 \times & (0.9 \\ 46.6 \times & (46.6 \end{array}$	·	$34.3 \times (1.1)$ $57.1 \times (38)$
	DNS	CNN-1 FC-0	$\begin{array}{ c c c c c }\hline 90.0 \times & (1.2) \\ 224 \times & (7.0) \\ \hline \end{array}$, ,	$89.1 \times (1.1)$ $347 \times (4.5)$
VGG-16	Magnitude	FC-0 FC-1	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$, ,	$112 \times (3502)$ $83.5 \times (767)$

Table 3: Network compression with Bloomier filter encoding. In Weightless, Bloomier encoded weights were compressed with arithmetic coding. Below are results of applying both Huffman and arithmetic coding to Bloomier encoded weights for all models considered. This was done to show the relative benefits of different compression techniques independent of the Bloomier encoding scheme.

Model	Pruning Method	Layer	Compression Factor (Size KB) Bloomier Huffman Arithmetic				
LeNet-300-100	Magnitude	FC-0	$45.8 \times (20.8)$	$.1) 50.3 \times$	(18.3)	$60.1 \times$	(15.3)
		FC-1	$56.0 \times (2.0)$	$(99) 40.3 \times$	(2.9)	$64.3 \times$	(1.82)
	DNS	FC-0	$152 \times (6.0)$	04) 145×	(6.3)	$174 \times$	(5.27)
		FC-1	$174 \times (0.6)$	$67)$ $125 \times$	(0.9)	$195 \times$	(0.60)
LeNet5	Magnitude	CNN-1	$46.2 \times (0.00)$.8) 31.4×	(1.1)	$51.6 \times$	(0.70)
		FC-0	$62.8 \times (34.8)$.6) $78.4 \times$	(27.9)	$74.2 \times$	(31.1)
	DNS	CNN-1	98× (1.	.2) 73.7×	(1.3)	114×	(0.86)
		FC-0	$445 \times (3.5)$	(52) 427×	(3.7)	$496 \times$	(3.16)
VGG-16	Magnitude	FC-0	$142 \times (275)$	50) 155×	(2530)	$157 \times$	(2500)
		FC-1	$74.6 \times (86)$	80) 82.8×	(774)	$85.8 \times$	(740)

Table 4: Weight reconstruction runtimes. Included in this table are the runtimes for Bloomier weight reconstruction using an Intel i7-6700K desktop CPU and a ARM A53 (600MHz clock) mobile class CPU. All numbers reported use only a single core.

Model	Pruning Method	Layer	Runtime Desktop	(Seconds) Mobile
LeNet-300-100	Magnitude	FC-0 FC-1	0.52 0.066	7.1 0.9
Ecret 900 100	DNS	FC-0 FC-1	0.52 0.067	7.0 0.91
LeNet5	Magnitude	CNN-1 FC-0	0.02	0.28 17.9
Ecross	DNS	CNN-1 FC-0	0.055 0.89	0.76 12.1
VGG-16	Magnitude	FC-0 FC-1	22.8 3.72	296 51.9