Differentiable Sorting Networks for Scalable Sorting and Ranking Supervision

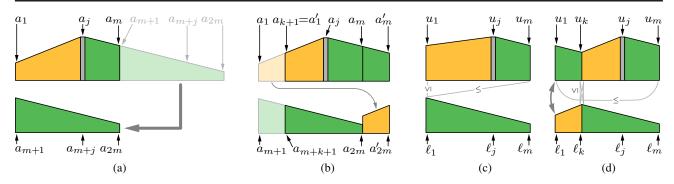


Figure 6: Bitonic merge turns a bitonic input sequence into two bitonic output sequences, with all elements in the one (upper, u_i) sequence larger than all elements in the other (lower, ℓ_i) sequence. The diagrams show the vertical alignment of elements to compare (a) and the invariance to cyclic permutations (b). Depending on the values in (a), no exchanges (c) or exchanges (d) are executed.

A. The Bitonic Sorting Network

In the following, we detail the bitonic sorting network and sketch a proof of why the bitonic sorting networks sorts:

A bitonic sequence is sorted by several bitonic merge blocks, shown in orange and green in Figure 2. Each block takes a bitonic input sequence $a_1a_2 \dots a_{2m}$ of length 2mand turns it into two bitonic output sequences $\ell_1\ell_2 \dots \ell_m$ and $u_1u_2 \dots u_m$ of length m that satisfy $\max_{i=1}^m l_i \leq \min_{i=1}^m u_i$. These subsequences are recursively processed by bitonic merge blocks, until the output sequences are of length 1. At this point, the initial bitonic sequence has been turned into a monotonic sequence due to the minimum/maximum conditions that hold between the output sequences (and thus elements).

A bitonic merge block computes its output as $\ell_i = \min(a_i, a_{m+i})$ and $u_i = \max(a_i, a_{m+i})$. This is depicted in Figure 2 by the arrows pointing from the minimum to the maximum. To demonstrate that bitonic merge works, we show that this operation indeed produces two bitonic output sequences for which the relationship $\max_{i=1}^{m} l_i \leq \min_{i=1}^{m} u_i$ holds.

Note that neither a cyclic permutation of the sequence $(a'_i = a_{(i+k-1 \mod 2m)+1}$ for some k, Figure 6b), nor a reversal, change the bitonic character of the sequence. As can be seen in Figure 6b, even under cyclic permutation still the same pairs of elements are considered for a potential swap. Thus, as a cyclic permutation or a reversal only causes the output sequences to be analogously cyclically permuted or reversed, this changes neither the bitonic character of these sequences nor the relationship between them. Therefore, it suffices to consider the special case shown in Figure 6a, with a monotonically increasing sequence (green) and the maximum element a_j (gray) in the first half. Note that in this case $\forall i; j \leq i \leq m : a_i \geq a_{m+i} \land u_i = a_i \land \ell_i = a_{m+i}$.

For this case, we have to distinguish two sub-cases: $a_1 \ge a_{m+1}$ and $a_1 < a_{m+1}$.

If, on one hand, $a_1 \ge a_{m+1}$, we have the situation shown in Figure 6c: the output sequence $u_1u_2 \ldots u_m$ is simply the first half of the sequence, the output sequence $\ell_1\ell_2 \ldots \ell_m$ is the second half. Thus, both output sequences are bitonic (since they are subsequences of a bitonic input sequence) and $\min_{i=1}^m u_i = \min(u_1, u_m) \ge \ell_1 = \max_{i=1}^m \ell_i$.

If, on the other hand, $a_1 < a_{m+1}$, we can infer $\exists k; 1 \leq k < j : a_k > a_{m+k} \land a_{k+1} \leq a_{m+k+1}$. This situation is depicted in Figure 6d. Thus, $\forall i; 1 \leq i \leq k : u_i = a_{m+i} \land \ell_i = a_i$ and $\forall i; k < i \leq m : u_i = a_i \land \ell_i = a_{m+i}$. Since $u_k = a_{m+k} > a_k = \ell_k, u_k = a_{m+k} \geq a_{m+k+1} = \ell_{k+1}, u_{k+1} = a_{k+1} \geq a_m + k + 1 = \ell_{k+1}, u_{k+1} = a_{k+1} \geq a_k = \ell_k$, we obtain $\max_{i=1}^m l_i \leq \min_{i=1}^m u_i$. Figure 6d shows that the two output sequences are bitonic and that all elements of the upper output sequence are greater than or equal to all elements of the lower output sequence.

B. Implementation Details

B.1. MNIST

For the MNIST based task, we use the same convolutional neural network architecture as in previous works (Grover et al., 2019; Cuturi et al., 2019). That is, two convolutional layers with a kernel size of 5×5 , 32 and 64 channels respectively, each followed by a ReLU and MaxPool layer; after flattening, this is followed by a fully connected layer with a size of 64, a ReLU layer, and a fully connected output layer mapping to a scalar.

B.2. SVHN

For the SVHN task, we use a network with four convolutional layers with a kernel size of 5×5 and (32, 64, 128, 256)filters, each followed by a ReLU and a max-pooling layer with stride 2×2 ; followed by a fully connected layer with size 64, a ReLU, and a layer with output size 1.

B.3. Fast Sort & Rank

To evaluate the fast sorting and ranking method by Blondel et al. (2020), we used the mean-squared-error loss between predicted and ground truth ranks as this method does not produce differentiable permutation matrices.

B.4. Top-k Supervision

For top-k supervision, we use ResNet18 as well as a Vanilla CNN with 4 convolutional and 2 fully connected lay-

ers. The vanilla CNN is has the following architecture: C16-BN-R-C32-BN-R-Max2-C64-BN-R-C128-BN -R-Max2-F256-Fc where Ck denotes a convolutional layer with k output channels, a 3×3 kernel, and padding of 1, BN denotes BatchNorm (Ioffe & Szegedy, 2015), R denotes ReLU, Max2 denotes MaxPool with a 2×2 kernel, and Fk denotes a fully connected layer with k outputs. This vanilla CNN is inspired from Blondel et al. (2020). We train each model using Adam (Kingma & Ba, 2015) for 500 epochs at a learning rate of 10^{-3} .

C. Standard Deviations of the Results

Tables 7, 8, 9, 10, and 11 display the standard deviations for the results in this work.

MNIST	n = 3	n = 5	n = 7	n = 8	n = 15
Fast Sort & Rank	90.6 93.5 73.5 $\pm 0.4 \pm 0.3 \pm 0.8$	$\begin{array}{c c} 71.5 & & 87.2 & & 71.5 \\ \pm 0.9 & & \pm 0.4 & & \pm 0.9 \end{array}$	$\begin{array}{c} 49.7 \mid 81.3 \mid 70.5 \\ \pm 0.6 \mid \pm 0.3 \mid \pm 0.4 \end{array}$	$\begin{array}{c} 29.0 \mid 75.2 \mid 69.2 \\ \pm 1.1 \mid \pm 0.6 \mid \pm 0.7 \end{array}$	$\begin{array}{c c} 2.8 & & 60.9 & & 67.4 \\ \pm 0.2 & & \pm 0.4 & & \pm 0.6 \end{array}$
Odd-Even	$\begin{array}{c} 95.2 \mid 96.7 \mid 86.1 \\ \pm 0.3 \mid \pm 0.2 \mid \pm 0.6 \end{array}$	$\begin{array}{c c} 86.3 & & 93.8 & & 86.3 \\ \pm 0.9 & & \pm 0.4 & & \pm 0.9 \end{array}$	$\begin{array}{c c} 75.4 & 91.2 & 86.4 \\ \pm 1.8 & \pm 0.6 & \pm 0.9 \end{array}$	$\begin{array}{c c} 64.3 & 89.0 & 86.7 \\ \pm 1.8 & \pm 0.6 & \pm 1.1 \end{array}$	$\begin{array}{c} 35.4 \mid 83.7 \mid 87.6 \\ \pm 1.8 \mid \pm 0.5 \mid \pm 0.5 \end{array}$
MNIST	n = 2	n = 4	n = 8	n = 16	n = 32
Odd-Even	$\begin{array}{c} 98.1 \mid 98.1 \mid 84.3 \\ \pm 0.3 \mid \pm 0.3 \mid \pm 0.9 \end{array}$	90.5 94.9 85.5 $\pm 1.2 \mid \pm 0.6 \mid \pm 1.5$	$\begin{array}{c} 63.6 \mid 87.9 \mid 83.6 \\ \pm 11.6 \mid \pm 4.2 \mid \pm 6.1 \end{array}$	$\begin{array}{c} 31.7 \mid 82.8 \mid 87.3 \\ \pm 1.5 \mid \pm 0.5 \mid \pm 0.5 \end{array}$	$\begin{array}{c c} 1.7 & 69.1 & 86.7 \\ \pm 0.5 & \pm 1.5 & \pm 1.0 \end{array}$
Bitonic	98.1 98.1 84.0 $\pm 0.2 \pm 0.2 \pm 1.2$	91.4 95.3 86.7 $\pm 0.6 \pm 0.3 \pm 0.4$	$70.6 90.3 86.9 \\ \pm 4.4 \pm 1.3 \pm 1.8$	$\begin{array}{c c} 30.5 & & 81.7 & & 86.6 \\ \pm 1.8 & & \pm 1.2 & & \pm 0.9 \end{array}$	$\begin{array}{c c} 2.7 & 67.3 & 85.4 \\ \pm 1.3 & \pm 2.7 & \pm 1.7 \end{array}$

Table 7: Same as Table 1 but with additional standard deviations.

Table 8: Same as Table 2 but with additional standard deviations.

SVHN	n = 2	n = 4	n = 8	n = 16	n = 32
Det. NeuralSort	90.1 90.1 39.9 $\pm 0.7 \pm 0.7 \pm 1.7$	$\begin{array}{c c} 61.4 & 78.1 & 45.4 \\ \pm 0.8 & \pm 0.3 & \pm 1.2 \end{array}$	$\begin{array}{c c} 15.7 & 62.3 & 48.5 \\ \pm 1.6 & \pm 1.2 & \pm 1.6 \end{array}$	$\begin{array}{c} 0.1 \mid 45.7 \mid 51.0 \\ \pm 0.1 \mid \pm 0.6 \mid \pm 1.2 \end{array}$	$\begin{array}{c c} 0.0 & & 29.9 & & 52.7 \\ \pm 0.0 & & \pm 1.4 & & \pm 1.5 \end{array}$
Optimal Transport	$\begin{array}{c} 85.5 \mid 85.5 \mid 25.9 \\ \pm 0.0 \mid \pm 0.0 \mid \pm 0.0 \end{array}$	$\begin{array}{c} 57.6 \mid 75.6 \mid 41.6 \\ \pm 1.1 \mid \ \pm \ 0.8 \mid \ \pm \ 1.8 \end{array}$	$\begin{array}{c} 19.9 \mid 64.5 \mid 51.7 \\ \pm 1.9 \mid \pm 1.1 \mid \pm 1.2 \end{array}$	$\begin{array}{c c} 0.3 & & 47.7 & & 53.8 \\ \pm 0.2 & & \pm 1.7 & & \pm 1.4 \end{array}$	$\begin{array}{c c} 0.0 & & 29.4 & & 53.3 \\ \pm 0.0 & & \pm 1.0 & & \pm 1.9 \end{array}$
Fast Sort & Rank	$\begin{array}{c} 93.4 \mid 93.4 \mid 57.6 \\ \pm 0.7 \mid \pm 0.7 \mid \pm 3.7 \end{array}$	$\begin{array}{c} 58.0 \mid 75.8 \mid 41.5 \\ \pm 1.1 \mid \ \pm \ 0.7 \mid \ \pm \ 1.0 \end{array}$	$\begin{array}{c} 8.6 \mid 52.7 \mid 34.4 \\ \pm 1.0 \mid \pm 0.6 \mid \pm 0.3 \end{array}$	$\begin{array}{c c} 0.3 & & 36.5 & & 41.6 \\ \pm 0.2 & & \pm 1.4 & & \pm 1.8 \end{array}$	$\begin{array}{c c} 0.0 & & 14.0 & & 27.5 \\ \pm 0.0 & & \pm 3.1 & & \pm 9.1 \end{array}$
Odd-Even	$\begin{array}{c} 93.4 \mid 93.4 \mid 58.0 \\ \pm 0.4 \mid \pm 0.4 \mid \pm 2.0 \end{array}$	$74.8 \mid 85.5 \mid 62.6 \\ \pm 1.2 \mid \pm 0.7 \mid \pm 1.1$	$\begin{array}{c} 35.2 \mid 73.5 \mid 63.9 \\ \pm 1.2 \mid \pm 0.5 \mid \pm 1.1 \end{array}$	$\begin{array}{c c} 1.8 & 54.4 & 62.3 \\ \pm 0.8 & \pm 1.6 & \pm 1.6 \end{array}$	$\begin{array}{c c} 0.0 & & 36.6 & & 62.6 \\ \pm 0.0 & & \pm 1.5 & & \pm 0.8 \end{array}$
Bitonic	$\begin{array}{c} 93.8 \mid 93.8 \mid 58.6 \\ \pm 0.3 \mid \pm 0.3 \mid \pm 0.8 \end{array}$	$\begin{array}{c} 74.4 \mid 85.3 \mid 62.1 \\ \pm 0.7 \mid \pm 0.3 \mid \pm 1.1 \end{array}$	$\begin{array}{c} 38.3 \mid 75.1 \mid 66.8 \\ \pm 2.4 \mid \pm 1.1 \mid \pm 1.4 \end{array}$	$\begin{array}{c c} 3.9 & 59.6 & 66.8 \\ \pm 0.3 & \pm 0.8 & \pm 1.4 \end{array}$	$\begin{array}{c c} 0.0 & & 42.4 & & 67.7 \\ \pm 0.0 & & \pm 3.5 & & \pm 3.6 \end{array}$

λ	0.25	0.4	0.4	0.4	0.4	0.4	0.4	0.4	0.4	0.4	0.4	0.4	0.4
n	32	32	64	128	256	512	1024	32	64	128	256	512	1024
batch size	128	128	64	32	16	8	4	4	4	4	4	4	4
s = 30	78.20 ± 2.35	$ \begin{array}{c} 79.89 \\ \pm 1.97 \end{array}$	81.25 ± 1.93	82.50 ± 1.09	82.05 ± 2.62	82.50 ± 1.75	82.80 ± 2.27	$ \begin{array}{c} 71.08 \\ \pm 1.67 \end{array} $	75.88 ± 2.30	79.43 ± 2.35	81.46 ± 1.47	82.98 ± 2.02	82.80 ± 2.27
s = 32.5	76.98 ± 0.86	$\begin{array}{ c c } 79.62 \\ \pm 3.62 \end{array}$	81.66 ± 2.42	80.15 ± 3.84	81.87 ± 2.19	82.64 ± 1.60	81.63 ± 6.22	$\begin{vmatrix} 72.31 \\ \pm 2.04 \end{vmatrix}$	75.59 ± 2.05	79.71 ± 1.57	81.36 ± 1.98	82.99 ± 1.67	81.63 ± 6.22
s = 35	77.45 ± 1.64	$ \begin{array}{c} 80.93 \\ \pm 2.75 \end{array}$	81.26 ± 2.41		81.42 ± 2.09	81.51 ± 2.12	81.15 ± 3.12	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	75.73 ± 2.46	78.81 ± 1.36	79.32 ± 4.85	82.30 ± 1.22	81.15 ± 3.12
s = 37.5	76.40 ± 3.90	$ \begin{array}{c} 80.02\\ \pm 1.74 \end{array}$		81.50 ± 2.03	80.05 ± 3.94	82.67 ± 2.21	$\begin{array}{c} 80.07 \\ \pm 3.67 \end{array}$	$ \begin{array}{c} 70.69 \\ \pm 2.26 \end{array} $	75.80 ± 1.22	79.11 ± 1.88	80.64 ± 2.18	82.70 ± 1.66	$\begin{array}{c} 80.07 \\ \pm 3.67 \end{array}$
s = 40	77.69 ± 1.54	$ \begin{array}{c} 80.97 \\ \pm 2.03 \end{array}$	80.23 ± 3.51	81.55 ± 1.97	79.75 ± 5.41	81.89 ± 2.51	81.15 ± 3.31	$ \begin{array}{c} 70.20 \\ \pm 2.06 \end{array} $	74.67 ± 2.45	78.14 ± 2.49		81.39 ± 1.67	81.15 ± 3.31
mean	77.35 ± 2.06	$ \begin{array}{c} 80.29 \\ \pm 2.48 \end{array}$	80.89 ± 2.48	81.28 ± 2.80	81.03 ± 3.48	82.24 ± 2.03	81.36 ± 3.97	$ \begin{array}{c} 71.09 \\ \pm 2.00 \end{array} $	75.53 ± 2.10	79.04 ± 1.97	80.57 ± 2.77	82.47 ± 1.71	81.36 ± 3.97
best s	78.20 ± 3.90	$ \begin{array}{c} 80.97 \\ \pm 3.62 \end{array}$	81.66 ± 3.51	82.50 ± 3.89	82.05 ± 5.41	82.67 ± 2.51	82.80 ± 6.22	$\begin{vmatrix} 72.31 \\ \pm 2.26 \end{vmatrix}$	75.88 ± 2.46	79.71 ± 2.49	81.46 ± 4.85	$82.99 \\ \pm 2.02$	82.80 ± 6.22
worst s	$\begin{array}{c} 76.40 \\ \pm 0.86 \end{array}$	$ \begin{array}{c} 79.62 \\ \pm 1.74 \end{array}$			79.75 ± 2.09	81.51 ± 1.60	80.07 ± 2.27	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	74.67 ± 1.22	78.14 ± 1.36	79.32 ± 1.47	81.39 ± 1.22	80.07 ± 2.27

Table 9: Same as Table 3 but with additional standard deviations.

Table 10: Same as Table 5 but with additional standard deviations.

	<i>n</i> =	n = 4		32
Setting / λ	0	0.25	0	0.25
Odd-Even (MNIST) Bitonic (MNIST)	$\begin{array}{c} 94.5 \pm 0.3 \\ 93.6 \pm 1.4 \end{array}$	94.9 ± 0.6 95.3 ± 0.3	61.5 ± 1.9 62.8 ± 15.5	$69.1 \pm 1.5 \\ 67.3 \pm 2.7$
Odd-Even (SVHN) Bitonic (SVHN)	77.3 ± 1.0 78.1 ± 0.2	85.5 ± 0.7 85.3 ± 0.3	28.5 ± 2.7 35.0 ± 0.8	$36.6 \pm 1.5 \\ 42.4 \pm 3.5$

Table 11: Same as Table 6 but with additional standard deviations.

Setting	Softmax CE	Diff. Top-k
CIFAR-10, Vanilla CNN CIFAR-10, ResNet18 CIFAR-100, Vanilla CNN CIFAR-100, ResNet18	$\begin{array}{c} 87.2\% \pm 0.2\% \\ 91.0\% \pm 0.3\% \\ 58.2\% \pm 0.3\% \\ 61.9\% \pm 0.4\% \end{array}$	$\begin{array}{c} 88.0\% \pm 0.4\% \\ 90.9\% \pm 0.2\% \\ 56.3\% \pm 0.5\% \\ 63.3\% \pm 0.6\% \end{array}$