Table 1: Performance metrics for different algorithms and models. Results with \dagger are run with our code; with \ast copied from related work. Models with a diamond (\diamondsuit) iterate for the correct amount of steps during train time (may differ between datapoints). Filled diamond (\blacklozenge) means the ground truth number of steps is also given at test time. HCS stands for hardcoded steps to 64 during test time; colour is used if on average 64 steps are more/less than ground truth. LT stands for learnt termination.

| Algorithm | $\mathbf{NAR}^{\dagger igodoldsymbol{ h}}$ | $\mathbf{NAR}^{\dagger \diamondsuit}_{(\mathrm{HCS})}$ | $\begin{array}{c} \mathbf{NAR}^{* \blacklozenge} \\ (\mathrm{Triplet-MPNN}) \end{array}$ | $\begin{array}{c} \mathbf{NAR}^{\dagger \blacklozenge} \\ (\mathrm{Triplet-MPNN}) \end{array}$ | $\mathbf{NAR}^{\dagger \diamondsuit}_{\mathrm{(LT)}}$ | $\mathbf{DEAR}^{\dagger}_{(\mathrm{ours})}$ |
|---|--|--|---|---|---|---|
| Bellman-F. Floyd-W. DSP MST Prim | $\begin{array}{c} 97.06\%\pm0.40\\ 52.53\%\pm0.98\\ 94.21\%\pm1.77\\ 93.56\%\pm0.77 \end{array}$ | $\begin{array}{c} 68.26\% \pm 13.58 \\ 52.53\% \pm 0.98 \\ 93.56\% \pm 1.16 \\ 93.56\% \pm 0.77 \end{array}$ | $\begin{array}{c} 93.26\% \pm 0.04 \\ 40.80\% \pm 2.90 \\ 97.62\% \pm 0.21 \\ 99.24\% \pm 0.21 \end{array}$ | $\begin{array}{c} 97.23\% \pm 0.15 \\ 61.86\% \pm 1.57 \\ 93.32\% \pm 1.60 \\ 92.01\% \pm 1.50 \end{array}$ | $\begin{array}{c} 95.39\% \pm 1.42 \\ 49.30\% \pm 0.53 \\ 88.30\% \pm 1.04 \\ 87.69\% \pm 1.17 \end{array}$ | $\begin{array}{c} 96.78\% \pm 0.43 \\ 55.75\% \pm 2.20 \\ 89.81\% \pm 0.14 \\ 88.67\% \pm 0.74 \end{array}$ |
| BFS DFS SCC | $\begin{array}{c} 99.85\% \pm 0.09 \\ 16.89\% \pm 5.73 \\ 40.70\% \pm 1.39 \end{array}$ | $\begin{array}{c} 87.42\% \pm 11.91 \\ 27.54\% \pm 0.52 \\ 41.79\% \pm 0.56 \end{array}$ | $\begin{array}{c} 99.89\% \pm 0.03 \\ 47.79\% \pm 4.19 \\ 57.63\% \pm 0.68 \end{array}$ | $\begin{array}{c} 99.69\% \pm 0.29 \\ 31.20\% \pm 4.02 \\ 46.84\% \pm 1.70 \end{array}$ | $\begin{array}{c} 99.51\% \pm 0.06 \\ 29.07\% \pm 2.32 \\ 39.33\% \pm 1.52 \end{array}$ | $\begin{array}{c} 98.73\% \pm 0.37 \\ 40.62\% \pm 0.44 \\ 43.63\% \pm 1.19 \end{array}$ |
| Search (Binary) Minimum Sort (Ins.) | $\begin{array}{c} 94.67\% \pm 2.31 \\ 97.67\% \pm 5.73 \\ 27.07\% \pm 10.3 \end{array}$ | $\begin{array}{c} 27.00\% \pm 24.25 \\ 97.67\% \pm 5.73 \\ 27.07\% \pm 10.3 \end{array}$ | $\begin{array}{c} 93.21\% \pm 1.10 \\ 99.24\% \pm 0.21 \\ 77.29\% \pm 7.42 \end{array}$ | $\begin{array}{c} 93.33\% \pm 2.31 \\ 96.67\% \pm 2.31 \\ 63.67\% \pm 39.97 \end{array}$ | $\begin{array}{c} 84.33\%\pm8.33\\ 94.00\%\pm2.00\\ 33.8\%\pm12.04 \end{array}$ | $\begin{array}{c} 59.00\% \pm 12.3 \\ 97.22\% \pm 3.82 \\ 86.93\% \pm 3.87 \end{array}$ |
| Overall | 71.42% | 61.64% | 80.60% | 77.58% | 70.07% | 75.42% |

Table 2: Fixing anomalies with CLRS-30's binary search further increases our overall score making our approach competitive to Triplet-MPNN. In the tables below, we will use the fixed version of search (which we reran for all models). Notation taken from Table 1.

| Algorithm | $\mathbf{NAR}^{\dagger \blacklozenge}$ | $\mathbf{NAR}^{\dagger \diamondsuit}_{(\mathrm{HCS})}$ | $\begin{array}{c} \mathbf{NAR}^{\dagger \blacklozenge} \\ (\mathrm{Triplet-MPNN}) \end{array}$ | $\mathbf{NAR}^{\dagger \diamondsuit}_{\mathrm{(LT)}}$ | $\mathbf{DEAR}^{\dagger}_{(\mathrm{ours})}$ |
|-------------|--|--|--|---|---|
| Search | $95.67\% \pm 0.58$ | $95.00\% \pm 1.73$ | $93.33\% \pm 0.58$ | $93.33\% \pm 3.05$ | $85.67\% \pm 0.58$ |
| New Overall | 71.52% | 68.44% | 77.58% | 70.97% | 78.38% |

Table 3: Extreme OOD testing of NAR and DEAR. Search may skew the results at larger scales, hence we report the overall without it too.

| Algorithm | Size 128 $(8\times)$ | | Size 250 | $6 (16 \times)$ | Size 512 $(32\times)$ | | |
|--|--|--|--|--|--|---|--|
| | $\mathbf{NAR}^{\dagger \blacklozenge}$ | $\mathbf{DEAR}^{\dagger}_{(\mathrm{ours})}$ | NAR [†] ♦ | $\mathbf{DEAR}^{\dagger}_{(\mathrm{ours})}$ | $\mathbf{NAR}^{\dagger \blacklozenge}$ | $\mathbf{DEAR}^{\dagger}_{(\mathrm{ours})}$ | |
| Bellman-F. Floyd-W. DSP MST Prim BFS DFS SCC Search Minimum Sort (Ins.) | $\begin{array}{c} 94.98\% \pm 0.52 \\ 28.16\% \pm 1.38 \\ 88.52\% \pm 3.44 \\ 86.95\% \pm 1.11 \\ 99.95\% \pm 0.00 \\ 8.14\% \pm 2.04 \\ 22.93\% \pm 1.48 \\ 87.67\% \pm 5.51 \\ 94.00\% \pm 0.00 \\ 11.59\% \pm 4.85 \end{array}$ | $\begin{array}{c} 95.65\% \pm 0.45\\ 31.37\% \pm 2.59\\ 81.21\% \pm 4.29\\ 80.59\% \pm 1.76\\ 98.87\% \pm 0.44\\ 21.67\% \pm 1.64\\ 28.56\% \pm 1.25\\ 71.67\% \pm 11.37\\ 94.00\% \pm 3.46\\ 48.90\% \pm 11.41 \end{array}$ | $\begin{array}{c} 91.83\% \pm 1.50 \\ 39.28\% \pm 0.71 \\ 78.65\% \pm 5.04 \\ 75.83\% \pm 2.69\% \\ 99.62\% \pm 0.12 \\ 3.65\% \pm 0.75 \\ \text{OOM} \\ 79.67\% \pm 10.02 \\ 92.33\% \pm 0.58 \\ 3.53\% \pm 1.11 \end{array}$ | $\begin{array}{c} 93.09\% \pm 1.33 \\ 41.63\% \pm 1.61 \\ 64.72\% \pm 8.73 \\ 69.59\% \pm 3.13 \\ 97.35\% \pm 0.05 \\ 11.53\% \pm 1.49 \\ OOM \\ 42.67\% \pm 9.71\% \\ 91.67\% \pm 3.06 \\ 16.37\% \pm 10.5 \end{array}$ | $\begin{array}{c} 87.29\% \pm 1.89\% \\ OOM \\ 54.38\% \pm 17.71\% \\ 66.18\% \pm 4.08\% \\ 99.71\% \pm 0.21\% \\ 1.24\% \pm 0.20\% \\ OOM \\ 72.00\% \pm 22.07\% \\ 87.00\% \pm 1.0\% \\ 1.63\% \pm 0.71\% \end{array}$ | $\begin{array}{c} 89.05\% \pm 2.11\% \\ OOM \\ 43.98\% \pm 9.75\% \\ 59.83\% \pm 3.23\% \\ 97.52\% \pm 0.87\% \\ 4.47\% \pm 0.63\% \\ OOM \\ 20.00\% \pm 5.29\% \\ 86.67\% \pm 2.89\% \\ 5.83\% \pm 4.80\% \end{array}$ | |
| Overall | 62.29% | 65.25% | 62.71% | 58.74% | 58.68% | 50.92% | |
| Overall w/o search | 59.47% | 64.54% | 60.59% | 60.75% | 56.78% | 55.34% | |

Table 4: Mean runtime per sample in seconds (A100 80GB GPU). OOM indicates out-of-memory. Rows and columns correspond to Table 3. Vertical bar | is used to separate the symbols (up/down arrows) across the three sizes. Double symbol is used for substantial differences.

| - • | (-1 | | , | | | * | |
|--------------------|---|----------------------|--------|------------------------|--------|------------------------|--------|
| Algorithm | | Size 128 $(8\times)$ | | Size 256 $(16 \times)$ | | Size 512 $(32 \times)$ | |
| | | NAR | DEAR | NAR | DEAR | NAR | DEAR |
| Bellman-F. | $\uparrow \uparrow \uparrow$ | 0.0340 | 0.0575 | 0.0413 | 0.0912 | 0.0710 | 0.1256 |
| Floyd-W. | $\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow -$ | 0.5856 | 0.1578 | 2.9694 | 0.3850 | OOM | OOM |
| DSP | | 0.6756 | 0.1339 | 1.5654 | 0.1632 | 8.2327 | 0.3279 |
| MST Prim | | 0.3883 | 0.0766 | 0.9941 | 0.1203 | 4.1825 | 0.1790 |
| BFS | $\uparrow \uparrow \uparrow$ | 0.0223 | 0.0310 | 0.0283 | 0.0380 | 0.0450 | 0.0765 |
| DFS | | 1.2965 | 0.1186 | 3.8432 | 0.1511 | 21.4690 | 0.3256 |
| SCC | $\downarrow\downarrow$ | 2.1008 | 0.0908 | OOM | OOM | OOM | OOM |
| Search | $\uparrow \approx \approx$ | 0.0365 | 0.0426 | 0.0725 | 0.0787 | 0.1809 | 0.1778 |
| Minimum | | 0.3451 | 0.0413 | 1.3009 | 0.0798 | 7.6732 | 0.2022 |
| Sort (Ins.) | $\downarrow\downarrow \downarrow\downarrow \downarrow\downarrow \downarrow\downarrow$ | 2.9813 | 0.0643 | 6.9552 | 0.1100 | 33.1898 | 0.2747 |

Table 5: DEAR is architecture invariant and can run with any type of processor. Notation taken from Table 1.

| | Floyd-W. | DFS | SCC | Search | Sort (Ins.) | $\begin{array}{c} \mathbf{Overall} \\ (\mathrm{subset}) \end{array}$ |
|--|--------------------------|--------------------------|--------------------------|--------------------------|-----------------------------|--|
| $\begin{array}{l} \mathbf{NAR}^{\dagger \blacklozenge} \\ (\text{Triplet-MPNN}) \end{array}$ | $61.86\% \pm 1.57$ | $31.20\% \pm 4.02$ | ${\bf 46.84\% \pm 1.70}$ | ${\bf 93.33\% \pm 0.58}$ | $63.67\% \pm 39.97$ | 59.18% |
| \mathbf{DEAR}^{\dagger} (ours; Triplet-MPNN) | ${\bf 62.29\% \pm 2.71}$ | ${\bf 42.73\% \pm 2.79}$ | $45.12\% \pm 1.52$ | $87.00\% \pm 5.57$ | $\mathbf{82.34\%} \pm 9.46$ | $\boldsymbol{63.90\%}$ |