

9 EXTENDED APPENDIX

9.1 REDS comparison to early exits

9.2 REDS performance on DNN and CNN architectures

In addition to the results on DS-CNN reported in the main paper, we show in Table 4 and Table 5 REDS performance on DNN and CNN architectures (with full fine-tuning) and compare to training model of each capacity from scratch and training REDS from scratch. Despite full fine-tuning, the results for S architecture show superior performance of the BU heuristic over TD.

MACs	Small (S) - Accuracy 83.82				Large (L) - Accuracy 86.87			
	Scratch	Knapsack BU	Knapsack TD	REDS training	Scratch	Knapsack BU	Knapsack TD	REDS training
100%	84.30 \pm 0.11	83.52 \pm 0.07	82.80 \pm 0.16	82.13 \pm 0.20	86.54 \pm 0.24	86.46 \pm 0.34	86.25 \pm 0.19	85.06 \pm 0.19
75%	83.77 \pm 0.23	82.29 \pm 0.35	81.88 \pm 0.30	81.23 \pm 0.21	85.96 \pm 0.13	86.38 \pm 0.65	86.09 \pm 0.03	84.93 \pm 0.20
50%	80.91 \pm 0.11	78.36 \pm 1.40	78.59 \pm 0.24	77.05 \pm 0.34	85.24 \pm 0.35	85.62 \pm 0.24	85.58 \pm 0.35	84.08 \pm 0.22
25%	69.77 \pm 0.67	64.42 \pm 1.99	61.43 \pm 3.35	63.69 \pm 3.26	84.22 \pm 0.13	82.61 \pm 0.61	83.00 \pm 0.62	82.03 \pm 0.59

Table 4: Test set accuracy [%] of training S and L fully-connected (DNN) architectures taken from [51]: training a network of each size from scratch ("Scratch"), conversion from a pre-trained network using two knapsack versions ("Knapsack BU" and "Knapsack TD"), and training REDS structure from scratch ("REDS training"). Reported results from three independent runs. The accuracy of each 100 % network reported in [51] is listed in the header row.

MACs	Small (S) - Accuracy 92.24				Large (L) - Accuracy 93.24			
	Scratch	Knapsack BU	Knapsack TD	REDS training	Scratch	Knapsack BU	Knapsack TD	REDS training
100%	91.10 \pm 0.23	91.60 \pm 0.39	91.20 \pm 0.35	88.89 \pm 0.26	92.94 \pm 0.20	92.83 \pm 0.26	92.97 \pm 0.15	90.97 \pm 0.21
75%	90.40 \pm 0.27	90.63 \pm 0.19	90.20 \pm 0.13	87.64 \pm 0.46	92.74 \pm 0.12	92.74 \pm 0.23	92.54 \pm 0.07	90.67 \pm 0.09
50%	89.07 \pm 0.24	88.39 \pm 0.31	88.39 \pm 0.52	85.99 \pm 0.19	92.44 \pm 0.30	91.95 \pm 0.22	91.93 \pm 0.28	90.36 \pm 0.30
25%	82.57 \pm 0.41	79.52 \pm 0.67	79.28 \pm 0.54	79.25 \pm 0.40	90.98 \pm 0.60	90.24 \pm 0.35	90.31 \pm 0.03	88.25 \pm 0.66

Table 5: The same as in Table 4 for S and L convolutional architectures (CNN) from [51].

Fig. 10 shows the impact of the architecture on REDS structure found by the knapsack BU solver. We present the results for DNN S, L and CNN S, L from left to right, respectively.

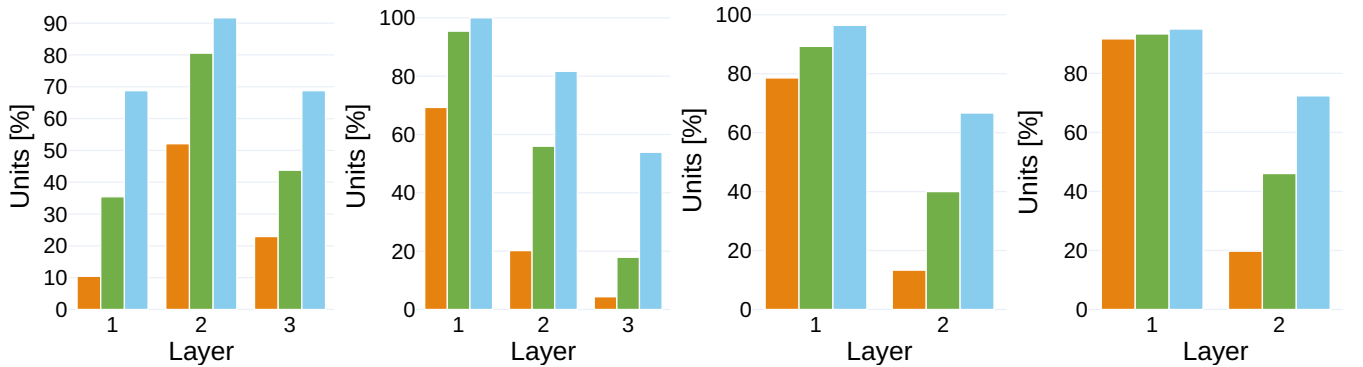


Figure 10: Analysis of the subnetwork architecture obtained by the knapsack BU heuristics. From left to right: DNN S, DNN L, CNN S and CNN L on GOOGLE SPEECH COMMANDS. The patterns as to which computational units constitute a child subnetwork are architecture-specific.

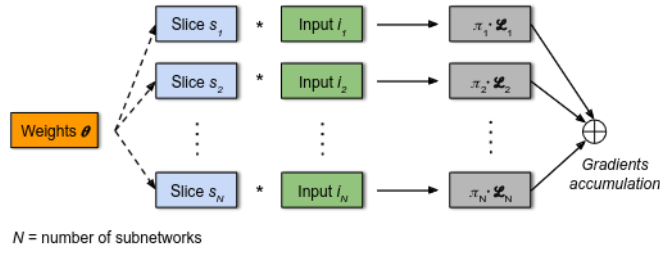


Figure 11: REDS fine-tuning. The parameters $\pi_{i=1}^N$ ensure the contribution of individual models to the loss aligns with the fraction of the shared weights.

9.3 REDS performance with 10 nested subnetworks

Table 6 and Fig. 12 show the performance of DNN, CNN and DS-CNN on GOOGLE SPEECH COMMANDS, when REDS structure comprises 10 subnetworks, compared to 4 subnetworks in the main paper. A larger number of subnetworks does not degrade model accuracy.

MACs	DNN		CNN		DS-CNN	
	Small	Large	Small	Large	Small	Large
100%	83.07 \pm 0.35	86.19 \pm 0.26	91.16 \pm 0.45	93.1 \pm 0.1	93.5 \pm 0.15	94.34 \pm 0.07
90%	82.93 \pm 0.4	86.17 \pm 0.36	90.4 \pm 0.47	92.84 \pm 0.27	93.33 \pm 0.11	94.32 \pm 0.1
80%	82.67 \pm 0.53	86.1 \pm 0.08	90.1 \pm 0.07	92.77 \pm 0.06	92.84 \pm 0.15	94.31 \pm 0.1
70%	81.67 \pm 0.53	85.78 \pm 0.11	89.43 \pm 0.41	92.25 \pm 0.47	92.64 \pm 0.24	94.21 \pm 0.14
60%	80.37 \pm 0.57	85.66 \pm 0.25	88.84 \pm 0.25	92.28 \pm 0.11	92.27 \pm 0.31	94.08 \pm 0.03
50%	78.26 \pm 0.53	85.33 \pm 0.09	88.4 \pm 0.2	92.03 \pm 0.23	91.04 \pm 0.51	93.97 \pm 0.04
40%	75.37 \pm 1.26	84.8 \pm 0.45	85.76 \pm 0.18	91.78 \pm 0.04	89.42 \pm 0.66	93.83 \pm 0.22
30%	67.76 \pm 1.59	82.66 \pm 0.18	81.92 \pm 0.49	90.52 \pm 0.54	87.63 \pm 1.03	93.59 \pm 0.14
20%	52.36 \pm 6.99	80.19 \pm 0.39	73.74 \pm 0.04	88.61 \pm 0.51	84.14 \pm 2.57	93.35 \pm 0.15
10%	23.93 \pm 5.88	50.7 \pm 7.5	58.87 \pm 5.27	81.15 \pm 1.39	58.46 \pm 3.35	90.38 \pm 0.17

Table 6: Test set accuracy [%] from Small (S) and Large (L) pretrained fully-connected (DNN), convolutional (CNN), and depth-wise separable convolutional (DS-CNN) networks taken from [51]. For each pre-trained architecture, REDS can support ten subnetworks obtained from the Knapsack BU formulation. The accuracies of the DS-CNN and CNN subnetworks do not degrade drastically until the lowest percentage of MACs considered. In contrast, the accuracies in the DNN subnetworks show a more pronounced drop from 30% MACs.

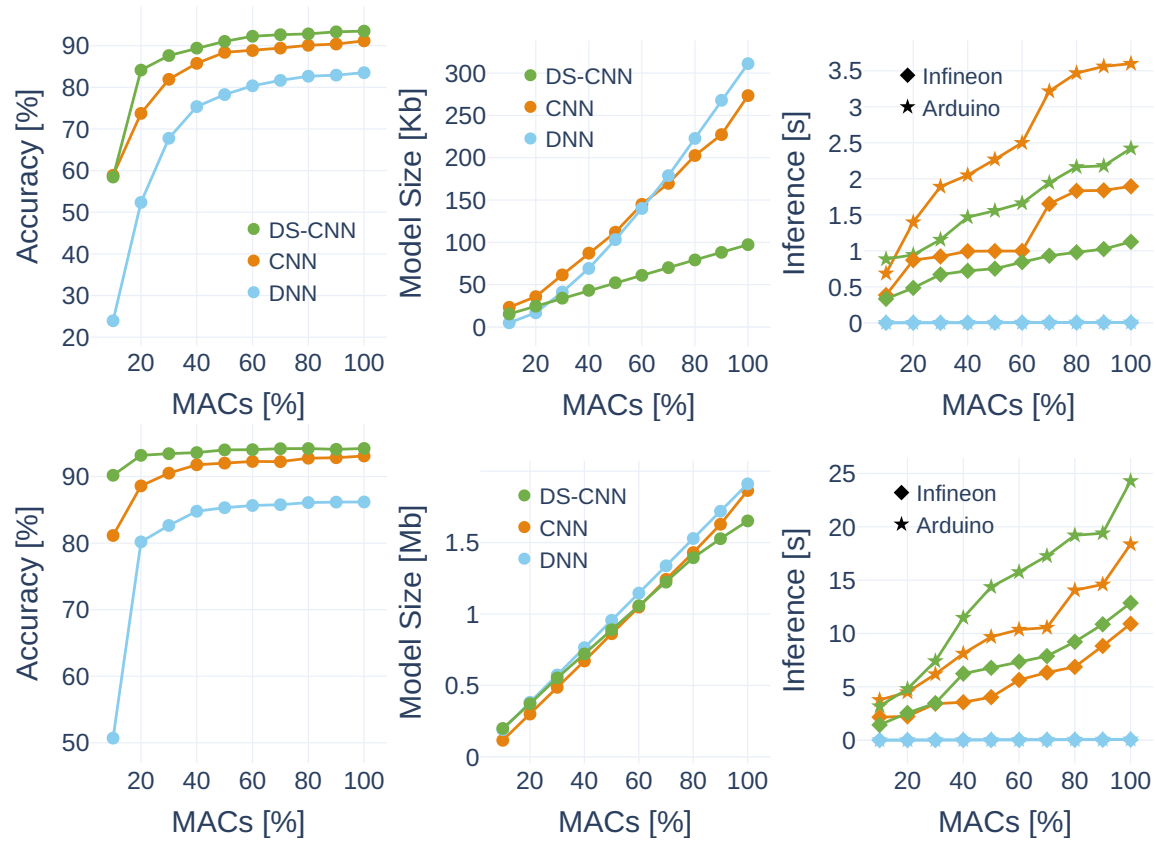


Figure 12: REDS size S (top row) and L (bottom row) architectures analysis finetuned on Google Speech Commands [49] with ten subnetworks. The plots from left to right show the subnetworks size, the subnetworks accuracy and the subnetworks inference time as a function of MAC percentage.

9.4 REDS on FASHION-MNIST and CIFAR10

The results in Table 7 and Table 8 show REDS performance using DS-CNN architecture of size S on FASHION-MNIST and CIFAR10. BU heuristic was used to obtain the results. REDS supports a different data domain without degrading the accuracy of the pre-trained model, reported in the header row. Compared to the state-of-the-art such as μ NAS [32], REDS demonstrates a faster architecture search time for both FASHION-MNIST and CIFAR10. In the former, REDS takes 19 minutes as opposed to 3 days; in the latter, REDS takes 90 minutes as opposed to 39 days while requiring less memory for model storage for both datasets. After finding and freezing the 25% MACs subnetwork architecture, the BU heuristic takes only a few seconds to find the other 50% and 75% MACs subnetworks architectures.

MACs	Acc (%) - Pre-trained 90.59	Model Size (Kb)	Time Taken (m)
100%	91.6 \pm 0.2	128.54	–
75%	91.51 \pm 0.28	107.73	1.58 [s]
50%	90.75 \pm 0.32	87.4	5.83 [s]
25%	89.22 \pm 0.45	66.63	19 [m]

Table 7: Analysis of the BU knapsack subnetworks obtained from a depth-wise separable convolutional (DS-CNN) S network, pre-trained on FASHION-MNIST [50]. REDS supports a different data domain without degrading the accuracy of the pre-trained model, reported in the header row.

MACs	Acc (%) - Pre-trained 79.36	Model Size (Kb)	Time Taken
100%	81.07 \pm 0.71	128.54	–
75%	80.17 \pm 0.69	109.41	2.89 [s]
50%	76.72 \pm 1.37	88.01	10.59 [s]
25%	68.66 \pm 1.65	69.63	90 [m]

Table 8: The same evaluation as in Table 7 for CIFAR10 [29].

9.5 REDS energy efficiency

Table 9 reports the energy consumption of inference time by different network architectures measured by Power Profiler Kit (PPK2) on Nordic nRF52840 (Arduino Nano 33 BLE Sense). In comparison, the model adaptation time takes less than 0.01 mJ.

MACs	DNN [mJ]	CNN [mJ]	DS-CNN [mJ]
100%	0.18	90.61	61
75%	0.15	86.01	44.03
50%	0.07	50.3	33.81
25%	0.05	44.81	20.44

Table 9: Energy consumption of REDS size S architectures measured by the Power Profiler Kit (PPK2) on Nordic nRF52840. The results are obtained by performing an inference pass for each subnetwork model and recording the inference current.