599 A Algorithm 1 explanation

We follow the standard Knapsack problem dynamic programming solution to break down the original 600 problem into sub-problems. Specifically, for each neuron j (or neuron group) in layer l, we can 601 choose to either include or not include it under the latency constraint c. When it is kept, the total 602 importance score increases \mathcal{I}_{l}^{j} while the latency constraint for the other neurons becomes $c - c_{l}^{j}$; If 603 the neuron is removed, the latency constraint for the other neurons remains c. We choose to keep or 604 remove the current neuron to maximize total importance. At the same time, we check whether the 605 more important neurons in the same layer are included to ensure the correctness of the latency. The 606 neuron selection from the remaining neurons is a sub-problem to solve. 607

Precisely, we use a vector max $V \in \mathbb{R}^{(C+1)}$ to store the maximum importance that we can achieve 608 under the latency constraint $c, 0 \le c \le C$ and $keep \in \mathbb{R}^{L \times (C+1)}$, a 2D vector where keep[l, c] denotes 609 the number of neuron groups we need to maintain in layer l to obtain the maximum importance 610 $\max V[c]$. We process the neurons according to their importance score in decreasing order. In this 611 way, all preceding neurons to the current one (i.e., neurons with a higher importance score in the same 612 layer) will be always considered first. To decide if we keep or remove the current neuron, we check 613 the total importance score and the inclusion status of its preceding neurons, so we can maximize the 614 total importance and ensure the latency cost correctness. 615

616 **B** Experimental settings

For image classification, in the main paper we focus on pruning networks on the large-scale ImageNet
ILSVRC2012 dataset [52] (1.3M images, 1000 classes). Each pruning process consumes a single
node with eight NVIDIA Tesla V100 GPUs. We use PyTorch [48] V1.4.0 model zoo for pretrained
weights for our pruning for a fair comparison with literature.

In our experiments we perform iterative pruning. Specifically, we prune every 320 minibatches after 621 loading the pretrained model with k = 30 pruning steps in total to satisfy the constraint. Unless 622 otherwise specified, we finetune the network for 90 epochs in total with an individual batch size 623 at 128 for each GPU. For finetuning, we follow NVIDIA's recipe [47] with mixed precision and 624 Distributed Data Parallel training. The learning rate is warmed up linearly in the first 8 epochs and 625 reaches the highest learning rate, then follows a cosine decay over the remaining epochs [37]. For the 626 result in Fig.1 that is trained with knowledge distillation, we use RegNetY-16GF (top1 82.9%) as the 627 teacher model when finetuning the pruned model. We use hard distillation on the logits and the final 628 training loss is calculated as $L = (1 - \alpha)L_{base} + \alpha L_{distil}$ where $\alpha = 0.5$ to balance between the 629 original loss L_{base} and the distillation loss L_{distil} . 630

For latency lookup table construction, we target a NVIDIA TITAN V GPU with batch size 256 for latency measurement to allow for highest throughput for inference, and target a Jetson TX2 with inference batch size 32. We pre-generate a layer latency look-up table on the platform by iteratively reducing of the number of neurons in a layer to characterize the latency with NVIDIA cuDNN [7] V7.6.5. We profile each latency measurement 100 times and take the average to avoid randomness.

We also provide pruning results on the small-scale CIFAR10 [29] dataset in appendix Sec. C. For
CIFAR10 experiments with ResNet-50/-56, we train the model on a single GPU for 200 epochs in
total where we perform pruning step very one epoch in the first 30 epochs and finetune the pruned
model during the remaining 170 epochs. The initial learning rate is set to 0.1 with batch size 128.
For DenseNet, we extend the finetuning epochs to 300 epochs

641 C More pruning results on CIFAR10 and ImageNet

We provide the pruning results of our method on CIFAR10 dataset in this section. We chose 3 network architectures for the experiment: ResNet50, ResNet56 and DenseNet40-12 [26]. As most of the prior methods perform pruning under the FLOPs constraint, in our CIFAR10 experiment we also use FLOPs constraint instead of the latency constraint. We also add some additional ImageNet-ResNet50 results comparison in the table. For a fair comparison, in the ImageNet50 experiment, we also provide the model accuracy under the same finetuning recipe and epochs as the methods to be compared to alleviate the potential impact of the different finetuning settings. Specifically, when compared to

1030 01	the unprunet	anetwork		
Dataset	Model	Method	FLOPs (%) \downarrow	Top1 (%) ↑
	ResNet50	ChipNet [58]	17.7	92.8
	Resiliences	LASP(Ours)	13.7	93.2
		CHIP [55]	27.7	92.05
CIFAR10	D N 5 (LASP (Ours)	26.8	93.22
	Residences	GDP [19]	34.36	93.55
ļ		LASP (Ours)	33.72	93.68
	D	QCQP [28]	29.2	93.80^{\dagger}
	DenselNet40-12	LASP (Ours)	29.2	93.15^{\dagger}
		GBN-60 [67]	59.46	76.19
	ResNet50	QCQP [28]	59.0	76.00
ImageNet		LASP (Ours)	58.12	$76.93 \\ (76.49^* / 76.53^{**})$
		GBN-50 [67]	44.94	75.18
		LASP (Ours)	42.05	$76.09(75.27^{**})$
Dataset	Model	Method	Inf speedup (%) ↑	Top1 drop (%) \downarrow
ImagaNat	DecNet50	QCQP [28]	$1.52 \times$	0.32
mageinet	Residet50	LASP (Ours)	1.60 imes	$-0.22(0.16^{**})$

Table 5: Additional pruning results and comparison on CIFAR10 and ImageNet dataset. FLOPS (%) are relative to those of the unpruned network

The baseline model used in QCQP has 95.01% top1 accuracy, while our pretrained model has 94.40% top1 accuracy. The accuracy drop is 1.21% vs. 1.25%, which is comparable use the same finetune recipe and epochs as GBN [67]

** use the same finetune recipe and epochs as QCQP [28]

Table 6: Pruning MobileNet-V1 and MobileNet-V2 on the ImageNet dataset with different targets.

Method	FLOPs (M)	Top1 (%)	Top5 (%)	FPS (im/s)	Speedup		Method	FLOPs (M)	Top1 (%)	Top5 (%)	FPS (im/s)	Speedup
No pruning	1 569	MobileN 72.64	et-V1 90.88	3415	$1 \times$:	No pruning	301	MobileN 72.10	et-V2 90.60	3080	1×
LASP-40% LASP-42% LASP-50% LASP-60% LASP-70% LASP-80% LASP-90%	$ 154 \\ 171 \\ 237 \\ 297 \\ 360 \\ 416 \\ 507 $	67.20 68.30 69.79 71.31 71.78 72.52 72.95	87.32 88.08 89.08 90.05 90.39 90.78 91.02	8293 7940 6887 5754 4870 4167 3765	$2.43 \times 2.32 \times 2.02 \times 1.68 \times 1.43 \times 1.22 \times 1.10 \times 100 \times 1000 \times 100 \times 100 \times 100 \times 100 \times 100 \times 1000 \times 100 \times 100 \times 100 \times 1$	-	LASP-60% LASP-65% LASP-70% LASP-75% LASP-90% LASP-95%	183 218 227 249 273 281	$\begin{array}{c} 70.42 \\ 71.41 \\ 71.88 \\ 72.16 \\ 72.45 \\ 72.55 \end{array}$	89.75 90.08 90.39 90.44 90.68 90.79	$5668 \\ 5003 \\ 4478 \\ 4109 \\ 3443 \\ 3265$	$1.84 \times 1.62 \times 1.45 \times 1.33 \times 1.12 \times 1.06 \times$

GBN [67], we finetune the pruned network for 60 epochs, where initial learning rate is set to 0.01649 with batch size 256. The learning rate is divided by 10 at epoch 36, 48 and 54. When compared to 650 QCQP [28], the pruned network is finetuned for 80 epochs with batch size 384 and the initial learning 651 rate of 0.015. Then, the learning rate is decayed at epoch 30 and 60 by deviding 10. As shown in 652 Tab. 5, LASP method consistently outperforms with lower FLOPs and higher Top1 accuracy. When 653 it comes to the actual inference speed comparison with QCQP [28], our method yields 0.16% less 654 accuracy drop while getting $0.08 \times$ faster speed. 655

Table. 6 and Fig. 5 provide additional prun-656 ing results for lightweight networks such as 657 MobileNet-V1 and MobileNet-V2. For the un-658 pruned models, we find that even MobileNet-V2 659 has significantly lower FLOPs, the inference 660 time is larger compared to MobileNet-V1.In 661 both cases, LASP yields inference speeds-ups 662 of $1.22 \times$ and $1.33 \times$ for MobileNet-V1 and 663 MobileNet-V2 respectively, while maintaining 664 the original top1 accuracy. 665



Figure 5: Pruning MobileNets on the ImageNet dataset

Efficacy of neuron grouping on MobileNet D 666

In this section, we show the benefits of latency-aware neuron grouping and the performance under 667 different group size settings on MobileNetV1. 668

Since MobileNet has group convolutional layers to speedup the inference, we take the group convolu-669 tional layer with its preceding connected convolutional layer together as coupled cross-layers [17] to 670

make sure the input channel number and out-671 put channel number of the group convolution 672 remain the same. All the 27 convolutional lay-673 ers can be divided into 14 coupled layers. In our 674 method, with the neuron grouping, we set the 675 individual group size of 1 coupled layer to 16, 676 of 3 coupled layers to 32 and 10 coupled layers 677 to 64. Also, for MobileNetV1 pruning, we add 678 the additional constraint that each layer has at 679 least one group of neurons remaining to make 680 sure that the pruned network is trainable. 681

We compare our latency-aware neuron grouping
with an heuristic option by setting a fixed group
size for all layers. Fig. 6 shows the comparison



Figure 6: Performance comparison of our latency-aware grouping to different fixed sizes for a MobielNetV1 pruned with different latency constraints on ImageNet. We compare to heuristic-based group selection studied by [66]. LG denotes the proposed latency-aware grouping in LASP that yields consistent latency benefits per final accuracy.

results between our neuron grouping method and various fixed group sizes for a MobielNet pruned
 with different latency constraints on ImageNet. As shown, similar to ResNet50, using small group
 sizes such as 8, 16 leads to worse performance; a large group size like 128 also harms the performance
 significantly. Our observations on ResNet50 pruning also hold in MobileNetV1 setting, further
 emphasizing the efficacy of our latency-aware neuron grouping.

690 E Ablation study of pruning step k

In this work, similar to many other prior meth-691 ods [2, 44, 68], we do iterative pruning with k 692 pruning steps in total. In this experiment, we 693 analyze the the accuracy of the final result as 694 a function of k. We set the value of k to 10, 695 20, 30 and 40 for iterative pruning, and also use 696 k = 1 to perform a single-shot pruning. The 697 result of this experiment is shown in Fig. 7. As 698 shown, we get similar results independent of 699 k. Imporantly, all these results outperform Ea-700 gleEye [31]. As expected, there is a drop in 701 accuracy for single-shot pruning (k = 1), espe-702



Figure 7: Performance comparison of different pruning steps k for ResNet50 pruning on ImageNet.

cially for large pruning ratios. The main reason is the neuron importance would change as we remove some other neurons and, in this setting, the value is not updated. Iterative pruning does not have this limitation as the importance score and the latency cost of the remaining neurons is updated after each pruning step to reflect any changes. In our experiments, we use k = 30 as it provides a good trade off between latency and accuracy.

708 F Comparison with EagleEye on ImageNet

We now use the same unpruned baseline model provided by EagleEye [31] to compare our proposed LASP method with EagleEye [31] varying the latency constraint. As shown in Fig. 8, our approach dominates EagleEye by consistently delivering a higher top-1 accuracy with a significantly faster inference time.

We then analyze the structure difference between our pruned model and the EagleEye
model. As mentioned in the main text that the
proposed LASP method tries to make the number of remaining neurons in each layer fall to the
right side of a step if the latency on the targeting

722 platform presenting a staircase pattern. Fig. 9



Figure 8: Pruning ResNet50 on the ImageNet dataset using the same baseline model as in EagleEye with a top-1 accuracy of 77.23%. The proposed LASP surpasses EagleEye ECCV20 [31] in accuracy and latency. Top-left is better.

shows two examples of pruned layers after pruning from LASP-45% and EagleEye-2G model. In the



Figure 9: Two examples of pruned layers from LASP model and EagleEye [31] model. The scattered black points are the locations of the layers fall to after pruning.

Table 7: Pruning ResNet50 on the ImageNet dataset (TITAN V) targeting on inference with batch size 1. LASP-X% indicates that X% latency to remain after pruning. The speedup is calculated as the ratio of FPS between the pruned network and the unpruned model.

Method	FLOPs (G)	Top1 Acc (%)	Top5 Acc (%)	FPS (imgs/s)	Speedup
No pruning	4.1	76.2	92.87	181	$1 \times$
0.75× ResNet50 [21]	2.3	74.8	-	192	$1.06 \times$
AutoSlim [68]	2.0	75.6	-	181	$1.00 \times$
MetaPruning [35]	2.0	75.4	-	190	$1.05 \times$
EagleEye-2G [31]	2.1	76.4	92.89	190	$1.05 \times$
GReg-2 [61]	1.8	75.4	-	196	$1.09 \times$
LASP-90% (Ours)	2.9	76.4	93.10	220	1.22 imes
0.50× ResNet50 [21]	1.1	72.0	-	193	$1.07 \times$
AutoSlim [68]	1.0	74.0	-	191	$1.06 \times$
MetaPruning [35]	1.0	73.4	-	196	$1.09 \times$
EagleEye-1G [31]	1.0	74.2	91.77	192	$1.06 \times$
GReg-2 [61]	1.3	73.9	-	206	$1.14 \times$
LASP-80% (Ours)	2.3	75.3	92.35	247	1.37 imes

left figure, we show that the layer in our pruned model has only 5 more neurons pruned than that in EagleEye model, the latency is reduced to a much lower level which is a 0.76ms drop while we have more input channels. In the right figure, we also show that sometimes we can remain a lot more neurons (30 neurons) in layer with only little latency (0.21ms) increase. These two examples both show the ability of method to fully exploit the latency traits and benefit the inference speed.

Our method benefits a lot from the non-linear latency characteristic since we are trying to keep as 729 many neurons as possible under the latency constraint. If the latency of the layer on the targeting 730 platform shows linear pattern, the advantage of our method becomes smaller. Fig. 9 shows the latency 731 behavior of the example layers on the targeting platform when reducing the number of input and 732 output channels. As we can see, the staircase pattern becomes less obvious as the number of input 733 channel reduces and the GPU has sufficient capacity for the reduced computation. This happens 734 during pruning, especially for large prune ratios. In such a case, the FLOP count reflects the latency 735 more accurately, and the performance gap between reducing FLOPs and reducing latency can possibly 736 become small. Nevertheless, our method can help avoid some latency peaks as shown in Fig. 9, which 737 could otherwise happen using other pruning methods. 738

739 G Pruning results for small batch size

In the main paper, we use a large batch 256 in the experiment to allow for highest throughput for inference, which also makes the latency of the convolution layers show apparent staircase pattern so that we can take full advantage of the latency characteristic. In this section, we show that with small batch size 1 that no obvious staircase pattern showing up in layer latency, our LASP algorithm still delivers better results compared to other methods.

When we use batch size 1 for inference, the layer latency of ResNet50 does not show obvious staircase pattern in most of the layers due to the insufficient usage of GPU. Therefore in this experiment, we use the latency lookup table granularity as a neuron grouping size, which in our case is 2, to fully exploit the hardware latency traits during pruning. We show our pruned results and the comparison

|--|

Model	mAP	FLOPs (G)	params (M)	FPS (BS=1)	FPS (BS=32)
SSD512-RN50, base model	77.98	65.56	21.97	68.24	103.48
SSD512-RN50-slim	75.83	46.09	16.33	76.49	114.80
SSD300-RN50	75.69	16.23	15.43	128.85	309.32
SSD300-VGG16 [34]	76.72	31.44	26.29	122.28	262.93
FasterRCNN-VGG16 [51]	70.10	91.23	137.08	29.21	-
RetinaNet-RN50 [33]	77.27	106.50	36.50	36.92	-
SSD512-RN50-LASP (Ours)	77.42	15.38	10.40	132.57	323.36

with other methods in Tab. 7. As shown in the table, while other methods reduce the total FLOPs of the network after pruning, they do not reduce the actual latency much, which is up to $1.09 \times$ faster than the original one at the cost of 2.8% top1 accuracy drop. Compared to these methods, although we get less FLOPs reduction using our proposed method, the pruned models are faster and get higher accuracy, which is $1.22 \times$ faster than the unpruned model while getting slightly higher accuracy and $1.37 \times$ faster with only 0.9% accuracy drop.

755 H Pruning results on object detection

In this section we show the detailed pruning results on objection detection task for Sec. 4.5. To prune
 the detector, we first train a SSD512 with ResNet50 as backbone. We also train some other popular
 models for performance comparison. The detailed numbers of Fig. 4 are shown in Tab. 8.

759 I Implementation details

Convert latency in float to int. Solving the neuron selection problem using the proposed augmented knapsack solver (Algo. 1 in the main paper), requires the neuron latency contribution and the latency constraint to be integers as shown in line 4 of the algorithm. To convert the measured latency from a full precision floating-point number to integer type, we multiply the latency by 1000 and perform rounding. Accordingly, we also scale and round the latency constraint value.

Deal with negative latency contribution. The neuron latency contribution in our augmented 765 knapsack solver must be a non-negative value since we have $dp_array \in \mathbb{R}^C$ and we need to visit 766 dp_array[$c - c_n$] as in line 5 of Algo. 1 in the main paper. However, by analyzing the layer latency 767 from the look-up table we find that for some layers the measured latency might even increase when 768 reducing some number of neurons. This means that the latency contribution could possibly be negative. 769 770 The simplest way to deal with the negative values is to directly set the negative latency contributions to be 0. This leads to the problem that the summed latency contribution would be larger than the 771 actual latency value, causing less neurons being selected. Thus, during our implementation, we keep 772 those negative latency values as they are, but update the vector size of dp_array to $\mathbb{R}^{C-\min(\min(\mathbf{c}),0)}$ 773 where $\min(\mathbf{c})$ is the minimum latency contribution. With such, the vector size of dp array would be 774 extended when there is negative latency contribution. This makes it possible to add one neuron with 775 negative latency contribution to a subset of neurons whose summed latency is larger than the latency 776 constraint. After the addition, the total latency will still remain under the constraint. 777

Pruning of the first layer. In our ImageNet experiments, we leave the first convolutional layer of ResNets unpruned to help maintain the top-1 classification accuracy. For MobileNet, the first convolutional layer is coupled with its following group convolutional layer. In our MobileNet experiments, we prune the first coupled layers at most to the half of neurons.

SSD for object detection. Our SSD model is based on [34]. When we train SSD-VGG16, we use the exactly same structure as described in the paper. When we train a SSD-ResNet50, the main difference between our model and the model described in the original paper is in the backbone, where the VGG is replaced by the ResNet50. Following [27], we apply the following enhancements in our backbone:

- The last stage of convolution layers, last avgpool and fc layers are removed from the original ResNet50 classification model.
- All strides in the 3rd stage of ResNet50 layers are set to 1×1 .

layer	SSD512	SSD512-slim	SSD300
layer1-conv1	(512, 3, 1, 1)	(256, 1, 1, 0)	(512, 3, 1, 1)
layer1-conv2	(512, 3, 2, 1)	(512, 3, 2, 1)	(512, 3, 2, 1)
layer2-conv1	(256, 1, 1, 0)	(256, 1, 1, 0)	(256, 1, 1, 0)
layer2-conv2	(512, 3, 2, 1)	(512, 3, 2, 1)	(512, 3, 2, 1)
layer3-conv1	(128, 1, 1, 0)	(128, 1, 1, 0)	(128, 1, 1, 0)
layer3-conv2	(256, 3, 2, 1)	(256, 3, 2, 1)	(256, 3, 2, 1)
layer4-conv1	(128, 1, 1, 0)	(128, 1, 1, 0)	(128, 1, 1, 0)
layer4-conv2	(256, 3, 2, 1)	(256, 3, 2, 1)	(256, 3, 1, 0)
layer5-conv1	(128, 1, 1, 0)	(128, 1, 1, 0)	(128, 1, 1, 0)
layer5-conv2	(256, 3, 2, 1)	(256, 3, 2, 1)	(256, 3, 1, 0)
layer6-conv1	(128, 1, 1, 0)	(128, 1, 1, 0)	-
layer6-conv2	(256, 4, 1, 1)	(256, 4, 1, 1)	-

Table 9: The additional convolution layers in SSD.

Table 10: Pruning ResNet50 on the ImageNet dataset with FLOPs constraint and comparison with state-of-the-art method EagleEye (ECCV'20) [31]. We remeasure the FLOPs, top1 and top5 accuracy of EagleEye to get results with two digits.

77 6 ---

											*	
Method	FLOPs (G)	Top1 Acc (%)	Top5 Acc (%)	77					^			
No pruning	4.1	77.23	93.70	76.5		/	A					
EagleEye-3G FLOP-T (Ours)	3.08 2.99	77.10 77.36	93.36 93.62) 90 V 10 75.5		//						
EagleEye-2G FLOP-T (Ours)	2.06 1.95	76.38 76.64	92.90 93.21	75 -								
EagleEye-1G FLOP-T (Ours)	1.03 0.96	74.18 74.84	91.78 92.26	74	5 1	1.5	2 F	2.5 LOPs (0	3	3.5	4	4.

The backbone is followed by 6 additional coupled convolution layers for input size 512×512 , or 5 for input size 300×300 . A BatchNorm layer is added after each convolution layer. The settings of these additional convolution layers are listed in Tab. 9, each layer is represented as (output channel,

⁷⁹² kernel size, stride, padding).

The detector heads are similar to the ones in the original paper. The first detection head is attached to the last layer of the backbone. The rest detection heads are attached to the corresponding additional layers. No additional BatchNorm layer in the detector heads.

796 J FLOPs-constrained pruning

Our implementation of latency-constrained pruning can be easily converted to be a FLOPs-constrained. 797 When constraining on FLOPs, $\Phi(\cdot)$ in the objective function (Eq.1 in the main paper) becomes the 798 FLOPs measurement function and C becomes the FLOPs constraint. Since the FLOPs of a layer 799 linearly decreases as the number of neurons decreases in the layer, we do not need to group neurons 800 in a layer any more. The problem can also be solved by original knapsack solver since each neuron's 801 FLOPs contribution in a layer is exactly the same and no preceding constraint is required. We conduct 802 some experiments by constraining the FLOPs and compare the results with EagleEye [31]. We name 803 the experiments using the same algorithm as LASP but targeting on optimizing the FLOPs as FLOP-T. 804 As shown in Tab. 10, with our pruning framework applying the knapsack solver, our results show 805 higher top-1 accuracy compared to the pruned networks of EagleEye with similar FLOPs remaining. 806 We also observe a larger gap between the methods when it comes to a more compact network. 807

808 K FLOPs vs. latency

FLOPs can be regarded as a proxy of inference latency; however, they are not equivalent [4, 33,
35, 40, 42]. We do global filter-wise pruning and have the same problem as NAS. The latency on a
GPU usually imposes staircase-shaped patterns for convolutional operators with varying channels
and requires pruning in groups. In contrast, FLOPs will change linearly. Depth-wise convolution,
compared to dense counterparts, has significantly fewer FLOPs but almost the same GPU latency

raore	11. 10001	eteo praim	5	of faconey constraint.
method	Top1(%)	FLOPs (G)	FPS (imgs/s)	FPS vs FLOPs
FLOP-T	74.84	0.962	2202	2500
LASP	74.92	1.210	2396	ିଙ୍କୁ mail ALASP
FLOP-T	76.64	1.949	1436	B2000
LASP	76.55	1.957	1672	S 1500
FLOP-T	77.36	2.988	1146	1000
LASP	77.45	2.988	1203	0.5 1 1.5 2 2.5 3 FLOPs (G)

Table 11: ResNet50 pruning with FLOPs/latency constrain

due to execution being memory-bounded¹. The discrepancy also holds for ResNets where the same amount of FLOPs impose more latency in earlier layers than later ones as the number of channels increases and feature map dimension shrinks – both increase compute parallelism. For example, the first 7×7 conv layer and the first bottleneck 3×3 conv in ResNet50 have nearly identical FLOPs but the former is 60% slower on-chip.

We compare our results of FLOPs-targeted (FLOP-T) showed in Sec. J and results using latencytargeted pruning (LASP) in Tab. 11. As shown in the table, using different optimization targets leads to quite different FPS vs FLOPs curves. In overall, with similar FLOPs remaining, using our LASP algorithm targeting on reducing the actual latency can get more efficient networks with more image being processed per second.

We also show a more straightforward relationship between the 824 actual latency of a layer and its FLOPs in Fig. 10. We use the 825 2nd convolution layer in the 1st residual block of ResNet50 as 826 an example. We vary the number of neurons of the layer from 827 0 to 128 and measure the actual latency on GPU (TITAN V) 828 as well as the FLOPs of the layer. We can see from the figure 829 that the actual latency does not strictly linearly decrease as the 830 FLOPs decreasing. 831



Figure 10: The measured latency vs. FLOPs of the 2nd convolution layer in the 1st residual block of ResNet50.

832 L Different choice of importance calculation

We use the first Taylor expansion [44] to estimate the loss change induced by pruning as the importance
score of the neurons. It is a gradient-based importance calculation and is shown to be given promising
results. In this section, we use the L2 norm of the neuron weights as the importance measurement
and apply LASP framework to ResNet50 ImageNet classification task. As shown in Tab. 12, Our
algorithm is generic applying to different importance measurements. As shown, using L2 norm of
weights as importance measurements leads to slightly lower accuracy.

Table 12: The results of LASP algorithm on ResNet50 ImageNet classification task with different choices of neuron importance measurements.

Method	First-Ta	ylor Expansi	ion (gradien	t-based)	L2	2 norm (mag	nitude-base	d)
	FLOPs(G)	Top1(%)	Top5(%)	FPS(im/s)	FLOPs(G)	Top1(%)	Top5(%)	FPS(im/s)
LASP-80%	3.0	77.5	93.60	1203	3.0	77.3	93.60	1196
LASP-55%	2.1	76.7	93.16	1672	2.0	75.7	92.66	1595

838

839 M Latency look-up table creation and calibration

In this section, we provide additional details to build the latency look-up table used in LASP, computational cost and the correlation between the estimated and the real ones. As mentioned in Appendix B, we pre-generate the layer latency look-up table on the platform with NVIDIA cuDNN [7] V7.6.5. For each layer, we iteratively reduce the number of neurons in the layer (each time reduce neurons) and characterize the corresponding latency. For each latency measurement, we use one profile for GPU warm up and another 3 profiles and take the average to avoid randomness. The average standard deviation of profiles for an operation is $8.67e^{-3}$.

¹https://tlkh.dev/depsep-convs-perf-investigations/

On a single TITAN V GPU, it takes around 5 847 hours to build the look-up table for ResNets fam-848 ily and 1 hour for MobileNets family. Note that 849 the LUT can be shared by the network architec-850 tures within the same family since they usually 851 have similar layer structures. We only need to 852 create the latency look-up table once for all the 853 possible latency targets. 854

There are some gaps between the predicted latency and the real latency of the model, because the latency look-up table is created layer-wise on convolution layers. There are additional costs in real inference such as pooling, non-linear activation etc. We plot the correlation between the



Figure 11: The correlation between the predicted latency reduction and the real latency reduction of the pruned models.

expected latency reduction from look-up table and the real latency reduction ratio of our pruned models in Fig. 11. We also calculate the Pearson Correlation Coefficient r for all the networks in the figure. We can clearly see a linear correlation between the predicted and real latency reduction from the figure, showing that the latency lookup table provides a good approximation and it is possible to calibrate the latency estimation using the linear coefficient to have a better estimation.

Limitation of layer-wise latency lookup table. We apply layer-wise latency lookup table in our 866 method, which does not consider the caching and parallelization among layers. For networks like 867 VGG without bypass paths, the layer operations will be executed sequentially, in which case, the 868 latency lookup table models the actual latency well. For models with parallel paths, it depends on 869 the hardware implementation whether there will be parallel execution in practice. When there is 870 parallelization, e.g., in accelerators, models like ResNets still work well because the skip connection 871 only takes small portion of computation; for models like InceptionNet, taking the parallelization 872 and into consideration when generating the lookup table would help better estimate the latency. For 873 wider applications in the future, the more domain knowledge we have about the GPU execution and 874 improve accordingly, the more accurate estimation we will obtain using the latency lookup table. 875

Comparison with quadratic model. Recent work QCQP [28] models latency using a quadratic 876 equation and solve a latency constrained optimization problem. The main difference between our 877 estimation from lookup table and the QCQP's estimation from quadratic modeling is that we use 878 different ways to model each layer's latency. QCQP[1] uses $\alpha_l + \beta_l ||r^{(l-1)}||_1 + \delta_l ||r^{(l-1)}||_1 ||r^{(l)}||_1$ to model the layer latency where $||r||_l$ is the number of remaining channels in the corresponding layer, 879 880 α , β and δ are the coefficients that need to be optimized for the targeting platform. Note that QCQP 881 also needs to profile latency for different samples of each layer, like what we do to create the lookup 882 table but will less samples, in order to optimize α , β and δ . It is also important to note that QCQP 883 uses a linear model between the layer latency and FLOPs (the quadratic part) and memory. Therefore, 884 QCQP fails to capture the latency staircase pattern (see Fig. 3a in the QCQP paper) which is the key 885 to maximize GPU utilization (see latency surface in Fig. 1). As shown in Fig.3a in the QCQP paper, 886 the quadratic modeling gives different latency estimations to layers with different number of input 887 and output channels. However, these layers have the same real inference time. As a result, the larger 888 the number of input and output channels, the larger the error in the estimation of QCQP. 889

890 N Pruning for INT8 quantization

We now focus on results when the target is INT8 891 inference which is a common requirement for 892 real-world applications. In particular, we use 893 NVIDIA Xavier as the target platform as, in this 894 platform, INT8 speedup is supported. We create 895 a INT8 latency look-up table for INT8 inference. 896 For comparison, we also create a FP32 latency 897 look-up table on the same platform and use both 898 look-up tables for pruning a ResNet50 model 899 on ImageNet classification. After convergence, 900 results are quantized into INT8 and the latency 901



Figure 12: The ResNet50 ImageNet pruning targeting INT8 inference on NVIDIA Xavier.

and accuracy is measured directly on the Xavier platform. We measure latency with a batch size 128, using TensorRT (V8.2.0.1) to get INT8 speedup. Results for this experiment are shown in Fig. 12.

As shown, even use a FP32 latency look-up table as the latency guidance, our method LASP outperforms the state-of-the-art method EagleEye [31]. These results are consistent with Sec. 4.2 where we show the LASP acceleration on GPUs with TensorRT. Pruning results of using a INT8 look-up table show that our method yields higher accuracy with lower latency on this platform. We obtain a $1.32 \times$ speedup while maintaining the original Top1 accuracy. Compared to EagleEye, LASP achieves up to $1.26 \times$ relative speedup and 0.15% higher accuracy.

910 O Breakdown of the algorithm execution time

We provide additional details of the algorithm process of different methods in this section. We estimate the time cost needed to get the pruned network structure for each method. The time of following finetuning is not taken into consideration. For a fair comparison, we set the number of pruning steps k for all iterative pruning methods to 30. All the values are approximated as all the methods are running on the same device (a NVIDIA V100 GPU) to get a pruned ResNet50. For AutoSlim [68], MetaPruning [35] and AMC [22], more GPU time is needed for additional training of the network.

Method	Evaluate proposals?	Auxiliary net training?	Sub-network selection	Additional time cost	Estimate time (RN50)
NetAdapt	Y	Ν	N candidates evaluation + finetune after each prune. Repeat $k \mbox{ times}$	Latency look-up table creation	$\sim 195 {\rm h~GPU}$
ThiNet	Y	Ν	1 or 2 train epochs after each pruning. Repeat k times	Additional forward pass to get neuron importance	$\sim 210 {\rm h~GPU}$
EagleEye	Y	Ν	1000 candidates evaluation	Monte Carlo sampling, prune to get 1000 candidates	30h GPU
AutoSlim	Y	Y	Train slimmable model k candidates evaluation		
MetaPruning	Y	Y	Train an auxiliary network k candidates evaluation		
AMC	Ν	Y	Train an RL agent		
LASP(Ours)	N	N	40 train iterations after each pruning. Repeat 30 times. (< 1 train epoch in total)	Augmented Knapsack solver (~ 30min in total) Latency look-up table creation	6.5h GPU 0.5h CPU

Table 13: The breakdown details of the execution process of different methods.

918 **P** Difference with prior work

KnapsackPruning (KP) [1] is one work that is mostly close to our method in the paper. While 919 both works look at similar problems from the same combinatorial perspective, there are several key 920 differences. First, KP focuses on constraining FLOPs and shows an instantiation on latency; in 921 contrast, we directly optimize the latency, which is more practical. Second, we show the latency 922 characterization on device and augment the original knapsack problem formulation accordingly, 923 Eq. 7, to accommodate to the latency traits - the neuron latency is dependent on the order of neuron 924 pruning in a layer, while KP uses a standard knapsack where the neuron FLOP cost is independent 925 of each other. Please note here that the formulation in KP assumes the independence thus can not 926 directly apply to the latency-targeted pruning, Third, we are the first ones to use the latency-aware 927 grouping which assigns different grouping sizes to each layer according to the latency traits rather 928 than predefined fixed values. 929

For a fair comparison of pruning to the KP method, we use the PyTorch baseline as unpruned model and both without knowledge distillation during training. The results for ResNet50 pruning on ImageNet are show as Tab. 16. As shown, our method performs significantly better leading to pruned model with higher accuracy but less FLOPs. On the other side, as they are not considering the actual latency, the resultant network structure is not GPU friendly and would fail to maximize the GPU utilization.

Q Detailed configuration of pruned models

⁹³⁷ We provide the detailed configuration of our pruned models of Tab. 1. For each model, we list the ⁹³⁸ number of neurons remaining in each convolution layer, starting from the input to the output. For

Table 14: ResNet50 pruning results on ImageNet comparison with Knapsack Pruning [1].

Method	FLOPs (G)	Top1 (%)	Acc Drop (%)
KP [1]	2.38	76.17	0.03
KP [1]	2.03	75.94	0.21
LASP (Ours)	2.01	76.46	-0.26

ResNets we use $[\cdot]$ to denote a residual block and (\cdot) to denote the neuron number of the residual bypass layer.

Note that for ResNets, the configuration is the "raw" configuration after the pruning. For each residual block $[x_1, x_2, x_3]$, because of the existing of bypass layer, we allow the entire layer pruning and it still leads to trainable networks. Whenever there is entire layer pruning $(x_1 == 0 \text{ or } x_2 == 0 \text{ or} x_3 == 0)$, all the other layers in the block can be removed and this branch can be replaced by a constant value. In such a case, the corresponding block in our final pruned model is cleaned to be [0, 0, 0]. We also provide the detailed structure plot of two pruned ResNet50 models in Fig. 13 for a better visualization.

Table 15: The detailed configuration of the LASP pruned models.

	ResNet50 - EagleEye [31] baseline
LASP-80%	$ \begin{array}{l} 64, [64, 32, 256](256), [32, 32, 256], [0, 32, 256], [128, 128, 512](512), [64, 96, 512], [64, 128, 512], [26, 128, 128, 128], [256, 128, 1024], [25$
LASP- 45%	$ \begin{array}{l} 64, [64, 32, 128](128), [32, 0, 128], [0, 32, 128], [64, 64, 384](384), [64, 96, 384], [32, 96, 384], [64, 128, 384], [256, 192, 1024](1024), [128, 96, 1024], [256, 64, 1024], [256, 96, 1024], [128, 32, 1024], [256, 96, 1024], [128, 32, 1024], [256, 96, 1024], [128, 32, 1024], [256, 96,$
LASP-30%	$ \begin{array}{l} 64, [0,0,64](64), [0,0,64], [0,32,64], [64,64,256](256), [64,64,256], [64,128,256], [64,96,256], [256,64,896](896), [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [0,0,896], [128,64,896], [128,66], [128,66], $
	ResNet101
LASP-60%	$ \begin{array}{l} 64, [64, 32, 192](192), [32, 32, 192], [64, 32, 192], [64, 128, 384](384), [64, 96, 384], [96, 96, 384], [128, 128, 384], [256, 256, 896](896), [256, 96, 896], [256, 128, 896], [256, 96, 896], [256, 64, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 96, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 896], [256, 128, 128, 128, 128, 128, 128, 128, 128$
LASP-50%	$ \begin{array}{l} 64, [64, 32, 128](128), [0, 32, 128], [64, 32, 128], [64, 128, 384](384), [32, 32, 384], [128, 96, 384], [128, 96, 384], [256, 256, 896](896), [256, 96, 896], [256, 9$
LASP-40%	$ \begin{array}{l} 64, [64, 32, 128](128), [0, 32, 128], [0, 0, 128], [64, 128, 384](384), [32, 64, 384], [64, 96, 384], [64, 128, 384], [256, 256, 896](896), [0, 64, 896], [128, 96, 896], [0, 64, 896], [0, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 64, 896], [128, 364, 896], [128, 64, 896], [126, 64, 896], [128, 64, 896], [128, 364, 896],$
LASP-30%	$ \begin{array}{l} 64, [64, 32, 128](128), [0, 0, 128], [0, 0, 128], [64, 128, 256](256), [0, 96, 256], [64, 96, 256], [64, 128, 256], [256, 192, 768](768), [0, 64, 768], [128, 64, 768], [0, 4, 768], [0, 4, 768], [0, 64, 768], [128, 64, 768], [0, 64, 768], [128, 64, 768], [0, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [0, 0, 768], [0, 0, 768], [0, 0, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [128, 64, 768], [0, 0, 768], [0, 0, 768], [256, 64, 768], [256, 64, 768], [256, 2048] \end{array} $
	MobileNet-V1
LASP-60%	896, 14, 14, 32, 32, 64, 64, 64, 64, 192, 192, 192, 192, 384, 384, 320, 320, 384, 384, 384, 384, 384, 384, 448, 448
LASP-42%	832, 16, 16, 32, 32, 32, 32, 32, 32, 32, 64, 64, 128, 128, 320, 320, 256, 256, 256, 256, 256, 256, 256, 320, 320, 320, 320, 896, 896, 896, 896, 896, 896, 896, 896
	MobileNet-V2
LASP-75%	16, 16, 16, 64, 64, 24, 64, 24, 112, 112, 32, 128, 128, 32, 128, 128, 32, 192, 192, 64, 384, 384, 64, 352, 352, 64, 352, 352, 64, 384, 384, 96, 512, 512, 96, 512, 512, 96, 512, 512, 96, 576, 576, 160, 960, 960, 160, 960, 960, 160, 960, 960, 320, 1280
LASP-60%	16, 16, 8, 32, 32, 16, 16, 16, 64, 64, 32, 32, 32, 32, 32, 64, 64, 32, 176, 176, 64, 288, 288, 64, 320, 320, 64, 320, 320, 64, 384, 384, 96, 448, 448, 96, 448, 448, 96, 576, 576, 160, 960, 960, 160, 960, 960, 160, 960, 960, 192, 1152

948 **R** Discussion on the augmented knapsack solver

Our augmented knapsack solver in Algo. 1 is modified based on the standard dynamic programming 949 solution for the 0-1 knapsack problem [3, 42, 43]. With the original 0-1 knapsack problem formulation, 950 each neuron can be selected to be removed or kept independently. However, the fact is that 1), the 951 layer always has the same latency with the same number of channels remaining no matter which 952 neurons are selected; 2). the neurons with higher importance are favored to maximize the accuracy. 953 So in this latency-aware pruning problem, we add an additional constraint, namely, each neuron can 954 be selected only when all the more importance neurons in the same layer are already kept, leading to 955 Eq. 7. The original solution is modified in lines 6-9 in Algo. 1 accordingly, and is converted into a 956 greedy approximation selection to retain the same time complexity as the original problem. We'll 957 discuss the non-greedy solution later. 958

In the augmented solver lines 6-9, every time when we decide whether to keep a neuron, we not only compare the total importance score can be reached under the constraint, but also check the inclusion of the "preceding" neurons. This is a greedy selection as it does not consider the potential importance value that the following neurons in the same layer would bring if the current neuron is kept. The



Figure 13: Visualization of the pruned ResNet50 structure.

overall time complexity of the solution is $O(N \times C)$ where $N = \sum_{l=1}^{L} N_l$ is the total number of neuron groups in the network and C is the latency constraint. We also provide the non-greedy solution 963 964 in Algo. 2. In this solution, for each neuron we add a process calculating and comparing the potential 965 importance score that the layer would further bring if the current neuron is selected to be kept. This 966 brings additional $O(N_l)$ complexity for each neuron in layer l. As a result, the total time complexity 967 of the solution increases to $O(\sum_{l=1}^{L} N_l^2 \times C)$. We test both of the solutions and observe similar performance in the ImageNet experiments as shown in Tab. 16. We hypothesize that the efficacy of 968 969 the greedy approach suffices from the already decreasingly ranked neurons feeding into the solver 970 and the iterative nature during pruning. Thus, the greedy approximation solution Algo. 1 is applied in 971 our method to have a better pruning efficiency. 972

0						
Method	1G model		2G model		3G model	
	Algo. 1	Algo. 2	Algo. 1	Algo. 2	Algo. 1	Algo. 2
Top1 (%)	77.45	77.51	76.56	76.51	74.45	74.51
FPS (img/s)	1203	1185	1672	1688	2597	2524

Table 16: ResNet50 pruning results on ImageNet with the greedy method Algo. 1 and non-greedy method Algo 2

Algorithm 2 Non-greedy solution for Eq. 7

Input: Importance score $\{\mathcal{I}_l \in \mathbb{R}^{N_l}\}_{l=1}^L$ where \mathcal{I}_l is sorted descendingly; Neuron latency contribution $\{c_l \in \mathbb{R}^{N_l}\}_{l=1}^L$; Latency constraint C. 1: maxV $\in \mathbb{R}^{(C+1)}$, keep $\in \mathbb{R}^{L \times (C+1)}$ \triangleright maxV[c]: max importance under constraint c; keep[l, c]: # neurons to keep in layer l to achieve maxI[c] 1. max $i \in \mathbb{R}^{d}$, keep $\in \mathbb{R}$ 2. for $l = 1, \dots, L$ do 3. for $j = 1, \dots, N_l$ do 4. for $c = 1, \dots, C$ do $\begin{array}{l} v_{keep} = \mathcal{I}_l^j + \max \mathbf{V}[c-c_l^j], \ v_{prune} = \max \mathbf{V}[c] \\ \text{flag = False} \end{array}$ $v_l - j + 1, \dots, N_l \, \mathbf{do}$ $v_{potential} = \sum_{j'=j}^{p_l} \mathcal{I}_l^{j'} + \max V[c - \sum_{j'=j}^{p_l} c_l^{j'}] \qquad \triangleright \text{ calculate the potential score this layer would bring if keep this neuron.}$ if $v_{potential} > v_{prune}$ and keep $[l, c - c_l^j] == j - 1$ then \triangleright check if leads to higher score and more important neurons in layer are kept flag = True break
end if 5: \triangleright total importance can achieve under constraint c with object n being kept or not 6: 7: for $p_l = j + 1, \dots, N_l$ do 8: 9: 10: 11: 12: 12: 13: 14: 15: end for if flag == True then keep[l, c] = j, update_maxV[c] = v_{keep} 16: else 17: 18: $\mathrm{keep}[l,c] = \mathrm{keep}[l,c-1], \, \mathrm{update_maxV}[c] = v_{prune}$ end if 19: end for $maxV \leftarrow update_maxV$ end for 20: 21: 22: end for 23: 24: keep_n = to save the kept neurons in model 25. for $l = l_1, \dots, l_d$ 26. $p_l = \text{keep}[l, C]$ 27. keep_n \leftarrow keep_n $\cup \{p_l \text{ top ranked neurons in layer } l\}$ ▷ retrieve the set of kept neurons 28: $C \leftarrow C - \sum_{j=1}^{p_l} c_l^j$ 29: end for Output: Kept important neurons (keep_n).