

## Appendix

### A Algorithm Procedure

A full workflow of channel number search with the proposed transitionaly APS is shown in Algorithm 1. The overall procedure consists of two stages. In the first stage, we fix the controller  $\pi(\theta)$  and initialize  $\Phi$  of affine parameter sharing with maximal value, so that meta weights  $\mathcal{W}$  can be efficiently optimized. During this stage architectures are uniformly sampled from the search space  $\mathcal{C}$  and are thus equally updated. In the second stage, we gradually anneal  $\Phi$  via Equation 5 so as to transit the sharing level  $\Phi$ . We alternatively update meta weights  $\mathcal{W}$  and controller parameter  $\theta$  based on architectures sampled from the controller.

---

#### Algorithm 1 RL-based CNS algorithm with Transitionaly APS

---

**Input:**

Training data  $\mathcal{D}_{tr}$ , validation data  $\mathcal{D}_{val}$ ;  
 Base network with meta weights  $\mathcal{W}$ , transformation matrices  $\mathcal{P}, \mathcal{Q}$   
 RL controller  $\pi(\theta)$  and channel search space  $\mathcal{C}$ .

**Output:**

Optimal channel configurations.

// Stage 1: fast optimization of meta-weights  $\mathcal{W}$

- 1: Initialize  $\mathcal{P}, \mathcal{Q}$  with maximal  $\Phi$ ;
  - 2: **for**  $t = 1, \dots, T_1$  **do**
  - 3:   Sample the architecture  $a$  uniformly from  $\mathcal{C}$ ;
  - 4:   Update  $\mathcal{W}$  via gradient descent with  $a$  on  $\mathcal{D}_{tr}$ ;
  - 5: **end for**
  - // Stage 2: transitionaly affine parameter sharing
  - 6: **for**  $t = 1, \dots, T_2$  **do**
  - 7:   Sample the architecture from controller  $a \sim \pi(\theta)$ ;
  - 8:   Update  $\mathcal{W}$  via gradient descent with  $a$  on  $\mathcal{D}_{tr}$ ;
  - 9:   Update  $\theta^{t+1} = \theta^t + \eta \mathbb{E}_a [\nabla_{\theta} \log p(a) \mathcal{R}]$  on  $\mathcal{D}_{val}$ ;
  - 10:   Anneal  $\Phi$  by updating  $\mathcal{P}, \mathcal{Q}$  in Equation (5);
  - 11: **end for**
- 

### B Proof for Theorem 3.1

**Theorem.** For  $\forall i \leq \tilde{i}$  and  $\forall o \leq \tilde{o}$ , the overall level  $\Phi$  of APS is maximized if  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}(\mathbf{Q}^{\tilde{i}})$  and  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}(\mathbf{P}^{\tilde{o}})$ .  $\Phi$  is minimized if  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}^{\perp}(\mathbf{Q}^{\tilde{i}})$  and  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}^{\perp}(\mathbf{P}^{\tilde{o}})$ .

To prove Theorem 3.1, we first show the case of two candidate decision  $(i, o)$  and  $(\tilde{i}, \tilde{o})$ , after which we can combine the pairwise optimal conditions together to yield Theorem 3.1. Without loss of generality, suppose  $c_{\tilde{i}} > c_i$  and  $c_{\tilde{o}} > c_o$ , we have the following lemma:

**Lemma 1.** Given candidate decisions  $(i, o)$  and  $(\tilde{i}, \tilde{o})$ ,  $\phi(i, o; \tilde{i}, \tilde{o})$  is maximized if  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}(\mathbf{Q}^{\tilde{i}})$  and  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}(\mathbf{P}^{\tilde{o}})$ ;  $\phi(i, o; \tilde{i}, \tilde{o})$  is minimized if  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}^{\perp}(\mathbf{Q}^{\tilde{i}})$  or  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}^{\perp}(\mathbf{P}^{\tilde{o}})$ .

*Proof.* As  $\phi(i, o; \tilde{i}, \tilde{o})$  is defined as the squared sum of pairwise correlation between  $\mathbf{W}^{i,o}$  and  $\mathbf{W}^{\tilde{i},\tilde{o}}$ , we can explicitly write it out as:

$$\begin{aligned} \phi(i, o; \tilde{i}, \tilde{o}) &= \sum_{x=1}^{c_i} \sum_{y=1}^{c_o} \sum_{\tilde{x}=1}^{c_{\tilde{i}}} \sum_{\tilde{y}=1}^{c_{\tilde{o}}} \left[ \text{Cov}_{x,y,\tilde{x},\tilde{y}}(\mathbf{W}^{i,o}, \mathbf{W}^{\tilde{i},\tilde{o}}) \right]^2 \\ &= \sum_{x=1}^{c_i} \sum_{y=1}^{c_o} \sum_{\tilde{x}=1}^{c_{\tilde{i}}} \sum_{\tilde{y}=1}^{c_{\tilde{o}}} \left[ \mathbb{E}(W_{x,y}^{i,o} W_{\tilde{x},\tilde{y}}^{\tilde{i},\tilde{o}}) - \mathbb{E}(W_{x,y}^{i,o}) \mathbb{E}(W_{\tilde{x},\tilde{y}}^{\tilde{i},\tilde{o}}) \right]^2. \end{aligned} \quad (6)$$

Note that the second term can be removed since  $\mathbb{E}(W_{x,y}^{i,o}) = \mathbb{E}((\mathbf{q}_x^i)^\top \mathbf{W} \mathbf{p}_y^o) = (\mathbf{q}_x^i)^\top \mathbb{E}(\mathbf{W}) \mathbf{p}_y^o = 0$  and similarly  $\mathbb{E}(W_{\tilde{x},\tilde{y}}^{\tilde{i},\tilde{o}}) = 0$ . Therefore  $\phi(i, o; \tilde{i}, \tilde{o})$  can be simplified as

$$\begin{aligned}
\phi(i, o; \tilde{i}, \tilde{o}) &= \sum_{x=1}^{c_i} \sum_{y=1}^{c_o} \sum_{\tilde{x}=1}^{c_{\tilde{i}}} \sum_{\tilde{y}=1}^{c_{\tilde{o}}} \mathbb{E}^2 \left[ W_{x,y}^{i,o} \cdot W_{\tilde{x},\tilde{y}}^{\tilde{i},\tilde{o}} \right] \\
&= \sum_{x,y} \sum_{\tilde{x},\tilde{y}} \mathbb{E}^2 \left[ (\mathbf{q}_x^i)^\top \mathbf{W} \mathbf{p}_y^o \cdot (\mathbf{q}_{\tilde{x}}^{\tilde{i}})^\top \mathbf{W} \mathbf{p}_{\tilde{y}}^{\tilde{o}} \right] \\
&= \sum_{x,y} \sum_{\tilde{x},\tilde{y}} \mathbb{E}^2 \left[ (\mathbf{q}_x^i)^\top \mathbf{W} \mathbf{p}_y^o \cdot (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{W}^\top \mathbf{q}_{\tilde{x}}^{\tilde{i}} \right] \\
&= \sum_{x,y} \sum_{\tilde{x},\tilde{y}} \left( (\mathbf{q}_x^i)^\top \mathbb{E} \left[ \mathbf{W} \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{W}^\top \right] \mathbf{q}_{\tilde{x}}^{\tilde{i}} \right)^2. \tag{7}
\end{aligned}$$

Expanding the expectation  $\mathbb{E}[\mathbf{W} \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{W}^\top]$  elementwisely, we have

$$\mathbb{E} \begin{bmatrix} \mathbf{w}_1 \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{w}_1^\top & \cdots & \mathbf{w}_1 \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{w}_c^\top \\ \vdots & \ddots & \vdots \\ \mathbf{w}_c \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{w}_1^\top & \cdots & \mathbf{w}_c \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{w}_c^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{p}_y^o)^\top \mathbf{p}_{\tilde{y}}^{\tilde{o}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & (\mathbf{p}_y^o)^\top \mathbf{p}_{\tilde{y}}^{\tilde{o}} \end{bmatrix} \tag{8}$$

where we have used the fact that  $\mathbb{E}[\mathbf{w}_m \mathbf{p}_y^o (\mathbf{p}_{\tilde{y}}^{\tilde{o}})^\top \mathbf{w}_n^\top] = (\mathbf{p}_y^o)^\top \mathbb{E}[\mathbf{w}_m^\top \mathbf{w}_n] \mathbf{p}_{\tilde{y}}^{\tilde{o}} = (\mathbf{p}_y^o)^\top \mathbf{p}_{\tilde{y}}^{\tilde{o}}$  if  $m = n$ , and 0 otherwise. With Equation 8, Equation 7 can be simplified to

$$\begin{aligned}
\phi(i, o; \tilde{i}, \tilde{o}) &= \sum_{x,\tilde{x}} \sum_{y,\tilde{y}} \left[ (\mathbf{q}_x^i)^\top \mathbf{q}_{\tilde{x}}^{\tilde{i}} \cdot (\mathbf{p}_y^o)^\top \mathbf{p}_{\tilde{y}}^{\tilde{o}} \right]^2 \\
&= \sum_{x,\tilde{x}} \sum_{y,\tilde{y}} \left[ \sum_{m=1}^c q_{m,x}^i q_{m,\tilde{x}}^{\tilde{i}} \cdot \sum_{n=1}^c p_{n,y}^o p_{n,\tilde{y}}^{\tilde{o}} \right]^2. \tag{9}
\end{aligned}$$

Without loss of generality, we take  $\mathbf{Q}^i$  and  $\mathbf{P}^{\tilde{o}}$  as standard orthogonal basis, i.e.  $q_{\tilde{x},\tilde{x}}^{\tilde{i}} = 1$  and  $q_{m,\tilde{x}}^{\tilde{i}} = 0$  for  $m \neq \tilde{x}$  and  $\tilde{x} \in \{1, \dots, c_{\tilde{i}}\}$ , and similarly for  $\mathbf{P}^{\tilde{o}}$ . Thus Equation 9 can be further reduced to

$$\sum_{x,\tilde{x}} \sum_{y,\tilde{y}} \left[ q_{\tilde{x},x}^i \cdot p_{\tilde{y},y}^o \right]^2 = \sum_{x=1}^{c_i} \sum_{y=1}^{c_o} \left( \sum_{\tilde{x}=1}^{c_{\tilde{i}}} (q_{\tilde{x},x}^i)^2 \right) \cdot \left( \sum_{\tilde{y}=1}^{c_{\tilde{o}}} (p_{\tilde{y},y}^o)^2 \right) \leq \sum_x \sum_y 1 = c_i c_o. \tag{10}$$

The equality holds if  $q_{m,x}^i = 0$  for  $m > c_{\tilde{i}}$  and  $p_{n,y}^o = 0$  for  $n > c_{\tilde{o}}$ . Therefore the maximum is attained when orthogonal basis of  $\mathbf{Q}^i$  and  $\mathbf{P}^o$  lie in the span of those in  $\mathbf{Q}^{\tilde{i}}$  and  $\mathbf{P}^{\tilde{o}}$  respectively, i.e.  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}(\mathbf{Q}^{\tilde{i}})$  and  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}(\mathbf{P}^{\tilde{o}})$ .

Conversely, minimum for Equation 9 is attained if  $q_{m,x}^i = 0$  for  $m \leq c_{\tilde{i}}$  or  $p_{n,y}^o = 0$  for  $n \leq c_{\tilde{o}}$ , which is equivalent to  $\mathcal{R}(\mathbf{Q}^i) \subseteq \mathcal{R}^\perp(\mathbf{Q}^{\tilde{i}})$  or  $\mathcal{R}(\mathbf{P}^o) \subseteq \mathcal{R}^\perp(\mathbf{P}^{\tilde{o}})$ .  $\square$

Finally, to prove Theorem 3.1, we only need to extend Lemma 1 to the case of multiple candidate decisions. The maximum and minimum of  $\Phi = \sum_{i,\tilde{i}} \sum_{o,\tilde{o}} \phi(i, o; \tilde{i}, \tilde{o})$  can be achieved when each  $\phi(i, o; \tilde{i}, \tilde{o})$  attains its maximum and minimum respectively. This corresponds to the intersection of optimal conditions in Lemma 1, which is exactly Theorem 3.1

## C Implementation Details and Hyper-parameters

**Design of RL controller** We follow ENAS [25] to take the RL-based algorithm to illustrate channel number search, and the overall formulation is given in Equation 1. Here we cover more details in the design. For the RL controller, we take a two-layer LSTM with 100 hidden units, and channel decisions are generated auto-regressively. The RL states contain the previous layer width decision, the one-hot layer index encoding, and available FLOPs left. Naively, one could simply set the accuracy

as the reward for controller training. However, this may be not reasonable under budget-constrained search. Instead we follow MNasNet [30] and design the reward  $\mathcal{R}$  as:

$$\mathcal{R} = Acc(a) \times \left[ \frac{\mathcal{B}(a)}{B} \right]^\gamma, \quad \text{where } \gamma = \begin{cases} \alpha & \text{if } \mathcal{B}(a) < B \\ \beta & \text{otherwise} \end{cases}, \quad (11)$$

where  $a \sim \pi(\theta)$  is the layerwise channel decision,  $Acc(\cdot)$  is the accuracy function, and  $\alpha, \beta$  are coefficients of FLOPs penalty. Such design is shown to approximate the pareto-optimal solutions [30]. We adopt policy gradient to maximize the reward function. To prevent the RL controller  $\pi(\theta)$  from getting stuck in local-minimal, we follow [25] to add an entropy regularization.

**Calculating the cosine similarity of gradients** Given two different candidates  $(i, o)$  and  $(\tilde{i}, \tilde{o})$  of one convolutional layer, we describe the steps to calculate the alignment of gradients ( $\cos(\mathbf{g}, \tilde{\mathbf{g}})$ ) on the meta-weight. We change the channel configuration in one layer and fix all the other layers. Specifically, for  $a = (a_1, a_2, \dots, a_L)$ , we vary  $a_l \in \mathcal{A}$  and fix all the rest  $a_j$  for  $j \neq l$ . Varying  $a_l \in \mathcal{A}$  gives  $A$  different gradients on  $\mathcal{W}$ . We thus compute and average the pairwise cosine similarities of these gradients at different values of  $\Phi$ , each of which involves  $A(A-1)/2$  combinations. Then these averaged cosine values are summed over all layers to produce Figure 3(a). During the transition of  $\Phi$ , we repeat the above procedure at each training epoch to plot the cosine values in Figure 4.

**Calculating the norm of coupled gradients** To collect the norm of coupled gradients, we maintain an accumulator for each candidate. During the searching process at time step  $t$ , each gradient term  $\mathbf{Q}^t (\nabla_{\mathcal{W}} \mathcal{L}(\mathbf{W}^t)) (\mathbf{P}^t)^\top$  from  $a_l$  is added to its own accumulator. The corresponding coupled gradients can be computed by summing over all candidate accumulators except for  $a_l$  and then multiply  $\mathbf{Q}^t$  and  $\mathbf{P}^t$ . Similar to the above subsection, the norm of coupled gradients across different layers are then averaged. The accumulators are reset every epoch to fully show the level of coupling change with respect to the parameter sharing transition. We follow such procedure at different values of  $\Phi$  to produce Figure 2(b). To plot the variation of coupled gradients in Figure 4, we repeat this process at every training epoch, and clean accumulators before the start of next epoch. We also empirically observe that the step-wise change of learning rates affects the observation of coupling. Thus we set constant learning rate as 1e-2 during the collection of coupled gradients.

**Efficient optimization of  $\Phi$**  The form of objective function  $\Phi$  w.r.t to  $\mathcal{P}, \mathcal{Q}$  can be readily obtained by vectorizing the element-wise summation in Equation 9. Direct optimization of the sharing level  $\Phi$  could be computationally expensive. Instead, we alternatively update  $\mathcal{P}, \mathcal{Q}$  to minimize  $\Phi$  with lower frequencies. For instance, we fix  $\mathcal{P}$  to optimize  $\mathcal{Q}$  for ten iterations, where the product terms of  $\mathcal{P}$  are computed in advance and can be reused them for multiple turns. In this way the computational burden can be effectively reduced. Furthermore, one can update  $\mathcal{P}, \mathcal{Q}$  separately for each layer, which is an equivalent implementation but with less memory overhead.

**Hyper-parameter settings** Finally, the detailed hyper-parameters for the RL based CNS algorithms on both CIFAR-10 and ImageNet datasets are shown in Table 3. Given the searched architecture, we follow default settings of the base models to train stand-alone models from scratch.

## D Visualization of Channel Configurations

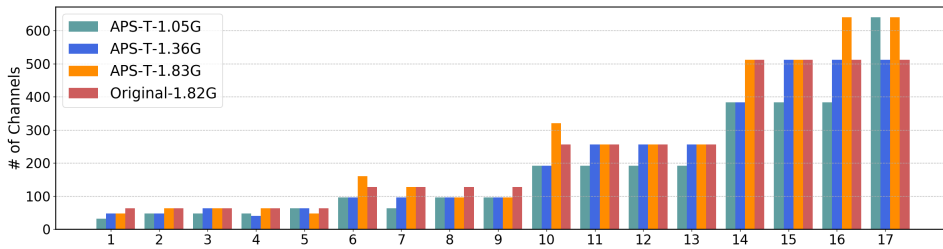
We visualize the searched channel configurations under different FLOPs constraint of ResNet-18 and MobileNet-v2 in Figure 10. Note that for MobileNet-v2, we omit the output channels of the depth-wise convolution in each block since it equals to the input channels. It can be observed that for both ResNet-18 and MobileNet-v2 under various FLOPs constraint, APS-T tends to find configurations with less channels in front layers and more channels in deep layers.

## E More Results on Unconstrained Search

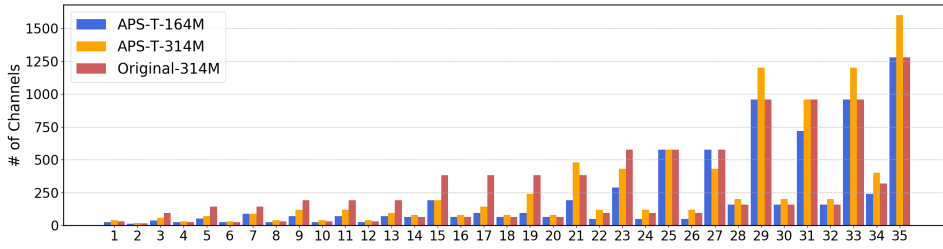
We provide more results on the architecture discrimination without FLOPs constraint in Figure 11. The hyper-parameters are consistent to those used Figure 5. It can be found that both APS-O and APS-T are relatively robust in the searched results. However, APS-I ends with large architecture discrimination but sometimes converges to sub-optimal configurations. We conjecture it is possibly due to early convergence of the controller as a result of insufficient training of meta weights  $\mathcal{W}$ .

Table 3: Hyper-parameters for different base models on CIFAR-10 and ImageNet.

Hyper-parameters	CIFAR-10		ImageNet	
	ResNet-20	ResNet-56	ResNet-18	MobileNet-v2
Channel Number $\mathcal{C}$	[4,8,16,32,64]	[4,8,16,32,64]	[32,48,64,80]	[8,12,16,20]
Width Multipliers	1	1	2	default
Max Channel Width $c$	208	208	128	32
Batch Size Per GPU	256	256	256	256
Init. Learning Rate of $\mathcal{W}$	1e-1	1e-1	1e-1	5e-2
Learning Rate Decay	Stepwise	Stepwise	Cosine	Cosine
Optimizer	SGD	SGD	SGD	SGD
Momentum	0.9	0.9	0.9	0.9
Nestrov	False	False	True	True
Learning Rate of $\mathcal{P}, \mathcal{Q}$	1e-3	1e-3	1e-3	1e-3
Learning Rate of $\theta$	1.6e-4	1.6e-4	1.6e-4	1.6e-4
FLOPs Penalty $\alpha, \beta$	0, -0.1	0, -0.06	0, -0.1	0, -0.1
Entropy Regularization	4e-3	4e-3	5e-1	4e-1
Weight Decay	2e-4	2e-4	1e-4	4e-5
Warmup Epochs	200	200	80	80
Max Epochs	600	600	160	160

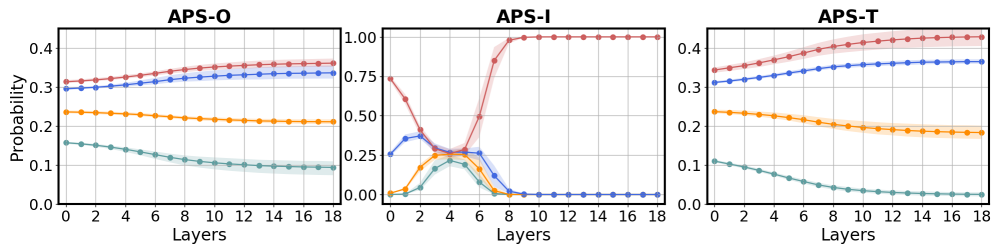


(a) ResNet-18



(b) MobileNet-v2

Figure 10: Channel configurations of ResNet-18 and MobileNet-v2 under different FLOPs constraint.



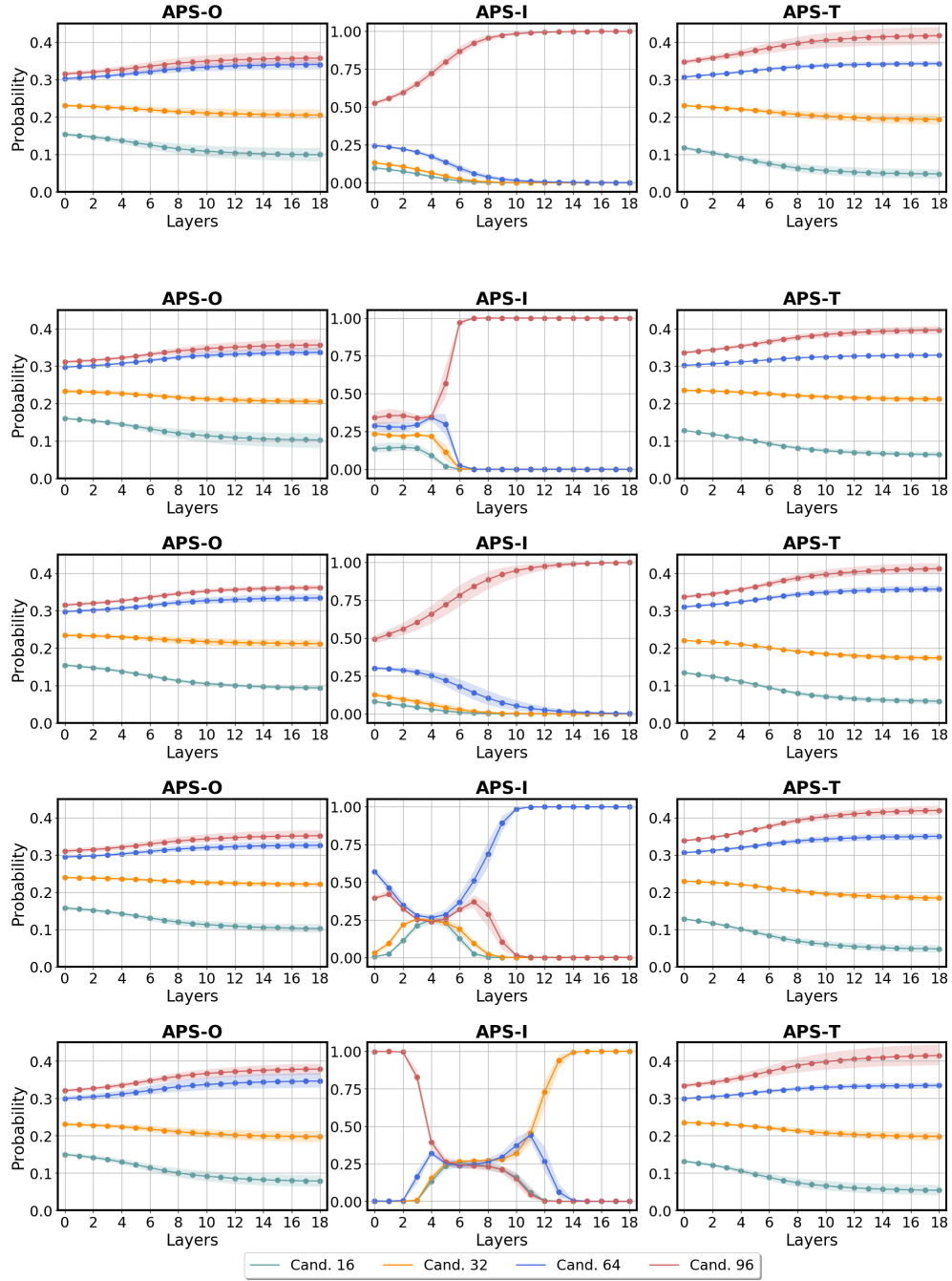


Figure 11: More searching results of APS-O, APS-I and APS-T without FLOPs constraint.