

Profiling Neural Blocks and Design Spaces for Mobile Neural Architecture Search

Keith G. Mills¹, Fred X. Han², Jialin Zhang³, Seyed Saeed Changiz Rezaei², Fabian Chudak², Wei Lu², Shuo Lian³, Shangling Jui³ and Di Niu¹

¹University of Alberta

²Huawei Technologies Canada Co., Ltd.

³Huawei Kirin Solution, Shanghai, China

Open-source repository: <https://github.com/Ascend-Research/BlockProfile>



Neural Architecture Search (NAS) automates neural network design and has achieved state-of-the-art results in numerous deep learning applications.

Three primary components: Design Space, Search Algorithm, Performance Estimation Strategy. Bulk of research done for the latter two.

We propose a simple methodology for profiling Design Spaces across valuable performance metrics, e.g., accuracy and inference latency, across different target hardware devices.

Devices considered: Huawei Kirin 9000 NPU, Nvidia RTX 2080 Ti GPU, AMD Threadripper 2990WX CPU and Samsung Note10.

Insights gleaned can be used to make pruned search spaces that outperform the original, including Once-for-All on MobileNetV3 (MBv3).

BLOCK-WISE SAMPLING AND PROFILING

Our method is rooted in the random sampling of architectures. To measure the impact of block b at layer l of unit u (denoted (u, l, b)), we sample many random architectures and affix b to (u, l) in each architecture, then measure the end-to-end response.

We can then iterate across all locations in the network, then average the results to get the general effect of block b on a desired metric.

		OFA MBConv Response																	
		Accuracy [%]			FLOPS [M]			RTX 2080 Ti [ms]			Huawei NPU-R224 [ms]			Huawei NPU-208 [ms]			Huawei NPU-192 [ms]		
Expansion Ratio	6	77.51	77.54	77.58	812	817	824	8.111	8.112	8.112	8.012	8.029	9.432	9.44	9.46	11.32	9.37	9.41	11.18
	4	77.44	77.48	77.49	797	800	805	8.111	8.110	8.111	7.961	7.999	8.539	9.38	9.42	10.04	9.32	9.36	10.01
	3	77.43	77.43	77.42	790	792	796	8.111	8.112	8.111	7.967	7.981	8.380	9.38	9.40	9.87	9.34	9.34	9.80
	Kernel Size	3	5	7	3	5	7	3	5	7	3	5	7	3	5	7	3	5	7
Expansion Ratio	6	413	413	413	405	406	405	380	380	379	42.12	42.66	43.41	37.23	37.69	38.34	31.06	31.45	31.99
	4	412	412	412	404	404	404	376	376	376	41.46	41.82	42.31	36.64	36.95	37.37	30.56	30.82	31.17
	3	409	409	409	401	401	401	376	376	376	41.14	41.41	41.81	36.35	36.59	36.92	30.33	30.53	30.80
	Kernel Size	3	5	7	3	5	7	3	5	7	3	5	7	3	5	7	3	5	7

RESULTS ON ONCE-FOR-ALL

Block-Wise Results on Once-for-All:

- Accuracy follows block size, FLOPS and visually correlated with Note10 latency.
- Latency on RTX 2080 Ti GPU close to constant.
- 2990WX CPU latency depends on expansion ratio (channels)
- NPU latency depends on kernel size.
- Lower resolution does not always mean lower latency.

EXPLOITING DISTRIBUTION DIFFERENCES

Pareto front performance depends on how different the accuracy and latency profiles are:

- NPU: Our insights outperform the original because Kernel Size 7 operations are not hardware friendly, but Kernel Size 5 is still very accuracy friendly.
- GPU: Negligible difference in block latency allows us to cut low-accuracy blocks with almost no additional latency.
- CPU: Latency differences exist but are very small.
- Note 10: Correlation means high latency operations are high accuracy operations. There is a trade-off.

WHAT IS A DESIGN SPACE?

The variable components of the network structure, primarily in the body. We abstract the Design Space to 3 levels:

- Units, u : Operate on unique tensor dimensions.
- Layers, l : Contain blocks, variable number within a unit.
- Blocks, b : The selectable operation sequences.
- Input image resolution size may also be considered.

We study three Design Spaces: Once-for-All, ProxylessNAS and ResNet50 and gain insights to block and layer sensitivity.

APPLICATION TO NAS

Derived insights can be used to prune a pre-defined space:

- Remove blocks from specific units or altogether.
- Limit the number of layers in a unit.
- Focus search on specific units.
- Optimize for one or more desired metrics.

Test insights using a simple random mutation algorithm.

- Compare pruned search space to the original.
- Accuracy-latency Pareto frontier optimization.
- Maximum accuracy search.

LATENCY CONSTRAINED PARETO FRONTIER SEARCH

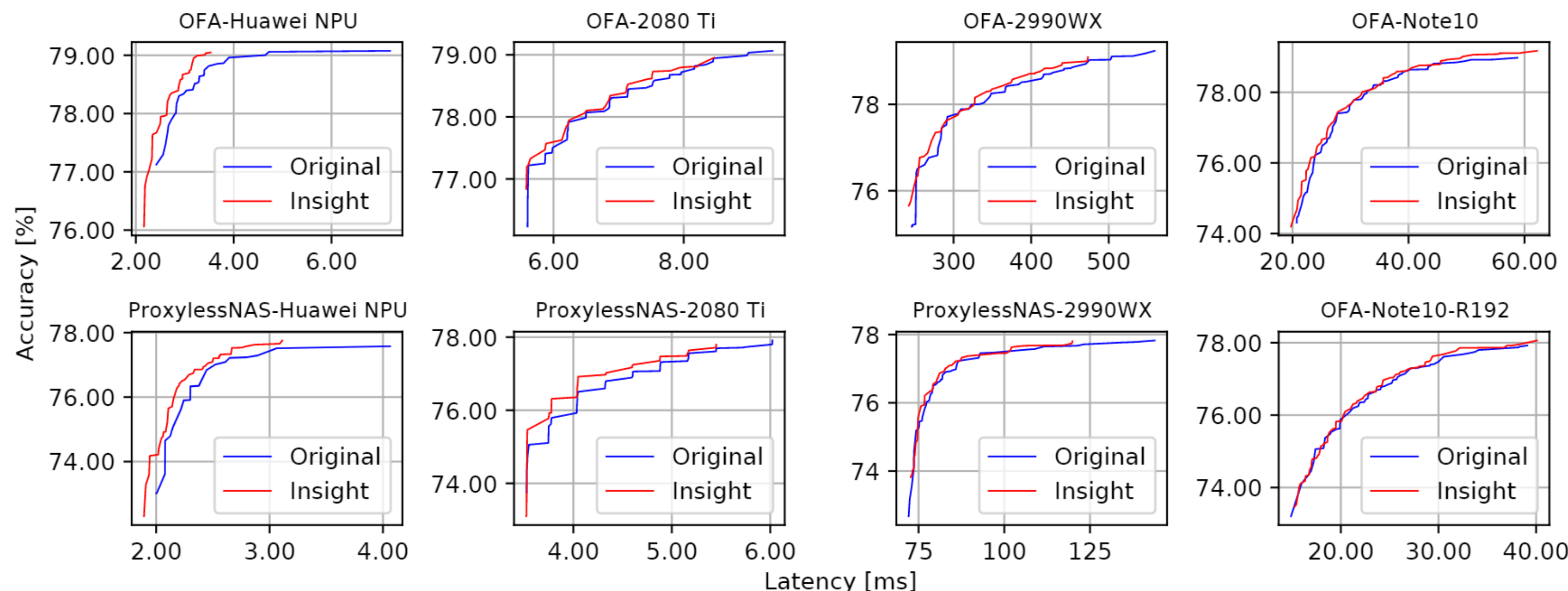


Figure 8: Pareto frontiers contrasting the original search spaces (blue) with our insight-based search spaces (red).

MAXIMUM ACCURACY SEARCH

Table 2: Maximum top-1 ImageNet accuracy search results on different design spaces, compared to existing works. We show averages over 5 random seeds for our experiments.

Model	Accuracy	MACs
MobileNetV2 [21]	72.0	300M
MobileNetV3-Large [10]	75.2	219M
OFA [2]	76.0	230M
OFA _{Large}	79.0	595M
OFA-insight	79.2 ± 0.04	342M
OFA-base	78.9 ± 0.07	292M
ProxylessNAS-insight	77.9 ± 0.04	417M
ProxylessNAS-base	77.6 ± 0.08	359M
ResNet50-insight	80.0 ± 0.03	2.81B
ResNet50-base	79.9 ± 0.09	2.64B

Without considering latency and just optimizing for accuracy, our insights on Once-for-All allow us to find a pruned search space that not only outperforms the original version, but the original OFA_{Large}!