



EfficientNet

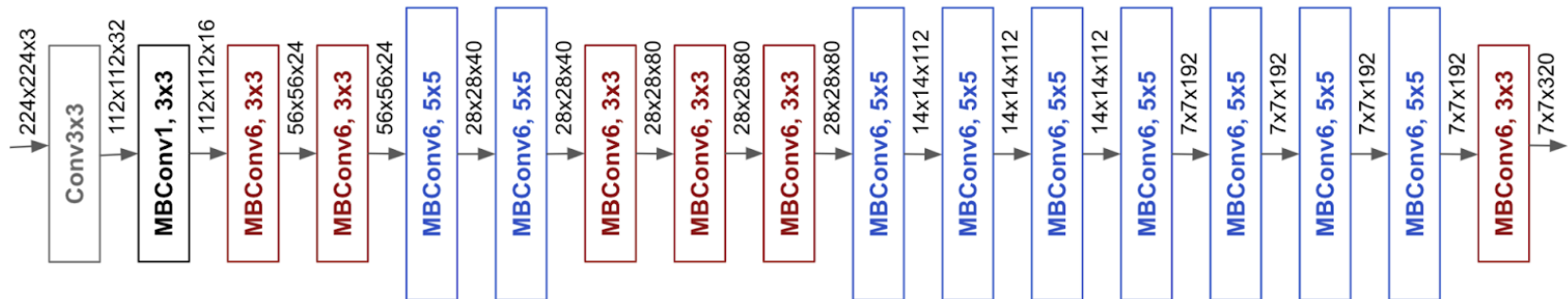
M.Tan et al, EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks, arXiv 2019

Scale Up

► Recent Trends of CNNs

► Repeating Base Blocks

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	28×28	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1



Scale Up

- ▶ **Scaling up ConvNets is widely used to achieve better accuracy.**
 - ▶ ResNet can be scaled from ResNet 18 to ResNet 200 by using more layers.
 - ▶ GPipe achieved 84.3% ImageNet top 1 accuracy by scaling up a baseline model 4 times larger.
- ▶ **What to Scale Up**
 - ▶ # of Layers: Depth
 - ▶ # of Channels: Width
 - ▶ Size of Input Images: Resolution

Scale Up

▶ What to Scale Up

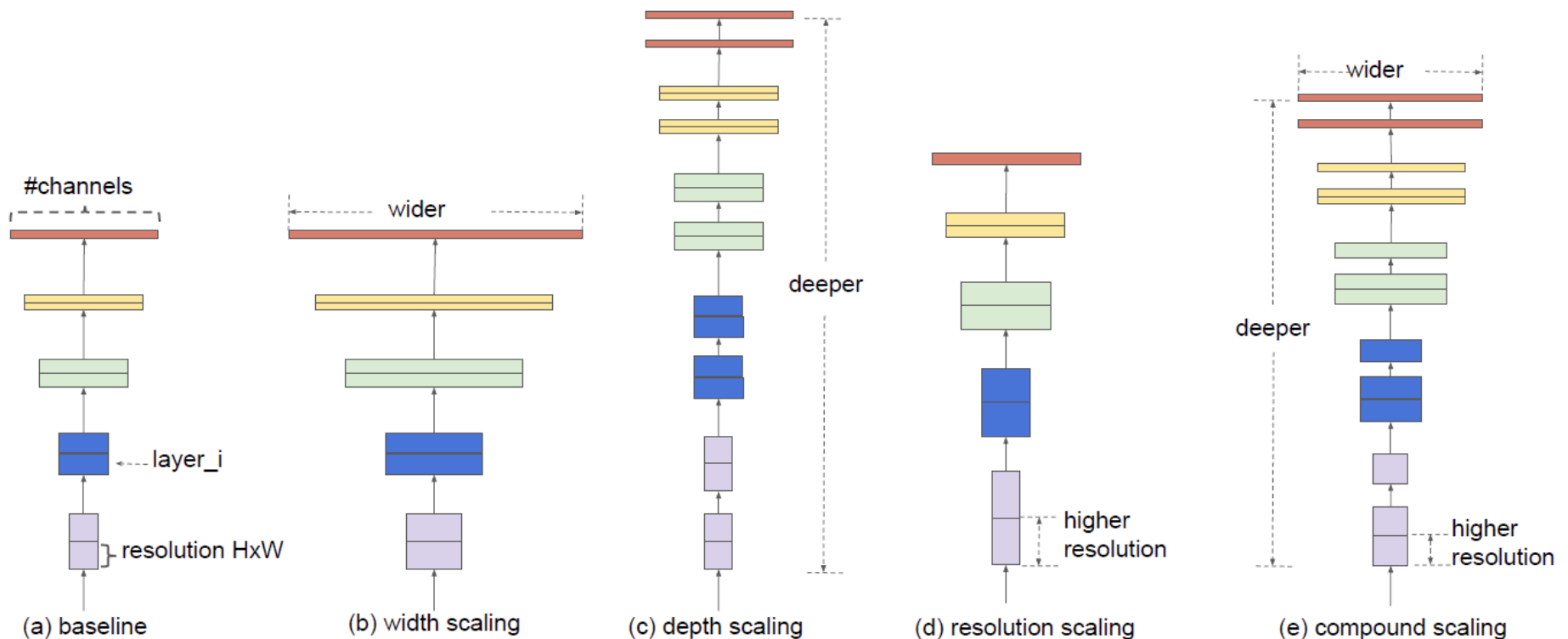
- ▶ # of Layers: Depth
- ▶ # of Channels: Width
- ▶ Size of Input Images: Resolution

▶ Why to Scale Up

- ▶ If a network has more layers
 - ▶ We can capture richer and more complex features
- ▶ If a network has more channels
 - ▶ We can have various patterns
- ▶ If an input image is bigger
 - ▶ We can use fine-grained patterns
 - ▶ Early networks used 224x224, but these days use 299x299, 331x331 or 480x480 (Gpipe)

Scale Up

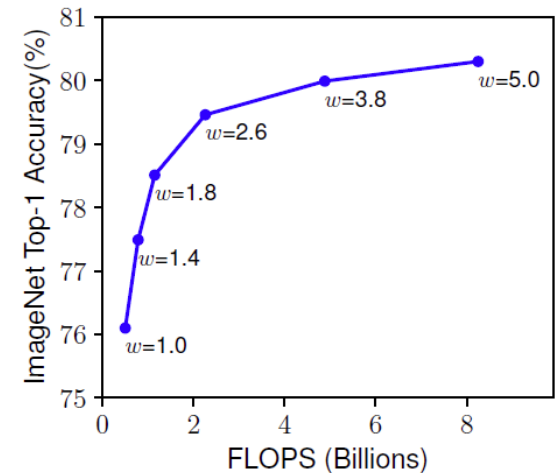
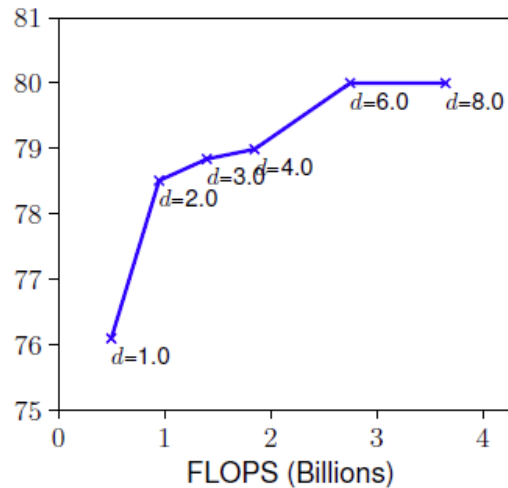
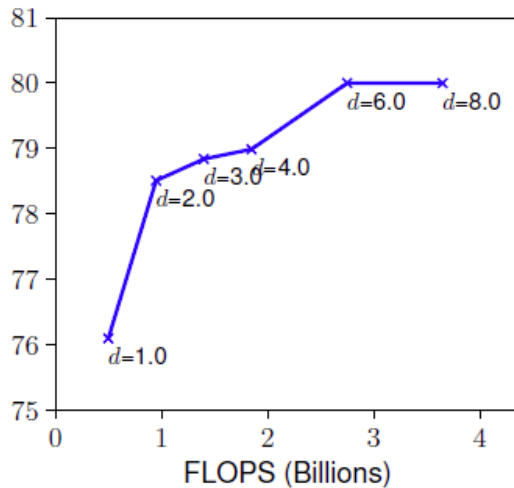
▶ Width, Depth, Resolution Scaling



Difficulties

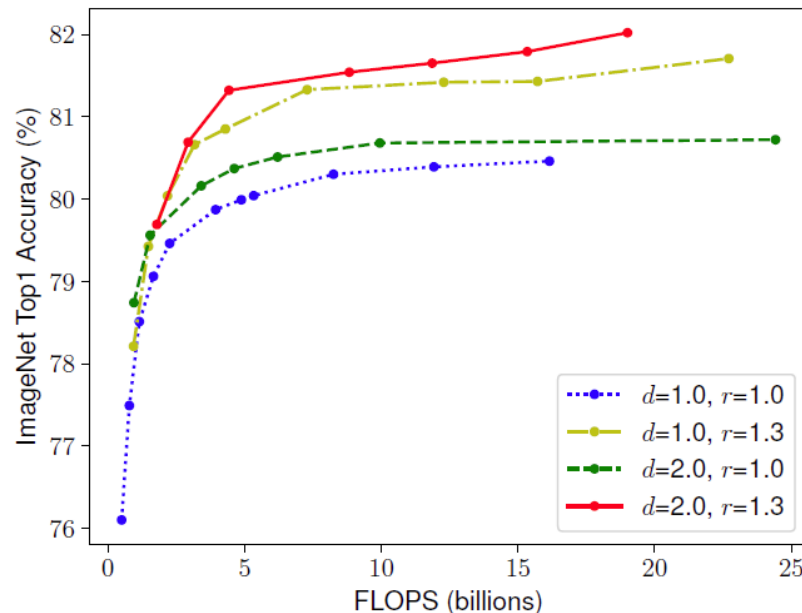
► Saturation

- ResNet 1000 has similar accuracy as ResNet 101 even though it has much more layers.
- Hard to capture good features if networks are shallow even though it is wider



Difficulties

- ▶ It is critical to balance width, depth and resolution
 - ▶ Scaling width w without changing depth ($d=1.0$) and resolution ($r=1.0$) results in quick saturation
 - ▶ With deeper ($d=2.0$) and higher resolution ($r=2.0$), width scaling achieves much better accuracy under the same FLOPS cost.



Idea for Best Compound Scaling

1. Find out a good baseline model
2. Find out the golden ratio of width, depth and resolution for scaling
3. Scaling up each dimension of the baseline model keeping the golden ratio of width, depth and resolution

Formulation

► Recap: CNN

$$\mathcal{N} = \bigodot_{i=1 \dots s} \mathcal{F}_i^{L_i} (X_{\langle H_i, W_i, C_i \rangle})$$

Example

s : stage,
 F_i : operation of stage i ,
 L_i : repetition of F_i ,
 $X_{\langle H_i, W_i, C_i \rangle}$: input

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	28×28	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1

Formulation

► Compound Scaling

$$\max_{d,w,r} \text{Accuracy}(\mathcal{N}(d, w, r))$$

$$s.t. \quad \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d \cdot \hat{L}_i} (X_{\langle r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i \rangle})$$

$$\text{Memory}(\mathcal{N}) \leq \text{target_memory}$$

$$\text{FLOPS}(\mathcal{N}) \leq \text{target_flops}$$

s : stage

F_i : operation of stage i ,

L_i : repetition of F_i ,

$X_{\langle H_i, W_i, C_i \rangle}$: input

d : scale factor of depth

w : scale factor of width

r : scale factor of resolution

► Assumption

- All stages and layers share the scaling factors to reduce the search space

Formulation

- ▶ **FLOPS of a CNN is proportional to d, w^2, r^2**
 - ▶ Doubling depth doubles FLOPS
 - ▶ Doubling width or resolution increases FLOPS by four times
- ▶ **So, following constraints are added to searching for d, w and r**

$$\text{depth: } d = \alpha^\phi$$

ϕ is a user specific parameter

$$\text{width: } w = \beta^\phi$$

$$\text{resolution: } r = \gamma^\phi$$

$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

Searching for Baseline Model

▶ By MNasNet

- ▶ a multi objective neural architecture search that optimizes both accuracy and FLOPS
- ▶ Optimization Goal :

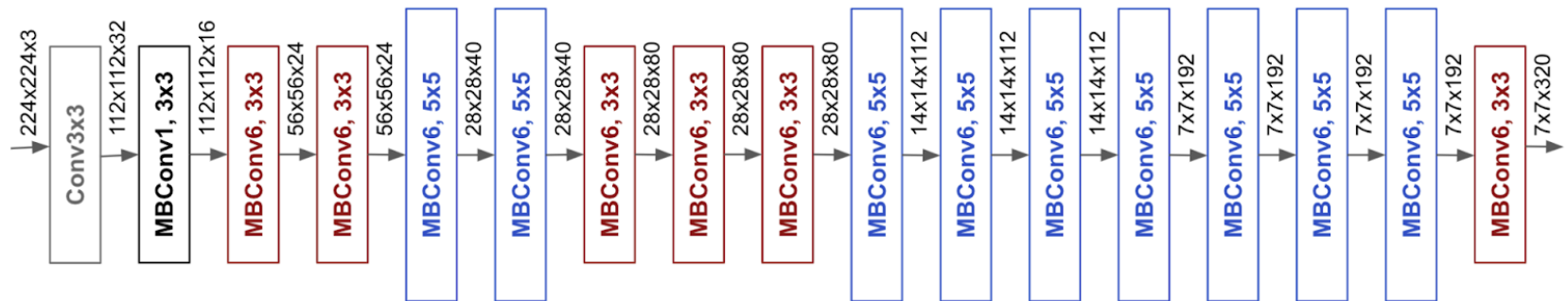
$$ACC(m) \times [FLOPS(m)/T]^w$$

- ▶ The found baseline model will be scaled up for better accuracy with less resources.

Searching for Baseline Model

► Found Baseline Model

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	28×28	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1



Golden Ratio and Model Scaling

- ▶ **Golden Ratio of α, β, γ by grid search**
 - ▶ Assuming twice more resources available, i.e., $\phi = 1$
 - ▶ The best values are $\alpha = 1.2, \beta = 1.1, \gamma = 1.15$
- ▶ **Scaling the Baseline Model**
 - ▶ Choosing ϕ considering available resources
 - ▶ Scale the base model by $w = \alpha^\phi, d = \beta^\phi, r = \gamma^\phi$
 - ▶ They chose 7 different values for ϕ , i.e., generated 7 different networks by scaling.

Performance

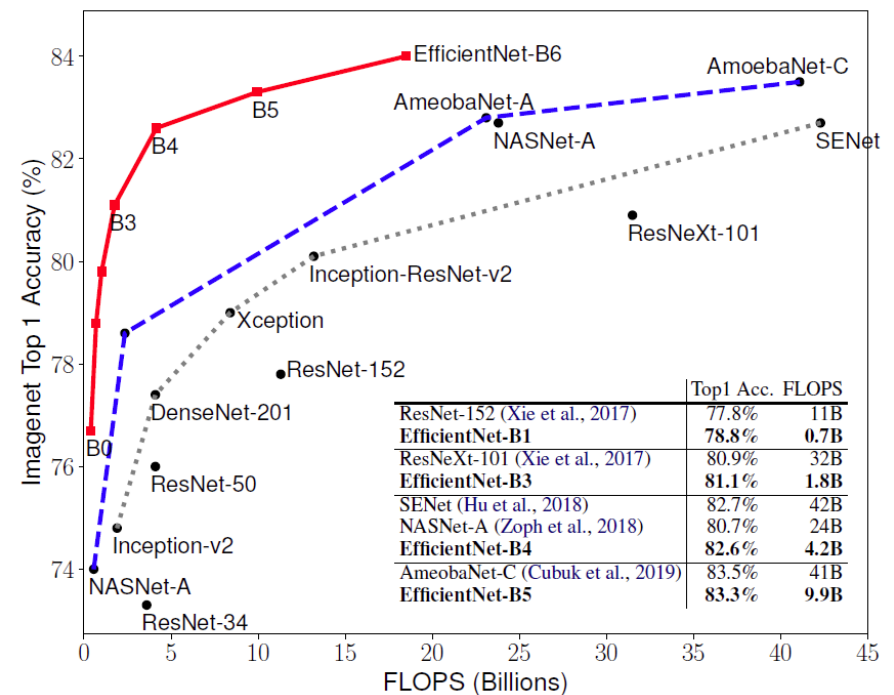
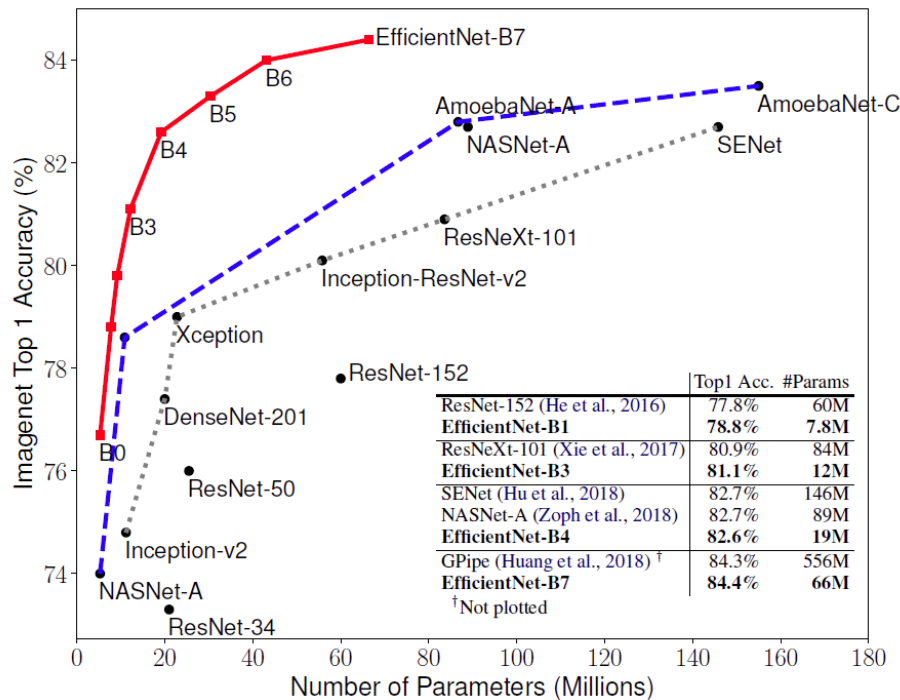
► ImageNet Classification

Model	Top-1 Acc.	Top-5 Acc.	#Params	Ratio-to-EfficientNet	#FLOPS	Ratio-to-EfficientNet
EfficientNet-B0	76.3%	93.2%	5.3M	1x	0.39B	1x
ResNet-50 (He et al., 2016)	76.0%	93.0%	26M	4.9x	4.1B	11x
DenseNet-169 (Huang et al., 2017)	76.2%	93.2%	14M	2.6x	3.5B	8.9x
EfficientNet-B1	78.8%	94.4%	7.8M	1x	0.70B	1x
ResNet-152 (He et al., 2016)	77.8%	93.8%	60M	7.6x	11B	16x
DenseNet-264 (Huang et al., 2017)	77.9%	93.9%	34M	4.3x	6.0B	8.6x
Inception-v3 (Szegedy et al., 2016)	78.8%	94.4%	24M	3.0x	5.7B	8.1x
Xception (Chollet, 2017)	79.0%	94.5%	23M	3.0x	8.4B	12x
EfficientNet-B2	79.8%	94.9%	9.2M	1x	1.0B	1x
Inception-v4 (Szegedy et al., 2017)	80.0%	95.0%	48M	5.2x	13B	13x
Inception-resnet-v2 (Szegedy et al., 2017)	80.1%	95.1%	56M	6.1x	13B	13x
EfficientNet-B3	81.1%	95.5%	12M	1x	1.8B	1x
ResNeXt-101 (Xie et al., 2017)	80.9%	95.6%	84M	7.0x	32B	18x
PolyNet (Zhang et al., 2017)	81.3%	95.8%	92M	7.7x	35B	19x
EfficientNet-B4	82.6%	96.3%	19M	1x	4.2B	1x
SENet (Hu et al., 2018)	82.7%	96.2%	146M	7.7x	42B	10x
NASNet-A (Zoph et al., 2018)	82.7%	96.2%	89M	4.7x	24B	5.7x
AmoebaNet-A (Real et al., 2019)	82.8%	96.1%	87M	4.6x	23B	5.5x
PNASNet (Liu et al., 2018)	82.9%	96.2%	86M	4.5x	23B	6.0x
EfficientNet-B5	83.3%	96.7%	30M	1x	9.9B	1x
AmoebaNet-C (Cubuk et al., 2019)	83.5%	96.5%	155M	5.2x	41B	4.1x
EfficientNet-B6	84.0%	96.9%	43M	1x	19B	1x
EfficientNet-B7	84.4%	97.1%	66M	1x	37B	1x
GPipe (Huang et al., 2018)	84.3%	97.0%	557M	8.4x	-	-

We omit ensemble and multi-crop models (Hu et al., 2018), or models pretrained on 3.5B Instagram images (Mahajan et al., 2018).

Performance

► ImageNet Classification



Performance

► ImageNet Classification

Table 4. Inference Latency Comparison – Latency is measured with batch size 1 on a single core of Intel Xeon CPU E5-2690.

Acc. @ Latency		Acc. @ Latency	
ResNet-152	77.8% @ 0.554s	GPipe	84.3% @ 19.0s
EfficientNet-B1	78.8% @ 0.098s	EfficientNet-B7	84.4% @ 3.1s
Speedup	5.7x	Speedup	6.1x

Performance

► Scaling Comparison

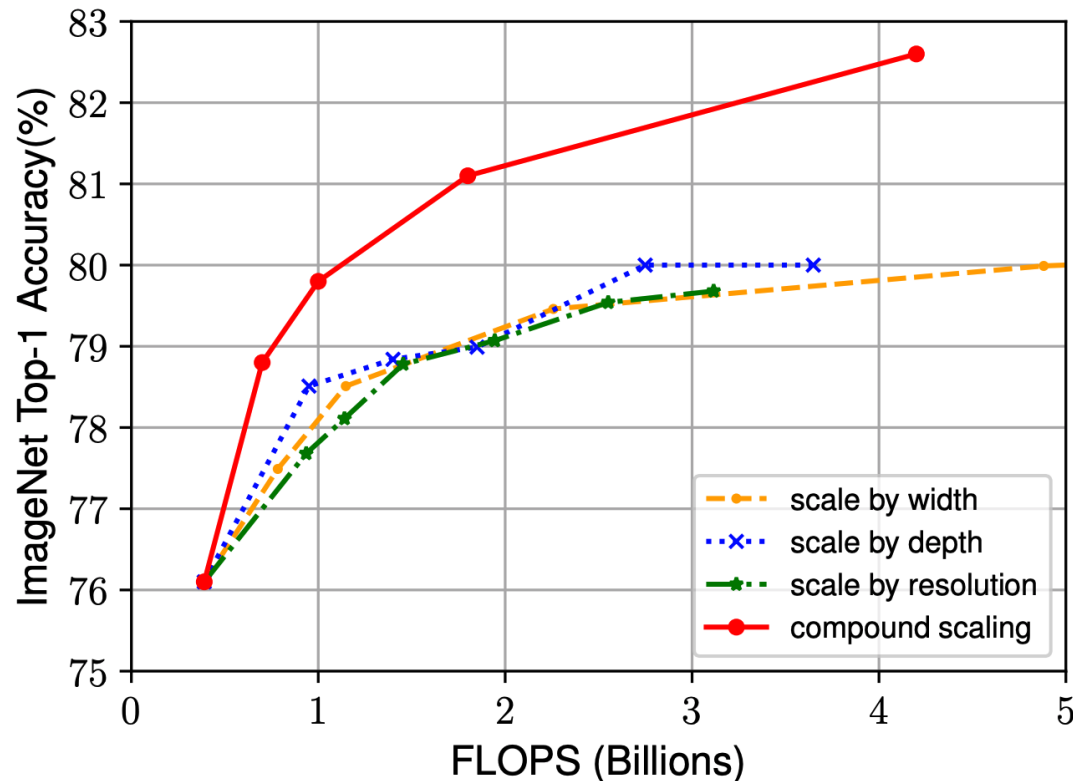


Figure 8. Scaling Up EfficientNet-B0 with Different Methods.

Performance

► Scaling Comparison

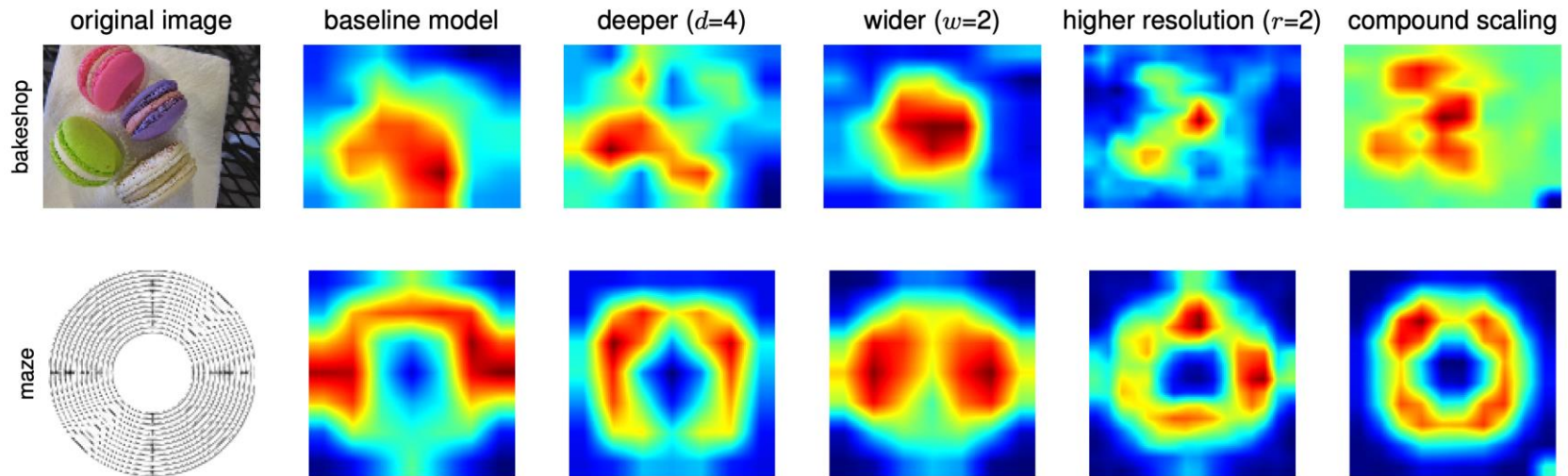


Figure 7. Class Activation Map (CAM) (Zhou et al., 2016) for Different Models in Table 7 - Our compound scaling method allows the scaled model (last column) to focus on more relevant regions with more object details. Model details are in Table 7.

Baseline: MobileNet and ResNet

- ▶ They have applied the same approach to scale MobileNet and ResNet

Model	FLOPS	Top-1 Acc.
Baseline MobileNetV1 (Howard et al., 2017)	0.6B	70.6%
Scale MobileNetV1 by width ($w=2$)	2.2B	74.2%
Scale MobileNetV1 by resolution ($r=2$)	2.2B	72.7%
compound scale ($d=1.4, w=1.2, r=1.3$)	2.3B	75.6%
Baseline MobileNetV2 (Sandler et al., 2018)	0.3B	72.0%
Scale MobileNetV2 by depth ($d=4$)	1.2B	76.8%
Scale MobileNetV2 by width ($w=2$)	1.1B	76.4%
Scale MobileNetV2 by resolution ($r=2$)	1.2B	74.8%
MobileNetV2 compound scale	1.3B	77.4%
Baseline ResNet-50 (He et al., 2016)	4.1B	76.0%
Scale ResNet-50 by depth ($d=4$)	16.2B	78.1%
Scale ResNet-50 by width ($w=2$)	14.7B	77.7%
Scale ResNet-50 by resolution ($r=2$)	16.4B	77.5%
ResNet-50 compound scale	16.7B	78.8%

Question and Answer