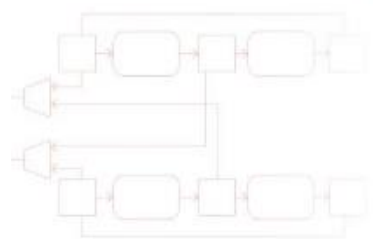
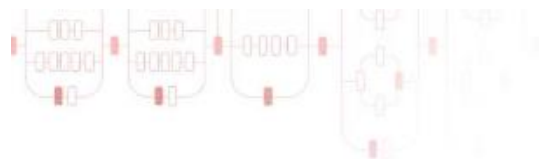


Efficient CNN

효율성





TensorFlow
Korea

가입함 ▾ 알림 [공유하기](#) [... 더 보기](#)

멤버 44,579명

멤버 찾기

관리자 및 댓글 관리자 9



이진원
관리자
최근 게시를 2개 보기
Samsung Electronics Staff Engineer



이지민
관리자
최근 게시를 1개 보기
서울대학교

메시지 보내기



Hwalsuk Lee
관리자
최근 게시를 1개 보기
Naver Corporation Research Engineer

메시지 보내기



pr12

공통 필터

PR12 Season 2

visionNoob • 업데이트: 2일 전

PR-101: Deep Feature Consistent Variational Autoencoder • 32:39

PR-102: Everybody Dance Now • 27:20

[모든 재생목록 보기](#)

PR12 Season 1

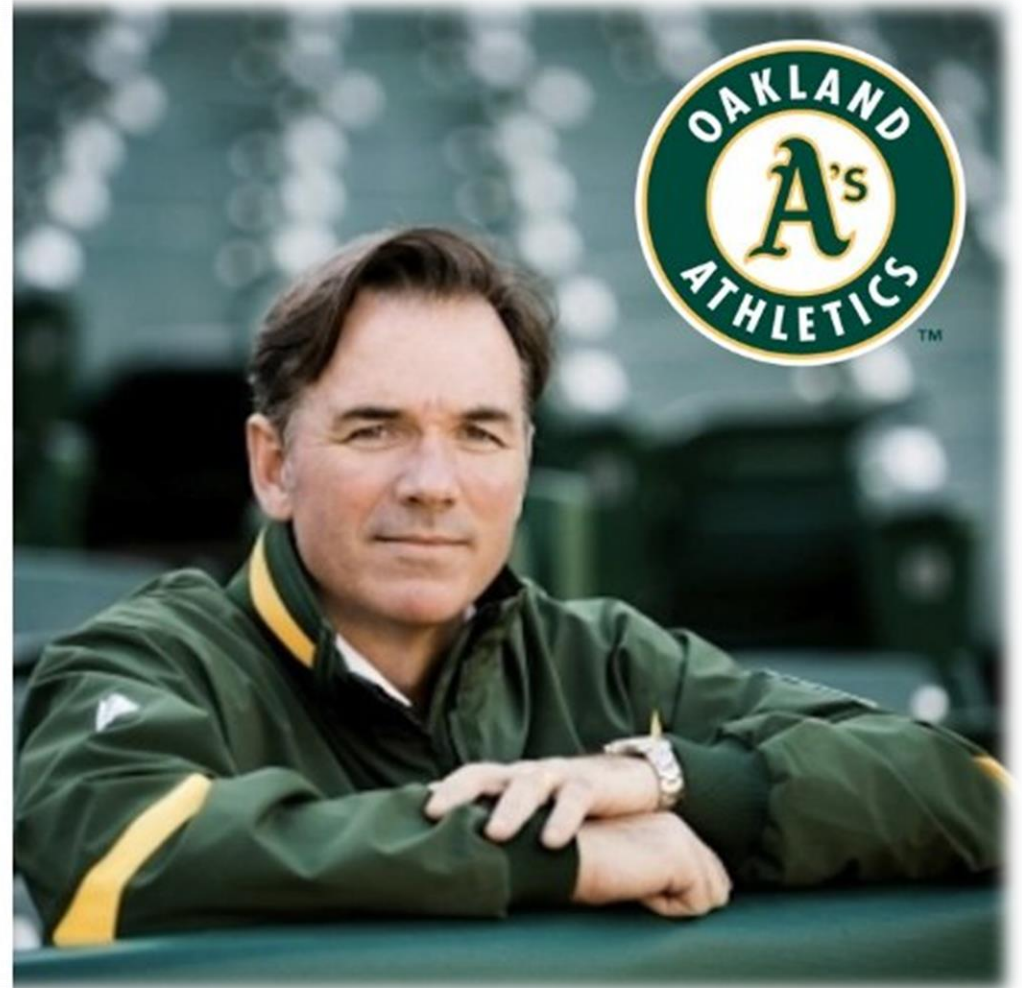
visionNoob

PR000: 논문 읽기 각오를 다집니다. • 9:16

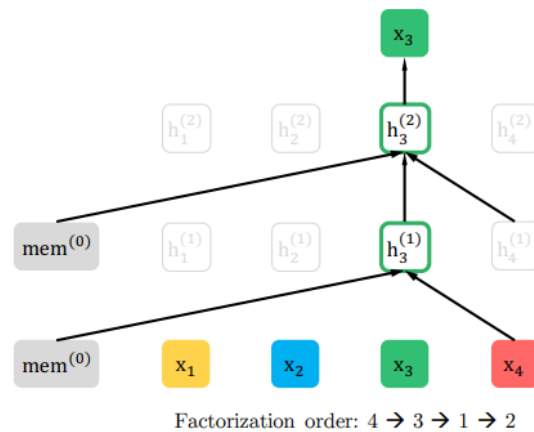
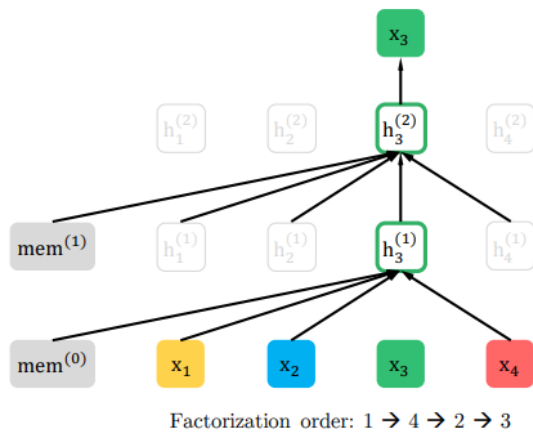
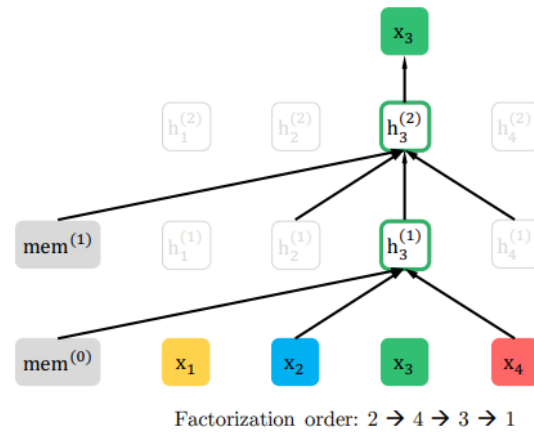
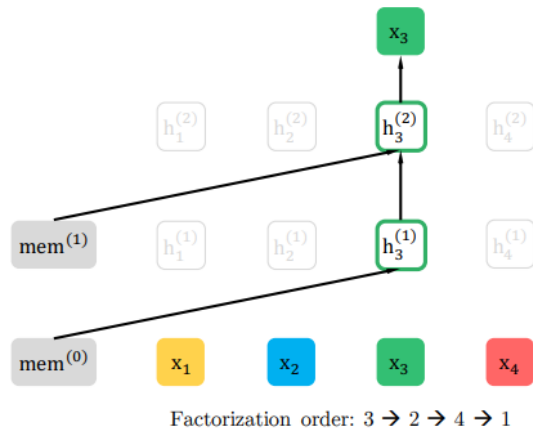
PR-001: Generative adversarial nets by Jaejun Yoo (2017/4/13) • 35:05

[모든 재생목록 보기](#)

MoneyBall

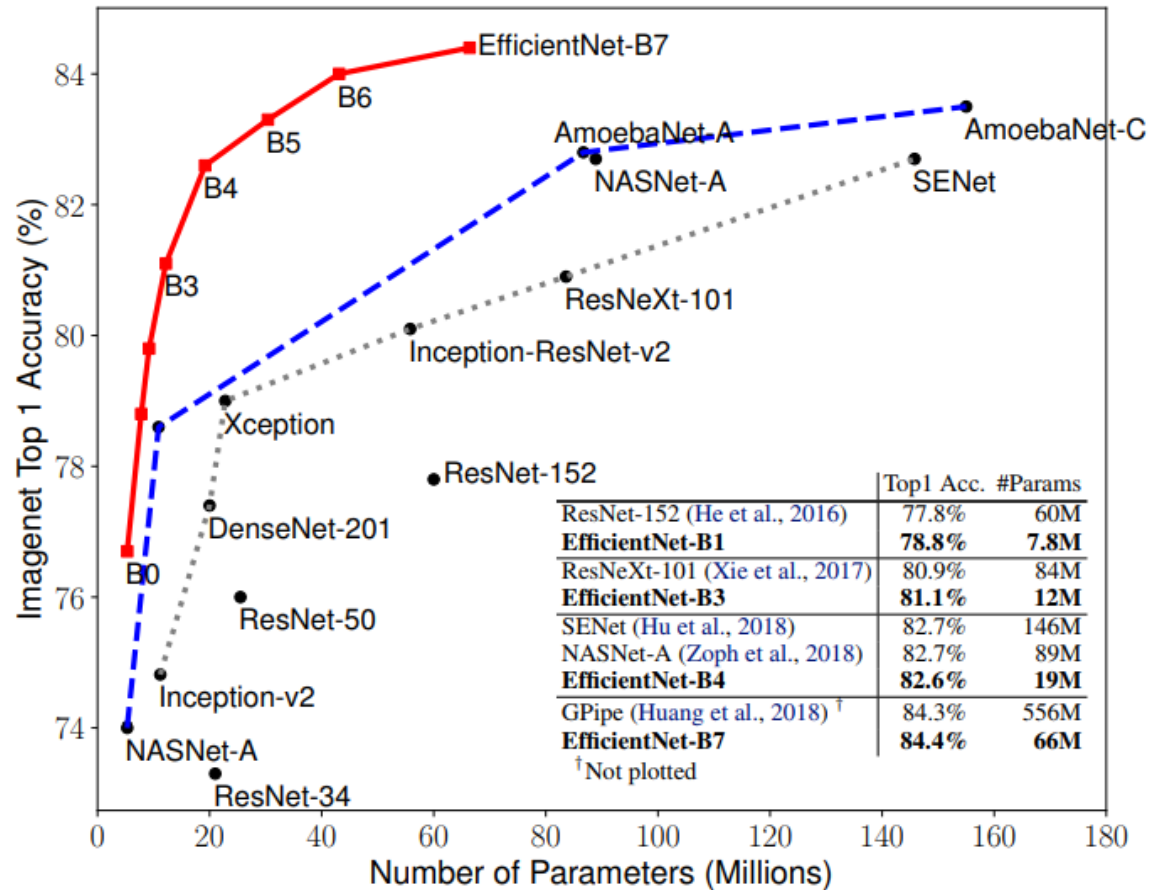


XLNet(Yang, arxiv 19 Jun 2019)



- Parameters
 - 340 million parameters
- Training
 - 512 TPU v3 chips for 500K steps
 - 2.5 days
- 512 TPU x 2.5 days x \$8
= \$245,000

Convolutional Neural Networks



- GPipe
 - 556 million parameters
- AmoebaNet
 - 450 K40 GPU, 7 days training
- NAS
 - 800 GPU, 28 days training

Toward More Accurate Models

- The impressive results of the current object category recognition systems are obtained on a limited number(20~1000) of basic-level object categories



Dog

Pomeranian



Bicycle

Road Bike



Car

Sports Car



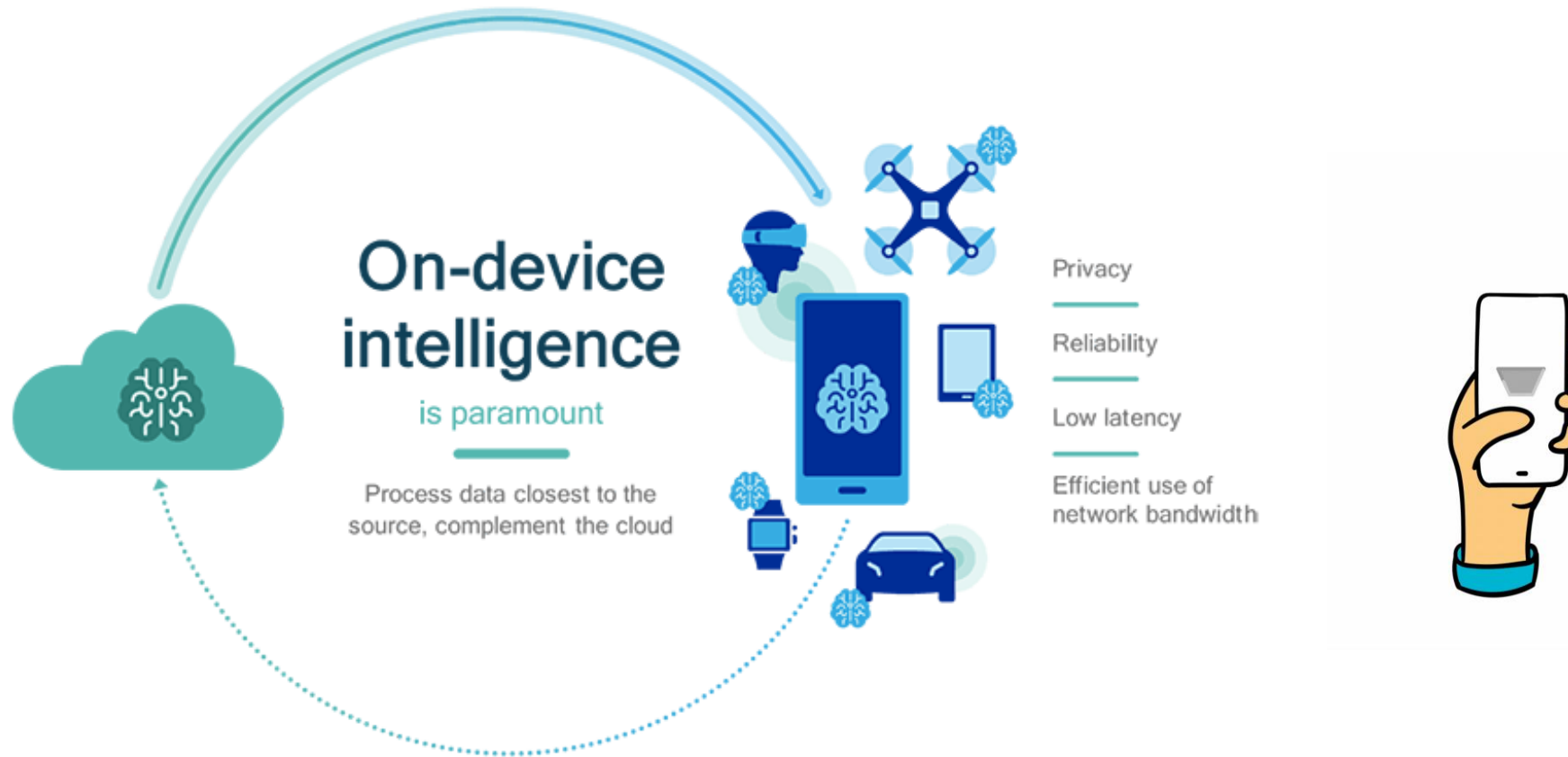
Bag

Shoulder bag

System	ImageNet 22,000 Top-1 Accuracy
[Le12]	13.6%
[Chilimbi14]	29.8%

- With large scale dataset with 22,000 categories, the performance drops to 30%

Efficient Models for On-device AI



그래서 오늘은...

Convolutional Neural Network

그 중에서도

Efficient CNN에는 어떤 model들이 있는지...

논문의 핵심 idea들 위주로
살펴보겠습니다

그리고 몇가지 bonus insight도...

오늘 다루지 않는 것들

- AutoML (TTT)
- Pruning
- Knowledge Distillation
- 오늘 소개할 논문들의 Experimental Results

오늘 살펴볼 논문들...

1. [Going Deeper with Convolutions](#)
2. [Very Deep Convolutional Networks for Large-Scale Image Recognition](#)
3. [SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and ...](#)
4. [SqueezeNext: Hardware-Aware Neural Network Design](#)
5. [MobileNets: Efficient Convolutional Neural Networks for Mobile Vision ...](#)
6. [MobileNetV2: Inverted Residuals and Linear Bottlenecks](#)
7. [Searching for MobileNetV3](#)
8. [ShuffleNet: An Extremely Efficient Convolutional Neural Network for ...](#)
9. [ShuffleNet V2: Practical Guidelines for Efficient CNN Architecture Design](#)
10. [Shift: A Zero FLOP, Zero Parameter Alternative to Spatial Convolutions](#)
11. [CondenseNet: An Efficient DenseNet using Learned Group ...](#)
12. [All You Need is a Few Shifts: Designing Efficient Convolutional Neural ...](#)
13. [Fully Learnable Group Convolution for Acceleration of Deep Neural ...](#)

깨알홍보 - 12PR



JinWon Lee

구독자 745명

채널 맞춤설정

YOUTUBE 스튜디오(베타)

홈

동영상

재생목록

채널

토론

정보



업로드한 동영상 모두 재생

정렬 기준

<p>Problem Formulation</p> <p>35:53</p>	<p>Randomly Wired Neural Networks</p> <p>39:04</p>	<p>The Fire Module</p> <p>39:18</p>	<p>Default Boxes and Aspect Ratios</p> <p>36:01</p>	<p>Analysis of the Runtime Performance (ShuffleNet v2 and MobileNet v2)</p> <p>33:34</p>	<p>Residual Blocks</p> <p>31:12</p>
---	--	-------------------------------------	---	--	-------------------------------------

PR-169: EfficientNet: Rethinking Model Scaling for Image Classification
조회수 613회 · 2주 전

PR-155: Exploring Randomly Wired Neural Networks for Efficient Image Classification
조회수 2.6천회 · 2개월 전

PR-144: SqueezeNext: Hardware-Aware Neural Architecture Search for Efficient Image Classification
조회수 766회 · 4개월 전

PR-132: SSD: Single Shot MultiBox Detector
조회수 2.6천회 · 5개월 전

PR-120: ShuffleNet V2: Practical Guidelines for Efficient Image Classification
조회수 1.2천회 · 7개월 전

PR-108: MobileNetV2: Inverted Residuals and Linear Bottlenecks
조회수 2천회 · 8개월 전

<p>Pentomino</p> <p>24:15</p>	<p>TPU Block Diagram</p> <p>38:29</p>	<p>NASNet-C</p> <p>32:14</p>	<p>Grouped Convolution</p> <p>33:13</p>	<p>Recap - Convolution Operation</p> <p>30:56</p>	<p>Recap - Faster R-CNN</p> <p>26:31</p>
-------------------------------	---------------------------------------	------------------------------	---	---	--

PR-095: Modularity Matters: Learning Invariant Relationships for Image Classification
조회수 458회 · 11개월 전

PR-085: In-Datacenter Performance Analysis of a Deep Neural Network
조회수 921회 · 1년 전

PR-069: Efficient Neural Architecture Search via Genetic Algorithms
조회수 2.6천회 · 1년 전

PR-054: ShuffleNet: An Extremely Efficient Convolutional Neural Network
조회수 3.8천회 · 1년 전

PR-044: MobileNet
조회수 8.7천회 · 1년 전

PR-033: PVANet: Lightweight Deep Neural Networks for Real-Time Object Detection
조회수 2.4천회 · 1년 전

Better

- Dimension Clusters
- Running & means clustering on the feature map to locate the final grid cells
- To prevent larger boxes from making more errors they use a different distance metric (L1/L2 distance)
- 4-5 is chosen as a good tradeoff between complexity and high recall

37:17

PR-023: YOLO9000: Better, Faster, Stronger
조회수 7.4천회 · 1년 전

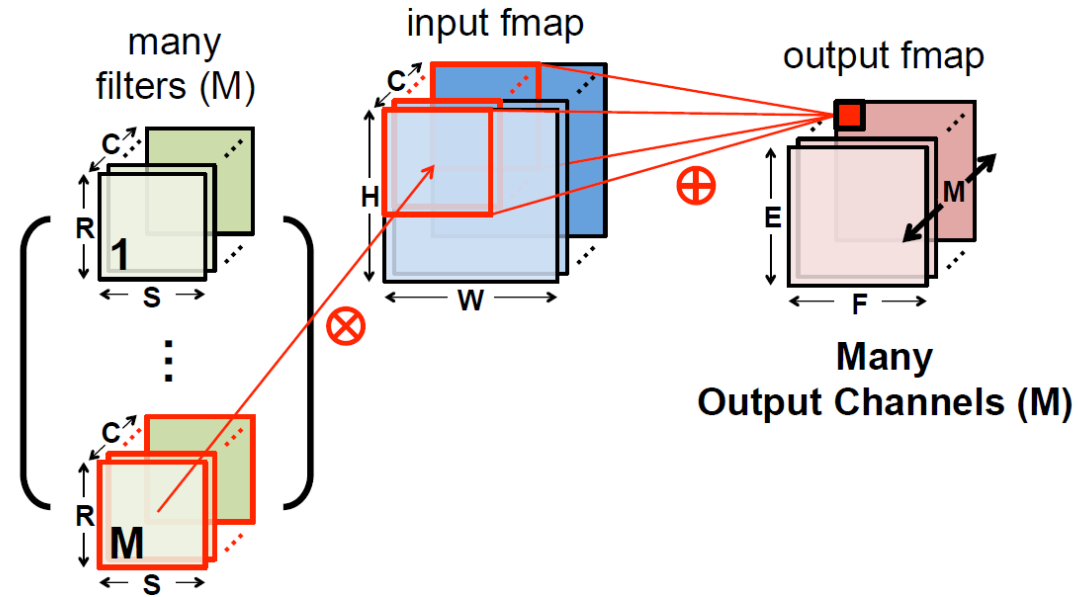
Region Proposals - Selective Search

- Bottom-up segmentation, merging regions at multiple scales

38:46

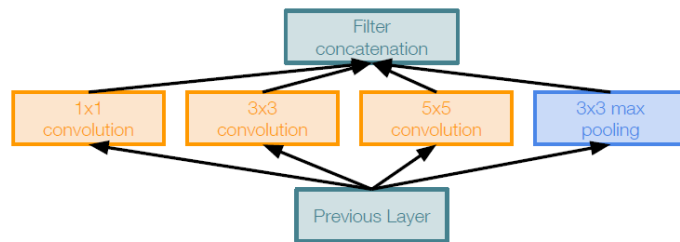
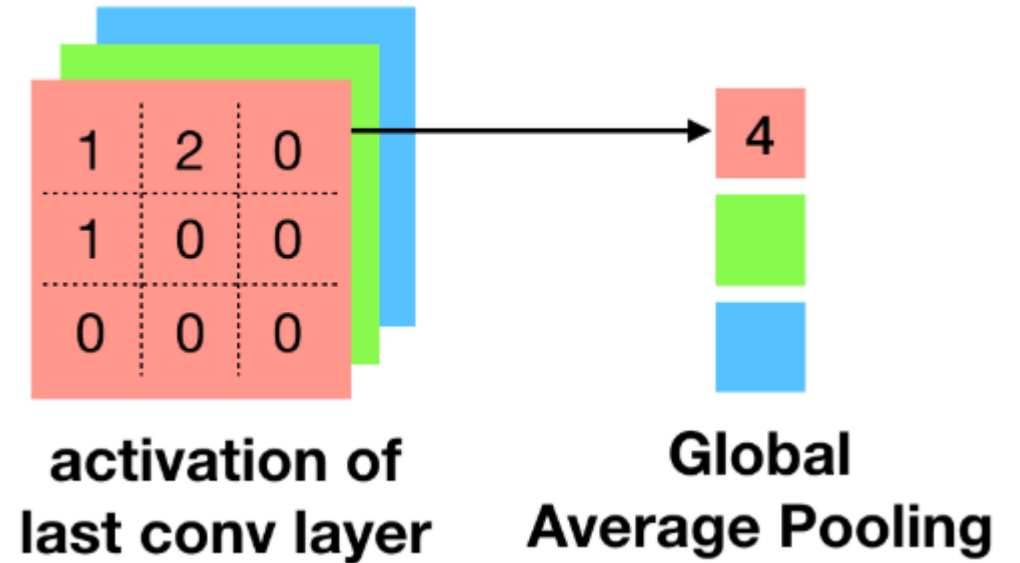
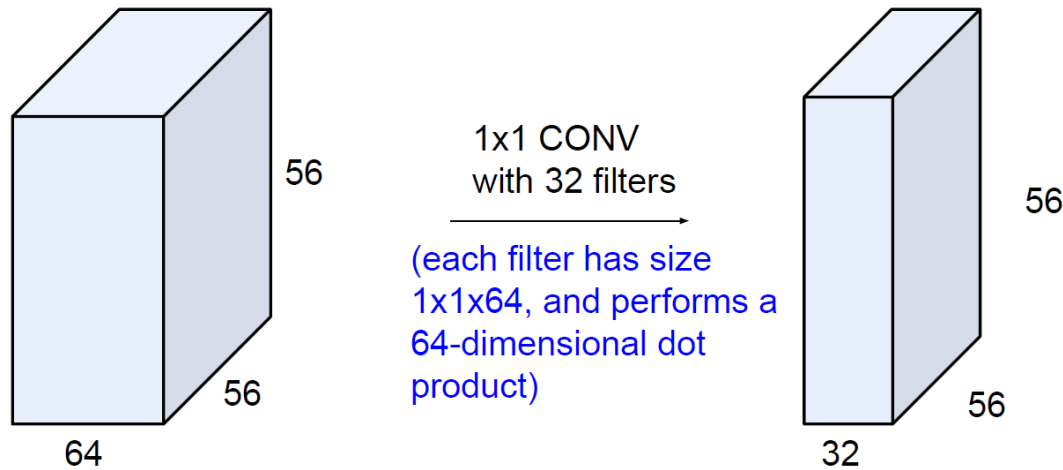
PR-012: Faster R-CNN: Towards Real-Time Object Detection
조회수 3만회 · 2년 전

Convolution

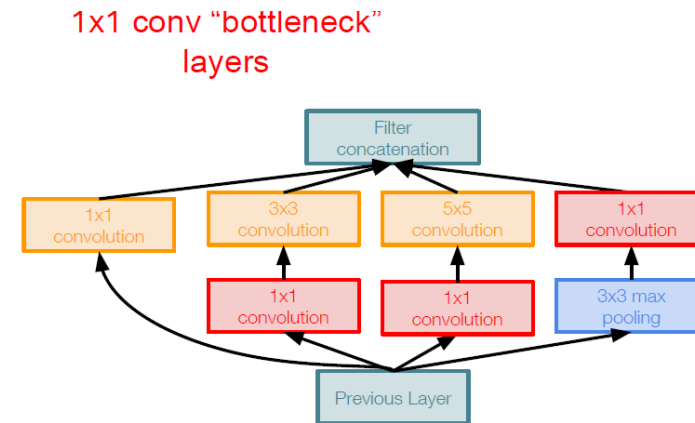


- $H=E, W=F$ 인 경우,
 - Parameter 수 : $R \times S \times C \times M$
 - 연산량 : $R \times S \times C \times H \times W \times M$
 - Memory Access Cost : $R \times S \times C \times M + (H \times W) \times (C + M)$

1. Bottleneck Layer, Global Average Pooling



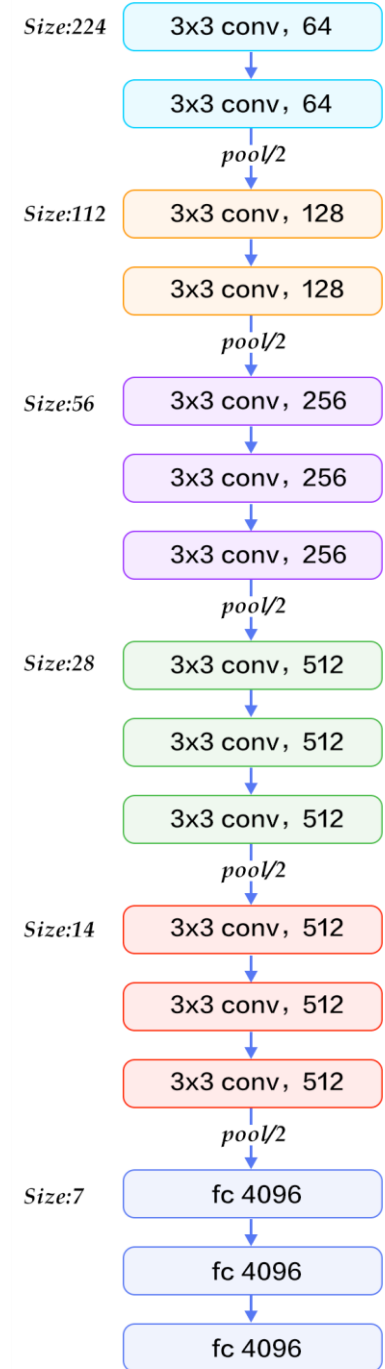
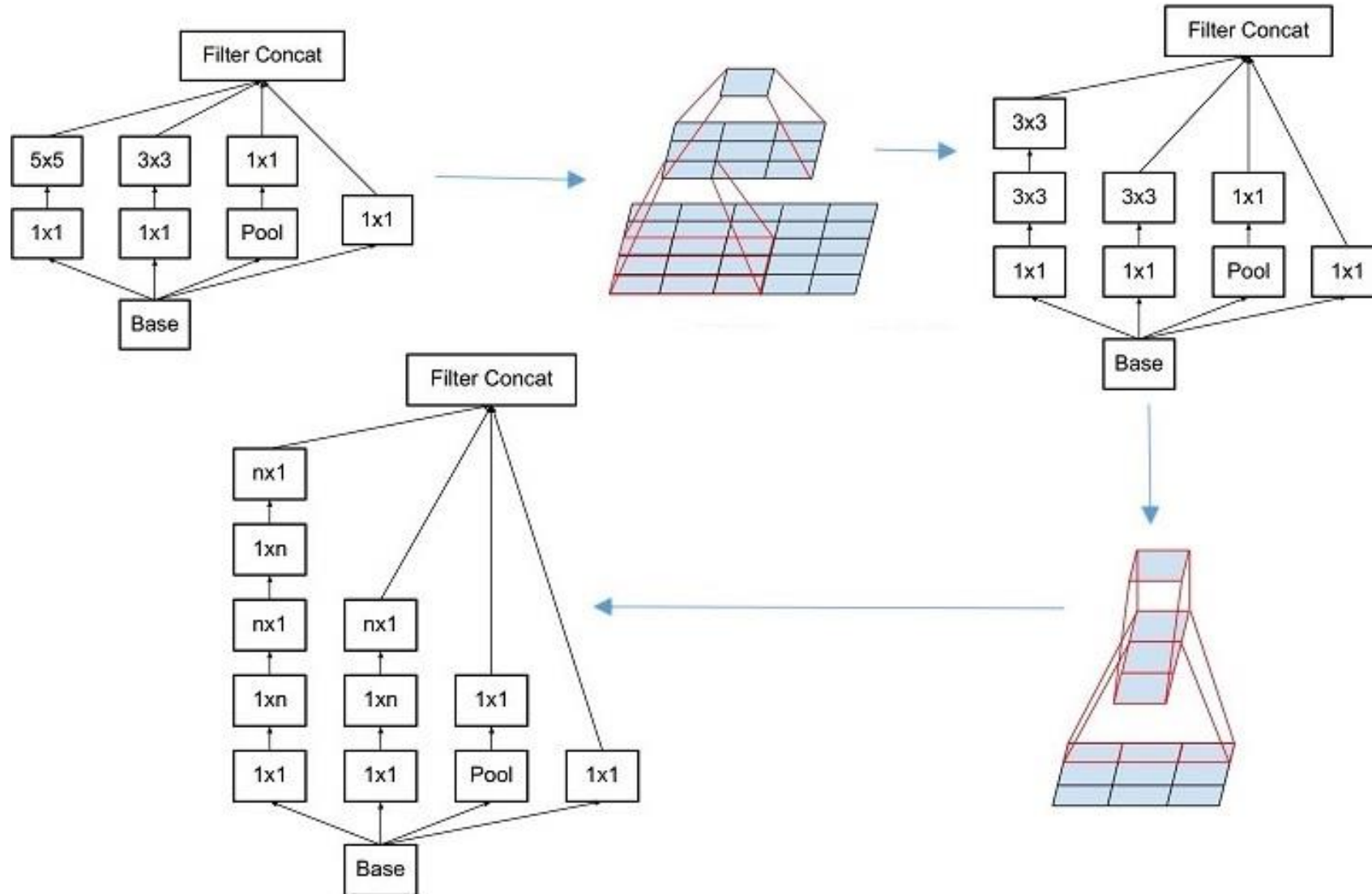
Naive Inception module



Inception module with dimension reduction

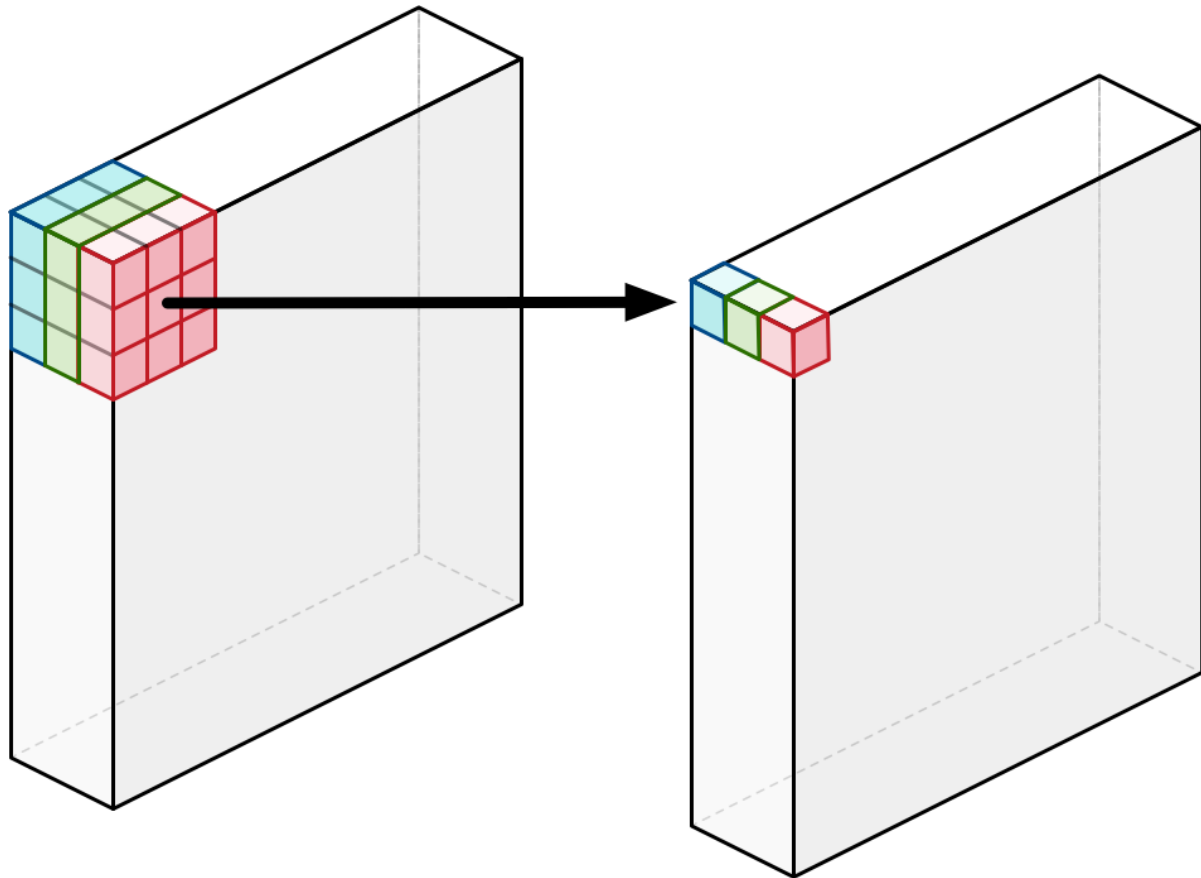
2. Filter Factorization

2. Filter Factorization

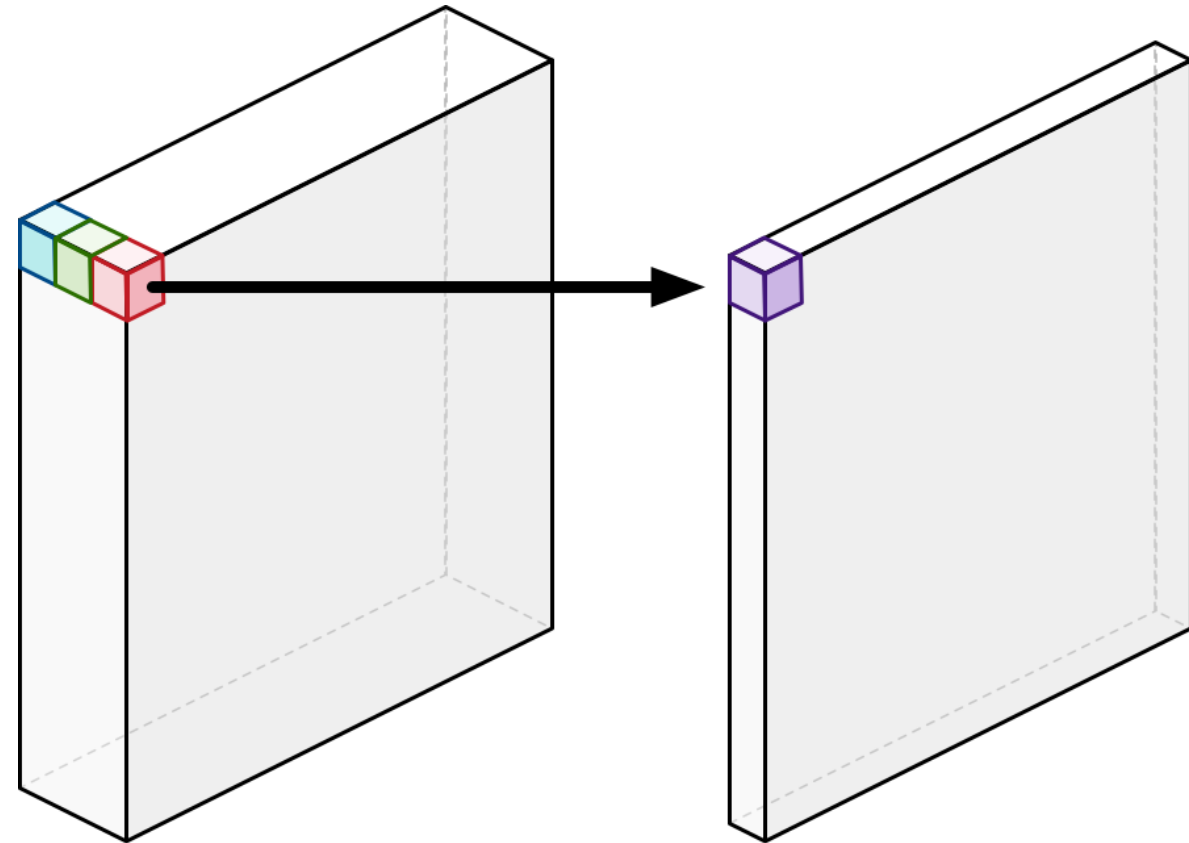


3. Depthwise Separable Convolution

3. Depthwise Separable Convolution

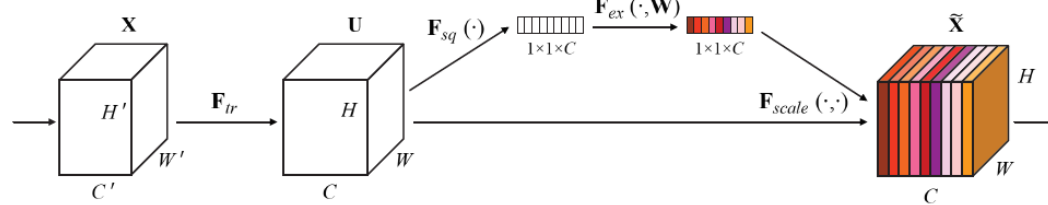
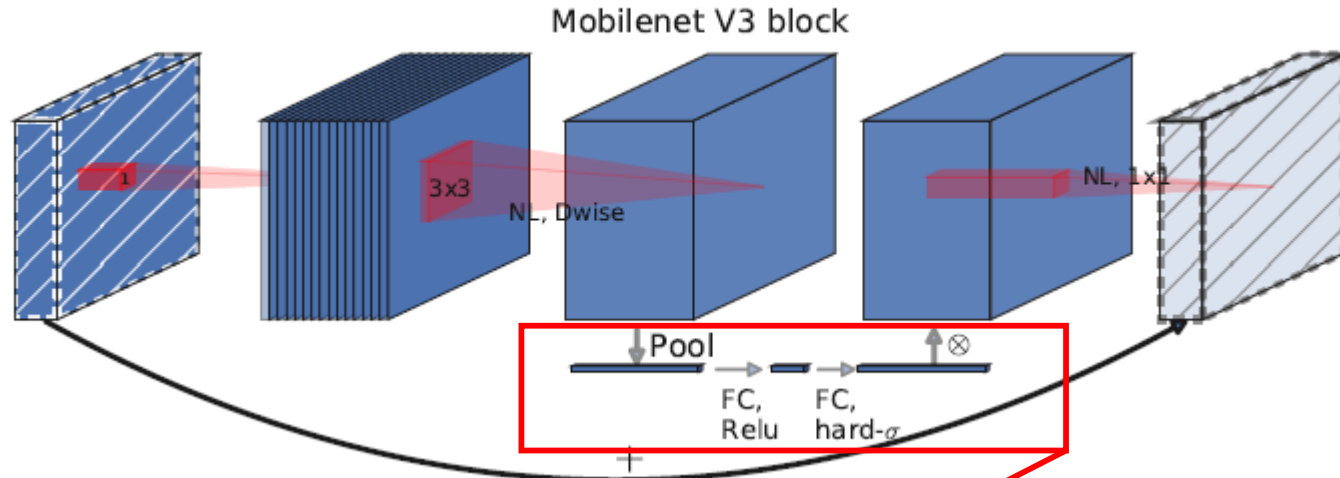


Depthwise convolution



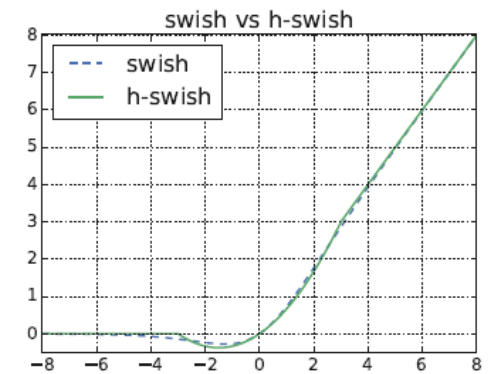
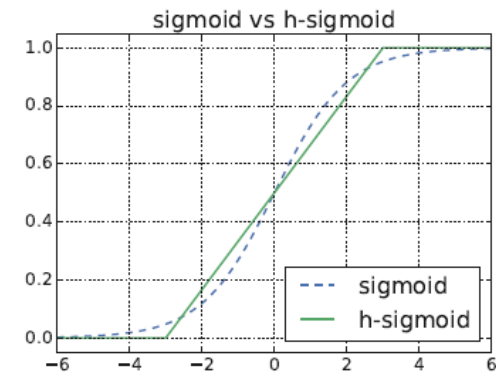
Pointwise convolution

MobileNetV3



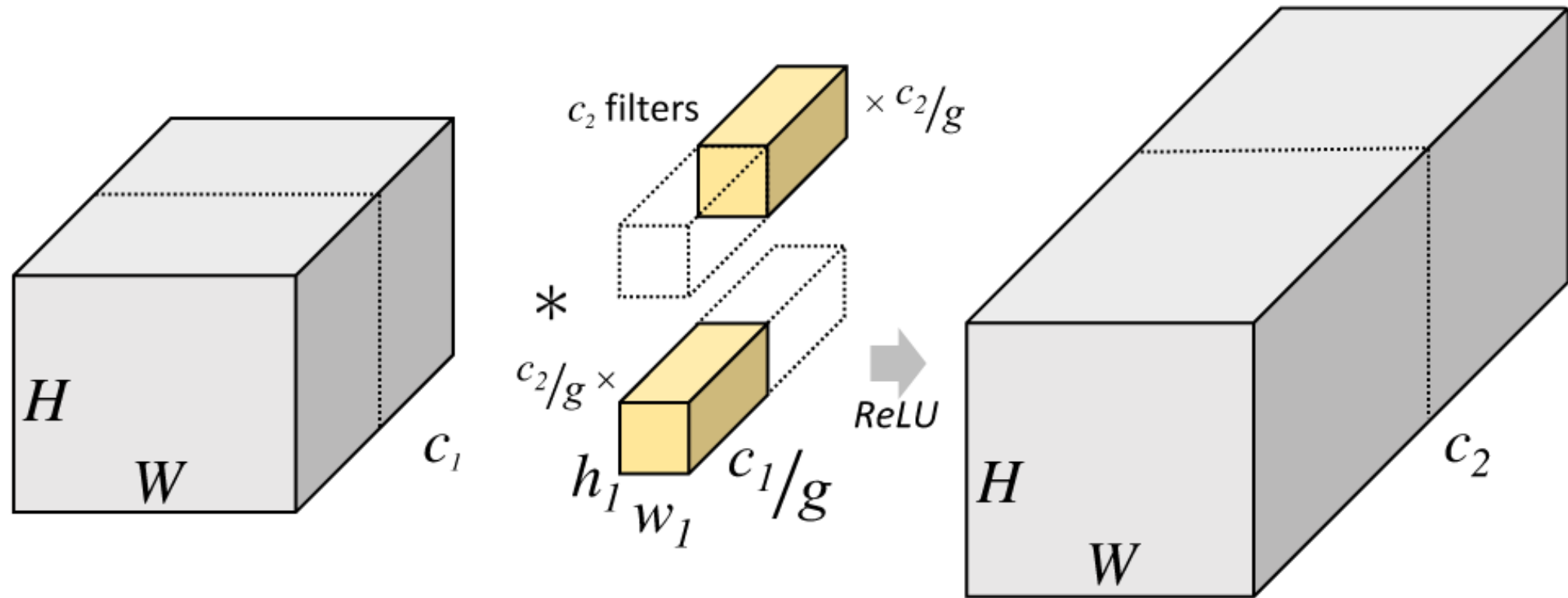
$$\text{swish } x = x \cdot \sigma(x)$$

$$\text{h-swish}[x] = x \frac{\text{ReLU6}(x + 3)}{6}$$



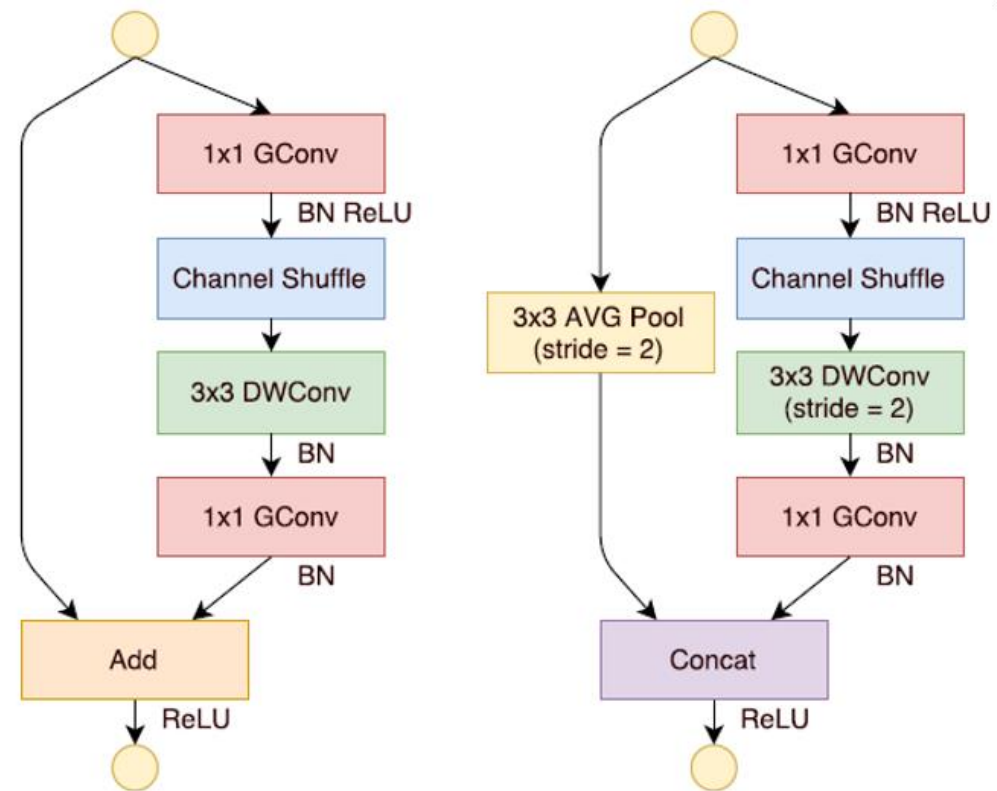
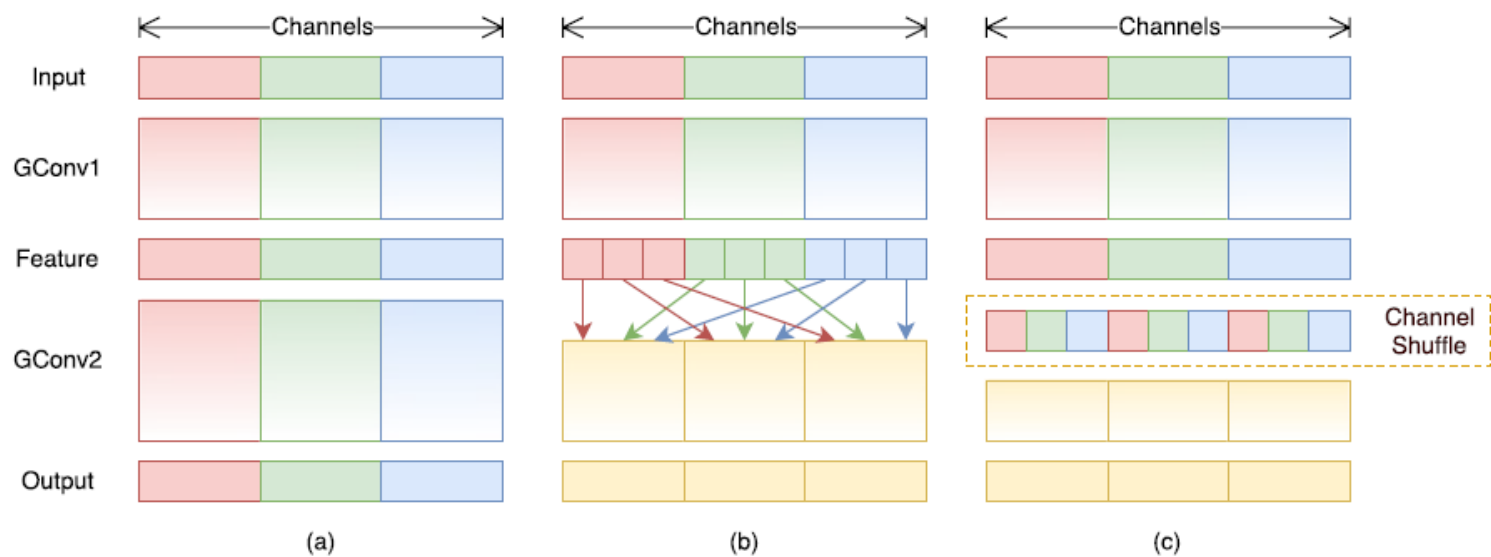
4. Group Convolution

Grouped Convolution

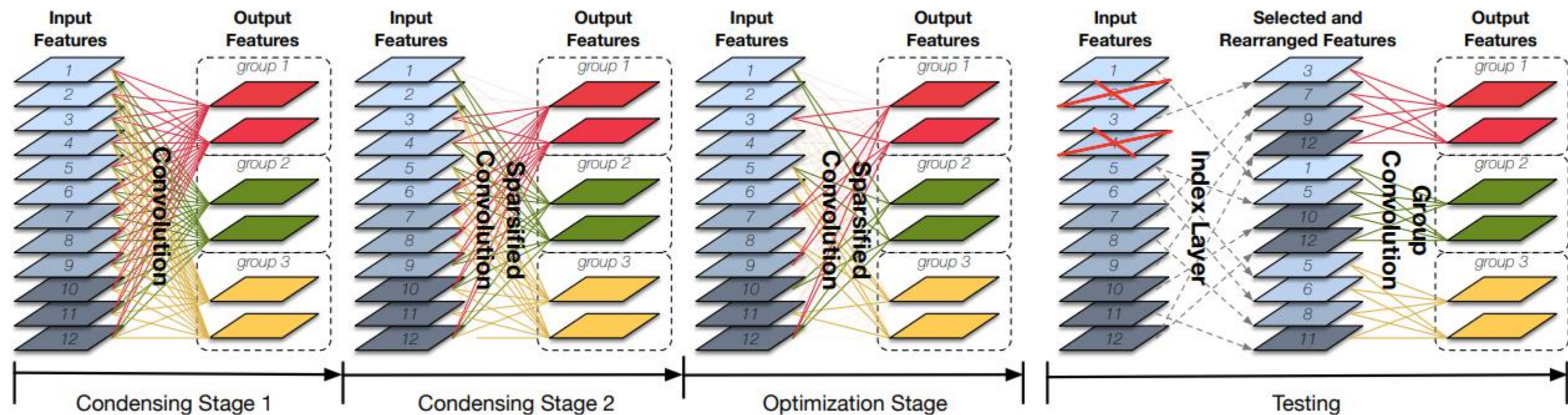


A convolutional layer with 2 filter groups. Note that each of the filters in the grouped convolutional layer is now exactly half the depth, i.e. half the parameters and half the compute as the original filter.

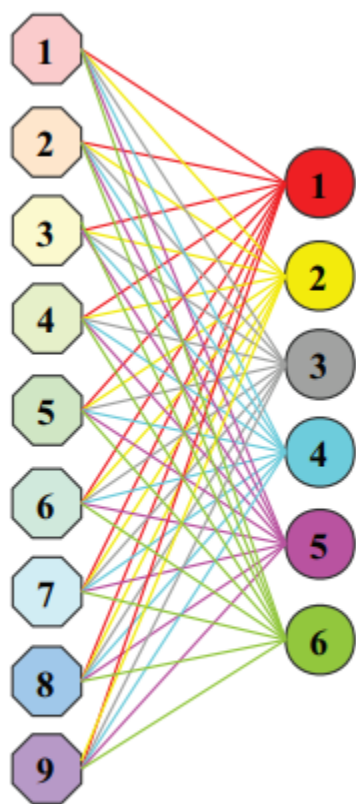
ShuffleNet



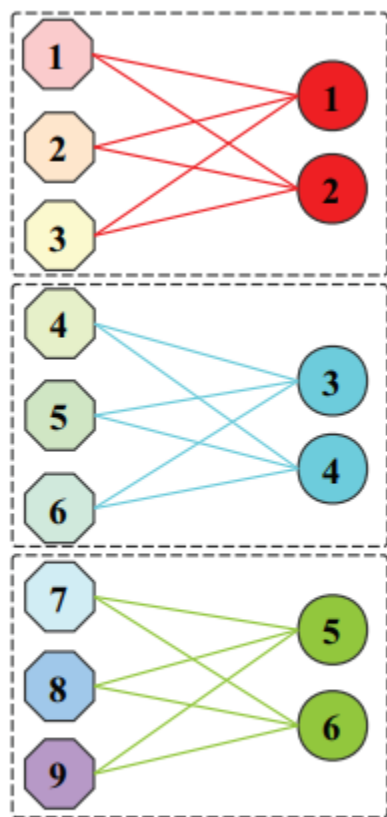
CondenseNet



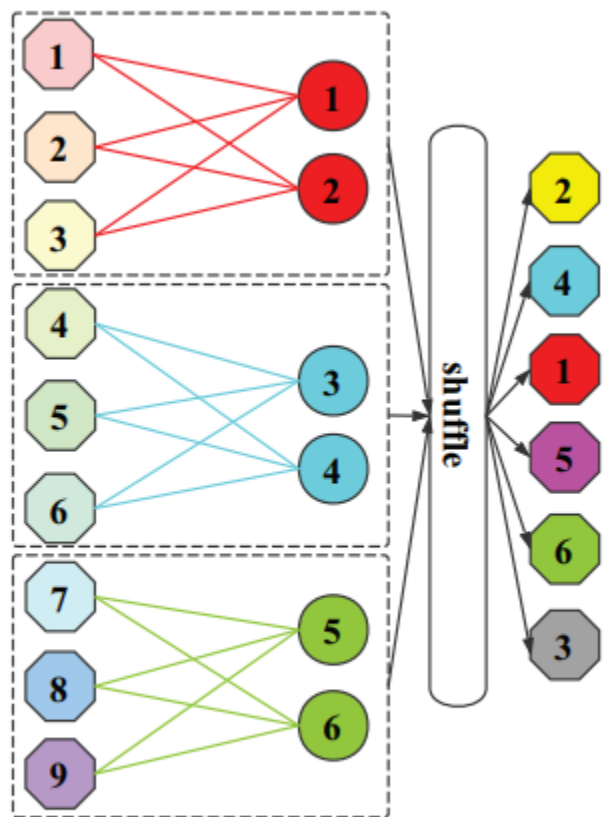
FLGC – CVPR 2019



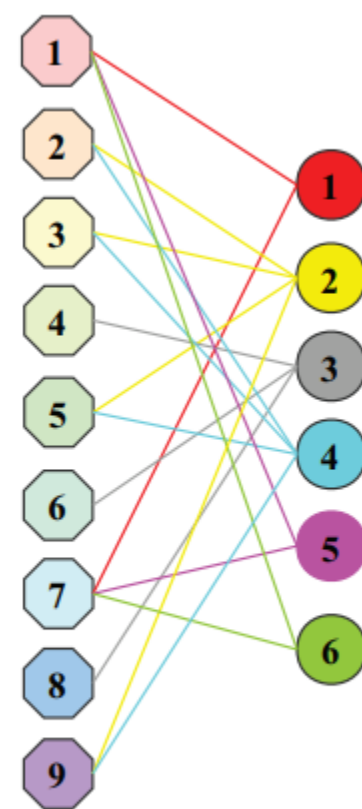
(a) Normal conv



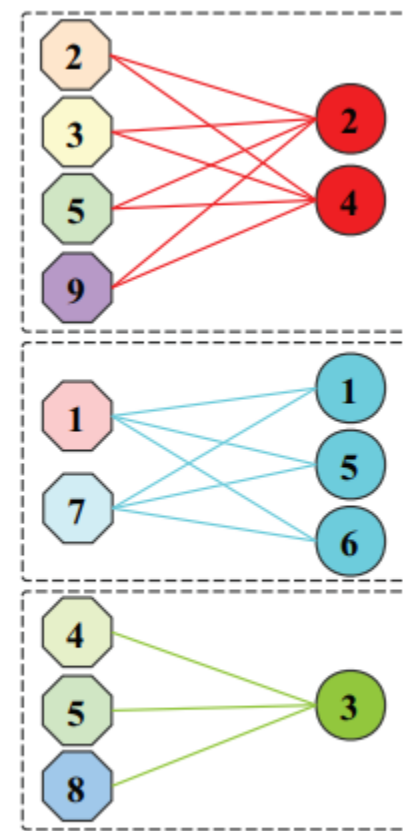
(b) Standard group conv



(c) ShuffleNet group conv unit



training

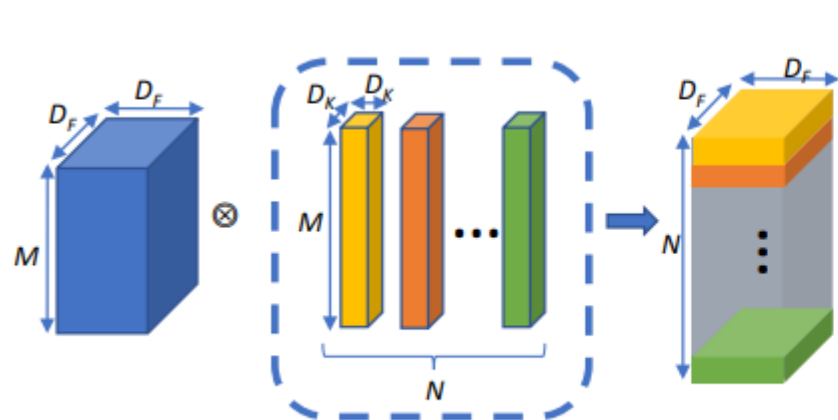


testing

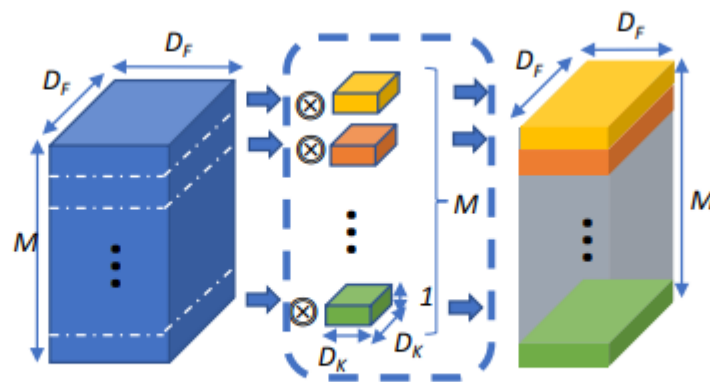
(d) FLGC

5. Shift

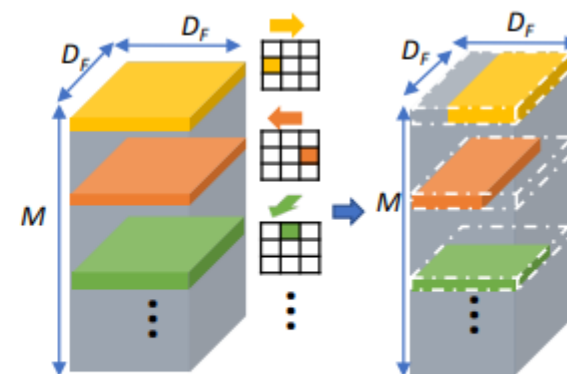
ShiftNet



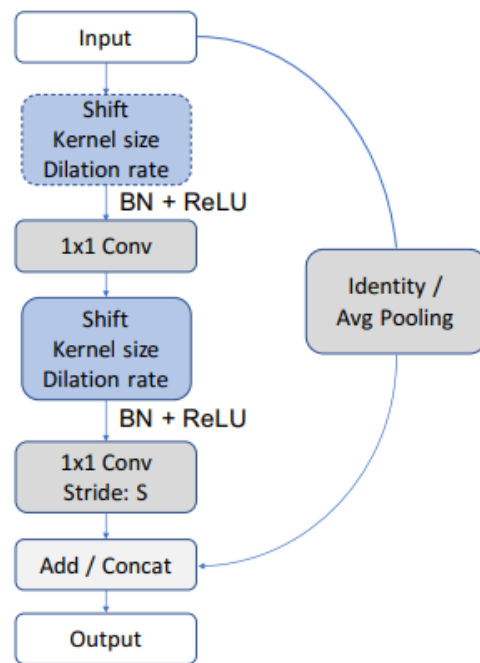
(a) Spatial Convolution



(b) Depth-wise convolution



(c) Shift



All You Need is a Few Shifts – CVPR 2019

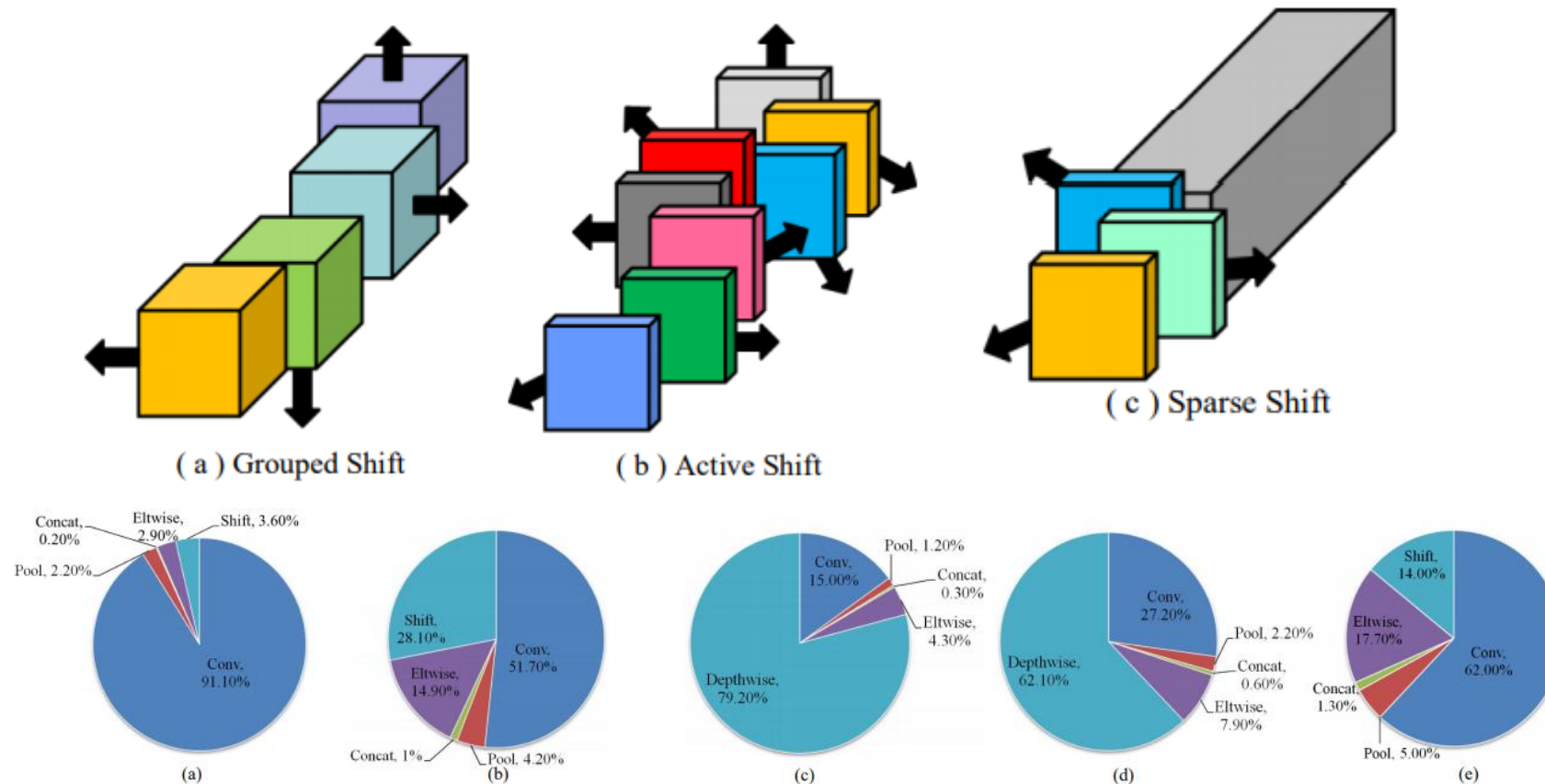
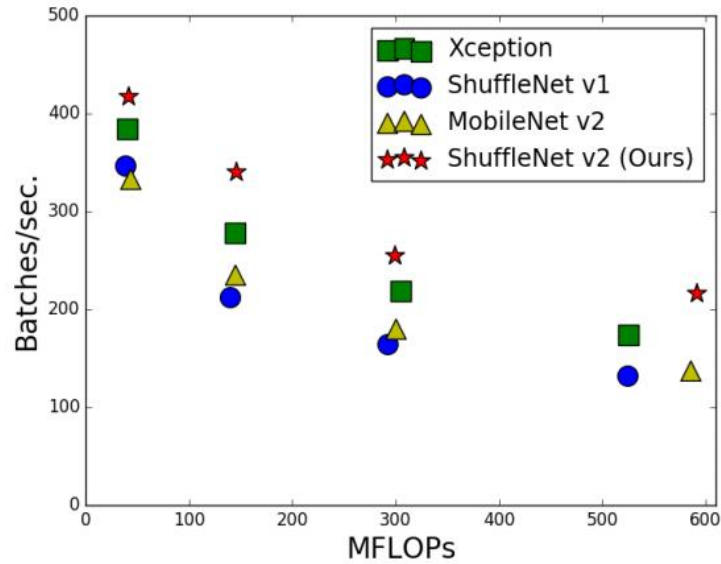
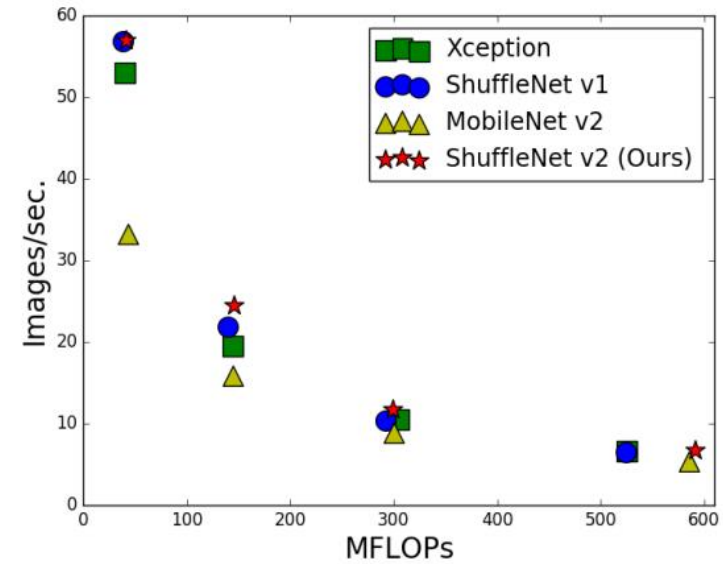


Figure 2. The practical runtime analysis. For clear comparison, both batch-normalization and ReLU layers are neglected since they can be merged into convolutional layer for inference. Also data feeding and preprocessing time are not considered here. Results are achieved under Caffe with mini-batch 32. They are averaged from 100 runs. (a) ShiftNet-A [37] on CPU (Intel Xeon E5-2650, atlas). (b) ShiftNet-A on GPU (TITAN X Pascal, CUDA8 and cuDNN5). (c) Shift layers in ShiftNet-A are replaced by depthwise separable convolution layers. (d) Depthwise separable convolution layers with kernel size 5 are replaced by the ones with kernel size 3. (e) ShiftNet-A with 80% shift sparsity on GPU (Shift sparsity denotes the ratio of unshifted feature maps).

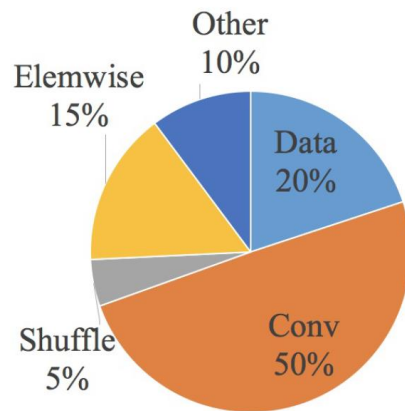
6. Use Direct Metric



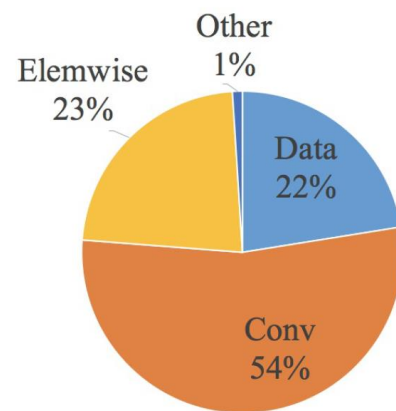
(c) GPU



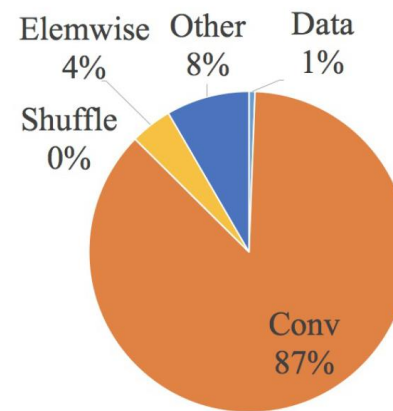
(d) ARM



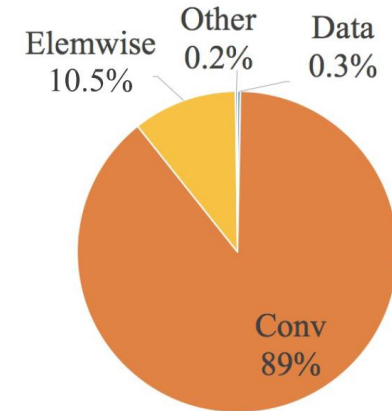
ShuffleNet V1 on GPU



MobileNet V2 on GPU



ShuffleNet V1 on ARM



MobileNet V2 on ARM

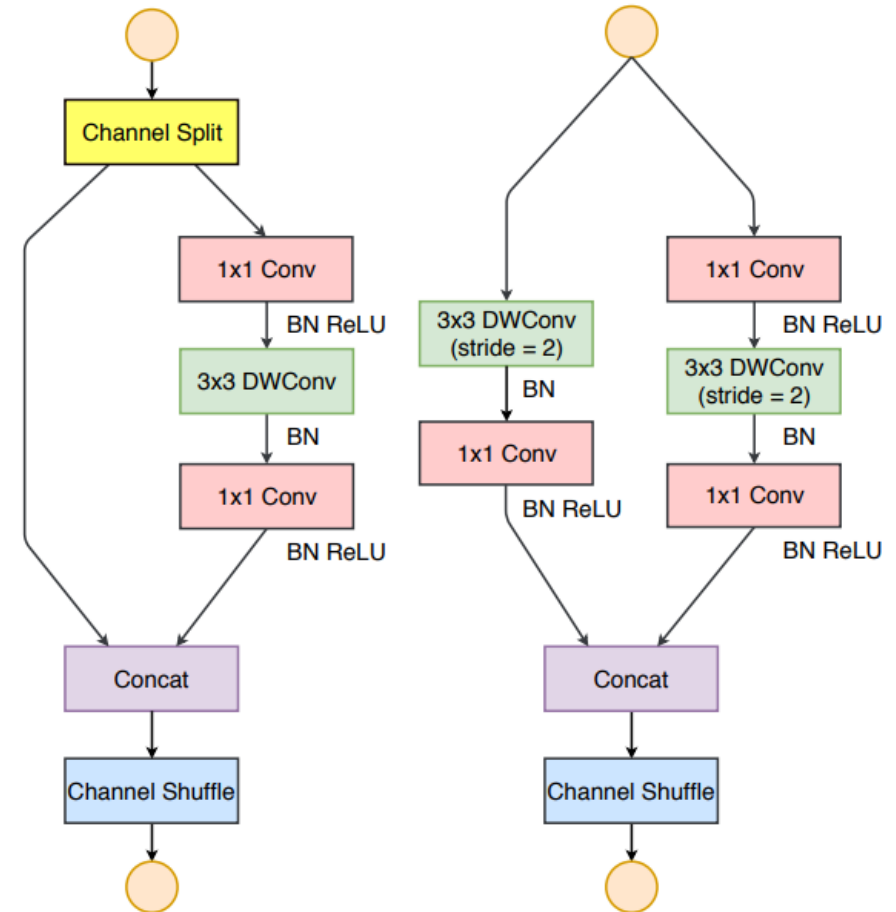
ShuffleNet V2 – Guide 1

- **Equal channel width minimizes memory access cost (MAC).**
- Let h and w be the spatial size of the feature map, the FLOPs of the 1×1 convolution is $B = hwc_1c_2$.
- The memory access cost (MAC), or the number of memory access operations, is $MAC = hw(c_1 + c_2) + c_1c_2$.
- MAC has a lower bound given by FLOPs. **It reaches the lower bound when the numbers of input and output channels are equal.**

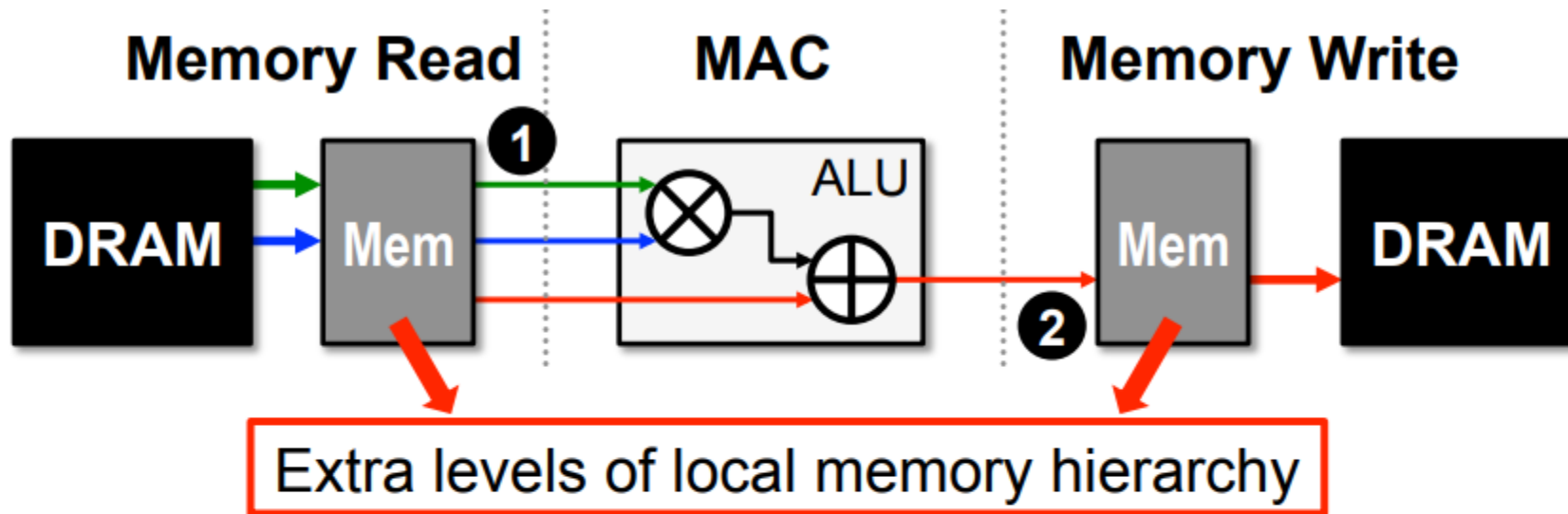
$$MAC \geq 2\sqrt{hwB} + \frac{B}{hw}$$

ShuffleNet V2

1. Use "balanced" convolutions (equal channel width).
2. Be aware of the cost of using group convolution.
3. Reduce the degree of fragmentation.
4. Reduce element-wise operations.



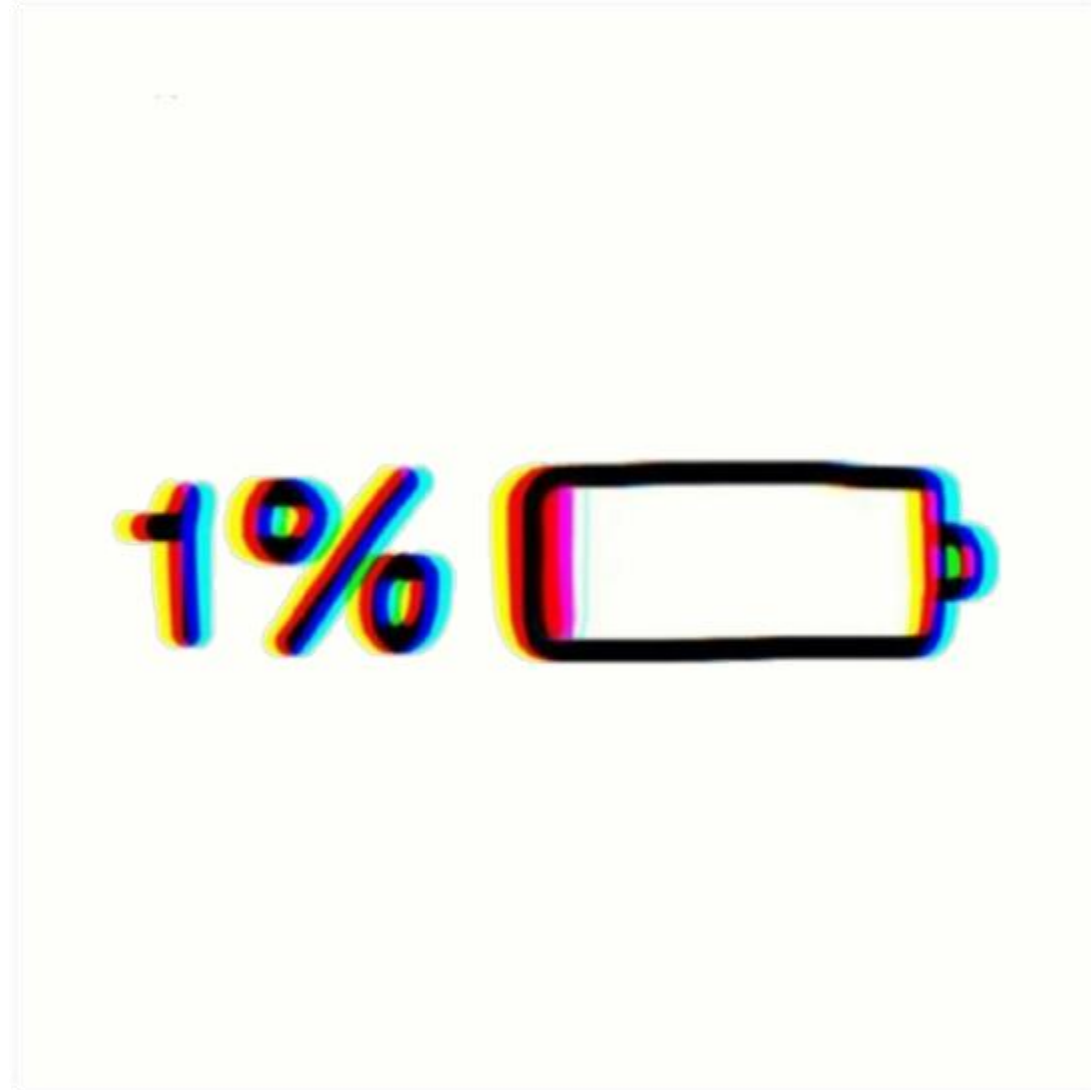
Most DNN Accelerators



Opportunities: ❶ **data reuse** ❷ **local accumulation**

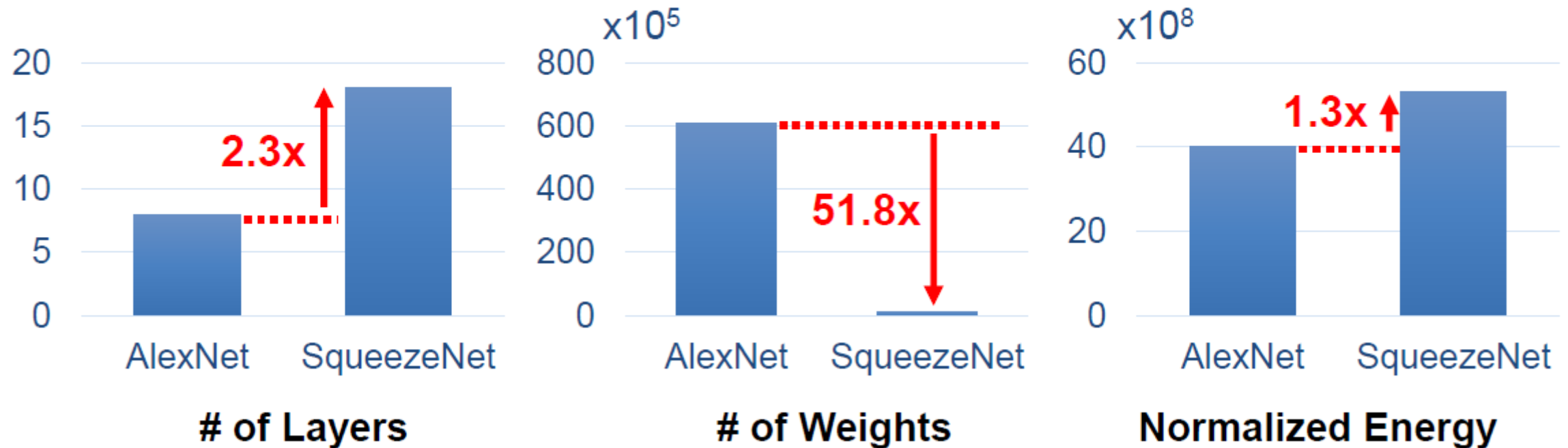
- ❶ Can reduce DRAM reads of **filter/fmap** by up to **500×**
- ❷ **Partial sum** accumulation does **NOT** have to access DRAM
 - Example: DRAM access in AlexNet can be reduced from **2896M** to **61M** (best case)

Energy is very Important



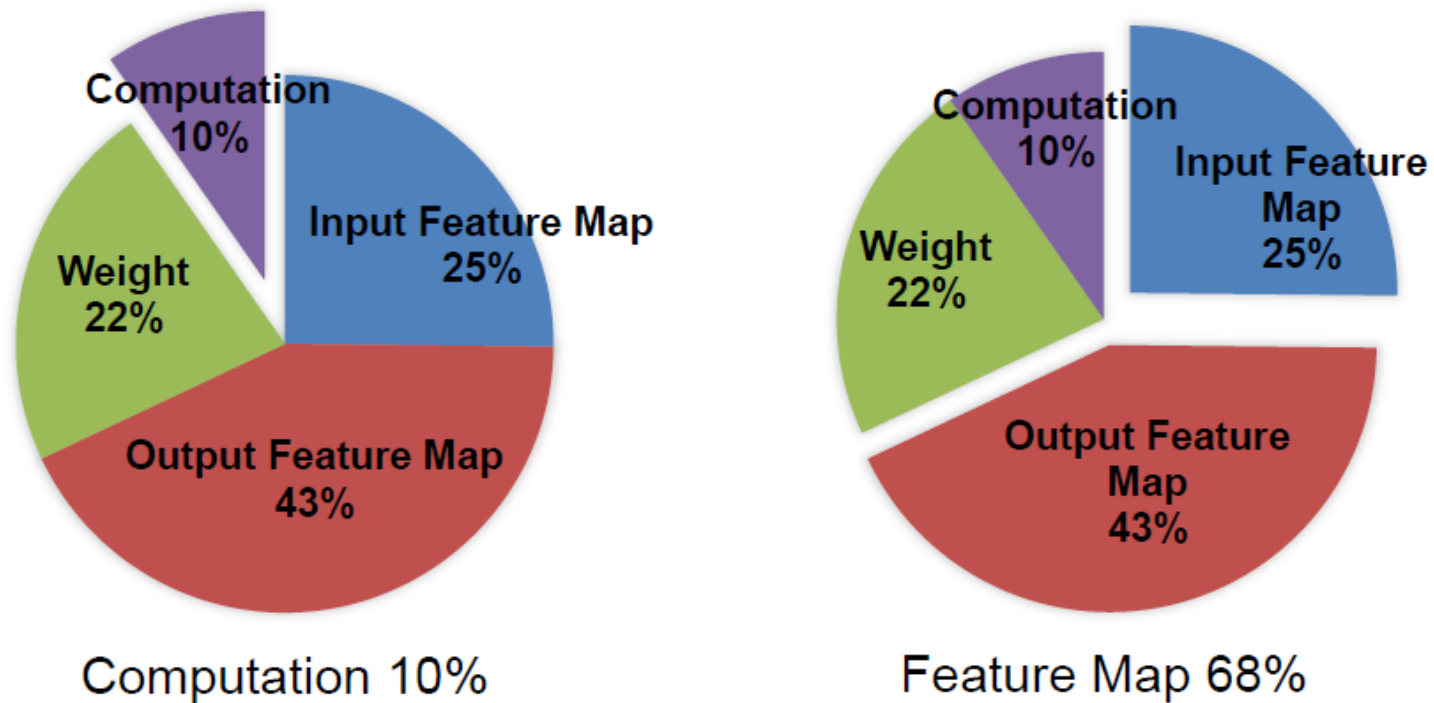
Energy?!

Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights



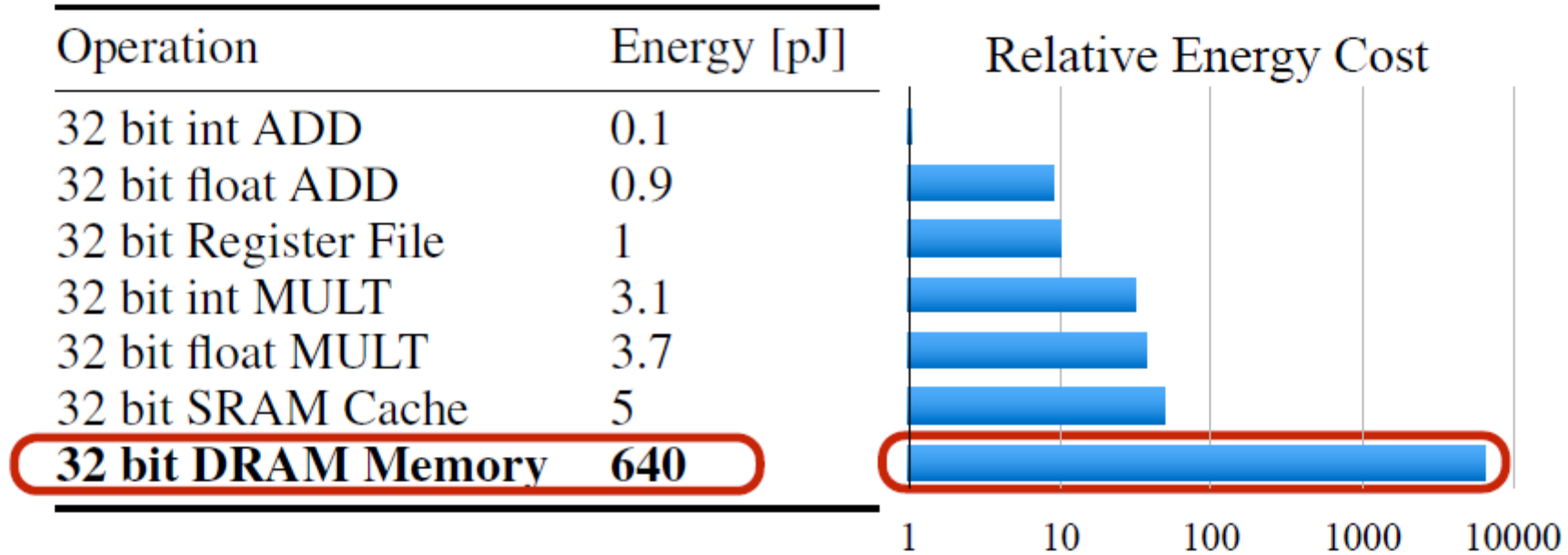
Key Insights

- Data movement is more expensive than computation
- Feature maps need to be taken into account

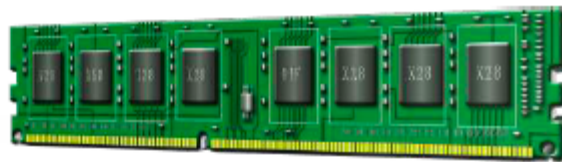


GoogLeNet Energy Breakdown

Where is the Energy Consumed?



1



[This image](#) is in the public domain

= 1000 × +

EfficientNet

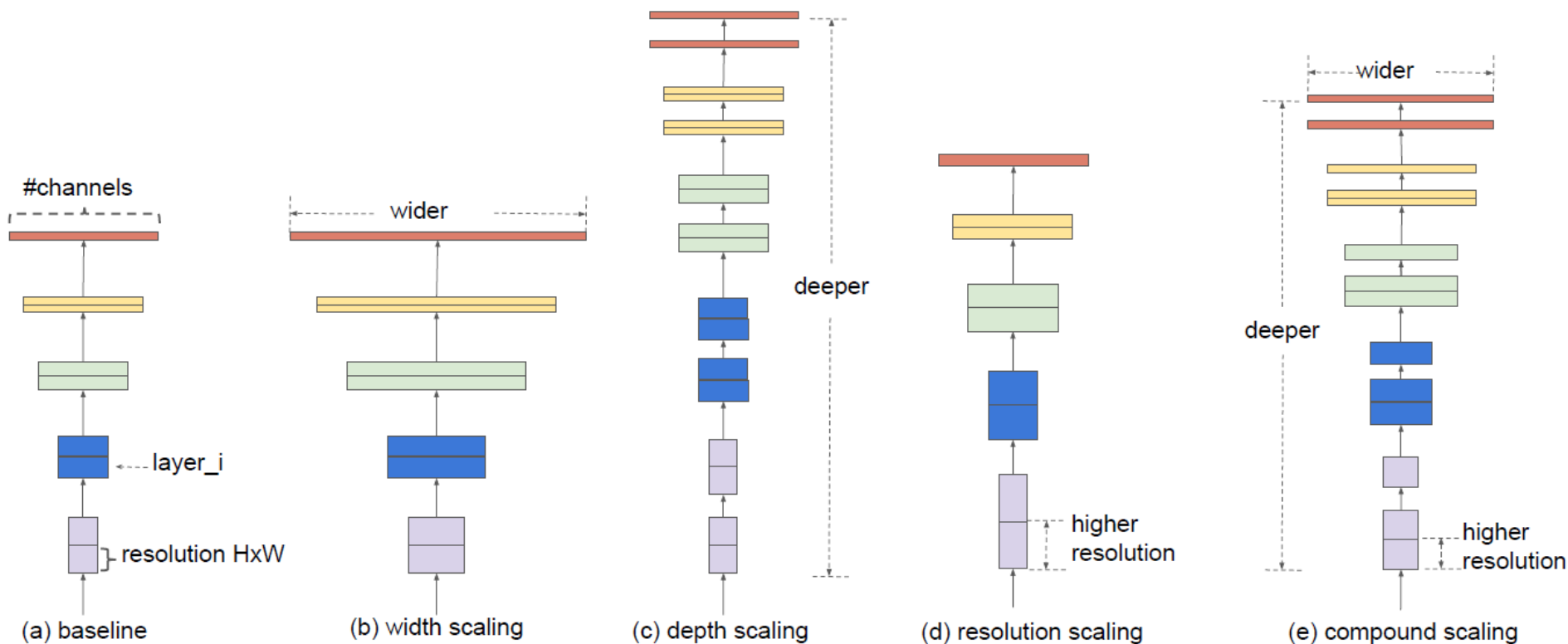
depth: $d = \alpha^\phi$

width: $w = \beta^\phi$

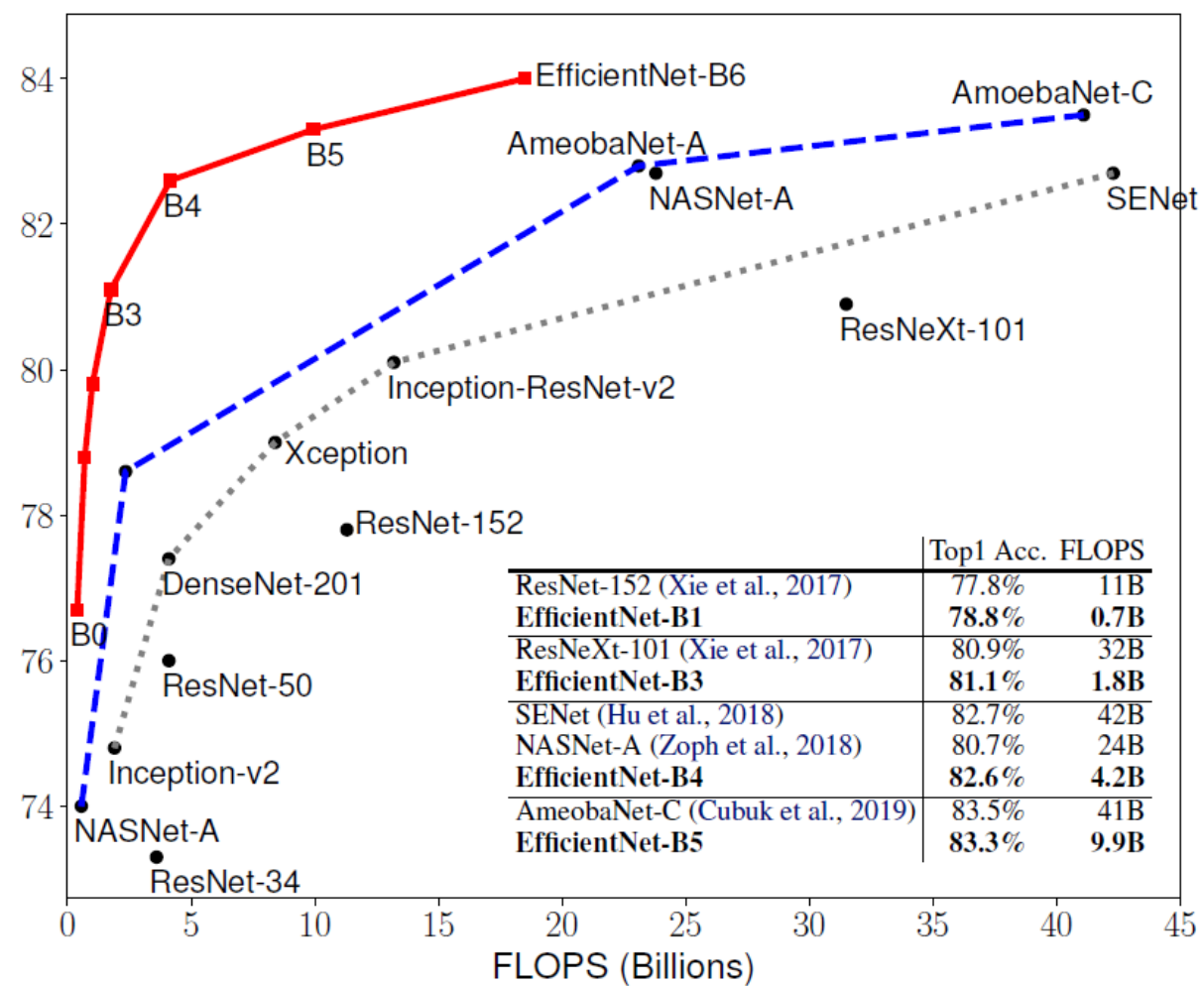
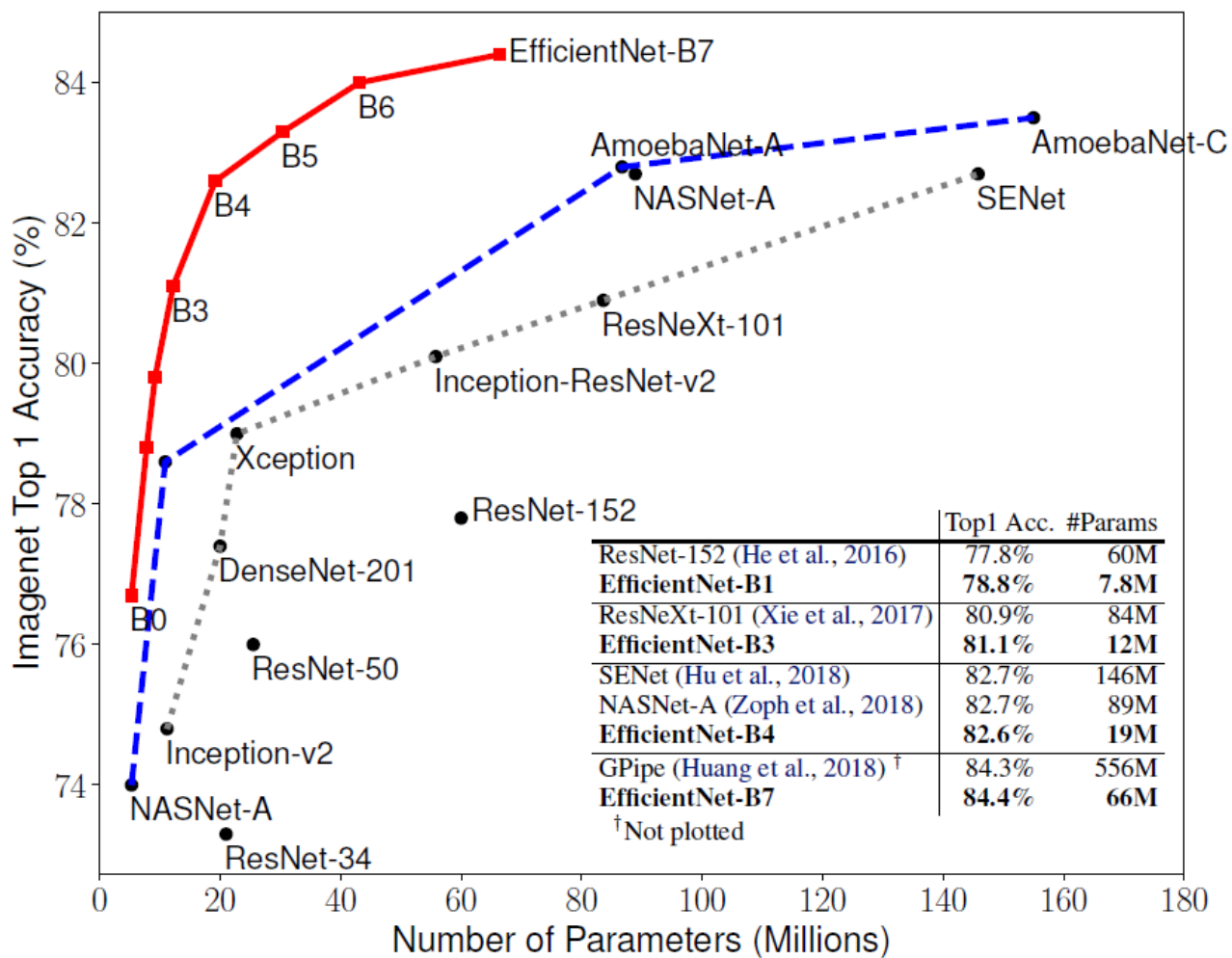
resolution: $r = \gamma^\phi$

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

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



EfficientNet



Thank you