

fmap in memory

fmap not in memory

 $i^{\text{th}}$  mobile inverted bottleneck block

learnable params

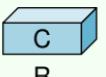
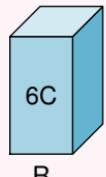
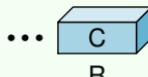
fixed params

weight

bias

 $1 \times 1$  Conv

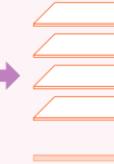
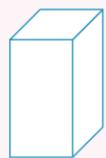
Depth-wise Conv

 $1 \times 1$  Conv

a) Fine-tune the full network (Conventional)

 $1 \times 1$  Conv

Depth-wise Conv

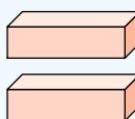
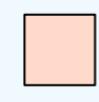
 $1 \times 1$  Conv

b) Fine-tune bias only

Downsample



Group Conv

 $1 \times 1$  Conv

Upsample



c) Lite residual learning