Legend: fmap in memory; fmap not in memory; $i^{th}$ mobile inverted bottleneck block; learnable params; fixed params; weight; bias
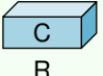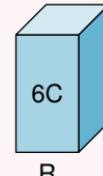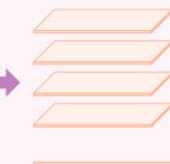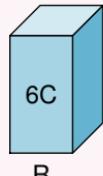
a) Fine-tune the full network (Conventional)

$1 \times 1$ Conv — Depth-wise Conv — $1 \times 1$ Conv

b) Fine-tune bias only

$1 \times 1$ Conv — Depth-wise Conv — $1 \times 1$ Conv

c) Lite residual learning

Downsample — Group Conv — $1 \times 1$ Conv — Upsample