# Lite Pose: Efficient Architecture Design 2D for Human Pose Estimation

Yihan Wang[1], Muyang Li[2], Han Cai[3], Wei-Ming Chen[3], Song Han[3]

[1]Tsinghua University   [2]CMU   [3]MIT

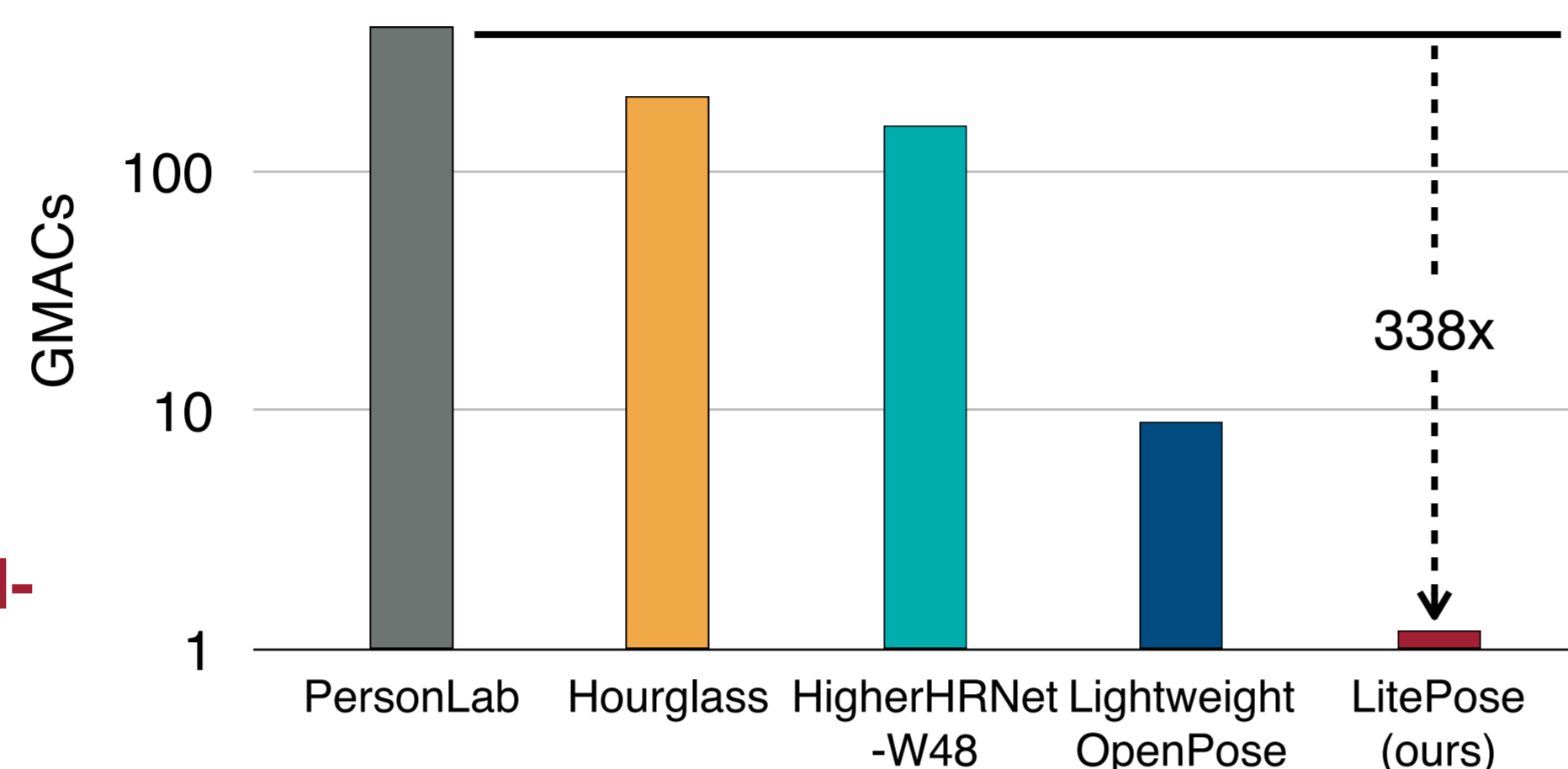Code: https://github.com/mit-han-lab/litepose

## Real-Time Multi-Person Pose Estimation on Edge
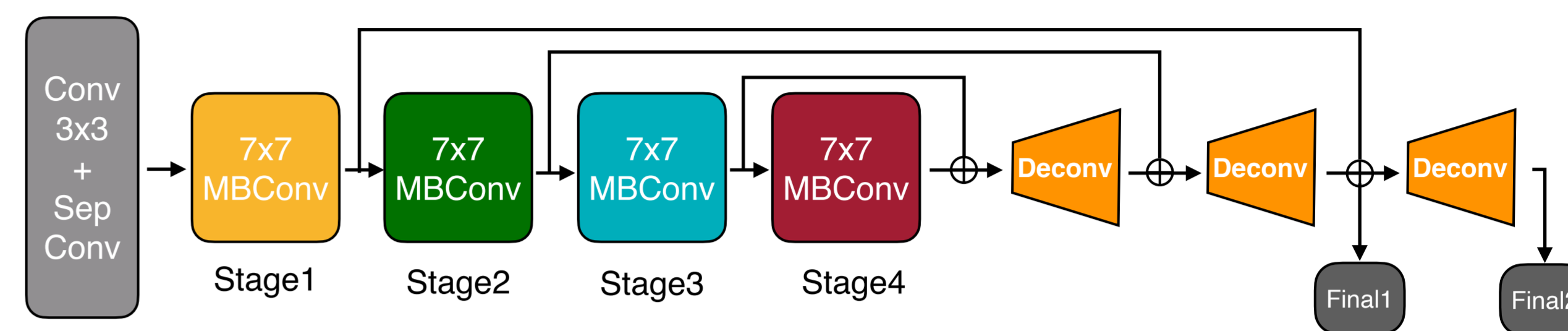


Multi-Person Pose Estimation → Edge Devices

Many human-centered vision applications rely on **real-time multi-person** pose estimation on **edge** devices, requiring **low-computation** pose estimation models.

However, current pose estimation models are too **heavy** for edge devices. We introduce **LitePose** to close the gap.
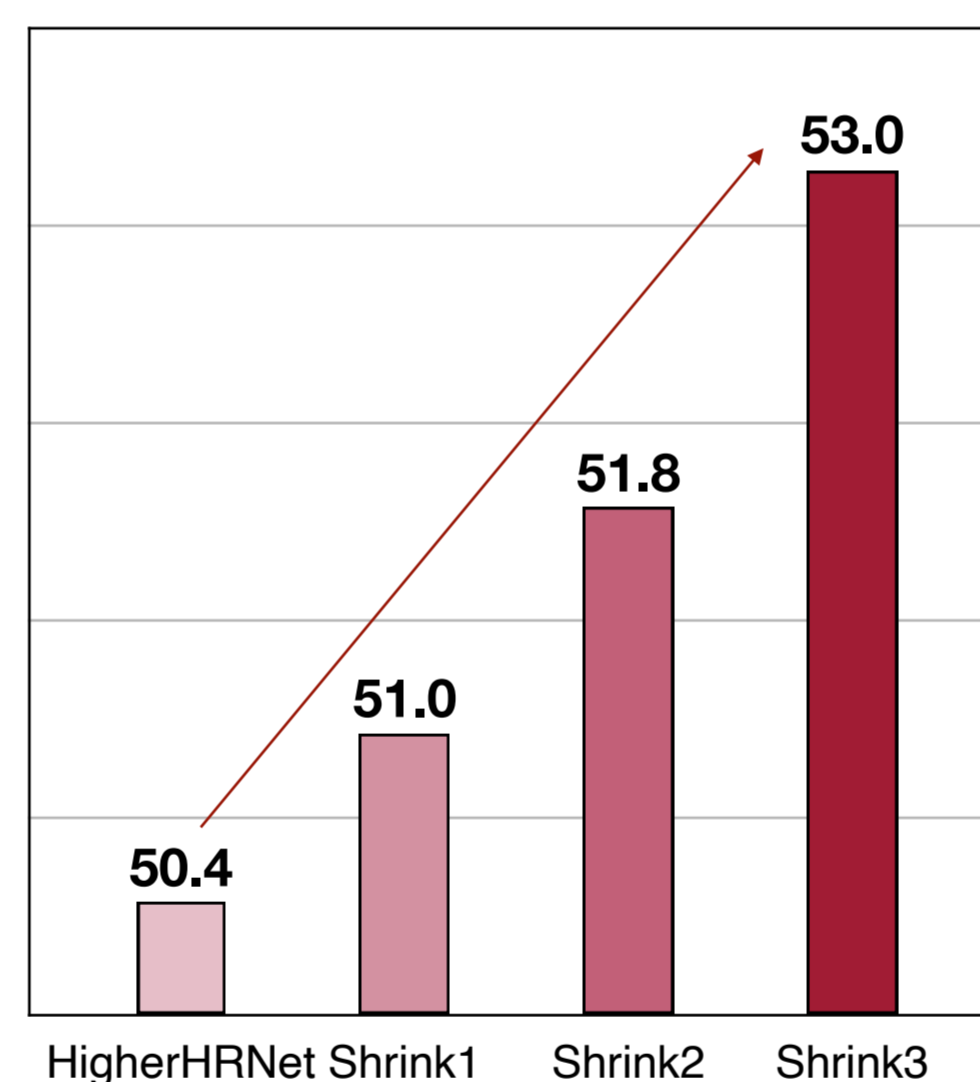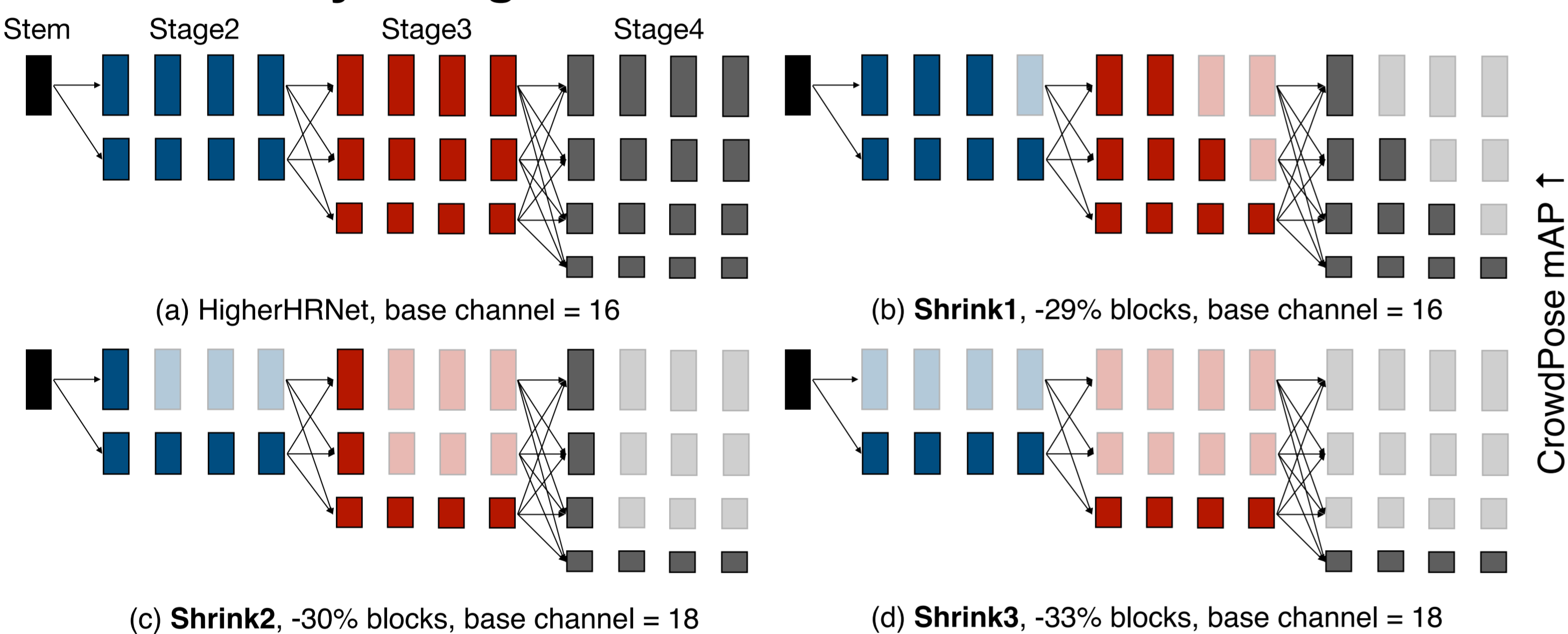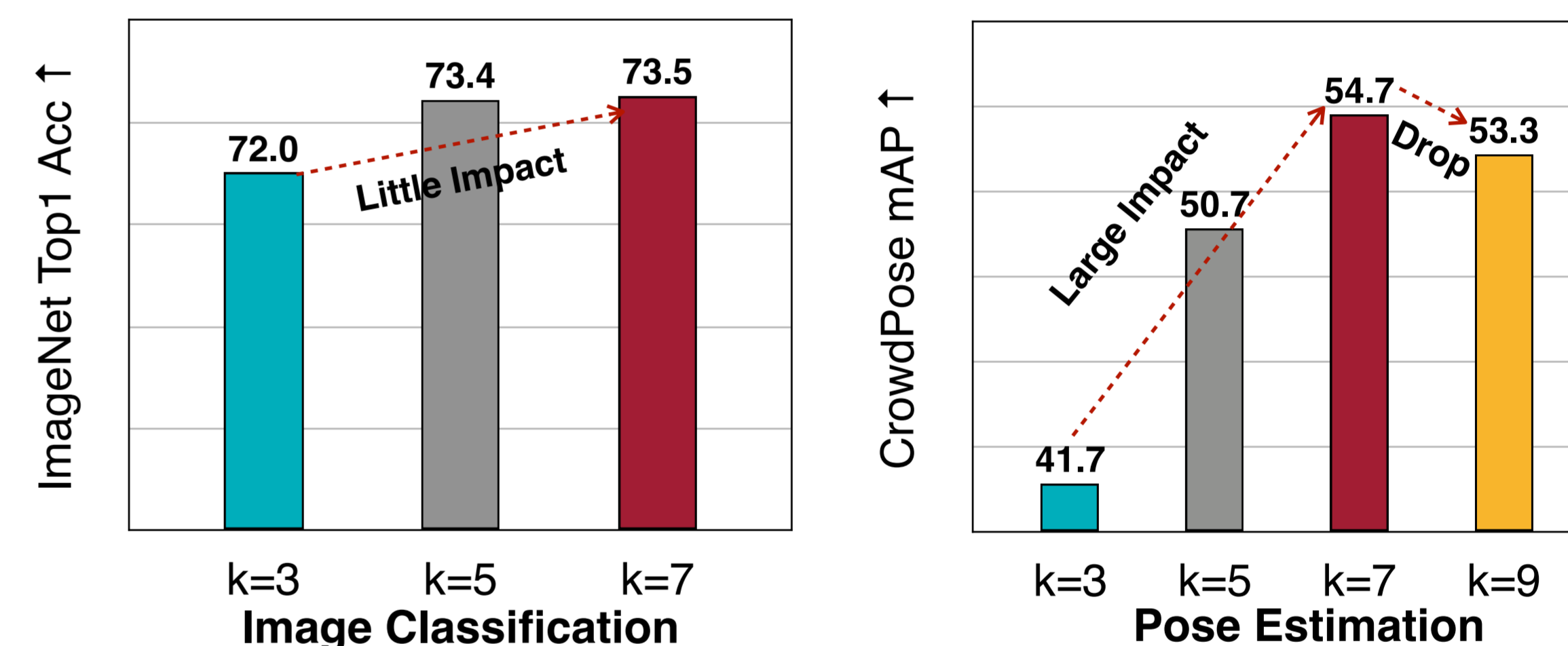


GMACs, PersonLab, Hourglass, HigherHRNet-W48, Lightweight OpenPose, LitePose (ours), 338x

## Overview of LitePose



Conv 3x3 + Sep Conv → 7x7 MBConv (Stage1) → 7x7 MBConv (Stage2) → 7x7 MBConv (Stage3) → 7x7 MBConv (Stage4) → Deconv → Deconv → Deconv, Final1, Final2

### Key insights:

1. **Single-branch architecture is efficient**
2. **Large kernel convolution is efficient.**
3. **Light-weight fusion deconv head.**

## Redundancy in High-Resolution Branches



(a) HigherHRNet, base channel = 16

(b) **Shrink1**, -29% blocks, base channel = 16

(c) **Shrink2**, -30% blocks, base channel = 18

(d) **Shrink3**, -33% blocks, base channel = 18

CrowdPose mAP ↑: HigherHRNet 50.4, Shrink1 51.0, Shrink2 51.8, Shrink3 53.0

We gradually remove blocks in high-resolution branches starting from HigherHRNet. Removed blocks are shown in transparent. The performance improves as we shrink the high-resolution branches.

## Large Kernel is Efficient



ImageNet Top1 Acc ↑: k=3 72.0, k=5 73.4, k=7 73.5, Little Impact — Image Classification

CrowdPose mAP ↑: k=3 41.7, k=5 50.7, k=7 54.7, k=9 53.3, Large Impact, Drop — Pose Estimation
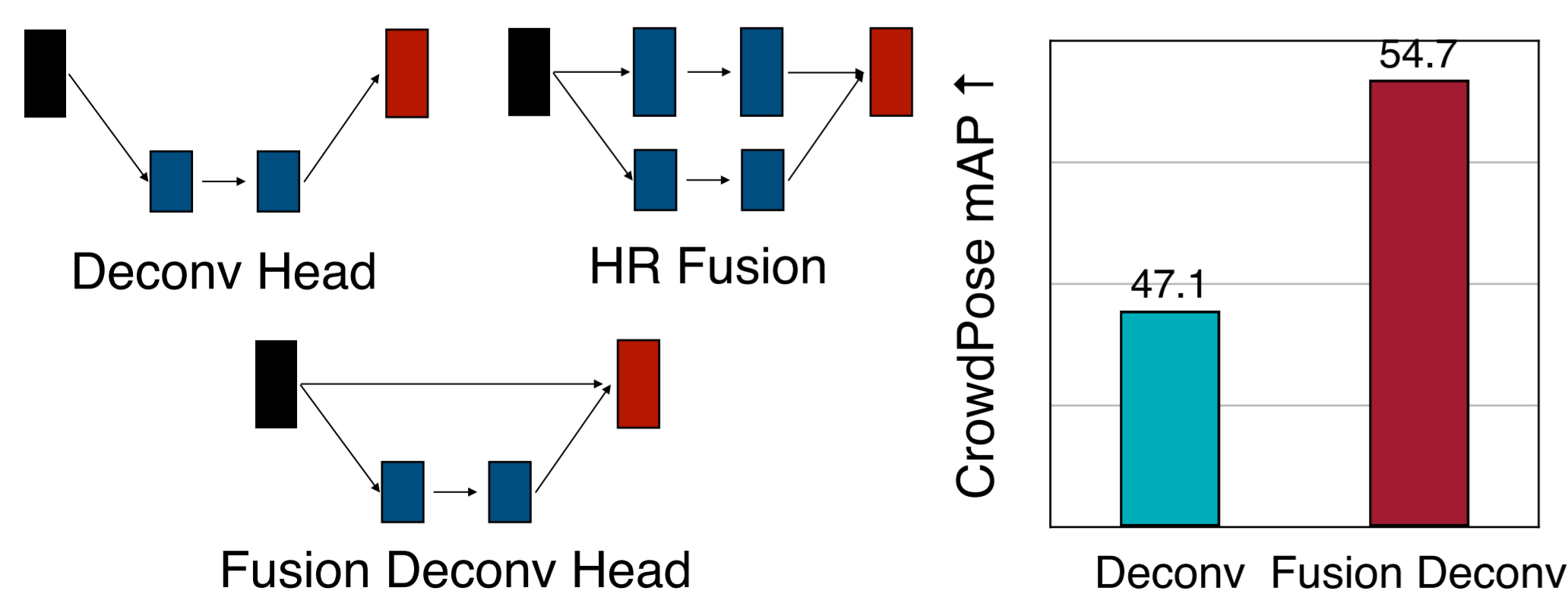
Unlike image classification, large kernel depthwise convolution plays a critical role in pose estimation. Increasing the kernel size from 3 to 7 improves the mAP by 13% mAP on the CrowdPose dataset with little overhead.

## Light-weight Fusion Deconv Head



Deconv Head   HR Fusion   Fusion Deconv Head

CrowdPose mAP ↑: Deconv 47.1, Fusion Deconv 54.7

We employ the lightweight fusion deconv head to enable multi-resolution feature fusion without heavy high-resolution branches.

## Compare with SOTA on the CrowdPose Dataset: 2.8-5x measured speedup



CrowdPose mAP ↑ vs GMACs — LitePose (49.5, 58.3, 59.9, 61.9), EfficientHRNet (46.1, 53.6, 56.3), HigherHRNet (50.4, 57.4). 2.8x reduction

CrowdPose mAP ↑ vs Raspberry Pi 4B+ (ms) — 5.0x faster

CrowdPose mAP ↑ vs Qualcomm Snapdragon 855 (ms) — 4.9x faster

CrowdPose mAP ↑ vs Latency on NVIDIA Jetson Nano (ms) — 5.0x faster