# MicroNet: Improving Image Recognition with Extremely Low FLOPs

Yunsheng Li[1], Yinpeng Chen[2], Xiyang Dai[2], Dongdong Chen[2], Mengchen Liu[2],
Lu Yuan[2], Zicheng Liu[2], Lei Zhang[2], Nuno Vasconcelos[1]

[1] University of California San Diego [2] Microsoft

{yul554,nvasconcelos}@ucsd.edu,
{yiche,xidai,dochen,mengcliu,luyuan,zliu,leizhang}@microsoft.com

## Abstract

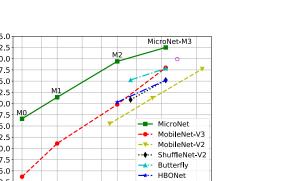This paper aims at addressing the problem of substantial performance degradation at extremely low computational cost (e.g. 5M FLOPs on ImageNet classification). We found that two factors, sparse connectivity and dynamic activation function, are effective to improve the accuracy. The former avoids the significant reduction of network width, while the latter mitigates the detriment of reduction in network depth. Technically, we propose micro-factorized convolution, which factorizes a convolution matrix into low rank matrices, to integrate sparse connectivity into convolution. We also present a new dynamic activation function, named Dynamic Shift-Max, to improve the non-linearity via maxing out multiple dynamic fusions between an input feature map and its circular channel shift. Building upon these two new operators, we develop a family of MicroNet models that achieve significant performance gains over the state of the art in the low FLOP regime.

## 1. Introduction

Recent progress in efficient CNN architectures [...] enables the design of reasonably accurate models under constrained computation.



Figure 1. **Computational Cost (MAdds) vs. ImageNet Accuracy.** MicroNet significantly outperforms the state-of-the-art efficient networks at very low FLOPs (from 6M to 21M MAdds).

## 2. Related Work

**Efficient CNNs:** MobileNet [...], [...], [...] decompose $k \times k$ convolution into a depthwise and a pointwise convolution.

## 3. Micro-Factorized Convolution

The goal of Micro-Factorized convolution is to optimize the trade-off between the number of channels and node connectivity.

## 4. Dynamic Shift-Max

## 5. MicroNet

## 6. Experiments

## 7. Conclusion