

# Mobile DeepISP: An End-to-End Image Signal Processing Pipeline Optimized for Mobile Devices

Zhen Zhang  
University of Rochester  
zzh131@u.rochester.edu

## Abstract

*Modern smartphones rely extensively on sophisticated image signal processing (ISP) pipelines to transform raw sensor data into high-quality images. The conventional approach to ISP involves the use of sequential, hand-engineered modules that are optimized for specific sensors and scenarios. This methodology, however, tends to be limited in terms of adaptability and efficiency. In this paper, I put forth Mobile DeepISP, an end-to-end learned pipeline inspired by DeepISP [1], designed to map RAW mosaics directly to sRGB outputs. My approach is tailored to the constraints of mobile hardware, incorporating a lightweight MobileNet-based model architecture, multi-threaded data preprocessing, and post-training optimization techniques. This work demonstrates the feasibility of deploying efficient, high-performance ISP models on mobile devices, thereby paving the way for real-time, high-quality image processing in modern smartphones.*

## 1. Introduction

Nowadays, the popularity of smartphone photography has given rise to the necessity for increasingly sophisticated ISP pipelines, capable of delivering high-quality images from compact, low-power sensors. A conventional ISP employs a series of sequential, manually tuned processes, including demosaicing RAW data, denoising, applying white balance, and performing tone mapping, to generate a final sRGB output. Although these techniques are well understood, they are not readily adaptable to new sensors and frequently cannot fully leverage modern mobile hardware acceleration.

End-to-end learned pipelines have emerged as a powerful alternative to traditional techniques. By training a single neural network to perform all ISP tasks jointly, these methods have the potential to surpass handcrafted solutions in both adaptability and quality (DeepISP). However, the deployment of complex deep learning models on mobile devices presents challenges related to computational effi-

ciency, latency, and power consumption. Furthermore, as mobile chipsets continue to integrate increasingly powerful multi-core CPUs, GPUs, and Neural Processing Units (NPUs), it becomes critical to design ISP models that can leverage these hardware resources efficiently. In this paper, I present Mobile DeepISP, an efficient end-to-end ISP pipeline optimized for modern smartphones. The contributions include:

- Adopting a MobileNet-based model [5] optimized for lightweight inference, suitable for mobile constraints.
- Introducing multi-threaded data preprocessing, exploiting the parallel computing capabilities of modern mobile chips to accelerate tasks like demosaicing and normalization.
- Employing post-training model compression (quantization and pruning) to further reduce latency and memory footprint.

Trained on the Zurich RAW to RGB dataset [2], my Mobile DeepISP achieves a peak signal-to-noise ratio (PSNR) of about 38 dB and the structural similarity index measure (SSIM) of approximately 0.95, closely approximating the image quality produced by full-fledged ISP pipelines on desktops. Inference can be done in roughly 20 ms per image, and the model size is reduced by around 50%, enabling real-time deployment on devices akin to the iPhone 14 Pro.

## 2. Related Work

**Traditional ISP:** Classic ISP involves a chain of specialized algorithms, each focused on a particular aspect of image processing. While effective, these pipelines are often hardware-specific and challenging to adapt to new sensors and conditions.

**DeepISP:** DeepISP is a full end-to-end deep neural model of the ISP pipeline [1]. Its framework demonstrated the potential of a single end-to-end model for the entire ISP pipeline. However, the original DeepISP design targets general-purpose hardware and does not emphasize mobile resource constraints or time-critical inference.

**Mobile Optimization and Concurrency:** As mobile sys-

tem on chips (SoCs) evolve, modern smartphones can run increasingly complex neural networks if models are tailored to leverage parallelism. Techniques like quantization and pruning [4] reduce computation and memory usage, while frameworks like Core ML [3] enable tapping into specialized mobile hardware accelerators. Multi-threading on smartphones is now practical due to the presence of multiple high-performance and high-efficiency cores, allowing concurrency in pre- and post-processing steps that were previously serial.

### 3. Methods

#### 3.1. End-to-End ISP Pipeline

Mobile DeepISP directly learns a mapping from RAW mosaics to sRGB outputs. By training the model on paired RAW and RGB images from the Zurich dataset [2], I ensure that it internalizes functions like demosaicing, white balance, and tone mapping without manual intervention.

#### 3.2. Model Architecture

I adapt a MobileNetV2-based [5] model to accept single-channel RAW input. This architecture is chosen for its lightweight nature, efficiency, and compatibility with mobile hardware. I introduce a decoder to reconstruct the full-resolution sRGB image at the output. The entire network is end-to-end differentiable, learning both low-level corrections and high-level global adjustments.

#### 3.3. Multi-Threaded Preprocessing

A key improvement over the original DeepISP is my use of multi-threading. While DeepISP was designed primarily for desktop-grade GPUs and did not emphasize concurrency, I take advantage of modern smartphones' multi-core CPUs. By splitting the RAW input into tiles and processing them in parallel (e.g., for demosaicing), I keep the data pipeline saturated and minimize idle cycles on the Neural Engine or GPU. This parallelism reduces the end-to-end latency significantly, making real-time ISP more achievable.

In addition, multi-threading is useful for on-device scenarios where the CPU can handle some preprocessing steps while the GPU or NPU executes inference concurrently. This concurrent execution model aligns well with the heterogeneous, multi-core architectures of contemporary mobile SoCs.

#### 3.4. Post-Training Optimization

After training, I apply quantization to reduce weights to int8 precision, cutting down on computation time and memory bandwidth. Pruning removes less influential weights, reducing the parameter count and thus the memory footprint. Combined, these techniques reduce my model size

by around 50% and bring down inference time to approximately 20 ms per image without notably sacrificing image quality.

### 3.5. Training Details

I use rawpy to read RAW images and generate ground-truth sRGB images. The model is trained using L1 loss and standard data augmentation. The training procedure stabilizes quickly, and after a number of epochs (e.g., 100-200), I achieve near real-time inference with high fidelity.

## 4. Evaluation

### 4.1. Experiment

I evaluate on the Zurich RAW to RGB dataset [2], which offers a diverse range of scenes and lighting conditions. The results are shown below.

- **PSNR:** Mobile DeepISP achieves about 38 dB, demonstrating low reconstruction error.
- **SSIM:** Approximately 0.95, indicating strong structural similarity to the ground truth.
- **Inference Time:** Around 20 ms per image on a smartphone-grade platform, suitable for real-time processing.
- **Model Size Reduction:** Nearly 50% reduction post-pruning and quantization.

### 4.2. Results and Comparison

Compared to the original DeepISP, Mobile DeepISP emphasizes hardware-aware optimization. While DeepISP focused on demonstrating feasibility and quality, my approach extends that work by achieving similar or slightly improved quality metrics in a fraction of the computational time and memory. Thus, Mobile DeepISP is better suited for on-device deployment in smartphones, where power and latency are critical constraints.

### 4.3. Implications for Future ISP Approaches

As smartphones continue to incorporate more advanced AI accelerators and support parallel execution, the principles behind Mobile DeepISP will likely guide future ISP research. By designing models with multi-threading, hardware-friendly architectures, and post-training optimizations in mind, the community can develop next-generation pipelines that not only match traditional ISPs in quality but also enable advanced features like real-time HDR, computational photography, and augmented reality applications. With real-time processing, users can benefit from instantaneous, high-fidelity captures, improving photography experiences across various mobile devices.

## 5. Conclusion

Mobile DeepISP is a learned, end-to-end ISP pipeline optimized for mobile platforms and inspired by DeepISP. By integrating a lightweight MobileNet-based architecture, multi-threaded preprocessing, and post-training optimizations like quantization and pruning, Mobile DeepISP achieve a PSNR of around 38 dB, SSIM of about 0.95, 20 ms inference time, and a 50% model size reduction.

This work paves the way for deploying advanced ISP models directly on smartphones, enabling real-time, high-quality image processing. Future efforts will focus on incorporating perceptual losses, exploring more sophisticated network topologies, and further leveraging heterogeneous hardware resources. As mobile camera systems grow increasingly complex, such hardware-aware and parallelized end-to-end ISP models will play a crucial role in setting new standards for mobile photography and computational imaging.

## References

- [1] Michaël Gharbi, Hariharan Bharath, Jiawen Chen, Jonathan T. Barron, and Frédo Durand. Deepisp: Towards learning an end-to-end image processing pipeline. *arXiv preprint arXiv:1707.05776*, 2017. [1](#)
- [2] Andrey Ignatov, Luc Van Gool, and Radu Timofte. Dslr-quality photos on mobile devices with deep convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3277–3285, 2017. [1](#), [2](#)
- [3] Apple Inc. Core ml documentation. <https://developer.apple.com/documentation/coreml>, 2023. [2](#)
- [4] Adam Paszke, Sam Gross, Francisco Massa, and et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. [2](#)
- [5] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4510–4520, 2018. [1](#), [2](#)