# Image Harmonization with Attention-based Foreground-background Feature Map Modulation
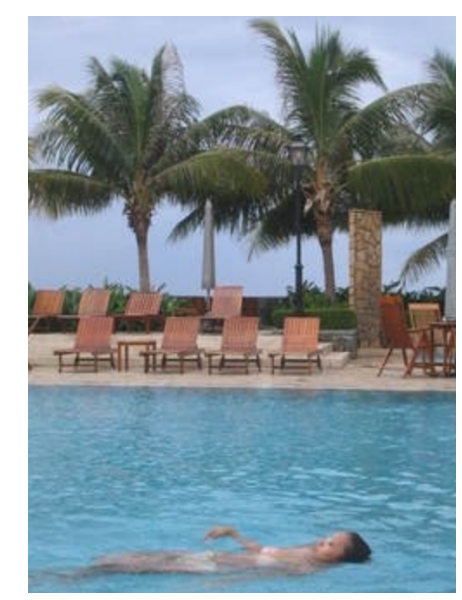
Guoqing Hao, Satoshi Iizuka, Kazuhiro Fukui

*Graduate School of Systems and Information Engineering, University of Tsukuba, Japan*

## Overview

1. We present an underlined attention-based deep feature modulation layer, which allows modulating the feature map of foreground according to those of similarity-weighted background, to improve realism of composites
2. Experimental results on the image harmonization dataset and real composite images show that our method outperforms existing methods both quantitatively and qualitatively

## Background

Foreground

Background

Composited

Harmonized

1. Inharmonious appearance between foreground and background degrades quality of composite images
2. Image harmonization is often conducted manually by experts and requires a significant amount of time
3. Our goal is to generate realistic composites automatically
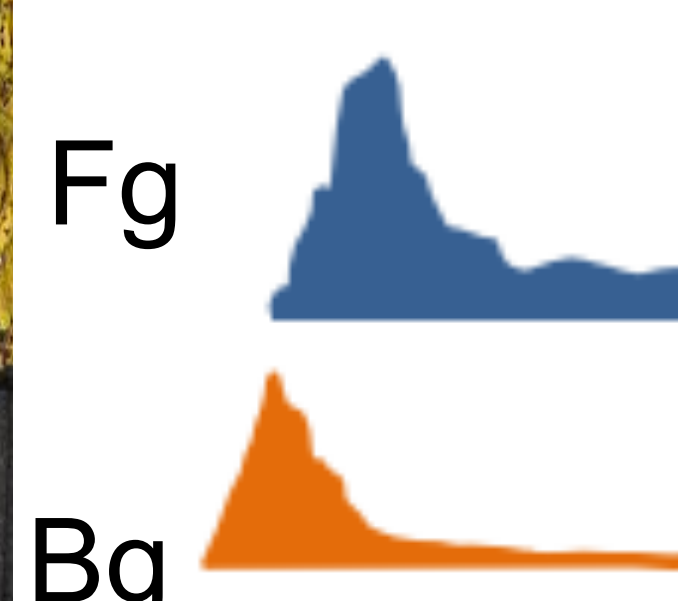
## Background
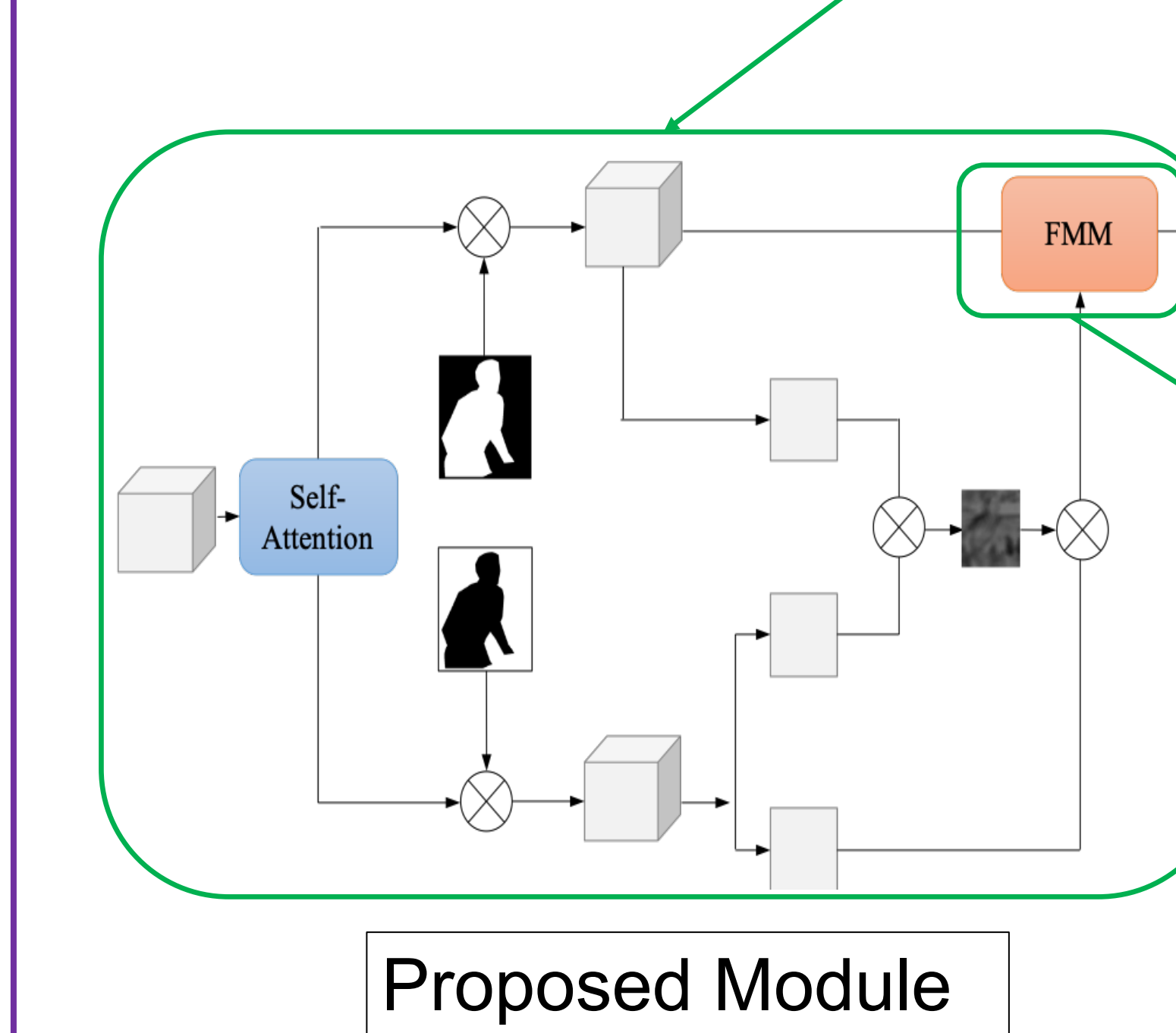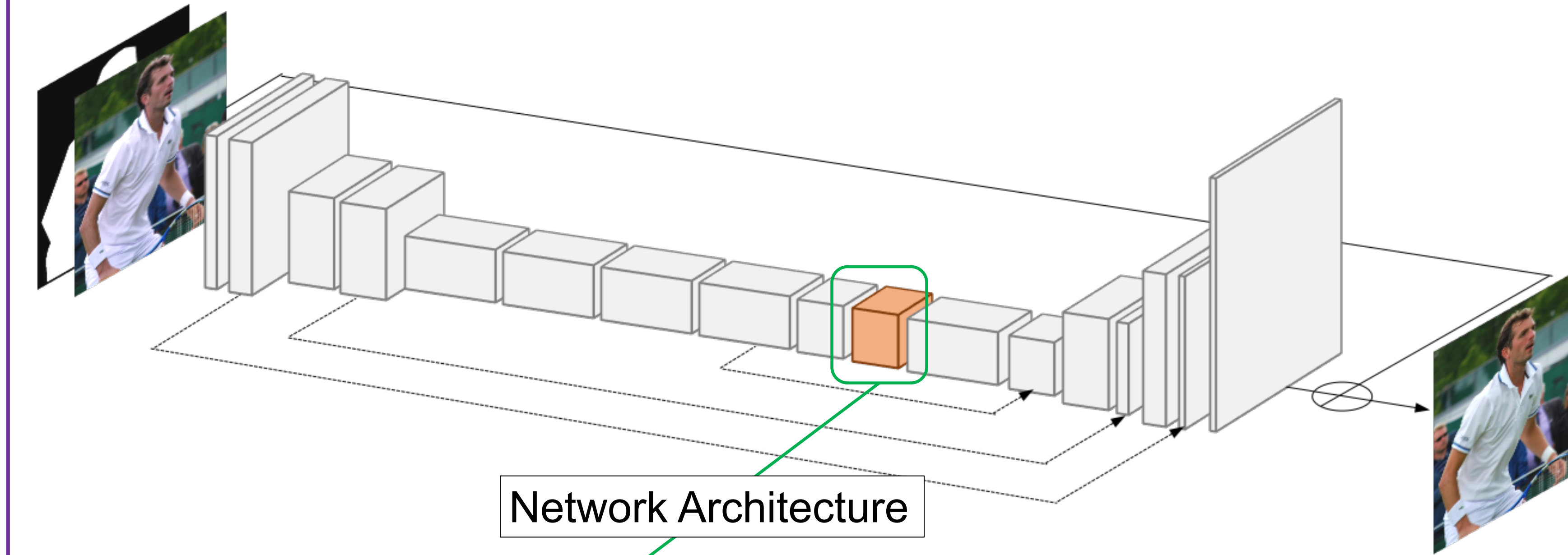
Composited

Histogram

Harmonized

Histogram

Fg

Bg

Our work is motivated by the fact that the specific image statistics such as histogram of luminance between foreground and background typically matches in realistic composite images

## Our Method

Network Architecture

Proposed Module

Self-Attention

FMM

$$FMM(h_f, h_b^a) = \gamma \times (h_f \times \frac{\sigma_{h_b^a}}{\sigma_{h_f}})$$

$\gamma$: learned parameter

$h_f$: foreground features

$h_b^a$: similarity-weighted background features

$\sigma$: standard deviation

**Network:**
- formed exclusively by convolutional layers

**Proposed Module:**
- adjust high-level feature statistics of foreground according to those of background

- capture non-local dependencies between foreground and background

- trained in an end-to-end fashion

- easily inserted into any convolutional neural networks with only a small additional computational cost

## Results

Comparisons between our method against existing methods

| Methods | PSNR | MSE |
|---|---|---|
| S²AM [1] | 34.35 | 59.67 |
| DoveNet [2] | 34.76 | 52.33 |
| Ours | **35.86** | **30.37** |

Ablation results of our proposed module. The "baseline" stands for the backbone net-work in our full method. The "A" stands for "remove self-attention block from our proposed module"

| Methods | PSNR |
|---|---|
| Baseline | 32.98 |
| Baseline + Self-attention [3] | 35.06 |
| Baseline + A | 35.17 |
| Baseline + proposed module | **35.86** |

## References

[1]. X. Cun and C. Pun. Improving the harmony of the composite image by spatial-separated attention module. IEEE Transactions on Image Processing. 2020

[2]. W. Cong, J. Zhang, L. Niu, L. Liu, Z. Ling, W. Li, and L. Zhang. Dovenet: Deep image harmonization via domain verification. In CVPR, 2020.

[3]. Han Zhang, Ian J. Goodfellow, Dimitris N. Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In ICML, 2019