



清華大學

Tsinghua University



西湖大學

WESTLAKE UNIVERSITY

ETH zürich



ICML

International Conference  
On Machine Learning

# Flow-Guided Sparse Transformer for Video Deblurring

ICML 2022

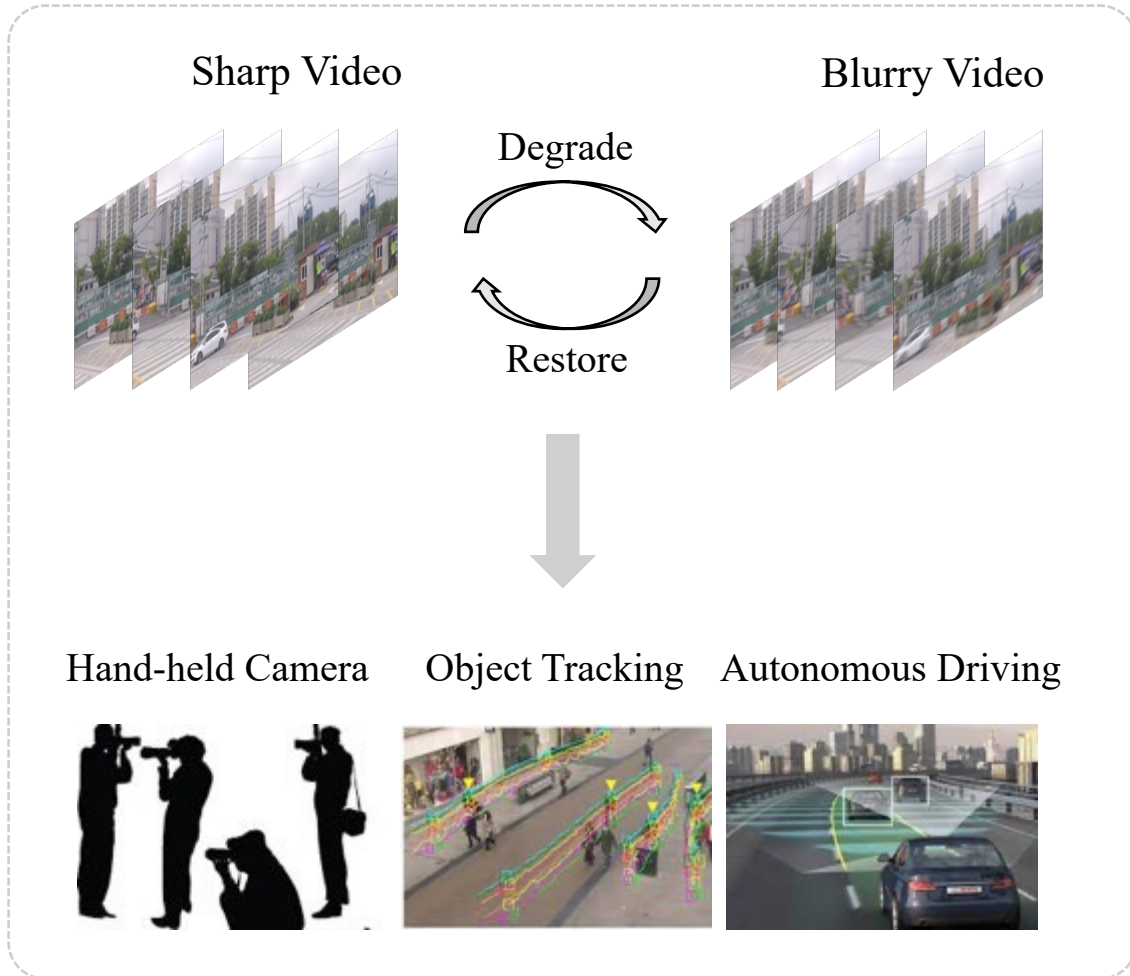
Jing Lin<sup>\*1</sup>, Yuanhao Cai<sup>\*1</sup>, Xiaowan Hu<sup>1</sup>, Haoqian Wang<sup>†1</sup>, Youliang Yan<sup>2</sup>  
Xueyi Zou<sup>†2</sup>, Henghui Ding<sup>3</sup>, Yulun Zhang<sup>3</sup>, Radu Timofte<sup>3</sup>, and Luc Van Gool<sup>3</sup>

The Shenzhen International Graduate School, Tsinghua University<sup>1</sup>  
Huawei Noah's Ark Lab<sup>2</sup>, ETH Zürich<sup>3</sup>

# Outline

- Background and Motivation
- The Proposed Flow-Guided Sparse Transformer
  - FGST: Overall framework
  - FGS-MSA
  - FGSW-MSA
  - RE
- Experiment Results

# Video Deblurring



## Existing Methods

- Conventional Methods: Based on hand-crafted prior, poor generalization ability, and limited representation capacity
- CNN-based Methods: Show limitations in capturing long-range dependencies and non-local self-similarity

Transformer ?

# Transformer

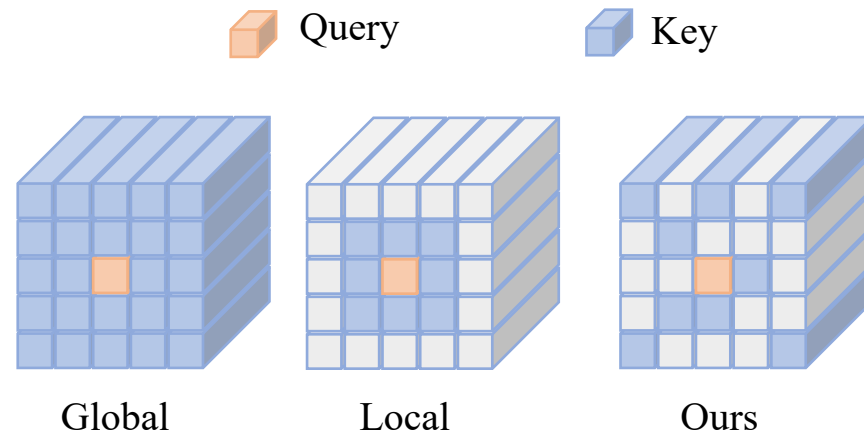
- Global Transformer: non-trivial computation cost
- Local Transformer: local receptive field, may miss some content-related tokens when fast motion exists



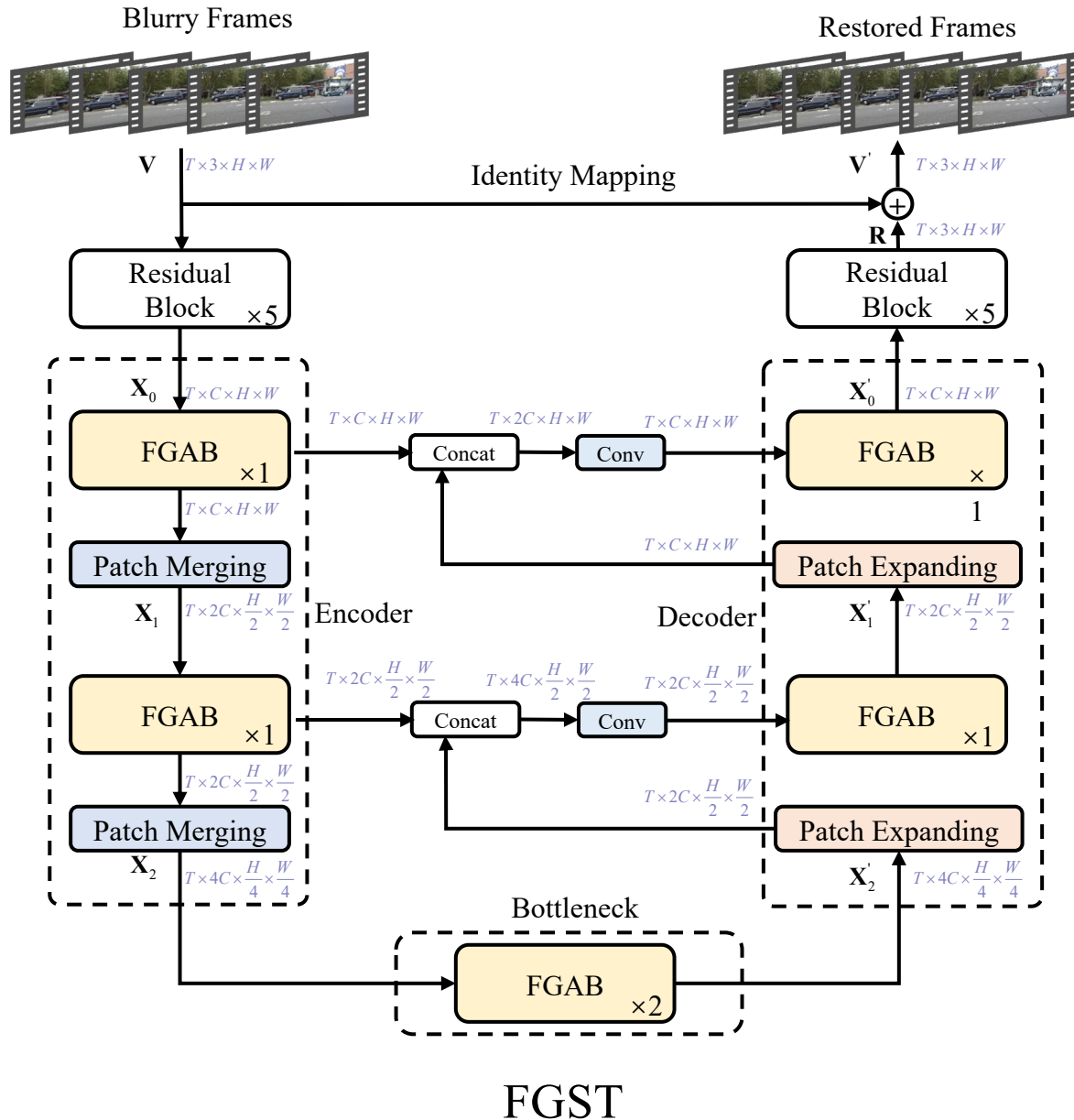
Previous Transformers lack the guidance of motion information when computing self-attention



Integrate optical flow into self-attention module

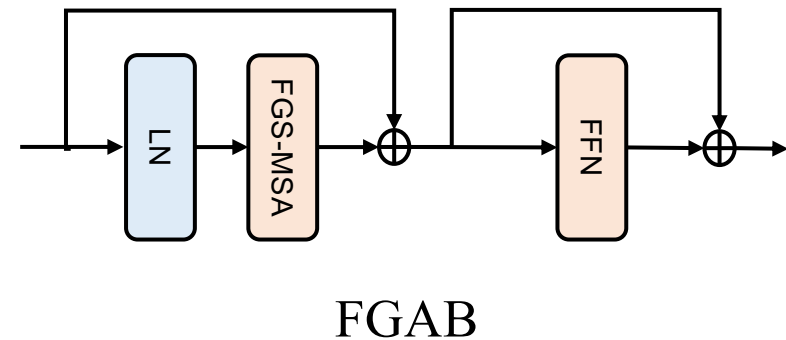


# Framework

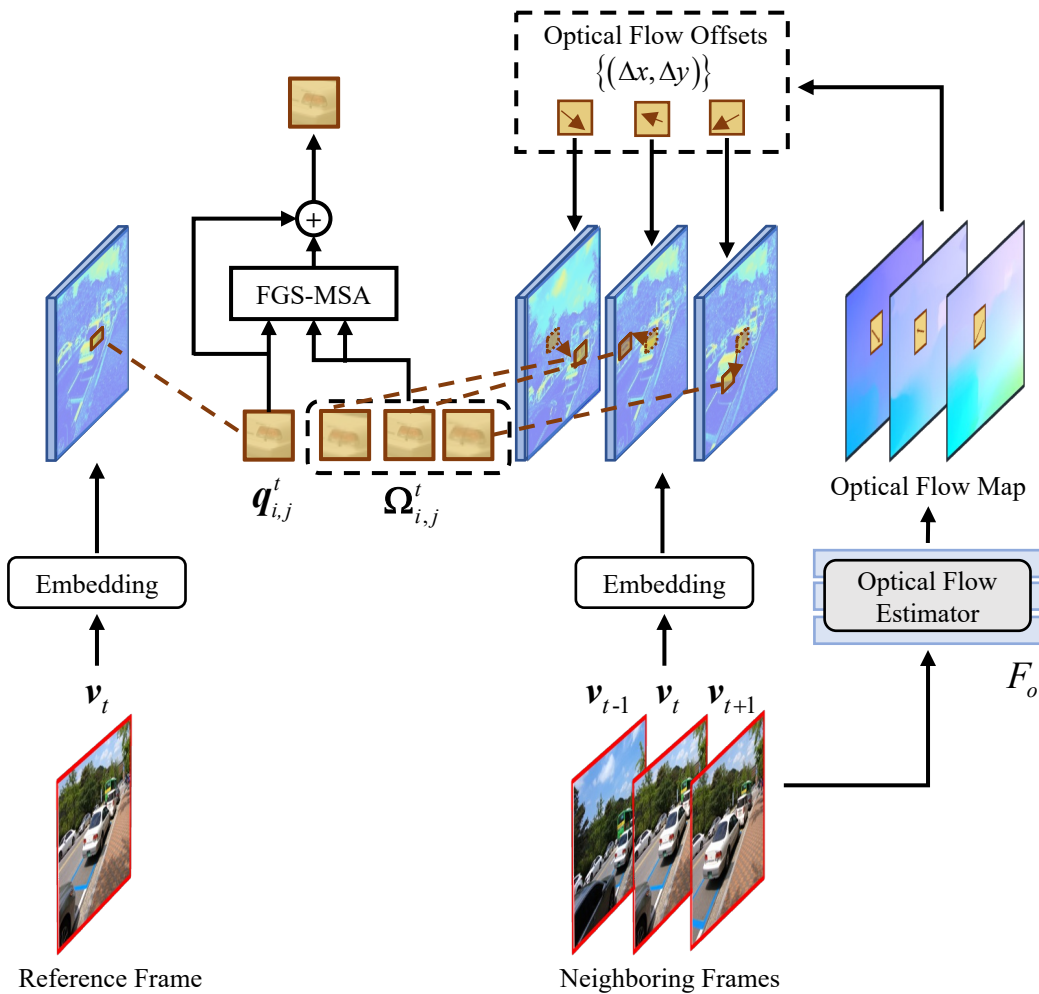


## Flow-Guided Sparse Transformer

- The first Transformer-based method for video deblurring
- Adopts a U-shaped structure consisting of an encoder, a bottleneck, and a decoder
- Built up by Flow-Guided Attention Blocks (FGABs)



# FGS-MSA



FGS-MSA

## Flow-Guided Sparse Multi-head Self-Attention

- Optical Flow Estimation

$$(\Delta x_f, \Delta y_f) = F_o(v_t, v_f) (i, j)$$

- Key Elements Sampling

$$\Omega_{i,j}^t = \{\mathbf{k}_{i+\Delta x_f, j+\Delta y_f}^f \mid |f - t| \leq r\}$$

- Self-Attention Calculation

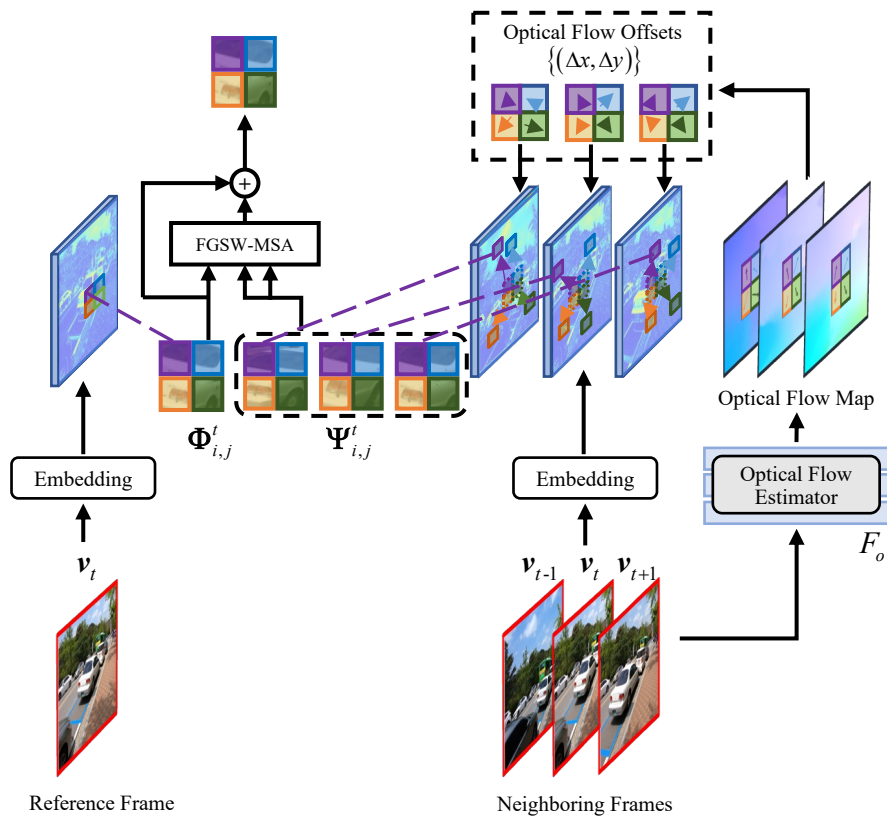
$$\text{FGS-MSA}(q_{i,j}^t, \Omega_{i,j}^t) = \sum_{n=1}^N \mathbf{W}_n \sum_{\mathbf{k} \in \Omega_{i,j}^t} \mathbf{A}_{nq_{i,j}^t, \mathbf{k}} \mathbf{W}'_n \mathbf{k},$$

Enjoy global receptive fields and linear computation complexity

$$O(\text{global MSA}) = 4(THW)C^2 + 2(THW)^2C,$$

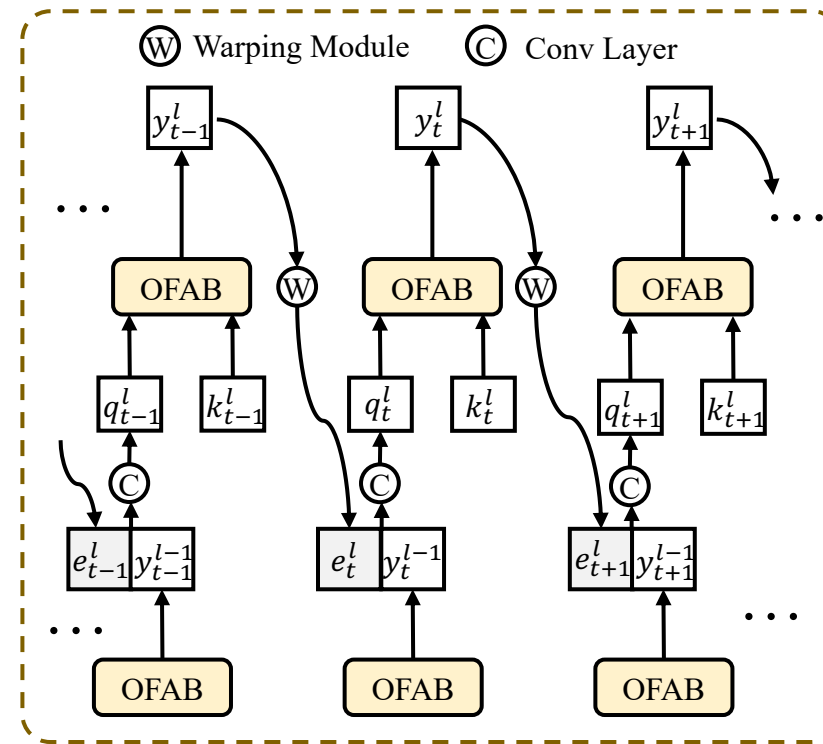
$$O(\text{FGS-MSA}) = 2(THW)C(2(r+1)C + 2r + 1).$$

# Improvements



FGSW-MSA

- FGSW-MSA: more robust to accommodate pixel-level flow offset prediction deviations



RE

- RE: Inspired by RNN, to establish long-range temporal dependencies

# Experiment

Method	EDVR (Wang et al. 2019)	Tao et al. (Tao et al. 2018)	Su et al. (Su et al. 2017)	DBLRNet (Zhang et al. 2018)	STFAN (Zhou et al. 2019)	Xiang et al. (Xiang et al. 2020)	TSP (Pan et al. 2020)	Suin et al. (Suin et al. 2021)	ARVo (Li et al. 2021)	<b>FGST (Ours)</b>
PSNR $\uparrow$	28.51	29.98	30.01	30.08	31.15	31.68	32.13	32.53	32.80	<b>33.36</b>
SSIM $\uparrow$	0.864	0.884	0.888	0.885	0.905	0.916	0.927	0.947	0.935	<b>0.950</b>

Tab. 1 Quantitative Comparison with SOTA methods on DVD dataset.

Method	RDN (Patrick et al. 2017)	Kim et al. (Kim et al. 2015)	EDVR (Wang et al. 2019)	Su et al. (Su et al. 2017)	STFAN (Zhou et al. 2019)	Nah et al. (Nah et al. 2019)	Tao et al. (Tao et al. 2018)	TSP (Pan et al. 2020)	Suin et al. (Suin et al. 2021)	<b>FGST (Ours)</b>
PSNR $\uparrow$	25.19	26.82	26.83	27.31	28.59	29.97	30.29	31.67	32.10	<b>32.90</b>
SSIM $\uparrow$	0.779	0.825	0.843	0.826	0.861	0.895	0.901	0.928	0.960	<b>0.961</b>

Tab. 2 Quantitative Comparison with SOTA methods on GOPRO dataset.

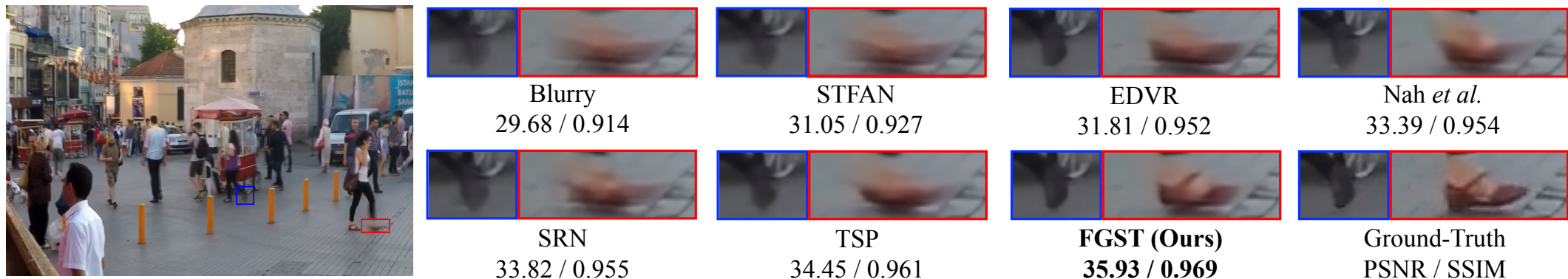


Fig. 1 Qualitative Comparison with SOTA methods.

Our FGST significantly outperforms SOTA methods quantitatively and qualitatively.



# Thanks



Code & Paper