

# Motion Estimation for Quanta Image Sensors Using Spatio-Temporal Priors

Hiroya Fukawa<sup>†</sup>, Kosuke Kurihara<sup>†</sup>, Yoshihiro Maeda<sup>‡</sup>, Shunichi Sato<sup>†</sup>, and Takayuki Hamamoto<sup>†</sup>

<sup>†</sup>Tokyo University of Science, Tokyo, 125-8585, Japan    <sup>‡</sup>Shibaura Institute of Technology, Tokyo, 135-8548, Japan

<sup>‡</sup>ymaeda@shibaura-it.ac.jp

**Abstract**—Quanta image sensors are a novel paradigm in image sensor technology. Their direct application to quanta image sensors-based imaging systems is challenging because a bit-plane image is a set of binary images. In this paper, we introduce spatio-temporal priors based on the intensity invariance and smoothness characteristics of the motion vector. Specifically, we model when the image sequences align with the correct motion vector, the spatiotemporal structure becomes more consistent. Moreover, the spatial smoothness prior is incorporated through the smoothing filtering of the evaluation metrics of motion vector candidates. The experimental results show that the proposed method is more effective than conventional methods.

**Index terms**— motion estimation, quanta image sensors, bit-plane images, spatio-temporal prior

## I. INTRODUCTION

Quanta image sensors (QISs) represent a novel paradigm in image sensor technology [1]–[3]. These sensors can detect individual photons while maintaining ultralow read noise [4]–[6]. They also feature a high temporal resolution of over 10,000 frames per second (fps), enabling the detection of individual photons. QIS-based imaging systems produce bit-plane images, which are time sequences of binary images that indicate whether a photon is incident [3]. These sensors outperformed traditional CMOS image sensors in low-light environments and for imaging high-speed moving objects. Therefore, QIS-based imaging systems can be employed in various applications, including high-dynamic-range (HDR) imaging [7]–[10] and low-light environment imaging [11], [12].

QIS-based imaging systems are particularly suitable for motion estimation tasks because they facilitate the capture of motion information from objects moving at ultra-high speeds. However, QIS-based imaging systems produce bit-plane images that fundamentally differ from multibit natural images. Consequently, it is difficult to estimate motion vectors with high accuracy for QIS-based imaging systems by applying

This work was supported by JSPS KAKENHI JP24K02964 and JP24K23903.

This article has been accepted for publication in IEEE International Conference on Visual Communications and Image Processing (VCIP) 2024. This is the author's version which has not been fully edited and content may change prior to final publication. Citation information: DOI 10.1109/VCIP63160.2024.10849878

© 2024 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

conventional motion estimation methods to multi-bit imaging systems. Therefore, we must develop a motion estimation algorithm adapted specifically for QIS-based imaging systems.

A few motion estimation methods for QIS-based imaging systems have been reported [13]–[15]. The conventional motion estimation methods for these systems can be categorized into two approaches: indirect estimation from bit-plane images [13] and direct estimation from bit-plane images [14], [15]. Jungerman *et al.* [13] proposed an indirect motion estimation method from bit-plane images. They reconstructed multiple multi-bit images from bit-plane images and estimated the motion vectors using the obtained reconstructed image sequences. However, this approach significantly degrades the time resolution because numerous bit-plane images are required to reconstruct a single multi-bit image. Moreover, motion within the bit-plane images causes motion blur in the reconstructed image; thus, the accuracy of motion estimation is degraded. These issues indicate that indirect motion estimation techniques do not fully leverage QIS-based imaging systems.

Unlike indirect methods, direct estimation methods directly estimate the motion from bit-plane images [14], [15]. These methods extract the time sequence of incident photons from bit-plane images along candidate motion vectors and statistically evaluate the temporal variation in incident photons caused by motion. To evaluate the temporal variation, the method proposed by Gyongy *et al.* [14] employs the confidence interval for the probability of incident photons in the spatio-temporally reduced cubicle of the bit-plane images. In addition, Iwabuchi *et al.* [15] used chi-square values to analyze the time sequence of the incident photons in a single pixel. This statistical approach is primarily motivated by the fact that the number of incident photons fluctuates over time, even without motion, because the behavior of the incident photons follows the Poisson process [7]–[10]. Although these methods utilize the temporal characteristics of incident photons, they do not effectively leverage the priors of motion vectors, such as the spatial intensity invariance between the corresponding pixels in adjacent frames and the smoothness of the motion vectors in neighboring pixels [16].

The incorporation of motion vector priors is a potential approach for accurate motion estimation in QIS-based imaging systems. These priors have contributed to the success of accurate motion estimation of multi-bit image sequences. However, their direct application to QIS-based systems is challenging because a bit-plane image is a set of binary images. Therefore,

these priors must be reconsidered within the framework of QIS-based imaging systems.

In this study, we propose a novel direct motion estimation method for QIS-based imaging systems that introduce *spatio-temporal priors* based on the intensity invariance and smoothness characteristics of the motion vector. Inspired by the intensity invariance characteristics used in the optical flow framework [17], we introduce a spatio-temporal consistency prior derived from the spatial gradient of the photon intensity. As known as intensity invariance characteristics, discriminative spatial structures, such as object boundaries and textures, provide useful cues for motion estimation [17]. We re-examined the aforementioned characteristics for spatio-temporal consistency. Specifically, we model when the bit-plane image sequences align with the correct motion vector, the spatio-temporal structure becomes more consistent. This leads to an increase in the spatial gradient of the photon intensity calculated from the sum of aligned bit-plane image sequences; thus, this metric can effectively assess the reliability of the motion estimation. Furthermore, the spatial smoothness prior was incorporated through the smoothing filtering of the evaluation metrics of the motion vector candidates. Considering these priors, we formulate the motion vector estimation as a weighted optimization problem based on a joint filtering framework [18].

## II. PROPOSED METHOD

### A. Overview

The proposed method estimates motion vectors by introducing spatio-temporal priors directly from bit-plane images. Figure 1 shows an overview of the proposed method. Inspired by [17], the objective function for the motion vector estimation is formulated as follows:

$$\mathbf{u}_q^* = \arg \min_{\mathbf{u} \in \psi} \sum_{\mathbf{r} \in \mathcal{N}} w_s w_t E_r^{(\mathbf{u})}, \quad (1)$$

where  $\mathbf{u}_q^*$  is the estimated motion vector at pixel position  $\mathbf{q}$ ,  $E_r^{(\mathbf{u})}$  is the error evaluation metric for motion vector candidate  $\mathbf{u}$  at  $\mathbf{r}$ , and  $\mathcal{N}$  is the set of pixel positions in the neighborhood centered on  $\mathbf{q}$ . In addition,  $w_s$  and  $w_t$  are weights based on the spatial and temporal priors, respectively.

Eq. (1) implies that the error evaluation metrics  $E_r^{(\mathbf{u})}$  are filtered using weights based on the spatial and temporal priors,  $w_s$  and  $w_t$ . The error evaluation metrics for bit-plane images in this method correspond to the representation of the intensity invariance prior, although this constraint is weak owing to the stochastic time fluctuations of the photons. Moreover, conventional methods do not incorporate a smoothing prior. Therefore, the proposed method introduces weights  $w_s$  and  $w_t$  as error evaluation metrics to effectively represent these priors.

### B. Calculation of Chi-Square Values as Error Evaluation Metrics

The proposed method employs chi-square values,  $(\chi^2)_q^{(\mathbf{u})}$ , as error evaluation metrics,  $E_r^{(\mathbf{u})} := (\chi^2)_q^{(\mathbf{u})}$ , similar to

the approach described in [15]. This is because the incident photons are known to have stochastic characteristics. Using chi-square values, we can stochastically evaluate the temporal variation in the number of photons incident on the region of interest.

The set of bit-plane images  $\{\mathbf{B}_t\}_{t=1}^T$  is spatially shifted according to the candidate motion vector  $\mathbf{u}$  to obtain a set of images  $\{\mathbf{B}_t^{(\mathbf{u})}\}_{t=1}^T$ . From this set, we extracted the spatial patches of  $K \times K$  pixels centered at position  $\mathbf{q}$  within the image-plane domain. The time sequence of these extracted patches is denoted as  $\{\{\mathbf{B}_{t,\mathbf{q}}^{(\mathbf{u})}\}_{\mathbf{q} \in \Omega}\}_{t=1}^T$ , where  $\Omega$  represents the set of  $K \times K$  pixel positions centered at  $\mathbf{q}$ .

We divide the time sequence of these extracted patches,  $\{\{\mathbf{B}_{t,\mathbf{q}}^{(\mathbf{u})}\}_{\mathbf{q} \in \Omega}\}_{t=1}^T$ , into  $G (= \lceil T/M \rceil)$  groups along the time direction, where  $M$  represents the number of frames in each group. Subsequently, the total number of photons incident in each group was calculated. The total photon count for the  $g$ -th group at position  $\mathbf{q}$ ,  $s_{g,\mathbf{q}}^{(\mathbf{u})}$ , is computed as follows:

$$s_{g,\mathbf{q}}^{(\mathbf{u})} = \sum_{t=M(g-1)}^{gM-1} \sum_{\mathbf{q} \in \Omega} \mathbf{B}_{t,\mathbf{q}}^{(\mathbf{u})}. \quad (2)$$

The total photon count,  $s_{g,\mathbf{q}}^{(\mathbf{u})}$ , is known to follow a binomial distribution, modeled by the  $MK^2$  Bernoulli trials of whether a photon was detected [15].

We compute the chi-square values,  $(\chi^2)_q^{(\mathbf{u})}$ , using the formula from  $\{s_{g,\mathbf{q}}^{(\mathbf{u})}\}_{g=1}^G$  as described in the method [15]:

$$(\chi^2)_q^{(\mathbf{u})} = \sum_{g=1}^G \left( Z_{g,\mathbf{q}}^{(\mathbf{u})} \right)^2, \quad (3)$$

where  $Z_{g,\mathbf{q}}^{(\mathbf{u})}$  is the standardized value of  $s_{g,\mathbf{q}}^{(\mathbf{u})}$  approximated to the standard normal distribution using the de Moivre-Laplace theorem [19]:

$$Z_{g,\mathbf{q}}^{(\mathbf{u})} = \frac{s_{g,\mathbf{q}}^{(\mathbf{u})} - MK^2 p_q^{(\mathbf{u})}}{\sqrt{MK^2 p_q^{(\mathbf{u})} (1 - p_q^{(\mathbf{u})})}}, \quad (4)$$

where  $p_q^{(\mathbf{u})}$  denotes the mean probability of detecting photons within the set of  $\{\{\mathbf{B}_{t,\mathbf{q}}^{(\mathbf{u})}\}_{\mathbf{q} \in \Omega}\}_{t=1}^T$  in the spatio-temporal domain, calculated as follows:

$$p_q^{(\mathbf{u})} = \frac{1}{TK^2} \sum_{t=1}^T \sum_{\mathbf{q} \in \Omega} \mathbf{B}_{t,\mathbf{q}}^{(\mathbf{u})}. \quad (5)$$

### C. Spatio-Temporal Priors to Chi-Square Values

We introduce spatio-temporal priors to the chi-square values obtained  $(\chi^2)_q^{(\mathbf{u})}$  by using a joint bilateral filtering scheme [20]. Spatio-temporal priors represent the spatial consistency of motion vector smoothness and the intensity values derived from the time sequence of incident photons.

We first reconstruct the guidance image,  $\mathbf{I}^{(\mathbf{u})}$ , for each motion vector candidate,  $\mathbf{u}$ , from  $\{\mathbf{B}_t^{(\mathbf{u})}\}_{t=1}^T$ , using the image reconstruction method [3], which reconstructs the multibit image from bit-plane images. Regions of the image  $\mathbf{I}^{(\mathbf{u})}$  where

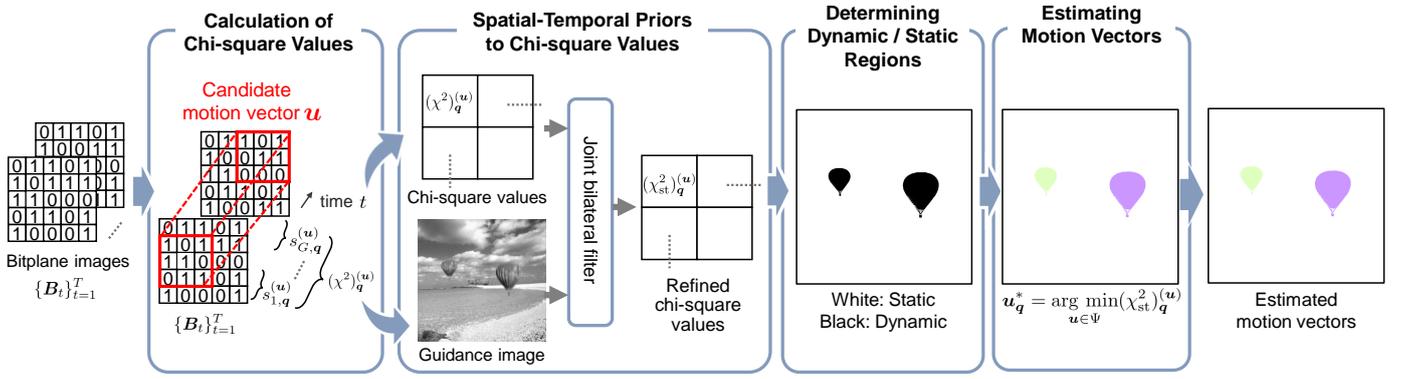


Fig. 1: Overview of the proposed method. The proposed method consists of four steps: calculating the error evaluation metrics using chi-square values, introducing spatio-temporal priors to the chi-square values, determining the dynamic or static region, and estimating the motion vectors.

the motion vector candidates match the correct motion vector without motion blur. This indicates that in regions without blurring, the image preserves its spatial structure; thus,  $\mathbf{I}^{(u)}$  is suitable as a guidance image. Furthermore, the intensity of the guidance image corresponds to the temporal features of the bit-plane images. Consequently, using the image as a guidance image in the joint filter for the chi-square values, we can refine these values while preserving the spatial smoothness of the motion vectors. Moreover, the intensity difference between neighboring pixels is similar to the intensity invariance prior in general motion estimation methods. This ensures that the refinement process improves the accuracy and reliability of the motion vector estimation.

Using  $\mathbf{I}^{(u)}$  as the guidance image, we then obtain refined chi-square values,  $(\chi_{st}^2)_q^{(u)}$ , applying the joint bilateral filtering to  $(\chi^2)_q^{(u)}$ :

$$(\chi_{st}^2)_q^{(u)} = \frac{1}{\eta} \sum_{\mathbf{r} \in \mathcal{N}} w_s w_t (\chi^2)_r^{(u)}, \quad (6)$$

$$w_s := \exp\left(\frac{\|\mathbf{q} - \mathbf{r}\|_2^2}{-2\sigma_s^2}\right), \quad (7)$$

$$w_t := \exp\left(\frac{\|\mathbf{I}_q^{(u)} - \mathbf{I}_r^{(u)}\|_2^2}{-2\sigma_t^2}\right), \quad (8)$$

where  $\sigma_s$  and  $\sigma_t$  are the smoothing parameters for the spatial and range spaces, respectively. In addition,  $\eta$  is the normalization term, and  $\mathcal{N}$  is the set of pixel positions in the neighborhood pixels centered on  $\mathbf{q}$ .

The weight  $w_t$  calculated from the intensity difference between neighborhood pixels in the guidance image represents the similarity of the time sequence of the incident photons. In addition, combined with the weights for spatial smoothness  $w_s$ , Eq. 6 favors spatially and temporally consistent neighborhoods. Thus, these weights effectively improved the estimation accuracy of the motion vectors.

#### D. Determining the Dynamic and Static Regions

We determined the dynamic and static regions to precisely estimate the motion vectors. The number of incident photons

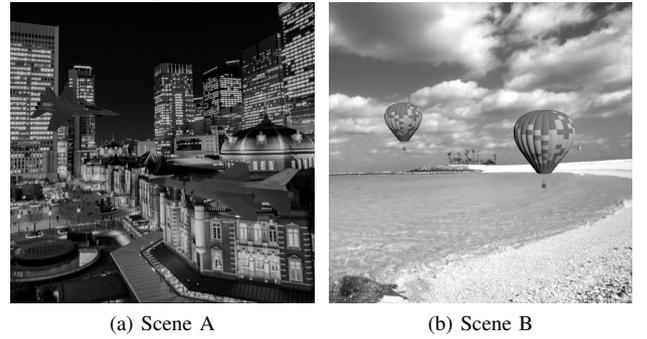


Fig. 2: Test scenes: In Figs. (a) and (b), there are two moving airplanes and balloons, respectively. In both figures, the moving objects move vertically and horizontally, respectively.

in the static regions was stable and its temporal variation was low. Therefore, the chi-square values for static regions are ideally small. However, in practice, the number of incident photons fluctuates because it follows a Poisson process [3]. These variations in the number of photons cause performance degradation in the motion vector estimation. Therefore, the proposed method determines the dynamic and static regions before estimating motion vectors.

To determine the dynamic or static regions, we applied thresholding to  $(\chi_{st}^2)_q^{(u=0)}$ , which is a refined chi-squared value when the non-motion condition is assumed (i.e.,  $\mathbf{u} = \mathbf{0}$ ). The set of pixel positions in the dynamic region  $\Phi$  is obtained as follows:

$$\Phi = \{\mathbf{q} \mid \alpha < (\chi_{st}^2)_q^{(u=0)}\}, \quad (9)$$

where  $\alpha$  is the threshold value for determining the dynamic or static regions. A high value of  $(\chi_{st}^2)_q^{(u=0)}$  suggests a significant temporal variation, indicating a dynamic region. This thresholding process functions as a chi-squared test for  $(\chi_{st}^2)_q^{(u=0)}$ , where  $\alpha$  is set according to the significance level of the chi-squared distribution with  $G - 1$  degrees of freedom [15].

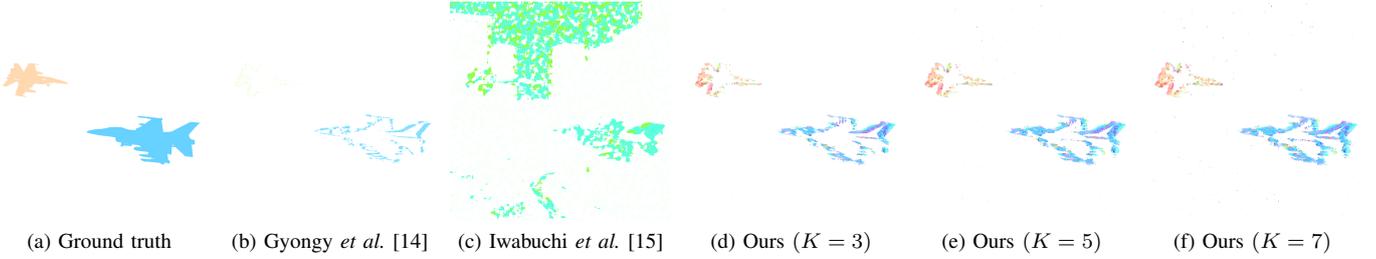


Fig. 3: Visual comparisons of estimated motion vectors (Scene A).

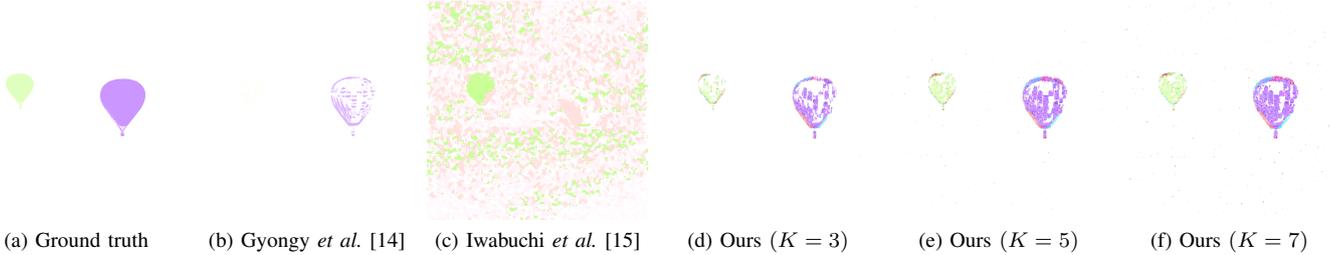


Fig. 4: Visual comparisons of estimated motion vectors (Scene B).

TABLE I: Quantitative results of EPE.

	[14]	[15]	The proposed method		
			$K = 3$	$K = 5$	$K = 7$
Scene A	2.324	9.598	2.378	2.302	<b>2.271</b>
Scene B	1.018	6.616	1.059	0.967	<b>0.956</b>

### E. Estimating Motion Vectors

The optimal motion vector minimizes the error evaluation metric of the incident photons. Thus, we explore the motion vector among all candidate motion vectors  $\Psi$  such that  $(\chi_{st}^2)^{(u)}_{\mathbf{q}}$  is minimized. In addition, the estimated motion vector in the static regions should be  $\mathbf{0}$ . The estimated motion vector at the pixel position  $\mathbf{q}$ ,  $\mathbf{u}_{\mathbf{q}}^*$ , is obtained as follows:

$$\mathbf{u}_{\mathbf{q}}^* = \begin{cases} \arg \min_{\mathbf{u} \in \Psi} (\chi_{st}^2)^{(u)}_{\mathbf{q}} & \text{if } \mathbf{q} \in \Phi \\ \mathbf{0} & \text{otherwise.} \end{cases} \quad (10)$$

## III. EXPERIMENTS

### A. Experimental setting

To demonstrate the effectiveness of the proposed method, we conducted experiments using synthesized bit-plane images. The bit-plane images were synthesized by simulating the QIS imaging process. We used the synthesized bit-plane images as input data for the motion estimation methods. We created 300 bit-plane images (i.e.,  $T = 300$ ) with a  $1024 \times 1024$  resolution. As shown in Fig. 2, two scenes containing moving objects were used for the evaluation. Scene A contains two airplanes, and Scene B contains two balloons as moving objects. These moving objects moved linearly with a constant velocity in the vertical and horizontal direction, respectively.

For comparison, we used methods [14], [15], that directly estimate motion from a set of bit-plane images. We ensured that the control parameters of the other methods were optimal.

Furthermore, we quantitatively evaluated the results using the endpoint error (EPE), which is denoted as the Euclidean distance between the ground truth and the estimated motion vectors. Lower EPE indicates better performance. We also visualized the direction and magnitude of the estimated motion vectors for each pixel using a color map.

Based on the preliminary experiments, we set the control parameters for the proposed method as follows. We set the kernel size of the joint bilateral filter to  $7 \times 7$  pixels and the smoothing parameters  $\sigma_s$  and  $\sigma_t$  to 1.0 and 7.0, respectively. We set the number of bit-plane images in each group  $M$  and the number of groups as 5 and 60, respectively. The threshold value  $\alpha$  was set to a significance level of 1% in the chi-square distribution with  $G - 1$  degrees of freedom. To investigate the impact of the patch size parameter  $K$ , we set  $K = 3, 5, 7$ . The experiments were run on a Windows OS with an Intel Core i7-10700KF and 64 GB RAM. We used MATLAB R2024a for the experiments.

### B. Experimental results

Table I shows the comparison results using the average EPEs of the estimated motion vectors of each pixel for scene A and B. The proposed method has a lower EPE than the comparison methods, indicating that the proposed method can accurately estimate the motion vector. We can also see that as the patch size parameter  $K$  increases, the EPE becomes lower. This indicates that incorporating an extensive spatio-temporal information into motion estimation contributes to improved performance.

Figures 3 and 4 present the results of visualizing the estimated motion vectors for each scene. The proposed method can estimate the motion vectors more accurately than comparison methods by incorporating spatio-temporal priors.

#### IV. CONCLUSION

In this paper, we propose a novel direct motion estimation method for QIS-based imaging systems that introduce spatio-temporal priors based on intensity invariance and smoothness priors. The proposed method achieves more accurate motion estimation than conventional methods by integrating the spatio-temporal information contained in bit-plane images and statistically evaluating its temporal variation.

#### REFERENCES

- [1] Jiaju Ma, Stanley Chan, and Eric R. Fossum, "Review of quanta image sensors for ultralow-light imaging," *IEEE Transactions on Electron Devices*, vol. 69, no. 6, pp. 2824–2839, 2022.
- [2] Mohammed A. Al-Rawhani, James Beeley, and David R. S. Cumming, "Wireless fluorescence capsule for endoscopy using single photon-based detection," *Scientific Reports*, vol. 5, 2015.
- [3] Feng Yang, Y. M. Lu, L. Sbaiz, and M. Vetterli, "Bits from photons: Oversampled image acquisition using binary poisson statistics," *IEEE Transactions on Image Processing*, vol. 21, no. 4, pp. 1421–1436, Apr. 2012.
- [4] Jing Gao, Xiaoxing Du, Kaiming Nie, and Jiangtao Xu, "Analysis and elimination method of the non-fixed column noise for quantum image sensor," *IEEE Sensors Journal*, vol. 20, no. 1, pp. 318–327, 2020.
- [5] Jiaju Ma and Eric R. Fossum, "Quanta image sensor jot with sub 0.3e-r.m.s. read noise and photon counting capability," *IEEE Electron Device Letters*, vol. 36, no. 9, pp. 926–928, 2015.
- [6] Shuichi Namiki, Shunichi Sato, Yusuke Kameda, and Takayuki Hamamoto, "Imaging method using multi-threshold pattern for photon detection of quanta image sensor," in *Proc. International Workshop on Advanced Imaging Technology (IWAIT)*, Apr. 2022, vol. 12177, p. 1217702.
- [7] Stanley H. Chan, Omar A. Elgendy, and Xiran Wang, "Images from bits: Non-iterative image reconstruction for quanta image sensors," *Sensors*, vol. 16, no. 11, 2016.
- [8] Abhiram Gnanasambandam and Stanley H. Chan, "Hdr imaging with quanta image sensors: Theoretical limits and optimal reconstruction," *IEEE Transactions on Computational Imaging*, vol. 6, pp. 1571–1585, 2020.
- [9] Omar A. Elgendy and Stanley H. Chan, "Image reconstruction and threshold design for quanta image sensors," in *Proc. IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 978–982.
- [10] Neale A.W. Dutton, Tarek Al Abbas, Istvan Gyongy, Francescopaolo Mattioli Della Rocca, and Robert K. Henderson, "High dynamic range imaging at the quantum limit with single photon avalanche diode-based image sensors," *Sensors*, vol. 18, no. 4, 2018.
- [11] Stanley H. Chan and Yue M. Lu, "Efficient image reconstruction for gigapixel quantum image sensors," in *Proc. IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, 2014, pp. 312–316.
- [12] Abhiram Gnanasambandam, Omar Elgendy, Jiaju Ma, and Stanley H. Chan, "Megapixel photon-counting color imaging using quanta image sensor," *Opt. Express*, vol. 27, no. 12, pp. 17298–17310, Jun 2019.
- [13] S. Jungerman, A. Ingle, and M. Gupta, "Panoramas from photons," in *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 10592–10602.
- [14] Istvan Gyongy, Neale Dutton, and Robert Henderson, "Single-photon tracking for high-speed vision," *Sensors*, vol. 18, no. 2, Jan. 2018.
- [15] Kiyotaka Iwabuchi, Yusuke Kameda, and Takayuki Hamamoto, "Image quality improvements based on motion-based deblurring for single-photon imaging," *IEEE Access*, vol. 9, pp. 30080–30094, 2021.
- [16] Mingliang Zhai, Xuezi Xiang, Ning Lv, and Xiangdong Kong, "Optical flow and scene flow estimation: A survey," *Pattern Recognition*, vol. 114, pp. 107861, 2021.
- [17] Michael Tao, Jiamin Bai, Pushmeet Kohli, and Sylvain Paris, "Simpleflow: A non-iterative, sublinear optical flow algorithm," *Computer Graphics Forum*, vol. 31, no. 2pt1, pp. 345–353, 2012.
- [18] Georg Petschnigg, Richard Szeliski, Maneesh Agrawala, Michael Cohen, Hugues Hoppe, and Kentaro Toyama, "Digital photography with flash and no-flash image pairs," *ACM Trans. Graph.*, vol. 23, no. 3, pp. 664–672, aug 2004.
- [19] Martin Raab and Angelika Steger, "'balls into bins' - a simple and tight analysis," in *Proc. the Second International Workshop on Randomization and Approximation Techniques in Computer Science*, Berlin, Heidelberg, 1998, p. 159–170, Springer-Verlag.
- [20] Chunxia Xiao and Jiajia Gan, "Fast image dehazing using guided joint bilateral filter," *The Visual Computer*, vol. 28, pp. 713–721, 2012.