# Pose Error Reduction for Focus Enhancement in Thermal Synthetic Aperture Visualization

Indrajit Kurmi, David C. Schedl, and Oliver Bimber, *Senior Member, IEEE*

*Abstract*—Airborne optical sectioning, an effective aerial synthetic aperture imaging technique for revealing artifacts occluded by forests, requires precise measurements of drone poses. In this article we present a new approach for reducing pose estimation errors beyond the possibilities of conventional Perspective-n-Point solutions by considering the underlying optimization as a focusing problem. We present an efficient image integration technique, which also reduces the parameter search space to achieve realistic processing times, and improves the quality of resulting synthetic integral images.

*Index Terms*—Image Processing and Computer Vision, Enhancement.

## I. Introduction

**A**IRBORNE optical sectioning (AOS) [1]–[6] is an effective aerial synthetic aperture imaging technique for revealing artifacts which are otherwise occluded through dense forest and remain invisible in individual images. For applications such as search and rescue (SAR) [2] and wild life observation [6], AOS acquires the fractional heat signal radiated from occluded targets as an unstructured thermal lightfield recorded with a camera drone (Fig. 2). Capturing over a wide synthetic aperture area (30–100 m diameter) support optical sectioning by image integration [1]. The integral images (representing the signal of a wide-aperture sensor measurement) are computed by registering and averaging all single image recordings with respect to the drone's three-dimensional poses and a defined synthetic focal plane. RGB images in visible waveband (Sony Alpha 6000 RGB camera) are acquired for drone pose estimation and thermal images (Flir Vue Pro; 7.5-13.5 $\mu m$ spectral band) for integration. Single image recordings (both thermal and RGB) are stored on the camera's internal memory cards and are downloaded for further processing after flights. After integration, targets at the focal plane will appear clear and sharp while occluders at greater distances will be strongly blurred and vanish. Thus, shifting focus by adjusting the focal plane interactively or automatically [5] towards its optimal position and orientation near the ground allows optical slicing through forest and the discovery and inspection of concealed targets, such as human bodies or animals.

AOS requires precise measurements of the drone's pose for each captured thermal image. Currently it is implemented

I. Kurmi, D. Schedl, and O. Bimber are with the Computer Science Department of the Johannes Kepler University, Linz, Austria, e-mail: firstname.lastname@jku.at
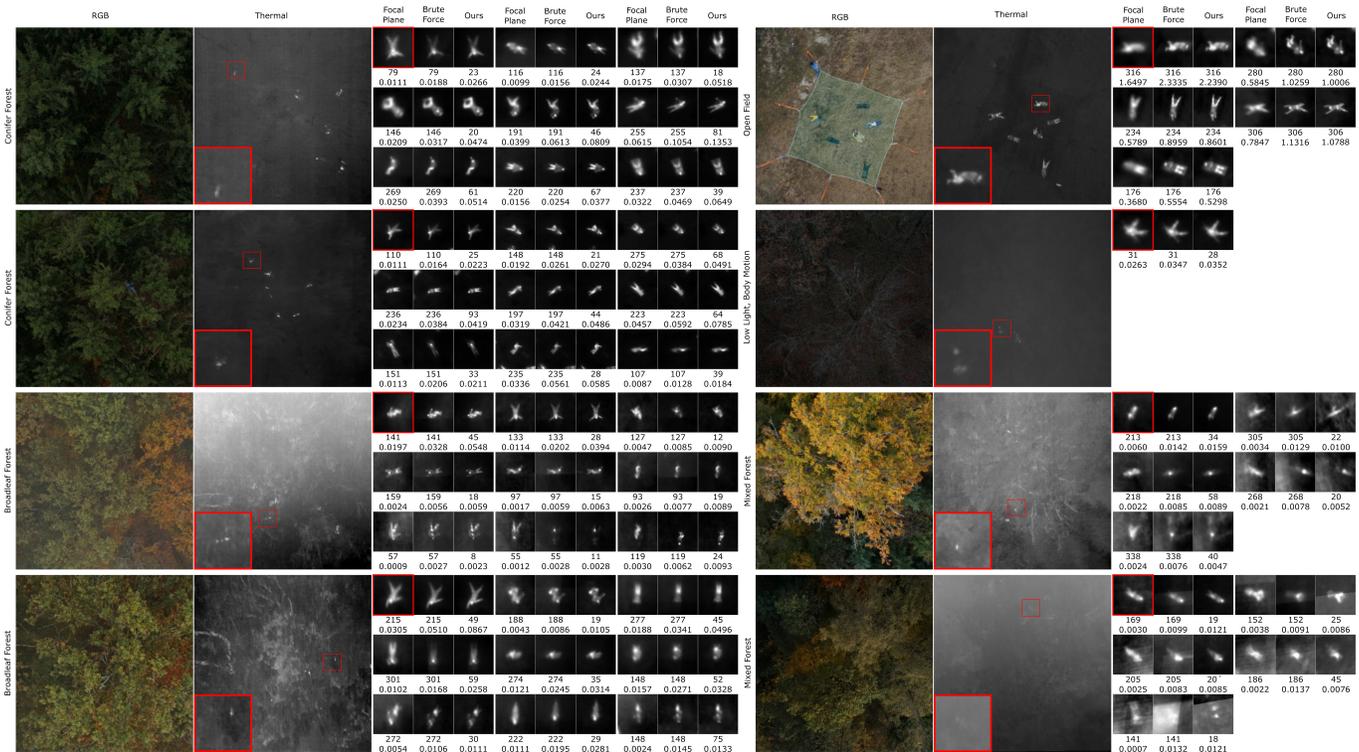
Fig. 1. Results of pose error reduction for 8 different scenes and varying recording conditions (exemplary RGB and thermal images are shown). We compare *focal plane* optimization [5], *brute force* search, and *our* approach for 52 test cases (i.e., humans). The number of integrated images ($N$) as well as the resulting normalized variance ($N \cdot \text{Var}[X]$) are displayed below each test case. Insets in the thermal images show positions and closeups for one exemplary human per scene (highlighted in red).

by applying state-of-the-art computer vision algorithms using general-purpose structure-from-motion and multi-view stereo pipeline [7], [8] to the high resolution RGB images. However, a recent study [9] has shown that primarily estimating or fitting a model to data with outliers is an NP-hard problem. Thus, pose estimation remains error prone and resulting misregistrations of single recordings leads to defocus in AOS integral images. Note that, due to their low resolution, high noise level, low contrast, and poor texture of natural scenes, thermal images do not provide usable image features for reliable pose estimation.

Various approaches for correcting the remaining pose-errors have been proposed, such as optimizations with either imposing structural constraints [10]–[12], by utilizing localized depth maps [13], [14], or by fusing different sensors [15], [16].

In this article, we present a different approach for reducing remaining pose errors by considering the underlying optimization as a focusing problem instead of solving the Perspective-n-Point (PnP) problem [9]. PnP solutions [10]–[16] require robust feature detection which is challenging in low-resolution thermal images and infeasible in the presence of strong occlusions. Thus, we solve the focusing problem by optimizing 3 or 6 pose parameters per image using simple gray level variance (GLV)[17] (which does not rely on any image features) as an error metric. We have already proven in [5], that GLV is invariant to occlusion and is therefore (in contrast to traditional gradient-, Laplacian-, and wavelet-based metrics reviewed in [18]) an ideal visibility metric for AOS thermal integral images.

Figure 1 illustrates AOS results when the synthetic focal plane is globally optimized as presented in [5] (focal plane), and the improved results when our new pose error reduction approach is applied to either 6D (brute force) or 3D (ours) pose parameters. For the latter, resulting integral images reveal more details as they are better focused and have a higher contrast.

In the following sections, we will discuss effects of errors in different pose parameters and various integration strategies leading to our final optimization approach with reduced parameter search space and early stopping. We will demonstrate that our approach improves the quality of integral images (compared to state-of-the-art pose estimation methods) while also being significantly faster and more efficient than the brute force optimization of all pose parameters.

## II. THE EFFECT OF POSE ERRORS

By applying a conventional geometric camera model, a scene point $P = (X, Y, Z, 1)^T$ in world coordinates is projected to the image point $p = (x, y)^T$ in coordinates of the camera's image plane by the projection matrix $M$ ($p = 1/zMP$), where $M = K(R\,t)$. While the camera's intrinsic matrix $K$ is constant for calibrated fixed focal length cameras, the extrinsic matrix $(R\,t)$ encodes the pose in form of a position column vector $t$ and an rotation matrix $R$ (see Fig. 2). Note that $(R\,t)$ denotes the column-wise concatenation of a matrix and vector. The coefficients of $R$ are composed of the three rotation parameters $\alpha$ (rotation around x-axis / pitch axis), $\beta$ (rotation around y-axis / roll axis), $\gamma$ (rotation around z-axis / yaw axis). To calculate $R$ we use $R =$
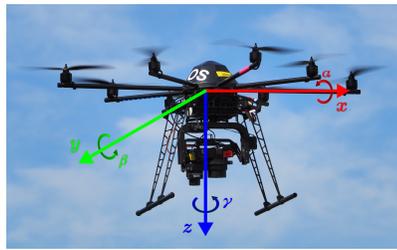


Fig. 2. Our MikroKopter OktoXL 6S12 octocopter equipped with a Flir Vue Pro thermal camera (9 mm fixed focal length lens, 14 bit tonal range covering a 7.5-13.5 $\mu m$ spectral band) and a Sony Alpha 6000 RGB camera (16-50 $mm$ lens, set to infinite focus) and its local coordinate system. The transformation to the global coordinate system is given by the extrinsic transformation parameters $(t_x, t_y, t_z, \alpha, \beta, \gamma)$.

$R_x(\alpha)R_y(\beta)R_z(\gamma)$, where $R_x, R_y, R_z$ are individual rotation matrices about axes x,y,z ($z$-$y$-$x$ extrinsic active rotation). The coefficients of $t$ are the three translation parameters $(t_x, t_y, t_z)$ along pitch, roll, and yaw axes.

Analyzing how much (on average) pose estimation error effects projection error on the image plane reveals that not all extrinsic parameters have the same impact. Figure 3 illustrate these dependencies under the assumption that the camera is located sufficiently far away from the recorded object ($t_z \gg Z$) to ignore perspective distortion. While a pose error ($\Delta t_z$) in $t_z$ leads to no significant error on the image plane, pose errors ($\Delta t_x, \Delta t_y$) in $t_x$ and $t_y$ result in notable projection errors (Fig. 3a). Rotation errors along pitch and roll axes ($\Delta\alpha$ and $\Delta\beta$) cause more projection error than a rotation error along the yaw axis $\Delta\gamma$ (Fig. 3b). By ignoring perspective distortions, the same image plane error caused by $\Delta\beta$ can also be caused by a corresponding (adjusted) $\Delta t_x$ (Fig. 3c). Thus, both errors in $\Delta\beta$ and $\Delta t_x$ can be compensated by a single correction either in $\Delta\beta$ or $\Delta t_x$. This also applies for $\Delta\alpha$ and $\Delta t_y$. The correspondence between correlating $\Delta t_x$, $\Delta t_y$ and $\Delta\alpha$, $\Delta\beta$ is unique and constant for all points. This does not apply for $\Delta\gamma$, which would require an individual $\Delta t_x, \Delta t_y$ compensation for each point (Fig. 3c). This analysis suggests, that three
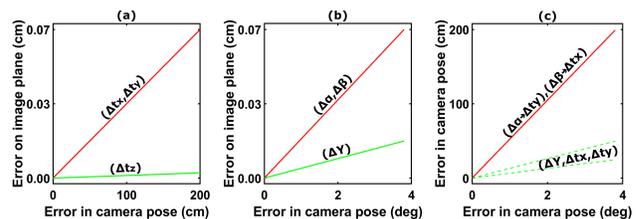


Fig. 3. Effect of pose estimation error on image plane error. In this simulation, we assume that $t_z \gg Z$ and that $t_z = 30$m. (a) Image plane error for $\Delta t_z$ (green) and for $\Delta t_x, \Delta t_y$ (red). (b) Image plane error for $\Delta\gamma$ (green) and for $\Delta\alpha, \Delta\beta$ (red). (c) Same image plane error caused by $\Delta t_x, \Delta t_y$ or by corresponding $\Delta\alpha, \Delta\beta$ (red). Image plane error of two different points caused by $\Delta t_x, \Delta t_y$ show that no constant corresponding $\Delta\gamma$ exists for all points (dashed green).

instead of six extrinsic parameters are sufficient for pose error reduction in our case. Since $\Delta t_z$ has little effect and $\Delta\alpha, \Delta\beta$ can be compensated, only $\Delta t_x$, $\Delta t_y$, $\Delta\gamma$ need to be estimated. Less parameters for pose error estimation leads to significantly faster and more efficient optimizations.

## III. EFFICIENT IMAGE INTEGRATION

In [5] we show that the GLV [17] of an AOS thermal integral image $X$ is an occlusion invariant measure for determining the
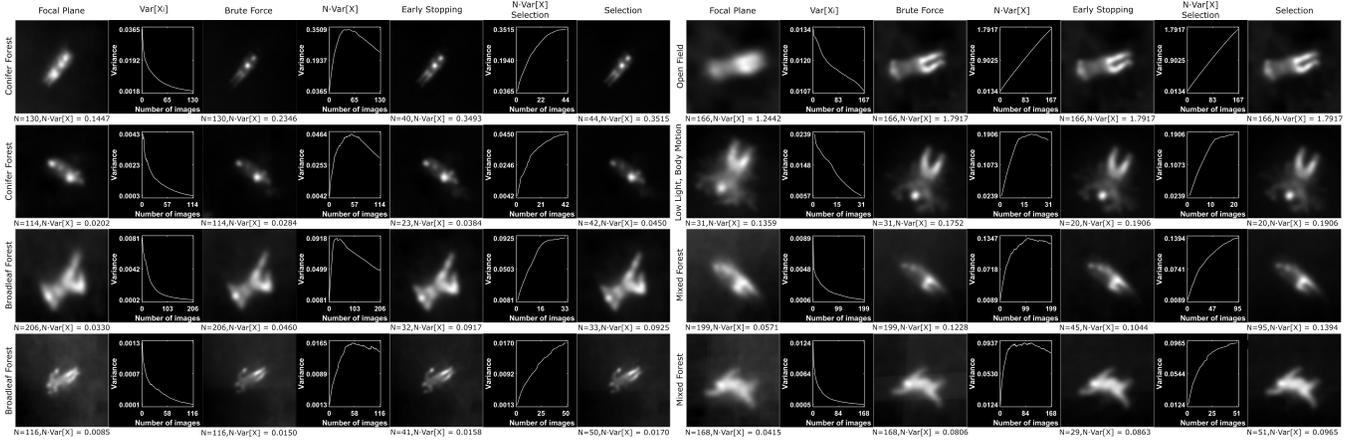
Fig. 4.  Integration strategies for a set of representative cases (different scene types shown in Fig. 1). From left to right: Results of *focal plane optimization* only (without additional pose error reduction), $\text{Var}[X_i]$-plot for all individual images in variance-decreasing order, results of *brute force* integration of all images in variance decreasing order, $N \cdot \text{Var}[X]$-plot for integral images computed from $N$ individual images in variance-decreasing order, results of *early stopping* the integration (in variance-decreasing order) at the maximum of $N \cdot \text{Var}[X]$, $N \cdot \text{Var}[X]$-plot for images that lead to an increase of $N \cdot \text{Var}[X]$ after integration in variance-decreasing order (*selection*), results of the *selection* strategy.

optimal synthetic focal plane pose, which minimizes occlusion and defocus. We prove (see Appendix) that the variance of an integral image is proportional to

$$\text{Var}[X] = \frac{\text{Var}[X_i]}{N} + (1-D)^2\left(1 - \frac{1}{N}\right)\sigma_s^2, \qquad (1)$$

where the variance of a single image is

$$\text{Var}[X_i] = D(1-D)((\mu_o - \mu_s)^2) + D\sigma_o^2 + (1-D)\sigma_s^2, \quad (2)$$

and $D$ is the probability of occlusion, while $\mu_o$, $\sigma_o^2$ and $\mu_s$, $\sigma_s^2$ are statistical properties of occlusion and the signal respectively. Note that the variance of a single image $\text{Var}[X_i]$ is only dependening on its content (signal and occlusion). Individual pose parameters $(\alpha, \beta, \gamma, t_x, t_y, t_z)$ do not influence $\text{Var}[X_i]$ but indirectly influence the focus/defocus of the integral image $X$. Thus, integrating $N$ individual images that are miss-registered due to remaining errors in pose estimation will result in defocus and thus change the variance of the integral image $\text{Var}(X)$—even if the optimal focal plane settings are found. We propose to optimize the initial pose estimation results of each image $X_i$ by maximizing $\text{Var}[X]$ while searching for better pose parameters $(\alpha, \beta, \gamma, t_x, t_y, t_z)_i$, where $X$ is the integral image of all $X_i$ $(i = 1 \ldots N)$. First, we determine the optimal synthetic focal plane, as explained in [5]. As explained earlier, we use general-purpose structure-from-motion and multi-view stereo pipeline [7], [8] applied to the high resolution RGB images for initial pose estimation. Second, we sort all thermal images $X_i$ in decreasing $\text{Var}[X_i]$ order, since higher variance indicates less occlusion. Finally, we sequentially register and integrate all images $X_i$ in decreasing variance order to an integral image $X$ by maximizing $\text{Var}[X]$ while searching for optimal pose parameters $(\alpha, \beta, \gamma, t_x, t_y, t_z)_i$ which register $X_i$ to the previous integral image (i.e., the one that integrates $X_1 \ldots X_{(i-1)}$ – starting initially with $X = X_1$. This way, we have to optimize only for six parameters per integration step, rather than for $6N$ parameters at once. We chose Nelder-Mead [19] (a commonly used direct search approach with complexity of $\mathcal{O}(n)$ per iteration where $n$ is the number of parameters) for parameter optimization. Note, that when integrating $N$ images, the

contrast of the resulting integral image drops proportionally to $N$ (due to the present occlusion on each individual image)[3]. Thus, $N \cdot \text{Var}[X]$ normalizes the resulting variance to become invariant to $N$.

As shown in Fig. 4, this clearly improves the focus of the integral images compared to a focal plane optimization only. However, such a brute force registration of all $N$ images by optimizing all six parameters is still time consuming and scales linearly with $N$, even if we search only $N$ times within a 6-dimensional parameter space rather than once within a $6N$-dimensional parameter space. Furthermore, focus can worsen if too many images are integrated in case of remaining optimization errors.

To overcome these problems, we investigated two alternatives to the above *brute force* strategy: First, an *early stopping* strategy stops the integration process as soon as $N \cdot \text{Var}[X]$ has reached its maximum and starts to decrease. Second, an *selection* strategy registers all images but integrates only those that lead to an increase in $N \cdot \text{Var}[X]$. Results of these three integration strategies are illustrated in Fig. 4.

Since early stopping is fastest and does not show a significant degradation in focus for a set of representative test cases, we chose this as a final integration strategy.

## IV. SEARCH SPACE REDUCTION

As discussed in section II, we could further reduce the number of pose parameters to be optimized for each image from six $(t_x, t_y, t_z, \alpha, \beta, \gamma)$ to three $(t_x, t_y, \gamma)$. Here, we evaluate this option for our representative test cases while applying early stopping for integration.

As illustrated in Fig. 5, there is no considerable degradation in focus even if the number of pose parameters is reduced to two $(t_x, t_y)$. In section II we explain that errors in $\alpha, \beta$ can be compensated by corrections in $t_x, t_y$, and that $t_z$ has no significant impact in case of $t_z \gg Z$. The reason why corrections of $\gamma$ seems also to have little impact might be that the remaining error resulting from the initial pose estimation was little for our test cases. However, since this cannot be assumed in general, we chose the set of three pose $(t_x, t_y, \gamma)$ as the best trade-off between quality and performance.
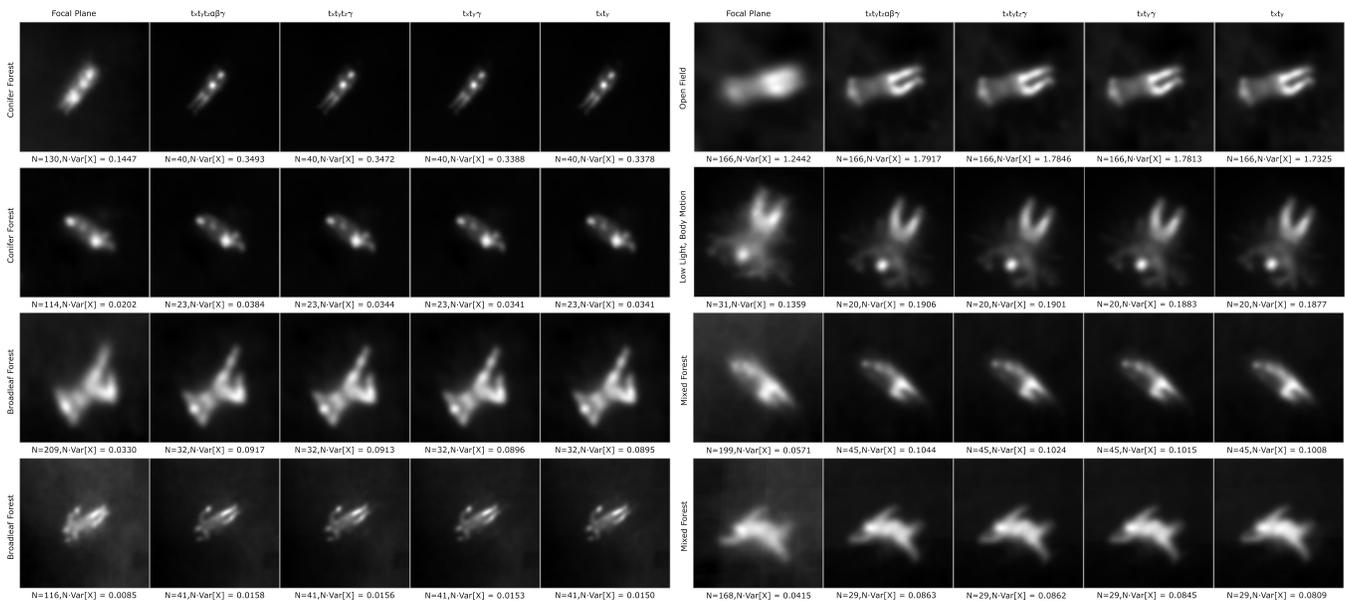
Fig. 5. Search space reduction by decreasing the number of pose parameters for a set of representative cases (same as in Fig. 4). From left to right: Results of *focal plane optimization* only (without additional pose error reduction), *early stopping* integration applied to all six parameters $(t_x, t_y, t_z, \alpha, \beta, \gamma)$, applied to four parameters $(t_x, t_y, t_z, \gamma)$, applied to three parameters $(t_x, t_y, \gamma)$, and applied to two parameters $(t_x, t_y)$.

## V. RESULTS

Figure 1 presents final results of our pose error reduction approach being applied to a total of 52 test cases captured under different conditions (conifer forest, broadleaf forest, mixed forest, open field, low-light and sunlight conditions, moving and still subjects, different densities). It achieves significant qualitative improvements under open field conditions (without occlusion), sparse and moderate dense occlusion conditions, and even under low-light conditions with body motion. Under conditions with severe occlusion, however, only minor improvements are made.

Our optimization framework (including search space reduction and early stopping) has, on an average, $89.92\%$ less parameters to optimize when compared to brute force search (only $3N$ instead of $6N$ parameters, where on average $N$ is lower due to early stopping). It significantly improves the focus of the integral images by $204.77\%$ (according to the measured $N \cdot \text{Var}[X]$ gain) when compared to only optimizing the focal plane as in [5].

When executed on a Nvidia GeForce GTX 1060, our GPU implementation requires on average 355s for the optimization of AOS integrals computed form 50 single images (before early stopping).

## VI. CONCLUSION AND FUTURE WORK

In this article we utilize GLV as an optimization metric for pose error reduction in thermal AOS integral images. We study effects of errors in different pose parameters and demonstrate that the number of pose parameters to be optimized can be reduced from six to three without any significant degradation of quality. We also investigate different image integration strategies to further reduce the search space, and showed that the contribution of a single image in the total integral is proportional to its GLV. Thus, single images are best integrated in decreasing GLV order – which supports early stopping.

Our results demonstrate that we achieved qualitative and quantitative improvements for representative test cases. The reported computation time of our implementation (on average 355s) can be reduced by an efficient hierarchical multi-scale processing approach (with a first test implementation we already achieved an average of 197s for each test case). Furthermore, advanced high-end GPUs (e.g., Nvidia Tesla V100 or Quadro RTX 8000) will lead to ×3-4 speed-ups when compared to the applied Nvidia GeForce GTX 1060. Overall, we estimate that processing times of less than 50s will be realistic and applicable as post-processing image-enhancement step. In comparison, conventional PnP solutions, such as in [7], [8], required on average 24 minutes for pose estimations with the full set of around 300 images.

In the future, we want to investigate if our pose-error-reduction approach (based on thermal images) can replace the conventional pose-estimation process (based on RGB images). The latter is not applicable at night as it relies on sufficient visible light. If the drone's GPS recordings are used as initial poses and our approach reduces the remaining pose errors sufficiently, then AOS night operations are possible.

We are currently also working on a people classification deep neural network that operates on unoptimized AOS integral images. This will allow the fully automatic selection of people regions to be optimized (which are manually selected at the moment).

## APPENDIX

Here, we present the derivation of the relation between variance $(\text{Var}[X])$ of the integral image $X$ and variance $(\text{Var}[X_i])$ of the individual image recordings $X_i$. We rely on the statistical model of [5], where the integral image $X$ is composed of $N$ single image recordings $X_i$ and each single image pixel in $X_i$ is either occlusion free $(S)$ or occluded $(O)$ determined by $Z$:

$$X_i = Z_i O_i + (1 - Z_i)S. \tag{3}$$

Similar to [5], all variables are independent and identically distributed with $Z_i$ following a Bernoulli distribution with success parameter $D$ (i.e., $\mathrm{E}[Z_i] = \mathrm{E}[Z_i^2] = D$; furthermore, note that $\mathrm{E}[Z_i(1 - Z_i)] = 0$ is true). The random variable $S$ follows a distribution with mean $\mathrm{E}[S] = \mu_s$ and $\mathrm{E}[S^2] = (\mu_s^2 + \sigma_s^2)$ and analogously $O_i$ follows a distribution with mean $\mathrm{E}[O_i] = \mu_o$ and $\mathrm{E}[O_i^2] = (\mu_o^2 + \sigma_o^2)$. For the mean and the variance of $X_i$, we determine the first and second moments of $X_i$:

$$\mathrm{E}[X_i] = D\mu_o + (1-D)\mu_s \qquad (4)$$

and

$$\mathrm{E}[X_i^2] = D(\mu_o^2 + \sigma_o^2) + (1-D)(\mu_s^2 + \sigma_s^2). \qquad (5)$$

Variance of single image recordings $X_i$ can be obtained as:

$$
\begin{aligned}
\mathrm{Var}[X_i] &= \mathrm{E}[X_i^2] - (\mathrm{E}[X_i])^2 \\
&= D(\mu_o^2 + \sigma_o^2) + (1-D)(\mu_s^2 + \sigma_s^2) - (D^2\mu_o^2 + (1-D)^2\mu_s^2 + 2D(1-D)\mu_o\mu_s) \\
&= D(1-D)((\mu_o - \mu_s)^2) + D\sigma_o^2 + (1-D)\sigma_s^2.
\end{aligned}
\qquad (6)
$$

Similarly, for $X$, we determine the first and second moments where the first moment of $X$ is given by:

$$
\begin{aligned}
\mathrm{E}[X] &= \mathrm{E}\left[\frac{1}{N}\sum_{i=1}^{N} Z_i O_i + (1-Z_i)S\right] \\
&= D\mu_o + (1-D)\mu_s
\end{aligned}
\qquad (7)
$$

and the second moment of $X$ is as derived in [5]:

$$
\begin{aligned}
\mathrm{E}[X^2] = \frac{1}{N^2}\Big( &N\big(D(\sigma_o^2 + \mu_o^2) + (1-D)(\sigma_s^2 + \mu_s^2)\big) \\
&+ N(N-1)\big(D^2\mu_o^2 + 2D(1-D)\mu_s\mu_o + (1-D)^2(\sigma_s^2 + \mu_s^2)\big)\Big).
\end{aligned}
\qquad (8)
$$

Consecutively, we calculate variance of the integral image as:

$$
\begin{aligned}
\mathrm{Var}[X] &= \mathrm{E}[X^2] - (\mathrm{E}[X])^2 \\
&= \frac{1}{N}\big(D(\sigma_o^2 + \mu_o^2) + (1-D)(\sigma_s^2 + \mu_s^2)\big) \\
&\quad + \big(D^2\mu_o^2 + 2D(1-D)\mu_s\mu_o + (1-D)^2(\sigma_s^2 + \mu_s^2)\big) \\
&\quad - \frac{1}{N}\big(D^2\mu_o^2 + 2D(1-D)\mu_s\mu_o + (1-D)^2(\sigma_s^2 + \mu_s^2)\big) \\
&\quad - \big(D^2\mu_o^2 + (1-D)^2\mu_s^2 + 2D(1-D)\mu_o\mu_s\big) \\
&= \frac{1}{N}\big(D(\sigma_o^2 + \mu_o^2) + (1-D)(\sigma_s^2 + \mu_s^2)\big) \\
&\quad + (1-D)^2\sigma_s^2 \\
&\quad - \frac{1}{N}\big(D^2\mu_o^2 + 2D(1-D)\mu_s\mu_o + (1-D)^2(\sigma_s^2 + \mu_s^2)\big) \\
&= \frac{1}{N}\big(D(1-D)((\mu_o - \mu_s)^2) + D\sigma_o^2 + (1-D)\sigma_s^2\big) \\
&\quad + (1-D)^2\big(1 - \frac{1}{N}\big)\sigma_s^2.
\end{aligned}
\qquad (9)
$$

Substituting (6) in (9) yields (1).

## REFERENCES

[1] I. Kurmi, D. C. Schedl, and O. Bimber, "Airborne optical sectioning," *Journal of Imaging*, vol. 4, no. 8, 2018. DOI: 10.3390/jimaging4080102.

[2] I. Kurmi, D. C. Schedl, and O. Bimber, "Thermal airborne optical sectioning," *Remote Sensing*, vol. 11, no. 14, 2019. DOI: 10.3390/rs11141668.

[3] I. Kurmi, D. C. Schedl, and O. Bimber, "A statistical view on synthetic aperture imaging for occlusion removal," *IEEE Sensors Journal*, pp. 1–1, 2019.

[4] O. Bimber, I. Kurmi, D. C. Schedl, and M. Potel, "Synthetic aperture imaging with drones," *IEEE Computer Graphics and Applications*, vol. 39, pp. 8–15, 2019.

[5] I. Kurmi, D. C. Schedl, and O. Bimber, "Fast Automatic Visibility Optimization for Thermal Synthetic Aperture Visualization," *IEEE Geoscience and Remote Sensing Letters*, vol. to appear, 2020.

[6] D. C. Schedl, I. Kurmi, and O. Bimber, "Airborne optical sectioning for nesting observation," *Scientific Reports*, vol. 10, p. 7254, 2020.

[7] J. L. Schönberger and J. Frahm, "Structure-from-motion revisited," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016, pp. 4104–4113. DOI: 10.1109/CVPR.2016.445.

[8] J. L. Schönberger, E. Zheng, M. Pollefeys, and J. Frahm, "Pixelwise view selection for unstructured multi-view stereo," in *European Conference on Computer Vision (ECCV)*, 2016.

[9] T. J. Chin, Z. Cai, and F. Neumann, "Robust fitting in computer vision: Easy or hard?" In *Computer Vision – ECCV 2018*, Cham: Springer International Publishing, 2018, pp. 715–730, ISBN: 978-3-030-01258-8.

[10] H. Li, J. Yao, J. Bazin, X. Lu, Y. Xing, and K. Liu, "A monocular slam system leveraging structural regularity in manhattan world," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 2518–2525.

[11] V. Ovechkin and V. Indelman, "Bafs: Bundle adjustment with feature scale constraints for enhanced estimation accuracy," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 804–810, 2018.

[12] M. Hsiao, E. Westman, and M. Kaess, "Dense planar-inertial slam with structural constraints," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 6521–6528.

[13] Y. Shin, Y. S. Park, and A. Kim, "Direct visual slam using sparse depth for camera-lidar system," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 5144–5151.

[14] E. J. Shamwell, K. Lindgren, S. Leung, and W. D. Nothwang, "Unsupervised deep visual-inertial odometry with online error correction for rgb-d imagery," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2019.

[15] R. Mur-Artal and J. D. Tardós, "Visual-inertial monocular slam with map reuse," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 796–803, 2017.

[16] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.

[17] L. Firestone, K. Cook, K. Culp, N. Talsania, and K. Preston, "Comparison of autofocus methods for automated microscopy," *Cytometry*, vol. 12, pp. 195–206, 1991.

[18] S. Pertuz, D. Puig, and M. A. Garcia, "Analysis of focus measure operators for shape-from-focus," *Pattern Recognition*, vol. 46, no. 5, pp. 1415–1432, 2013.

[19] J. Nelder and R. Mead, "A Simplex Method for Function Minimization," *Comput. J.*, vol. 7, pp. 308–313, 1965.