

D2NT: A High-Performing Depth-to-Normal Translator

Yi Feng, Bohuan Xue, Ming Liu, Qijun Chen, and Rui Fan[✉]

Abstract—Surface normal holds significant importance in visual environmental perception, serving as a source of rich geometric information. However, the state-of-the-art (SoTA) surface normal estimators (SNEs) generally suffer from an unsatisfactory trade-off between efficiency and accuracy. To resolve this dilemma, this paper first presents a superfast depth-to-normal translator (D2NT), which can directly translate depth images into surface normal maps without calculating 3D coordinates. We then propose a discontinuity-aware gradient (DAG) filter, which adaptively generates gradient convolution kernels to improve depth gradient estimation. Finally, we propose a surface normal refinement module that can easily be integrated into any depth-to-normal SNEs, substantially improving the surface normal estimation accuracy. Our proposed algorithm demonstrates the best accuracy among all other existing real-time SNEs and achieves the SoTA trade-off between efficiency and accuracy.

SOURCE CODE, DEMO VIDEO, & SUPPLEMENT

Our source code, demo video, and supplement are publicly available at mias.group/D2NT.

I. INTRODUCTION

Surface normal is an informative visual feature that has been widely used in a variety of robot environmental perception tasks, *e.g.*, visual odometry [1], [2], scene parsing [3]–[7], and depth estimation [8], [9]. Due to the requirement for real-time execution in such tasks, surface normal estimators (SNEs) should be both accurate and computationally efficient [10].

Early geometry-based SNEs compute surface normals via either plane fitting (solvable with energy minimization techniques) or weighted neighboring normal aggregation. However, these SNEs typically have an imbalance between accuracy and speed (see Fig. 1). In 2015, Nakagawa *et al.* [11] proposed an efficient SNE, which computes surface

This work was supported by the National Key R&D Program of China under Grant 2020AAA0108100, the National Natural Science Foundation of China under Grant 62233013, the Science and Technology Commission of Shanghai Municipal under Grant 22511104500, the Fundamental Research Funds for the Central Universities under Grants 22120220184 and 22120220214, and the Shanghai Municipal Science and Technology Major Project under Grant 2021SHZDZX0100. (Yi Feng and Bohuan Xue contributed equally to this work.) (Corresponding author: Rui Fan.)

Yi Feng, Qijun Chen, and Rui Fan are with the Robotics & Artificial Intelligence Laboratory (RAIL), the College of Electronic & Information Engineering, the State Key Laboratory of Intelligent Autonomous Systems, Tongji University, Shanghai 201804, P. R. China. (e-mails: fengyi@ieee.org, qjchen@tongji.edu.cn, rui.fan@ieee.org)

Bohuan Xue is with the Department of Computer Science & Engineering, the Hong Kong University of Science and Technology, Hong Kong SAR, P. R. China. (e-mail: bxueaa@ust.hk)

Ming Liu is with the Robotics & Autonomous Systems Thrust of the Systems Hub, the Hong Kong University of Science and Technology (Guangzhou), Nansha, Guangzhou 511400, P. R. China. (e-mail: eelium@ust.hk)

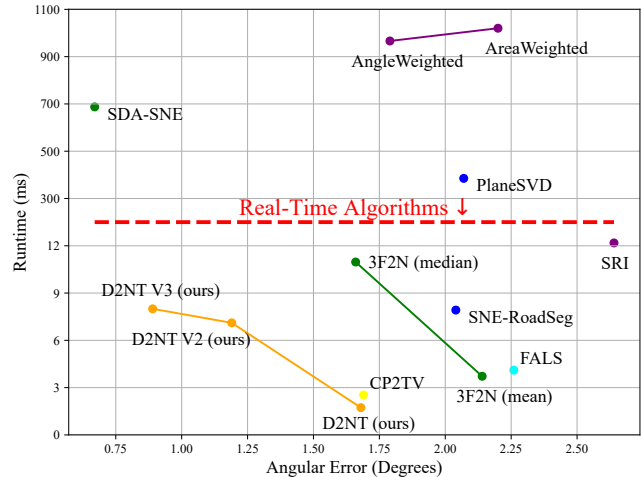


Fig. 1. Efficiency versus accuracy trade-off comparison among all SoTA geometry-based SNEs (on the 3F2N Easy dataset). D2NT has the highest computational efficiency, and D2NT V3 achieves the best trade-off between speed and accuracy.

normals via the cross product of two orthogonal tangent vectors (hereafter called CP2TV). However, its performance on spatial discontinuities is unsatisfactory as a result of inaccurate observed tangent vectors. Recently, Fan *et al.* [10] introduced an efficient and accurate SNE, referred to as three-filters-to-normal (3F2N). Although 3F2N achieves state-of-the-art (SoTA) performance, the aggregation of neighboring surface normals with a mean or median filter is still computationally intensive.

Therefore, there is a strong necessity to develop an SNE that achieves a balance between rapid computation and high accuracy. In this paper, we present a high-performing depth-to-normal translator (D2NT), which significantly improves the efficiency and accuracy trade-off, and significantly refines the estimation results in and around discontinuities. The contributions of our work are summarized as follows:

- 1) **D2NT**, a cutting-edge **Depth-to-Normal Translator**. In comparison to other geometry-based SNEs, D2NT computes surface normals directly from depth maps, demonstrating remarkable computational efficiency. Compared to existing SoTA SNEs, D2NT establishes the most direct relationship between depth and surface normal.
- 2) **Discontinuity-Aware Gradient (DAG) filter**, a depth gradient filter that selectively identifies discontinuities and eliminates outliers (non-coplanar points in relation to the reference point). Compared to traditional finite difference (FD) operators, our proposed DAG filter provides a significant improvement in terms of depth gradient estimation accuracy.

3) **Markov random field-based Normal Refinement (MNR) module**, which dramatically reduces surface normal estimation errors. It can also be integrated with any depth-to-normal SNEs to further enhance the quality of their estimated surface normals.

II. RELATED WORK

This section provides an overview of geometry-based surface normal estimators. As shown in Table I, the existing SNEs can be divided into three categories: energy minimization-based, averaging-based, and depth-to-normal.

Let $\mathbf{P} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n\}$ be the given 3D point set. For an arbitrary 3D point $\mathbf{p}_i \in \mathbf{P}$, its surface normal is represented as $\mathbf{n}_i = [n_{ix}, n_{iy}, n_{iz}]^\top$. To find the optimal \mathbf{n}_i , $\mathbf{Q}_i = \{\mathbf{q}_{i1}, \mathbf{q}_{i2}, \dots, \mathbf{q}_{ik} \mid \mathbf{q}_{ik} \in \mathbf{P}\}$, the neighboring points of \mathbf{p}_i are typically considered.

A. Energy Minimization-Based Methods

This category of methods computes surface normals by finding a best-fit plane from the augmented neighboring point set $\mathbf{Q}_i^+ = \{\mathbf{Q}_i, \mathbf{p}_i\}$ as follows:

$$\hat{\mathbf{n}}_i = \arg \min_{\mathbf{n}_i} E(\mathbf{Q}_i^+, \mathbf{n}_i), \quad (1)$$

where $\hat{\mathbf{n}}_i$ is obtained by minimizing the energy function E .

PlaneSVD [12] fits a local planar surface to \mathbf{Q}_i^+ by minimizing the least squares of the distances from the points to the surface using SVD. Similarly, PlanePCA [13] finds the minimum variance of \mathbf{Q}_i^+ with respect to the centroid $\bar{\mathbf{p}}_i = \frac{1}{k+1}(\mathbf{p}_i + \sum_{j=1}^k \mathbf{q}_{ij})$. VectorSVD [14] fits the local planar surface by minimizing the sum of the squared dot products between the surface normal and tangent vectors.

Recently, Ming *et al.* [15] proposed SDA-SNE, a highly accurate surface normal estimator based on multi-directional dynamic programming and iterative polynomial interpolation. Nevertheless, its demanding computational requirements and iterative nature result in subpar real-time performance. The computation-intensive nature of energy minimization and the calculation of 3D coordinates make these SNEs suffer from slow processing speed and noise.

B. Averaging-Based Methods

This category of methods estimates surface normals by averaging the normal vectors of the surrounding triangles:

$$\mathbf{n}_i = \frac{1}{k} \sum_{j=1}^k w_j \frac{\mathbf{r}_{ij} \times \mathbf{r}_{ij+1}}{\|\mathbf{r}_{ij} \times \mathbf{r}_{ij+1}\|_2}, \quad (2)$$

where $\mathbf{r}_{ij} = \mathbf{q}_{ij} - \mathbf{p}_i$, $\mathbf{r}_{ik+1} = \mathbf{r}_{i1}$, and w_j is the weight calculated based on either the area (AreaWeighted [16]) or the angle (AngleWeighted [17]) of the triangles. Nonetheless, both of these methods necessitate an initial estimation of the normals and can only be utilized as a back-end optimization technique.

TABLE I
TAXONOMY OF THE SOTA GEOMETRY-BASED SNEs.

| Category | Algorithm | Expression |
|---------------------------|--------------------|---|
| Energy Minimization-Based | PlaneSVD [12] | $\min \left\ \left[\mathbf{Q}_i^+, \mathbf{1}_k \right] \mathbf{b}_i \right\ _2$ |
| | PlanePCA [13] | $\min \left\ \left[\mathbf{Q}_i^+ - \bar{\mathbf{p}} \right] \mathbf{n}_i \right\ _2$ |
| | VectorSVD [14] | $\min \left\ \left[\mathbf{Q}_i - \mathbf{1}_k \mathbf{p}_i^\top \right] \mathbf{n}_i \right\ _2$ |
| | SDA-SNE [15] | $\min \left\{ \mathcal{T} \left(\mathbf{E}^{(k-1)}, \mathbf{S} \right) \right\}$ |
| Averaging-Based | AreaWeighted [16] | $w_j = \frac{1}{2} \left\ \mathbf{r}_{ij} \times \mathbf{r}_{ij+1} \right\ _2$ |
| | AngleWeighted [17] | $w_j = \cos^{-1} \left(\frac{\langle \mathbf{r}_{ij}, \mathbf{r}_{ij+1} \rangle}{\ \mathbf{r}_{ij}\ _2 \ \mathbf{r}_{ij+1}\ _2} \right)$ |
| Depth-to-Normal | 3F2N [10] | $n_x = f_x \frac{\partial 1/z}{\partial u}, \quad n_y = f_y \frac{\partial 1/z}{\partial v}$ |
| | | $\hat{n}_z = -\Phi \left\{ \frac{\Delta x_{ij} n_x + \Delta y_{ij} n_y}{\Delta z_{ij}} \right\}$ |
| | CP2TV [11] | $\mathbf{n}_i = \mathbf{t}_u \times \mathbf{t}_v$ |

C. Depth-to-Normal Methods

Fan *et al.* [10] proposed 3F2N, a fast and accurate surface normal estimator, which directly converts the structured range sensor data, such as depth or disparity images, into surface normal maps using two gradient filters and a mean/median filter. This category of methods typically assumes that the range sensor is a pinhole camera model as follows:

$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{K} \mathbf{p}_i = \begin{bmatrix} f_x & 0 & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix}, \quad (3)$$

where \mathbf{K} represents the camera intrinsic matrix, $\mathbf{p}_0 = [u_0, v_0]^\top$ is the principal point in pixels, and f_x and f_y denote the camera's focal lengths in the x and y directions, respectively. This method achieves fast computational speed and high accuracy, but it still involves the calculation of 3D coordinates, which is redundant and computationally demanding.

Nakagawa *et al.* [11] presented CP2TV, an SNE that utilizes the cross-products of tangent vectors of local planar surfaces to directly estimate surface normals from depth maps. However, the accuracy of this method is inadequate in and around discontinuities, as it adopts a finite difference operator to estimate depth gradients. Inaccurate tangent vectors generated in these regions lead to substantial calculation errors in the estimated surface normals.

III. METHODOLOGY

In this section, we first introduce a highly efficient method for estimating surface normals from structured range sensor data in an end-to-end manner. Then, we present a novel approach to improve the accuracy of depth gradient estimation. Additionally, we propose an optimization strategy to refine surface normal estimation in and around discontinuities, which can be well embedded into any existing depth-to-normal SNEs. The pipeline of our algorithm is illustrated in Fig. 2.

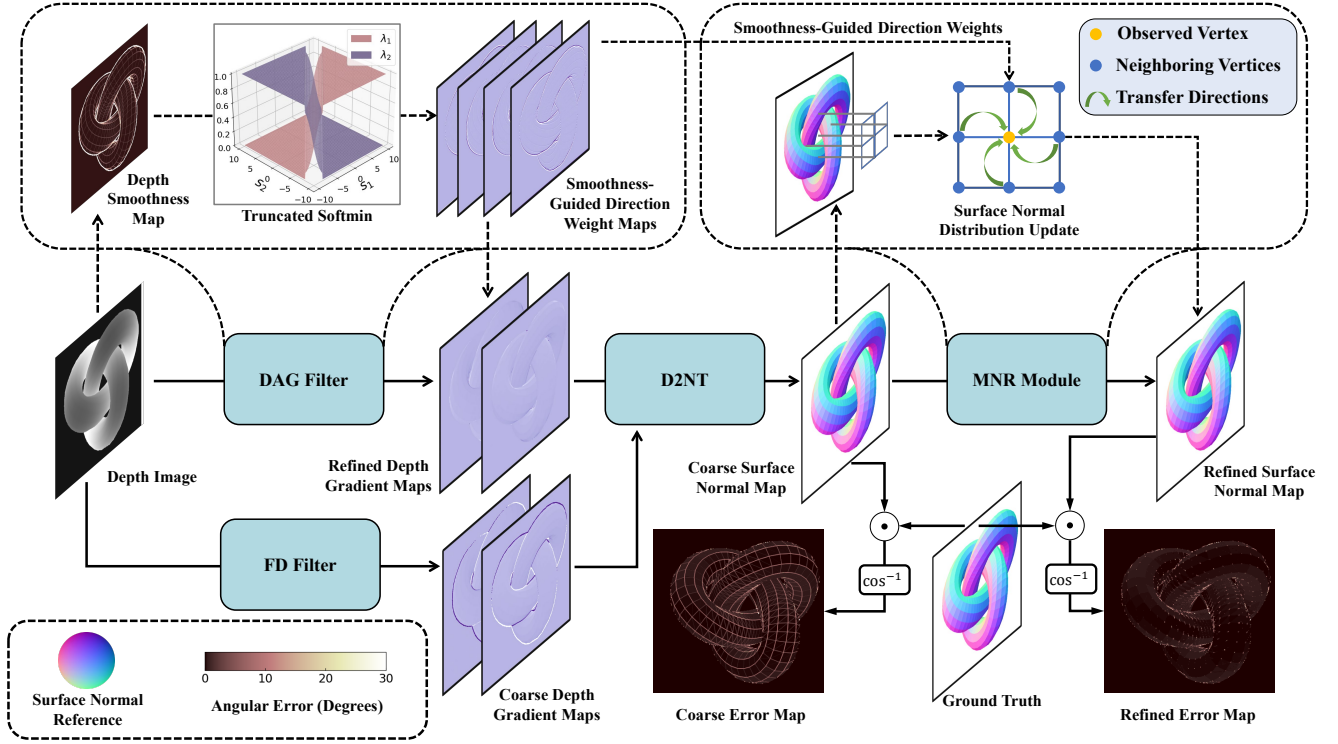


Fig. 2. The illustration of our proposed D2NT, DAG filter, and MNR module. D2NT translates depth images into surface normal maps in an end-to-end fashion; DAG filter adaptively generates smoothness-guided direction weights for improved depth gradient estimation in and around discontinuities; MNR module further refines the estimated surface normals based on the smoothness of neighboring pixels.

A. Depth-to-Normal Translation

An observed 3D point $\mathbf{p} = [x, y, z]^\top$ and its surface normal $\mathbf{n} = [n_x, n_y, n_z]^\top$ have the following relation:

$$n_x x + n_y y + n_z z + d = 0, \quad (4)$$

where d is the distance between the origin and the tangent plane. Combining (3) with (4) results in the following expression:

$$\frac{z(u - u_0)}{f_x} n_x + \frac{z(v - v_0)}{f_y} n_y + n_z z + d = 0. \quad (5)$$

(5) contains an implicit function $z(u, v)$. We compute the partial derivatives of z with respect to u and v , as follows:

$$\begin{aligned} z_u \left(\frac{u - u_0}{f_x} n_x + \frac{v - v_0}{f_y} n_y + n_z \right) + \frac{z}{f_x} n_x &= 0, \\ z_v \left(\frac{u - u_0}{f_x} n_x + \frac{v - v_0}{f_y} n_y + n_z \right) + \frac{z}{f_y} n_y &= 0, \end{aligned} \quad (6)$$

where $z_u = \frac{\partial z}{\partial u}$ and $z_v = \frac{\partial z}{\partial v}$. n_x and n_y can then be obtained by plugging (5) into (6):

$$n_x = \frac{f_x d}{z^2} z_u, \quad n_y = \frac{f_y d}{z^2} z_v. \quad (7)$$

n_z can therefore be computed by plugging (7) into (5):

$$n_z = -\frac{d}{z^2} (z + (u - u_0) z_u + (v - v_0) z_v). \quad (8)$$

Removing the common factor $\frac{d}{z^2}$ results in a simplified expression for surface normal:

$$\mathbf{n} = \begin{bmatrix} -f_x & 0 & 0 \\ 0 & -f_y & 0 \\ u - u_0 & v - v_0 & z \end{bmatrix} \begin{bmatrix} z_u \\ z_v \\ 1 \end{bmatrix}. \quad (9)$$

(9) describes an end-to-end translation from a given depth image to its surface normal map. Compared with other SoTA SNEs, such as 3F2N [10] and SNE-RoadSeg [3], D2NT eliminates the need to calculate 3D coordinates by leveraging the explicit relationship between depth and normal, demonstrating remarkable computational efficiency.

B. Discontinuity-Aware Gradient Filtering

As (9) demonstrates, surface normals can be directly calculated from structured range sensor data when the camera parameters are known. The accuracy of the partial derivatives directly affects the accuracy of the surface normal estimation. This subsection introduces an improved depth gradient computation approach.

The existing depth-to-normal methods generally utilize regular image gradient filters, such as FD¹, to approximate depth gradients. However, these filters tend to yield poor results on discontinuities, such as ridges, ditches, and edges, as outliers (non-coplanar adjacent points) are involved in depth difference computation.

To address this issue, we define a horizontal gradient filter G_h and a vertical gradient filter G_v as follows:

$$G_h = \lambda_l \Delta_b + \lambda_r \Delta_f, \quad G_v = \lambda_u \Delta_b^\top + \lambda_d \Delta_f^\top, \quad (10)$$

¹Horizontal FD kernel: $\Delta = [-1, 0, 1]$.

where $\Delta_b = [-1, 1, 0]$ and $\Delta_f = [0, -1, 1]$ are the backward and forward difference operators, respectively. λ denotes the weight distribution along four different directions. Hereafter the subscripts l , r , u , and d denote left, right, up, and down directions, respectively.

To obtain more accurate G_u and G_v in areas with discontinuities, we must distinguish between distinct continuous surfaces and assign appropriate weights to Δ_f and Δ_b based on the smoothness of the neighboring pixels' surfaces. The local surface smoothness s_p can be reflected by:

$$s_p = |\nabla^2 z_p|, \quad (11)$$

where ∇^2 is a second-order Discrete Laplace Filter (DLF), and z_p is the depth at p . The weights of the difference operator along four directions can then be assigned as follows:

$$\lambda_l, \lambda_r = \mathcal{M}(s_l, s_r), \quad \lambda_u, \lambda_d = \mathcal{M}(s_u, s_d), \quad (12)$$

where s_l , s_r , s_u , and s_d represent the smoothness of four neighboring points, respectively, and \mathcal{M} is the softmin function:

$$\mathcal{M}(s_i) = \frac{e^{-s_i/\tau}}{\sum_{j=1}^n e^{-s_j/\tau}}, \quad i = 1 \text{ or } 2, \quad (13)$$

where τ is the coefficient that regulates the "softness" of the softmin function. Additionally, we observe that significant estimation errors generally occur on the boundaries of surfaces which have small depth gradients (*i.e.*, surfaces that are nearly parallel to the XOY plane). This is due to the fact that the adjacent plane typically exhibits a much larger difference in depth gradient magnitude when compared to the reference surface (*i.e.*, the plane where the reference point is located). As a result, the calculated depth gradient is bound to differ from the depth gradient of the reference plane, even if the weight assigned to the reference plane's depth gradient is high, according to (13). To tackle this problem, when the smoothness of neighboring points differs greatly from each other, the weights for the depth gradients of adjacent and reference planes should be automatically assigned to 0 and 1, respectively. Therefore, we introduce truncated softmin

$$\mathcal{M}_t(s_i) = \begin{cases} \mathcal{M}(s_i) & (|s_2 - s_1| \leq 1) \\ \mathbb{1}_{\mathbb{R}^+}(s_i - 1) & (|s_2 - s_1| > 1) \end{cases} \quad (14)$$

to further improve surface normal accuracy, where $\mathbb{1}_{\mathbb{R}^+}(\cdot)$ is the indicator function mapping weight to either 0 or 1 based on the difference in adjacent pixel's surface smoothness. Our proposed DAG filter adaptively generates gradient filters based on surface smoothness, resulting in more accurate estimations of depth gradients, as outliers are effectively filtered out. In summary, the gradient filter of a given point p_i can be represented by the following expression²:

$$\mathbf{G}_i = \begin{bmatrix} G_h \\ G_v^\top \end{bmatrix}_i = \mathcal{M}_t(|\nabla^2 \mathbf{z}_i|) \begin{bmatrix} \Delta_b \\ \Delta_f \end{bmatrix}, \quad (15)$$

where $\mathbf{z}_i = \begin{bmatrix} z_l & z_r \\ z_u & z_d \end{bmatrix}$ is the neighborhood depth matrix of p_i .

²Here $\nabla^2(\cdot)$, $|\cdot|$, and $\mathcal{M}_t(\cdot)$ are element-wise operators.

C. MRF-Based Surface Normal Refinement

Our observation reveals that the surface normals of pixels near/on discontinuities are generally incorrect. This is due to the fact that non-coplanar points are used for local planar surface fitting, causing incorrect depth gradients (discussed in the previous subsection). To resolve this issue, we propose a fast and effective MRF-based optimization (post-processing) method, which significantly improves surface normal accuracy while having minimal impact on the processing speed.

The depth image can be modeled as an undirected graph $\mathcal{G} = (\mathcal{P}, \mathcal{E})$, where each node represents a pixel in the depth map, and each edge describes the connection between adjacent pixels. Let $\mathbf{N} = \{\mathbf{n}_p \mid p \in \mathcal{P}\}$ be a random variable set, where \mathbf{n}_p represents the estimated surface normal of point p . \mathbf{n}_p is conditionally independent of all other variables in \mathbf{N} :

$$\mathbf{n}_p \perp \mathbf{n}_{\mathcal{P} \setminus Q_p^+} \mid \mathbf{n}_{Q_p}. \quad (16)$$

Let $P(\mathbf{N} = n)$ be the joint probability distribution of \mathbf{N} , representing the probability of a particular field configuration n in surface normal field \mathbf{N} . Specifically, we set the size of the maximum clique to 2 to model our pairwise MRF. According to the Hammersley-Clifford theorem [18], $P(\mathbf{N} = n)$ is represented as follows:

$$P(\mathbf{N} = n) = \frac{1}{Z} \prod_{p \in \mathcal{P}} \phi(p) \prod_{(p, q_i) \in \mathcal{E}} \psi(p, q_i), \quad (17)$$

where Z is the partition function. (17) is mathematically equivalent to the energy minimization problem as follows:

$$\mathbf{E} = \sum_{p \in \mathcal{P}} \Phi(p) + \sum_{(p, q_i) \in \mathcal{E}} \Psi(p, q_i), \quad (18)$$

where the data term

$$\Phi(p) = \|\hat{\mathbf{n}}_p - \mathbf{n}_p\|_2, \quad (19)$$

enforces the consistency between the estimated surface normals $\hat{\mathbf{n}}_p$ and the observed surface normals \mathbf{n}_p , and the smoothness term

$$\Psi(p, q) = s_p \sum_i \mathcal{M}(s_i) \|\hat{\mathbf{n}}_p - \mathbf{n}_{q_i}\|_2, \quad (20)$$

smoothens the surface normal distribution between reference point p and its neighboring point q_i , where $\mathcal{M}(s_i)$ is the weight of the neighboring point q_i generated by the softmin function (13), \mathbf{n}_{q_i} is the observed surface normal of q_i , and s_p is the local surface smoothness that decides the weight between data term and smoothness term.

(19) suggests that the difference between the observed and estimated surface normals should be insignificant, while (20) implies that adjacent points on the same local planar surface should have consistent normal distributions.

IV. EXPERIMENTS

This section evaluates the performance of our proposed surface normal estimator and compares it with SoTA geometry-based SNEs. To simplify the presentation, we refer

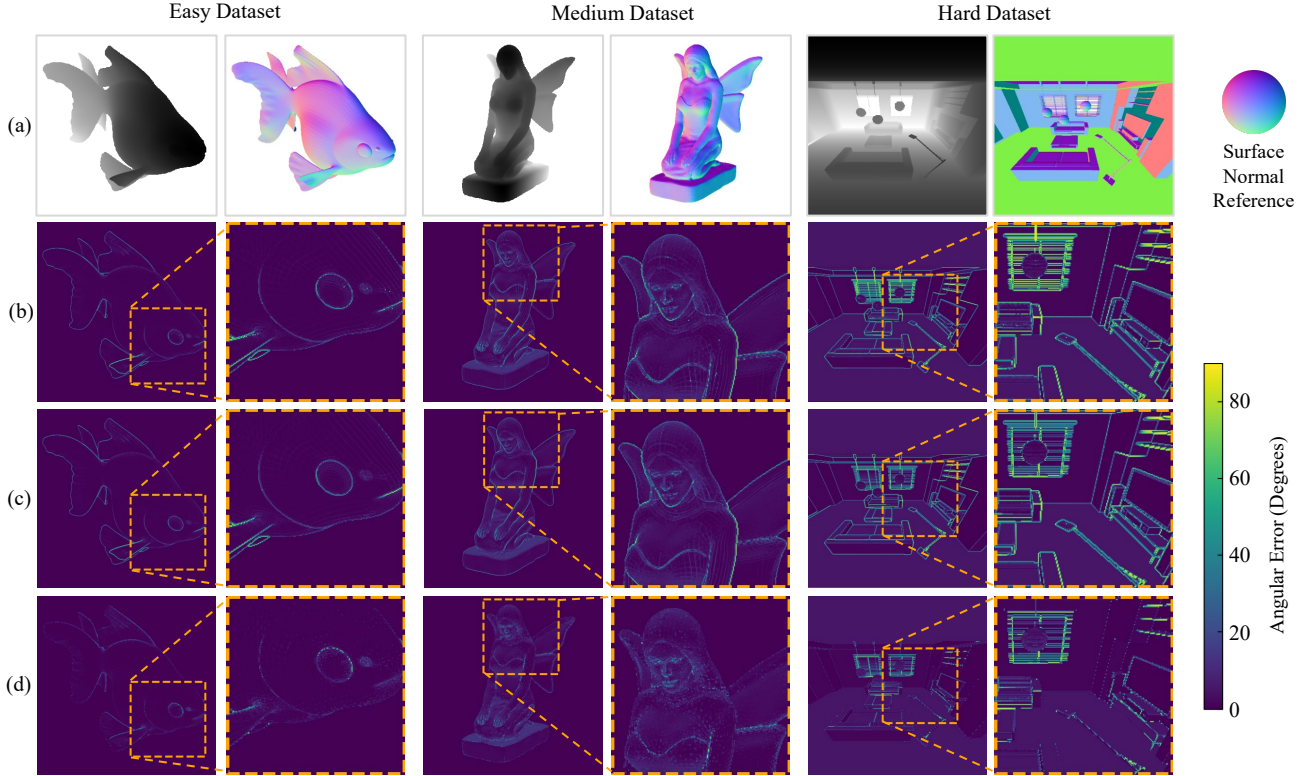


Fig. 3. Comparison of our proposed SNE with other SoTA geometry-based SNEs on the 3F2N [10] dataset: (a) depth maps and ground-truth surface normal maps; (b) error maps obtained using 3F2N (median filter); (c) error maps obtained using CP2TV; (d) error maps obtained using our proposed D2NT V3.

TABLE II
EVALUATION OF THE PROPOSED SNE USING FOUR DISCRETE LAPLACIAN FILTERS ON THE 3F2N DATASETS.

| Filter Config | D2NT+DAG | | | D2NT+MNR | | |
|------------------|-------------|-------------|--------------|-------------|-------------|-------------|
| | Easy | Medium | Hard | Easy | Medium | Hard |
| 1D DLF | 1.19 | 4.87 | 12.84 | 1.19 | 5.08 | 11.87 |
| DLF- α | 1.40 | 5.17 | 12.85 | 0.79 | 4.80 | 9.86 |
| DLF- β | 1.30 | 4.99 | 12.46 | 0.93 | 4.84 | 10.44 |
| DLF- γ | 1.36 | 5.05 | 13.03 | 1.36 | 5.05 | 13.03 |

to the basic depth-to-normal translator introduced in Sec. III-A as **D2NT**, the version that includes the DAG filter only as **D2NT V2**, and the version that includes both the DAG filter and the MNR module as **D2NT V3**.

Accurately determining surface normals from real-world range sensor data is infeasible due to the presence of noise. Although public datasets, such as NYUv2 [19] and DIODE [20], provide surface normal “ground truth”, it is often obtained through the interpolation of point sets into local planar surfaces, making the evaluation of SNEs with such “ground truth” unreliable. As a result, we conduct experiments on our previously published synthetic dataset [10].

A. Implementation Details and Evaluation Metrics

As discussed in Sec. III, local surface smoothness is computed through the convolution of the depth map with Laplacian kernels. To find the best convolution kernel,

four DLFs are used, including 1D DLF (horizontal kernel: $[1, -2, 1]$, and vertical kernel: $[1, -2, 1]^T$), DLF- α , DLF- β , and DLF- γ ³ [21]. The execution time of the four DLFs is comparable, as the optimization only occupies a minor portion of the overall process. As demonstrated in Table II, the best results on the 3F2N easy and medium datasets are achieved when using the 1D DLF for D2NT+DAG. Additionally, D2NT+MNR shows the best performance across all three 3F2N datasets when using the DLF- α for computing local surface smoothness. Therefore, we use the 1D DLF for the DAG filter and the DLF- β for the MNR module.

Moreover, to meet the real-time requirement, we simplified the implementation of our proposed MNR module. Specifically, when the level of discontinuity in the local surface is assessed to be low according to (11), we exclude the smoothness term in (18). Similarly, if a point is identified to be in and around discontinuities, we omit the data term and instead use the surface normal of the neighboring point with the highest smoothness.

Following [10], we use the average angular error e_A to quantitatively evaluate the performance of SNEs:

$$e_A = \frac{1}{N} \sum_{k=1}^N \cos^{-1} \left(\frac{\langle \mathbf{n}_k, \hat{\mathbf{n}}_k \rangle}{\|\mathbf{n}_k\|_2 \|\hat{\mathbf{n}}_k\|_2} \right), \quad (21)$$

$${}^3\text{DLF-}\alpha: \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \text{DLF-}\beta: \begin{bmatrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \text{DLF-}\gamma: \begin{bmatrix} 1 & 2 & 1 \\ 2 & -12 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

TABLE III

SPEED, ACCURACY, AND TRADE-OFF COMPARISONS AMONG SOTA GEOMETRY-BASED SNEs ON THE 3F2N DATASET.

| Real-Time | Method | t (ms) ↓ | e_A (degrees) ↓ | | | π (degrees/kHz) ↓ | | |
|-----------|--------------------|---------------|-------------------|-------------|-------------|-----------------------|----------------|----------------|
| | | | Easy | Medium | Hard | Easy | Medium | Hard |
| N | PlaneSVD [13] | 393.69 | 2.07 | 6.07 | 17.59 | 813.87 | 2389.73 | 6923.18 |
| | PlanePCA [14] | 631.88 | 2.07 | 6.07 | 17.59 | 1306.29 | 3835.59 | 11111.92 |
| | VectorSVD [16] | 563.21 | 2.13 | 6.27 | 18.01 | 1199.63 | 3529.11 | 10142.34 |
| | AreaWeighted [16] | 1092.24 | 2.20 | 6.27 | 17.03 | 2407.74 | 6843.56 | 18600.68 |
| | AngleWeighted [16] | 1032.88 | 1.79 | 5.67 | 13.26 | 1850.00 | 5855.62 | 13693.24 |
| | SDA-SNE [15] | 726.18 | 0.68 | 4.38 | 8.10 | 493.8 | 3180.67 | 5882.06 |
| Y | SNE-RoadSeg [3] | 7.92 | 2.04 | 6.28 | 16.37 | 16.16 | 49.74 | 129.65 |
| | 3F2N [10] | 10.97 | 1.66 | 5.69 | 15.31 | 18.18 | 62.38 | 168.03 |
| | CP2TV [11] | 2.23 | 1.69 | 6.01 | 13.82 | 3.75 | 13.39 | 30.76 |
| | D2NT (ours) | 1.82 | 1.54 | 5.64 | 15.32 | 3.05 | 10.25 | 27.84 |
| | D2NT V2 (ours) | 7.80 | 1.19 | 4.87 | 12.84 | 8.44 | 34.67 | 91.33 |
| | D2NT V3 (ours) | 7.99 | 0.89 | 4.78 | 9.86 | 7.09 | 38.28 | 78.91 |

TABLE IV

COMPARISON BETWEEN 3F2N AND CP2TV WITH AND WITHOUT OUR PROPOSED MNR MODULE EMBEDDED.

| Module Config | 3F2N | | | CP2TV | | |
|------------------|-------------|-------------|--------------|-------------|-------------|-------------|
| | Easy | Medium | Hard | Easy | Medium | Hard |
| w/o MNR | 1.66 | 5.69 | 15.32 | 1.69 | 6.02 | 13.82 |
| w/ MNR | 0.82 | 4.89 | 10.33 | 0.91 | 4.80 | 9.86 |
| Improvement | 50.7% | 14.0% | 32.5% | 40.8% | 15.0% | 35.6% |

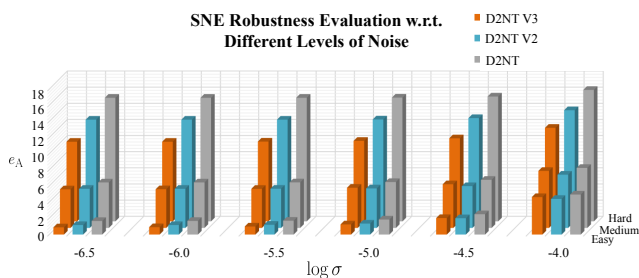


Fig. 4. Comparison among the three D2NT versions on the 3F2N datasets with different levels of Gaussian noise added.

where N is the number of valid pixels, and \mathbf{n}_k and $\hat{\mathbf{n}}_k$ are the ground truth and estimated surface normals, respectively.

In addition to accuracy evaluation, we adopt the metric

$$\pi = e_A t \text{ (degrees/kHz)} \quad (22)$$

proposed in [10] to quantify the trade-off between efficiency and accuracy of a given SNE. A high-performing (fast and accurate) SNE achieves a low π score.

B. Performance Comparison

As shown in Table III, our proposed surface normal estimators demonstrate superior performance compared to all other SoTA SNEs. D2NT achieves the highest computational efficiency and the optimum trade-off between speed and accuracy, while D2NT V3 achieves the highest accuracy (the e_A scores achieved by D2NT are less than 1° , 5° , and 9°

on the 3F2N easy, medium, and hard datasets, respectively). Furthermore, as illustrated in Fig. 3, our D2NT outperforms 3F2N and CP2TV, particularly in and around discontinuities.

We also conducted supplementary experiments to demonstrate the compatibility of our proposed MNR module with other depth-to-normal SNEs, as shown in Table IV. When incorporating the MNR module with 3F2N and CP2TV, the quality of their estimated surface normals is greatly improved, with a drop in 3F2N's e_A scores by 51%, 14%, and 33% on the 3F2N easy, medium, and hard datasets respectively and a decrease in CP2TV's e_A scores by 41%, 15%, and 36% on the same datasets. These results suggest that our proposed MNR module can be utilized in conjunction with other depth-to-normal SNEs and serves as an effective back-end optimization technique to enhance surface normal estimation in and around discontinuities.

As the used synthetic datasets are clean, we further evaluate the robustness of our methods in the presence of random Gaussian noise on the same datasets. As shown in Fig. 4, all three D2NT versions are stable with respect to different levels of Gaussian noise, and D2NT V3 is the most robust compared to the other two versions. In addition, it can be observed that our methods exhibit greater stability on the 3F2N medium and hard datasets, as compared to the 3F2N easy dataset. This is likely due to the added discontinuities caused by the Gaussian noise on the 3F2N easy dataset.

V. CONCLUSION

This paper presented an end-to-end depth-to-normal translator, a discontinuity-aware gradient filter, and an MRF-based surface normal refinement module. Extensive experimental results demonstrate that 1) our proposed depth-to-normal translator achieves the fastest execution speed and the best balance between computational efficiency and accuracy, and 2) the discontinuity-aware gradient filter and MRF-based surface normal refinement module can further improve its performance in and around discontinuities. Furthermore, our proposed MRF-based surface normal refinement module is also compatible with other depth-to-normal SNEs.

REFERENCES

- [1] Y. Li *et al.*, “Structure-SLAM: Low-drift monocular SLAM in indoor environments,” *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6583–6590, 2020.
- [2] Z. Liu *et al.*, “LPD-Net: 3D point cloud learning for large-scale place recognition and environment analysis,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 2831–2840.
- [3] R. Fan *et al.*, “SNE-RoadSeg: Incorporating surface normal information into semantic segmentation for accurate freespace detection,” in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 340–356.
- [4] H. Wang *et al.*, “Dynamic fusion module evolves drivable area and road anomaly detection: A benchmark and algorithms,” *IEEE Transactions on Cybernetics*, vol. 52, no. 10, pp. 10750–10760, 2022.
- [5] H. Wang *et al.*, “Applying surface normal information in drivable area and road anomaly detection for ground mobile robots,” in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 2706–2711.
- [6] R. Fan *et al.*, “Pothole detection based on disparity transformation and road surface modeling,” *IEEE Transactions on Image Processing*, vol. 29, pp. 897–908, 2019.
- [7] H. Wang *et al.*, “SNE-RoadSeg+: Rethinking depth-normal translation and deep supervision for freespace detection,” in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 1140–1145.
- [8] X. Qi *et al.*, “GeoNet: Geometric neural network for joint depth and surface normal estimation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 283–291.
- [9] X. Qi *et al.*, “GeoNet++: Iterative geometric neural network with edge-aware refinement for joint depth and surface normal estimation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 969–984, 2020.
- [10] R. Fan *et al.*, “Three-filters-to-normal: An accurate and ultrafast surface normal estimator,” *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 5405–5412, 2021.
- [11] Y. Nakagawa *et al.*, “Estimating surface normals with depth image gradients for fast and accurate registration,” in *2015 International Conference on 3D Vision (3DV)*. IEEE, 2015, pp. 640–647.
- [12] C. Wang *et al.*, “Comparison of local plane fitting methods for range data,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1. IEEE, 2001.
- [13] K. Klasing *et al.*, “Realtime segmentation of range data using continuous nearest neighbors,” in *2009 International Conference on Robotics and Automation (ICRA)*. IEEE, 2009, pp. 2431–2436.
- [14] K. Jordan *et al.*, “A quantitative evaluation of surface normal estimation in point clouds,” in *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2014, pp. 4220–4226.
- [15] N. Ming *et al.*, “SDA-SNE: Spatial discontinuity-aware surface normal estimation via multi-directional dynamic programming,” in *2022 International Conference on 3D Vision (3DV)*, 2022, pp. 486–494.
- [16] K. Klasing *et al.*, “Comparison of surface normal estimation methods for range sensing applications,” in *2009 International Conference on Robotics and Automation (ICRA)*. IEEE, 2009, pp. 3206–3211.
- [17] S. Jin *et al.*, “A comparison of algorithms for vertex normal computation,” *The Visual Computer*, vol. 21, no. 1, pp. 71–82, 2005.
- [18] J. M. Hammersley and P. Clifford, “Markov fields on finite graphs and lattices,” *Unpublished manuscript*, vol. 46, 1971.
- [19] N. Silberman *et al.*, “Indoor segmentation and support inference from RGBD images,” in *European Conference on Computer Vision (ECCV)*. Springer, 2012, pp. 746–760.
- [20] Vasiljevic *et al.*, “DIODE: A Dense Indoor and Outdoor D_Epth dataset,” *CoRR*, 2019.
- [21] M. Wardetzky, “Discrete laplace operators,” *An Excursion Through Discrete Differential Geometry: AMS Short Course, Discrete Differential Geometry, January 8-9, 2018, San Diego, California*, vol. 76, p. 1, 2020.