

Plant Leaf Meshes from Time-of-Flight RGB-D Sensors

Anonymous 3DV submission

Paper ID ****

Abstract

1. Introduction

This paper addresses the problem of automatically building 3D mesh models of plant leafs using inexpensive time-of-flight RGB-D sensors. Plant researchers, seeking to understand genetic underpinnings of plant growth [REF] and seeking to develop new varieties [REF], need automated ways to non-invasively measure plant phenotypes including growth, leaf distributions, orientations, photosynthesis and productivity [REF]. An important step in estimating all of these properties is obtaining 3D shape and pose for all the plant leafs. Plants cannot be moved or disturbed in growth chambers, and so our concept is to mount close-range RGB-D sensors in the chambers and acquire 3D mesh models of the leafs.

Time-of-flight RGB-D sensors have both advantages and drawbacks compared to other 3D sensors. [Expand the below]

- Dense depth image without scanning or rotating the target
- Depth even on non-textured regions
- Closely-aligned high resolution color image
- Near IR reflectance image
- Significantly higher depth error than laser scanners
- Short range and sensitive to specular artifacts
- Occlusion boundary artifacts (as with most ranging devices)

Here we address a key obstacle in obtaining accurate leaf models: the large noise in the depth images. [Expand the below:]

- There has been research in merging depth maps from multiple views and in reducing noise this way. This is

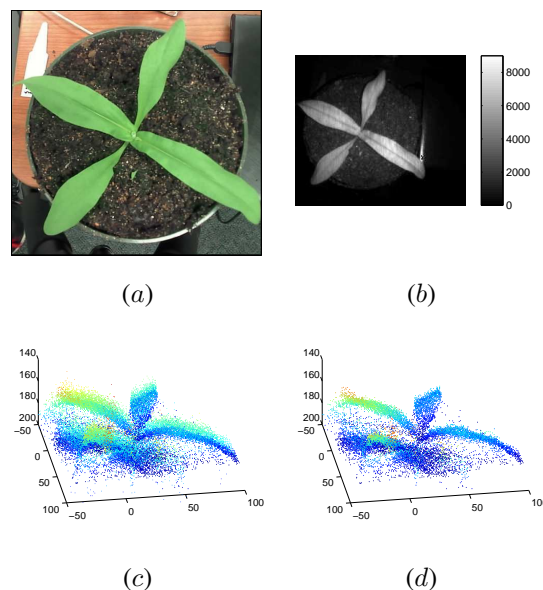


Figure 1. Illustration of sensor data. (a) Portion of color image. (b) IR reflectance image with reflectance values. (c) Portion of a single depth image surrounding plant projected into 3D showing significant depth noise. (d) Average of 60 depth images projected into 3D, with σ_S being the dominant source of noise. Units of 3D plots are mm. The goal of this paper is to combine these data elements into

not so helpful here as the sensor is fixed, and the size of the noise compared to the features we want to observe is large in our application.

- More ...

Our solution is a new mesh generation algorithm that leverages the high resolution color image, the dense depth estimates and the near-IR reflection image to overcome large depth errors. The paper is organized as follows. [Explain]

2. Related Work

3. Sensor Data Characterization

We explore using a new class of RGB-D sensor such as the Creative Sens3D [4] to build 3D mesh models of plant leaves. The sensor contains a high resolution color camera (1280×720 pixels) adjacent and parallel to a lower resolution depth camera (320×240 pixels). A flash IR emitter adjacent to these cameras illuminates the scene and depth sensor measures the time-of-travel for the reflected light as well as its reflectance over its pixel grid. These modalities are illustrated in Figure 1. It operates in a similar way to the Kinect 2 but is designed for closer range targets.

While the sensor produces dense depth measurements over target leaf surfaces, the noise in depth measurements is significantly larger than the features we are seeking to recover as illustrated in Figure 1(c). A key goal in this paper is to overcome this noise to maximize the accuracy of leaf shape estimates. We start in this section by modeling and quantifying the measurement noise.

3.1. Noise Characterization

Since the depth camera returns an IR reflectance in addition to a depth value at each pixel, both it and the color camera are initially calibrated using Zhang’s method [6] to obtain intrinsic and extrinsic parameters. Thus each pixel in each camera defines a ray from its camera origin. Depth noise is modeled as a one dimensional random variable, ε , along the ray for each pixel along its ray direction.

The depth noise, ε , is modeled as the sum of an image-varying term, ε_I , and a scene-varying term, ε_S :

$$\varepsilon = \varepsilon_I + \varepsilon_S. \quad (1)$$

The term ε_I models the random change in depth for camera pixels of subsequent images of a static scene from a static camera. To quantify this term we measured the standard deviation σ_I in depth of each pixel for a batch of 300 images of a fixed scene containing a flat matte surface. We repeated this at different poses and depths, and with different surface albedos. While target depth, inclination, albedo, and pixel position are all correlated with σ_I , we found that the best predictor for σ_I was the IR reflectance intensity, as shown in Figure 2. For typical scenes the single measurement noise in depth is roughly 5mm, although for low reflectivity objects or objects at long range this noise can increase significantly. Fortunately plant leaves are good IR reflectors.

Averaging depth measurements of a fixed scene will reduce the noise from ε_I , but will not reduce the noise from ε_S . This latter scene-varying term is constant for a static scene, but changes when the scene changes. To characterize this noise we first eliminated (approximately) the image-varying noise contribution by averaging over a large number of images (300). Then assuming ε_S is independent

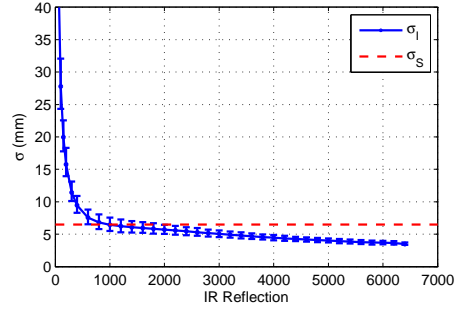


Figure 2. Image-varying noise, σ_I , is predicted well by the IR reflectance in raw units returned by the camera, see Figure 1(b). The scene-varying noise, σ_S , is plotted for comparison.

and identically distributed between pixels, we measured the variance of the pixel depth errors between a known flat surface and the estimated surface. In our experiments we obtained $\sigma_S = 6.5\text{mm}$, and found that it was insensitive to changes in depth.

The total pixel noise can be estimated assuming independence of ε_I and ε_S , and is given by:

$$\sigma^2 = \frac{\sigma_I^2}{N} + \sigma_S^2, \quad (2)$$

where N is the number of images averaged over. When averaging 5 or more depth images the scene-varying contribution, σ_S^2 , will dominate. There are additional sources of noise not modeled by this. These include object specularities, and mixed-depth pixels on object edges. These tend to produce very large image-varying noise, σ_I , and we can filter these points by discarding depth pixels with $\sigma_I > 20\text{mm}$.

4. Mesh Fitting

We pose mesh fitting to 3D point data as finding the most likely surface that would have generated those points. By incorporating prior surface assumptions, the fitting process estimates a continuous surface from discrete points that can eliminate much of the measurement noise. Methods that fit mesh models to 3D points often minimize the perpendicular distance of points to facets [REF]. This makes sense when point-cloud noise is equal in all directions or else the point noise is small compared to the facets. For our data the measurement noise is large and is not equal in all directions, but rather is along the depth camera rays. Hence the focus of this section is to develop a mesh fitting method that minimizes these pixel depth errors along the pixel rays.

In this paper we define a mesh in a 2D image space and project it into 3D. This is more limiting than full 3D meshes as it models only the surface portions visible from the sensor, but it also provides a number of advantages. Compared

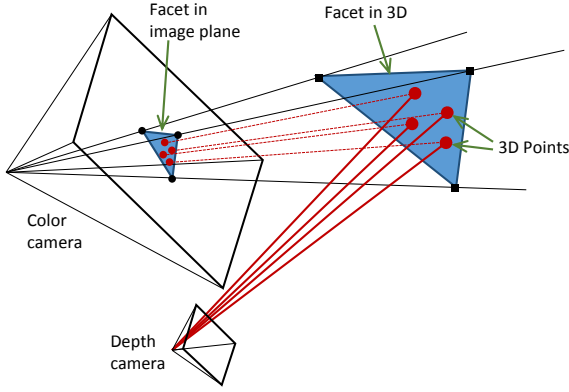


Figure 3. The parallel and adjacent color and depth cameras are shown as pyramids denoting their fields of view, and their size difference illustrates their relative resolutions. Three vertices in a color image define the rays on which the vertices of the corresponding 3D object facet must lie. This facet is fit using the 3D points projected out from the depth camera.

to methods that fit prior surface models to depth maps [ref], need to search of the space of poses, scales and distortions of the model with the chance of finding local minima. Compared to voxel-based models with implicit surfaces [ref], our method can better incorporate pixels uncertainties and surface priors, as well as having fewer discretization artifacts. In addition our method can naturally incorporate detailed features from the high-resolution color camera, and reflectivity information from the IR reflectance image.

4.1. Notation

A vertex, v , is a vector in 3D. In a given camera coordinate system, it projects onto a pixel on the unit focal-length image-plane $\tilde{v} = (u, v, 1)^\top$, where the “ \sim ” indicates a homogeneous vector, and u and v are the coordinates in this plane. Now \tilde{v} defines a ray from the camera origin, and the original vertex is obtained by scaling the image-plane vertex by its depth, λ , along the ray, namely: $v = \lambda \tilde{v}$.

4.2. Facet Model

Mesh fitting for an individual facet is illustrated in Figure 3. The 2D vertices and edge connections are determined in an image, in this case the color image although it could be the depth image, as described in section 5. If these vertices lie on a feature of the target leaf, such as its edge, we know that those 3D features like somewhere along the rays emanating from camera origin through those vertices. Hence a triangular facet approximation to the object surface will have vertices on these three rays.

The next step is to associate depth measurements with the facet. Pixels in the depth camera are projected along their rays out into 3D. The resulting 3D points that lie within

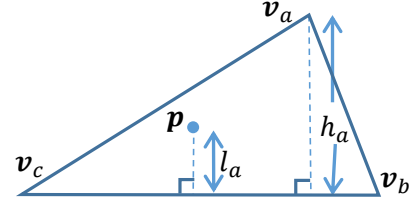


Figure 4. The coordinates of a point on a facet described by Eq. (3) are the weighted linear sum of the three vertex coordinates. The weight, α_a , for vertex v_a is given by $\alpha_a = \frac{l_a}{h_a}$, the ratio of its perpendicular distance l_a to the opposite edge to the vertex perpendicular distance h_a . Analogous expressions describe α_b and α_c .

the facet pyramid (defined by these rays through its vertices) will project into the 2D facet in the color image. Hence it is straightforward to associate 3D points with mesh facets.

To estimate the facet parameters from depth measurements we will express the depth points as a linear function of the vertices of its facet. We make a local orthographic approximation for the projection of a facet. This will be a good approximation as long as the facet size is small compared to is depth from the camera, which is true for most applications. Given this assumption, the coordinates of a point p lying on a facet can be expressed as a linear combination of the three vertex coordinates v_a , v_b and v_c as follows:

$$p = \alpha_a v_a + \alpha_b v_b + \alpha_c v_c. \quad (3)$$

The coefficients α_a , α_b and α_c are defined in Figure 4. Substituting in depth-scaled homogeneous vectors, and taking the third row, we obtain an equation for the point depth, λ_p :

$$\lambda_p = \alpha_a \lambda_a + \alpha_b \lambda_b + \alpha_c \lambda_c, \quad (4)$$

where λ_a is the depth of v_a and so forth.

4.3. Least Squares Depth

Equation (4) gives the depth of one point in terms of its facet vertices. For mesh with many facets and a measurement with many depth points, a vector of pixel depths, λ_d , and vector of vertex depths, λ_v , are related with a coefficient matrix, A , containing the appropriate α 's:

$$\lambda_d = A \lambda_v. \quad (5)$$

Given a measurement vector of point depths, $\bar{\lambda}_d$, expressed in the color camera coordinates, the error vector between these depths and the corresponding mesh points is: $\bar{\lambda}_d - A \lambda_v$. Notice that this error is along the color camera rays which to a good approximation are parallel to the depth camera rays, and thus the noise model in Eq. (2) applies. This justifies the following weighted squared error formulation:

$$E_{\text{depth}}(\lambda_v) = \|W \bar{\lambda}_d - W A \lambda_v\|^2, \quad (6)$$

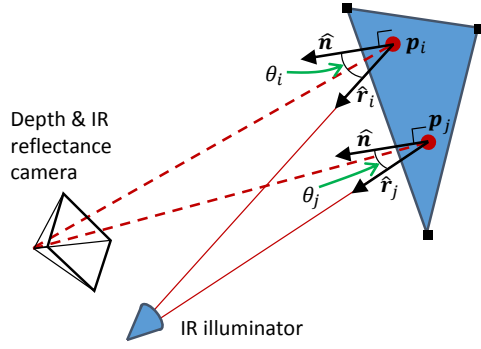


Figure 5. Our geometric model for the reflectance image. At each depth point the intensity of the reflectance image will depend on a number of factors including the angle between the facet normal and the ray to the illuminator shown for two points as θ_a and θ_b .

where W is a diagonal matrix containing the inverse standard deviation, σ^{-1} , from Eq. (2).

4.4. Shape from Shading

The depth camera measures IR reflectance in addition to depth at each pixel. Since the reflectance depends on the angle between the surface normal and the incident ray from the IR illuminator, the reflectance image can provide useful cues on the object surface. Shape from shading techniques model this dependence on the surface normal, along with additional surface assumptions such as smoothness, to estimate the normals and integrate an object surface [REF]. However the real world practicality of these methods has been limited since they generally require a single known light source position illuminating a Lambertian surface, the integration is sensitive to noise, and shape is obtained only up to a scale factor. Fortunately our application satisfies the key requirements of Shape from Shading, (we have a known light source and leafs are modeled well as Lambertian surfaces [1]), and our mesh model provides additional information that removes the need to integrate noisy surface normals. This section describes our use of Shape from Shading to improve shape recovery.

4.4.1 Reflectance Modeling

We build a simplified bidirectional reflectance model to explain the pixel values, R_i , of the IR reflectance image, shown in Figure 1(b). This model for pixel i is:

$$R_i = \frac{I_i \rho \hat{\mathbf{r}}_i \cdot \hat{\mathbf{n}}_i s_i}{r_i^h}. \quad (7)$$

Here I_i is the intensity of the ray from the IR illuminator assumed to be a point source and decreasing with the

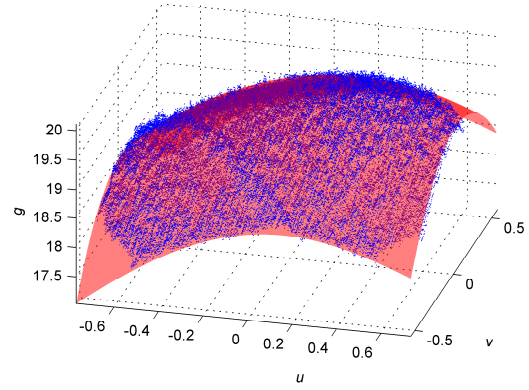


Figure 6. Our model for gain $g(u, v)$ in Eq. (9) is fit to a log reflectance image adjusted with known range and ray angles from Eq. (8). The dots are the measured data and the surface is the parameterized gain $g(u, v)$.

inverse square of the distance to target, r . Under a Lambertian assumption the reflected beam is decreased by the albedo, ρ , and the inner product of the ray direction $\hat{\mathbf{r}}_i$ and the normal, $\hat{\mathbf{n}}_i$, which is equal to the cosine of the angle between them, θ_i . Finally the sensor scales the incoming beam with a factor s_i . This model can be simplified further by approximating the outgoing ray I_i as being the same for a given pixel regardless of target depth. We define a pixel gain $g_i = \log(I_i s_i)$, and the resulting model is:

$$R_i = \frac{\exp(g_i) \rho \cos(\theta_i)}{r_i^2}. \quad (8)$$

The gain values can be calculated from a single depth image of a known surface with constant albedo up to a scale factor (of the unknown albedo). We observe large gain near the center of the image with tapering towards the edges, along with inter-pixel variations. We chose to treat the pixel variations as noise and model the gain as a polynomial in normalized pixel coordinates, u and v . That is our modeled gain is:

$$g(u, v) = \beta_1 + \beta_2 u + \beta_3 v + \beta_4 u^2 + \beta_5 v^2 \quad (9)$$

Taking the log of Eq. 8 we can do a least squares fit of the parameters β_i and so characterize the gain. We used a flat target with constant albedo and took data at various inclinations and depths. Figure 6 shows the resulting gain model, along with data from one depth image. We note that gain is primarily useful only up to a scale factor typically do not know the albedo of the target object.

4.4.2 Mesh Normals Estimation

Given a reflectance image providing, R_i for each pixel, a depth image from which we can calculate range to each

pixel, r_i , and our gain model, $g(u, v)$, we can rearrange and Eq. (eq:reflectance) to obtain an inclination prediction, η_i at each pixel:

$$\eta_i \equiv \rho_t \cos(\theta_i) = R_i r_i^2 \exp(-g(u, v)). \quad (10)$$

The scale factor is the unknown target albedo, ρ_t . If the target such as a leaf has a uniform albedo, this is the same for all target pixels.

Now each mesh facet has a unit normal which can be encoded in terms of the depths of its three vertices: $\hat{\mathbf{n}}(\lambda_a, \lambda_b, \lambda_c)$. The inner product of this normal with the unit ray to the illuminator is $\cos(\theta)$ and should agree with the values predicted in Eq. (10) for each reflectance pixel it contains. This leads to a Shape from Shading cost of

$$E_{SfS}(\lambda_v, \rho_t) = \sum_{i \in \mathcal{T}} \left\| \frac{\eta}{\rho_t} - \hat{\mathbf{n}} \cdot \hat{\mathbf{r}}_i \right\|^2. \quad (11)$$

The sum is over all target pixels, \mathcal{T} , projecting into the image mesh. One additional parameter is introduced, the unknown target albedo, ρ_t . This cost is nonlinear, but given a good initial parameters from the least squares solution to Eq. (6), it is readily minimized as a function of vertex depths and albedo.

4.5. Regularization

Prior models on surface properties can be incorporated into the mesh via regularization and in so doing reduce the impact of noise. Membrane energy is a well-used function in mesh optimization [3] and can be minimized using the discrete Laplacian operator. Here we use Laplacian smoothing due to its simplicity and good performance [3, 5, 2], although we modify it to accommodate image-based edge information. In addition, rather than apply Laplacian smoothing after the fact, entailing an iterative optimization [3], we show that Laplacian smoothing can be incorporated directly into the least squares mesh estimation. This has a number of advantages. First the smoothing penalty is traded off against measurement error rather than vertex offset. Second, the due to our ray constraints on the vertices we are able to derive a linear least squares solution to the Laplacian avoiding the need to iterate as in methods such as [3]. Finally when the regularization components are added to Eq. (6) they ensure that the solution is well-posed even when some of the facets have no depth points in them.

Laplacian smoothing uses an umbrella-operator [3] on a vertex, v , and its neighbors $v_i \in \mathcal{N}(v)$,

$$\mathbf{u}(v) = \frac{1}{n} \sum_{v_i \in \mathcal{N}} \mathbf{v}_i - \mathbf{v}, \quad (12)$$

for n neighbors as illustrated in Figure 7. In our model the vertices lie along known rays and so this operator can

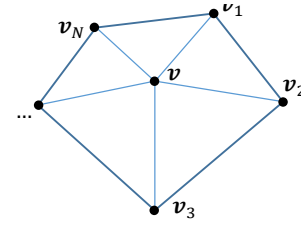


Figure 7. In discrete form the Laplacian is implemented as an umbrella operator, Eq. (12), over a vertex v and its first neighbors.

be expressed as a function of the vertex depth: $\mathbf{u}(\lambda) = \frac{1}{n} \sum_{i \in \mathcal{N}} \lambda_i \tilde{\mathbf{v}}_i - \lambda \tilde{\mathbf{v}}$. The squared magnitude $\|\mathbf{u}(\lambda)\|^2$ is a natural penalty term as it captures a discrete form of the membrane energy. Summing this over all vertices and arranging the known $\tilde{\mathbf{v}}$ components into a single matrix U , we obtain

$$E_{reg}(\lambda_v) = \|U\lambda\|^2. \quad (13)$$

The cost in Eq. (13) penalizes high curvature regions, and so provides a way incorporate smoothness priors into the mesh estimation. Now we may have image cues for creases or sharp edges on portions of a mesh, as we describe in section 5. We can modify the umbrella operator from Eq. (12) so that the neighbors are only other vertices on the crease. In this way the Laplacian acts along the crease, and not across it, allowing sharper folding along these edges.

5. Color-Based Mesh Cues

6. Results

7. Conclusion

References

- [1] M. Chelle. Could plant leaves be treated as lambertian surfaces in dense crop canopies to estimate light absorption? *Ecological Modelling*, 198(1):219 – 228, 2006. 4
- [2] C.-Y. Chen and K.-Y. Cheng. A sharpness dependent filter for mesh smoothing. *Computer Aided Geometric Design*, 22(5):376 – 391, 2005. Geometry Processing. 5
- [3] L. Kobbelt, S. Campagna, J. Vorsatz, and H.-P. Seidel. Interactive multi-resolution modeling on arbitrary meshes. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '98, pages 105–114, New York, NY, USA, 1998. ACM. 5
- [4] V. Nguyen, M. Chew, and S. Demidenko. Vietnamese sign language reader using intel creative senz3d. In *IEEE International Conference on Automation, Robotics and Applications (ICARA)*, pages 77–82, 2015. 2
- [5] Y. Ohtake, A. Belyaev, and I. Bogaevski. Mesh regularization and adaptive smoothing. *Computer-Aided Design*, 33(11):789 – 800, 2001. 5

- [6] Z. Zhang. A flexible new technique for camera calibration.
Pattern Analysis and Machine Intelligence, IEEE Transactions on, 22(11):1330–1334, Nov 2000. 2