

# Statistical Formulation of Region-Based Segmentation (2007)

$P(\Omega)$ : optimal partition of the image plane  $\Omega$   
 切割

By maximizing a posteriori probability  $P(P(\Omega)|I)$  for a given image  $I$ .

$$\therefore P(P(\Omega), I) = P(P(\Omega)|I) P(I) = P(I|P(\Omega)) P(P(\Omega))$$

$$\therefore P(P(\Omega)|I) \propto \underbrace{P(I|P(\Omega))}_{\text{Image Based cues observation given a model state.}} \underbrace{P(P(\Omega))}_{\text{geometric properties of the partition.}}$$

The image partition to be composed of  $N$  regions without correlation between labellings

$$P(I|P(\Omega)) = P(I|\{\Omega_1, \dots, \Omega_N\}) = \prod_{i=1}^N \underbrace{P(I|\Omega_i)}_{\substack{\text{the prob of observing an} \\ \text{image } I \text{ when } \Omega_i \text{ is a} \\ \text{region of interest.}}}$$

Assume that the observation  $I$  consists of a set of feature values  $f(x)$  associated with each image location. (image intensity, pixel color, gradient ....)

We make the assumption that the value of  $f$  at different locations of the same region can be modeled as independent and identically distributed realizations of the same random process.

Let  $p_i$  be the PDF of this random process in  $\Omega_i$ .

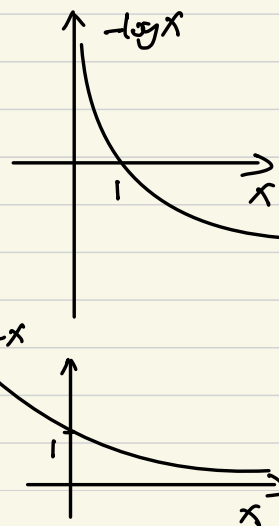
$$\therefore P(I|P(\Omega)) = \prod_{i=1}^N \prod_{x \in \Omega_i} \underbrace{(p_i(f(x)))^{dx}}_{\text{特征}}$$

$dx$  is the bin volume.

★ Maximization of the  $P(P(\Omega)|I)$  is equivalent to minimizing its negative logarithm

$$\begin{aligned} \therefore & -\log\left(\prod_{i=1}^N \prod_{x \in \Omega_i} (p_i(f(x)))^{dx}\right) \\ & = -\sum_i \int_{\Omega_i} \log(p_i(f(x))) dx \end{aligned}$$

If  $P(P(\Omega)) \propto e^{-\nu|C|}$ ,  $\nu > 0$ , favors a short length  $|C|$  of partition boundary.



$$\therefore E(\{\Omega_1, \dots, \Omega_N\}) = -\sum_i \int_{\Omega_i} \log p_i(f(x)) dx + v|\Omega|$$

$$\therefore \min(E\{\Omega_i\}) = \min(-\sum_i \int_{\Omega_i} p(f(x)) dx)$$

### Two-Phase Level Set

Assume that a partitioning of the domain  $\Omega$  such that each pixel is ascribed to one of two possible phases.

$$\therefore E(\phi) = \int_{\Omega} -H(\phi) \log p(f) - (1-H(\phi)) \log p(f) + v|\nabla H(\phi)| dx$$

$$\text{where } H(\phi) = \begin{cases} 1 & \text{if } \phi \geq 0, \\ 0 & \text{else.} \end{cases}$$



# Robust Real-Time Visual Tracking Using Pixel-Wise Posteriors (2018)

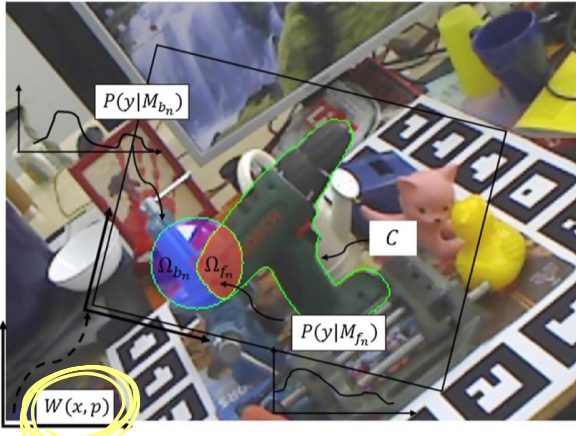


Fig. 12 Local region model

$W(x, p)$ : take a pixel location  $x$  in the object frame and warps it into the image frame according to parameters  $p$ .

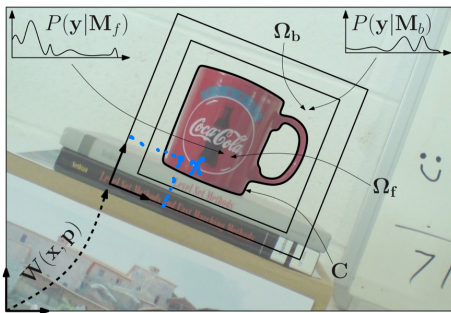
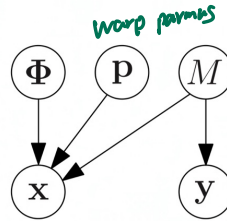


Fig. 1. (Left): Representation of object showing: the contour  $C$ , the set of foreground pixels  $\Omega_f$ , the set of background pixels  $\Omega_b$ , the foreground model  $P(y|M_f)$ , the background model  $P(y|M_b)$  and the warp  $W(x, p)$ ; (Right): Graphical representation of our generative model representing the image as a bag-of-pixels



- $x$ : A pixel location in the object coordinate frame.
- $y$ : A pixel value (in our experiments this is a YUV value).
- $I$ : Image.
- $W(x, p)$ : Warp with parameters  $p$ .
- $M = \{M_f, M_b\}$ : Model parameter either foreground or background.
- $P(y|M_f)$ : Foreground model over pixel values  $y$ .
- $P(y|M_b)$ : Background model over pixel values  $y$ .
- $C$ : The contour that segments the foreground from background.
- $\Phi(x)$ : Shape kernel (in our case the level-set embedding function).
- $\Omega = \{\Omega_f, \Omega_b\}$ : Pixels in the object frame  $[\{x_0, y_0\}, \dots, \{x_N, y_N\}]$ , which is partitioned into foreground pixels  $\Omega_f$  and background pixels  $\Omega_b$ .
- $H_e(z)$ : Smoothed Heaviside step function.
- $\delta_e(z)$ : Smoothed Dirac delta function.

The joint distribution for a single pixel given the model:

$$\begin{aligned}
 P(x, y, \Phi, p, M) &= P(x, y | \Phi, p, M) P(\Phi, p, M) \\
 &= P(x | \Phi, p, M) P(y | \Phi, p, M) P(\Phi) P(p) P(M) \\
 &= P(x | \Phi, p, M) P(y | M) P(\Phi) P(p) P(M) \\
 &= P(x | \Phi, p, M) P(y, M) P(\Phi) P(p)
 \end{aligned}$$

$$\therefore P(x, \Phi, p, M | y) P(y) = P(x | \Phi, p, M) P(M | y) P(y) P(\Phi) P(p)$$

$$\begin{aligned}
 \therefore P(x, \Phi, p, M | y) &= P(x | \Phi, p, M) P(M | y) P(\Phi) P(p) \\
 &= P(\Phi) P(p) \sum_{j=f,b} P(x | \Phi, p, M_j) P(M_j | y)
 \end{aligned}$$

$$P(M_j | y) = \frac{P(y | M_j) P(M_j)}{P(y)}$$

$$= \frac{P(y | M_j) P(M_j)}{\sum_{i=f,b} P(y | M_i) P(M_i)}, j \in \{f, b\} = P(x, \Phi, p | y)$$

pixel-wise posterior of the model  $M$  given a pixel values  $y$ .

$$P(x, \Phi, p | y) = P(\Phi, p | x, y) P(x | y) = P(\Phi, p | x, y) P(x)$$

$$\Rightarrow P(\Phi, p | x, y) = P(x, \Phi, p | y) \frac{1}{P(x)}$$

The prob of the shape and the location  $p$  given pixel  $\{x, y\}$

$$= P(\Phi) P(p) \left( \frac{1}{P(x)} \right) \sum_{j=f,b} P(x | \Phi, p, M_j) P(M_j | y)$$

it can be drop as it is constant.

$$P(x | \Phi, p, y)$$

The pixel wise posteriors:

$$P(\Phi, p | \Omega) \propto P(\Phi) P(p) \prod_{i=1}^N \left\{ \sum_{j=f,b} P(X_i | \Phi, p, M_j) P(M_j | y_i) \right\}$$

$\Phi$  is known from map1.  
constant for all pixels

$$\therefore P(\Phi, p | \Omega) \propto \prod_{i=1}^N \left\{ \sum_{j=f,b} P(X_i | \Phi, p, M_j) P(M_j | y_i) \right\}$$

$$P(X_i | \Phi, p, M_f) = \begin{cases} \frac{He(\Phi(x_i))}{n_f}, & j=f \\ \frac{1 - He(\Phi(x_i))}{n_b}, & j=b \end{cases}$$

如果  $x_i$  是在 foreground, 概率为  $\frac{1}{n_f}$ , 否则为 0.  
如果  $x_i$  是在 background, 概率为  $\frac{1}{n_b}$ , 否则为 0.

$$P(M_j | y_i) = \frac{P(y_i | M_j) P(M_j)}{\sum_{i=f,b} P(y_i | M_i) P(M_i)}$$

$$\text{and } P(M_j) = \begin{cases} \frac{n_f}{n}, & j=f \\ \frac{n_b}{n}, & j=b \end{cases}, \quad \eta = n_f + n_b = K$$

$$n_f = \sum_{i=1}^K He(\Phi(x_i))$$

$$n_b = \sum_{i=1}^K (1 - He(\Phi(x_i)))$$

$$\therefore P(\Phi, p | \Omega) \propto \prod_{i=1}^N \left\{ \sum_{j=f,b} P(X_i | \Phi, p, M_j) P(M_j | y_i) \right\}$$

$$= \prod_{i=1}^N \left\{ P(X_i | \Phi, p, M_f) P(M_f | y_i) + P(X_i | \Phi, p, M_b) P(M_b | y_i) \right\}$$

$$= \prod_{i=1}^N \left\{ \frac{He(\Phi(x_i))}{n_f} \frac{P(y_i | M_f) \frac{n_f}{n}}{P(y_i | M_f) \frac{n_f}{n} + P(y_i | M_b) \frac{n_b}{n}} \right.$$

$$\left. + \frac{1 - He(\Phi(x_i))}{n_b} \frac{P(y_i | M_b) \frac{n_b}{n}}{P(y_i | M_f) \frac{n_f}{n} + P(y_i | M_b) \frac{n_b}{n}} \right\}$$

$$= \prod_{i=1}^N \left\{ \frac{He(\Phi(x_i)) P(y_i | M_f)}{P(y_i | M_f) n_f + P(y_i | M_b) n_b} + \frac{[1 - He(\Phi(x_i))] P(y_i | M_b)}{P(y_i | M_f) n_f + P(y_i | M_b) n_b} \right\}$$

"  $p_f$  "  $p_b$

## 2D-3D Pose Estimation of Heterogeneous Objects Using a Region Based Approach (2016)

It is common in region-based segmentation for a closed curve to be evolved such that discrepancy (相差) between the statistics of the foreground and those of the background region is maximised:

$$E = \int_{\Omega_f} r_f(I(x), C) d\Omega + \int_{\Omega_b} r_b(I(x), C) d\Omega$$

$\Omega_f$ : foreground of image  $\subset \mathbb{R}^2$ ,  $I\{\{x_0, y_0\}, \dots, \{x_N, y_N\}\}$ ,  $N$  is number of pixels.

$\Omega_b$ : background

$x = [x, y]$

$I(x) = y$ , pixel value

$C = \{(x, y) \in \mathbb{R}^2 \mid \Phi(x) = 0\}$ : the contour, the zero level-set of the function  $\Phi(x)$

where  $r_f$  and  $r_b$  are two monotonically decreasing functions, measuring the matching quality of image pixels with respect to the foreground and background models.

① Rewrite the energy function:

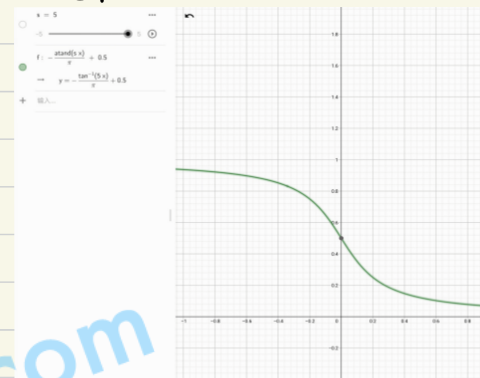
$$E(\Phi) = \int_{\Omega} H_\epsilon(\Phi) r_f(x) + (1 - H_\epsilon(\Phi)) r_b(x) d\Omega$$

$$\Phi(x) = \begin{cases} -d(x), & \forall x \in \Omega_f \\ d(x), & \forall x \in \Omega_b \end{cases}, \quad d(x) = \min_{x_c \in C} |x - x_c|$$

$H_\epsilon$  is smoothed Heaviside step,

$$H_\epsilon(\Phi(x)) = -\frac{1}{\pi} \arctan\left(\frac{\Phi(x)}{\epsilon}\right) + \frac{1}{2},$$

with  $\epsilon$  determining the pitch of the smoothed transition



$r_f$  and  $r_b$  are given by the likelihood of a pixel property (such as color) under a given model, i.e.  $r(x) = P(y|M)$ , with  $M \in \{M_f, M_b\}$ .

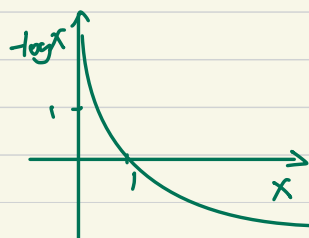
PDF of a pixel value  $y$  belonging to the foreground or background.

② In contrast, assuming pixel-wise independence, and replacing integration with summation, the energy given by the negative  $\log$ (posterior) probability of the shape of the contour  $(\Phi)$  given the image data:

$$P(\Phi | I) \propto \prod_{x \in \Omega} (H_\epsilon(\Phi(x)) P_f(y) + (1 - H_\epsilon(\Phi(x))) P_b(y))$$

$$\therefore E(\Phi) = -\log(P(\Phi | I))$$

$$= -\sum_{x \in \Omega} \log(H_\epsilon(\Phi) P_f + (1 - H_\epsilon(\Phi)) P_b)$$



$y = I(x)$ . pixel value

$$\therefore P(y) = P(y|M_f)P(M_f) + P(y|M_b)P(M_b)$$

$$\therefore \underline{P(M_j|y)} P(y) = P(y|M_j) P(M_j), \quad j = \{f, b\}$$

$$\Rightarrow \underline{P(M_j|y)} = \frac{P(y|M_j) P(M_j)}{\underbrace{\sum_{i=f,b} P(y|M_i) P(M_i)}_{P(y)}}$$

Let  $\eta_f = \sum_{i=1}^N H_e(\Phi(x_i))$ ,  $\eta_b = \sum_{i=1}^N (1 - H_e(\Phi(x_i)))$ ,  $\eta = \eta_f + \eta_b$

foreground pixels number      background pixels number      all pixels number

$$\therefore \underline{P(M_f)} = \frac{\eta_f}{\eta}, \quad \underline{P(M_b)} = \frac{\eta_b}{\eta}$$

$$\therefore \underline{P_j(y)} = \frac{P(M_j|y)}{\eta_j} = \frac{1}{\eta_j} \frac{P(y|M_j) \eta_j}{\sum_{i=f,b} P(y|M_i) \frac{\eta_i}{\eta}} = \frac{P(y|M_j)}{P(y|M_f) \eta_f + P(y|M_b) \eta_b}$$

The  $P(y|M_j)$  are estimated by calculating the histograms (256 bins) of each color space and smooth using a Gaussian kernel. For simplicity, we assume the RGB channels are independent.

$$\therefore \underline{P(y_{RGB}|M)} = P(y_R|M) P(y_G|M) P(y_B|M)$$

1. Initialize pose parameters.  $\lambda = \lambda_0$
2. For each iteration:
  - (a) Apply 3D transformation to object.
  - (b) Project 3D model on to image plane.
  - (c) For each local region:
    - (i) Estimate local region statistics.
    - (ii) Calculate local energy gradient with respect to pose parameters,  $\nabla E_n$ .
  - (d) Fuse local region gradients,  $\nabla E = f(\nabla E_n)$ .
  - (e) Find optimal step size,  $s$ .
  - (f) Update pose parameters  $\lambda = \lambda - s \nabla E$ .

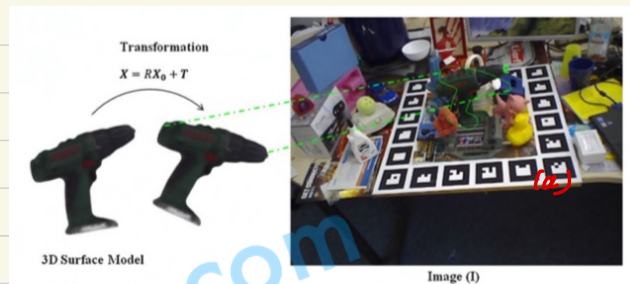


Fig. 5 Driller object projection onto the image plane

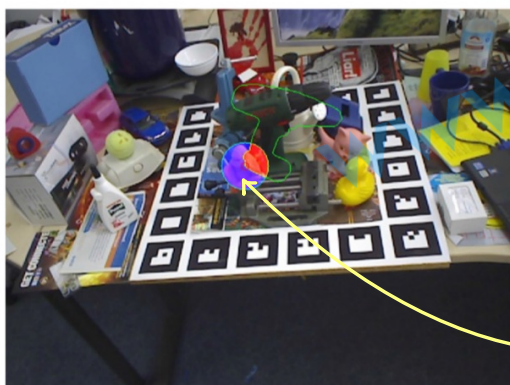
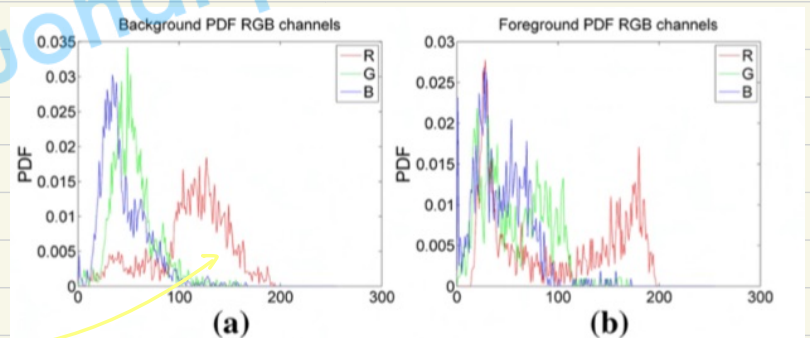


Fig. 8 Example of a local region extraction, divided into the local foreground (red) and local background (blue) (Color figure online)





Let  $\lambda = [\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5, \lambda_6]^T$  the pose parameters.

$$t = [\lambda_1, \lambda_2, \lambda_3]^T$$

$$R = \exp(\Lambda), \quad \Lambda = \begin{bmatrix} 0 & -\lambda_6 & \lambda_5 \\ \lambda_6 & 0 & -\lambda_4 \\ -\lambda_5 & \lambda_4 & 0 \end{bmatrix}$$

$$\therefore \frac{\partial F}{\partial \lambda_i} = \frac{\partial (-\sum_{x \in \Omega} \log(\text{He}(\Phi)P_f + (1 - \text{He}(\Phi))P_b))}{\partial \lambda_i}$$

$$\frac{d \ln x}{dx} = \frac{1}{x}$$

$$= -\sum_{x \in \Omega} \frac{1}{\underbrace{\text{He}(\Phi)P_f + (1 - \text{He}(\Phi))P_b}_{\text{He}(\Phi)(P_f - P_b) + P_b}} \cdot (P_f - P_b) \cdot \frac{\partial \text{He}(\Phi)}{\partial \lambda_i}$$

rely on statistical properties est  
weight applied to the geo differentials.

geo differentials of the  
object with respect to the  
pose parameters.

$$\therefore \frac{\partial \text{He}(\Phi)}{\partial \lambda_i} = \frac{\partial \text{He}(\Phi)}{\partial \Phi} \left[ \underbrace{\frac{\partial \Phi}{\partial u}}_{\substack{\uparrow \\ \Phi(u,v) \\ \text{image pixel}}} \frac{\partial u}{\partial \lambda_i} + \frac{\partial \Phi}{\partial v} \frac{\partial v}{\partial \lambda_i} \right] = J_{\Phi}(\Phi) \left[ \frac{\partial \Phi}{\partial u} \frac{\partial u}{\partial \lambda_i} + \frac{\partial \Phi}{\partial v} \frac{\partial v}{\partial \lambda_i} \right]$$

$$\therefore \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \frac{X}{Z} f_x + u_0 \\ \frac{Y}{Z} f_y + v_0 \end{bmatrix}$$

$$\therefore \frac{\partial u}{\partial \lambda_i} = f_x \frac{\partial (X/Z)}{\partial \lambda_i} = f_x \left( \frac{1}{Z} \frac{\partial X}{\partial \lambda_i} - X \frac{1}{Z^2} \frac{\partial Z}{\partial \lambda_i} \right)$$

$$\therefore \frac{\partial v}{\partial \lambda_i} = f_y \left( \frac{1}{Z} \frac{\partial Y}{\partial \lambda_i} - Y \frac{1}{Z^2} \frac{\partial Z}{\partial \lambda_i} \right)$$

## Local Region Based Pose Estimation

Defined a mask function  $B_n(x_i, x_c)$  which centered at  $x_c$ , with a radius  $d$ :

$$B_n(x_i, x_c) = \begin{cases} 1 & |x_i - x_c| < d \\ 0 & \text{else} \end{cases}$$

$$\frac{\partial E}{\partial \lambda_i} = - \sum_{x \in \Omega} \frac{(P_f - P_b)}{He(\Phi)P_f + (1 - He(\Phi))P_b} \cdot \frac{\partial He(\Phi)}{\partial \lambda_i}$$

$\therefore He(\Phi)$  value only change on the counter! So we define a local regions around every point along the object's contour.

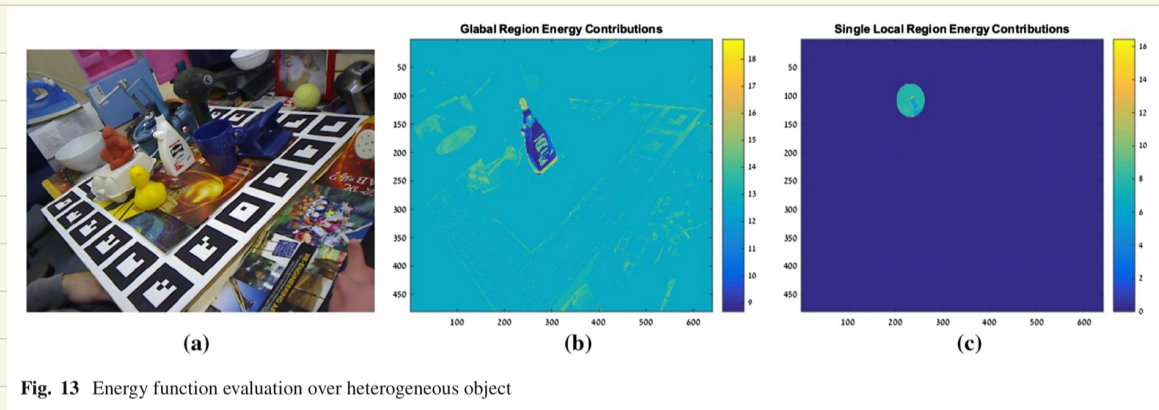


Fig. 13 Energy function evaluation over heterogeneous object

Using large sized regions, the region statistics will be more robust to the initialization of the object's pose. Change in the objects pose will have a lower impact on the statistical models. Using small sized regions, the ability to capture variations in the region's statistics increases.

A small object or a cluttered (~~slowly~~) background require a small radius to correctly capture the variation between regions, whereas for a large object and a slowly varying background, a large radius is preferred.

# <<Real-Time Monocular Pose Estimation of 3D Objects Using Temporally Consistent>> (2017)

Result evaluate :  $\frac{E}{|\Omega|} > t$  ,  $t \in [0.5, 0.6]$  , lost

$$E = - \sum_{x \in \Omega} \log (He(\Phi(x)) \bar{p}_f(x, y) + (1 - He(\Phi(x))) \bar{p}_b(x, y))$$

we want get the min of E

$\Omega \subset \mathbb{R}^2$  being the image domain,  $\Omega \rightarrow \{0, 1\}$

$|\Omega|$ : foreground region pixel number background  $\nwarrow$  foreground

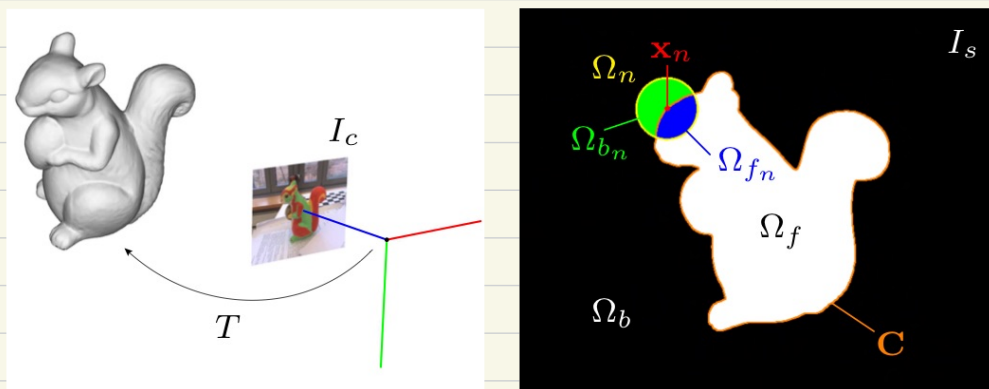


Figure 2. Overview of our pose estimation setting. Left: The object pose  $T$  relative to a camera based on color image  $I_c$  and a 3D model of the object. Right: Silhouette  $I_s$  generated by projecting the surface model into the image plane using an estimated pose  $T$ .

tlc-histograms :

那这里又存在一个local采样中心点的时序一致的问题，就是说前后帧的匹配问题，必须得找到  $x_i(t_l) \leftrightarrow x_i(t_{l-1})$  间的匹配关系。Real-Time Monocular Pose Estimation of 3D Objects using Temporally Consistent Local Color Histograms 提出 temporally consistent local color histograms (tlc-histograms) 试图解决这一问题。

大致思想是，由于三维顶点是唯一且有序的，其投影到二维平面有像素点与之对应。根据不同角度和视口，每个三维模型顶点  $X_i$ （都有可能投影成为二维图像边缘上的像素点，也就意味着）都可以有一个关联着的所谓的 local foreground & background histograms，文中也称呼为 tlc-histogram。

那也就是说齐次坐标的  $\tilde{x}_i$  其在二维投影平面的前后帧像素关系：

$$\pi(K(T(t_l)\tilde{X}_i)_{3 \times 1}) \leftrightarrow \pi(K(T(t_{l-1})\tilde{X}_i)_{3 \times 1})$$

对于第一次投影到二维轮廓线上的三维顶点，都在当前帧，以该点为圆心进行一次tlc-histogram初始化。而对于已保留有前帧histogram信息的点，即说明该像素  $x_i = \pi(K(T\tilde{X}_i)_{3 \times 1})$  其前后帧都落在轮廓线上，即  $x_i(t_l), x_i(t_{l-1}) \in C$ ，其更新方法就是前文提到的：

$$P(y|M_{f_i}) = (1 - \alpha_f)P^{t_{l-1}}(y|M_{f_i}) + \alpha_f P^{t_l}(y|M_{f_i}),$$

$$P(y|M_{b_i}) = (1 - \alpha_b)P^{t_{l-1}}(y|M_{b_i}) + \alpha_b P^{t_l}(y|M_{b_i}),$$

## Pose Detection

In order to cover different scales, each of these 48 base orientations is used to generate a template at a close, an intermediate and a far distance to the camera, resulting in overall 144 base templates.

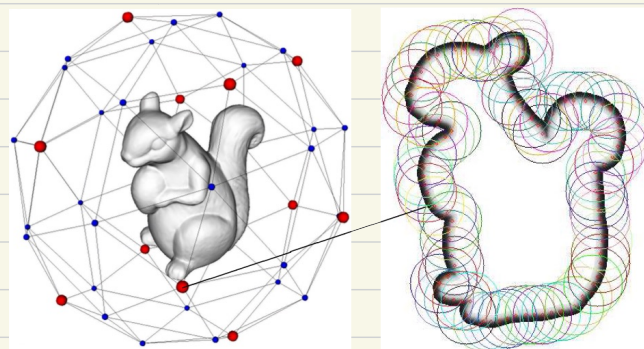


Figure 5. The template views used for pose detection. Left: A subdivided icosahedron generates the outer plane rotation of the template views. Red dots indicate vertices corresponding to the base templates, blue dots indicate those provided by subdivision used for the neighboring templates. Right: An example base template visualized by the corresponding  $H_e(\Phi)$  (grey pixels) with the local histogram regions depicted (colored circles) along the contour.



## 3.4. Approximating for real-time tracking

Computing the SDF from Eq. 3.1 has already three costly steps. We need a silhouette rendering  $\Omega_f$  of the current model pose, an extraction of the contour  $C$  and lastly, a subsequent distance transform embedding  $\phi$ . While [29] perform GPU rendering and couple computation of the SDF and its gradient in the same pass to be faster, [17] perform hierarchical ray-tracing on the CPU and extract the contour via Scharr operators. We make two key observations:

1. Only the actual contour points are required
2. Neighboring points provide superfluous information because of similar curvature

We thus propose a cheap yet very effective approximation of the model render space that avoids both online rendering and contour extraction. In an offline stage, we equidistantly

sample viewpoints  $V_i$  on a unit sphere around the object model, render from each and extract the 3D contour points to store view-dependent **sparse** 3D sampling sets in **local object space** (see Figure 4). Since we will utilize these points in 3D space, we neither need to sample in scale nor for different inplane rotations. Finally, we store for each contour point its 2D gradient orientation and sample a set of interior surface points with their normals (see Figure 5).

In a naive approach, all involved terms from Eq. 7 would be computed densely, *i.e.*  $\forall x \in \Omega$ , which is prohibitively costly for real-time scenarios. The related work evaluates the energy only in a narrow band around the contour since the residuals decay quickly when leaving the interface. We therefore propose to compute Eq. 8 in a narrow band along a sparse set of selected contour points where we compute  $\phi$  along rays. Each projected contour point shoots a positive and negative ray perpendicularly to the contour, *i.e.* along its normal. Building on that, we introduce the idea of ray integration for 3D contour points such that we do not create pixel-wise but **ray-wise** Jacobians which leads to a smaller reduction step and a better conditioning of the normal system in Eq. 9 than [17] and their approach.

To formalize, we have a model pose  $[R, t]$  during tracking and avoid rendering by computing the camera position in object space  $O := -R^T t$ . We normalize to unit length and find the closest viewpoint  $V^*$  quickly via dot products:

$$V^* := \underset{V_i}{\operatorname{argmax}} \langle V_i, O / \|O\| \rangle. \quad (16)$$

Each local 3D sample point of the contour  $X_i$  from  $V^*$  is then transformed and projected to a 2D contour sample point  $x_i = \pi(RX_i + t)$  which is then used to shoot rays into the object interior and into the opposite direction.

To get the orientation of each ray, we cannot rely anymore on the value during pre-rendering since the current model pose might have an inplane rotation not accounted for. Given a contour point with 2D rotation angle  $\theta$  during pre-rendering, we could embed it into 3D space via  $v = (\cos \theta, \sin \theta, 0)$  and later multiply it with the current model rotation  $R$ . Although this works in practice, the projection of  $R \cdot v$  onto the image plane can be off at times. We thus propose a new approximation of the inplane rotation where we seek to decompose  $R = R_{\text{inplane}} \cdot R_{\text{canonical}}$  *s.t.* one part describes a general rotation around the object center in a canonical frame and the other a rotation around the view direction of the camera (*i.e.* inplane). Although ill-posed in general, we exploit our knowledge about the closest viewpoint by assuming  $R_{\text{canonical}} \approx R_{V^*}$  and propose to approximate a rotation  $\tilde{R}$  on the xy-plane via

$$\tilde{R} := R \cdot R_{V^*}^T. \quad (17)$$

We then extract the angle  $\theta = \operatorname{acos}(\tilde{R}_{1,1})$  via the first element. With larger viewpoint deviation  $\|V^* - \frac{O}{\|O\|}\|$ , this approximation worsens but our sphere sampling is dense enough to alleviate this in practice. We re-orient each contour gradient  $\tilde{g}_i := (g_i + \theta) \bmod 2\pi$  and shoot rays to compute the residuals and  $\frac{\partial H_\phi}{\partial \phi}$  from Eq. 7 (see Figure 5 to compare the orientations and the bottom row in Figure 2 for the SDF rays).

The final missing building block is the derivative of the SDF  $\frac{\partial \phi}{\partial x}$  which cannot be computed numerically since we are missing dense information. We thus compute it geometrically, similar to [17]. Whereas their computation is exact when assuming local planarity by projections onto the principal ray, our approach is faster while incurring a small error which is negligible in practice. Given a ray  $r = (r_x, r_y)$  from contour point  $p = (p_x, p_y)$  we compute the horizontal derivative at  $\phi(p_x + r_x, p_y + r_y)$  as central difference

$$\frac{\| (p_x + r_x + 1, p_y + r_y) \| - \| (p_x + r_x - 1, p_y + r_y) \|}{2}.$$

The vertical derivative is computed analogously. Like the related work, we perform all computations on three pyramid levels in a coarse-to-fine manner and shoot the rays in a band of 8 steps on each level. Since we shoot two rays per

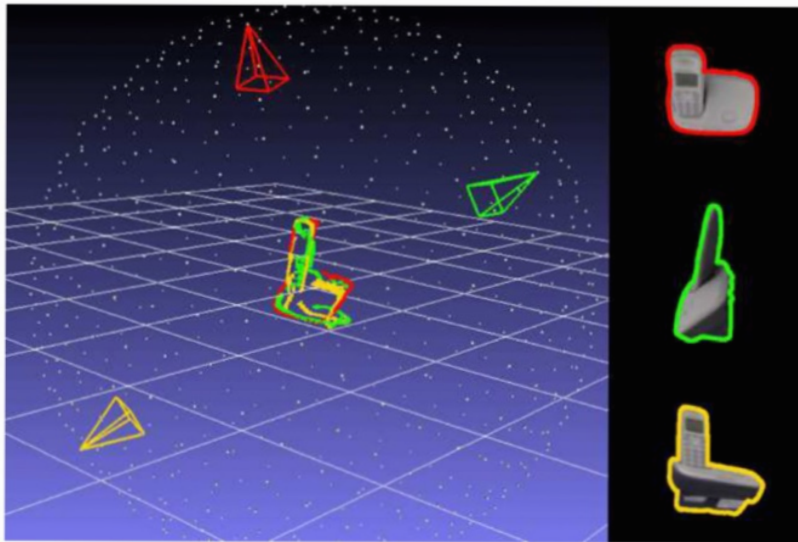


Figure 4. Object-local 3D contour points visualized for three view-points on the unit sphere. Each view captures a different contour which is used during tracking to circumvent costly renderings.

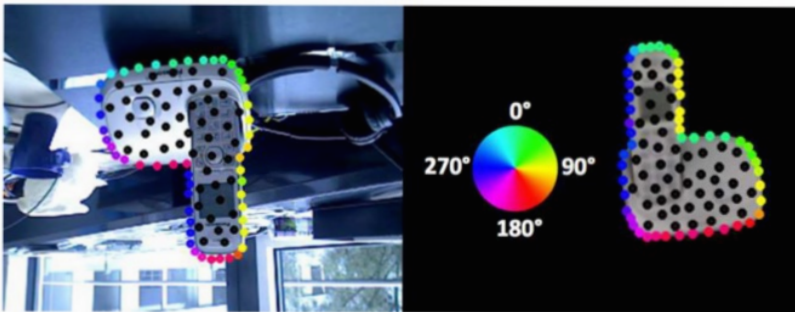
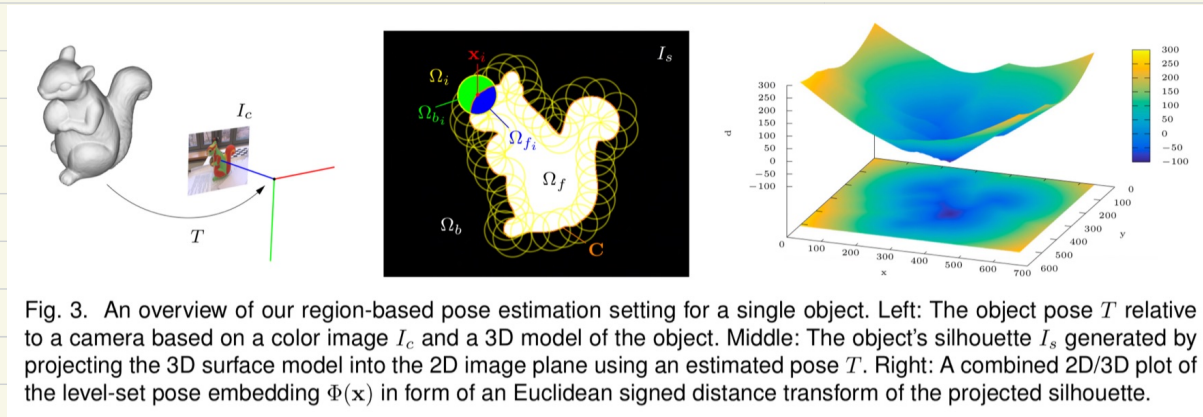


Figure 5. Current tracking and closest pre-rendered viewpoint augmented with contour and interior sampling points. The hue represents the normal orientation for each contour point. Note how we rotate the orientation of each contour point by our approximation of the inplane rotation such that the SDF computation is proper.

## A Region-Based Gauss-Newton Approach to Real-Time Monocular Multiple Object Tracking (2019)



## In 2D-3D Pose Estimation of Heterogeneous Objects Using a Region Based Approach (2016)

The local region centers  $x_i$  were radiated as arbitrary sets of pixel locations along  $C$  for each image. Thus, this approach in general does not allow to establish correspondences of centers across multiple frames ( $x_i(t_k) \leftrightarrow x_i(t_{k-1})$ ) which is required in order to update the respective histograms.

**Tdc-histograms**: projecting all model vertices into image plane,  $x_i = \pi(K(T\tilde{x}_i))_{3 \times 1}$  and selecting the subset of all  $x_i \in C$ . Center correspondences between frames:

$$\pi(K(T(t_k)\tilde{x}_i))_{3 \times 1} \leftrightarrow \pi(K(T(t_{k-1})\tilde{x}_i))_{3 \times 1}$$

This ensures to keep the individual histograms temporally consistent.

- ① If a vertex projects onto the contour for the first time, its corresponding histograms are initialed from the local region around its center in current frame.
- ② Otherwise:

$$P(y|M_{f_i}) = (1 - a_f) P^{t-1}(y|M_{f_i}) + a_f P^t(y|M_{f_i})$$

$$P(y|M_{b_i}) = (1 - a_b) P^{t-1}(y|M_{b_i}) + a_b P^t(y|M_{b_i})$$

More robust results can be obtained by computing the average posteriors from all local histograms instead:

$$\bar{P}_f(x) = \frac{1}{\sum_{i=1}^n B_i(x)} \sum_{i=1}^n P_{f_i}(x) B_i(x), \quad B_i(x) = \begin{cases} 1, & \forall x \in \Omega_i \\ 0, & \forall x \notin \Omega_i \end{cases}$$

$$\bar{P}_b(x) = \frac{1}{\sum_{i=1}^n B_i(x)} \sum_{i=1}^n P_{b_i}(x) B_i(x)$$



$$\therefore E(\xi) = - \sum_{x \in \Omega} \log(\text{He}(\Phi(x(\xi))) \bar{p}_f(x) + (1 - \text{He}(\Phi(x(\xi)))) \bar{p}_b(x)).$$

### Gauss-Newton Strategy

$$E(\xi) = \sum_{x \in \Omega} F(x, \xi), \quad F(x, \xi) = -\log(\text{He}(\Phi(x(\xi))) \bar{p}_f(x) + (1 - \text{He}(\Phi(x(\xi)))) \bar{p}_b(x))$$

rewrite it as a nonlinear weighted least-squares problem of the form :

$$E(\xi) = \frac{1}{2} \sum_{x \in \Omega} \frac{1}{F(x, \xi)} F^2(x, \xi) = \sum_{x \in \Omega} \psi(x) F^2(x, \xi)$$

fixed weights  $\psi(x)$  by means of Gauss-Newton opt and alternatively using the refined pose for updating the weights  $\psi(x)$ .

$$\therefore \frac{\partial E(\xi)}{\partial \xi} = \frac{1}{2} \sum_{x \in \Omega} \psi(x) \frac{\partial F^2(x, \xi)}{\partial \xi} = \sum_{x \in \Omega} \psi(x) F \frac{\partial F}{\partial \xi} = \sum_{x \in \Omega} \frac{\partial F}{\partial \xi}$$

with  $\psi(x) F(x, \xi) = 1$ .

$$\therefore \frac{\partial^2 E(\xi)}{\partial \xi^2} = \sum_{x \in \Omega} \psi(x) \left( \left( \frac{\partial F}{\partial \xi} \right)^T \frac{\partial F}{\partial \xi} + F \frac{\partial^2 F}{\partial \xi^2} \right)$$

$$\text{Let } \frac{\partial F(x, \xi)}{\partial \xi} = J$$

$$\therefore E(\xi + \Delta \xi) \approx E(\xi) + J \Delta \xi$$

$$\therefore \frac{1}{2} \| E(\xi) + J \Delta \xi \|^2 = \frac{1}{2} (\| E(\xi) \|^2 + 2 E(\xi)^T \sum_{x \in \Omega} J \Delta \xi + \Delta \xi^T \sum_{x \in \Omega} \psi(x) \Delta \xi^T J^T J \Delta \xi)$$

$$\text{Let } \frac{\partial (\frac{1}{2} \| E(\xi) + J \Delta \xi \|^2)}{\partial \Delta \xi} = 0$$

$$\therefore J^T J \Delta \xi = - J^T E(\xi)$$

$$\therefore \Delta \xi = (J^T J)^{-1} J^T E(\xi)$$

#### 4.1 Rendering Engine

We use the standard rasterization pipeline of OpenGL in order to obtain the silhouette masks  $I_s$ . Since we want to process the rendered images on the CPU, we perform offscreen rendering into a *FrameBufferObject*, which is afterwards downloaded to host memory. To generate synthetic views that match the real images, the intrinsic parameters (2) of the camera need to be included. For this, we model the transformation from 3D model coordinates  $\mathbf{X}$  to homogeneous coordinates within the canonical view volume of OpenGL as  $\tilde{\mathbf{V}} = P(K)LT\tilde{\mathbf{X}} \in \mathbb{R}^4$ . Here,  $L$  is the so-called *look-at* matrix

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4}, \quad (36)$$

that aligns the principal axes of the real camera coordinate frame with those of the virtual OpenGL camera and  $P(K)$  is a homogeneous projection matrix

$$P(K) = \begin{bmatrix} \frac{2f_x}{w} & 0 & 1 - \frac{2c_x}{w} & 0 \\ 0 & -\frac{2f_y}{h} & 1 - \frac{2c_y}{h} & 0 \\ 0 & 0 & -\frac{Z_f + Z_n}{Z_f - Z_n} & -\frac{2Z_f Z_n}{Z_f - Z_n} \\ 0 & 0 & -1 & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4}, \quad (37)$$

with respect to the camera matrix  $K$ . The scalars  $w$ ,  $h$  are the width and height of the real image  $I_c$  and  $Z_n$ ,  $Z_f$  are the near- and far-plane of the view frustum described by  $P(K)$ .

In case of tracking multiple objects, all 3D models are rendered in the same scene. Each mesh is rendered with a constant and unique color that corresponds to its model index  $j$ . This allows to separate the individual foreground regions

and identify their contours  $\mathbf{C}^j$  as required for computing the different level-sets. Here, mutual occlusions are natively handled by the OpenGL *Z-Buffer*.

As seen in (34), the derivatives used for pose optimization involve the coordinates of the 3D surface point  $\mathbf{X}'$  in the camera's frame of reference, corresponding to each pixel  $\mathbf{x}$ . In addition to the silhouette mask, we therefore also download the *Z-buffer* into a per pixel depth map  $I_d : \Omega \rightarrow [0, 1] \subset \mathbb{R}$ . Given  $I_d$ , the required coordinates are efficiently determined via backprojection as  $\mathbf{X}'(\mathbf{x}, I_d) = D(\mathbf{x}, I_d)K^{-1}\tilde{\mathbf{x}}$ , with

$$D(\mathbf{x}, I_d) = \frac{Z_n Z_f}{Z_f - I_d(\mathbf{x})(Z_f - Z_n)}, \quad \forall \mathbf{x} \in \Omega_f, \quad (38)$$

where  $\tilde{\mathbf{x}} = [x, y, 1]^\top$  is the homogeneous representation of an image point  $\mathbf{x} = [x, y]^\top$ .

In (26) it has been shown that it is beneficial not only to consider the points on the surface closest to the camera but also the most distant ones (on the backside of the object) for pose optimization. In order to obtain the respective coordinates for each pixel, we compute an additional reverse depth map  $I_d^r$ , for which we simply invert the OpenGL depth check used to compute the corresponding *Z-buffer* (see Figure 6). Given  $I_d^r$ , the farthest surface point  $\mathbf{X}'(\mathbf{x}, I_d^r)$  corresponding to a pixel  $\mathbf{x}$  is also recovered as  $\mathbf{X}'(\mathbf{x}, I_d^r) = D(\mathbf{x}, I_d^r)K^{-1}\tilde{\mathbf{x}}$ .

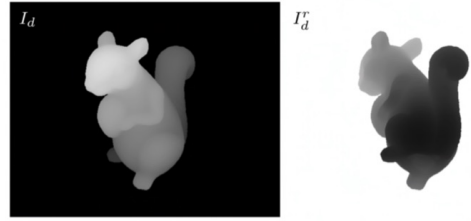


Fig. 6. The two depth map types used within our approach, where brighter pixels are closer to the camera. Left: The usual depth map  $I_d$  corresponding to the closest surface points. Right: The reverse depth map  $I_d^r$  corresponding to the most distant surface points.

## Pixel-Wise Weighted Region-Based 3D Object Tracking using Contour Constraints (2021)

By projecting 3D model into image plane with a given pose, we obtain the silhouette mask  $I_s$ , in which the projected contour segments the image into the foreground region  $\Omega_f$  and the background region  $\Omega_b$ . At each projected contour point  $m_i$ , the search line  $l_i$  is sampled along the direction vector  $n_i$  which is perpendicular to the projected contour. The  $j$ -th sampling pixel of the search line  $l_i$  is denoted by  $x_{ij}$ .  $s_i$  is the object contour point. The object pose that transforms a 3D model point in the object coordinate to the camera coordinate is represented by a  $4 \times 4$  homogeneous matrix

$${}_{model}^{cam}T = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}$$

$$x = \pi(K(T^{model}X))$$

$\Delta T$  is parameterized by  $p = [w_1, w_2, w_3, v_1, v_2, v_3]^T$

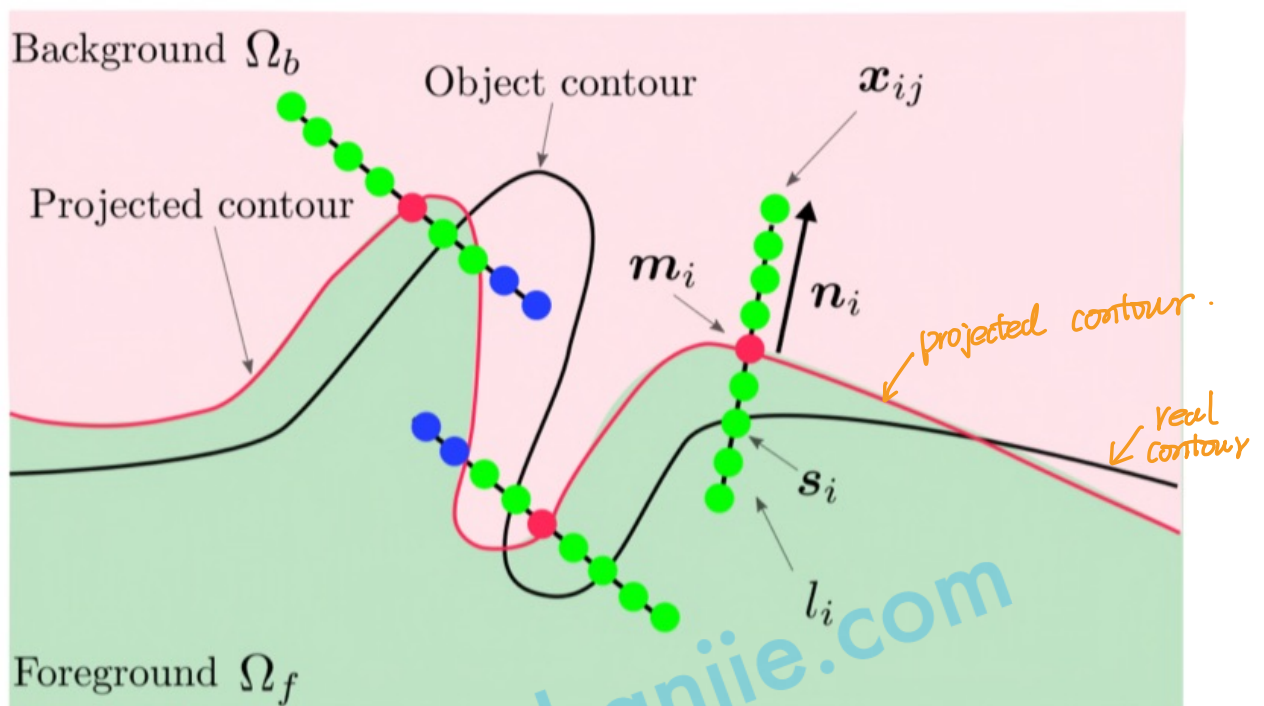


Fig. 2. The silhouette mask  $I_s$  generated by projecting the 3D model into the image plane.

$$E(p) = - \sum_{x_{ij} \in L} w(x_{ij}) \log [He(d(x_{ij})) P_f(x_{ij}) + (1 - He(d(x_{ij}))) P_b(x_{ij})]$$

$$\textcircled{1} p = [w_1, w_2, w_3, v_1, v_2, v_3]^T$$

$\textcircled{2} x_{ij}$  : The  $j$ -th sampling pixel of the search line  $l_i$

$\textcircled{3} L$  : all sampling pixels  $x_{ij}$

$$\textcircled{4} d(x_{ij}) = n_i^T (x_{ij} - m_i),$$

$m_i$  : projected contour point.

$n_i$  : perpendicular to the projected contour, pointing to background.

$$d(x_{ij}) \begin{cases} > 0, & x_{ij} \in \Omega_b \\ < 0, & x_{ij} \in \Omega_f \end{cases}$$

$$\textcircled{5} He(d(x_{ij})) = \frac{1}{\pi} (-\arctan(d(x_{ij})s) + \frac{\pi}{2})$$

$s$  controls the pitch of the smoothed transition.

$\textcircled{6} P_f(x_{ij})$  and  $P_b(x_{ij})$  : **representing the foreground and background region membership posterior probability of each pixel's color, which are calculated by temporally consistent local color histograms**

Candidate Contour Point:

**we build a bundle image  $I_b$  by simply stacking each search line. As shown in Fig. 3, each row of  $I_b$  is the sampling pixels  $x_{ij}$  of the search line  $l_i$ .**

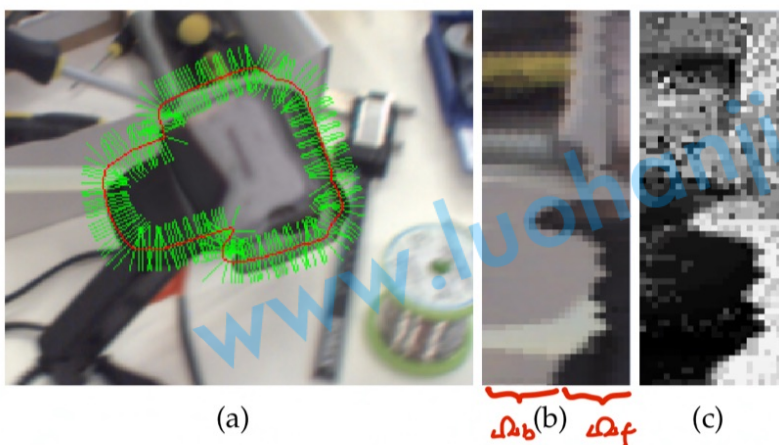


Fig. 3. Example of the search lines and the bundle image. (a) The search lines (green lines) are sampled around the projected contour (red curve). (b) Part of the bundle image  $I_b$ , which is built by stacking the search lines. (c) Foreground probability map  $I_p$  of the bundle image  $I_b$ .

we assume that the object contour point  $s_i$  is in the location where the foreground probability changes rapidly. Then the candidate point  $h_{ij}$  can be computed by 1D convolution of a  $1 \times 3$  filter kernel  $f = [-1, 0, 1]$  at each row of  $I_p$ .

$$h_{ij} = (f * I_p)(l_{ij}) > \epsilon$$

$l_{ij}$  is the location in  $I_p$ ,

The candidate contour point  $h_{ij}$  can be extracted if the convolution value is larger than the threshold.

### Probability of Contour Point

For a sampling point  $x_{ij}$ , the left area  $\Phi_{x_{ij}}^-$  is defined as  $(x_{i,j-1}, x_{i,j-2}, x_{i,j-3})$   
 right  $\Phi_{x_{ij}}^+$   $(x_{i,j+1}, x_{i,j+2}, x_{i,j+3})$

$$P(h_{ij} | C) = \prod_{x_{ij} \in \Phi_{h_{ij}}^-} P_b(x_{ij}) \prod_{x_{ij} \in \Phi_{h_{ij}}^+} P_f(x_{ij})$$

the probability that the candidate point  $h_{ij}$  belongs to the contour

$$P(h_{ij} | F) = \prod_{x_{ij} \in \Phi_{h_{ij}}^-} P_f(x_{ij}) \prod_{x_{ij} \in \Phi_{h_{ij}}^+} P_f(x_{ij})$$

$h_{ij}$  belongs to the foreground clutter point

$$P(h_{ij} | B) = \prod_{x_{ij} \in \Phi_{h_{ij}}^-} P_b(x_{ij}) \prod_{x_{ij} \in \Phi_{h_{ij}}^+} P_b(x_{ij})$$

$h_{ij}$  belongs to the background clutter point

The normalized probability that the candidate point  $h_{ij}$  belongs to the object contour point:

$$P_c(h_{ij}) = \frac{P(h_{ij} | C)}{P(h_{ij} | C) + P(h_{ij} | F) + P(h_{ij} | B)}$$

### Weight Function using Contour Constraints

$$W_c(x_{ij}) = \begin{cases} \exp(k_1 (1 - P(s_i | C))), & \text{if } s_i \neq \phi \\ \exp(k_1), & \text{otherwise} \end{cases}$$

If  $P(s_i | C)$  is small, the search line  $l_i$  is likely to be influenced by partial occlusions or ambiguous colors, then we reduce the weights of all sampling pixels  $x_{ij}$  on the search line  $l_i$ .



where  $k_1$  is a negative constant that controls the rate of decay, so that the weight function  $w_c(x_{ij})$  decays exponentially with  $1 - P(s_i | C)$ .

$$w_d(x_{ij}) = \begin{cases} \exp(k_2 D(x_{ij}, s_i)) & , \text{if } s_i \neq \emptyset \\ \exp(k_2) & , \text{otherwise} \end{cases}, \quad D(x_{ij}, s_i) = \|x_{ij} - s_i\| / N_c$$

length of search line  $l_i$

Since sampling pixels far from the object contour are more likely to be influenced by partial occlusions and cluttered backgrounds, we reduce the weights of the sampling pixels far from the object contour.

$$w(x_{ij}) = w_c(x_{ij}) w_d(x_{ij})$$

Pose Optimization

$$\begin{aligned} E(p) &= - \sum_{x_{ij} \in L} w(x_{ij}) \log [He(d(x_{ij})) P_f(x_{ij}) + (1 - He(d(x_{ij}))) P_b(x_{ij})] \\ &= \frac{1}{2} \sum_{x_{ij} \in L} w(x_{ij}) \psi(x_{ij}) F^2(x_{ij}, p) \end{aligned}$$

where  $\psi(x_{ij}) = \frac{1}{F(x_{ij}, p)}$ ,  $F(x_{ij}, p) = -\log[He(d(x_{ij})) P_f(x_{ij}) + (1 - He(d(x_{ij}))) P_b(x_{ij})]$

$$J = \frac{\partial F(x_{ij}, p)}{\partial He(d(x_{ij}))} \frac{\partial He(d(x_{ij}))}{\partial p}$$

$$\frac{\partial F(x_{ij}, p)}{\partial He(d(x_{ij}))} = \frac{P_f - P_b}{He(d(x_{ij})) P_f + (1 - He(d(x_{ij}))) P_b}$$

$$\frac{\partial He(d(x_{ij}))}{\partial p} = \frac{\partial He}{\partial d(x_{ij})} \frac{\partial d(x_{ij})}{\partial m_i} \frac{\partial m_i}{\partial p} = \delta_d(d(x_{ij})) n_i^T \frac{\partial m_i}{\partial p}$$

$$\delta_d(d(x_{ij})) = \frac{\partial He}{\partial d(x_{ij})} = \frac{S}{\pi (H d(x_{ij})^2 S^2)}$$

$$\therefore \Delta p = - \left( \sum_{x_{ij} \in L} w(x_{ij}) \psi(x_{ij}) J^T J \right)^T \sum_{x_{ij} \in L} w(x_{ij}) J^T$$

TABLE 3

Average runtime (in ms) of each step for one frame on RBOT dataset.  
Please see the text for more explanations.

Method	RD	PP	SL	CS	DT	PO	HU	Total
[25]	18.5	4.6	-	-	2.1	2.2	6.2	33.6
Ours	18.6	4.4	1.5	0.7	-	1.9	6.2	33.3

RD: rendering

PP: pixel-wise posterior probability calculation

SL: search lines sampling

CS: object contour points searching

DT: signed distance transform

PO: pose optimization

HU: histogram up-date

# A Sparse Gaussian Approach to Region-Based 6DOF object Tracking (2021)

$X_i = [X_i, Y_i, Z_i]^T$ , 3D model points

$\tilde{X}_i = [X_i, Y_i, Z_i, 1]^T$

$I: \Omega \rightarrow \{0, \dots, 255\}^3$ , color image  $I$  with image domain  $\Omega \subset \mathbb{R}^2$ .

$y_i = I(x_i)$ , color values  $y_i$  at image  $x_i = [x_i, y_i]^T$

$$x_i = \pi(X_i) = \begin{bmatrix} X_i f_x / Z_i + p_x \\ Y_i f_y / Z_i + p_y \end{bmatrix}$$

$$\overset{\text{camera ref frame}}{\uparrow} {}^c \tilde{X}_i = {}^c_M \tilde{X}_i = \begin{bmatrix} {}^c_M R & {}^c_M t \\ 0 & 1 \end{bmatrix} \overset{\text{model ref frame}}{\uparrow} {}^m \tilde{X}_i$$

$$\therefore R = \exp([R]_x) = I + [R]_x + \frac{1}{2!} [R]_x^2 + \frac{1}{3!} [R]_x^3 + \dots$$

$$[R]_x = \begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix}_x = \begin{bmatrix} 0 & r_3 & -r_2 \\ r_3 & 0 & -r_1 \\ -r_2 & r_1 & 0 \end{bmatrix}$$

The linear variation of a 3D model point represented in Camera reference frame  $C$ :

$${}^c \tilde{X}_i^+ = \begin{bmatrix} {}^c_M R & {}^c_M t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I + [R]_x & \theta_t \\ 0 & 1 \end{bmatrix} {}^m \tilde{X}_i$$

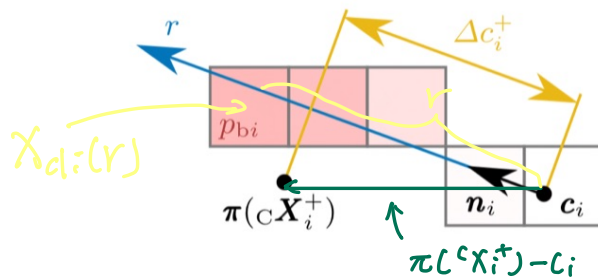
$$\theta^T = [\theta_r^T \quad \theta_t^T]$$

It is more natural to rotate a model point (右乘) in the model frame since object is typically moved significantly more than the camera.

A **correspondence line** is described by a center  $c_i = [c_x, c_y]^T \in \mathbb{R}^2$  and a normal vector  $n_i = [n_x, n_y]^T \in \mathbb{R}^2$ , with  $\|n_i\|_2 = 1$  in image.

3D colour point  $X_i$  and associated vector normal to Contour  $N_i$   $\xrightarrow[\text{image}]{\text{Project to}}$   $c_i, n_i$

pixels on the line :  $x_{cl_i}(r) = \lfloor c_i + rn_i + 0.5 \rfloor$   
 $\uparrow$   
distance from center  $c_i$



**Fig. 2.** Correspondence line defined by a center  $c_i$  and a normal vector  $n_i$ , with selected pixels and the projected difference  $\Delta c_i^+$  from  $c_i$  to the varied point  ${}_cX_i^+$ . The color intensity in red indicates the magnitude of the pixel-wise posterior  $p_{bi}$  for each pixel.

During the pose variation in bDof, the projected difference  $\Delta c_i^+$  from the **unmoved** center  $c_i$  to the **varriated** model point  ${}_cX_i^+$  is calculated as:

$$\Delta c_i^+ = n_i^T (\pi({}_cX_i^+) - c_i)$$

The color on the line with distance  $r$ :

$$y_i(r) = I(x_{cl_i}(r)) \quad x_{cl_i}(r) = \lfloor c_i + rn_i + 0.5 \rfloor$$

The pixel-wise posterior, of the model  $m_j$ , given a pixel value  $y_i(r)$ :

$$p_j(r) = P(m_j | y_i(r)) = \frac{P(y_i(r) | m_j) P(m_j)}{P(y_i(r))} = \frac{P(y_i(r) | m_j) P(m_j)}{P(y_i(r) | m_f) P(m_f) + P(y_i(r) | m_b) P(m_b)}, \quad j \in \{f, b\}$$

$\uparrow$   
 $m_f$ : foreground model  
background model

$$= \frac{P(y_i(r) | m_j)}{P(y_i(r) | m_f) + P(y_i(r) | m_b)}$$

Using the knowledge that foreground and background are equally likely along the correspondence line, i.e.  $p(m_f) = p(m_b)$ ,

$$p_j(y) = \frac{P(m_j | y)}{\eta_j} = \frac{P(y | m_j) \frac{\eta_j}{\sum_{i=f,b} P(y | m_i) \eta_i}}{P(y | m_f) \eta_f + P(y | m_b) \eta_b} = \frac{P(y | m_j) \eta_j}{P(y | m_f) \eta_f + P(y | m_b) \eta_b}$$

$\uparrow$  foreground pixels number       $\uparrow$  background pixels number

Define  $p_j(r) = P(m_j | y_i(r)) = \frac{P(y_i(r) | m_j) \eta_j}{P(y_i(r) | m_f) \eta_f + P(y_i(r) | m_b) \eta_b}$

$\uparrow$  model       $\uparrow$  pixel

Remap  $r \in \mathbb{R}_i \longrightarrow r_s$  in a discrete space:

$$\Delta C_i^+ = n_i^T (\pi(C_i^+) - c_i) \longrightarrow \Delta C_s^+$$

$$r_s = (r - \Delta r_i) \frac{\bar{n}_i}{S}, \quad \Delta C_s^+ = (\Delta C_i^+ - \Delta r_i) \frac{\bar{n}_i}{S}$$

$S \in \mathbb{N}^+$  the scale describes the number of pixels combined into a segment.

$\bar{n}_i = \max(|n_{x_i}|, |n_{y_i}|)$  projects a normal vector  $n_x$  to the closest horizontal or vertical coordinate.

$\Delta r_i \in \mathbb{R}$  distance from center  $c_i$  to defined segment location (chosen such that the center  $r_s = 0$  lies on the border between two segments).

The likelihood function in scale-space:

$$P(D_i | \Delta \tilde{C}_i) \propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{C}_i) P_{f_i}(r_s) + h_b(r_s - \Delta \tilde{C}_i) P_{b_i}(r_s))$$

$\uparrow$   
 $\Delta \tilde{C}_i$  is a function of  $\theta$  (pose)

$$P(D_i | \theta) \propto \prod_{r \in \mathbb{R}_i} (h_f(r - \Delta C_i^+) P_{f_i}(r) + h_b(r - \Delta C_i^+) P_{b_i}(r))$$

$R_s$  a set of distance to segment centers,  $\Delta \tilde{C}_i$  has to be chosen such that precalculated values aligned with segment centers.

For  $h_f(x), h_b(x)$ , 10 precomputed values  $x \in \{-4.5, -3.5, -2.5, -1.5, -0.5, 0.5, 1.5, 2.5, 3.5, 4.5\}$  are used. X function.

先将  $r \rightarrow r_s$ , 然后通过 offset  $\Delta \tilde{C}_i$  使得 precalculated values 对齐于 segment centers.

$P_{f_i}$  and  $P_{b_i}$  segment-wise posteriors. Assuming pixel-wise independence:

$r_s \xrightarrow{S(r_s)} \{r\}$

$$P_{sji}(r_s) = \frac{\prod_{r \in S(r_s)} P(y_i(r) | m_j)}{\prod_{r \in S(r_s)} P(y_i(r) | m_f) + \prod_{r \in S(r_s)} P(y_i(r) | m_b)}, \quad j \in \{f, b\}$$

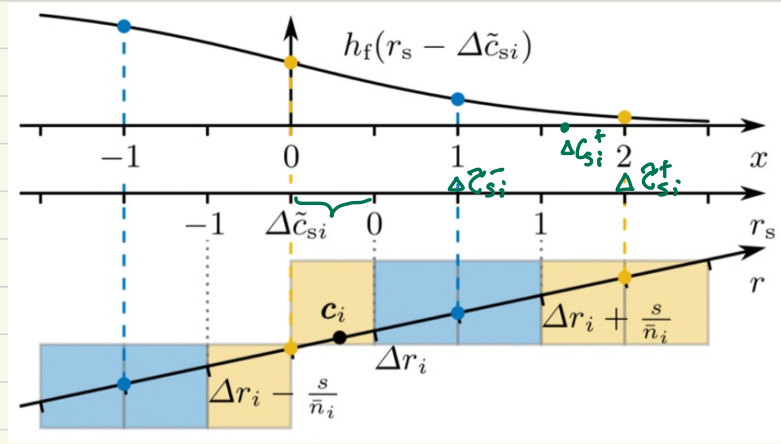
all  $r$  belongs to  $r_s$

$$P_{ji}(r) = \frac{P(m_j | y_i(r))}{\eta_j} = \frac{P(y_i(r) | m_j)}{P(y_i(r) | m_f) \eta_f + P(y_i(r) | m_b) \eta_b}$$

where  $S$  is a set-valued function that maps  $r_s$  to a set of values  $r$  that describe the  $S$  (scale) closest pixel centers of a segment.

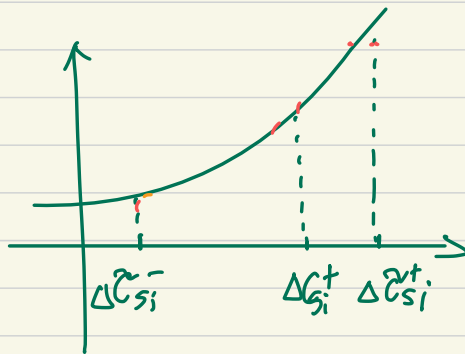


Due to a limited number of precalculated values for  $h_f$  and  $h_b$ ,  $P(D_i | \Delta \tilde{c}_{si})$  can only be evaluated at discrete values  $\Delta \tilde{c}_{si}$ , the upper and lower neighboring discretized values  $\Delta \tilde{c}_{si}^+$  and  $\Delta \tilde{c}_{si}^-$  are used to linearly interpolate:  
of  $\Delta c_{si}$



$$P(D_i | \theta) \propto (\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i | \Delta \tilde{c}_{si}^-) + (\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i | \Delta \tilde{c}_{si}^+)$$

$$P(D_i | \Delta \tilde{c}_{si}) \propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{c}_{si}) P_{sf_i}(r_s) + h_b(r_s - \Delta \tilde{c}_{si}) P_{sb_i}(r_s))$$



$$\Delta \tilde{c}_{si}^+ = (\Delta c_{si}^+ - \Delta r_i) \frac{\bar{n}_i}{S}$$

$$P(D_i | \Delta c_{si}^+) = P(D_i | \Delta \tilde{c}_{si}^-) + (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-) \frac{(P(D_i | \Delta \tilde{c}_{si}^+) - P(D_i | \Delta \tilde{c}_{si}^-))}{(\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-)}$$

$$= \frac{(\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-) - (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-)}{(\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-)} P(D_i | \Delta \tilde{c}_{si}^-)$$

$$+ \frac{(\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-)}{(\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-)} P(D_i | \Delta \tilde{c}_{si}^+)$$

$$= \frac{\Delta \tilde{c}_{si}^+ - \Delta c_{si}^+}{\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-} P(D_i | \Delta \tilde{c}_{si}^-) + \frac{\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-}{\Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-} P(D_i | \Delta \tilde{c}_{si}^+)$$

$\therefore \Delta \tilde{c}_{si}^+ - \Delta \tilde{c}_{si}^-$  is a fixed value, it can be dropped.

## Gaussian Equivalence

Newton optimization yields good results for Gaussian distributions.

The goal is to find smoothed step function  $h_f$  and  $h_b$  that ensure that the likelihood follows a Gaussian distribution. we want the first-order derivative of the logarithm of the likelihood has to be equal to that of a normal distribution.

$$P(\mathcal{D}_i | \theta) \propto \prod_{r \in \mathcal{R}_i} (h_f(r - \Delta C_i^+) P_{fi}(r) + h_b(r - \Delta C_i^+) P_{bi}(r))$$

$$P(\mathcal{D}, p | \Omega) \quad P(x_i | \mathcal{D}, p, M_f) \quad P(M_f | y_i)$$

$$\theta^T = [\theta_r^T \quad \theta_b^T]$$

①  $\mathcal{D}_i$ : the data specific to a single correspondence line,

②  $\mathcal{R}_i$ : a set distance  $r$  from the line center to pixel centers.

③  $P_{fi}$  and  $P_{bi}$ : the pixel-wise posteriors for foreground and background.

④  $\Delta C_i^+$ : the projected difference from correspondence line center  $C_i$  to the warped model point  ${}^c\tilde{X}_i^+$

$${}^c\tilde{X}_i^+ = \begin{bmatrix} {}^cR & {}^c t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + [\theta_r]_x & \theta_b \\ 0 & 1 \end{bmatrix} {}^M\tilde{X}_i, \quad C_i = \pi {}^cX_i$$

$$\Delta C_i^+ = n_i^T (\pi ({}^c\tilde{X}_i^+) - C_i)$$

For  $P_{fi}$  and  $P_{bi}$ , perfect segmentation (分割, 还不算 pose) and a contour at the correspondence line center assumed.

$$P_{fi}(r) = \begin{cases} 1 & \text{if } r \leq 0 \\ 0 & \text{else} \end{cases}, \quad P_{bi}(r) = \begin{cases} 0 & \text{if } r \leq 0 \\ 1 & \text{else} \end{cases}$$

Suppose:

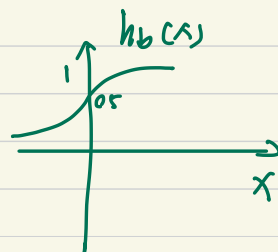
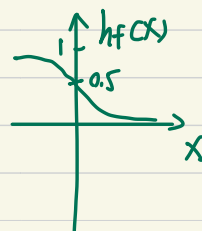
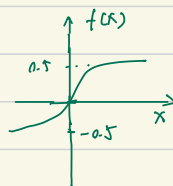
$$\textcircled{1} \quad h_f + h_b = 1$$

$$\textcircled{2} \quad h_f \text{ and } h_b \text{ be symmetric.}$$

$$\therefore \quad h_f(x) = 0.5 - f(x), \quad h_b(x) = 0.5 + f(x)$$

where  $f(x)$  is an odd function ( $-f(x) = f(-x)$ ) that lies within the interval  $[-0.5, 0.5]$  and

$$\lim_{x \rightarrow \infty} f(x) = 0.5 \quad \text{and} \quad \lim_{x \rightarrow -\infty} f(x) = -0.5$$





Infinitesimally small pixels are assumed to write the likelihood from in continuous form (Product Integral)

$$P(D_i|\theta) \propto \prod_{r=-\infty}^{\infty} (h_f(r - \Delta C_i^+) P_{fi}(r) + h_b(r - \Delta C_i^+) P_{bi}(r))^{dr}$$

Convert the product integral to Riemann integral in the perfect pixel-wise posteriors assumption:

$$\prod_a^b f(x) dx = \lim_{\Delta x \rightarrow 0} \prod f(x_i) \Delta x = \exp\left(\int_a^b \ln f(x) dx\right)$$

$$\begin{aligned} P(D_i|\theta) &\propto \exp\left(\int_{r=-\infty}^{\infty} \ln(h_f(r - \Delta C_i^+) P_{fi}(r) + h_b(r - \Delta C_i^+) P_{bi}(r)) dr\right) \\ &= \exp\left(\int_{r=-\infty}^{\infty} \ln(h_f(r - \Delta C_i^+) P_{fi}(r)) dr + \int_{r=-\infty}^{\infty} \ln(h_b(r - \Delta C_i^+) P_{bi}(r)) dr\right) \end{aligned}$$

$$P_{fi}(r) = \begin{cases} 1 & \text{if } r \leq 0 \\ 0 & \text{else} \end{cases}, \quad P_{bi}(r) = \begin{cases} 0 & \text{if } r \leq 0 \\ 1 & \text{else} \end{cases}$$

$$= \exp\left(\int_{r=-\infty}^0 \ln(h_f(r - \Delta C_i^+)) dr + \int_{r=0}^{\infty} \ln(h_b(r - \Delta C_i^+)) dr\right)$$

$$\text{let } r - \Delta C_i^+ = x \quad \therefore \int_{r=-\infty}^0 \Rightarrow \int_{x=-\infty}^{-\Delta C_i^+}, \quad dr = d(x + \Delta C_i^+) = dx$$

$$= \exp\left(\int_{x=-\infty}^{-\Delta C_i^+} \ln(h_f(x)) dx + \int_{x=-\Delta C_i^+}^{\infty} \ln(h_b(x)) dx\right)$$

To eliminate both constant scaling factors and the integral :

$$\frac{\partial \ln(P(D_i|\theta))}{\partial \Delta C_i^+} = \frac{\partial \ln \left[ \exp\left(\int_{x=-\infty}^{-\Delta C_i^+} \ln(h_f(x)) dx + \int_{x=-\Delta C_i^+}^{\infty} \ln(h_b(x)) dx\right) \right]}{\partial \Delta C_i^+}$$

$$= \frac{\partial \left( \int_{x=-\infty}^{-\Delta C_i^+} \ln(h_f(x)) dx + \int_{x=-\Delta C_i^+}^{\infty} \ln(h_b(x)) dx \right)}{\partial \Delta C_i^+}$$

$$\Rightarrow \frac{\partial \left( \int_{x=-\infty}^{-\Delta C_i^+} \ln(h_f(x)) dx \right)}{\partial \Delta C_i^+}$$

**Leibniz Rule:** If  $\Phi(t) = \int_{a(t)}^{b(t)} f(x, t) dx$ , then

$$\frac{d\Phi}{dt} = \int_{a(t)}^{b(t)} \frac{\partial f(x, t)}{\partial t} dx + f(b(t), t) \frac{db}{dt} - f(a(t), t) \frac{da}{dt}$$

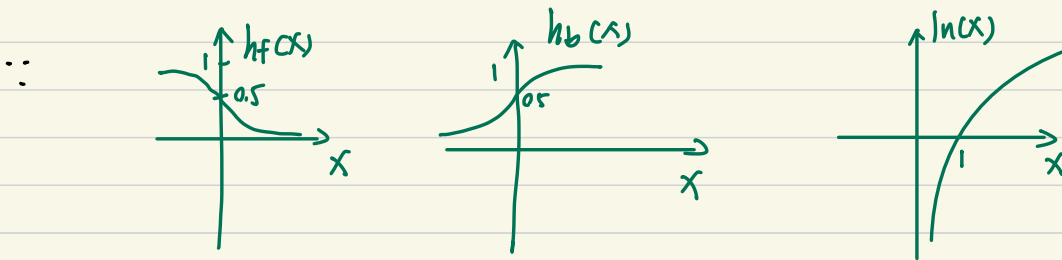
$$= \int_{x=-\infty}^{-\Delta C_i^+} \frac{\partial \ln(h_f(x))}{\partial \Delta C_i^+} dx - \ln(h_f(-\Delta C_i^+))$$

$$= -\ln(h_f(-\Delta C_i^+))$$

$$\frac{\partial \left( \int_{x=-\Delta C_i^+}^{\infty} \ln(h_b(x)) dx \right)}{\partial \Delta C_i^+}$$

$$= 0 + 0 + \ln(h_b(-\Delta C_i^+))$$

$$\therefore \frac{\partial \ln(P(D; \theta))}{\partial \Delta C_i^+} = -\ln(h_f(-\Delta C_i^+)) + \ln(h_b(-\Delta C_i^+))$$



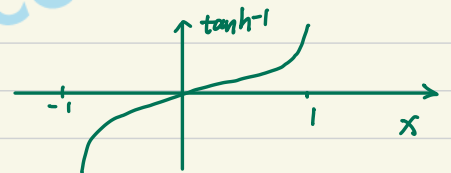
$$\therefore \lim_{x \rightarrow -\infty} \ln(h_f(x)) = 0, \quad \lim_{x \rightarrow \infty} \ln(h_b(x)) = 0$$

$$\therefore h_f(x) = 0.5 - f(x), \quad h_b(x) = 0.5 + f(x)$$

$$\begin{aligned} \therefore \frac{\partial \ln(P(D; \theta))}{\partial \Delta C_i^+} &= -\ln(0.5 - f(-\Delta C_i^+)) + \ln(0.5 + f(-\Delta C_i^+)) \\ &= -\ln(0.5(1 - 2f(-\Delta C_i^+))) + \ln(0.5(1 + 2f(-\Delta C_i^+))) \\ &= -\ln(1 - 2f(-\Delta C_i^+)) - \ln 0.5 + \ln(1 + 2f(-\Delta C_i^+)) + \ln 0.5 \\ &= -\ln(1 - 2f(-\Delta C_i^+)) + \ln(1 + 2f(-\Delta C_i^+)) \end{aligned}$$

$$\therefore \text{The inverse hyperbolic tangent: } 2 \tanh^{-1}(z) = -\ln(1-z) + \ln(1+z)$$

$$\frac{\partial \ln(P(D; \theta))}{\partial \Delta C_i^+} = 2 \tanh^{-1}(2f(-\Delta C_i^+))$$



We want to enforce our likelihood function follows a normal distribution  $X \sim N(\mu, \sigma^2)$

$$\text{The PDF: } f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\therefore \text{let } \Delta C_i^+ \sim N(0, S_h) \propto \exp\left(-\frac{1}{2S_h} (\Delta C_i^+)^2\right)$$

We attain a probabilistic model that follows a Gaussian distribution. Although assumptions, such as perfect pixel-wise posteriors and infinitesimally small pixels, are not an exact description of reality. Experiments demonstrate that we achieve excellent convergence for used regularized Newton optimization.

## Regularized Newton Method

To maximize the likelihood, we estimate the variation vector  $\theta$  and iteratively update the pose. We use Newton optimization with Tikhonov regularization:

$$\hat{\theta} = [-H + \begin{bmatrix} \lambda_r I_3 & 0 \\ 0 & \lambda_t I_3 \end{bmatrix}]^{-1} g$$

$$\text{where } g^T = \frac{\partial \ln(P(D|\theta))}{\partial \theta} \Big|_{\theta=0}, \quad H = \frac{\partial^2 \ln(P(D|\theta))}{\partial \theta^2} \Big|_{\theta=0}$$

$\lambda_r$  and  $\lambda_t$  the regularization parameters for rotation and translation.

$$\arg\min_x (\|Ax - b\|^2 + \|\Gamma x\|^2), \quad \Gamma \text{ is chosen as a } \Gamma = \alpha I, \quad \hat{x} = (A^T A + \Gamma^T \Gamma)^{-1} A^T b \\ = (A^T A + (\alpha I)^T (\alpha I))^{-1} A^T b \\ = (A^T A + \alpha^2 I) A^T b$$

Finally, the pose can be updated according to:

$$\hat{M}T = \hat{M}T \begin{bmatrix} \exp([\hat{\theta}_r]_x) & \hat{\theta}_t \\ 0 & 1 \end{bmatrix}$$

www.luohanjie.com

The full likelihood:  $n_{ci}$  is the number of correspondence lines.

$$P(D|\theta) \propto \prod_{i=1}^{n_{ci}} P(D_i|\theta)$$

$$\therefore \ln P(D|\theta) \propto \sum_{i=1}^{n_{ci}} \ln P(D_i|\theta)$$

$$\Delta \tilde{c}_{si}^- \leq \Delta c_{si}^+ \leq \Delta \tilde{c}_{si}^+$$

$$P(D_i|\theta) \propto (\Delta \tilde{c}_{si}^+ - \Delta c_{si}^+) \underbrace{P(D_i|\Delta \tilde{c}_{si}^-)} + (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i|\Delta \tilde{c}_{si}^+)$$

$$P(D_i|\Delta \tilde{c}_{si}) \propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{c}_{si}) p_{sf_i}(r_s) + h_b(r_s - \Delta \tilde{c}_{si}) p_{sf_i}(r_s))$$

$$\Delta c_i^+ = n_i^T (\pi(c_i^+) - c_i)$$

$$c_i^+ = \begin{bmatrix} c_m^R & c_m^t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + [\theta_r]_x & \theta_b \\ 0 & 1 \end{bmatrix}^M \tilde{x}_i$$

$$\Delta c_{si}^+ = (\Delta c_i^+ - \Delta r_i) \frac{\bar{n}_i}{S}$$

$$g_T = \frac{\partial \ln(P(D|\theta))}{\partial \theta} \Big|_{\theta=0} = \sum_{i=1}^{n_{ci}} \frac{\partial \ln(P(D_i|\theta))}{\partial \Delta c_{si}^+} \frac{\partial \Delta c_{si}^+}{\partial c_i^+} \frac{\partial c_i^+}{\partial \theta} \Big|_{\theta=0}$$

$$\begin{aligned} \frac{\partial \ln(P(D_i|\theta))}{\partial \Delta c_{si}^+} \Big|_{\theta=0} &\stackrel{\text{linear interpolate}}{\approx} \frac{\partial (\ln[(\Delta \tilde{c}_{si}^+ - \Delta c_{si}^+) P(D_i|\Delta \tilde{c}_{si}^-) + (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i|\Delta \tilde{c}_{si}^+)])}{\partial \Delta c_{si}^+} \Big|_{\theta=0} \\ &= \frac{-P(D_i|\Delta \tilde{c}_{si}^-) + P(D_i|\Delta \tilde{c}_{si}^+)}{(\Delta \tilde{c}_{si}^+ - \Delta c_{si}^+) P(D_i|\Delta \tilde{c}_{si}^-) + (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i|\Delta \tilde{c}_{si}^+)} \end{aligned}$$

$$\therefore P(D_i|\Delta \tilde{c}_{si}) \propto \mathcal{N}(\Delta \tilde{c}_{si} | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}})$$

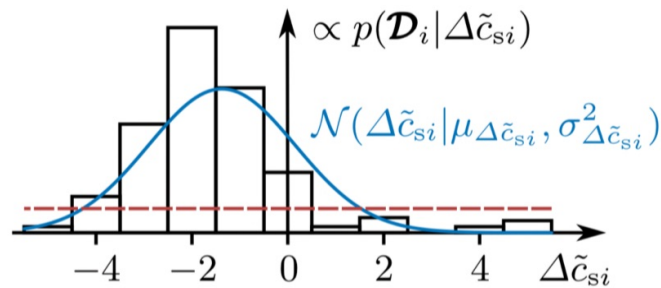
$$\therefore P(D_i|\theta) \propto (\Delta \tilde{c}_{si}^+ - \Delta c_{si}^+) P(D_i|\Delta \tilde{c}_{si}^-) + (\Delta c_{si}^+ - \Delta \tilde{c}_{si}^-) P(D_i|\Delta \tilde{c}_{si}^+)$$

The linear interpolation of two Gaussian distribution  $P(D_i|\Delta \tilde{c}_{si}^-)$  and  $P(D_i|\Delta \tilde{c}_{si}^+)$  is also a Gaussian distribution, and

$$P(D_i|\Delta \tilde{c}_{si}^-) \sim \mathcal{N}(\Delta c_{si}^+ | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}}) \sim P(D_i|\Delta \tilde{c}_{si}^+)$$

$$\therefore P(D_i|\theta) \sim \mathcal{N}(\Delta c_{si}^+ | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}})$$

$$\begin{aligned}
\frac{\partial \ln(P(\mathcal{D}_i | \Delta \tilde{c}_{si}))}{\partial \Delta c_{si}^+} &= \frac{\partial \ln(\mathcal{N}(\Delta c_{si}^+ | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}}^2))}{\partial \Delta c_{si}^+} & \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\
&= \frac{\partial \left( \ln \frac{1}{\sigma \sqrt{2\pi}} + \frac{-(\Delta c_{si}^+ - \mu_{\Delta \tilde{c}_{si}})^2}{2\sigma_{\Delta \tilde{c}_{si}}^2} \right)}{\partial \Delta c_{si}^+} \\
&= \frac{-2(\Delta c_{si}^+ - \mu_{\Delta \tilde{c}_{si}})}{2\sigma_{\Delta \tilde{c}_{si}}^2} = -\frac{(\Delta c_{si}^+ - \mu_{\Delta \tilde{c}_{si}})}{\sigma_{\Delta \tilde{c}_{si}}^2}
\end{aligned}$$



**Fig. 5.** Example showing normalized values of a noisy discrete likelihood  $p(\mathcal{D}_i | \Delta \tilde{c}_{si})$  and the normal distribution  $\mathcal{N}(\Delta \tilde{c}_{si} | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}}^2)$  that approximates that likelihood. The red line indicates a threshold for the probability values. To avoid errors from image noise and invalid pixel-wise posteriors, for values below this threshold, the normal distribution is used in the calculation of partial derivatives of the log-likelihood.

$$\begin{aligned}
P(\mathcal{D}_i | \Delta \tilde{c}_{si}) &\propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{c}_{si}) p_{sf_i}(r_s) + h_b(r_s - \Delta \tilde{c}_{si}) p_{st_i}(r_s)) \\
&\propto \mathcal{N}(\Delta \tilde{c}_{si} | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}}^2)
\end{aligned}$$

$$\text{and } \frac{\partial^2 \ln(P(\mathcal{D}_i | \Delta \tilde{c}_{si}))}{\partial \Delta c_{si}^{+2}} = -\frac{1}{\sigma_{\Delta \tilde{c}_{si}}^2}$$

$${}^c \tilde{\chi}_i^+ = \begin{bmatrix} {}^c \tilde{m}_R & {}^c \tilde{m}_t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + [\theta_r]_x & \theta_t \\ 0 & 1 \end{bmatrix} {}^M \tilde{\chi}_i$$

$$\left. \frac{\partial {}^c \tilde{\chi}_i^+}{\partial \theta} \right|_{\theta=0} = {}^c \tilde{m}_R \begin{bmatrix} -[{}^M \tilde{\chi}_i]_x & 1 \end{bmatrix}$$

$$\Delta C_i^+ = n_i^T (\pi(cX_i^+) - c_i)$$

$$\Delta C_{Si}^+ = (\Delta C_i^+ - \Delta r_i) \frac{\bar{n}_i}{S}$$

$$x_i = \pi(X_i) = \begin{bmatrix} X_i f_x / Z_i + p_x \\ Y_i f_y / Z_i + p_y \end{bmatrix}$$

$$\Delta C_{Si}^+ = (\Delta C_i^+ - \Delta r_i) \frac{\bar{n}_i}{S} = (n_i^T (\pi(cX_i^+) - c_i) - \Delta r_i) \frac{\bar{n}_i}{S}$$

$$= \underbrace{\frac{\bar{n}_i}{S}}_{\text{Scalar}} ([n_{ix} \ n_{iy}] \begin{bmatrix} {}^c X_i f_x / {}^c Z_i + p_x - c_{ix} \\ {}^c Y_i f_y / {}^c Z_i + p_y - c_{iy} \end{bmatrix} - \underbrace{\Delta r_i}_{\text{Scalar}})$$

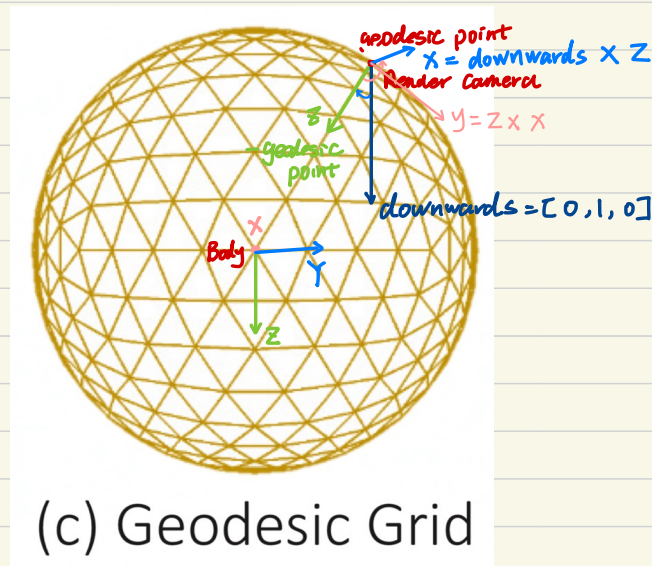
$$= \frac{\bar{n}_i}{S} (n_{ix} ({}^c X_i f_x / {}^c Z_i + p_x - c_{ix}) + n_{iy} ({}^c Y_i f_y / {}^c Z_i + p_y - c_{iy}) - \Delta r_i)$$

$$\frac{\partial \Delta C_{Si}^+}{\partial {}^c X_i^+} = \frac{\bar{n}_i}{S} \left[ \frac{n_{ix} f_x}{{}^c Z_i} \quad \frac{n_{iy} f_y}{{}^c Z_i} - \frac{n_{ix} {}^c X_i f_x + n_{iy} {}^c Y_i f_y}{({}^c Z_i)^2} \right]$$

$$H \propto \sum_{i=1}^{n_{cl}} \frac{\partial^2 \ln \mathcal{L}(\mathcal{D}_i | \theta)}{\partial \Delta C_{Si}^{+2}} \left( \frac{\partial \Delta C_{Si}^+}{\partial {}^c X_i^+} \quad \frac{\partial {}^c X_i^+}{\partial \theta} \right)^T \left( \frac{\partial \Delta C_{Si}^+}{\partial {}^c X_i^+} \quad \frac{\partial {}^c X_i^+}{\partial \theta} \right) \Big|_{\theta=0}$$

## Geodesic Grid

The 3D model of an object is thereby rendered from  $2^{16}$  different viewpoints that are placed on the vertices of a geodesic grid with a distance of 0.8m to the object center. For each rendering, we randomly sample  $N_{CL}=200$  points from the object contour.



# Code

## Init

1. ReadPosesRBOTDataset
2. Init Tracker()

NormalImageViewer {  
Render Geometry: holds rendering data for all assigned bodies as well as the glfw context for the render  
make render camera has the same intrinsics as real camera  
Dataset RBOTCamera: load images from the RBOT dataset.  
Renders all bodies into a normal image that is displayed on top of a color camera image

Body: holds all body info such as name, geometry data, pose and occlusion ids

Model: stores a model of a body that consists of template views with multiple contour points. It includes all functionality to generate, save and load the model.

RegionModality: {  
Body  
Model  
Dataset RBOTCamera

correspondence search, calculation of the gradient vector and hessian matrix. pose optimization, color histograms, visualization. occlusions using an occlusion mask

Tracker: {  
NormalImageViewer {  
Render Geometry  
Dataset RBOTCamera  
RegionModality: {  
Body  
Model  
Dataset RBOTCamera

## Evaluate

1. Init Bodies

① Model → Generate Model

a. Generate camera poses (GenerateGeodesicPoses) See Geodesic Grid

b. NormalImageRender: Render both a depth image and an image where the normal vector of the surface is encoded in the color of each pixel

setUpRender → RenderGeometry.AddBody → LoadMeshIntoVertices()



### C. Generate template views

For each view:

- ① Render image by GL
- ② FetchNormalImage (BGRA, SUC4)
- ③ FetchDepthImage (16U)

normal-image: CV-8UC4  
depth-image: CV-16U

#### ④ GeneratePointData

1. get A channel of BGRA image.

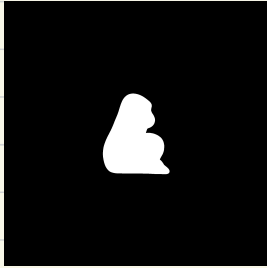
2. findContours

3. Calculate data for contour points

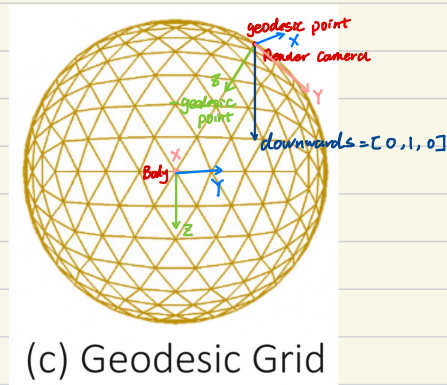
a. Randomly sample point on contour

b. Compute the point 3D position in the model body frame

2D pixel  $\rightarrow$  3D in Cam frame  $\rightarrow$  3D in Body frame. The 3D point of model of the contour.

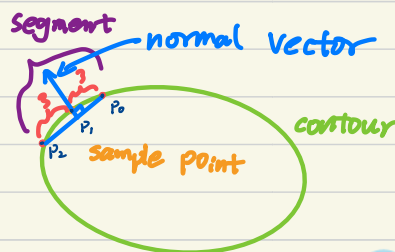


A channel of BGRA image.



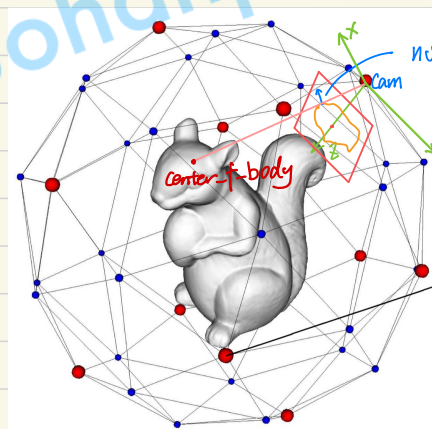
(c) Geodesic Grid

### C. Calculate contour segment and approximate normal vector



$$\begin{aligned}\vec{n} \cdot (P_0 - P_1) &= 0 \\ &= \vec{n} \cdot \vec{v} = 0 \\ &= n_x v_x + n_y v_y = 0 \\ \therefore \vec{n} &= [-v_y, v_x]\end{aligned}$$

### D. normal vector of the sample contour in Body frame

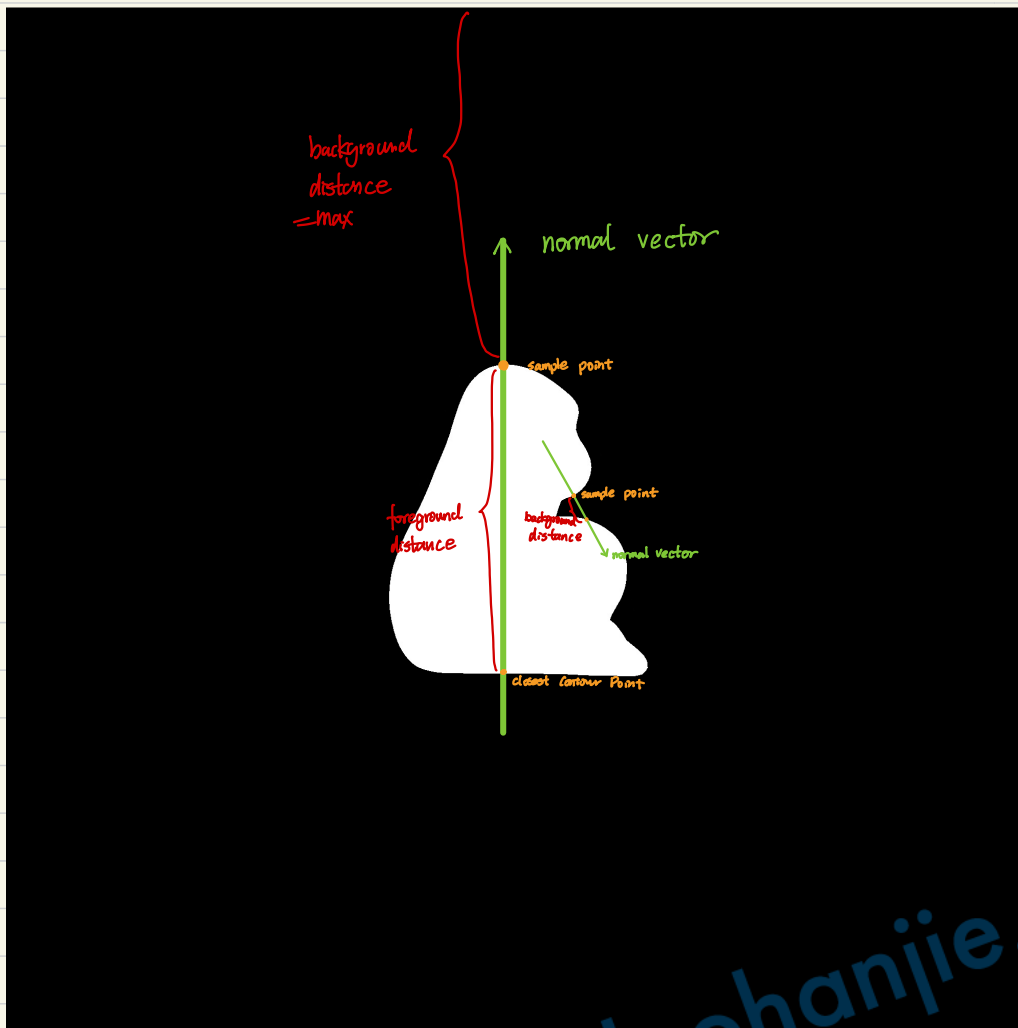


E. compute continuous distances for the foreground and background where the corresponding regions are not interrupted by the other. Starting from a coordinate, those distance are measured along the normal vector.

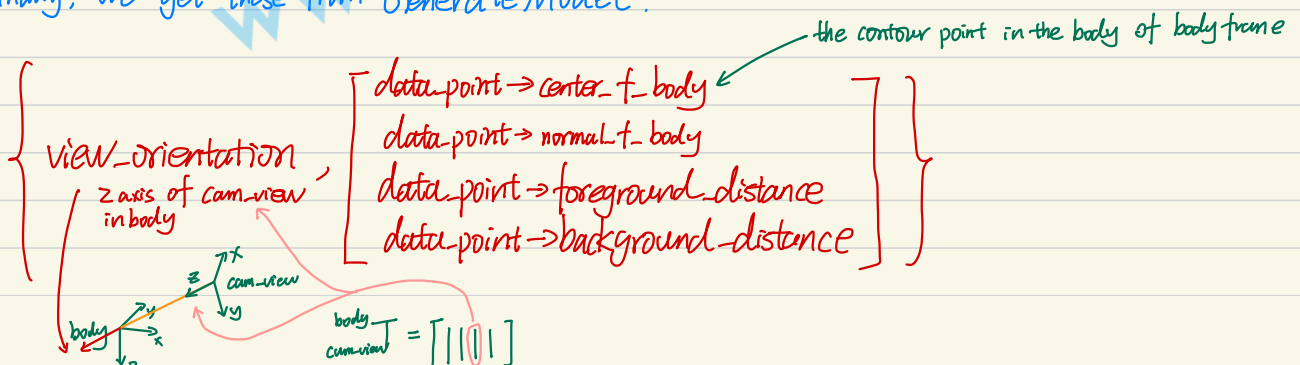
$$\begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} f_x x + c_x z \\ f_y y + c_y z \\ z \end{bmatrix} = \begin{bmatrix} f_x \frac{x}{z} + c_x \\ f_y \frac{y}{z} + c_y \\ 1 \end{bmatrix} = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

$$f_x \frac{x}{z} + c_x = u \Rightarrow \frac{z}{f_x} = \frac{x}{c_x - u}, \quad \times \frac{f_x}{z} = c_x - u$$

$$\therefore \text{pixel-to-meter} = \frac{z}{f_x}$$



Finally, we get these from GenerateModel:



```

version_id
1
sphere_radius
0.8
n_divides
4
n_points
200
image_size
2000
n_template_views
2562
geometry_path
/media/luohanjie/Hanjie/RBOT_dataset/apc/apc.obj

geometry_unit_in_meter
0.001
geometry_defined_counterclockwise
1
geometry_enable_culling
0
geometry2body_pose
1 0 0 0
0 1 0 0
0 0 1 0
0 0 0 1
maximum_body_diameter
0.3

```

Camera 2 body - pose :  $\text{view\_cam}^{\text{body}} T$  from Generate Point Data  
 body 2 world - pose :  $\text{world}^{\text{body}} T = \text{cam}^{\text{body}} T$  from model to real camera.  
 RegionModality  $\rightarrow$  body 2 camera - pose :  $\text{camera} \rightarrow \text{world} \rightarrow \text{camera} \cdot \text{body} \rightarrow \text{body} \rightarrow \text{world} = I$   $\text{cam}^{\text{body}} T$   
 geometry2body - pose :

Normal Image Renders  $\rightarrow$  camera 2 world - pose :  $\text{view\_cam}^{\text{body}} T$ , = Camera 2 body - pose

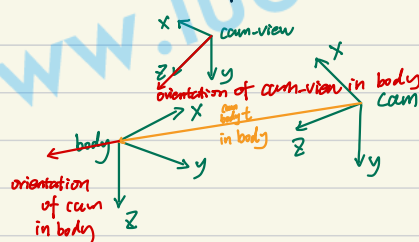
## 2. Reset Body (0)

① region\_modality\_ptr  $\rightarrow$  Start Modality C)

a. Init histograms

AddLinePixelColorsToTempHistograms()

① GetClosestTemplateView ( $\text{cam}^{\text{body}} T$ )



$$(\text{cam}^{\text{body}} R)^{-1} \text{cam}^{\text{body}} t = \text{body}^{\text{cam}} t, \quad \text{orientation} = \text{body}^{\text{cam}} t / \|\text{body}^{\text{cam}} t\|$$

get template\_view which ori of  $\text{cam}^T$  ori of cam-view is maximum.

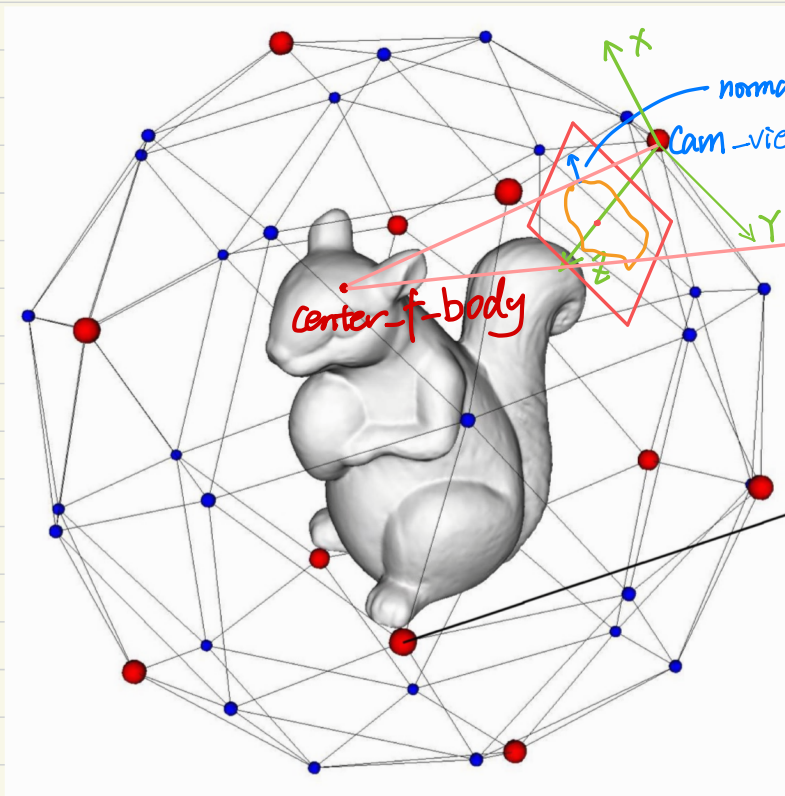
$n\_histogram\_bins = 32$

histogram - f ( $32^3$ )

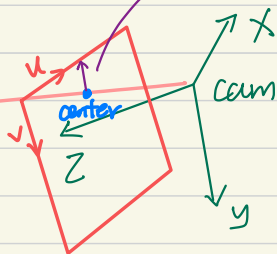
histogram - b ( $32^3$ )

For each datapoints of the matched template-view

body 2 Camera pose . center-f-body = center-f-camera .  
 $\text{cam body}^T$  , 3D point in the body project point in camera frame  
 which is projected become  
 contour .



normal (normal-f-body in com frame)



compute: data in cam frame.

center : in image

Normal: in cam frame

foreground-distance : in image

background distance = in image

compute the foreground and background histogram

along the line, which maximum is 10 pixels.



template\_view)

$$\text{Histogram} = [B \Rightarrow 3 \times 32^2 + G \Rightarrow 3 \times 32 \times 32 + R \Rightarrow 3]$$

white: [255, 255, 255]

$$= \frac{255}{8} \cdot 32 \times 32 + \frac{255}{8} 32 + \frac{255}{8}$$

$$= 255 \cdot 4 \times 32 + 255 \cdot 4 +$$

$$[\underbrace{B}_{32} \quad \underbrace{G}_{32} \quad \underbrace{R}_{32}]$$

`n_histogram_bins_:` 32

`temp_histogram_f_:`  $n\_histogram\_bins\_^3$

`temp_histogram_b_:`  $n\_histogram\_bins\_^3$

### 3. Iterate over all frames

#### 3.1 Execute Measured Tracking Cycle

To find the final pose, the process is repeated 7 times, starting each time with the retrieval of a new template view. The first and second iterations use a scale of  $s = 5$  and  $s = 2$ , respectively. For all other iterations, the scale is set to  $s = 1$ . Examples of different scales are shown in Fig. 1. The specified scales have the effect that in the first iteration a large area with low resolution is considered while short lines with high resolution are used in later iterations.

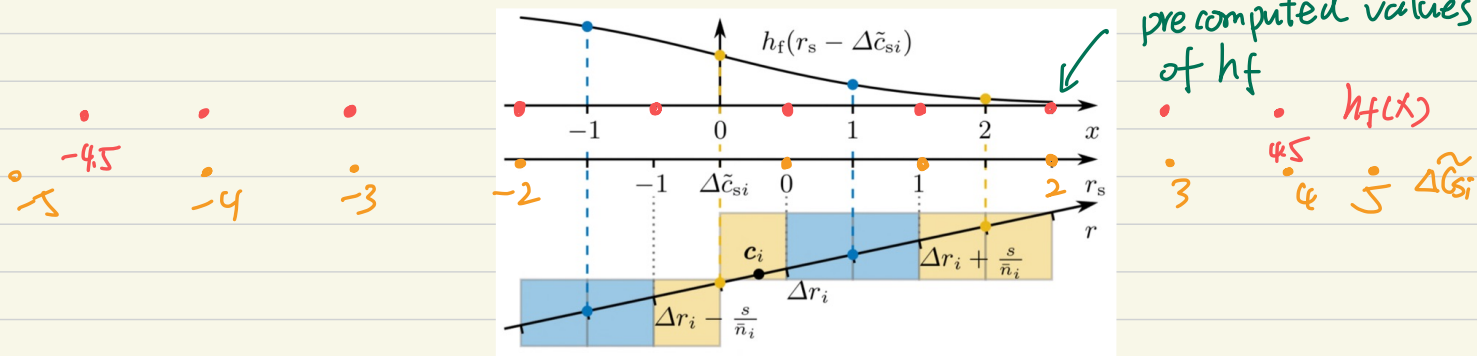
For each iteration:  $\text{Scale} = [5, 2, 1, 1, \dots]$

#### 3.1.1 Calculate Correspondence

$$\text{body} \rightarrow \text{camera-pose} = \text{camera} \rightarrow \text{world} \rightarrow \text{camera} \cdot \text{body} \rightarrow \text{body} \rightarrow \text{world} = \mathbf{I} \begin{matrix} \text{cam} \\ \text{body} \end{matrix}^T$$

$$\text{line-length-in-segments} = \underbrace{\text{function-length}}_{10} + \underbrace{\text{distribution-length}}_{11} - 1$$

20 segments



①  $P(D_i | \Delta \tilde{c}_{si})$  is calculated for 11 discrete values  $\Delta \tilde{c}_{si} \in \{-5, -4, \dots, 5\}$

② For  $h_f(x)$ ,  $h_b(x)$ , 10 precomputed values  $x \in \{-4.5, -3.5, -2.5, -1.5, -0.5, 0.5, 1.5, 2.5, 3.5, 4.5\}$  are used.

function.

$$\textcircled{3} \quad h_f(x) = \frac{1}{2} - \frac{1}{2} \tanh\left(\frac{x}{2s_h}\right), \quad h_b(x) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{x}{2s_h}\right)$$

$$s_h = 1.3$$

Precalculate Function Lookup: ① function length = 10,  $x \in \{-4.5, -3.5, -2.5, -1.5, -0.5, 0.5, 1.5, 2.5, 3.5, 4.5\}$

$$h_f(x) = \frac{1}{2} - \frac{1}{2} \tanh\left(\frac{x}{2 \cdot 1.3}\right)$$

$$h_b(x) = 1 - h_f(x)$$

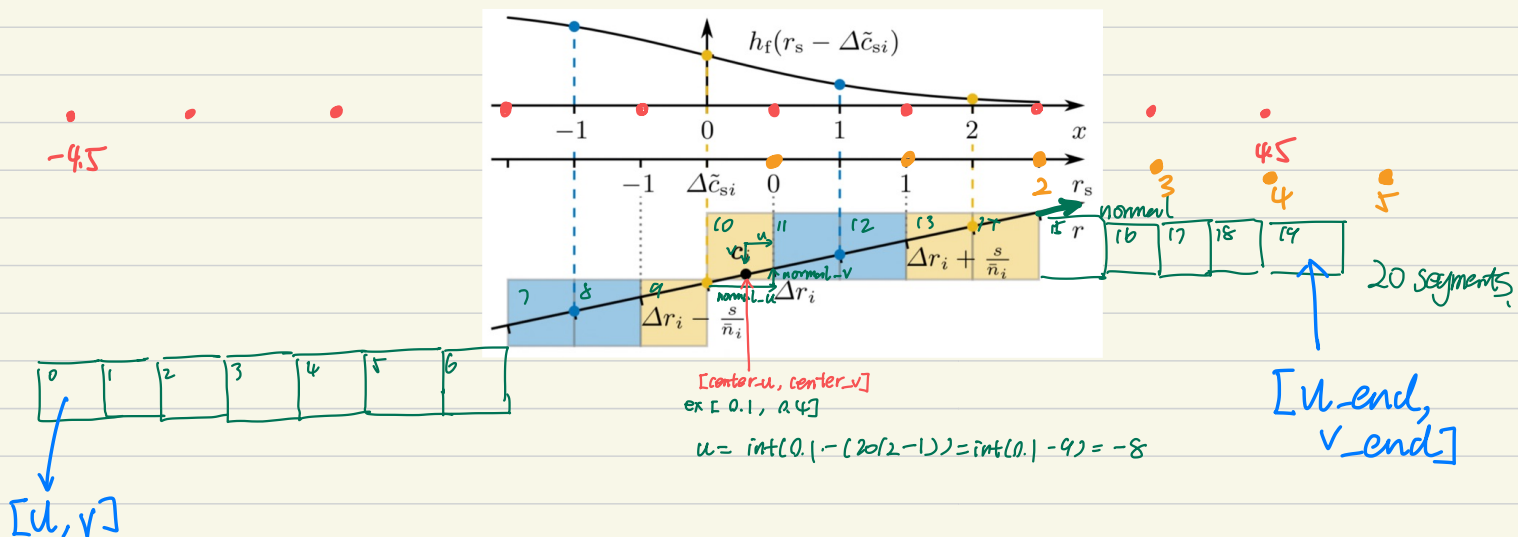
$$\text{line\_length} = \text{line\_length\_in\_segments} \times \text{scale}$$

Search closest template view

For each contour point of matched template-view :

compute : center-f-body  
center-f-camera  
center-u, center-v  
normal-u, normal-v  
continuous: in pixels  $\times$  scale

① Calculate Segment Probabilities :



$$\text{if } \text{normal\_u} > \text{normal\_v} \rightarrow \text{v\_step} = \frac{\text{normal\_v}}{\text{normal\_u}} \quad \text{假设 Scale} = 1$$

$$\text{starting point } u = \text{int}(\text{center\_u} - (\frac{\text{line\_length}}{2} - 0.5) + 0.5)$$

$$= \text{int}(\text{center\_u} - (\text{line\_length}/2 - 1))$$

$$u\_end = u + \text{line\_length} - 1$$



$$\Delta \tilde{c}_{s_i}^- \leq \Delta c_{s_i}^+ \leq \Delta \tilde{c}_{s_i}^+$$

$$P(D_i | \theta) \propto (\Delta \tilde{c}_{s_i}^+ - \Delta c_{s_i}^+) \underbrace{P(D_i | \Delta \tilde{c}_{s_i}^-)} + (\Delta c_{s_i}^+ - \Delta \tilde{c}_{s_i}^-) P(D_i | \Delta \tilde{c}_{s_i}^+)$$

$$P(D_i | \Delta \tilde{c}_{s_i}) \propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{c}_{s_i}) p_{sf_i}(r_s) + h_b(r_s - \Delta \tilde{c}_{s_i}) p_{sb_i}(r_s))$$

$P(D_i | \Delta \tilde{c}_{s_i})$  is calculated for 11 discrete values  $\Delta \tilde{c}_{s_i} \in \{-5, -4, \dots, 5\}$   
distribution-length

$R_s$  a set of distance to segment centers

Iterate over all pixels of line and calculate probabilities  $p_{sf_i}(r_s)$  and  $p_{sb_i}(r_s)$

for  $x \in [u, u\_end]$ ,  $y \in [v, v\_end]$ ,

get ① pixel-color-probability-f = histogram-f-[x,y] :  $p(y_i(r) | m_f)$

pixel-color-probability-b = histogram-b-[x,y] :  $p(y_i(r) | m_b)$

$$\tilde{p}(y_i(r) | m_f) = \frac{p(y_i(r) | m_f)}{p(y_i(r) | m_f) + p(y_i(r) | m_b)}$$

$$\tilde{p}(y_i(r) | m_b) = \frac{p(y_i(r) | m_b)}{p(y_i(r) | m_f) + p(y_i(r) | m_b)}$$

Then:

$$\text{segment\_probability\_f}[\text{segment\_id}, x] = \prod_{r \in S(r_s)} \tilde{p}(y_i(r) | m_f)$$

$$\text{segment\_probability\_b}[\text{segment\_id}, x] = \prod_{r \in S(r_s)} \tilde{p}(y_i(r) | m_b)$$

where  $S$  is a set-valued function that maps  $r_s$  to a set of values  $r$  that describe the  $S(\text{scale})$  closest pixel centers of a segment.

$$\text{normal-component-to-scale} = \bar{n}_i / s$$

$$\text{delta-r} = \Delta r_i$$

$\bar{n}_i = \max(|n_{x_i}|, |n_{y_i}|)$  projects a normal vector  $n_x$  to the closest horizontal or vertical coordinate.

$\Delta r_i \in \mathbb{R}$  distance from center  $c_i$  to defined segment location (chosen such that the center  $r_s = 0$  lies on the boarder between two segments).



(2) Calculate Distribution

$$\text{distribution\_lengths} = \|\Delta \vec{C}_i \in \{-5, -4, \dots, 5\}\|, \quad \|\Delta \vec{C}_i\|$$

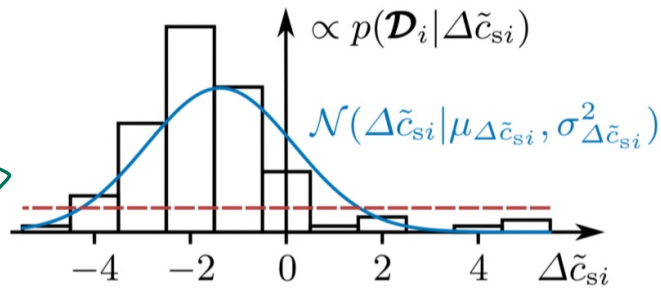
Loop over entire distribution and start values of segment probabilities:

Loop over values of segment probabilities and corresponding lookup values:

Lookup { functionLength = 10,  $x \in \{-4.5, -3.5, -2.5, -1.5, -0.5, 0.5, 1.5, 2.5, 3.5, 4.5\}$

$h_f(x) = \frac{1}{2} - \frac{1}{2} \tanh\left(\frac{x}{2 \cdot 1.3}\right)$

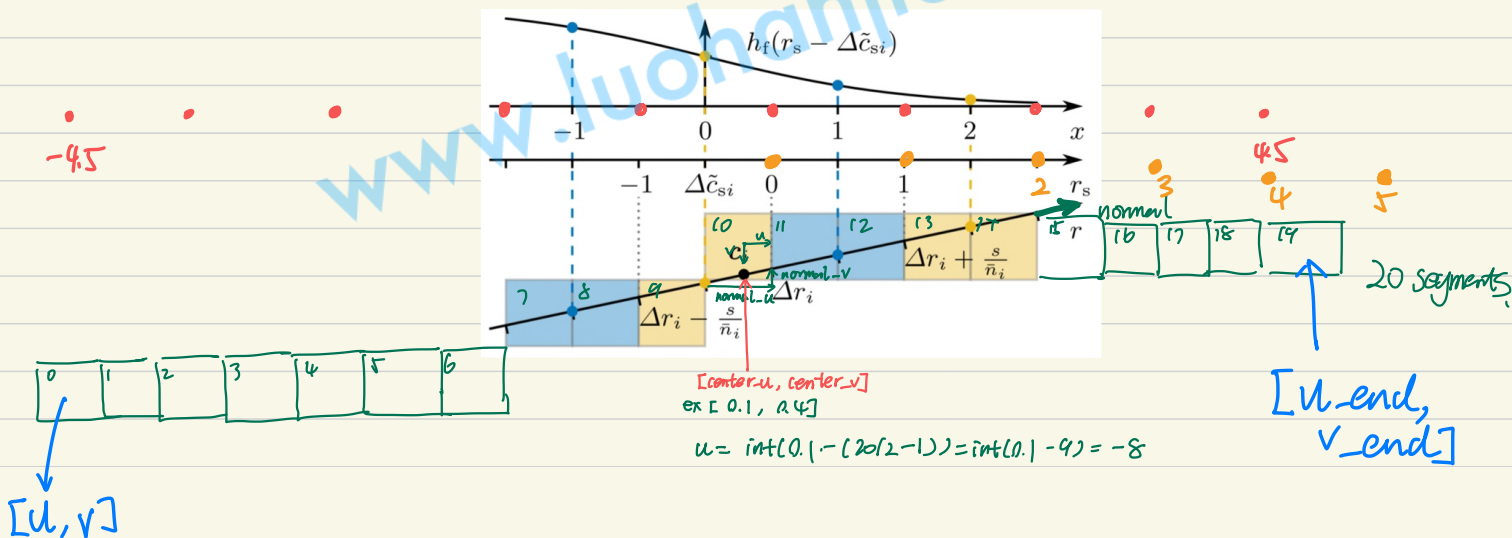
$h_b(x) = 1 - h_f(x)$



**Fig. 5.** Example showing normalized values of a noisy discrete likelihood  $p(\mathcal{D}_i | \Delta \tilde{c}_{si})$  and the normal distribution  $\mathcal{N}(\Delta \tilde{c}_{si} | \mu_{\Delta \tilde{c}_{si}}, \sigma_{\Delta \tilde{c}_{si}}^2)$  that approximates that likelihood. The red line indicates a threshold for the probability values. To avoid errors from image noise and invalid pixel-wise posteriors, for values below this threshold, the normal distribution is used in the calculation of partial derivatives of the log-likelihood.

$$p(D_i | \Delta \tilde{C}_{si}) \propto \prod_{r_s \in R_s} (h_f(r_s - \Delta \tilde{C}_{si}) p_{sf_i}(r_s) + h_b(r_s - \Delta \tilde{C}_{si}) p_{sf_i}(r_s))$$

compute  $P(D_i | \Delta \tilde{C}_i)$  for each  $\Delta \tilde{C}_i$



③ Calculate mean and var of  $p(D_i | \Delta \tilde{C}_{s_i})$

3.1.2 For each update iteration : 2

① Calculate PoseUpdate()

body2camera\_pose\_ : <sup>cam</sup>body<sup>T</sup>

$$\Delta C_i^+ = n_i^T (\pi(C_i^+) - c_i) \quad \Delta C_{s_i}^+ = (\Delta C_i^+ - \Delta r_i) \frac{\bar{n}_i}{S}$$

$$\frac{\partial \ln(P(D_i | \Delta \tilde{C}_{s_i}))}{\partial \Delta C_{s_i}^+} = \begin{cases} -\frac{(C_{s_i}^+ - M \Delta \tilde{C}_{s_i})}{6 \Delta \tilde{C}_{s_i}^2}, & P(D_i | \Delta \tilde{C}_{s_i}) \text{ and } P(D_i | \Delta \tilde{C}_{s_i}^+) < 0.01 \\ \frac{-P(D_i | \Delta \tilde{C}_{s_i}) + P(D_i | \Delta \tilde{C}_{s_i}^+)}{(C_{s_i}^+ - \Delta \tilde{C}_{s_i}) P(D_i | \Delta \tilde{C}_{s_i}) + (C_{s_i}^+ - \Delta \tilde{C}_{s_i}^+) P(D_i | \Delta \tilde{C}_{s_i}^+)} \end{cases}$$

$$\frac{\partial \Delta C_{s_i}^+}{\partial C_i^+} = \frac{\bar{n}_i}{S} \left[ \frac{n_{ix} f_x}{c_{z_i}} \quad \frac{n_{iy} f_y}{c_{z_i}} - \frac{n_{ix} c_{x_i} f_x + n_{iy} c_{y_i} f_y}{c_{z_i}^2} \right]$$

$$\left. \frac{\partial C_i^+}{\partial \theta} \right|_{\theta=0} = {}^c_m R [-[{}^M X_i]_x \quad I]$$

$$\frac{\partial^2 \ln(P(D_i | \Delta \tilde{C}_{s_i}))}{\partial \Delta C_{s_i}^{+2}} = -\frac{1}{6 \Delta \tilde{C}_{s_i}^2}$$

$$g^T = \left. \frac{\partial \ln(P(D|\theta))}{\partial \theta} \right|_{\theta=0} = \sum_{i=1}^{N_{cl}} \left. \frac{\partial \ln(P(D_i|\theta))}{\partial \Delta C_{s_i}^+} \frac{\partial \Delta C_{s_i}^+}{\partial C_i^+} \frac{\partial C_i^+}{\partial \theta} \right|_{\theta=0}$$

$$H = \left. \frac{\partial^2 \ln(P(D|\theta))}{\partial \theta^2} \right|_{\theta=0}$$

$$\approx \sum_{i=1}^{N_{cl}} \frac{\partial^2 \ln(P(D_i|\theta))}{\partial \Delta C_{s_i}^{+2}} \left( \frac{\partial \Delta C_{s_i}^+}{\partial C_i^+} \frac{\partial C_i^+}{\partial \theta} \right)^T \left( \frac{\partial \Delta C_{s_i}^+}{\partial C_i^+} \frac{\partial C_i^+}{\partial \theta} \right) \Big|_{\theta=0}$$

$$\delta = (-H + [{}^{\lambda_r} l_3 \quad 0; 0 \quad \lambda_{t13}])^{-1} g$$

$$\text{delta-cs} : \Delta C_{si}^+$$

$$\text{dloglikelihood\_ddelta-cs} : \frac{\partial \ln(P(D_i | \Delta \tilde{C}_{si}))}{\partial \Delta C_{si}^+}$$

$$\text{ddelta-cs\_dcenter} : \frac{\partial \Delta C_{si}^+}{\partial C_{xi}^+}$$

$$\text{ddelta-cs\_dtheta} : \frac{\partial \Delta C_{si}^+}{\partial \theta}$$

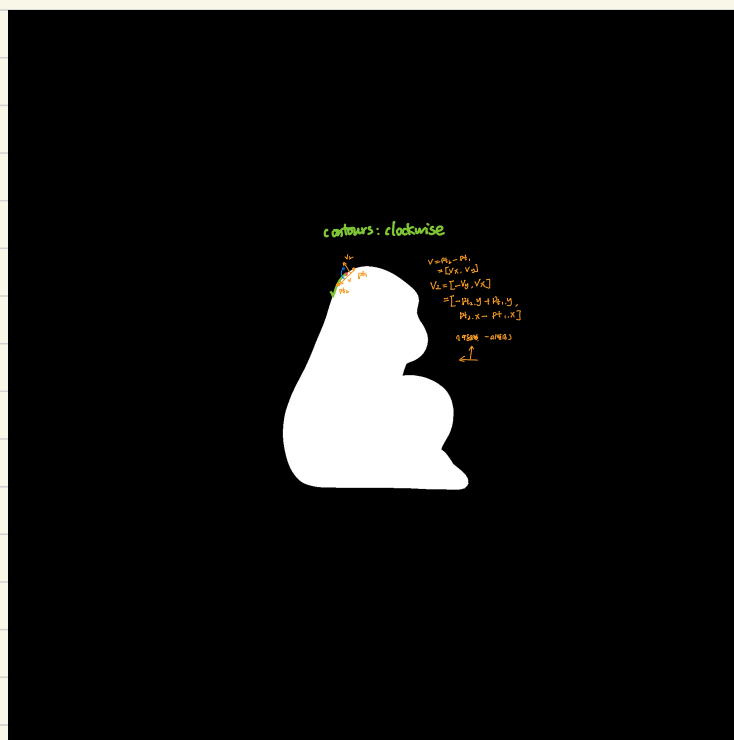
$$\text{gradient} : \sum_{i=1}^{n_{c1}} \frac{\partial \ln(P(D_i | \theta))}{\partial \Delta C_{si}^+} \frac{\partial \Delta C_{si}^+}{\partial C_{xi}^+} \frac{\partial C_{xi}^+}{\partial \theta} \bigg|_{\theta=0}$$

$$\text{hessian} :$$

$$\sum_{i=1}^{n_{c1}} - \frac{1}{6\Delta C_{si}^2} \left( \frac{\partial \Delta C_{si}^+}{\partial C_{xi}^+} \frac{\partial C_{xi}^+}{\partial \theta} \right)^T \left( \frac{\partial \Delta C_{si}^+}{\partial C_{xi}^+} \frac{\partial C_{xi}^+}{\partial \theta} \right) \bigg|_{\theta=0}$$

$$\hat{\theta} = (-H + \begin{bmatrix} \lambda_r I_3 & 0 \\ 0 & \lambda_{t13} \end{bmatrix})^{-1} g$$

$${}^{cam}_{body} T = {}^{cam}_{body} T \cdot \Delta T$$



# SRT3D: A Sparse Region-Based 3D Object Tracking Approach for the Real World (2021)

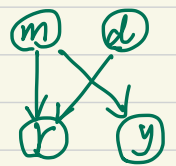
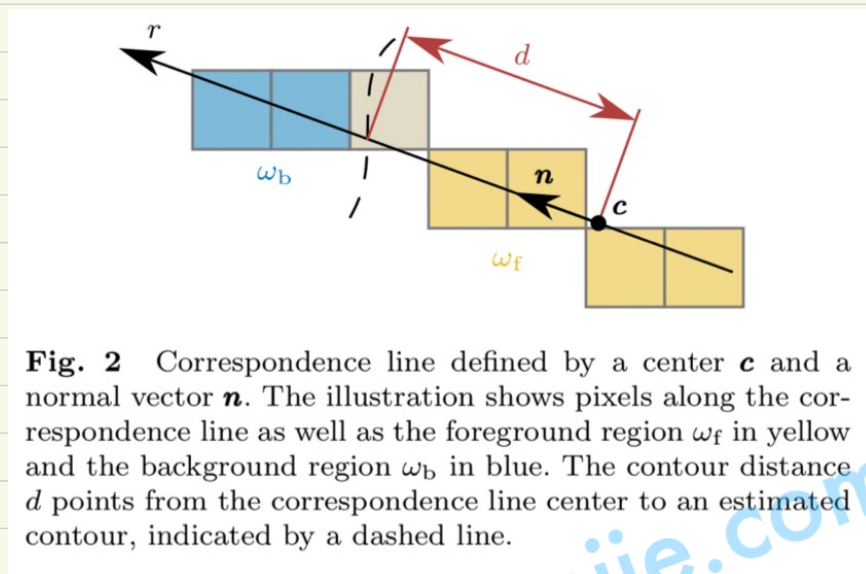
Images :  $I: \Omega \rightarrow \{0, \dots, 255\}^3$   
 $\uparrow$   
 $\Omega \subset \mathbb{R}^2$  image domain

correspondence line :  $l: W \rightarrow \{0, \dots, 255\}^3$   
 $\uparrow$   
 $W \subset \mathbb{R}$ , line domain

image value :  $y = I(x)$ ,  $x = [x \ y]^T$

$y = l(r)$ ,  $r \in \mathbb{R}$  in line.  
 $= I(c + rn)$

$c = [c_x, c_y]^T$  center,  $n = [n_x, n_y]^T$  normal vector  
 $\|n\|_2 = 1$



$\{r, y\}$  is a sample

1. For a single pixel

$$P(r, y, d, m) = P(r, y | d, m) P(d, m) = P(r | d, m) P(y | d, m) P(d, m)$$

$$\uparrow \quad \uparrow \quad \uparrow$$

$$\text{color} \quad \text{model}$$

$$= P(r | d, m) P(y | m) P(d) P(m)$$

distance from the correspondence line  
center  $c$  to an estimated contour.

$$\begin{aligned}
 P(r, y, d, m) &= P(r|d, m) P(y|m) P(d) P(m) \\
 P(r, d, m|y) P(y) &= P(r|d, m) P(y|m) P(d) P(m) \\
 P(r, d, m|y) &= P(r, |d, m) \left[ \frac{P(y|m) P(m)}{P(y)} \right] P(d) \\
 &\quad \uparrow \\
 &\quad P(m|y) P(y) = P(y|m) P(m) \\
 &= \underbrace{P(r, |d, m)}_{\text{conditional line coordinate probabilities}} P(m|y) P(d)
 \end{aligned}$$

$$\therefore P(m_i|y) = \frac{P(y|m_i) P(m_i)}{\sum_{j \in \{t, b\}} P(y|m_j) P(m_j)}, \quad i \in \{t, b\}$$

$P(m_i|y)$ : the probability distributions that describe how likely it is that a specific color value is part of the foreground / background region.

Foreground and background are equally likely along the correspondence line,  $P(m_f) = P(m_b)$  locally

$$\therefore P(m_i|y) = \frac{P(y|m_i)}{P(y|m_f) + P(y|m_b)}, \quad i \in \{t, b\}$$

$$\begin{aligned}
 P(X) &= \sum_y P(X, Y=y) \\
 &= \sum_y P(X|Y=y) P(Y=y)
 \end{aligned}$$

$$P(d|r, y) P(r, y) = P(d, r, y) = P(d|r, y) P(r) P(y)$$

$$P(d|r, y) = \frac{1}{P(r)} \frac{P(d, r, y)}{P(y)} = \frac{1}{P(r)} \frac{P(d, r|y) P(y)}{P(y)}$$

$$= \frac{1}{P(r)} P(d, r|y) = \frac{1}{P(r)} \sum_{i \in \{t, b\}} P(d, r, m_i|y)$$

$$= \frac{1}{P(r)} \sum_{i \in \{t, b\}} \underbrace{P(r, |d, m_i)}_{\substack{d, m_i \text{ and } r \text{ is probability}}} P(m_i|y) P(d)$$

Posterior probability over the entire correspondence line domain  $w$ :

$$P(d|w, l) = \prod_{r \in w} \sum_{i \in \{t, b\}} \frac{1}{P(r)} P(r|d, m_i) P(m_i|y) P(d)$$

The probability of the contour distance  $d$  given data from a correspondence line.

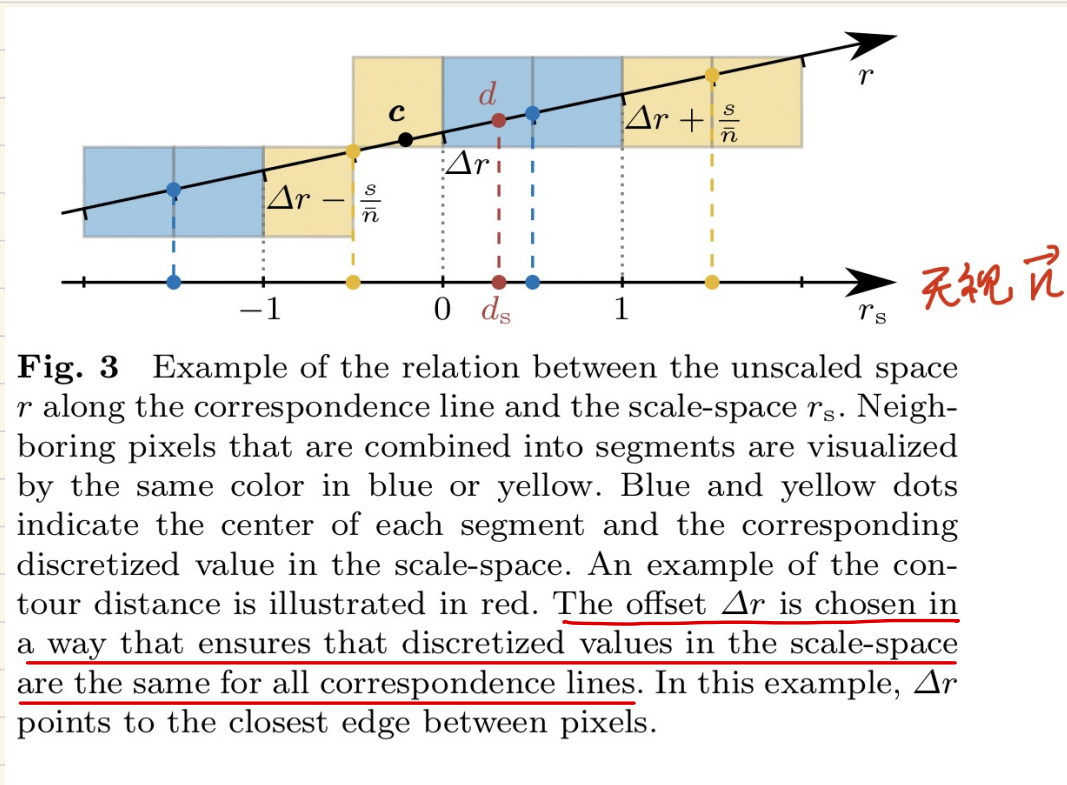
$$\propto \prod_{r \in w} \sum_{i \in \{t, b\}} P(r|d, m_i) P(m_i|l(r))$$

$l(r) = l(c+rn) = y$ ,  $P(r)$  and  $P(d)$  are considered to be uniform and constant and are thus dropped.



Estimating the  $P(d|w, l)$  is computationally expensive since, for each distance  $d$ , the product of  $P(d|w, l)$  has to be computed over the entire domain  $w$ .

2. Discrete scaled space formulation combines multiple pixels into segments and projects from continuous space along the correspondence line into a **discrete** space that is independent of a correspondence line's location and orientation.



$$r_3 = (r - \Delta r) \frac{\bar{n}}{5}$$

$$ds = (d - \Delta r) \frac{\bar{n}}{3}$$

$S \in \mathbb{N}^+$ : the number of pixels combined into a segment.  
 $\bar{n} = \max(|n_x|, |n_y|)$

$\therefore$  The posterior probability in the discrete scale-space:

$$p(d|w, l) \propto \prod_{r \in w} \sum_{i \in \{t, b\}} p(r|d, m_i) p(m_i | l(r))$$

$$\Rightarrow P(ds | w_s, l_s) \propto \prod_{r \in w_s} \sum_{i \in \{1, b\}} p(r_s | ds, m_i) p(m_i | l_s, r_s)$$

scaled correspondence  
line domain

$S = \{s(r_s)\}$  is a set-valued function that maps from the scaled line coordinate  $r_s$  to the segment  $S$ , which is a set of the closest  $s$  pixel values  $y$ .

$$\therefore P(m_i | y) = \frac{P(y | m_i)}{\sum_{j \in \{t, b\}} P(y | m_j)}, \quad i \in \{t, b\}$$

$$\therefore P(m_b) = P(m_t)$$

$$P(m_i | y_1, y_2, y_3) = \frac{P(y_1, y_2, y_3 | m_i) P(m_i)}{\sum_{j \in \{t, b\}} P(y_1, y_2, y_3 | m_j) P(m_j)}$$

$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1, y_2, y_3 | m_t) + P(y_1, y_2, y_3 | m_b)}$$

$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1, y_2, y_3 | m_t) + P(y_1, y_2, y_3 | m_b)}$$

$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1 | m_t) P(y_2 | m_t) P(y_3 | m_t) + P(y_1 | m_b) P(y_2 | m_b) P(y_3 | m_b)}$$

$$\therefore P(m_i | S) = \frac{\prod_{y \in S} P(y | m_i)}{\prod_{y \in S} P(y | m_t) + \prod_{y \in S} P(y | m_b)}$$

↑  
a set of the closest spread values y

$$P(ds | ws, l_s) \propto \prod_{r \in ws} \sum_{i \in \{t, b\}} P(r_s | ds, m_i) P(m_i | l_s(r_s))$$

$$\textcircled{1} \quad P(m_i | y_j) = \frac{P(y_j | m_i) P(m_i)}{P(y_j | m_t) P(m_t) + P(y_j | m_b) P(m_b)}$$

If color  $y_j$  do not sample,  $P(m_i | y_j) = 0.5$

$$\textcircled{2} \quad P_{m_i} = P(m_i | y_1) P(m_i | y_2) P(m_i | y_3) \dots$$

$$\textcircled{3} \quad P(m_i | S) = \frac{P_{m_i}}{P_{m_t} + P_{m_b}}$$

$$= \frac{P(y_1 | m_i)}{[P(y_1 | m_t) + P(y_1 | m_b)]} \frac{P(y_2 | m_i)}{[P(y_2 | m_t) + P(y_2 | m_b)]} \frac{P(y_3 | m_i)}{[P(y_3 | m_t) + P(y_3 | m_b)]}$$

$$\frac{P(y_1 | m_t)}{[P(y_1 | m_t) + P(y_1 | m_b)]} \frac{P(y_2 | m_t)}{[P(y_2 | m_t) + P(y_2 | m_b)]} \frac{P(y_3 | m_t)}{[P(y_3 | m_t) + P(y_3 | m_b)]} + \frac{P(y_1 | m_b)}{[P(y_1 | m_t) + P(y_1 | m_b)]} \frac{P(y_2 | m_b)}{[P(y_2 | m_t) + P(y_2 | m_b)]} \frac{P(y_3 | m_b)}{[P(y_3 | m_t) + P(y_3 | m_b)]}$$

$$= \frac{\prod_{y \in S} P(y | m_i)}{\prod_{y \in S} P(y | m_t) + \prod_{y \in S} P(y | m_b)}$$

$$\textcircled{1} \begin{array}{ccc} P(m_b|y_1) & P(m_b|y_2) & P(m_b|y_3) \\ 0.5 & 1 & 0.1 \\ P(m_f|y_1) & P(m_f|y_2) & P(m_f|y_3) \\ 0.5 & 0 & 0.9 \end{array}$$

$$P(m_b|S) = \frac{0.5 \cdot 1 \cdot 0.1}{0.5 \cdot 1 \cdot 0.1 + 0} = 1$$

$$\begin{aligned} \textcircled{2} \quad P(m_b|S) &= \frac{0.001 \times 0.1 \times 0.1}{0.001 \times 0.1 \times 0.1 + 0.001 \times 0.001 \times 0.2} \\ &= \frac{0.1}{0.1 + 0.001 \times 0.2} \\ &= \frac{0.1}{0.1 + 1e^{-4} \times 0.2} \approx 1 \end{aligned}$$

## Laplace smoothing

Laplace smoothing is a technique to smooth categorical data. Given a set of observation counts  $X = \{X_1, X_2, \dots, X_d\}$  from a  $d$ -dim multinomial distribution with  $N$  trials, a "smoothed" version of counts gives the est:

$$\theta_i^j = \frac{X_i + \alpha}{N + \alpha d} \quad (i = 1, \dots, d)$$

It can address the problem when we encounter a unseen data which the posterior probability will become 0.  
ex.  $X_i$  is the number of color  $y$  observation times.

$$\therefore p(y|m_i) = \frac{\text{the number of pixel which value } y \text{ (} X_i \text{)}}{\text{total sample pixel numbers (} N \text{)}}$$

$$\rightarrow \text{smooth } p(y|m_i) = \frac{\text{the number of pixel which value } y \text{ (} X_i \text{)} + 1}{\text{total sample pixel numbers (} N \text{)} + \text{the total number of } y \text{ values (} 3 \times 2768 = (256/8)^3 \text{)}}$$

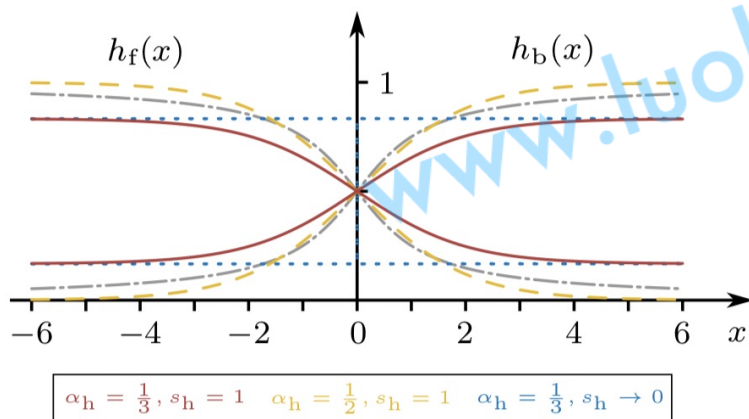
## Underflow

If  $P(y|m_i)$  is small,  $\prod_{y \in S} P(y|m_i)$  can produce numbers too small for floating point storage.

$$P(m_i|S) = \frac{\prod_{y \in S} P(y|m_i)}{\prod_{y \in S} P(y|m_f) + \prod_{y \in S} P(y|m_b)}$$

$$\prod_{y \in S} P(y|m_i) = \exp(\log(\prod_{y \in S} P(y|m_i))) = \exp(\sum_{y \in S} \log(P(y|m_i)))$$

3.  $p(r|d, m_f)$ ,  $p(r|d, m_b)$  models a local uncertainty with respect to the exact location of the foreground and background transition.



**Fig. 4** Smoothed step functions  $h_f$  and  $h_b$  that model the conditional line coordinate probabilities  $p(r | d, m_f)$  and  $p(r | d, m_b)$ . The functions  $h_f(x) = \frac{1}{2} - \frac{1}{\pi} \tan^{-1}(\frac{x}{s_h})$  and  $h_b(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}(\frac{x}{s_h})$  used by Zhong et al. (2020b) and Tjaden et al. (2018) are illustrated by dash-dotted gray lines. The definitions  $h_f(x) = \frac{1}{2} - \frac{1}{2} \tanh(\frac{x}{2s_h})$  and  $h_b(x) = \frac{1}{2} + \frac{1}{2} \tanh(\frac{x}{2s_h})$  from our previous work (Stoiber et al., 2020) are shown as dashed yellow lines. In both plots, the slope parameter  $s_h = 1$  is used. For the proposed functions in Eq. (12) and (13), solid red lines correspond to  $\alpha_h = \frac{1}{3}$  and  $s_h = 1$ , while dotted blue lines show the functions for  $\alpha_h = \frac{1}{3}$  and  $s_h \rightarrow 0$ . In addition to visualizing the definitions from our previous work, the dashed yellow lines illustrate the proposed functions for  $\alpha_h = \frac{1}{2}$  and  $s_h = 1$ .

$$h_i(x) = P(x | d, m_i)$$

$$\begin{cases} h_f(x) = \frac{1}{2} - a_n \tanh(x/2s_n) \\ h_b(x) = \frac{1}{2} + a_n \tanh(x/2s_n) \end{cases}$$

$$a_n \in [0, 0.5] \\ s_n \in \mathbb{R}^+$$

$$P(m_f) = P(m_b) = a_n \quad \star a_n \text{ 含义}$$

$$P(m_n) = 1 - 2a_n$$

Suppose  $m \in \{m_f, m_b, m_n\}$

$P(m_f) = P(m_b) = a_n$ , and the probability of a noise  $m_n$   
 $\uparrow$  noise  
 $[0, 0.5]$

$$P(m_n) = 1 - 2a_n$$

$$\therefore P(m_i | y) = \frac{P(y | m_i) P(m_i)}{\sum_{j \in \{f, b, n\}} P(y | m_j) P(m_j)} \quad i \in \{f, b, n\}$$

We define the conditional color probability given the noise model:

$$P(y | m_n) = \frac{1}{2} (P(y | m_f) + P(y | m_b))$$

$$\therefore P(m_i | y) = \frac{P(y | m_i) P(m_i)}{P(y | m_n) P(m_n) + P(y | m_f) P(m_f) + P(y | m_b) P(m_b)}$$

$$= \frac{P(y | m_i) P(m_i)}{P(y | m_n) (1 - 2a_n) + P(y | m_f) a_n + P(y | m_b) a_n}$$

$$= \frac{P(y | m_i) P(m_i)}{P(y | m_n) (1 - 2a_n) + P(y | m_f) a_n + P(y | m_b) a_n}$$

$$= \frac{P(y | m_i) P(m_i)}{\frac{1}{2} (P(y | m_f) + P(y | m_b)) (1 - 2a_n) + a_n (P(y | m_f) + P(y | m_b))}$$

$$= \frac{2 P(y | m_i) P(m_i)}{P(y | m_f) + P(y | m_b)}$$

For foreground and background:

$$P(m_i | y) = \frac{2 a_n P(y | m_i)}{P(y | m_f) + P(y | m_b)}, \quad i \in \{f, b\}$$

$$= 2 a_n P_i(r)$$

For noise:

$$P(m_n | y) = \frac{2 P(y | m_n) P(m_n)}{P(y | m_f) + P(y | m_b)} = \frac{2 \cdot \frac{1}{2} (P(y | m_f) + P(y | m_b)) P(m_n)}{P(y | m_f) + P(y | m_b)}$$

$$= P(m_n)$$

$$= 1 - 2a_n$$

$$\therefore P(d | r, y) \propto \sum_{i \in \{f, b\}} P(r | d, m_i) P(m_i | y), \quad h_i(x) = P(r | d, m_i)$$

$$P_i(r) = \frac{P(y | m_i)}{P(y | m_f) + P(y | m_b)}$$

we define  $P(r | d, m_n) = \frac{1}{2}$  for noise.

$$\therefore P(d | r, y) \propto \sum_{i \in \{f, b, n\}} P(r | d, m_i) P(m_i | y)$$

$$= 2a_n h_f(r-d) P_f(r) + 2a_n h_b(r-d) P_b(r) + \frac{1}{2} (1 - 2a_n)$$

Let

- ①  $h_f(x) = \frac{1}{2} - f(x)$
- ②  $h_b(x) = \frac{1}{2} + f(x)$
- ③  $P_f(r) + P_b(r) = 1$

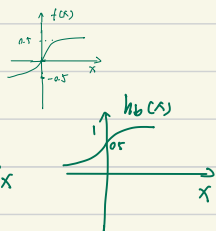
Suppose:

- ①  $h_f + h_b = 1$
- ②  $h_f$  and  $h_b$  be symmetric.

$$\therefore h_f(x) = 0.5 - f(x), \quad h_b(x) = 0.5 + f(x)$$

where  $f(x)$  is an odd function ( $-f(x) = f(-x)$ ) that lies within the interval  $[-0.5, 0.5]$  and

$$\lim_{x \rightarrow 0} f(x) = 0.5 \quad \text{and} \quad \lim_{x \rightarrow 0} f(x) = -0.5$$



$$\therefore P(d | r, y) \propto 2a_n \left( \frac{1}{2} - f(r-d) \right) P_f(r) + 2a_n \left( \frac{1}{2} + f(r-d) \right) P_b(r) + \frac{1}{2} (1 - 2a_n) (P_f(r) + P_b(r))$$

$$= a_n P_f(r) - 2a_n f(r-d) P_f(r) + a_n P_b(r) + 2a_n f(r-d) P_b(r) + \frac{1}{2} (P_f(r) + P_b(r)) - a_n (P_f(r) + P_b(r))$$

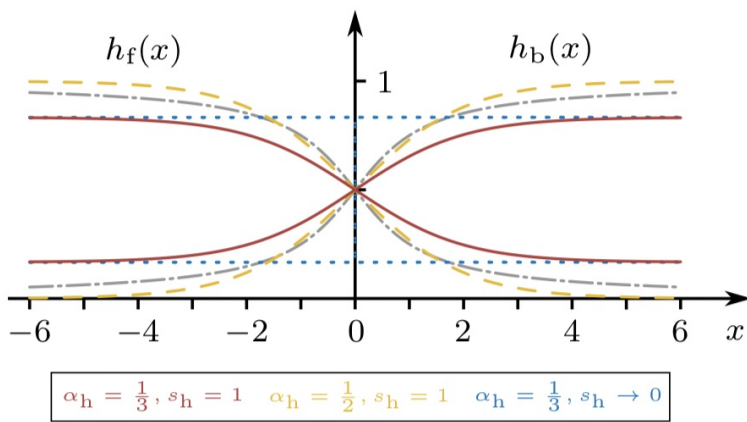
$$= \left( \frac{1}{2} - 2a_n f(r-d) \right) P_f(r) + \left( \frac{1}{2} + 2a_n f(r-d) \right) P_b(r)$$



$$\therefore f(x) = \frac{1}{2} \tanh\left(\frac{x}{2s_h}\right)$$

$$\therefore p(d|r, y) \propto \left(\frac{1}{2} - \alpha_h \tanh\left(\frac{r-d}{2s_h}\right)\right) p_f(r) + \left(\frac{1}{2} + \alpha_h \tanh\left(\frac{r-d}{2s_h}\right)\right) p_b(r)$$

Extending the probabilistic model with a noise model  $m_n$  with a defined constant uncertainty and using the foreground and background probabilities  $p(m_f) = p(m_b) = \alpha_h$  is equivalent to the introduction of a simple amplitude parameter  $\alpha_h$  into the smoothed step functions.



**Fig. 4** Smoothed step functions  $h_f$  and  $h_b$  that model the conditional line coordinate probabilities  $p(r | d, m_f)$  and  $p(r | d, m_b)$ . The functions  $h_f(x) = \frac{1}{2} - \frac{1}{\pi} \tan^{-1}\left(\frac{x}{s_h}\right)$  and  $h_b(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}\left(\frac{x}{s_h}\right)$  used by Zhong et al. (2020b) and Tjaden et al. (2018) are illustrated by dash-dotted gray lines. The definitions  $h_f(x) = \frac{1}{2} - \frac{1}{2} \tanh\left(\frac{x}{2s_h}\right)$  and  $h_b(x) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{x}{2s_h}\right)$  from our previous work (Stoiber et al., 2020) are shown as dashed yellow lines. In both plots, the slope parameter  $s_h = 1$  is used. For the proposed functions in Eq. (12) and (13), solid red lines correspond to  $\alpha_h = \frac{1}{3}$  and  $s_h = 1$ , while dotted blue lines show the functions for  $\alpha_h = \frac{1}{3}$  and  $s_h \rightarrow 0$ . In addition to visualizing the definitions from our previous work, the dashed yellow lines illustrate the proposed functions for  $\alpha_h = \frac{1}{2}$  and  $s_h = 1$ .

$$4. \quad \textcircled{1} p(d|w, l) \propto \prod_{r \in w} \sum_{i \in \{f, b\}} p(r|d, m_i) p(m_i | l(r))$$

$$\textcircled{2} h_i(r-d) = p(r|d, m_i) \quad \begin{cases} h_f(x) = \frac{1}{2} - \alpha_h \tanh\left(\frac{x}{2s_h}\right) \\ h_b(x) = \frac{1}{2} + \alpha_h \tanh\left(\frac{x}{2s_h}\right) \end{cases}$$

$$\textcircled{3} p_i(r) = p(m_i | l(r)) = p(m_i | y) = \frac{p(y | m_i)}{p(y | m_f) + p(y | m_b)}, \quad i \in \{f, b\}$$

$l(r) = l(c+rn) = y$

$$\therefore p(d|w, l) \propto \prod_{r \in w} h_f(r-d) p_f(r) + h_b(r-d) p_b(r)$$

To analyze the posterior probability  $P(d|w, l)$  distribution,

we assume a contour at the line center  $c$  and perfect step function  $p_i(x)$  for the pixel-wise posterior probabilities defined by:

$$\begin{cases} p_f(r) = \frac{1}{2} - \frac{1}{2} \operatorname{sgn}(r) \\ p_b(r) = \frac{1}{2} + \frac{1}{2} \operatorname{sgn}(r) \end{cases}$$

we consider infinitesimally small pixels and write the posterior probability distribution  $p(d|w, l)$  in continuous form for an infinite correspondence line:

The classical Riemann integral of a function  $f: [a, b] \rightarrow \mathbb{R}$ :

$$\int_a^b f(x) dx = \lim_{\Delta x \rightarrow 0} \sum f(x_i) \Delta x$$

Product integrals are similar, but take the limit of a product instead of the limit of sum. They can be thought of as "continuous" version of "discrete" products.

One kind of Product Integrals is Geometric Integral:

$$\prod_a^b f(x) dx = \lim_{\Delta x \rightarrow 0} \prod f(x_i) \Delta x = \exp\left(\int_a^b \ln f(x) dx\right)$$

$$\begin{aligned} P(d|w, l) &\propto \prod_{r=-\infty}^{\infty} (h_f(r-d)p_f(r) + h_b(r-d)p_b(r)) dr \\ &= \exp\left(\int_{r=-\infty}^{\infty} \ln(h_f(r-d)p_f(r) + h_b(r-d)p_b(r)) dr\right) \\ &= \exp\left(\int_{r=-\infty}^0 \ln(h_f(r-d)p_f(r) + h_b(r-d)p_b(r)) dr \right. \\ &\quad \left. + \int_{r=0}^{\infty} \ln(h_f(r-d)p_f(r) + h_b(r-d)p_b(r)) dr\right) \\ &= \exp\left(\int_{r=-\infty}^0 \ln(h_f(r-d)) dr + \int_{r=0}^{\infty} \ln(h_b(r-d)) dr\right) \end{aligned}$$

Let  $x = r - d$  ↖  $b(t)$ ,  $t =$

$$\therefore p(d|w, L) \propto \exp \left( \int_{x=-\infty}^{-d} \ln(h_f(x)) dx + \int_{x=-d}^{\infty} \ln(h_b(x)) dx \right)$$

↗  $a(t)$

$$\therefore \frac{\partial \ln(P(d|w, L))}{\partial d} = \frac{\partial \left( \int_{x=-\infty}^{-d} \ln(h_f(x)) dx \right)}{\partial d} + \frac{\partial \left( \int_{x=-d}^{\infty} \ln(h_b(x)) dx \right)}{\partial d}$$

**Leibniz Rule:** If  $\Phi(t) = \int_{a(t)}^{b(t)} f(x, t) dx$ , then

$$\frac{d\Phi}{dt} = \int_{a(t)}^{b(t)} \frac{\partial f(x, t)}{\partial t} dx + f(b(t), t) \frac{db}{dt} - f(a(t), t) \frac{da}{dt}$$

Let  $d \leftarrow t$ ,  $-d \leftarrow b(t)$ ,  $x = -\infty \leftarrow a(t)$

$\ln(h_f(x)) \leftarrow f(x, t)$

$$\therefore \frac{\partial \left( \int_{x=-\infty}^{-d} \ln(h_f(x)) dx \right)}{\partial d}$$

$$= \int_{x=-\infty}^{-d} \frac{\partial \ln(h_f(x))}{\partial d} dx - \ln(h_f(-d)) + 0$$

$$= -\ln(h_f(-d))$$

$$\therefore \frac{\partial \left( \int_{x=-d}^{\infty} \ln(h_b(x)) dx \right)}{\partial d}$$

$$= 0 + 0 + \ln(h_b(-d))$$

$$\frac{\partial \ln(P(d|w, L))}{\partial d} = -\ln(h_f(-d)) + \ln(h_b(-d))$$

$$\therefore h_f(x) = \frac{1}{2} - a_h \tanh(x/2s_h)$$

$$h_b(x) = \frac{1}{2} + a_h \tanh(x/2s_h)$$

$$\therefore \frac{\partial \ln(P(d|\omega, \mathbf{l}))}{\partial d} = -\ln\left(\frac{1}{2} - a_h \tanh\left(\frac{-d}{2s_h}\right)\right) + \ln\left(\frac{1}{2} + a_h \tanh\left(\frac{-d}{2s_h}\right)\right)$$

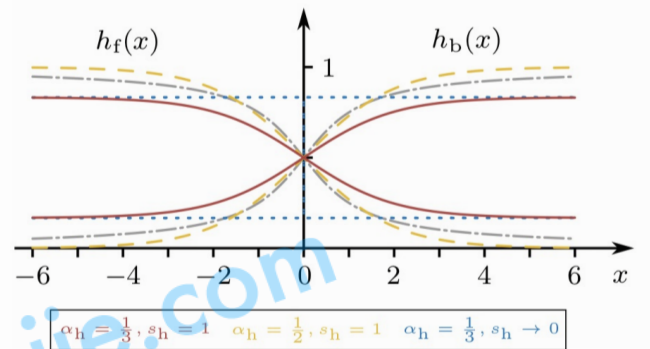
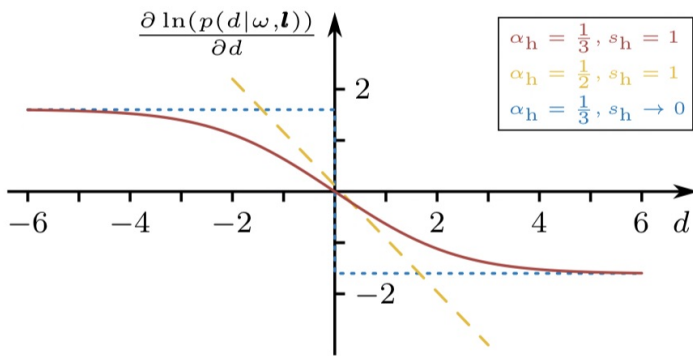
$$\star \tanh^{-1} = \operatorname{arctanh}$$

$$\therefore 2 \tanh^{-1}(x) = -\ln(1-x) + \ln(1+x)$$

$$\tanh(-x) = -\tanh(x)$$

$$\begin{aligned} \therefore \frac{\partial \ln(P(d|\omega, \mathbf{l}))}{\partial d} &= -\ln\left(\frac{1}{2}\left(1 - 2a_h \tanh\left(\frac{-d}{2s_h}\right)\right)\right) \\ &\quad + \ln\left(\frac{1}{2}\left(1 + 2a_h \tanh\left(\frac{-d}{2s_h}\right)\right)\right) \\ &= -\ln\left(1 - 2a_h \tanh\left(\frac{-d}{2s_h}\right)\right) + \ln\left(1 + 2a_h \tanh\left(\frac{-d}{2s_h}\right)\right) \\ &= 2 \tanh^{-1}\left(2a_h \tanh\left(\frac{-d}{2s_h}\right)\right) \end{aligned}$$

$$\frac{\partial \ln(P(d|\omega, \mathbf{l}))}{\partial d} = -2 \tanh^{-1}\left(2a_h \tanh\left(\frac{d}{2s_h}\right)\right)$$



**Fig. 5** First-order derivatives of the log-posterior with respect to the contour distance  $d$  for different slope and amplitude parameters  $s_h$  and  $\alpha_h$ . The solid red line shows the derivative for  $\alpha_h = \frac{1}{3}$  and  $s_h = 1$ , which yields a function with a smooth transition from an upper bound to a lower bound. The dashed yellow line shows the function for  $\alpha_h = \frac{1}{2}$  and  $s_h = 1$ . This produces a linear first-order derivative. Finally, using  $\alpha_h = \frac{1}{3}$  and  $s_h \rightarrow 0$  results in a perfect step function illustrated by the dotted blue line.

① If  $a_n = \frac{1}{2}$

PDF:  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\frac{\partial \ln(P(d|w, L))}{\partial d} = -\frac{d}{s_n} = \frac{\partial \ln(N(d|0, s_n))}{\partial d}$$

$$N(d|0, s_n) = \frac{1}{\sqrt{2\pi s_n}} \exp\left(-\frac{d^2}{2s_n}\right)$$

The first-order derivative of the log-likelihood of  $P(d|w, L)$  is equal to the first-order derivative of the logarithm of the normal distribution  $N(d|0, s_n)$ .

$$\therefore P(d|w, L) \propto \exp\left(-\frac{d^2}{2s_n}\right)$$

Because the PDF of  $P(d|w, L)$  that integrates to one, the final solution can only be the Gaussian distribution:

$$\text{If } a_n = \frac{1}{2}, \ln P(d|w, L) = \frac{1}{\sqrt{2\pi s_n}} \exp\left(-\frac{d^2}{2s_n}\right)$$

② If  $s_n \rightarrow 0$

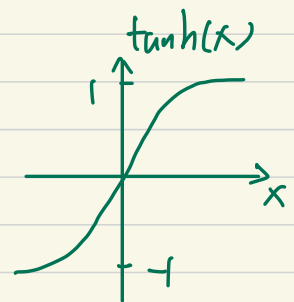
$$\frac{\partial \ln(P(d|w, L))}{\partial d} = -2 \tanh^{-1}\left(2a_n \tanh\left(\frac{d}{2s_n}\right)\right)$$

$$\Rightarrow = -2 \tanh^{-1}(2a_n) \operatorname{sgn}(d)$$

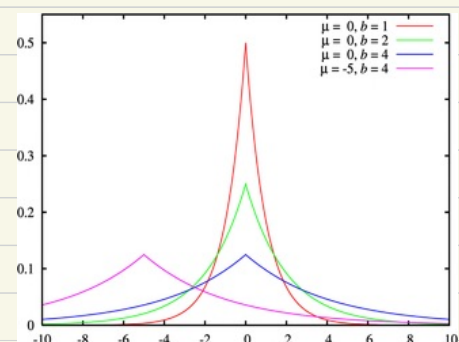
$$\begin{aligned} P(d|w, L) &\propto \exp\left(\int \frac{\partial \ln(P(d|w, L))}{\partial d} \Delta d\right) \\ &= \exp\left(\int -2 \tanh^{-1}(2a_n) \operatorname{sgn}(d) \Delta d\right) \\ &= \exp(-2 \tanh^{-1}(2a_n) |d|) + C \\ &= \exp\left(-\frac{|d|}{\frac{1}{2 \tanh^{-1}(2a_n)}}\right) + C \end{aligned}$$

Laplace distribution:

$$f(x|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right) = \frac{1}{2b} \begin{cases} \exp\left(-\frac{\mu-x}{b}\right), & \text{if } x < \mu \\ \exp\left(-\frac{x-\mu}{b}\right), & \text{if } x \geq \mu \end{cases}$$



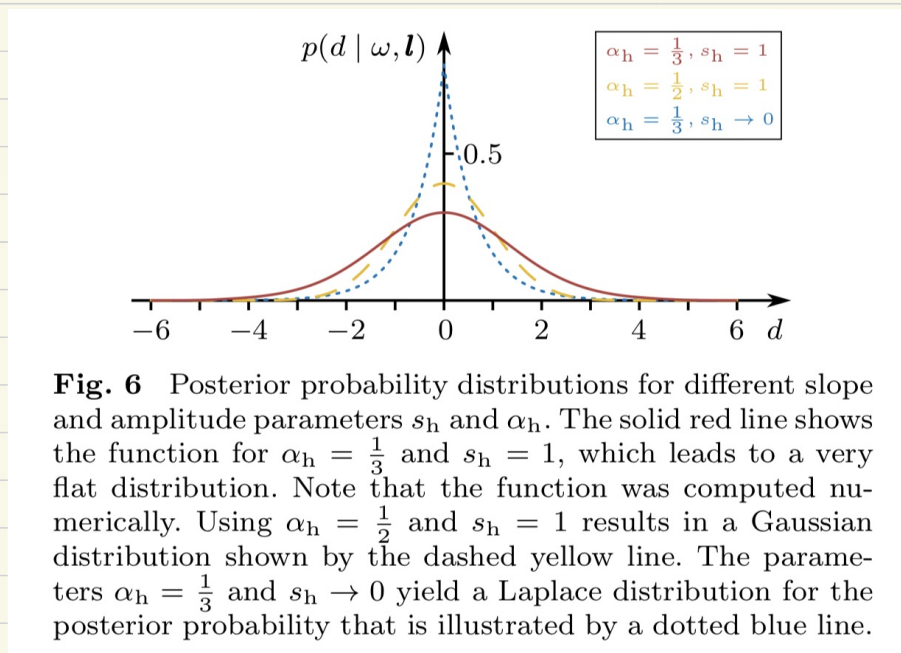
$$\frac{d|x|}{dx} = \operatorname{sgn}(x) \text{ for } x \neq 0$$



with the exception of a constant scaling factor, the function is equal to a Laplace distribution, which again is a valid PDF that integrates to one.

$$\text{If } s_h \rightarrow 0, \quad p(d|w, l) = \frac{1}{2b} \left( -\frac{|d|}{b} \right), \quad b = \frac{1}{2 \tanh'(2a_h)}$$

$$\text{The variance of Laplace distribution: } 2b^2 = \frac{1}{2 (\tanh'(2a_h))^2}$$



The slope parameter  $s_h$  controls local uncertainty, allowing multiple values  $d$  to be almost equally likely.

The amplitude parameter  $a_h$  controls the size of the peak compared to its surroundings, thereby controlling global uncertainty.

In the perfect situation,

$$\text{If } a_h = \frac{1}{2}, \quad \ln P(d|w, l) = \mathcal{N}(d|0, s_h) = \frac{1}{\sqrt{2\pi}s_h} \exp\left(-\frac{d^2}{2s_h}\right)$$

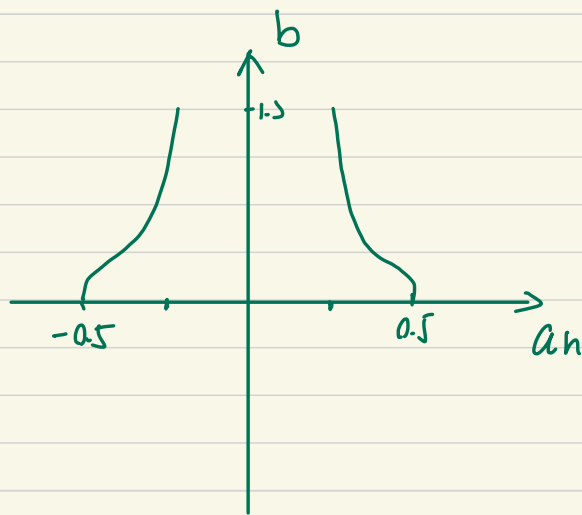
$$G_{\text{gauss}}^2 = s_h$$

$$\text{If } s_h \rightarrow 0, \quad p(d|w, l) = \frac{1}{2b} \left( -\frac{|d|}{b} \right), \quad b = \frac{1}{2 \tanh'(2a_h)}$$

$$G_{\text{Laplace}}^2 = 2b^2 = \frac{1}{2} \tanh^2(2a_h)$$

$\therefore$  In the normal situation, the variance of  $\mathcal{N}(d_i | \mu_i, G_i^2)$ ,  $G_i^2$  should be larger than  $G_{\text{gauss}}^2$  and  $G_{\text{Laplace}}^2$ .





5. 3D model points :  $X = [X \ Y \ Z]^T$   
 $\tilde{X} = [X \ Y \ Z \ 1]^T$

$$x = \pi(X) = \begin{bmatrix} \frac{X}{Z} f_x + p_x \\ \frac{Y}{Z} f_y + p_y \end{bmatrix}$$

$$X = \pi^{-1}(x, dz) = dz \begin{bmatrix} \frac{x - p_x}{f_x} \\ \frac{y - p_y}{f_y} \\ 1 \end{bmatrix}$$

$${}_c \tilde{X} = {}_c T_m \quad {}_m \tilde{X} = \begin{bmatrix} {}_c R_m & {}_c t_m \\ 0 & 1 \end{bmatrix} {}_m \tilde{X}$$

$$R = \exp([r]_x) = I + [r]_x + \frac{1}{2!} [r]_x^2 + \frac{1}{3!} [r]_x^3 + \dots$$

$${}_c \tilde{X}(\theta) = \begin{bmatrix} {}_c R_m & {}_c t_m \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I + [\theta_r] & \theta_t \\ 0 & 1 \end{bmatrix} {}_m \tilde{X}$$

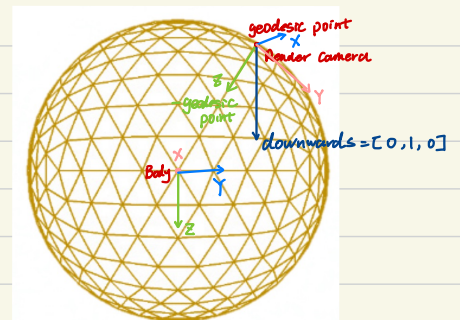
$$\theta_r \in \mathbb{R}^3, \theta_t \in \mathbb{R}^3, \theta^T = [\theta_r^T \theta_t^T]$$

6. The 3D geometry is rendered from a number of  $n_v$  viewpoints all around the object. Virtual cameras are placed on the vertices of a geodesic grid. For each rendering,  $n_c$  points  $x_i \in \mathbb{R}^2$  are randomly sampled from the contour.  $n_i \in \mathbb{R}^2$  is the normal is computed,  $\|n_i\|_2 = 1$ . 3D vectors with respect to the model frame:

$${}_m \tilde{X}_i = {}_m T_c \pi^{-1}(x_i, d(z_i))$$

$${}_m N_i = {}_m R_c \begin{bmatrix} n_i \\ 0 \end{bmatrix}$$

$$\text{orientation vector } M^v = {}_m R_c \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



(c) Geodesic Grid

Given a specific pose  ${}_m R_c, {}_c t_m$ . the closest precomputed view  $i_v$ :

$$i_v = \arg \max_{i \in \{1, \dots, n_v\}} (M^v)^T {}_m R_c {}_c t_m)$$

7. 3D model points and normal vectors from the closest view of the sparse view are projected into the image:

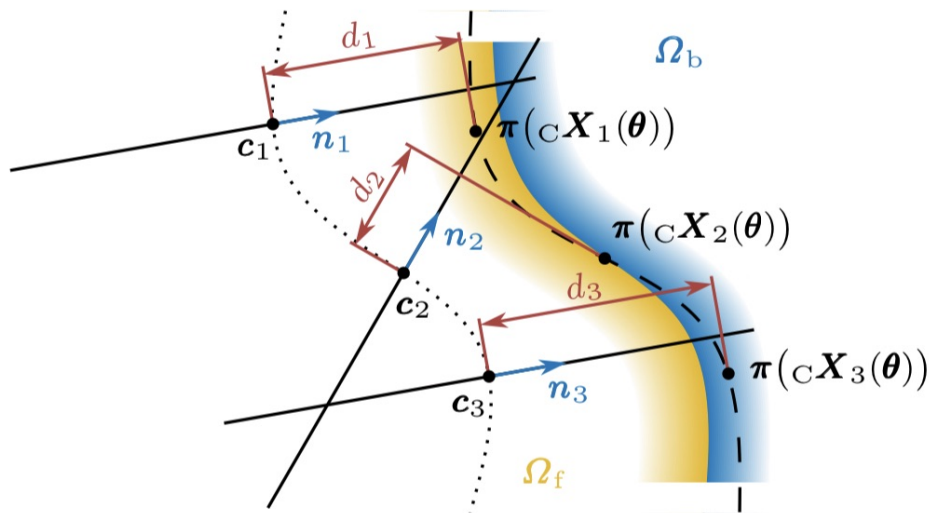
$$c_i = \pi({}_c T_m {}_m \tilde{X}_i)$$

$$n_i \propto ({}_c R_m {}_m N_i)_{2 \times 1}, \quad \|n_i\|_2 = 1$$

Contour distances are calculated as the distances along normal vectors  $n_i$  from correspondence line centers  $c_i$  to projected 3D model points  $X_i$ :

$$d_i(\theta) = n_i^T (\pi({}_c \tilde{X}_i(\theta)) - c_i)$$

$$\leftarrow {}_c \tilde{X}(\theta) = \begin{bmatrix} {}_c R_m & {}_c t_m \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I + [\theta_r] & \theta_t \\ 0 & 1 \end{bmatrix} {}_m \tilde{X}$$



**Fig. 8** Correspondence lines defined by a center  $\mathbf{c}_i$  and a normal vector  $\mathbf{n}_i$ . Variated contour distances  $d_i$  are measured along the correspondence lines from the centers  $\mathbf{c}_i$  to the projected 3D model points  $\pi(\mathbf{C}\mathbf{X}_i(\theta))$  that depend all on the same pose variation  $\theta$ . The object contour of the original pose estimate, which was used to define the correspondence lines, is indicated by a dotted line. The current estimate of the contour that depends on the pose variation vector  $\theta$  is shown by a dashed line. The ground truth segmentation that we try to estimate is given by the foreground region  $\Omega_f$  in yellow and the background region  $\Omega_b$  in blue. Note that while contours are illustrated as continuous lines, in our method, they are represented by points and normal vectors from the closest view of the sparse viewpoint model.

Assuming a number of  $n_c$  independent correspondence lines.

$$P(\theta | \mathcal{D}) \propto \prod_{i=1}^{n_c} P(d_{si}(\theta) | w_{si}, l_s) \quad (29), \quad \mathcal{D} : \text{the data from all correspondence lines}$$

$$\textcircled{1} P(d_s | w_s, l_s) \propto \prod_{r_s \in w_s} \sum_{i \in \{f, b\}} P(r_s | d_s, m_i) P(m_i | l_s(r_s)) \quad \textcircled{7} \textcircled{10}$$

$$\textcircled{2} P(r_s | d_s, m_i) = h_i(r_s - d_s) = \begin{cases} \frac{1}{2} - a_n \tanh h[(r_s - d_s)/2s_h] & , i = f \\ \frac{1}{2} + a_n \tanh h[(r_s - d_s)/2s_h] & , i = b \end{cases}$$

$$\textcircled{3} r_s = (r - \Delta r) \frac{\bar{n}}{s}, \quad \text{unknown, } l(r) = l(c + m)$$

$s \in \mathbb{N}^+$ : the number of pixels combined into a segment.  
 $\bar{n} = \max(|n_x|, |n_y|)$

$$d_s = (d - \Delta r) \frac{\bar{n}}{s} \quad \textcircled{4}$$

$$(4) d(\theta) = n^T (\underbrace{\pi(cX(\theta))}_{(5)}) - c$$

$$(5) cX(\theta) = \begin{bmatrix} cR_m & ct_m \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + [\theta_r] & \theta_t \\ 0 & 1 \end{bmatrix} m\tilde{x}$$

$$(3) p(m_i | \underbrace{S}_{(6)}) = \frac{\prod_{y \in S} p(y | m_i)}{\prod_{y \in S} p(y | m_f) + \prod_{y \in S} p(y | m_b)}$$

$\uparrow$   
 a set of the  
 closest  $s$  pixel values  $y$

$\uparrow$   
 (7)

$$(6) p(m_i | y) = \frac{p(y | m_i)}{p(y | m_f) + p(y | m_b)}, \quad i \in \{f, b\}$$

(7)  $S = \{s(r_s)\}$  is a set-valued function that maps from the scaled line coordinate  $r_s$  to the segment  $S$ , which is a set of the closest  $s$  pixel values  $y$ .

$P(\theta | D)$  describes how well the current pose  $cX(\theta)$  estimate explains the segmentation of the image into a foreground region and background region.

Pose  $cX(\theta) \longrightarrow$  contour distance:  $d(\theta) \longrightarrow P(\theta | D)$

$$\begin{aligned} \ln(P(\theta | D)) &\propto \sum_i^{\text{enc}} \ln[p(d_{s_i}(\theta) | w_{s_i}, l_s)] \\ &= \sum_i^{\text{enc}} \sum_{r_s}^{r_s=w_s} \ln \left[ \sum_{j \in \{f, b\}} p(r_s | d_s, m_j) p(m_j | \{s(r_s)\}) \right] \end{aligned}$$

www.luohanjie.com

8. To maximize the joint posterior probability  $\ln(P(\theta | D))$ , we estimate the variation vector  $\theta$  and iteratively update the pose.

$$\hat{\theta} = (-H + \begin{bmatrix} \lambda_r I_3 & 0 \\ 0 & \lambda_t I_3 \end{bmatrix})^{-1} g$$

$$g = \frac{\partial \ln(P(D|\theta))}{\partial \theta} \bigg|_{\theta=0}, \quad H = \frac{\partial^2 \ln(P(D|\theta))}{\partial \theta^2} \bigg|_{\theta=0}$$

$$\therefore cT_M = cT_M \begin{bmatrix} \exp([\hat{\theta}_r]_{\times}) & \hat{\theta}_t \\ 0 & 1 \end{bmatrix}$$

$$ds = (d - \Delta r) \frac{\bar{n}}{s}, \quad d(\theta) = n^T (\pi(c \chi(\theta)) - c)$$

$$\textcircled{1} \quad g^T = \sum_{i=1}^{n_c} \frac{\partial \ln(p(ds_i | w_i, l_i))}{\partial ds_i} \quad \begin{matrix} \downarrow \\ \frac{\partial ds_i}{\partial c \chi_i} \end{matrix} \quad \begin{matrix} \frac{\partial c \chi_i}{\partial \theta} \end{matrix} \bigg|_{\theta=0}$$

$\begin{matrix} 6 \times 1 & 1 \times 3 & 3 \times 6 \end{matrix}$

$$\therefore c \chi(\theta) = \begin{bmatrix} c R_m & c t_m \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + [\theta_r] & \theta_t \\ 0 & 1 \end{bmatrix} m \tilde{x}$$

$$\theta_r \in \mathbb{R}^3, \quad \theta_t \in \mathbb{R}^3, \quad \theta^T = [\theta_r^T \quad \theta_t^T]$$

$$\delta \xi = [\delta \rho, \delta \phi]^T$$

$$\frac{\partial \langle TP \rangle}{\partial \delta \xi} = \lim_{\delta \xi \rightarrow 0} \frac{\exp(\xi^T) \exp(\delta \xi^T) P - \exp(\xi^T) P}{\delta \xi}$$

$$\approx \lim_{\delta \xi \rightarrow 0} \frac{\exp(\xi^T) (1 + \delta \xi^T) P - \exp(\xi^T) P}{\delta \xi}$$

$$= \lim_{\delta \xi \rightarrow 0} \frac{\exp(\xi^T) \delta \xi^T P}{\delta \xi} \quad [s \phi^T P + s \rho]$$

$$= \lim_{\delta \xi \rightarrow 0} \frac{\begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} s \phi^T & s \rho \\ 0 & 1 \end{bmatrix} \begin{bmatrix} P \\ 1 \end{bmatrix}}{\delta \xi}$$

$$= \lim_{\delta \xi \rightarrow 0} \frac{\begin{bmatrix} R s \phi^T P + R s \rho + t \\ 1 \end{bmatrix}}{\delta \xi} = A$$

$$= \begin{bmatrix} \frac{A}{s \rho} & \frac{A}{s \phi} \end{bmatrix} \quad R s \phi^T P = -R P^T s \phi$$

$$= \begin{bmatrix} R & -R P^T \\ 0 & 0 \end{bmatrix}$$

$$\frac{\partial c \chi_i}{\partial \theta} = c R_m \begin{bmatrix} -[m \chi_i]_x & I_3 \end{bmatrix} = \begin{bmatrix} m \chi_i \times c R_m, & c R_m \end{bmatrix}$$

$\begin{matrix} 3 \times 3 & 3 \times 6 \end{matrix}$

$$\begin{matrix} [1 \times 3] \\ [x^T] \end{matrix} \begin{bmatrix} -[m \chi_i]_x & I_3 \end{bmatrix} \begin{matrix} 3 \times 3 \\ 3 \times 5 \end{matrix}$$



$$ds = (d - \Delta r) \frac{\bar{n}}{3}, \quad d(\theta) = n^T (\pi(cX(\theta)) - c)$$

$$x = \pi(X) = \begin{bmatrix} \frac{X}{Z} f_x + p_x \\ \frac{Y}{Z} f_y + p_y \end{bmatrix}$$

$$d_s = \left[ \left( c \frac{X}{Z} f_x + p_x - (x) n_x + \left( c \frac{Y}{Z} f_y + p_y - (y) n_y \right) - \Delta r \right) \frac{\bar{n}}{3} \right. \\ \left. + \frac{1}{Z} (cX + x n_x + cY + y n_y) \right]$$

$$\therefore \frac{\partial ds_i}{\partial c x_i} = \frac{\bar{n}_i}{3} \left[ \frac{f_x n_{xi}}{c Z_i} + \frac{f_y n_{yi}}{c Z_i} - \frac{c x_i f_x n_{xi} + c y_i f_y n_{yi}}{c Z_i^2} \right]$$

1x3

② Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   
 $g: \mathbb{R} \rightarrow \mathbb{R}$   
 and  $h(x) = g(f(x))$

$$\begin{aligned} \therefore \nabla^2 h(x) &= \nabla [g'(f(x)) \nabla f(x)] \\ &= \nabla [g'(f(x))] \nabla f(x) + g'(f(x)) \nabla [\nabla f(x)] \\ &= g''(f(x)) \nabla f(x) \nabla f(x)^T + g'(f(x)) \nabla^2 f(x) \end{aligned}$$

$n \times 1 \quad 1 \times n \quad n \times n$

Let  $g: \mathbb{R} \rightarrow \mathbb{R}, \quad \ln(p(ds_i | w_{si}, l_{si})) = g(ds_i)$   
 $f: \mathbb{R}^n \rightarrow \mathbb{R}, \quad f(\theta) = ds_i$

$b \times 1$

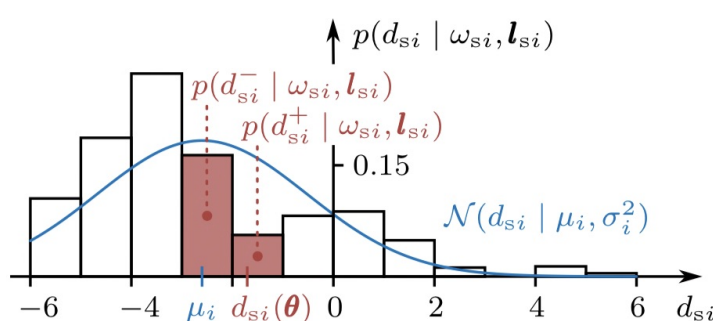
$$\begin{aligned} \therefore H &= \sum_{i=1}^{nc} \frac{\partial^2 \ln(p(ds_i | w_{si}, l_{si}))}{\partial ds_i^2} \underbrace{\left( \frac{\partial ds_i}{\partial \theta} \right)^T}_{b \times 1} \left( \frac{\partial ds_i}{\partial \theta} \right) \\ &\quad + \frac{\partial \ln(p(ds_i | w_{si}, l_{si}))}{\partial ds_i} \left( \frac{\partial^2 ds_i}{\partial \theta^2} \right) \\ &\approx \sum_{i=1}^{nc} \frac{\partial^2 \ln(p(ds_i | w_{si}, l_{si}))}{\partial ds_i^2} \left( \frac{\partial ds_i}{\partial \theta} \right)^T \left( \frac{\partial ds_i}{\partial \theta} \right) \end{aligned}$$

① For global optimization :  $\ln(p(d_{si} | w_{si}, l_{si})) \sim \mathcal{N}(d_{si} | \mu_i, \sigma_i^2)$

The required mean  $\mu_i$  and standard deviation  $\sigma_i$  are thereby estimated from a set of discretized contour distances  $d_{si}$  and their corresponding probability values.

$$\therefore \frac{\partial \ln(p(d_{si} | w_{si}, l_{si}))}{\partial d_{si}} \approx -\frac{1}{\sigma_i^2} (d_{si} - \mu_i) \quad \text{PDF: } f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\frac{\partial^2 \ln(p(d_{si} | w_{si}, l_{si}))}{\partial d_{si}^2} \approx -\frac{1}{\sigma_i^2} \quad \text{If } a_n = \frac{1}{n}, \ln P(d | w, l) = \frac{1}{N \sigma_{sh}} \exp(-\frac{d^2}{2\sigma_{sh}^2})$$



**Fig. 9** Discrete posterior probability distribution with noisy measurements. For global optimization, the distribution is approximated by a normal distribution  $\mathcal{N}(d_{si} | \mu_i, \sigma_i^2)$ . The normal distribution and its mean  $\mu_i$  are illustrated in blue. In the case of local optimization, only two discrete probability values that are closest to the current estimate of the contour distance  $d_{si}(\theta)$  are considered. The two discrete probability values  $p(d_{si}^- | w_{si}, l_{si})$  and  $p(d_{si}^+ | w_{si}, l_{si})$ , which are used to approximate the first-order derivative, are colored in red.

The approximated derivatives direct the optimization towards the mean  $\mu_i$  using the variance  $\sigma_i^2$  to consider uncertainty. In the real world, the mean does not exactly coincide with the maximum. Using the approximation has the advantage of fast convergence and that the opt avoids local minima resulting from invalid pixel-wise posteriors and image noise.

② Once the opt is closer to the maximum, the algorithm switches to local opt.

We use the probability  $d_{si}^-$  and  $d_{si}^+$  that are closest to the current estimate  $d_{si}(\theta)$  and:

$$\frac{\partial \ln(p(d_{si} | w_{si}, l_{si}))}{\partial d_{si}} \approx \frac{d_{si}}{\sigma_i^2} \ln \left( \frac{p(d_{si}^+ | w_{si}, l_{si})}{p(d_{si}^- | w_{si}, l_{si})} \right)$$

$$\frac{\partial^2 \ln(P(d_{si}|w_{si}, l_{si}))}{\partial d_{si}^2} \approx -\frac{1}{\sigma_i^2}$$

$$\text{If } E(\theta) = \sum_{i=1}^{n_s} \frac{\sigma_i}{\sigma_i^2} \ln(P(d_{si}|\theta) | w_{si}, l_{si})$$

The weighting  $\frac{1}{\sigma_i^2}$  improves robustness because correspondence lines with high uncertainty are considered less important.

The step size  $\sigma_s$  helps to balance the weight and specifies how far the opt proceeds, directly scaling the  $\theta$ .

$$\therefore g^T = \sum_{i=1}^{n_s} \frac{\sigma_i}{\sigma_i^2} \frac{\partial \ln(P(d_{si}|w_{si}, l_{si}))}{\partial d_{si}} \frac{\partial d_{si}}{\partial \theta} \Big|_{\theta=0}$$

$$H \approx \sum_{i=1}^{n_s} \frac{\sigma_i}{\sigma_i^2} \frac{\partial^2 \ln(P(d_{si}|w_{si}, l_{si}))}{\partial d_{si}^2} \left( \frac{\partial d_{si}}{\partial \theta} \right)^T \left( \frac{\partial d_{si}}{\partial \theta} \right) \Big|_{\theta=0}$$

Central difference:

$$\begin{aligned} \therefore \frac{\partial \ln(P(d_{si}|w_{si}, l_{si}))}{\partial d_{si}} &\approx \frac{\ln(P(d_{si}^+ | w_{si}, l_{si})) - \ln(P(d_{si}^- | w_{si}, l_{si}))}{\Delta d_{si} = 1} \\ &= \ln \left( \frac{P(d_{si}^+ | w_{si}, l_{si})}{P(d_{si}^- | w_{si}, l_{si})} \right) \end{aligned}$$

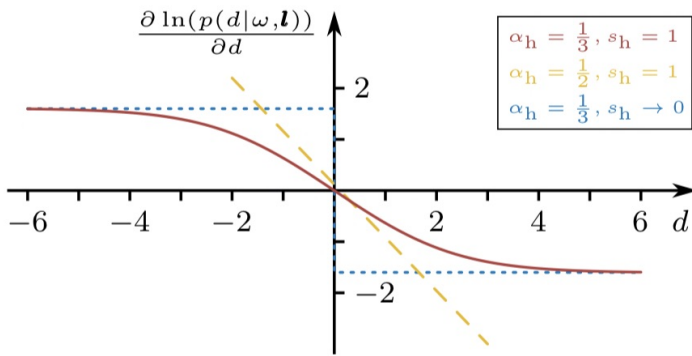
$$\therefore \frac{\partial E(\theta)}{\partial \theta} = \frac{1}{\sigma_i^2} \ln \left( \frac{P(d_{si}^+ | w_{si}, l_{si})}{P(d_{si}^- | w_{si}, l_{si})} \right) \left( \frac{\partial d_{si}}{\partial \theta} \right)^T \left( \frac{\partial d_{si}}{\partial \theta} \right)$$

$$\text{Let } \frac{\partial^2 \ln(P(d_{si}|w_{si}, l_{si}))}{\partial d_{si}^2} \approx \frac{1}{\sigma_i^2}$$

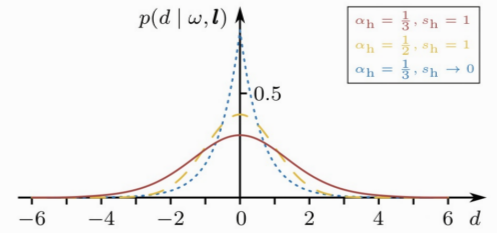
$$\therefore H \approx \sum_{i=1}^{n_s} \left( \frac{1}{\sigma_i^2} \right) \left( \frac{\partial d_{si}}{\partial \theta} \right)^T \left( \frac{\partial d_{si}}{\partial \theta} \right) \Big|_{\theta=0}$$



same as global opt



**Fig. 5** First-order derivatives of the log-posterior with respect to the contour distance  $d$  for different slope and amplitude parameters  $s_h$  and  $\alpha_h$ . The solid red line shows the derivative for  $\alpha_h = \frac{1}{3}$  and  $s_h = 1$ , which yields a function with a smooth transition from an upper bound to a lower bound. The dashed yellow line shows the function for  $\alpha_h = \frac{1}{2}$  and  $s_h = 1$ . This produces a linear first-order derivative. Finally, using  $\alpha_h = \frac{1}{3}$  and  $s_h \rightarrow 0$  results in a perfect step function illustrated by the dotted blue line.



**Fig. 6** Posterior probability distributions for different slope and amplitude parameters  $s_h$  and  $\alpha_h$ . The solid red line shows the function for  $\alpha_h = \frac{1}{3}$  and  $s_h = 1$ , which leads to a very flat distribution. Note that the function was computed numerically. Using  $\alpha_h = \frac{1}{2}$  and  $s_h = 1$  results in a Gaussian distribution shown by the dashed yellow line. The parameters  $\alpha_h = \frac{1}{3}$  and  $s_h \rightarrow 0$  yield a Laplace distribution for the posterior probability that is illustrated by a dotted blue line.

① For  $\alpha_h = \frac{1}{2}, s_h = 1$  where  $\frac{\partial \ln(p(d|\omega, l))}{\partial d}$  is a linear first-order derivative. It leads to the estimation of the weighted mean over all correspondence lines.

$$g^T = \sum_{i=1}^{n_c} \frac{\partial \ln(p(d_i | w_i, l_i))}{\partial d_i} \frac{\partial d_i}{\partial c x_i} \frac{\partial c x_i}{\partial \theta} \bigg|_{\theta=0}$$

② For  $\alpha_h = \frac{1}{3}, s_h \rightarrow 0$ ,  $\frac{\partial \ln(p(d|\omega, l))}{\partial d}$  is a binary derivatives guide the opt toward the weighted median.

9.

### Algorithm 1 Tracking Step

```

1: Update camera image
2: for  $i = 1, 2, \dots, 7$  do  $S = [5, 2, 1, 1, 1]$ 
3:   Optional: Render occlusion mask
4:   Find closest view of the sparse viewpoint model
5:   Define correspondence lines in the image
6:   Compute discrete distributions  $p(d_{si} | \omega_{si}, l_{si})$ 
7:   for  $j = 1, 2$  do
8:     Calculate gradient  $g$  and Hessian  $H$ 
9:     Estimate variation  $\hat{\theta}$  and update pose  ${}^cT_M$ 
10:  end for
11: end for
12: Update color histograms  $p(y | m_f)$  and  $p(y | m_b)$ 

```

1. continuous distances from the sparse viewpoint model are used to reject correspondence lines with distances that are below 6 segments.

2. an offset of one pixel, the first 18 pixels are considered in both the positive and negative direction of the normal vector.

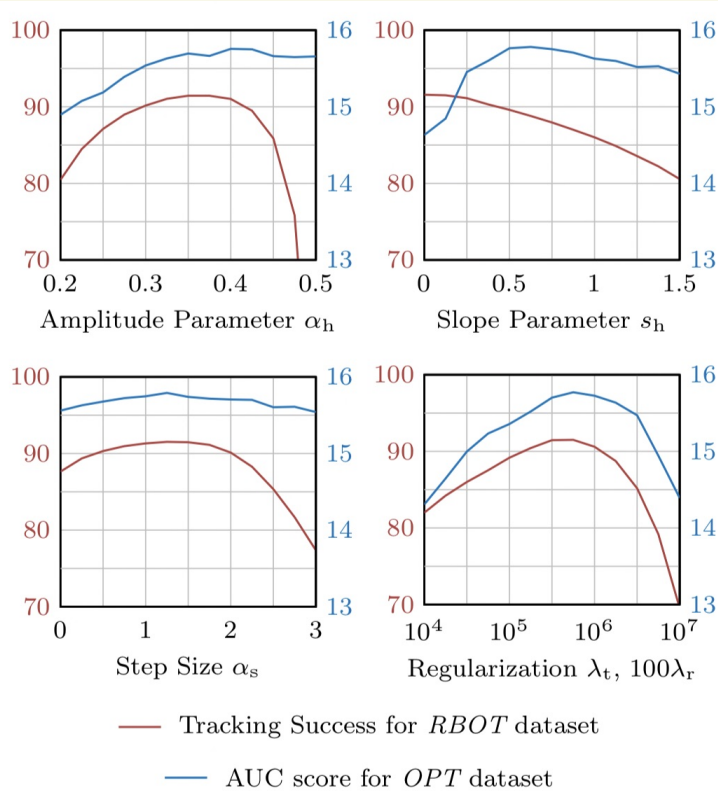
3. the posterior probability distribution  $p(d_{si} | \omega_{si}, l_{si})$  is evaluated at 12 discrete values  $d_{si} \in \{-5.5, -4.5, \dots, 5.5\}$ . *distribution-num = 12*

4. In the calculation, we use 8 precomputed values for the smoothed step functions  $h_f$  and  $h_b$ , corresponding to  $x \in \{-3.5, -2.5, \dots, 3.5\}$ . *function-num = 8*  
 $x = v_{si} - d_{si}$

5. Also, a minimal offset  $\Delta r_i$  is chosen such that the line coordinates  $r_i$  point to pixel centers while the scaled line coordinates  $r_{si}$  ensure matching values for  $x = r_{si} - d_{si}$ . In our case, this means that  $r_{si} \in \mathbb{Z}$ .

6. For the first iteration, the global optimization is used to quickly converge towards a rough pose estimate. In the second iteration, the local optimization is employed to refine this pose, using a step size of  $\alpha_s = 1.3$ . As regularization parameters, we use  $\lambda_r = 5000$  and  $\lambda_t = 500000$ .

## 10. Parameter Analysis



**Fig. 13** Average tracking success for the RBOT dataset and average AUC score for the OPT dataset over different values of the amplitude parameter  $\alpha_h$ , slope parameter  $s_h$ , step size  $\alpha_s$ , and the rotational and translational regularization parameters  $\lambda_r$  and  $\lambda_t$ . For the evaluation of the regularization parameters, we set  $\lambda_t = 100\lambda_r$ .

$a_n$ : a constant level of noise

$s_n$ : local uncertainty

RBOT:  $a_n = 0.36$   
 $s_n \rightarrow 0$

OPT:  $a_n = 0.42$   
 $s_n = 0.5$   
 $\lambda_r = 500000$  for soda

$a_s = 1.3$

$\lambda_t = 100\lambda_r$ ,  
 while  $\lambda_r = \lambda_t$  for soda

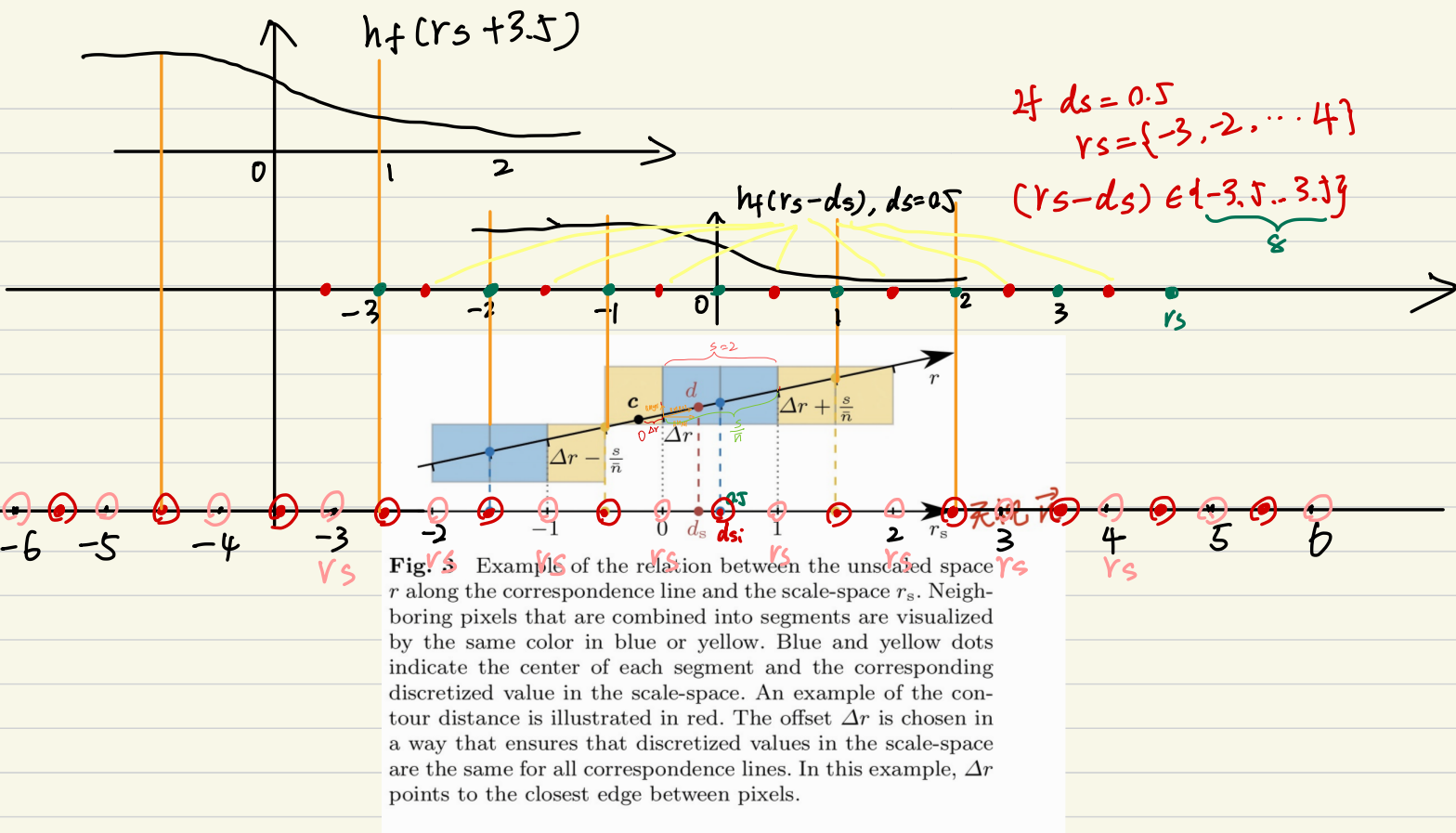
If  $\lambda_r, \lambda_t$  too small, the opt is unstable for directions in which no or very little information is available.

If  $\lambda_r, \lambda_t$  too large, the opt is slowed down and the final pose might not be reached.

Use regularization parameters that are in the same order of magnitude as the maximum rotational and translational diagonal elements of the Hessian matrix.

www.luohanjin.com

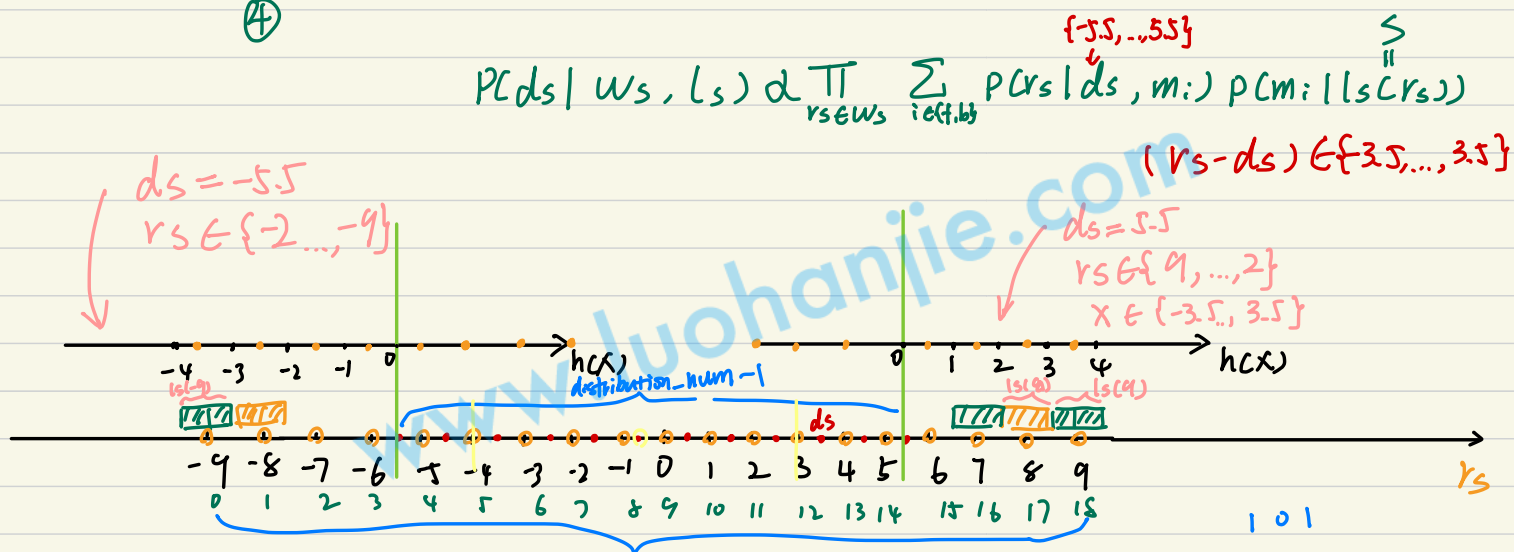




$$\textcircled{2} P(r_s | d_s, m_i) = h_i(r_s - d_s) = \begin{cases} \frac{1}{2} - a_n \tanh h[(r_s - d_s)/2s_n] \\ \frac{1}{2} + a_n \tanh h[(r_s - d_s)/2s_n] \end{cases}$$

$$\textcircled{3} r_s = (r - \Delta r) \frac{\bar{n}}{s}, \quad \begin{matrix} \swarrow \text{different from } s \\ s \in \mathbb{N}^+: \text{the number of pixels combined into a segment.} \\ \bar{n} = \max(|n_x|, |n_y|) \end{matrix}$$

$$d_s = (\underbrace{d - \Delta r}_{\textcircled{4}}) \frac{\bar{n}}{s}$$

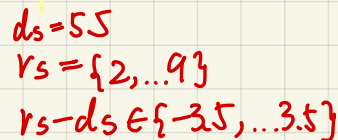


$$\text{line\_length\_segment} = (\text{function\_num} - 1) + (\text{distribution\_num} - 1)$$

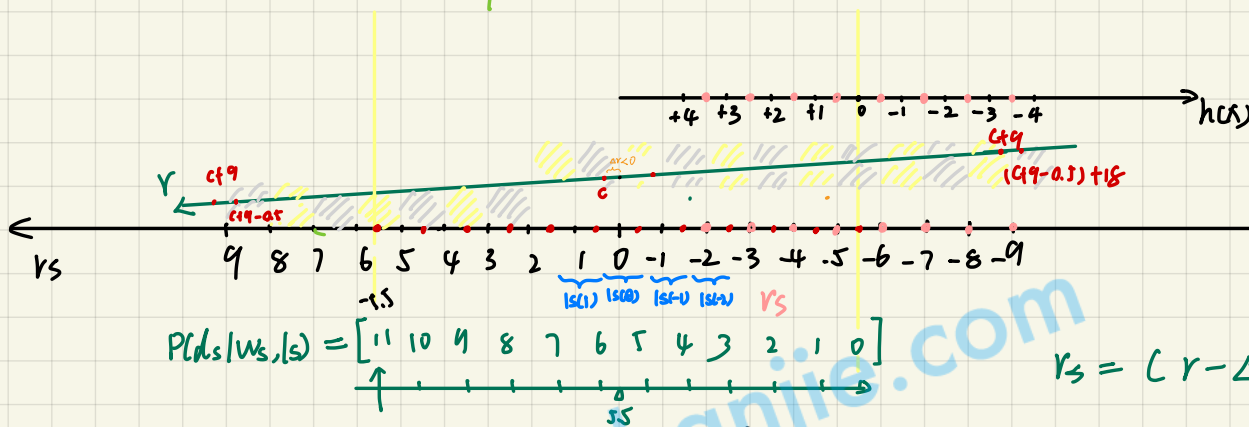
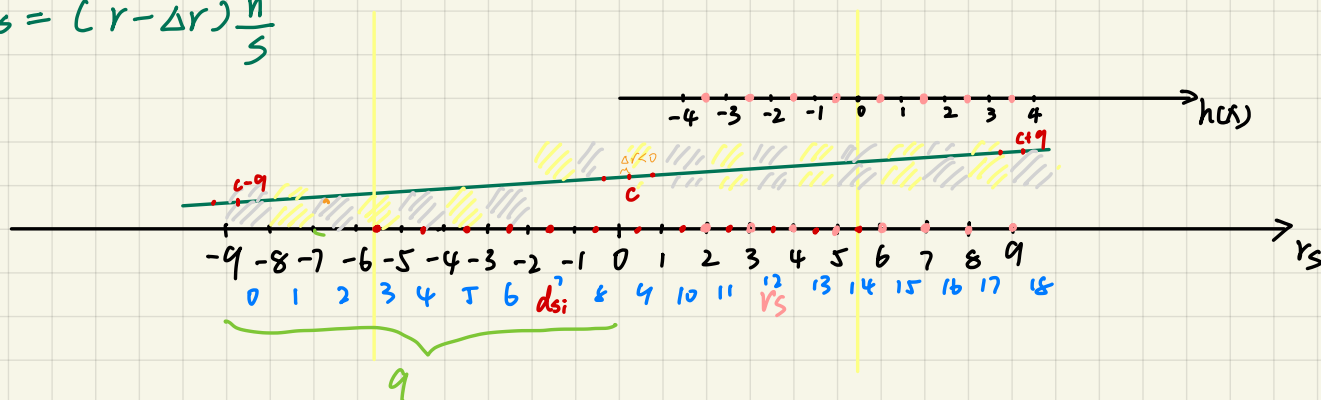
$$19 = \underbrace{X \in \{-3.5, 3.5\}}_8 + \underbrace{s_i \in \{-5.5, 5.5\}}_{12}$$

$$\text{line\_length\_pixel} = \text{line\_length\_segment} \times \text{scale}.$$

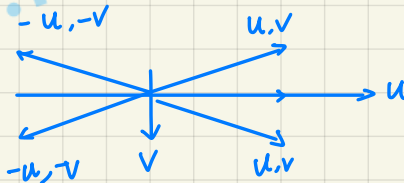
$$\Delta r = \Delta r_x / n_x$$



$$r_3 = (r - \Delta r) \frac{\bar{n}}{5}$$



$$r_3 = (r - \Delta r) \frac{\bar{n}}{5}$$

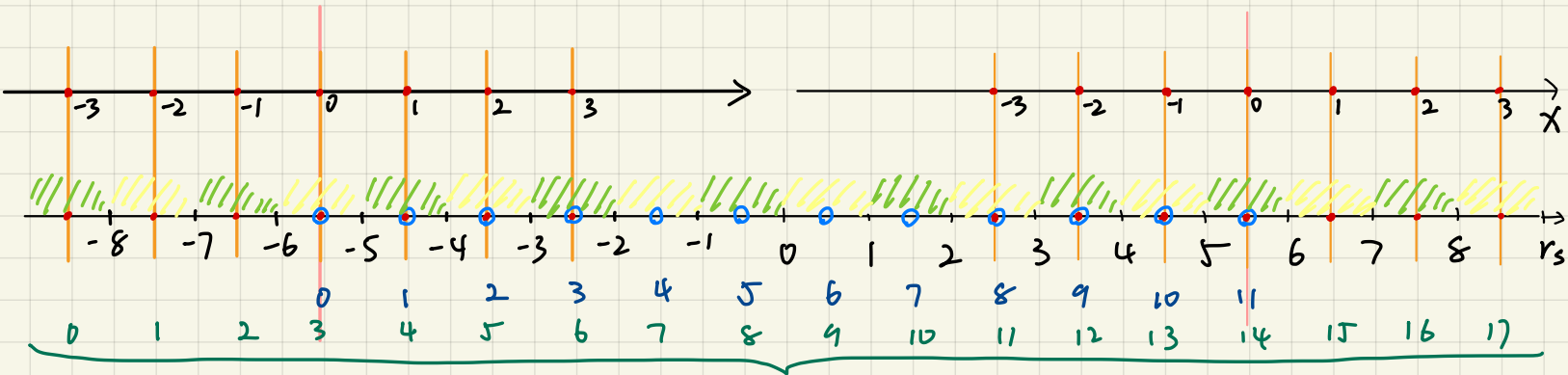


My Version

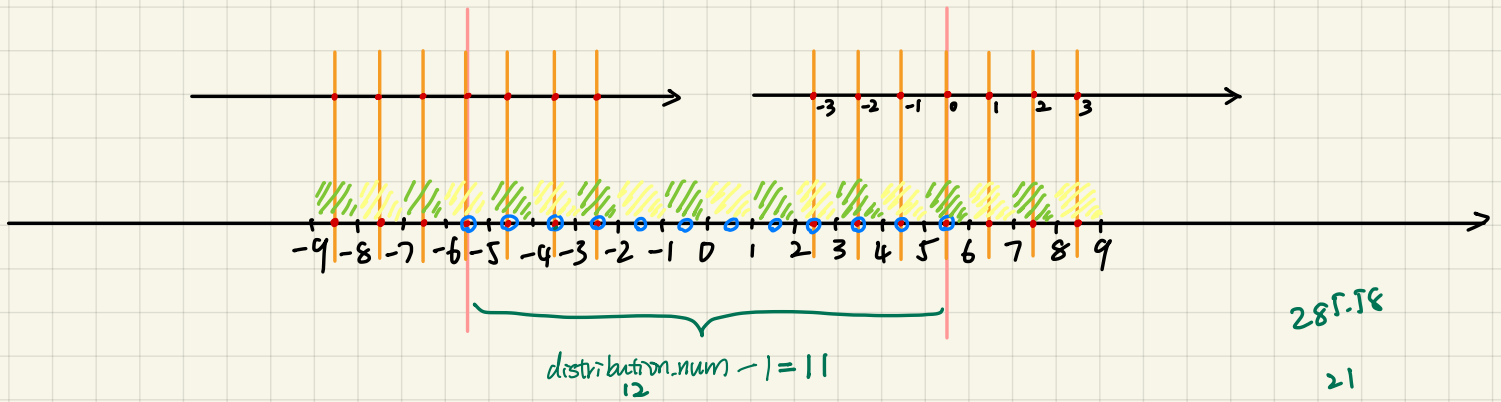
$$d_s = 5.5$$

$$r_s = \{2.5, 3.5, \dots, 7.5, 8.5\}$$

$$x = r_s - d_s = \{-3, -2, \dots, 3\}$$

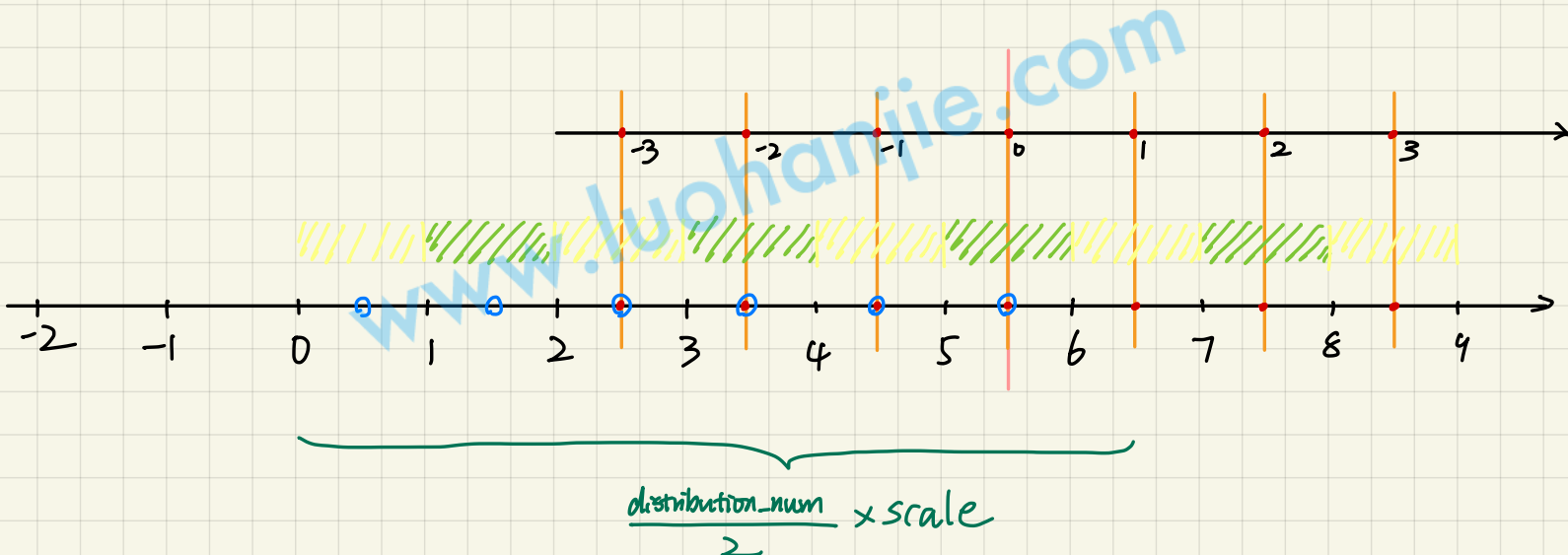


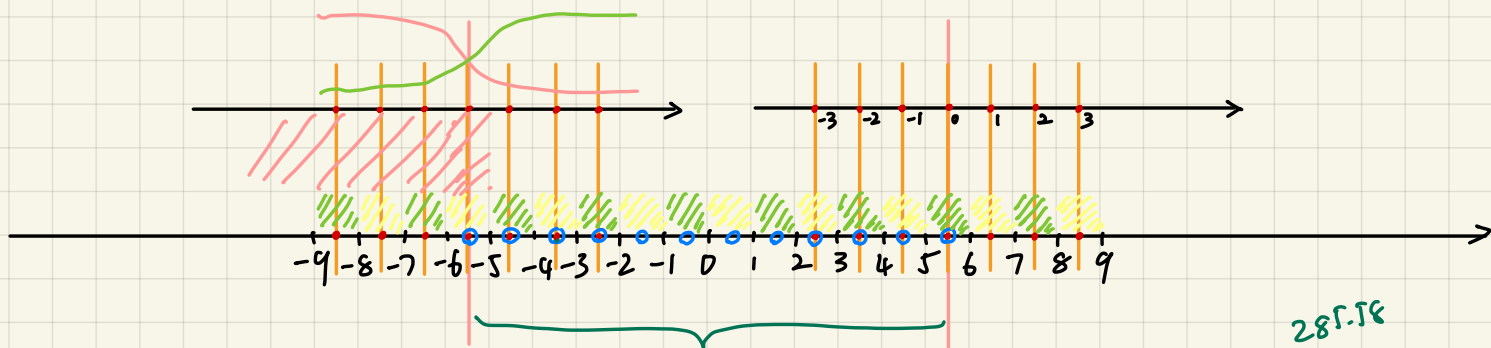
$$\begin{aligned} \text{line\_length\_pixel} &= (\text{line\_length\_segment} + 1) \times \text{scale} \\ &= 2 \times 18 = 36 \end{aligned}$$



$$\text{distribution\_num}_{12} - 1 = 11$$

$$\text{line\_length\_segment} = 17 = \text{function\_num}_7 + \text{distribution\_num}_{12} - 2$$





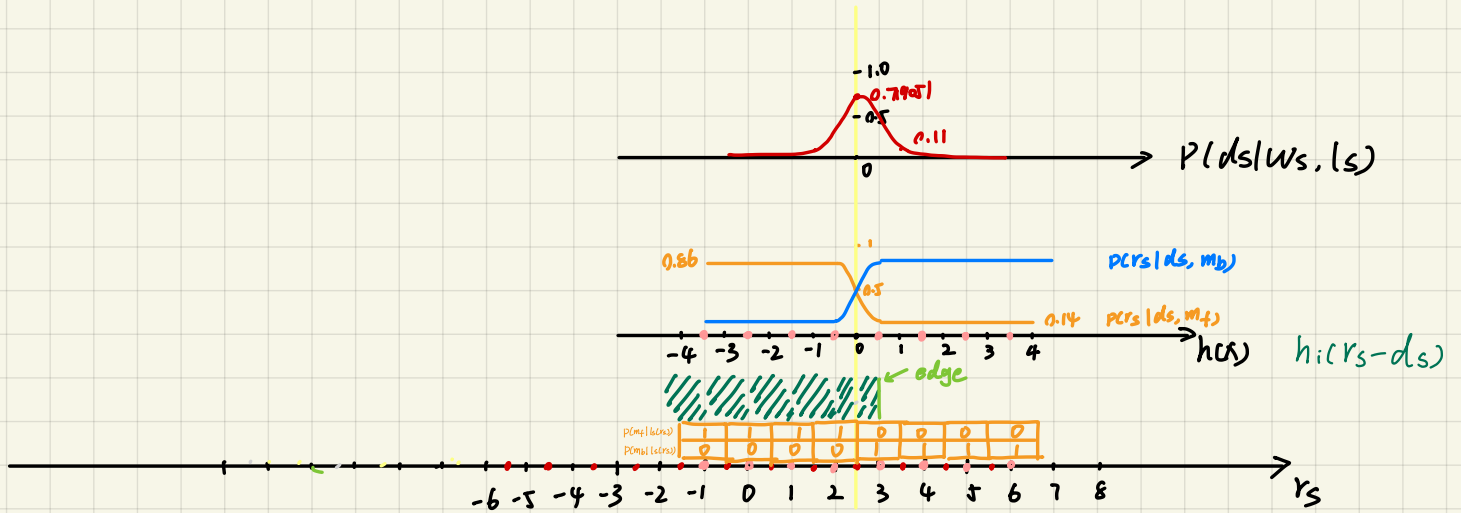
$$\text{distribution\_num} - 1 = 11$$

285.58

21

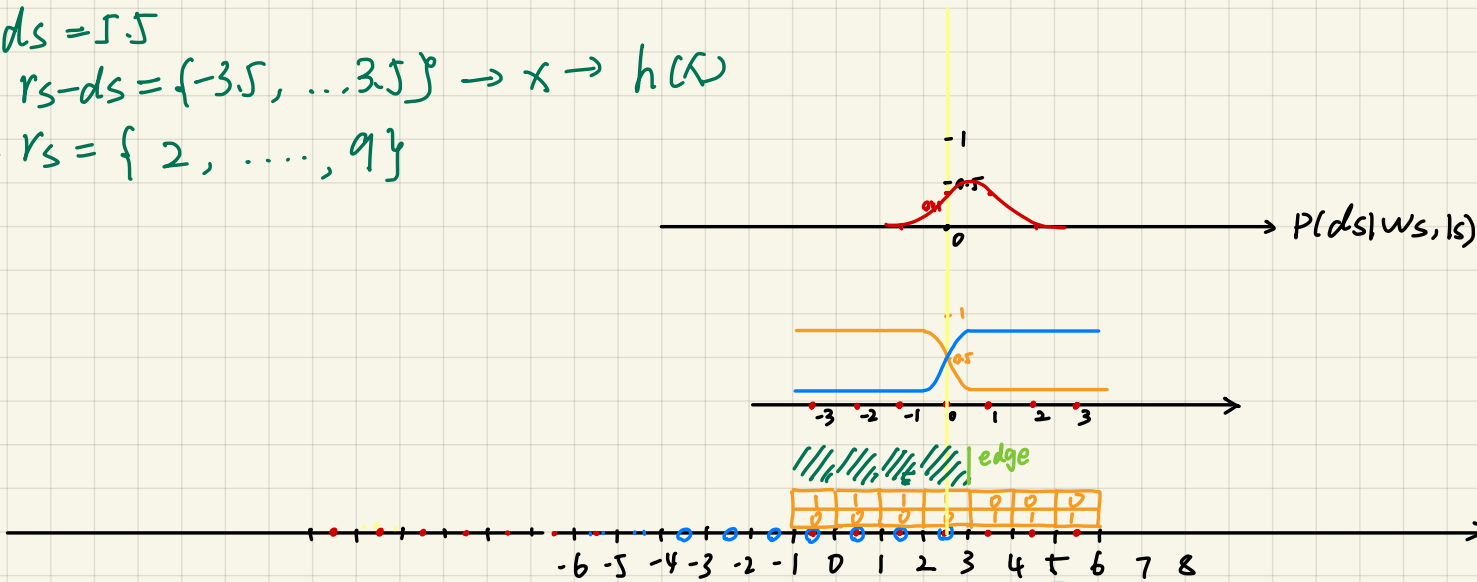
$$\text{line\_length\_segment} = 17 = \text{function\_num} + \text{distribution\_num} - 2$$

★ my version !



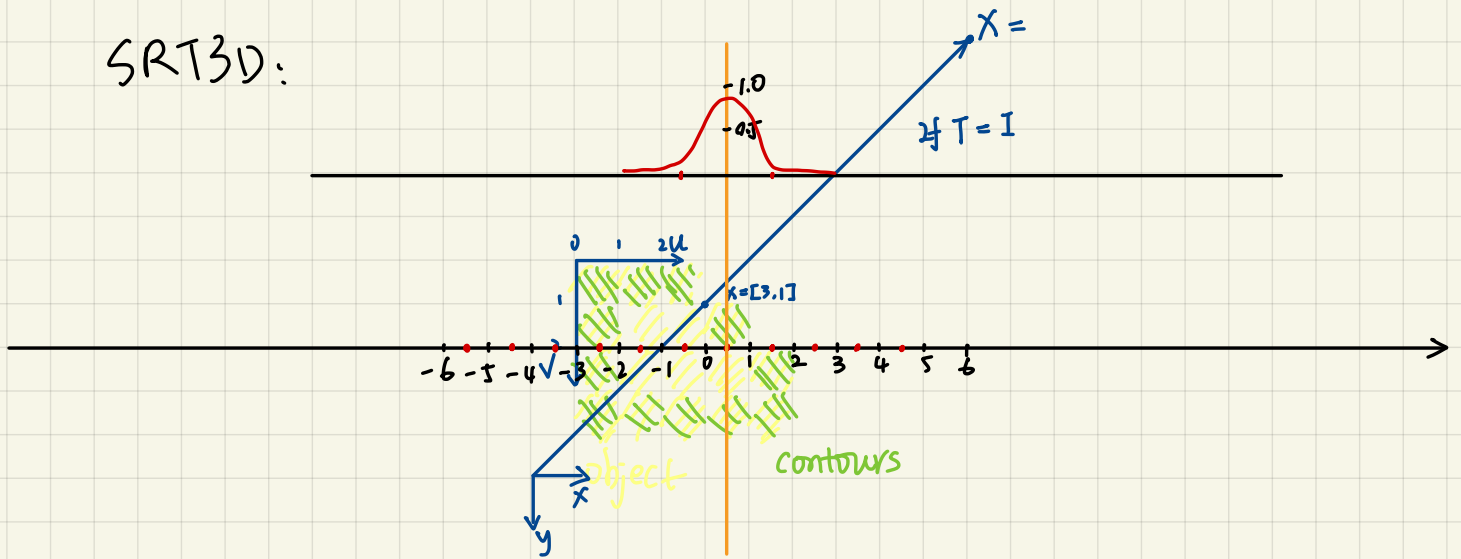
$$ds = \{-5.5, -4.5, -3.5, -2.5, -1.5, -0.5, 0.5, 1.5, 2.5, 3.5, 4.5, 5.5\}$$

If  $ds = 5.5$   
 $\therefore r_s - ds = \{-3.5, \dots, 3.5\} \rightarrow x \rightarrow h(x)$   
 $\therefore r_s = \{2, \dots, 9\}$



www.luohanjie.com

SRT3D:

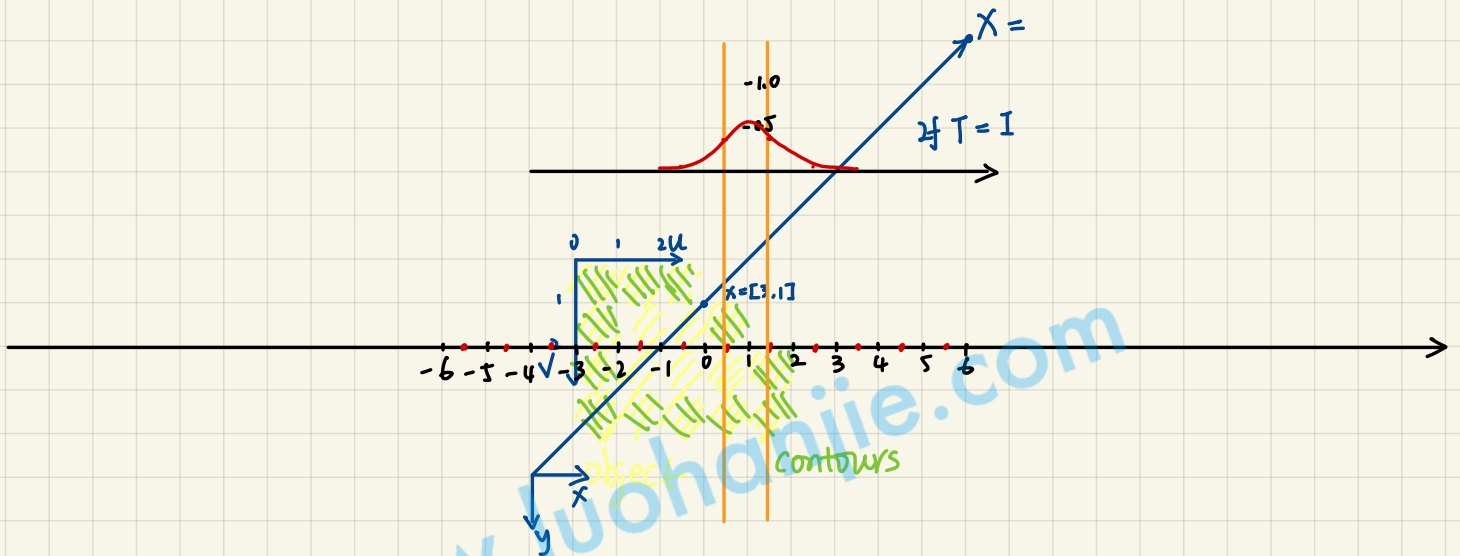


$$d(\theta) = n^T (\pi(c, X(\theta)) - c)$$

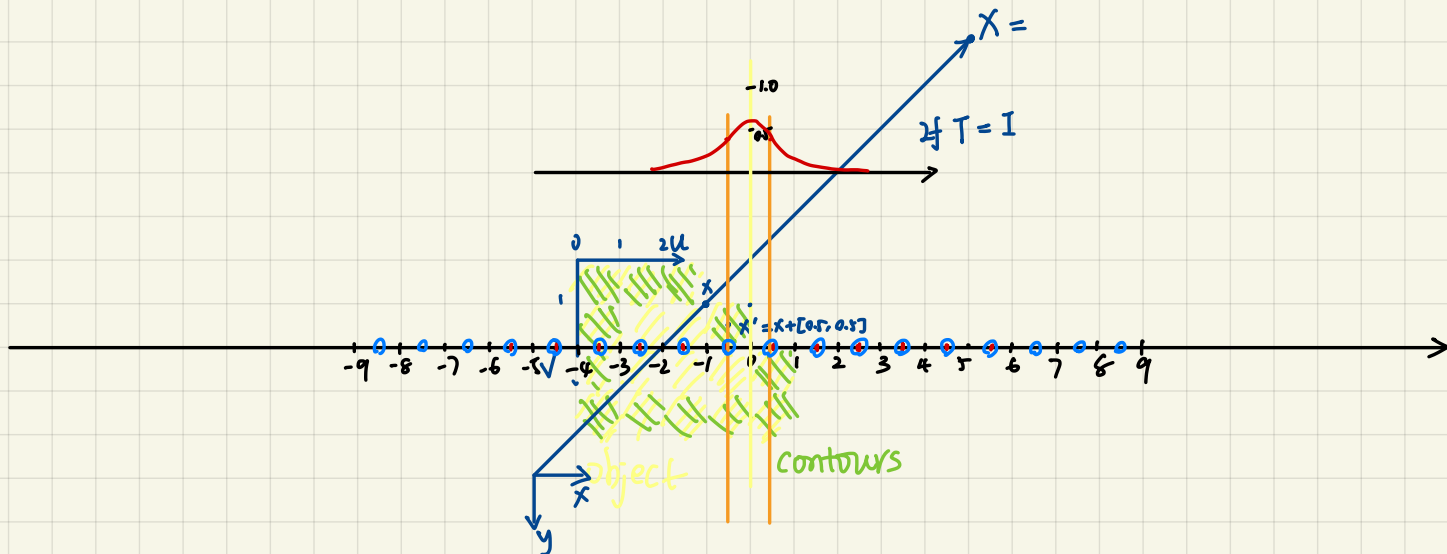
$$r_s = (r - \Delta r) \frac{\bar{n}}{s}$$

$$d_s = (d - \Delta r) \frac{\bar{n}}{s}$$

My:

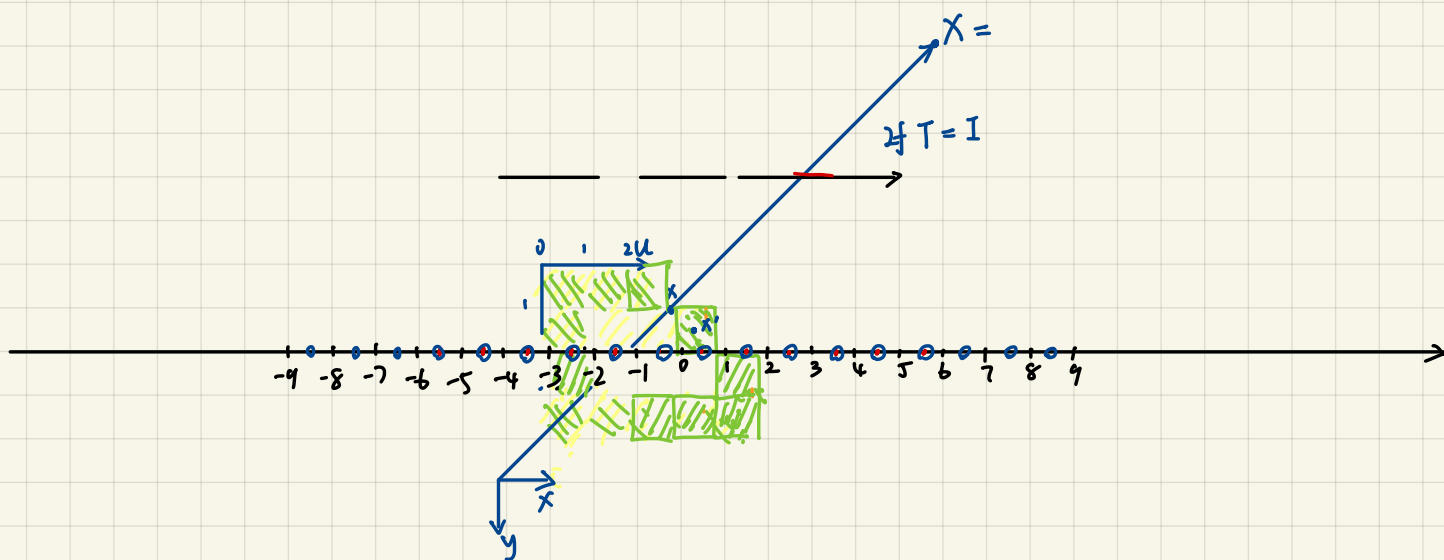






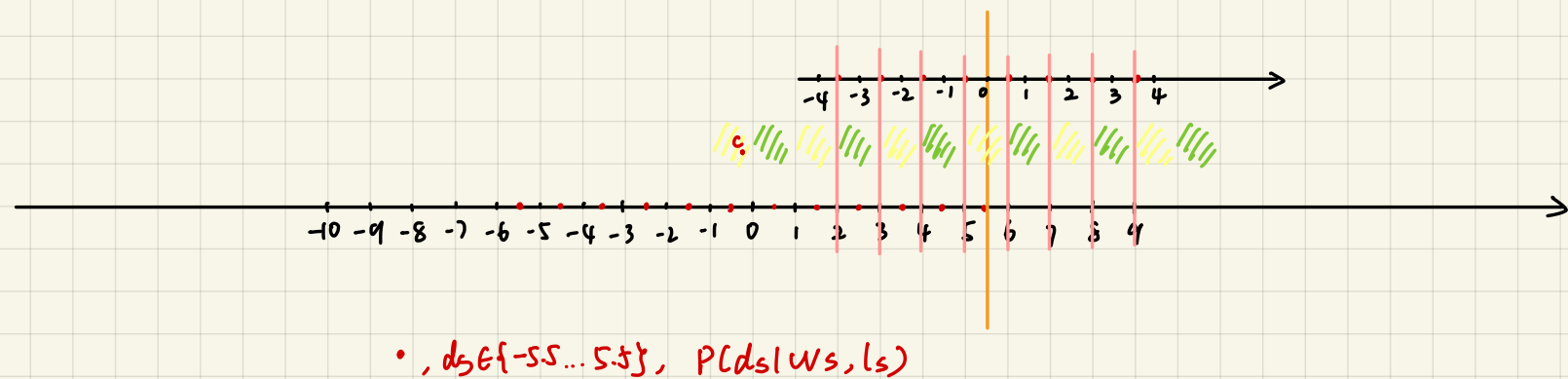
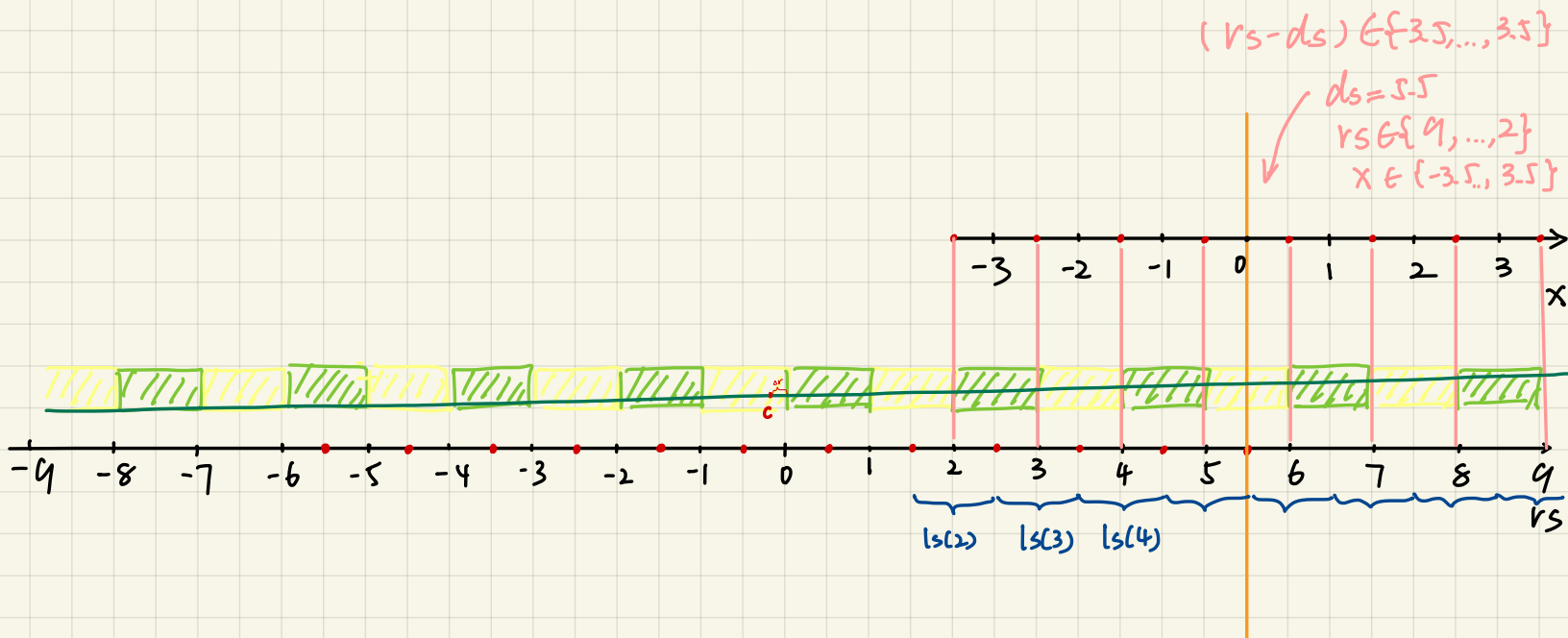
$$d(\theta) = n^T (\pi(c, X(\theta)) - c)$$

$$ds = (d - \Delta r) \frac{\bar{n}}{S}$$



www.luohanjie.com

Draft



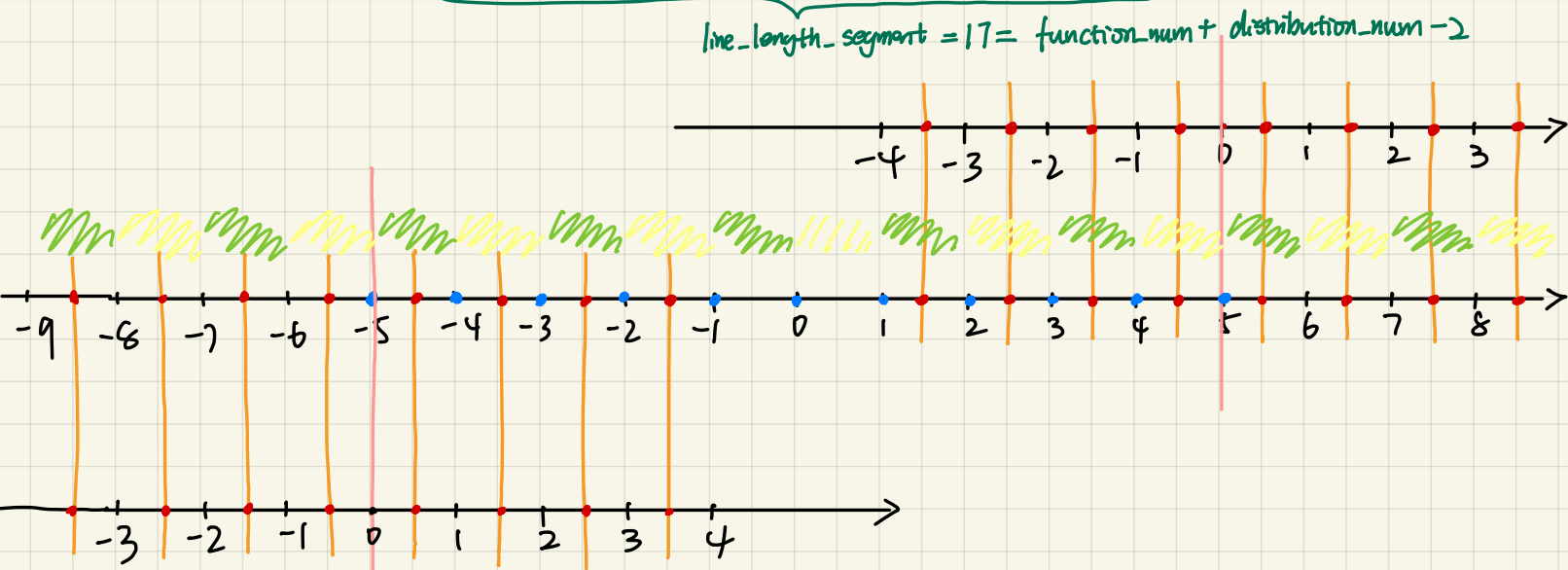
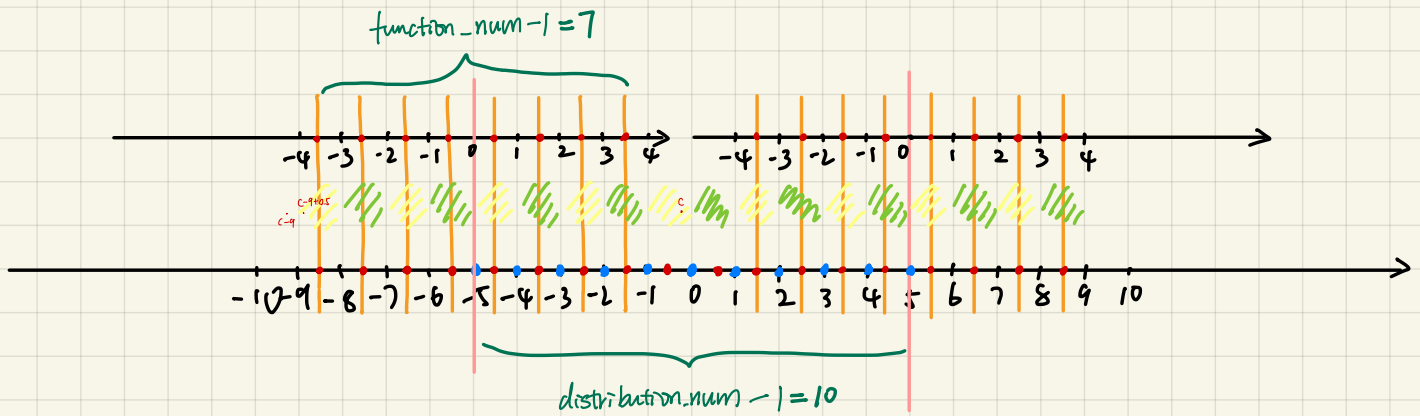
$$P(r_s | d_s, m_i) = h_i(r_s - d_s) = \begin{cases} \frac{1}{2} - a_n \tanh[(r_s - d_s)/2s_n] \\ \frac{1}{2} + a_n \tanh[(r_s - d_s)/2s_n] \end{cases}$$

$$P(d_s | w_s, l_s) \propto \prod_{r_s \in w_s} \sum_{i \in \{1, 2\}} P(r_s | d_s, m_i) p(m_i | l_s(r_s))$$

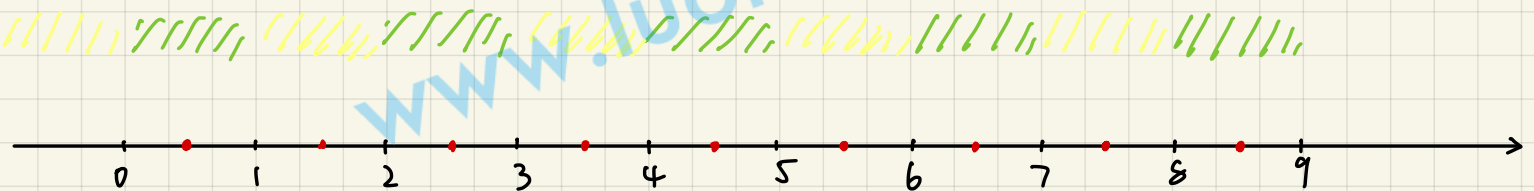
$\{ -5.5, \dots, 5.5 \}$

$$r_s = (r - \Delta r) \frac{\bar{n}}{5} \quad d_s = (\underline{d} - \Delta r) \frac{\bar{n}}{5}$$

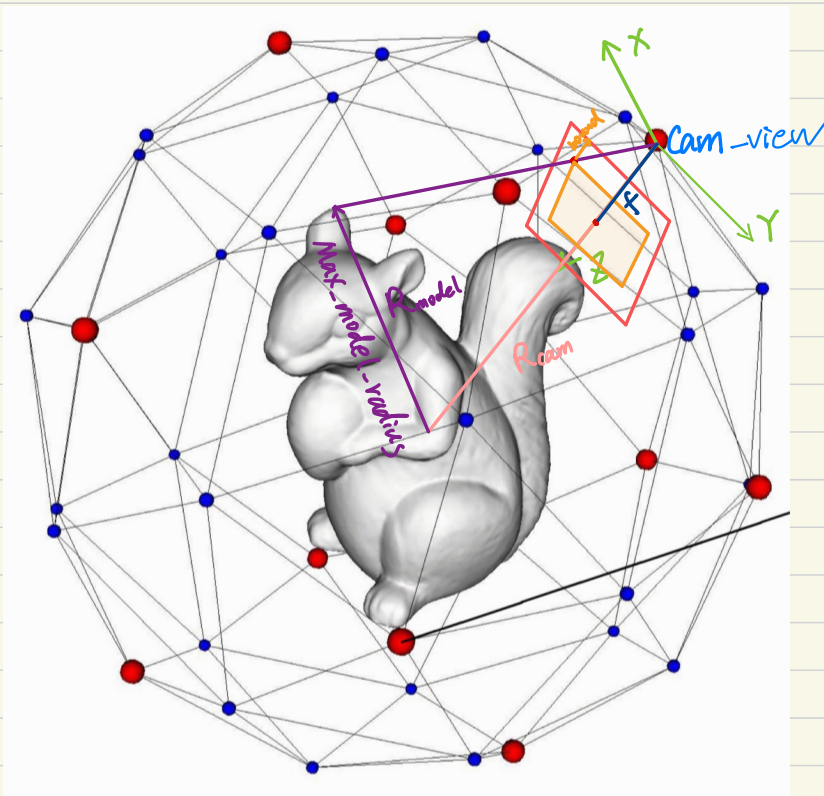
# Draft



$$\begin{aligned} \text{line\_length\_pixel\_} &= (\text{line\_length\_segpart} + 1) \times \text{scale} \\ &= 2 \times 18 = 36 \end{aligned}$$



www.lvohanjie.com



$$\therefore f_x \frac{X}{Z} + C_x = u$$

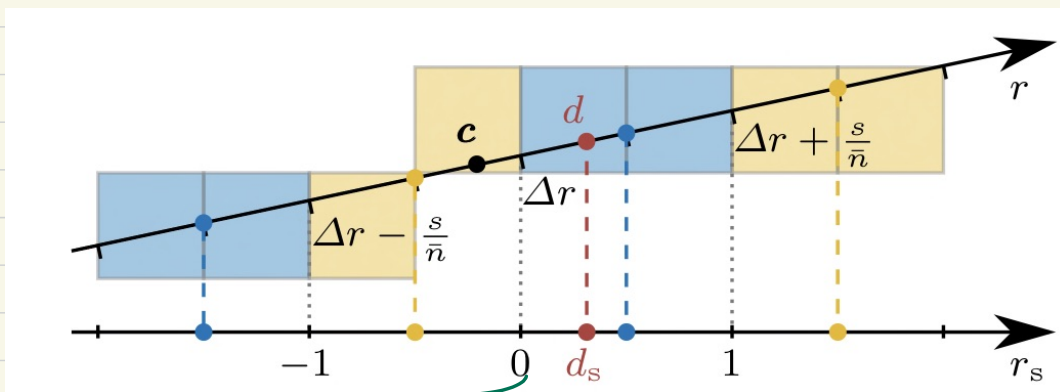
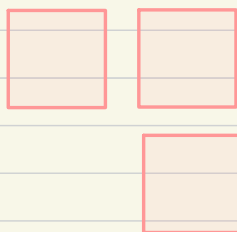
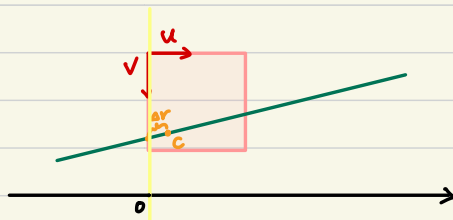
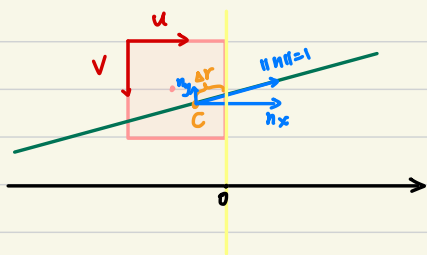
One Point in the max model sphere  
 $X = [R_{model}, 0, R_{cam}]$

$$f \frac{R_{model}}{R_{cam}} + C_x = \text{image-size} - \text{bound}$$

$$\text{make } C_x = \frac{\text{image-size}}{2}$$

$$\therefore f \frac{R_{model}}{R_{cam}} = \frac{\text{image-size} - \text{bound}}{2}$$

$$f = \left( \frac{\text{image-size} - \text{bound}}{2} \right) \frac{R_{cam}}{R_{model}}$$



$$\frac{(\text{line-length} - 1)}{2}$$

if  $\text{Scale} = 2$   
 $\text{line-length} =$

www.luohanjie.com

$$P(m_i | y_1, y_2, y_3) = \frac{P(y_1 | m_i)}{P(y_1 | m_f) + P(y_1 | m_b)} \cdot \frac{P(y_2 | m_i)}{P(y_2 | m_f) + P(y_2 | m_b)} \cdot \frac{P(y_3 | m_i)}{P(y_3 | m_f) + P(y_3 | m_b)}$$

$$= P(m_i | y_1) P(m_i | y_2) P(m_i | y_3).$$

$$\therefore P(m_b) = P(m_f)$$

$$P(m_i | y_1, y_2, y_3) = \frac{P(y_1, y_2, y_3 | m_i) P(m_i)}{P(y_1, y_2, y_3)} = \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i) P(m_i)}{\sum_{j \in \{f, b\}} P(y_1, y_2, y_3 | m_j) P(m_j)}$$

$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1, y_2, y_3 | m_f) + P(y_1, y_2, y_3 | m_b)}$$

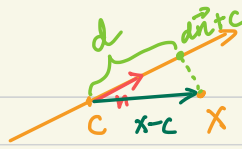
$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1, y_2, y_3 | m_f) + P(y_1, y_2, y_3 | m_b)}$$

$$= \frac{P(y_1 | m_i) P(y_2 | m_i) P(y_3 | m_i)}{P(y_1 | m_f) P(y_2 | m_f) P(y_3 | m_f) + P(y_1 | m_b) P(y_2 | m_b) P(y_3 | m_b)}$$

$$P(m_f | s) = \frac{1}{1 + \prod_{y \in s} \frac{P(y | m_b)}{P(y | m_f)}}$$



$$① \quad d = n^T (\overset{\text{an}}{\parallel} \vec{x} - c)$$



$$② \quad ds = (d - \Delta r) \frac{\bar{n}}{s}$$

we know  $ds$ , and want to know the pixel  $dn + c$

$$d = ds \frac{s}{\bar{n}} + \Delta r$$