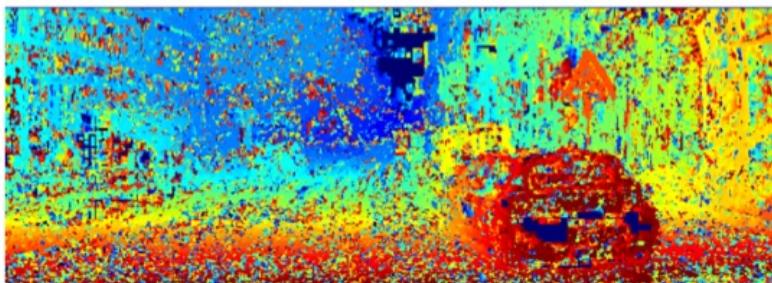


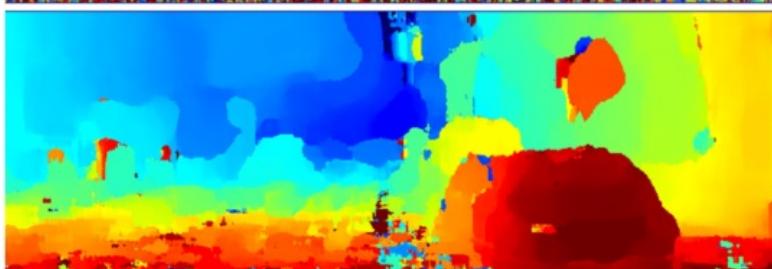
# LECTURE 10 : DEPTH FROM STEREO & VIDEO

## [1] Disparity Map from Rectified Images

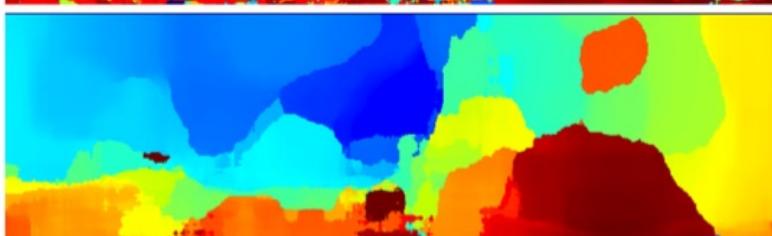
Pixel or patch based operations :



patch size = 5



patch size = 35



patch size = 85

Problem : small patches → more details but more noise  
large patch → less noise but less details

Solution: MRF



$$D^* = \underset{\{D\}}{\operatorname{argmin}} \sum_x \left( f_d(x, D(x), I, I') + \lambda \sum_{y \in N_x} f_p(D(x), D(y)) \right)$$

↓  
disparity map )      ↓  
                        pixel location      ↓  
                        input images  
                        ↓  
                        disparity to estimate  
 $\{D\} = \{d_{\min} \dots d_{\max}\}$   
the range of disparity values

Data term:

$$f_d(x, D(x), I, I') = [I(x) - I'(x + D(x))]^2$$

Photo-consistency constraint

If  $\{d_{\min} \dots d_{\max}\}$  has 10 values:  $d_1, d_2, \dots d_{10}$ ;  
then pixel  $x$  has 10 possible values of  $f_d$ , where each  
of them is the data cost of  $d_i$ .

Prior term:

$$f_p(D(x), D(y)) = \frac{(D(x) - D(y))^2}{\text{smoothness constraint}}$$

If we have 10 values of  $\{d_{\min} \dots d_{\max}\}$  for each pixel, we  
will have 100 values of  $f_p$ :

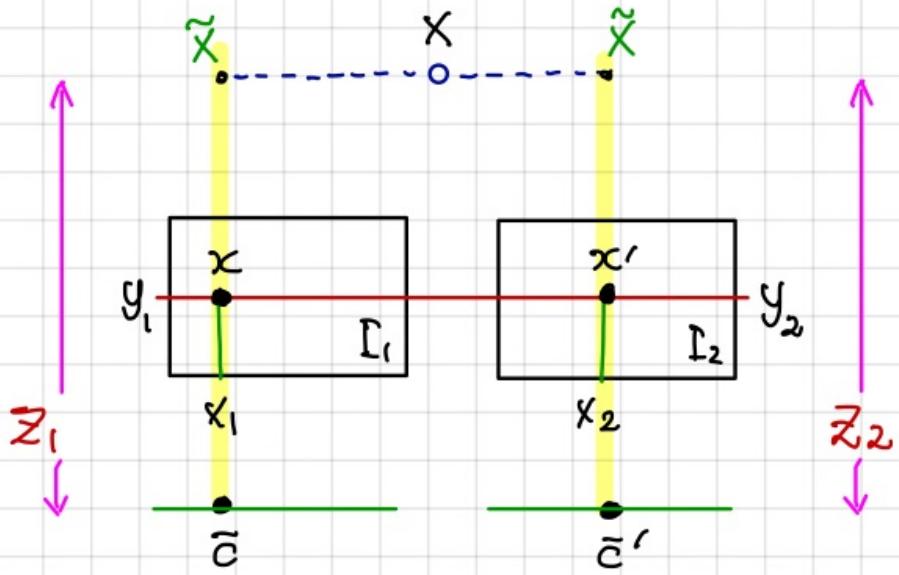
		$D(x)$	$d_1$	$d_2$	$\dots$	$d_{10}$
		$d_1$	$(d_1-d_1)^2$	$(d_2-d_1)^2$	$\vdots$	$(d_{10}-d_1)^2$
$D(y)$	$d_2$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	$d_{10}$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

These values are applied to all pixels, exactly in the same way.

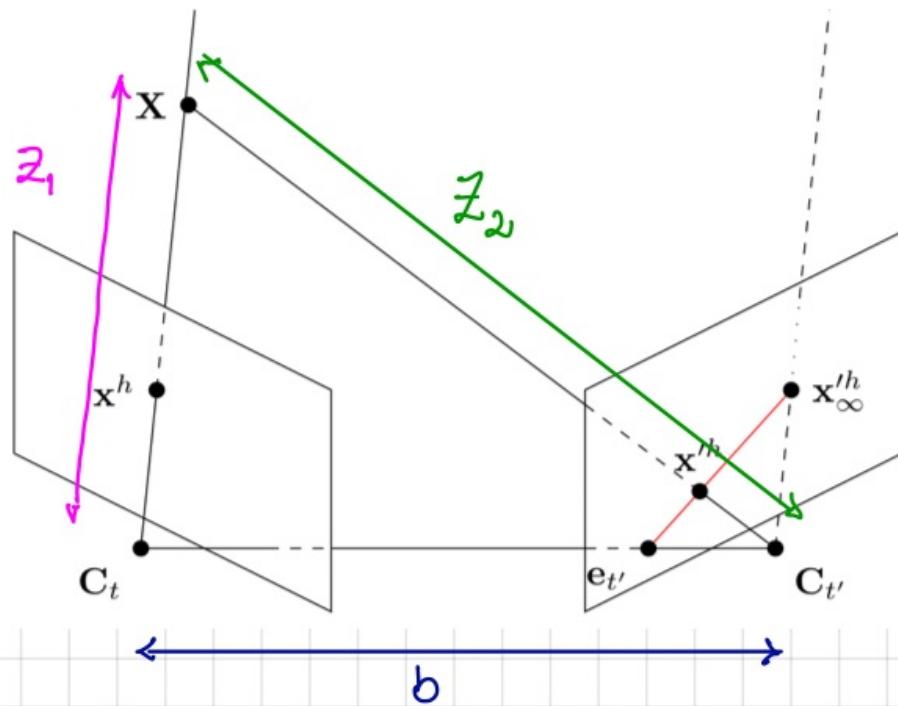


Having set the data and prior term for every pixel of the left  
image, we optimize the graph using graphcuts.

For a pair of rectified images, there is only one disparity or depth map:



Though, before being rectified, there are 2 depth maps:



Q: How to compute disparity of non-rectified images? #6

A: Two ways :

- (1) Rectify the images  $\rightarrow$  Epipolar rectification
- (2) Using epipolar geometry



Epipolar-geometry implicit rectification is preferable,  
since less steps are required in the computation.



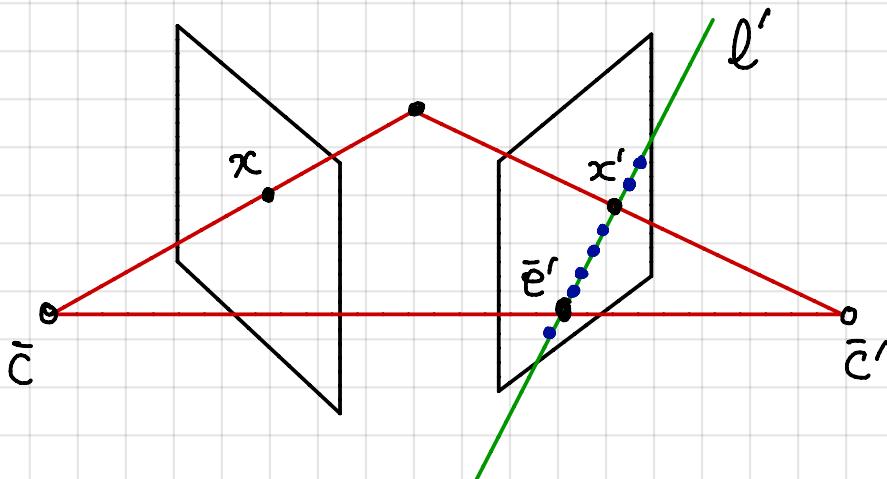
The basic formulation:

$$d = \frac{1}{z} \quad ; \quad \begin{array}{l} d = \text{disparity} \\ z = \text{depth} \end{array}$$

$$x' \sim K' R' R^T K^{-1} x + d K' R' (\bar{c} - \bar{c}')$$

Given  $x$ ,  $P = K[R|E]$  and  $P' = K'[R'|E']$ , we want to  
find  $d$  using the formula, such that  $I(x) = I'(x')$ .

The formula can generate all points on the green line below,  
when we change the value of  $d$ :

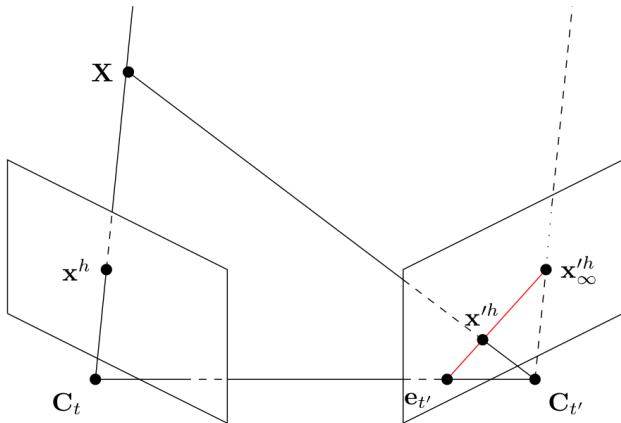


## [4] Derivation of Epipolar Rectification

To understand how we can get this equation:

$$x' \sim K' R' R^T K^{-1} x + d K' R' (\bar{c} - \bar{c}'),$$

consider:



$$x^h = P X = K [R | \bar{e}] X$$

Assume the existence of  $X_\infty$  along the line that passes through  $x$ , where  $X_\infty$  is located at infinity:

$$X_\infty = \begin{bmatrix} x \\ y \\ z \\ 0 \end{bmatrix} = \begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix}$$

$$\underset{3 \times 1}{x^h} = \underset{3 \times 3}{K} \underset{3 \times 3}{[R | \bar{e}]} \underset{3 \times 1}{X_\infty} = \underset{3 \times 3}{K} \underset{3 \times 3}{[R | \bar{e}]} \underset{\cancel{3 \times 4}}{\begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix}} = \underset{3 \times 3}{K R \hat{x}_\infty} \underset{3 \times 1}{\cancel{4 \times 1}}$$

Thus:

$$\hat{x}_\infty = R^T K^{-1} x^h$$

Projecting  $X_\infty$  onto the right image:

$$\begin{aligned} x'^h &= P' X_\infty = K' [R' | \bar{e}'] \begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix} = K' R' \hat{x}_\infty \\ &= K' R' R^T K^{-1} x^h \end{aligned}$$

What we want to find is  $x'^h$ , and not  $x_\infty^h$ .

Basic idea of finding  $x^h$ :

To find  $x^h$ , we can start the search from  $x_\infty^h$ . Then from  $x_\infty^h$ , we can trace along the epipolar line in the direction of the epipolar point,  $e_{t'}$ :

$$\begin{aligned}\bar{e}_{t'} &= P' \bar{C}_t = K' [R' | \bar{E}'] \bar{C}_t \\ &= K' R' \left[ \underbrace{I}_{3 \times 3} \underbrace{[(R')^{-1} \bar{E}']}_{3 \times 1} \right] \bar{C}_t \\ &= K' R' [\bar{C}_t - \bar{C}_{t'}]\end{aligned}$$

Recall:

$$\bar{C}_{t'} = -(R')^{-1} \bar{E}'$$

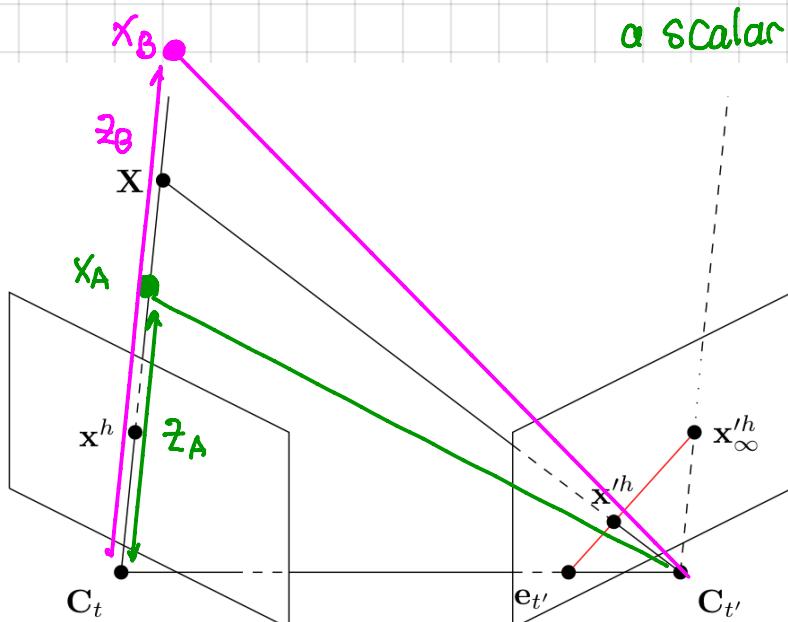
Therefore:

$$\underset{3 \times 1}{x^h} \sim \underset{3 \times 1}{x_\infty^h} + d \underset{3 \times 1}{\bar{e}_{t'}} \rightarrow \text{the line parametric equation}$$

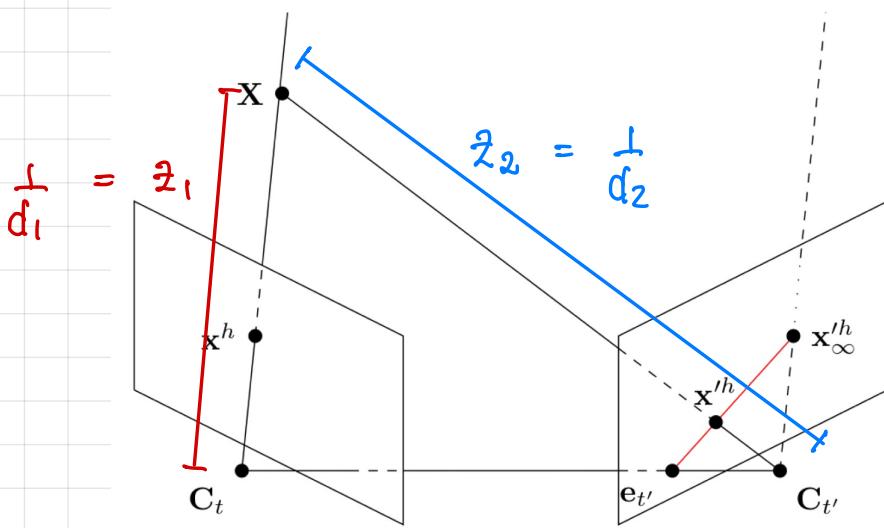
$$x^h \sim K' R' R^T K^{-1} x^h + d \bar{e}_{t'}$$

$$\underset{3 \times 1}{x^h} \sim \underset{3 \times 3}{K'} \underset{3 \times 3}{R'} \underset{3 \times 3}{R^T} \underset{3 \times 3}{K^{-1}} \underset{3 \times 1}{x^h} + d \underset{3 \times 3}{K'} \underset{3 \times 3}{R'} [\bar{C}_t - \bar{C}_{t'}]$$

a scalar value



The shorter the depth ( $z_A$ ), the larger the value of  $d$ , vice versa. This also indicates  $d$  determines the depth of point  $X$ .



$d$  = disparity  
 $z$  = depth

In the figure, with respect to  $X$ , there are 2 depths :

- $z_1 = 1/d_1$  is the depth from  $C_t$  to  $X$
- $z_2 = 1/d_2$  is the depth from  $C_{t'}$  to  $X$

$$x^h \sim x'^h + d, e_{t'} \quad \rightarrow \text{You need to normalize } x^h \text{ so that: } \begin{pmatrix} x \\ y \\ z \end{pmatrix} \rightarrow \begin{pmatrix} x/z \\ y/z \\ 1 \end{pmatrix}$$

Q: How to justify that this equation is correct ?

A: Consider the following cases :

If  $d=0$ , meaning  $z=\infty$  :

$$x'^h = x'^h_\infty$$

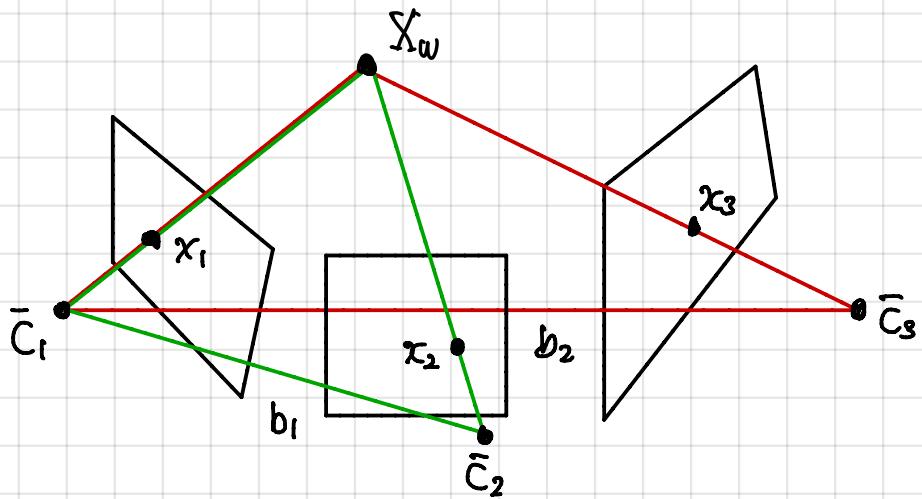
If  $d=\infty$ , meaning  $z=0$  (the 3D point is at  $\bar{C}_t$ ) :

$$\left. \begin{aligned} x'^h &= x'^h_\infty + \infty e' \\ &= e' \end{aligned} \right\} \text{see the illustration above}$$

If  $0 < d < \infty$ , then  $x'^h$  must lie on the epipolar line !

## [5] Epipolar Rectification for Multiple Images

#10



[0] For image 1 & 2 :

$$x_2 = x_{200} + d_{12} \quad \bar{e}_2 = K_2 R_2 R_1^T K_1^{-1} x_1 + d_{12} K_2 R_2 (\bar{c}_1 - \bar{c}_2)$$

[0] For image 1 & 3 :

$$x_3 = x_{300} + d_{13} \quad \bar{e}_3 = K_3 R_3 R_1^T K_1^{-1} x_1 + d_{13} K_3 R_3 (\bar{c}_1 - \bar{c}_3)$$



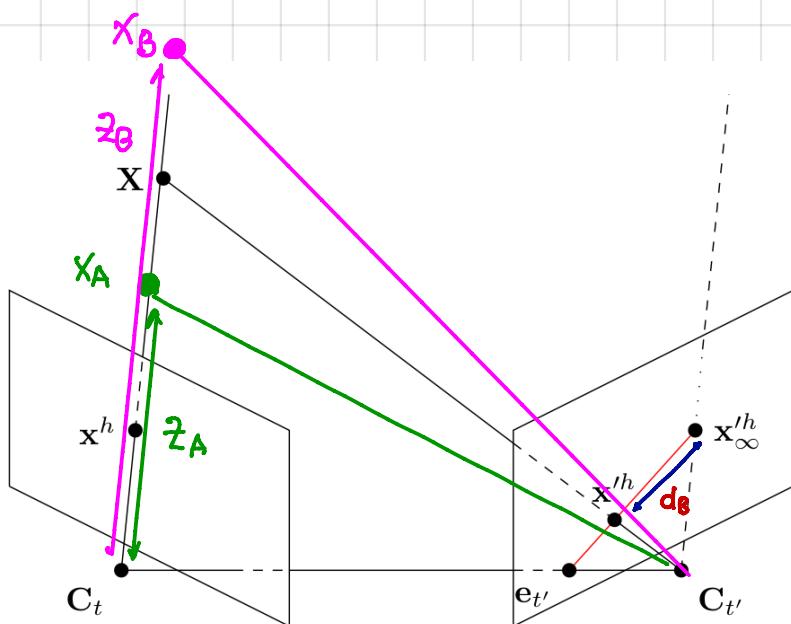
$d_{12} = d_{13}$

Implying :  $I_2(x_2) = I_1(x_1)$

$I_3(x_3) = I_1(x_1)$

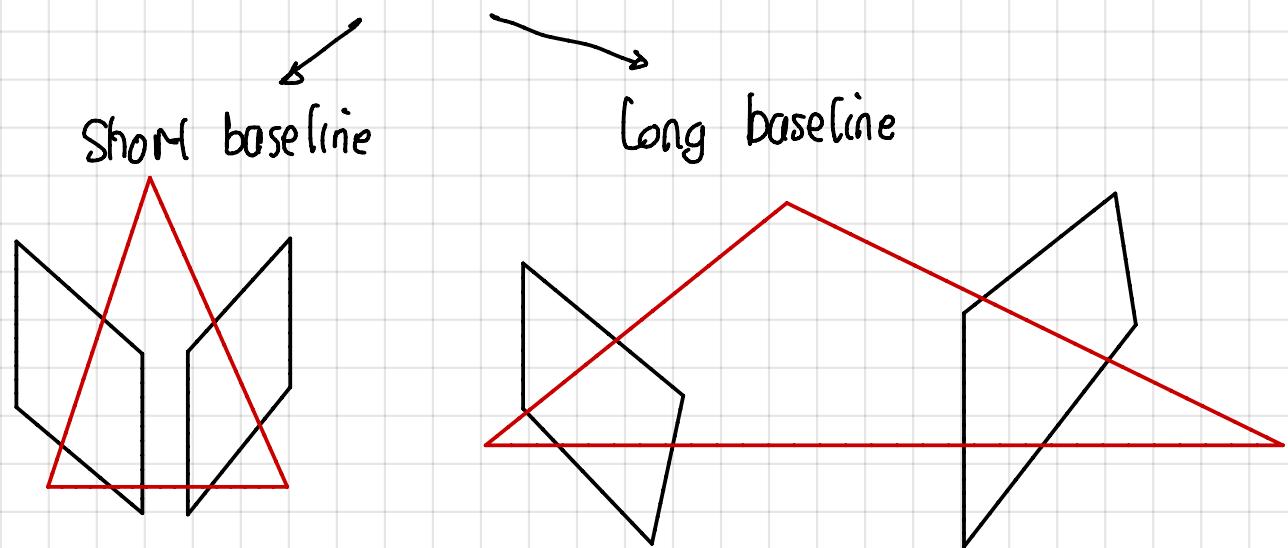
Q : Why  $d_{12} = d_{13}$  ?

A : Because  $d_{12}$  &  $d_{13}$  indicate the same depth of  $X$  from  $C_1$ .



## [6] Multiple Baseline Stereo = Depth from Video #1

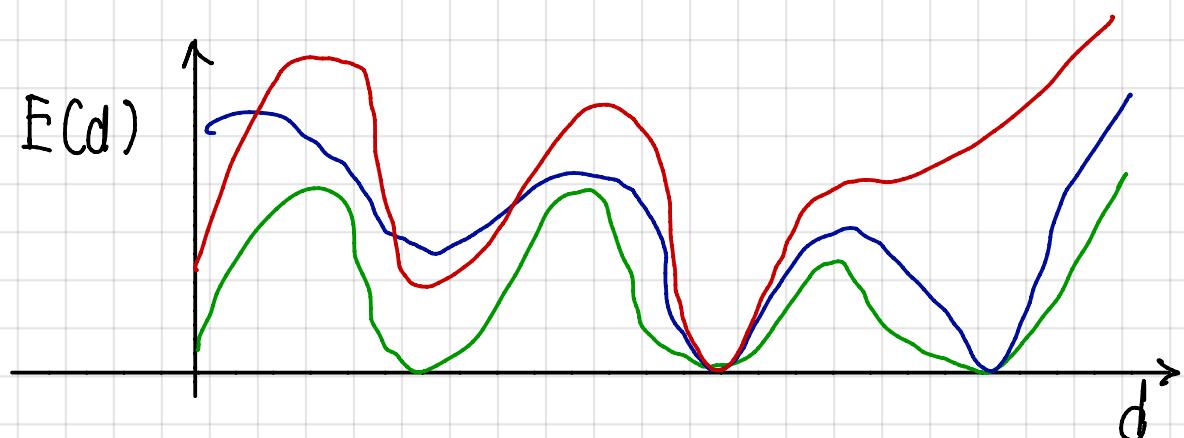
Disparity values depends on the baseline



For estimating disparity, a longer baseline is generally better, since  $b$  in the equation acts like a magnifier. However, longer baselines tend to suffer from occlusions.



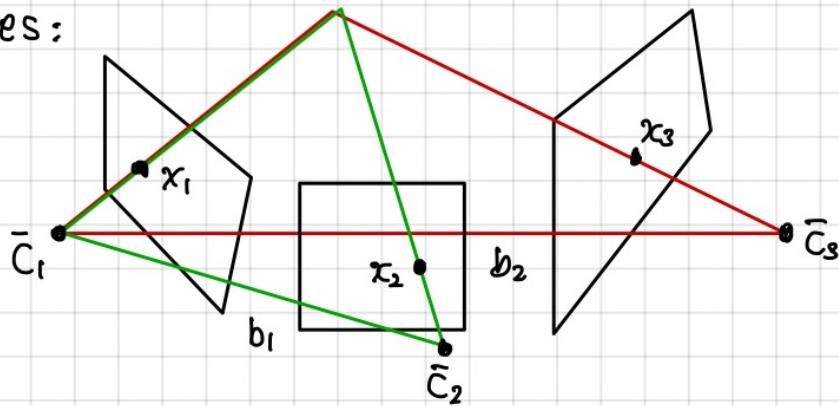
Multiple baselines allow us to get the benefits of short and long baselines:



The green line indicates that there are many candidates of  $d$  that make  $E(d)=0 \Rightarrow$  ambiguity!

# [7] Depth from Video (Multiple Baselines) #12

Multiple baselines:



Two main steps

(1) Disparity Initialization

Purpose: To have an initial depth map of each frame.

Input:  $\{I_t, P_t\}_{t=1}^N$ ;  $t = 1 \dots N$

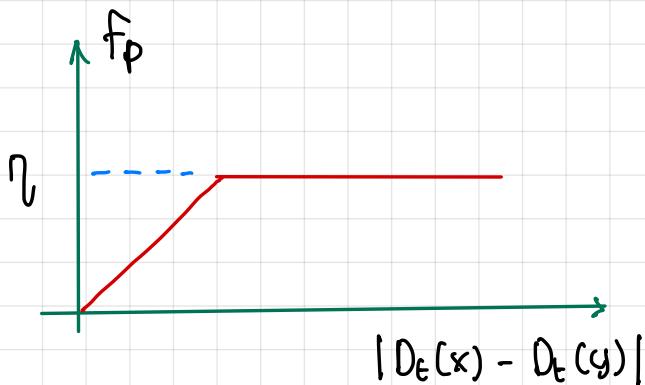
Output:  $\{D_t\}_{t=1}^N$  (the disparity maps)

$N$  = the number frames.

MAF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_x \left[ l - U(x) f_d^{\text{init}}(x, D_E(x)) - \sum_{y \in Nx} \lambda(x, y) f_p(D_t(x), D_t(y)) \right]$$

- Prior term:  $f_p(D_t(x), D_t(y)) = \min(|D_t(x) - D_t(y)|, \eta)$



To prevent  $f_p$  from being too large, since the difference between two disparity values shouldn't be that large.

Main reference: "Consistent depth map recovery from a video sequence."  
TPAMI, 2009.

MRF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_x \left[ l - U(x) f_d^{\text{init}}(x, D_t(x)) - \sum_{y \in N_x} x(y) f_p(D_t(x), D_t(y)) \right]$$

• Data term:

Use all other frames (prev. &amp; subsequent frames)

$$f_d^{\text{init}}(x, D_t(x)) = \sum_{t'}^N f_c(x, D_t(x), I_t, I_{t'})$$

e.g.:  $t = 1$   
 $t' = 2, 3, \dots, N$ 

where:

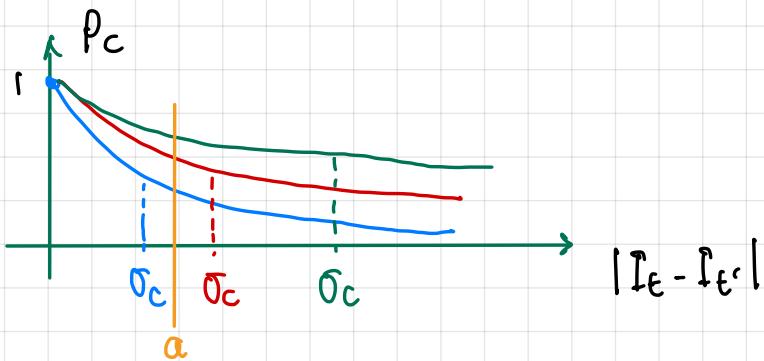
$$f_c(x, d, I_t, I_{t'}) = \frac{\sigma_c}{\sigma_c + |I_t(x) - I_{t'}(\ell_{t,t'}(x, d))|}$$

photo consistency

 $\ell_{t,t'}(x_t, d)$  means:

$$x_{t'} = K_{t'} R_{t'} R_t^T K_t^{-1} x_t + d K_{t'} R_{t'} [C_t - C_{t'}]$$

$3 \times 1 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 1 \quad 1 \times 1 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 1$

 $\sigma_c$  implies a variable to control the tolerance of the photo consistency:If  $|I_t - I_{t'}| = a$ , then different values of  $\sigma_c$  will produce different  $f_c$ .

The smaller  $\sigma_c$  (blue line) penalize  $|I_t - I_{t'}|$  more. Meaning, even though  $|I_t - I_{t'}|$  is relatively small,  $f_c$  will be smaller compared with the other values of  $\sigma_c$  (red and green lines).

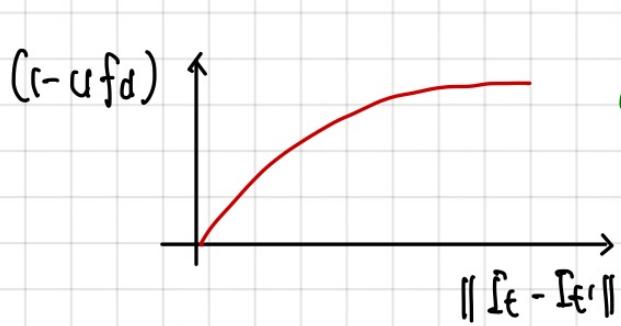
MRF :

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_x \left[ (1 - U(x)) f_d^{\text{init}}(x, D_t(x)) - \sum_{y \in N_x} \lambda(x,y) f_p(D_t(x), D_t(y)) \right]$$

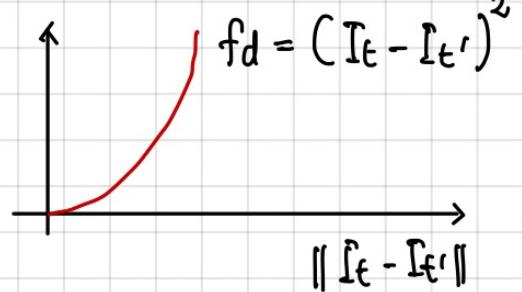
$$U(x) = \frac{1}{\max_{D_t(x)} f_d(x, D_t(x))}$$

→ the normalization factor.  
To ensure the values ranging from 0 to 1.

The graph of the overall data term:



more robust than  
our prev. cost  
function



If there is noise  $\|I_t - I_{t'}\|$  will be large, hence  $(I_t - I_{t'})^2$  will be also large, though the actual cost shouldn't be that large.

Using  $(1 - U f_d)$  will reduce the effect of noise in the optimization.

- Weighting factor  $\lambda$ :

MRF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_x \left[ I - U(x) f_d^{\text{init}}(x, D_E(x)) - \sum_{y \in N_x} \lambda(x, y) f_p(D_t(x), D_t(y)) \right]$$

Goal: To preserve the discontinuities.  $\lambda$  is defined to encourage the disparity discontinuities to be consistent with intensity / color discontinuities.

$$\lambda(x, y) = w_s \frac{U_\lambda(x)}{\|I_E(x) - I_E(y)\| + \epsilon}$$

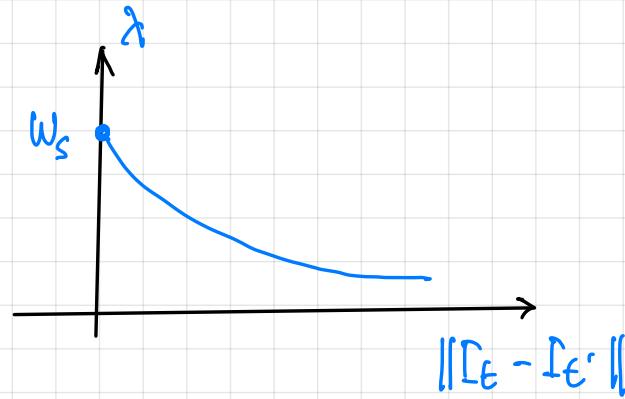
where:

$$U_\lambda(x) = \frac{|N_x|}{\sum_{y' \in N_x} \frac{1}{\|I_E(x) - I_E(y')\| + \epsilon}}$$

$w_s$  = the smoothness strength

$$\lambda(x, y) = w_s \frac{|N_x|}{(\|I_E(x) - I_E(y)\| + \epsilon) \sum_{y' \in N_x} \frac{1}{\|I_E(x) - I_E(y')\| + \epsilon}}$$

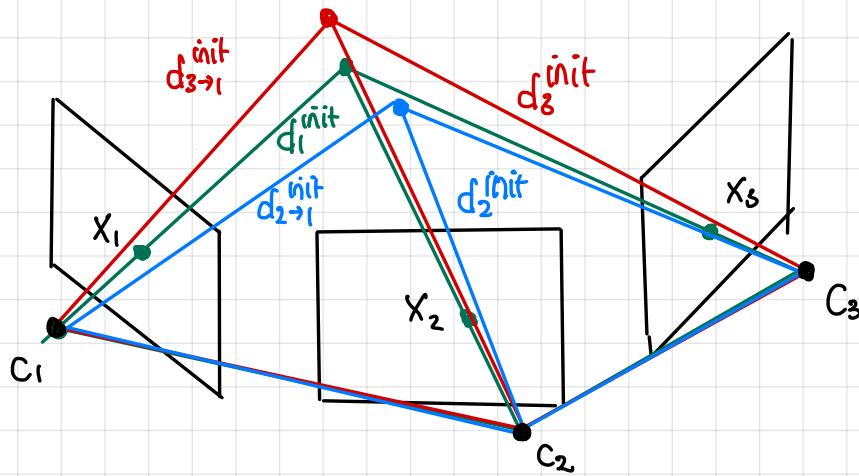
1. The smaller the intensity/color difference ( $\|I_E(x) - I_E(y)\|$ ) between two neighboring pixels, the larger  $\lambda$  is. A larger  $\lambda$  will encourage the smoothness constraint more.
2.  $U_\lambda$  is to normalize  $\|I_E(x) - I_E(y)\|$  so that  $\lambda$  should be within a certain range depending on  $w_s$ .



## [8] Bundle Optimization

#16

Problem :

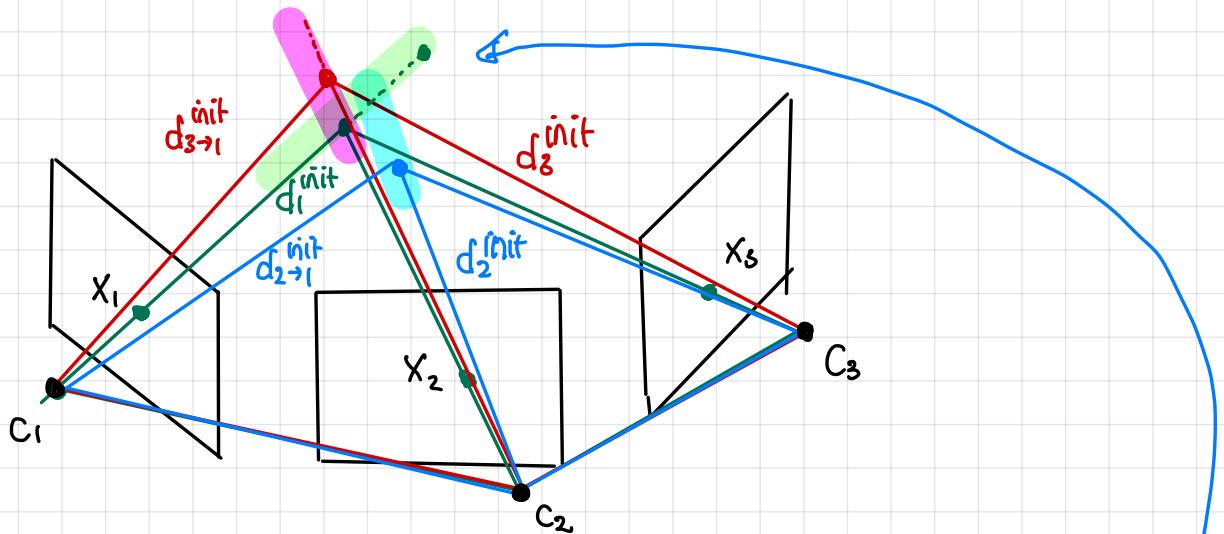


$$\text{ Ideally: } d_1^{\text{init}} = d_{2 \rightarrow 1}^{\text{init}} = d_{3 \rightarrow 1}^{\text{init}}$$

Reality :  $d_1^{\text{init}} \neq d_{2 \rightarrow 1}^{\text{init}} \neq d_{3 \rightarrow 1}^{\text{init}}$   $\rightarrow$  This is the cause of the flickering problem

To address the problem we need to make :

$d_1 \approx d_{2 \rightarrow 1} \approx d_{3 \rightarrow 1}$   $\rightarrow$  Enforce them to be as close as possible .



Basic idea :

To find other disparity value other than  $d_c^{\text{init}}$  that close the gap in  $d_1$  &  $d_{2 \rightarrow 1}$  &  $d_{3 \rightarrow 1}$

Candidates for  $d_t^{\text{init}}$  are :  $\{d_t^{\text{init}} - N\epsilon, \dots, d_t^{\text{init}}, d_t^{\text{init}} + \epsilon, \dots, d_t^{\text{init}} + N\epsilon\}$

$N$  &  $\epsilon$  are variables that need to be decided.

• Input:  $\{I_t, IP_t, D_t^{\text{init}}\}_{t=1}^N$

Output:  $D_t$

$$D_t^{\text{init}} = \underset{\{d_t^{\text{init}} + n\epsilon\}_{n=-N}^{+N}}{\operatorname{argmin}} \sum_x \left[ I - U(x) f_d^{\text{init}}(x, D_t(x)) - \sum_{y \in Nx} \lambda(x, y) f_p(D_t(x), D_t(y)) \right]$$

candidates:  $\{d_{t-NE}^{\text{init}}, \dots, d_{t-\epsilon}^{\text{init}}, d_t^{\text{init}}, d_{t+\epsilon}^{\text{init}}, \dots, d_{t+NE}^{\text{init}}\}$

where:

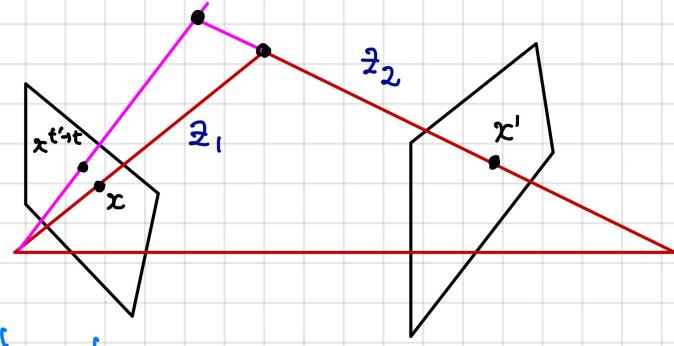
$$f_d(x, d) = \sum_{t'}^N f_c(x, d, I_t, I_{t'}) p_v(x, d, D_{t'})$$

photo consistency      geometric coherence

$\Downarrow$   
the same as in the previous step.

Geometric Coherence:

Goal: to measure how consistent  $z_1$  (or  $d_1$ ) and  $z_2$  (or  $d_2$ ) geometrically.



$$p_v(x, d, D_{t'}) = \exp \left( - \frac{\|x - l_{t', t}(x', D_{t'}(x))\|^2}{2\sigma_d^2} \right)$$

where:  $x' = l_{t', t}(x, d)$

$$x^{t' \rightarrow t} = l_{t', t}(x', D_{t'}(x))$$

$$\|x - x^{t' \rightarrow t}\| = \|x - l_{t', t}(x', D_{t'}(x))\|$$

similar to  $\sigma_c$ : to control how tolerant we are with the error



;  $\sigma_d$  = standard deviation