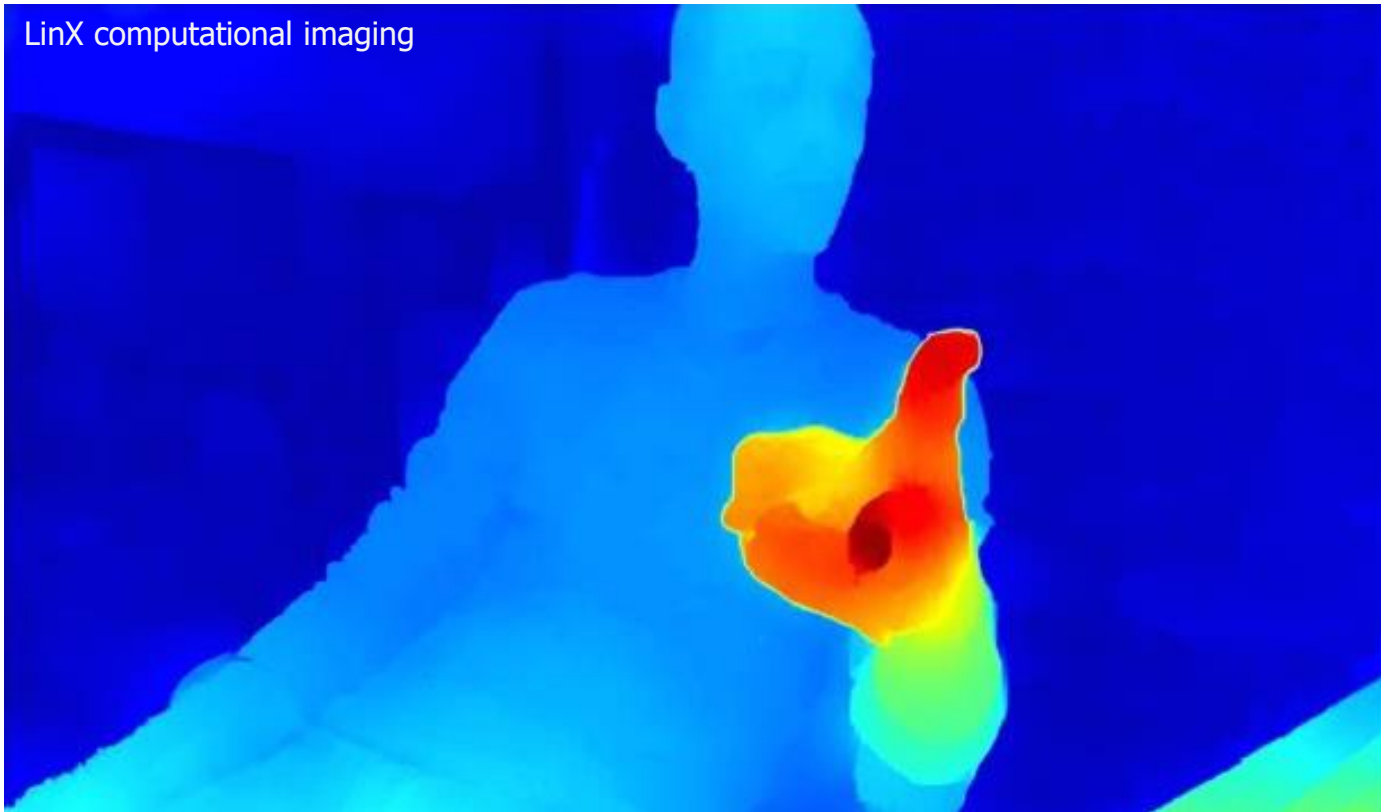


CS484/684 Computational Vision

Dense Stereo

LinX computational imaging



Dense Stereo

towards **dense** 3D reconstruction

assumption: known camera motion, i.e. known epipolar lines

- (dense) stereo is an example of **dense correspondence**
- another example is dense motion estimation (*optical flow*)

But, **stereo is simpler** since the search for correspondences is restricted to 1D epipolar lines (versus 2D search for non-rigid motion)

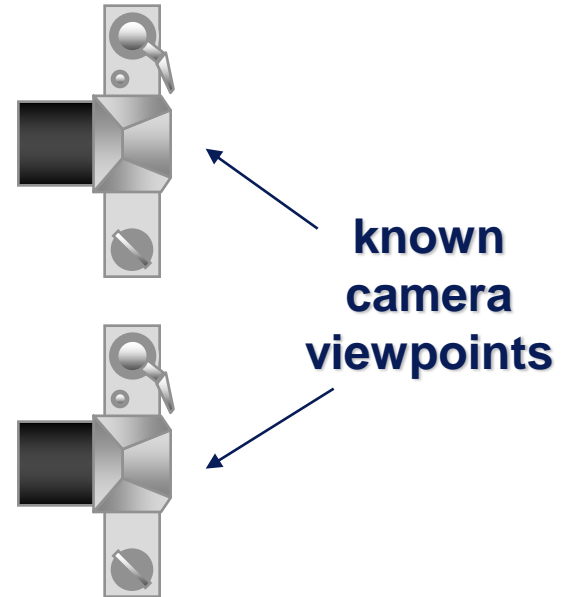
Dense Stereo Correspondence

- camera rectification for stereo pairs
- window-based (local) stereo
- scan-line stereo correspondence
 - optimization via DP, Viterbi, Dijkstra
- image-grid (global) stereo
 - optimization via graph cuts

examples of loss functions with spatial/geometric **regularization** of image labels

Szeliski, Chapter 11

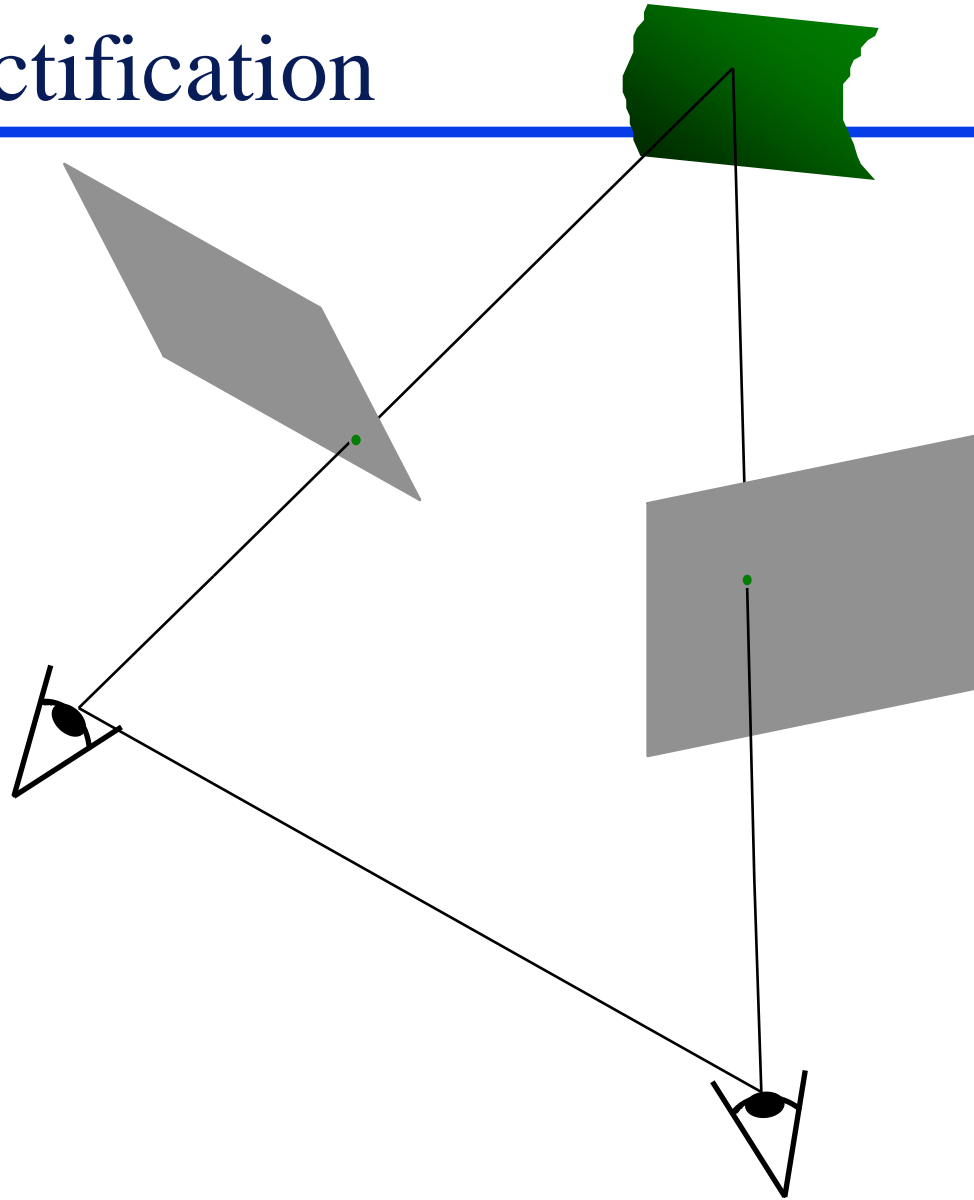
Stereo vision



Two views of the same scene from slightly different point of view
Also called, narrow baseline stereo.

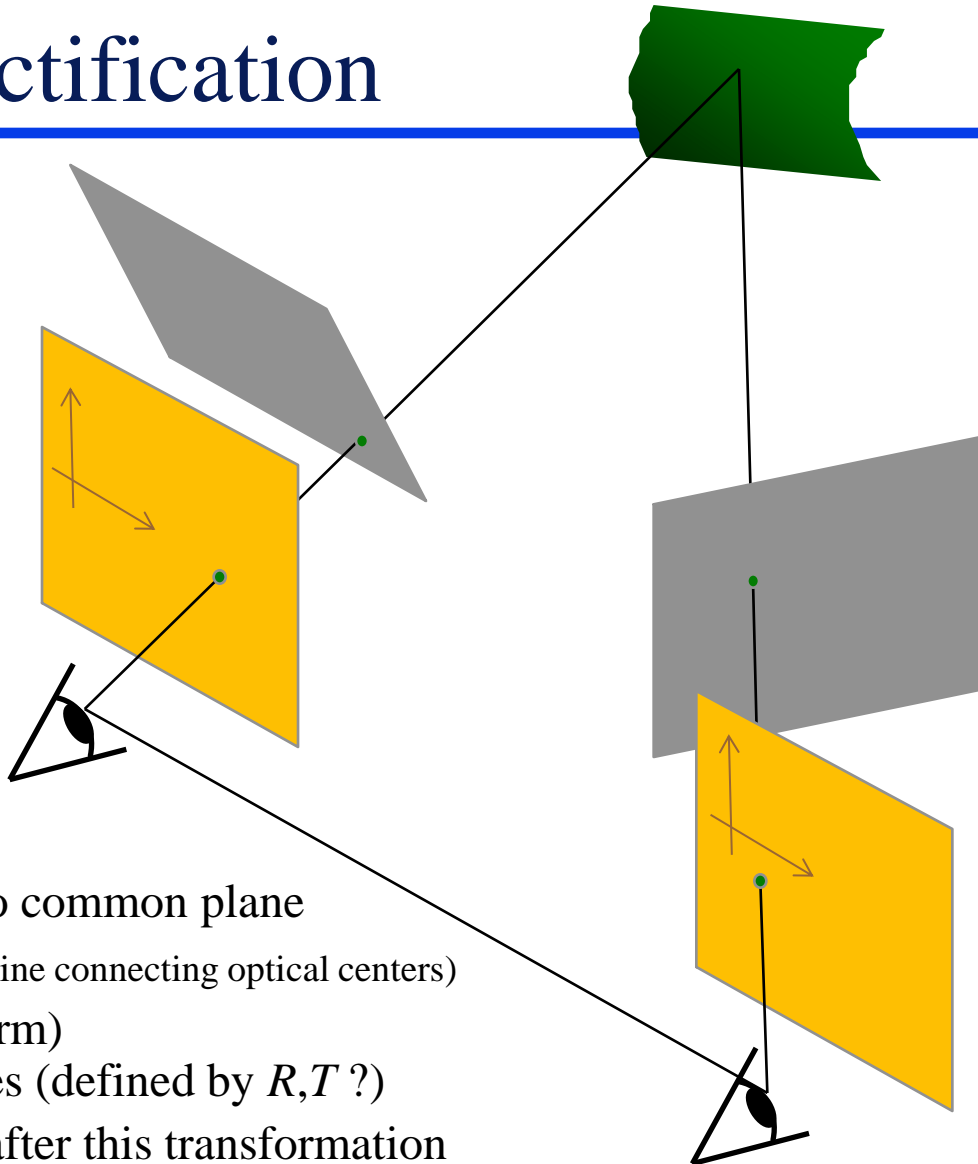
Motivation: - smaller difference in views allows to find **more matches** (Why?)
- scene reconstruction is simply represented via **depth map**

Stereo image rectification



Stereo image rectification

the corresponding (rectified)
camera geometry is analogous to
“panning motion”



□ Image Reprojection

- reproject image planes onto common plane parallel to the baseline (i.e. line connecting optical centers)
- homographies (3x3 transform) applied to both input images (defined by R, T ?)
- pixel motion is horizontal after this transformation
- C. Loop and Z. Zhang. [Computing Rectifying Homographies for Stereo Vision](#). IEEE Conf. Computer Vision and Pattern Recognition, 1999.

Stereo image rectification

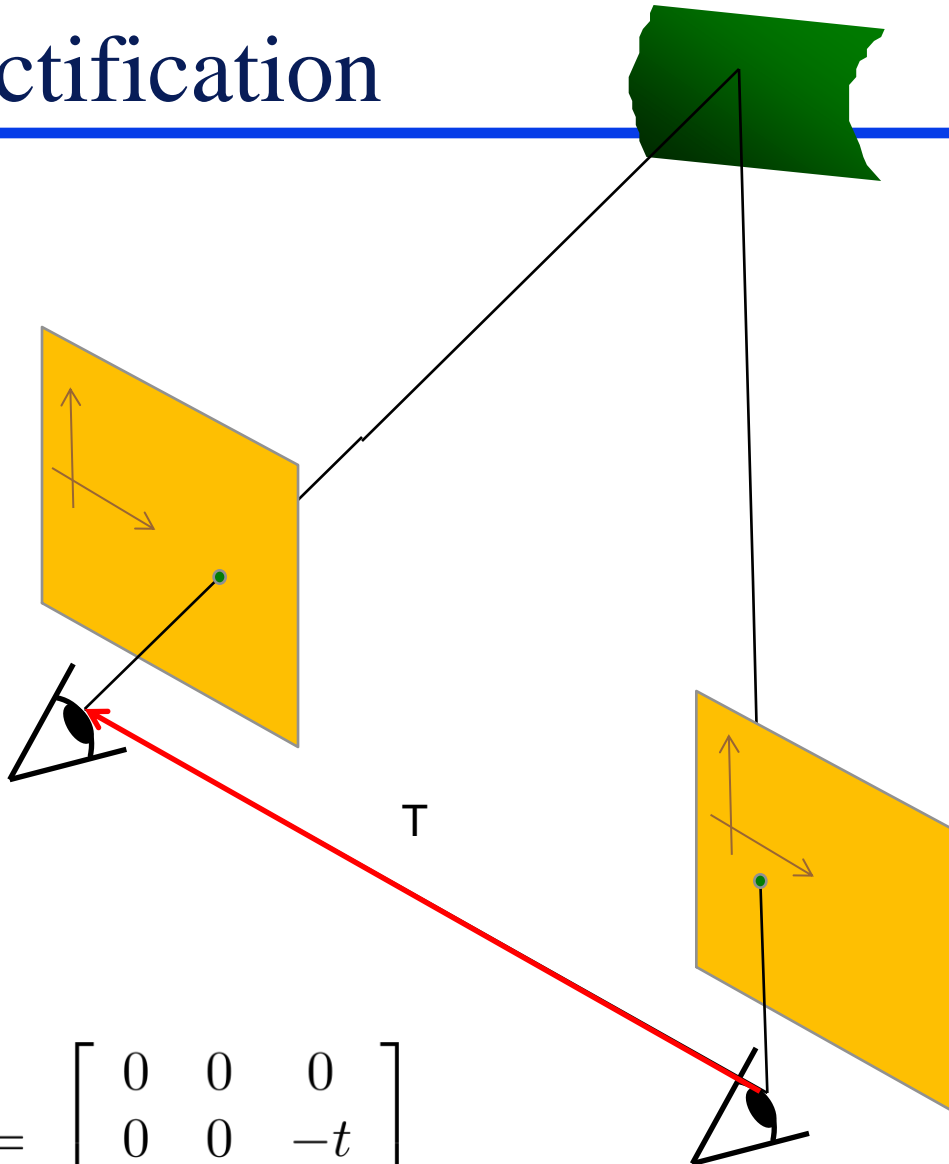
the corresponding (rectified)
camera geometry is analogous to
“panning motion”

□ Epipolar constraint:

$$R = I$$

$$T = \begin{bmatrix} t \\ 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow E = [T]_{\times} R = [T]_{\times} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -t \\ 0 & t & 0 \end{bmatrix}$$



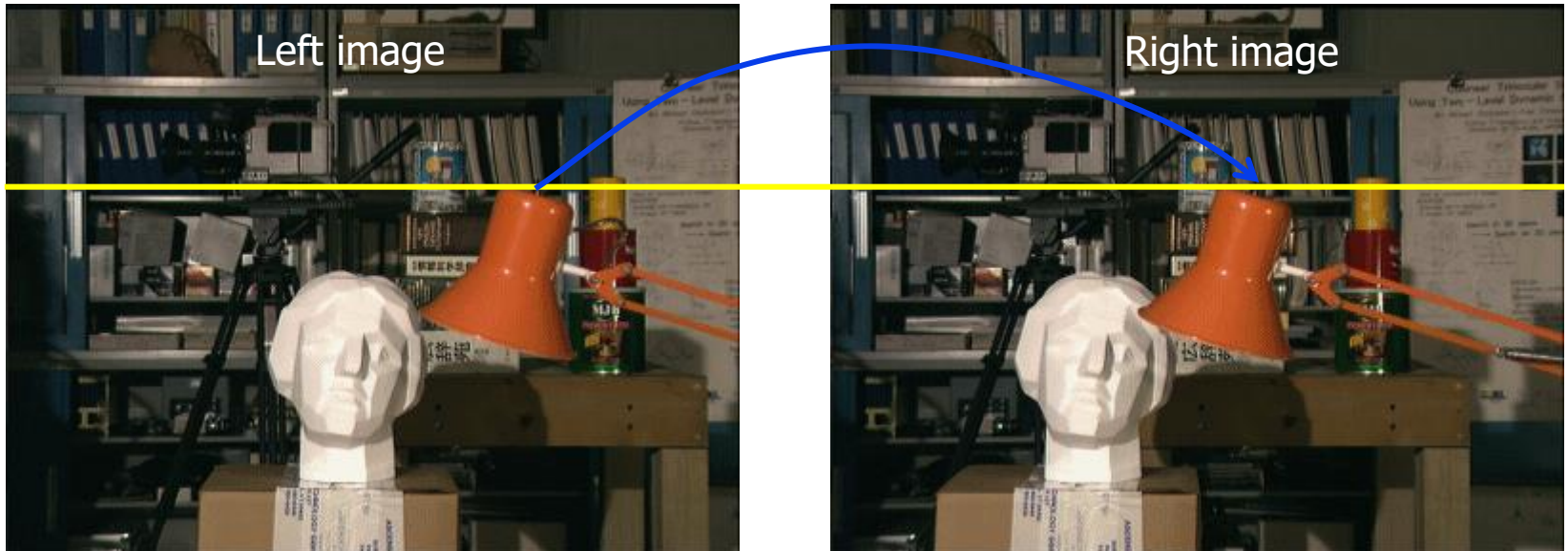
Stereo Rectification



Note projective distortion.
It will be much bigger
if images are taken from
very different viewpoints
(large baseline).

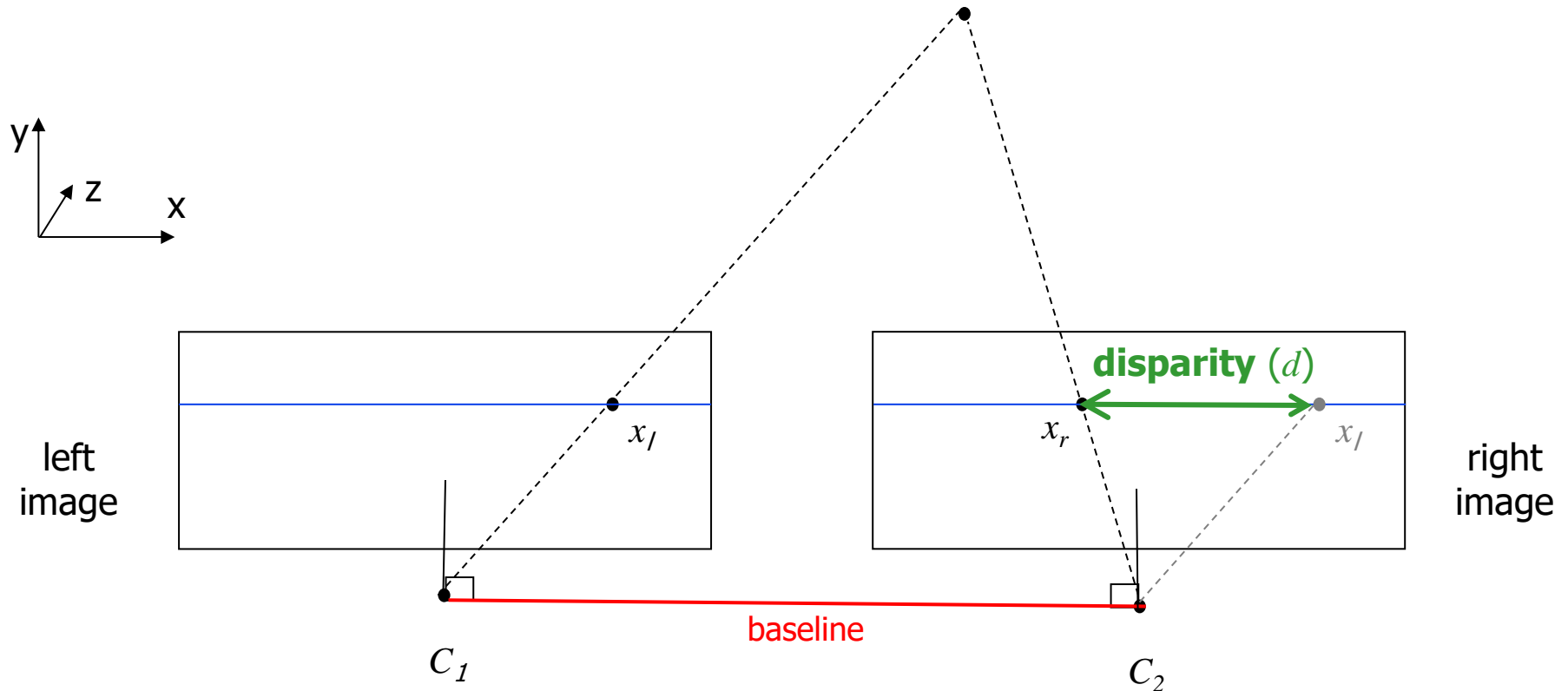
in this example the base line C_1C_2 is parallel to cube edges.

Stereo as a *correspondence* problem



(After rectification) all correspondences are along the same horizontal scan lines
(epipolar lines)

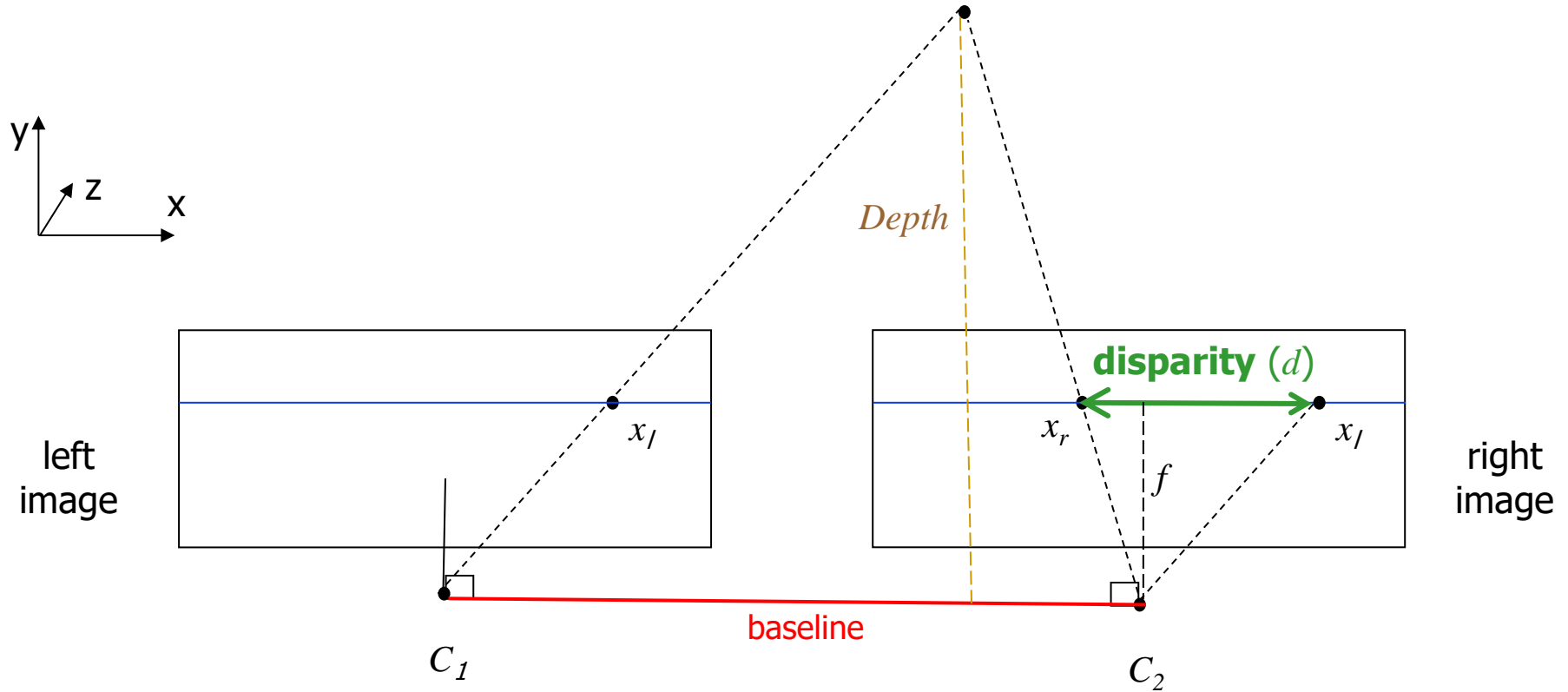
Rectified Cameras



epipolar lines are parallel to the x axis

difference between the x -coordinates of x_l and x_r is called the disparity

Rectified Cameras



$$\text{Depth} = |C_1 C_2| \cdot f / d$$

Stereo

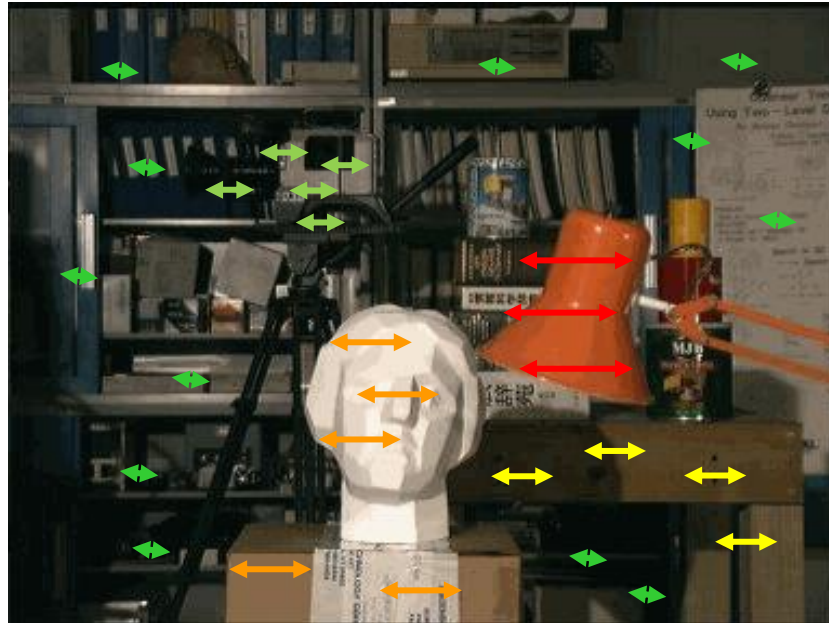


Correspondences are described by shifts along horizontal scan lines (epipolar lines)

which can be represented by scalars (**disparities**)

Stereo

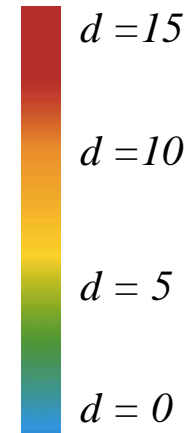
closer objects (smaller depths) correspond to **larger disparities**



Correspondences are described by shifts along horizontal scan lines (epipolar lines)

which can be represented by scalars (**disparities**)

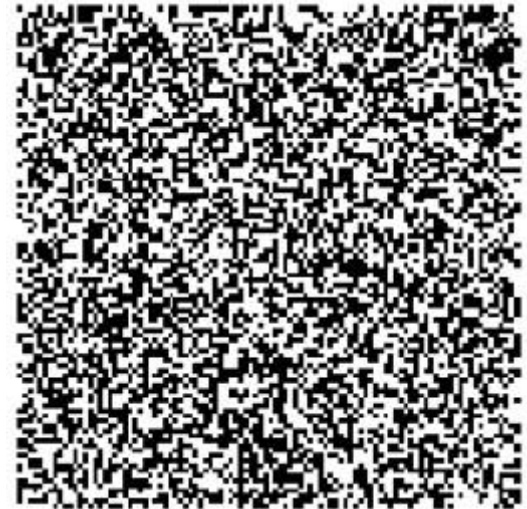
Stereo



- If x-shifts (disparities) are known for all pixels in the left (or right) image then we can visualize them as a **disparity map** – scalar valued function $d(p)$
- larger disparities correspond to closer objects

Stereo Correspondence problem

- Human vision can solve it
(even for “random dot” stereograms)

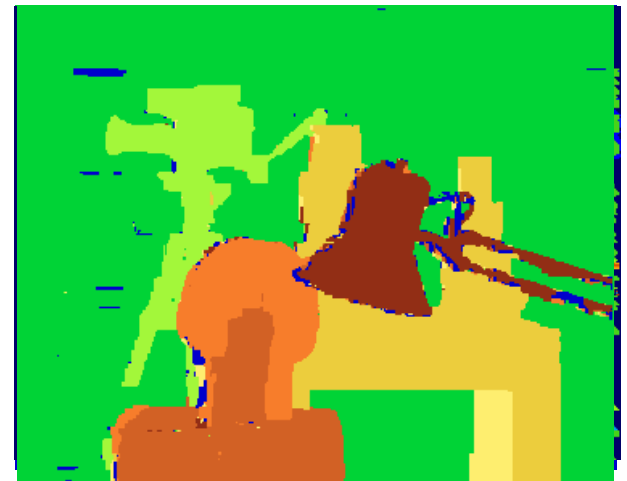


- Can computer vision solve it?

Maybe

see *Middlebury Stereo Database*
for the state-of-the art results

<http://cat.middlebury.edu/stereo/>

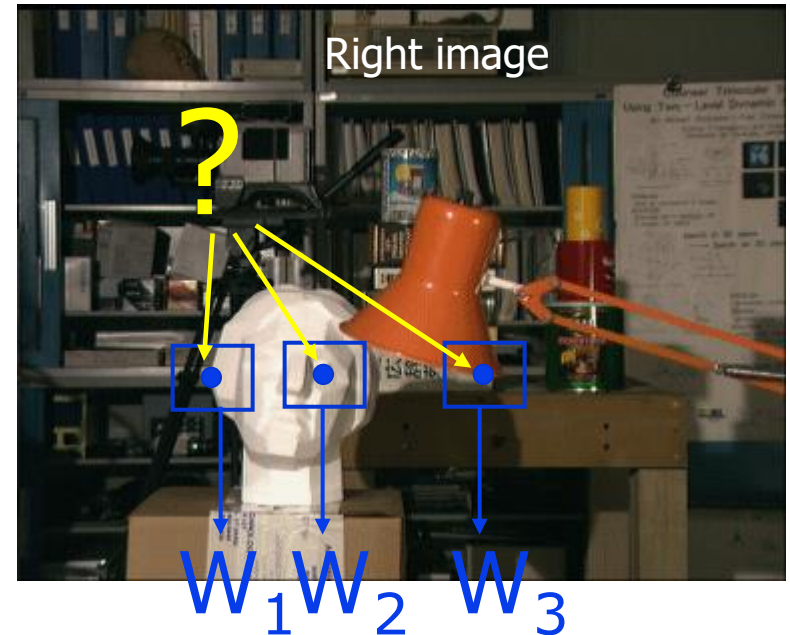
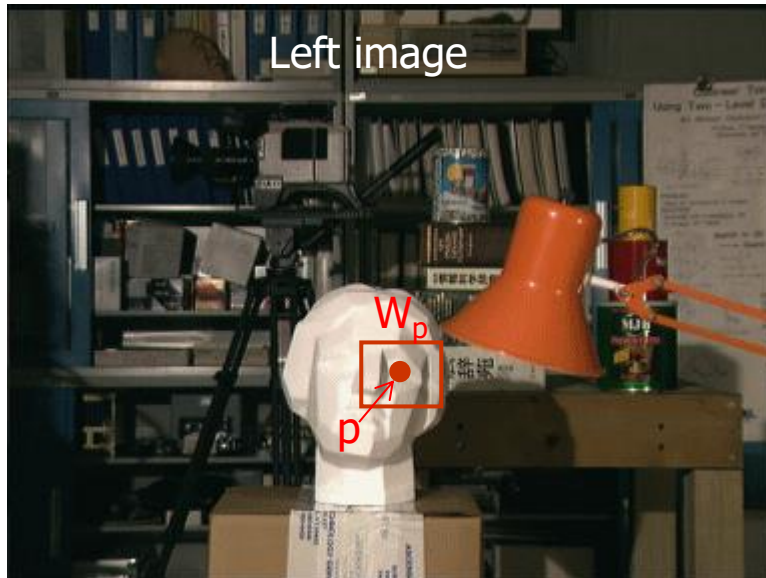


Stereo

- Window based
 - Matching rigid windows around each pixel
 - Each window is matched independently
- Scan-line based approach
 - Finding coherent correspondences for each scan-line
 - Scan-lines are independent
 - DP, *shortest paths*
- Global (muti-scan-line) approach
 - Finding coherent correspondences for all pixels (jointly)
 - spatial regularization over \mathbb{R}^2 (grid), e.g. *graph cuts*

Stereo Correspondence problem

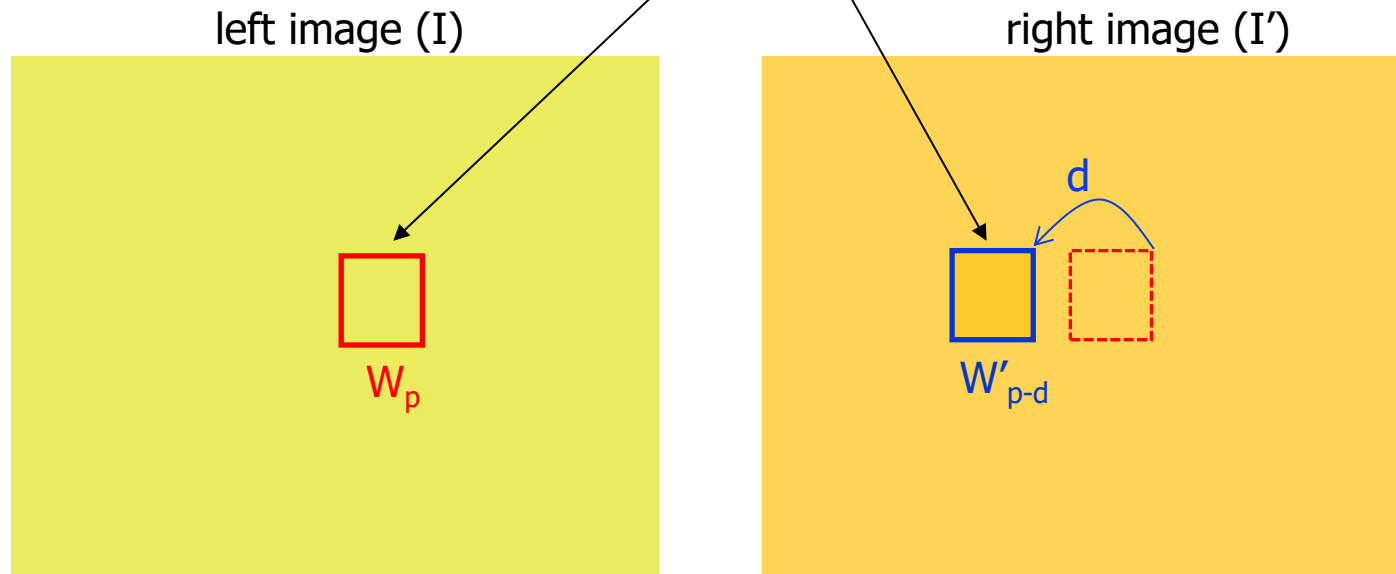
Window based approach



- For any given point p in left image consider window (or image patch) W_p around it
- Find **the best** matching window W_q on the same scan line in the right image that looks most similar to W_p

SSD (sum of squared differences) approach

$$\text{computing } SSD(p, d) = \sum_{(x, y) \in W_p} (I(x, y) - I'(x - d, y))^2$$



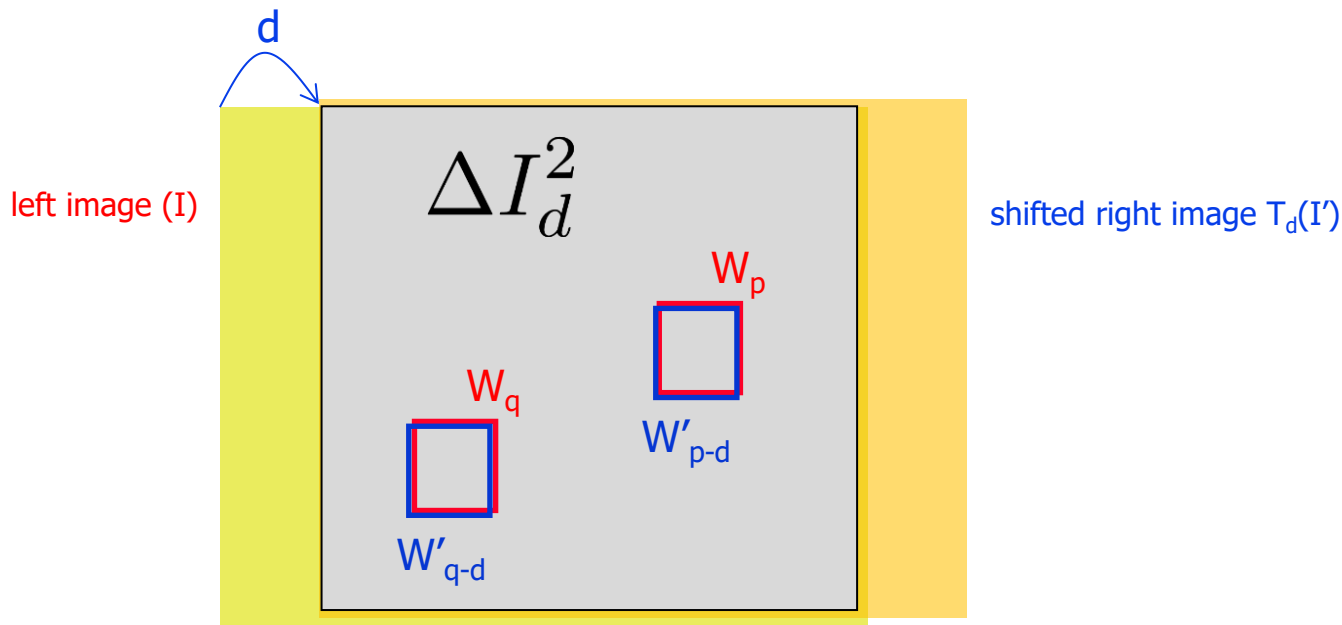
for any given pixel p compute SSD between windows W_p and W'_{p-d}
for all disparities d (in some interval $[min_d, max_d]$)

the best disparity for p can be defined as

$$\hat{d}_p = \arg \min_d SSD(p, d)$$

computing SSD efficiently

- For each fixed d , compute sq. differences image ΔI_d^2



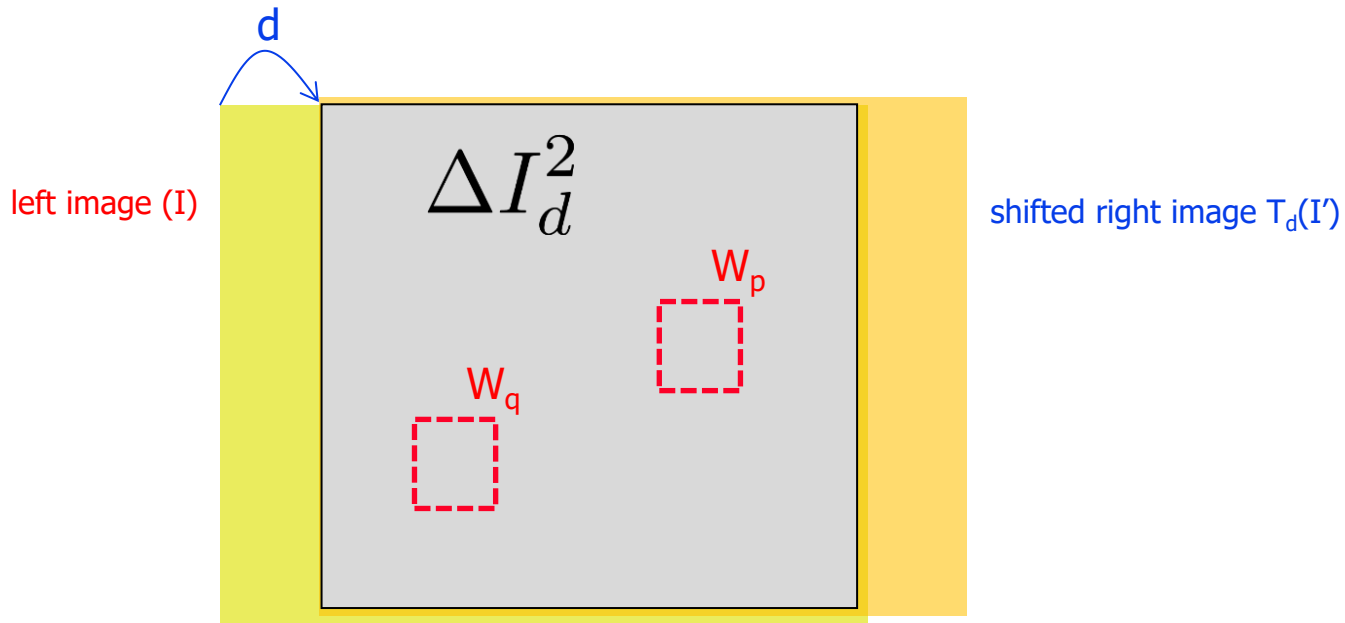
squared differences between the left image I and the shifted right image $T_d(I')$

$$\Delta I_d^2(x, y) := (I(x, y) - I'(x - d, y))^2$$

Then, $\text{SSD}(p, d)$ between W_p and W'_{p-d} is $\sum_{(x, y) \in W_p} \Delta I_d^2(x, y)$

computing SSD efficiently

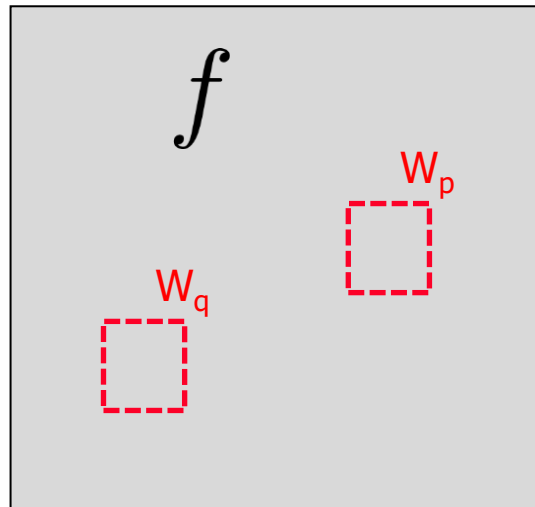
Need to sum pixel values at all possible windows



Then, SSD(p,d) between W_p and W'_{p-d} is $\sum_{(x,y) \in W_p} \Delta I_d^2(x,y)$

computing SSD efficiently

How to sum values at all possible rectangular windows □
in any given image f efficiently ?



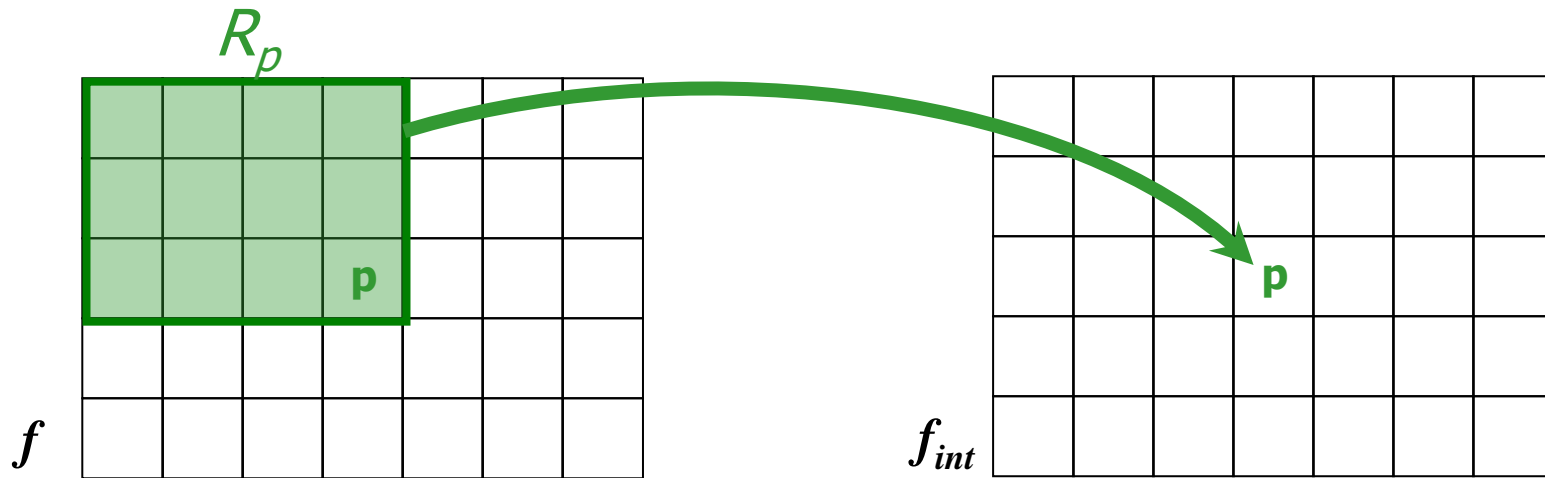
For SSD we just have special case $f(x, y) \equiv \Delta I_d^2(x, y)$

standard general trick:

“Integral Images”

$$f_{int}(p) := \sum_{q \in R_p} f(q)$$

- Define **integral image** $f_{int}(p)$ as the sum (integral) of image f over pixels in rectangle $R_p := \{q \mid “q \leq p”\}$



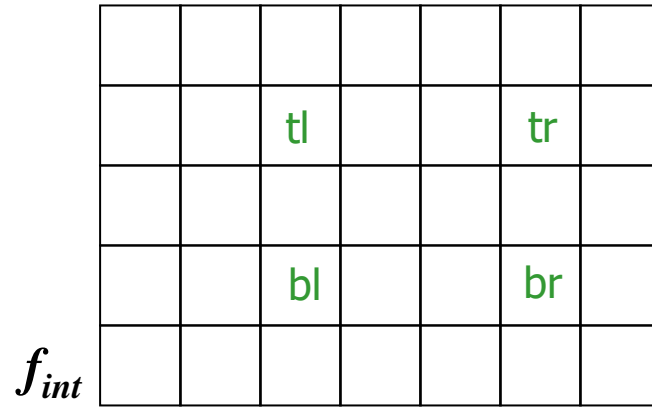
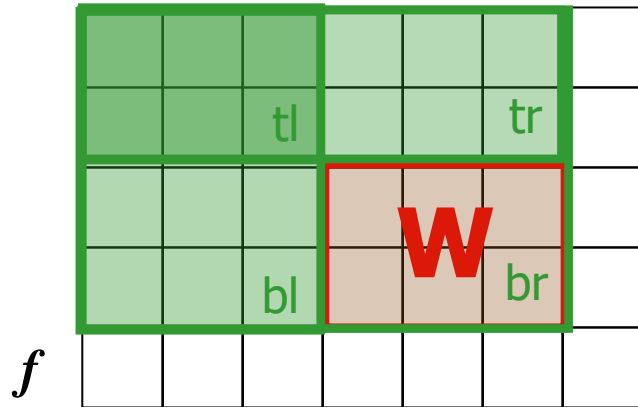
- Can compute $f_{int}(p)$ for all p in one or **two passes** over image f (How?)

standard general trick:

“Integral Images”

$$f_{int}(p) := \sum_{q \in R_p} f(q)$$

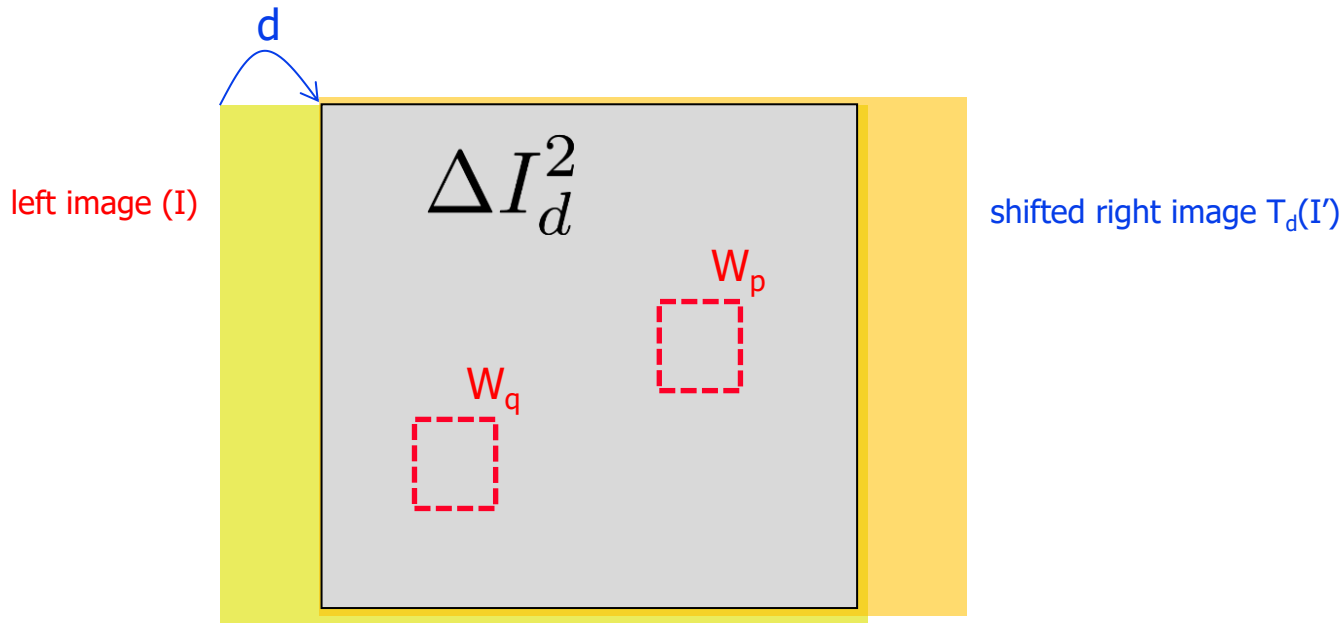
- Define **integral image** $f_{int}(p)$ as the sum (integral) of image f over pixels in rectangle $R_p := \{q \mid “q \leq p”\}$



- Now, for any W the sum (integral) of f inside that window can be computed as $\sum_{q \in W} f(q) = f_{int}(\text{br}) - f_{int}(\text{bl}) - f_{int}(\text{tr}) + f_{int}(\text{tl})$

computing SSD efficiently

For SSD we have special case $f(x, y) \equiv \Delta I_d^2(x, y)$



Now, the sum of ΔI_d^2 at any window takes 4 operations independently of window size

=> $O(|I| * |d|)$ window-based stereo algorithm

Problems with Fixed Windows

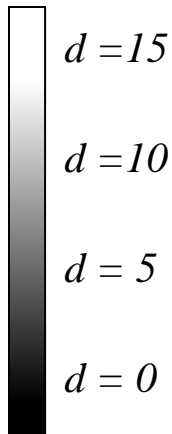
disparity maps $\hat{d}_p = \arg \min_d SSD(p, d)$ for:

small window



- better at boundaries
- noisy in low texture areas

large window



- better in low texture areas
- blurred boundaries

Q: what do we implicitly assume when using low $SSD(d, p)$ at a window around pixel p as a criteria for “good” disparity d ?

window algorithms

- Maybe variable window size (pixel specific)?
 - What is the right window size?
 - Correspondences are still found independently at each pixel (no coherence)

- All window-based solutions can be thought of as “local” solutions - but very fast!

- How to go to “global” solutions?
 - use *objectives*, a.k.a. *energy* or *loss* functions
 - *surface regularization* or *spatial coherence*
 - optimization

need priors
to compensate
for local data ambiguity

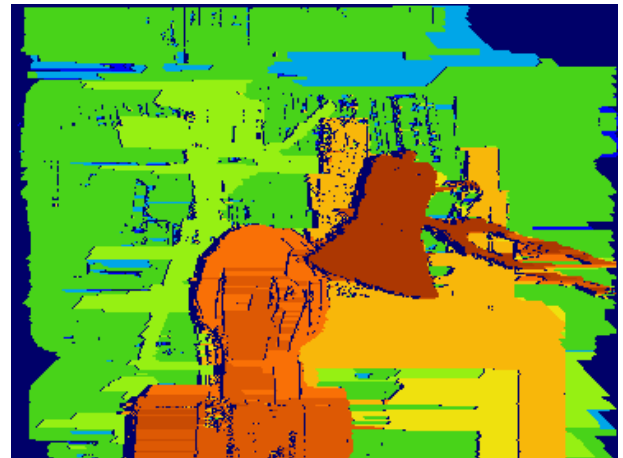
Stereo Correspondence problem

Scan-line approach

□ Scan-line stereo

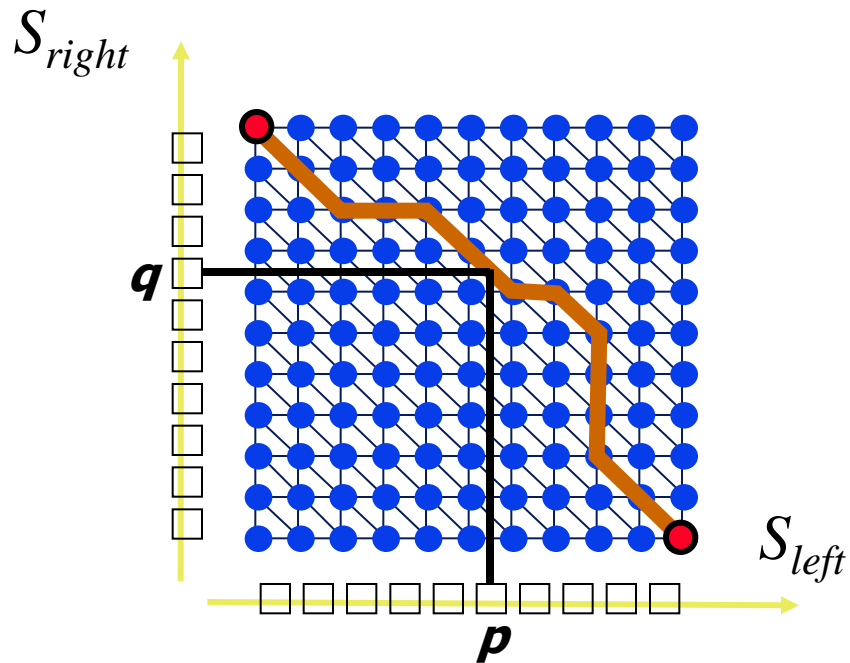
- coherently match pixels in each scan line
- DP or shortest paths work (easy 1D optimization)
- Note: scan lines are still matched independently

– streaking artifacts

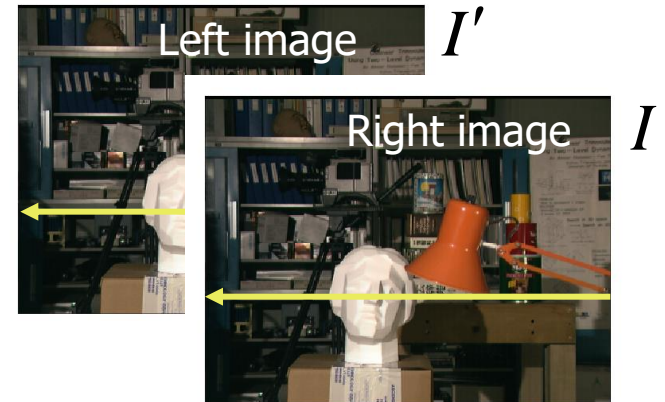


“Shortest paths” for Scan-line stereo

e.g. Ohta&Kanade’85, Cox et.al.’96

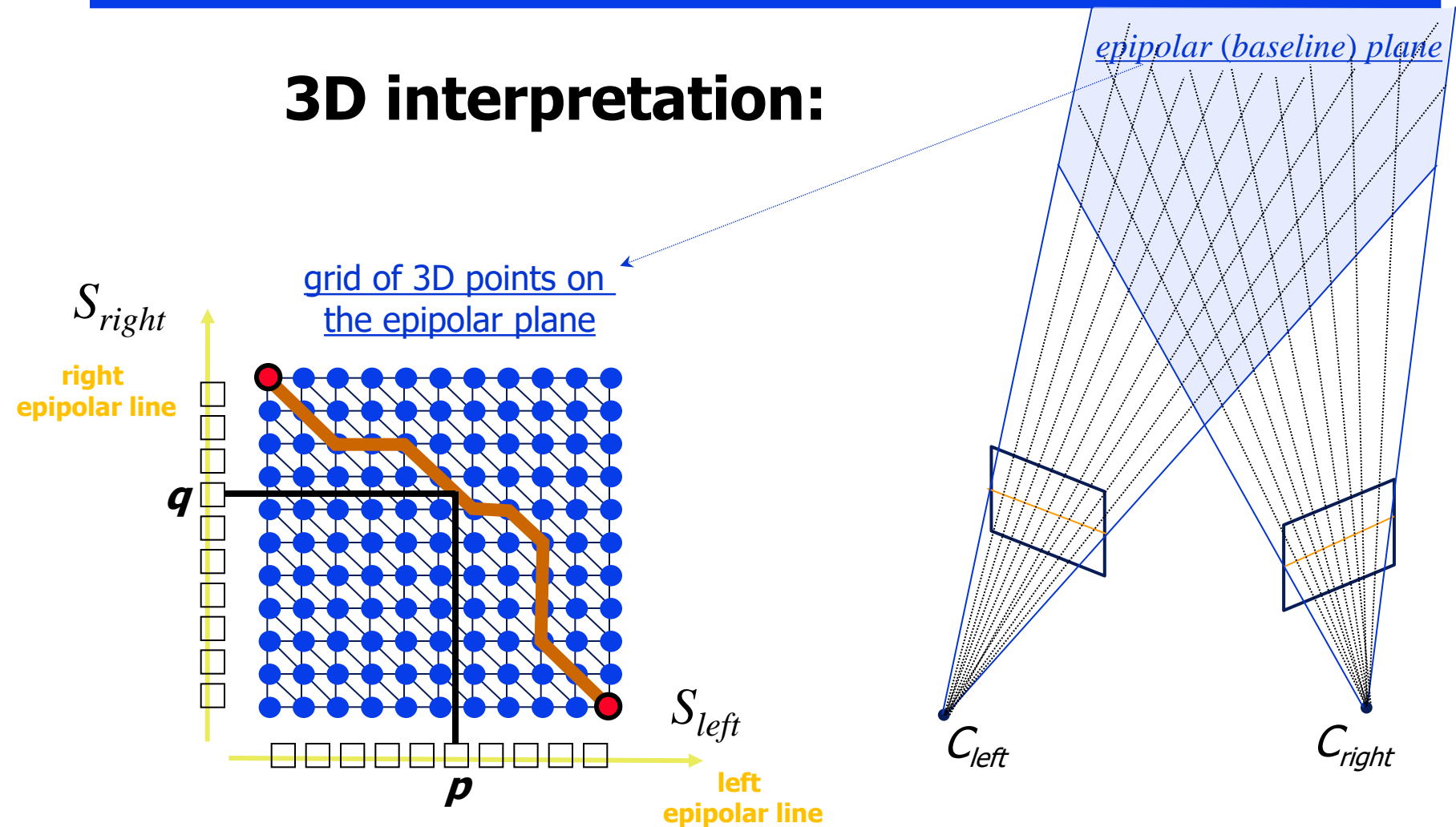


a **path** on this graph represents a matching function



“Shortest paths” for Scan-line stereo

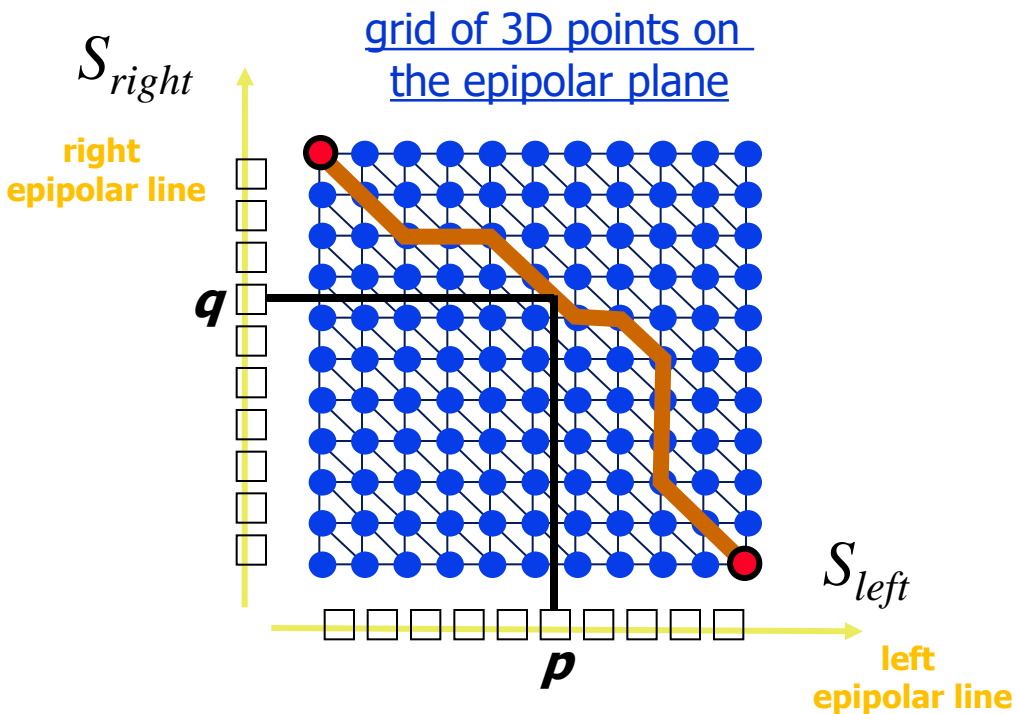
3D interpretation:



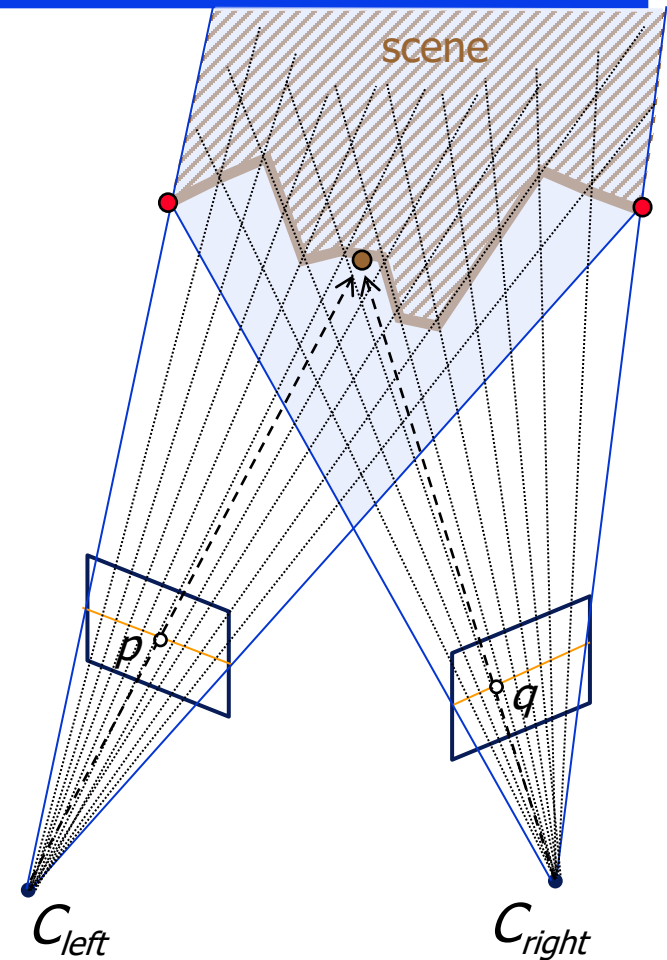
a **path** on this graph represents a matching function

“Shortest paths” for Scan-line stereo

3D interpretation:



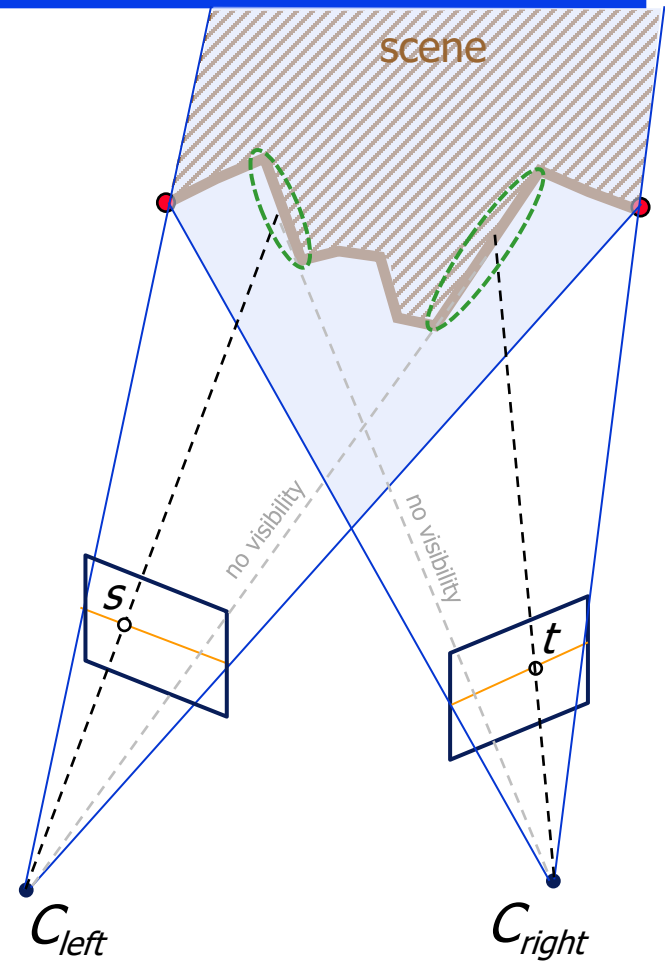
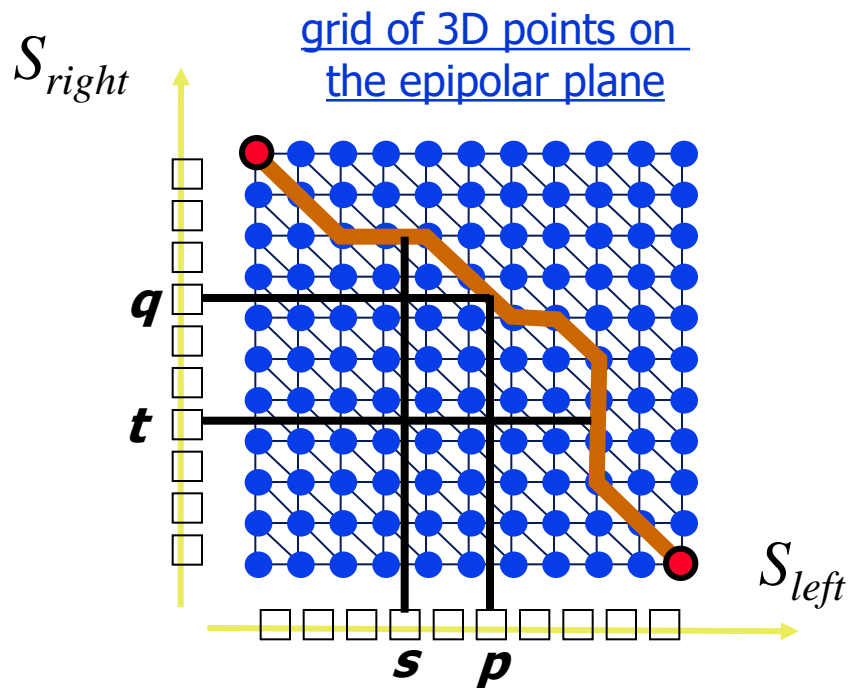
a **path** on this graph represents a matching function



This **path** corresponds to
an intersection of **epipolar plane**
with **3D scene surface**

“Shortest paths” for Scan-line stereo

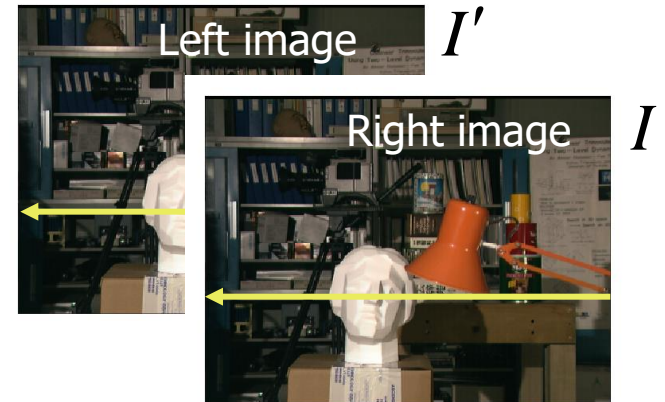
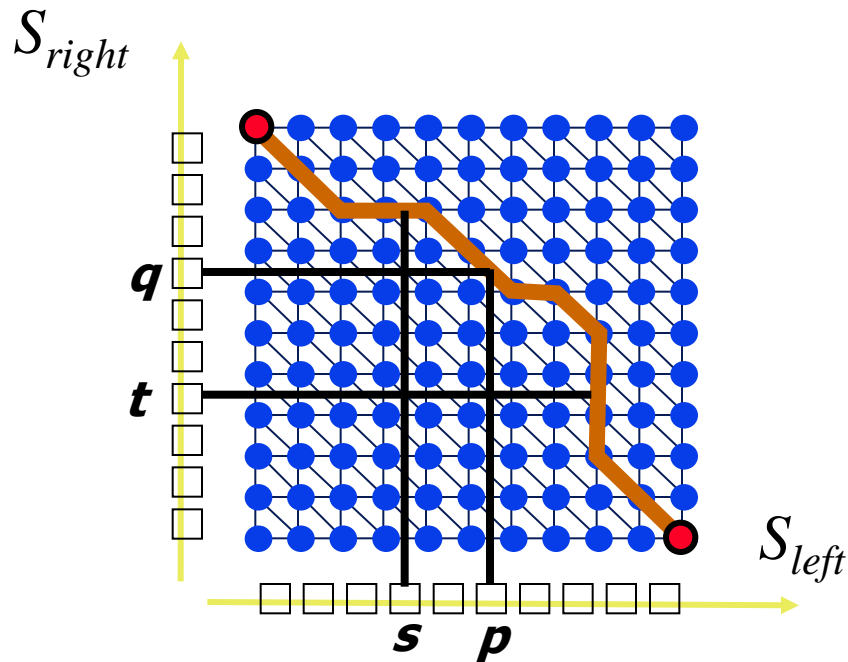
3D interpretation:



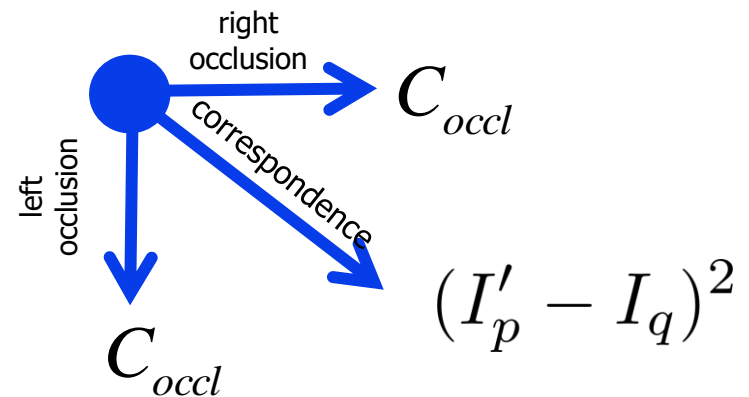
horizontal and vertical edges on the path imply “no correspondence” (*occlusion*)

“Shortest paths” for Scan-line stereo

e.g. Ohta&Kanade’85, Cox et.al.’96

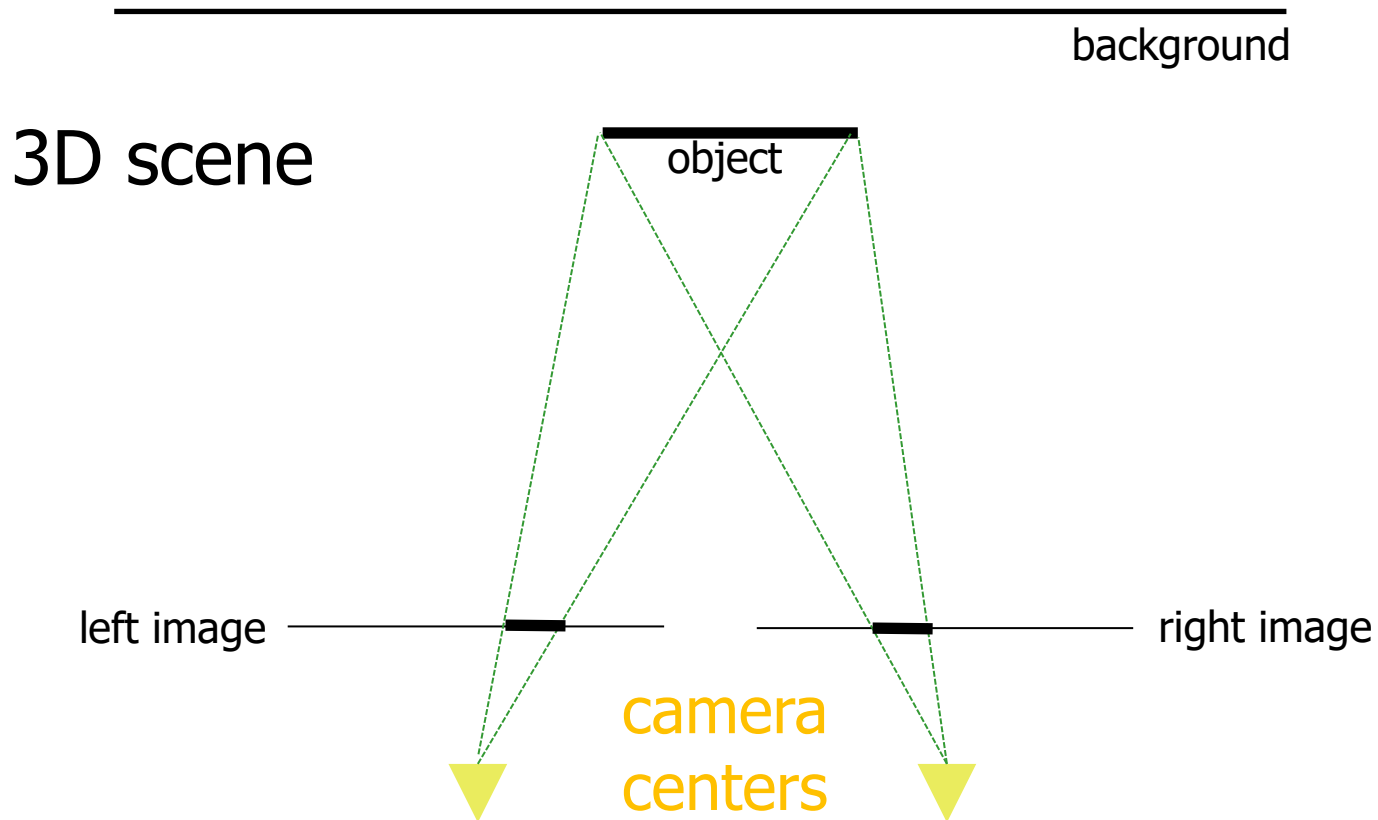


Edge weights:

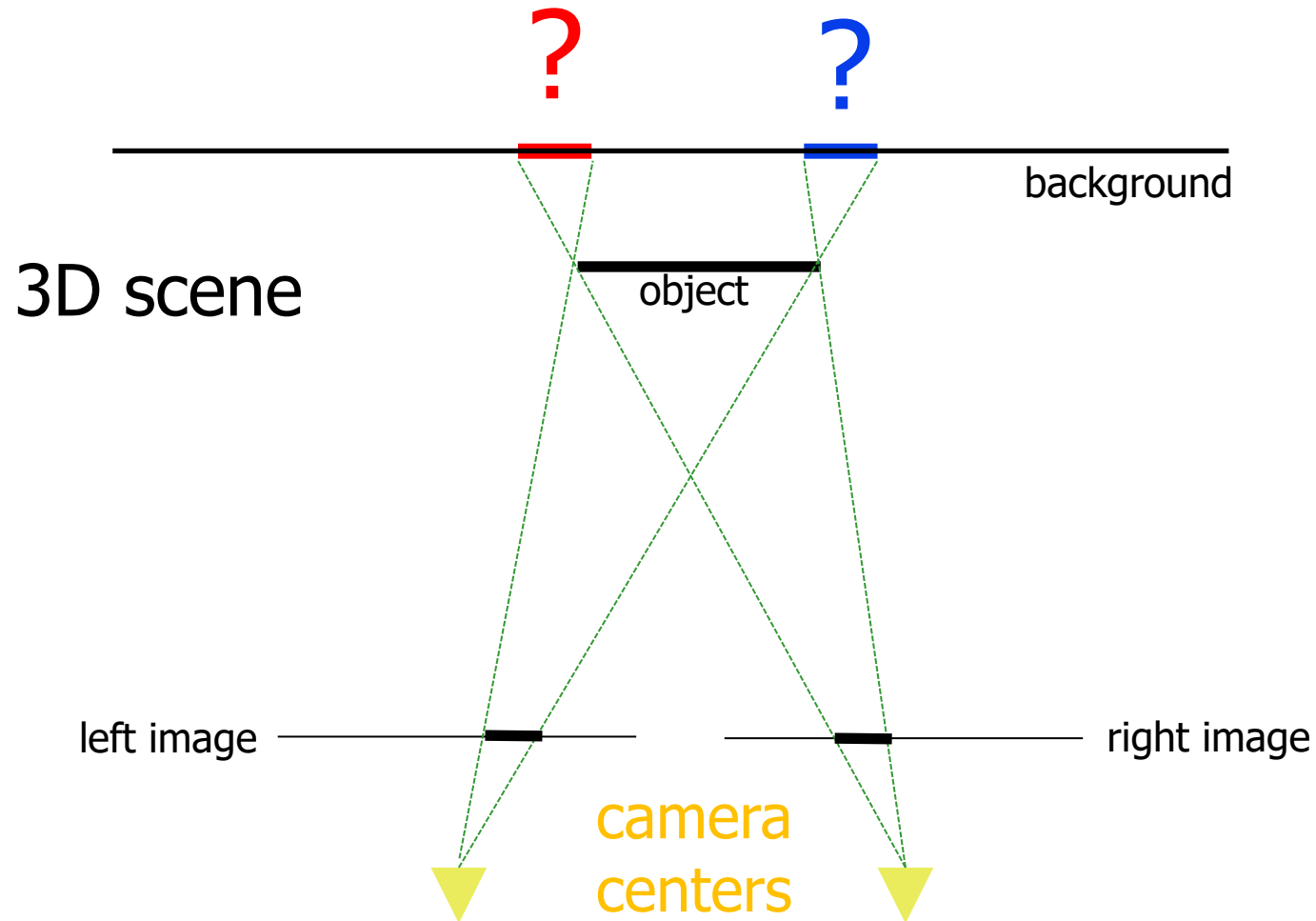


What is “occlusion” in general ?

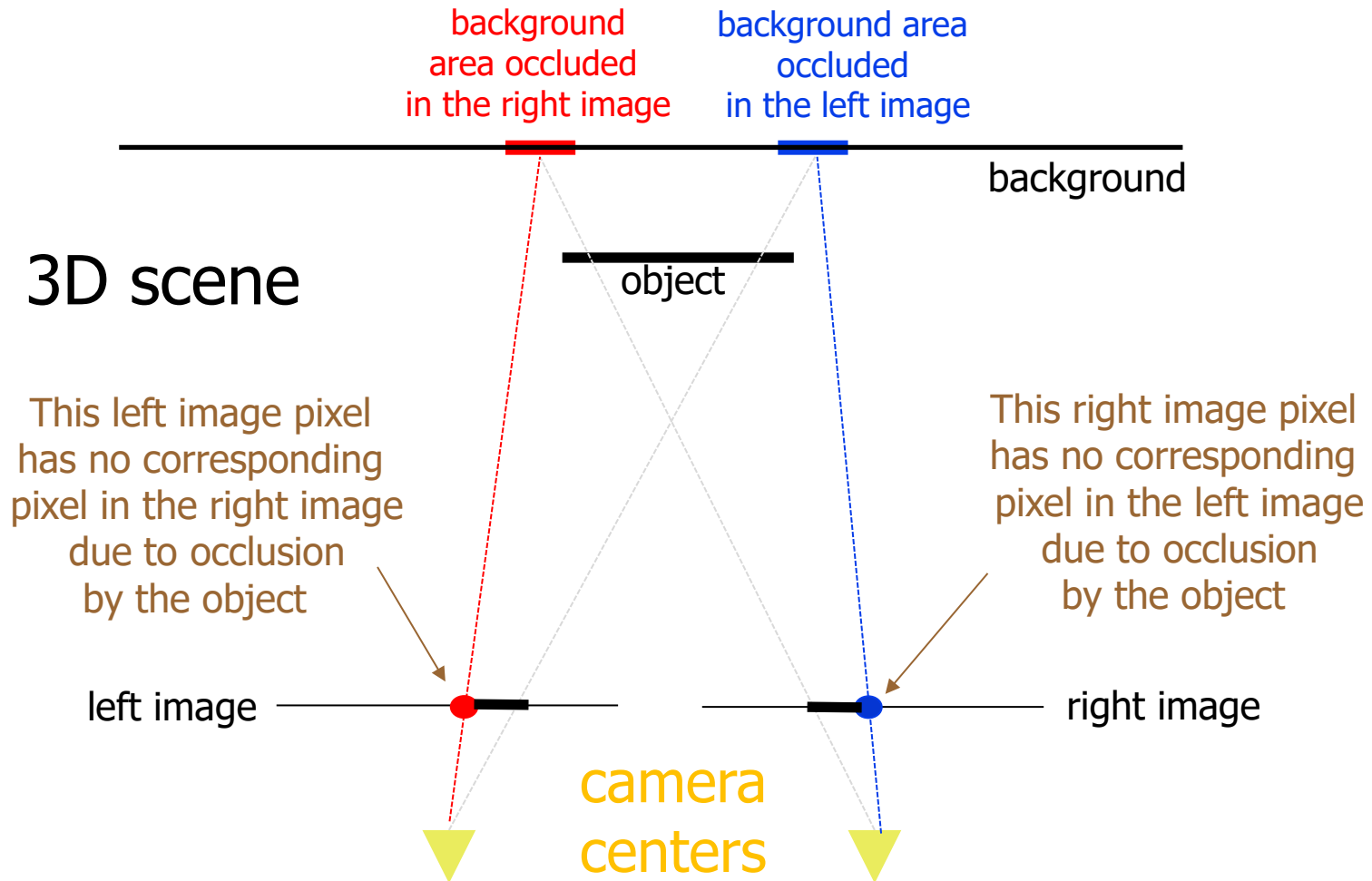
Occlusion in stereo



Occlusion in stereo

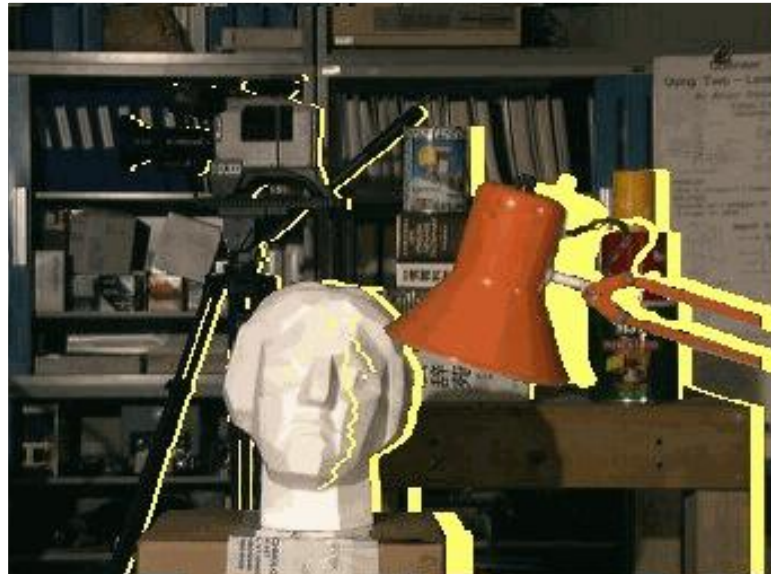


Occlusion in stereo



Note: occlusions occur at depth discontinuities/jumps

Stereo

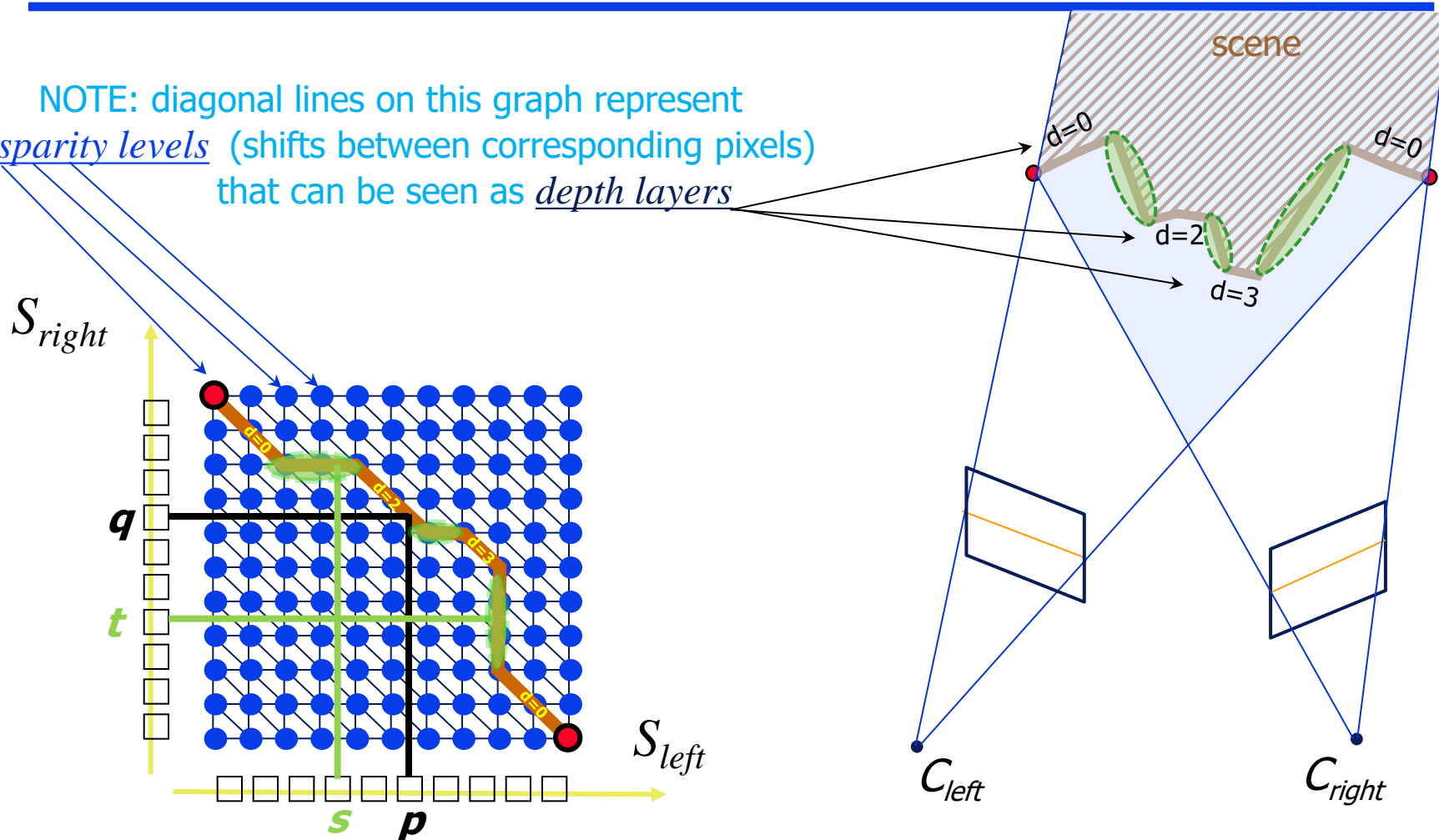


yellow marks occluded points in different viewpoints
(points not visible from the central/base viewpoint).

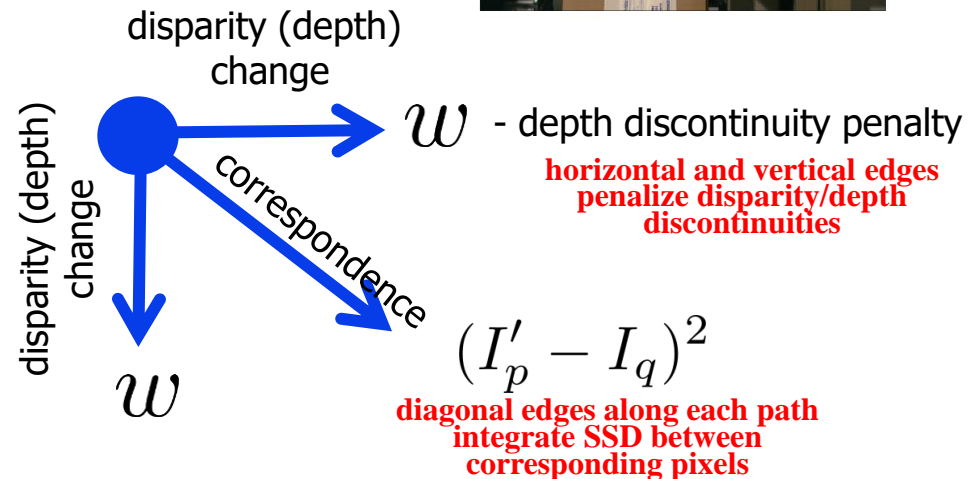
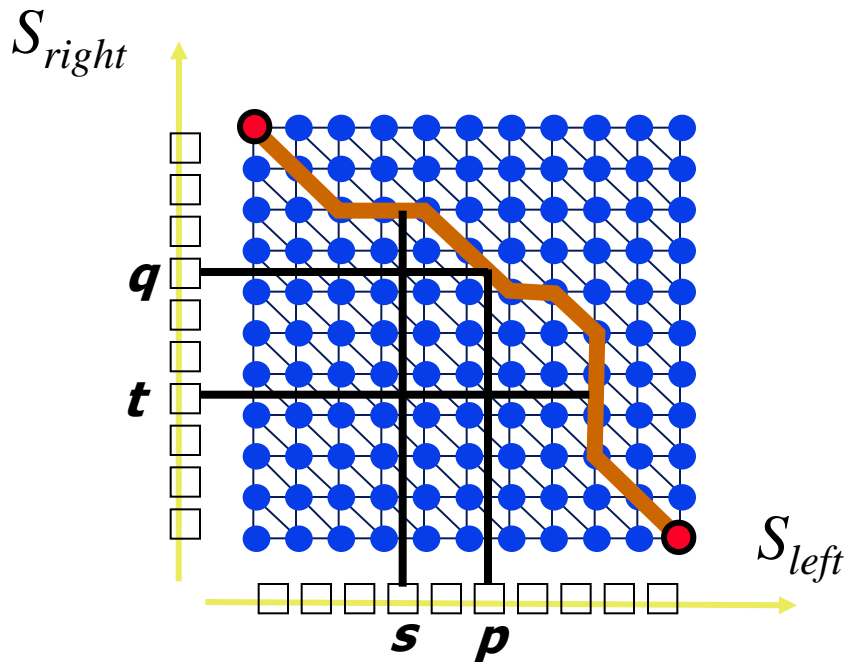
Note: occlusions occur at depth discontinuities/jumps

Occlusions vs disparity/depth jumps

NOTE: diagonal lines on this graph represent disparity levels (shifts between corresponding pixels) that can be seen as depth layers



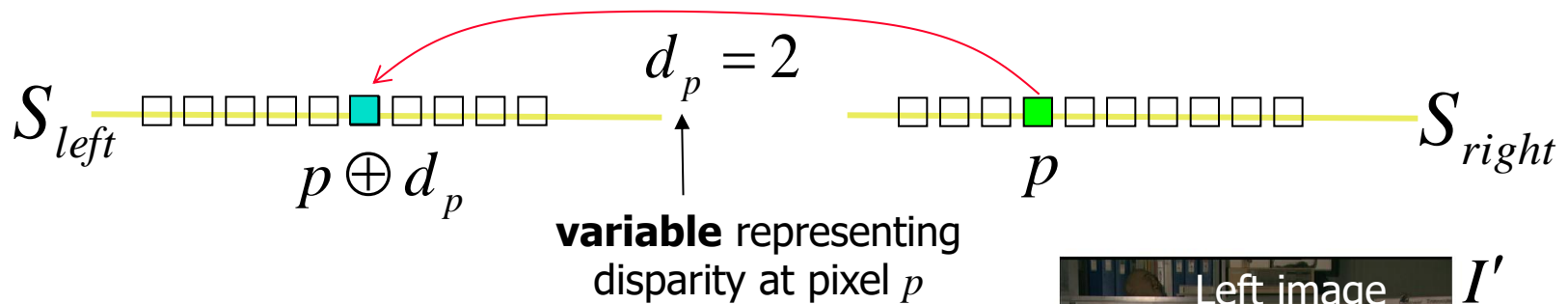
horizontal and vertical edges on this graph describe **occlusions**,
as well as **disparity jumps** or **depth discontinuities**



But, the actual implementation in OK'85 and C'96 uses *Viterbi* algorithm (DP)

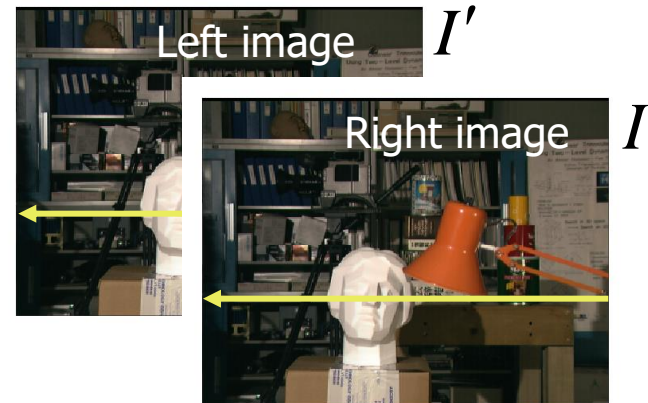
explicitly assigning “optimal” disparity labels d_n to all pixels p as follows...

More common representation of disparity map in stereo algorithms

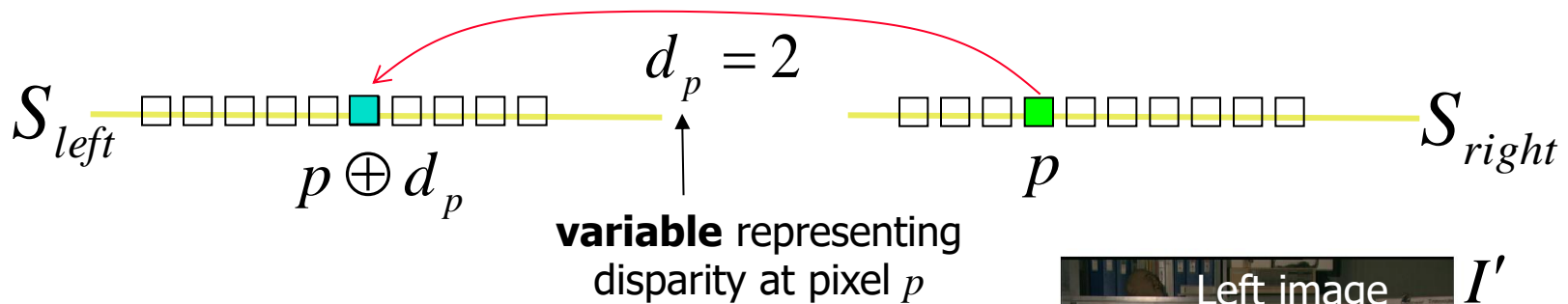


$$\mathbf{d} = \{d_p \mid p \in G\}$$

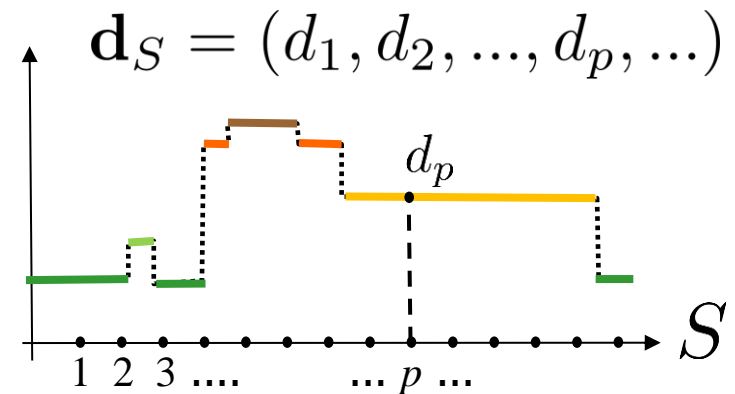
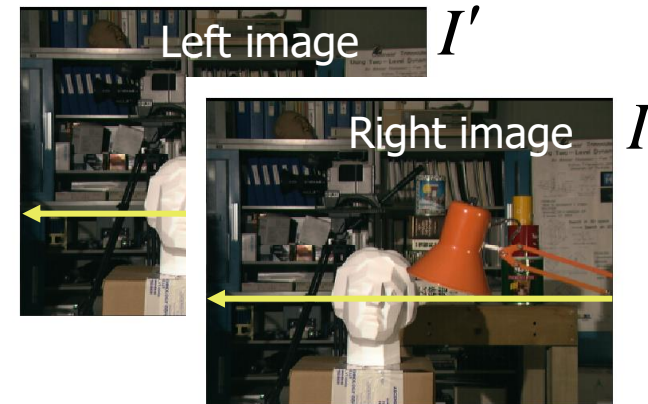
configuration (map) of disparities
for pixels p on **image grid**



More common representation of disparity map in stereo algorithms

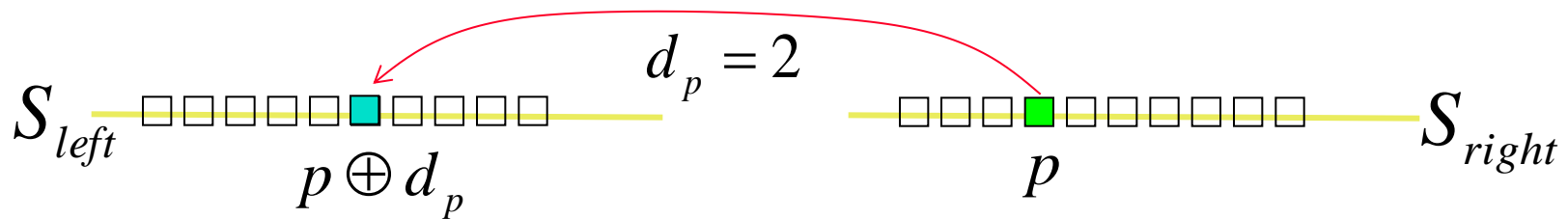


$\mathbf{d}_S = \{d_p \mid p \in S\}$
 configuration (map) of disparities
 for pixels p on a **scan line**

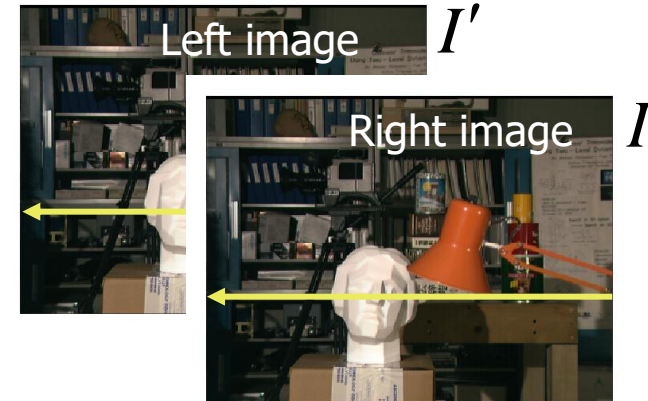


NOTE: *Loss, Energy, Cost, and Objective function* mean the same thing and are used indiscriminately in different contexts

DP for scan-line stereo



Viterbi algorithm can find an optimal *disparity configuration* $\mathbf{d}_S = (d_1, d_2, \dots, d_p, \dots)$ for pixels p on a given scan-line S_{right} minimizing the following **loss function**



$$E(\mathbf{d}_S) = \sum_{p \in S} \underbrace{D_p(d_p)}_{\substack{|I_p - I'_{p \oplus d_p}| \\ \text{photo consistency}}} + \sum_{p \in S} \underbrace{V(d_p, d_{p+1})}_{\substack{w |d_p - d_{p+1}| \\ \text{spatial coherence}}} = \sum_{\{p, q\} \in N} E(d_p, d_q)$$

Such pairwise loss can be optimized in $O(nm^2)$ on non-loopy graphs (e.g., **chains**)

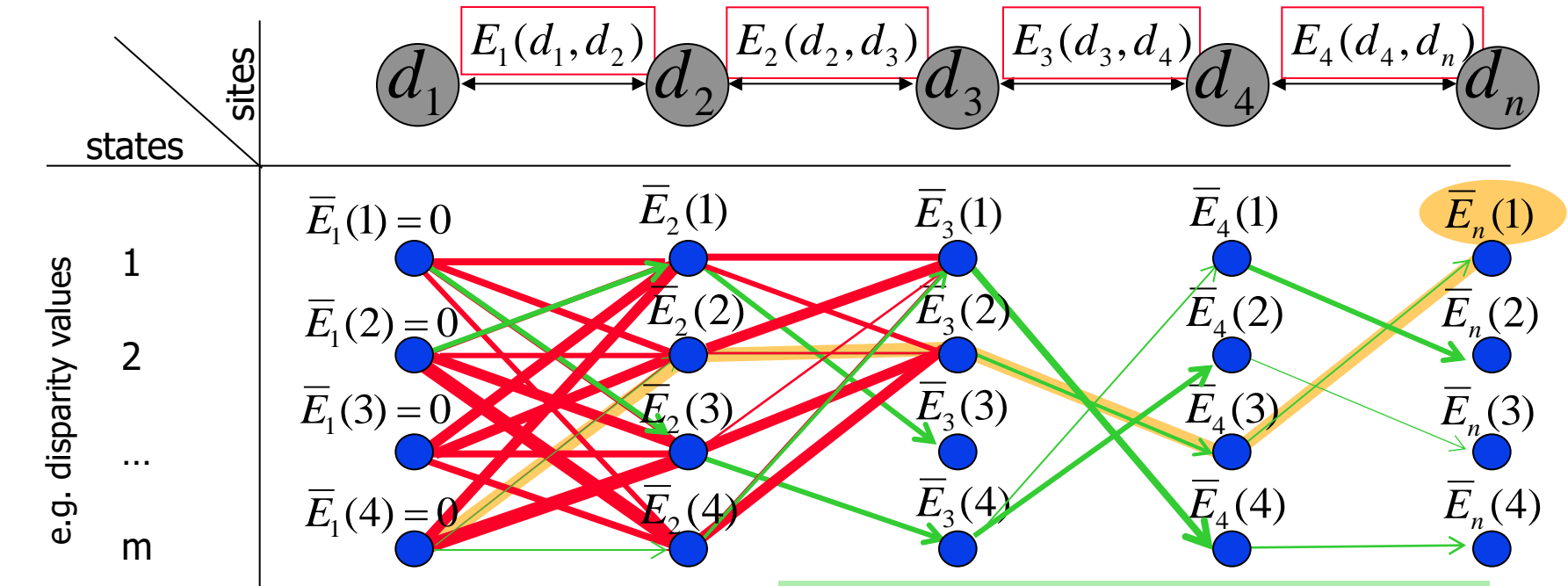
Dynamic Programming (DP)

Viterbi Algorithm

$$d_i \in \{0, 1, \dots, m-1\}$$

Consider **pair-wise interactions** between sites (pixels) on a **chain** (scan-line)

$$E(d_1, \dots, d_n) = E_1(d_1, d_2) + E_2(d_2, d_3) + \dots + E_{n-1}(d_{n-1}, d_n)$$



$\bar{E}_p(k)$ - internal "loss counter" at "site" p and "state" k

$$\bar{E}_{p+1}(k) = \min_i (\bar{E}_p(i) + E_p(i, k))$$

Complexity: $O(nm^2)$, worst case = best case

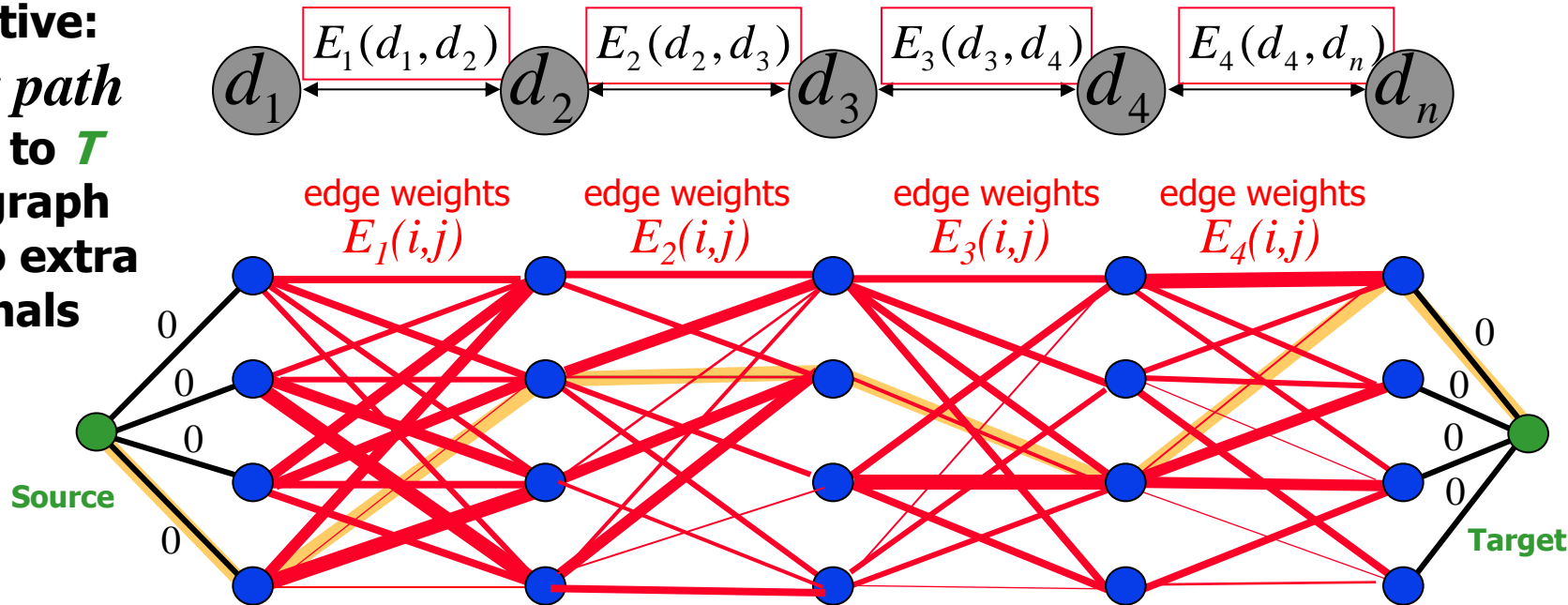
Dynamic Programming (DP)

Shortest paths Algorithm

Consider **pair-wise interactions** between sites (pixels) on a **chain** (scan-line)

$$E_1(d_1, d_2) + E_2(d_2, d_3) + \dots + E_{n-1}(d_{n-1}, d_n)$$

Alternative:
shortest path
from **S** to **T**
on the graph
with two extra
terminals



Complexity: $O(nm^2 + nm \log(nm))$ - worst case

But, the best case could be better than Viterbi. Why?

Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\parallel} + \sum_{\{p, q\} \in \underbrace{N}_{\parallel}} \underbrace{V(d_p, d_q)}_{\parallel}$$

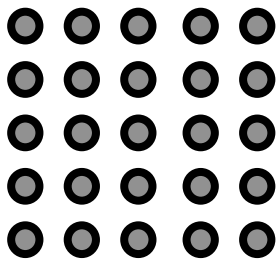
$$\boxed{|I_p - I'_{p \oplus d_p}|}$$

photo consistency

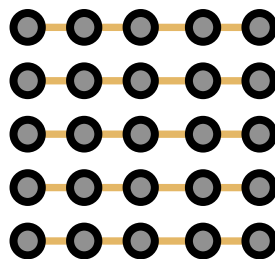
$$\boxed{w |d_p - d_q|}$$

spatial coherence

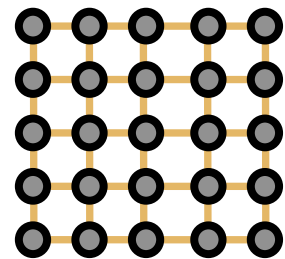
Consider three different neighborhood systems N :



$$N = \emptyset$$



$$N = \{\{p, p \pm 1\} : p \in G\}$$



$$N = \{\{p, q\} \subset G : |pq| \leq 1\}$$

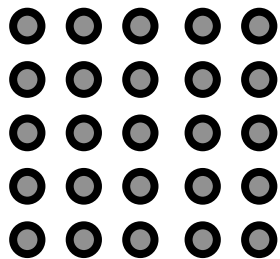
Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\parallel \begin{array}{c} |I_p - I'_{p \oplus d_p}| \\ \text{photo consistency} \end{array}} + \sum_{\{p, q\} \in \underbrace{N}_{\parallel \begin{array}{c} w |d_p - d_q| \\ \text{spatial coherence} \end{array}}} \underbrace{V(d_p, d_q)}_{\text{smoothness term disappears}}$$

CASE 1



$N = \emptyset$

Q: how to optimize $E(\mathbf{d})$ in this case?

$$\forall p \in G \quad \hat{d}_p = \arg \min_d D_p(d) \quad O(nm)$$

Q: How does this relate to window-based stereo?

Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

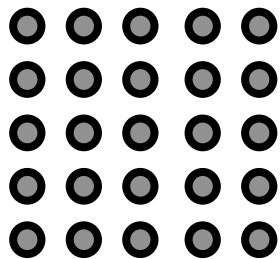
Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\substack{|| \\ |I_p - I'_{p \oplus d_p}| \\ \text{photo consistency}}} + \sum_{\{p, q\} \in \underbrace{N}_{\substack{|| \\ w |d_p - d_q| \\ \text{spatial coherence}}}} V(d_p, d_q)$$

smoothness term disappears

CASE 1



$N = \emptyset$

Nodes/pixels do not interact (are independent).
Optimization of the sum of **unary terms**,

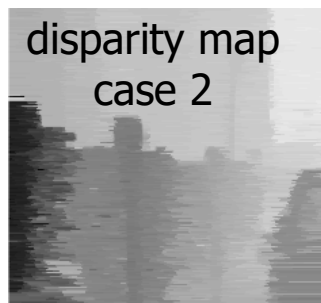
e.g. $\sum_{p \in G} D_p(d_p)$, is trivial: $O(nm)$

Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

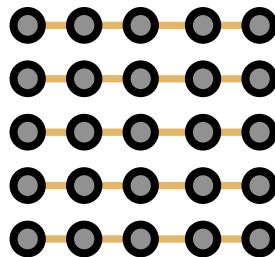
Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\substack{|I_p - I'_{p \oplus d_p}| \\ \text{photo consistency}}} + \underbrace{\sum_{\{p,q\} \in N} \underbrace{V(d_p, d_q)}_{\substack{w |d_p - d_q| \\ \text{spatial coherence}}}}_{\text{spatial coherence}}$$



CASE 2



Pairwise coherence is enforced,
but only between pixels on
the same scan line.

Q: how do we optimize such $E(\mathbf{d})$?

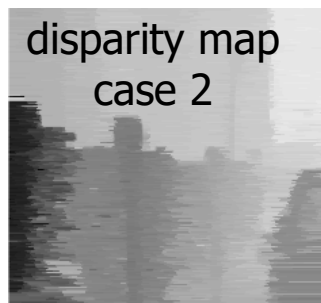
$$N = \{\{p, p \pm l\} : p \in G\}$$

$$O(nm^2)$$

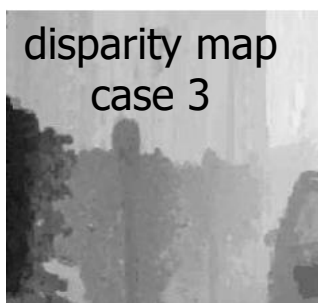
Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

Consider energy (loss) function over disparities
 $\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\substack{|I_p - I'_{p \oplus d_p}| \\ \text{photo consistency}}} + \underbrace{\sum_{\{p,q\} \in N} \underbrace{V(d_p, d_q)}_{\substack{w |d_p - d_q| \\ \text{spatial coherence}}}}_{\text{spatial coherence}}$$



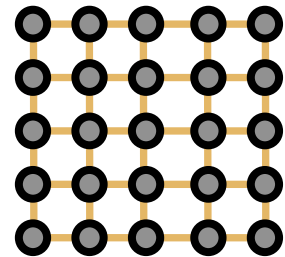
disparity map
case 2



disparity map
case 3

Pairwise smoothness
of the disparity map
is enforced both
horizontally and
vertically.

CASE 3



$$N = \{\{p, q\} \subset G : |pq| \leq 1\}$$

NOTE: *depth map* regularity/smoothness should be isotropic since
 3D scene surface is independent of scan-lines (epipolar lines) orientation.

Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\parallel} + \sum_{\{p,q\} \in \underbrace{N}_{\parallel}} \underbrace{V(d_p, d_q)}_{\parallel}$$

$|I_p - I'_{p \oplus d_p}|$
photo consistency

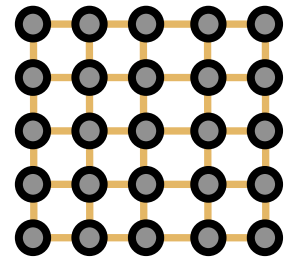
$w |d_p - d_q|$
spatial coherence

How to optimize “pairwise” loss on loopy graphs?

NOTE 1: **Viterbi does not apply**, but its extensions (e.g. *message passing*) provide approximate solutions on loopy graphs.

NOTE 2: “*Gradient descent*” can find only local minima for a continuous relaxation of $E(\mathbf{d})$ combining non-convex photo-consistency (1st term) and convex *total variation of \mathbf{d}* (2nd term).

CASE 3



$$N = \{\{p, q\} \subset G : |pq| \leq 1\}$$

Estimating (optimizing) disparities: over **points** vs. **scan-lines** vs. **grid**

Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\parallel} + \sum_{\{p,q\} \in \underbrace{N}_{\parallel}} \underbrace{V(d_p, d_q)}_{\parallel}$$

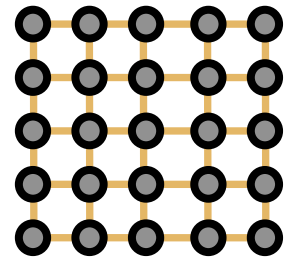
$|I_p - I'_{p \oplus d_p}|$
photo consistency

$w |d_p - d_q|$
spatial coherence

Can **globally minimize**
such pairwise losses
over any neighborhood N ,
e.g. *graph cut* algorithms

(**convex potentials V only** [Ishikawa et al 1998])

CASE 3



$$N = \{\{p, q\} \subset G : |pq| \leq 1\}$$

Useful extension:

local affinities w_{pq} and “edge alignment”

Consider energy (loss) function over disparities

$\mathbf{d} = \{d_p \mid p \in G\}$ for pixels p on grid G

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\parallel} + \sum_{\{p, q\} \in N} \underbrace{V(d_p, d_q)}_{\parallel}$$

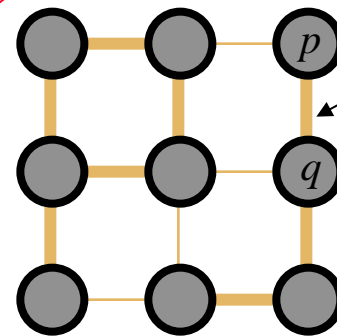
$|I_p - I'_{p \oplus d_p}|$

photo consistency

$w_{pq} \cdot |d_p - d_q|$

spatial coherence

In general, one can use
pairwise affinities w_{pq}
specific to each pair of
neighboring pixels p and q



w_{pq} - weights of
neighborhood edges
(e.g. may be assigned
according to local
intensity contrast in
the **reference image**)

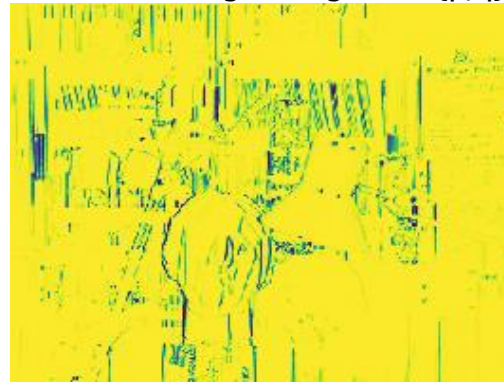
Useful extension: local affinities w_{pq} and “edge alignment”

visualization of weights w_{pq}

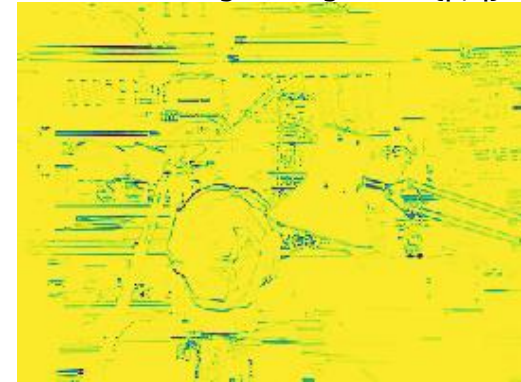
reference (e.g. right) image I



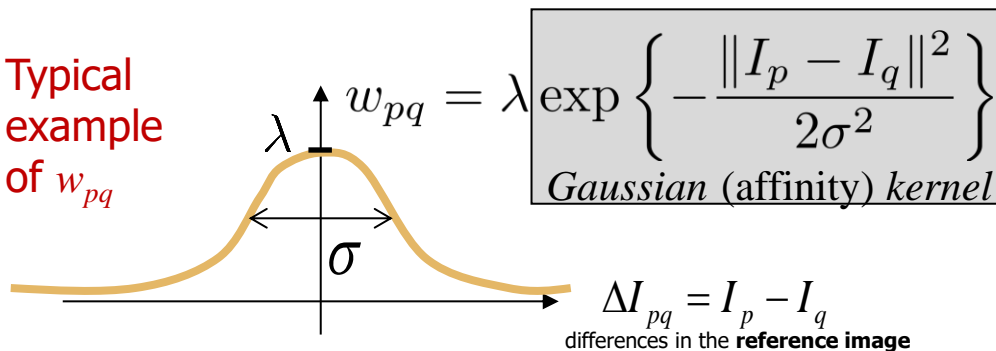
for horizontal edges/neighbors $\{p, q\}$



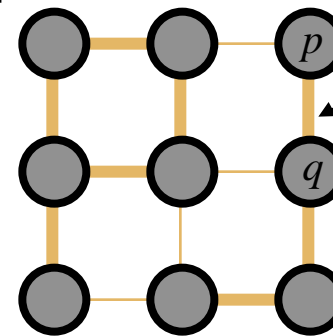
for vertical edges/neighbors $\{p, q\}$



Typical
example
of w_{pq}



Similar “edge aligning” affinity kernels are common in low-level segmentation (see topic 9). “Deep features” f_p can replace I_p (topic 12).



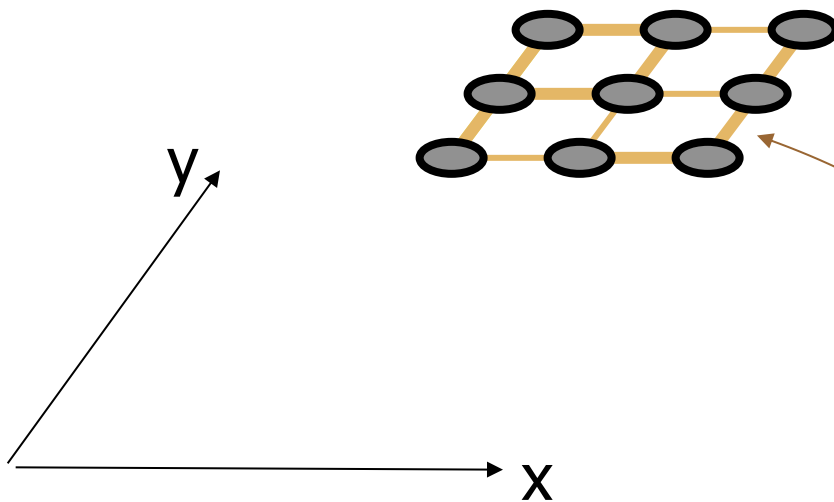
w_{pq} - weights of neighborhood edges
(e.g. may be assigned according to local intensity contrast in the reference image)

Motivation: such *static cues* (in ref. image) help to **align depth boundaries to high contrast edges** since the loss function gets lower when disparity changes $|d_p - d_q|$ happen near edges where $|\Delta I_{pq}| > \sigma$.

NOTE: **parameter σ** is important, it works as (soft) **edge detection threshold**.

Multi-scan-line stereo with s - t graph cuts [Roy&Cox'98, Ishikawa'98]

optional slides
(grad students)



$$E(\mathbf{d}) = \sum_{p \in G} D_p(d_p) + \sum_{\{pq\} \in N} w_{pq} |d_p - d_q|$$

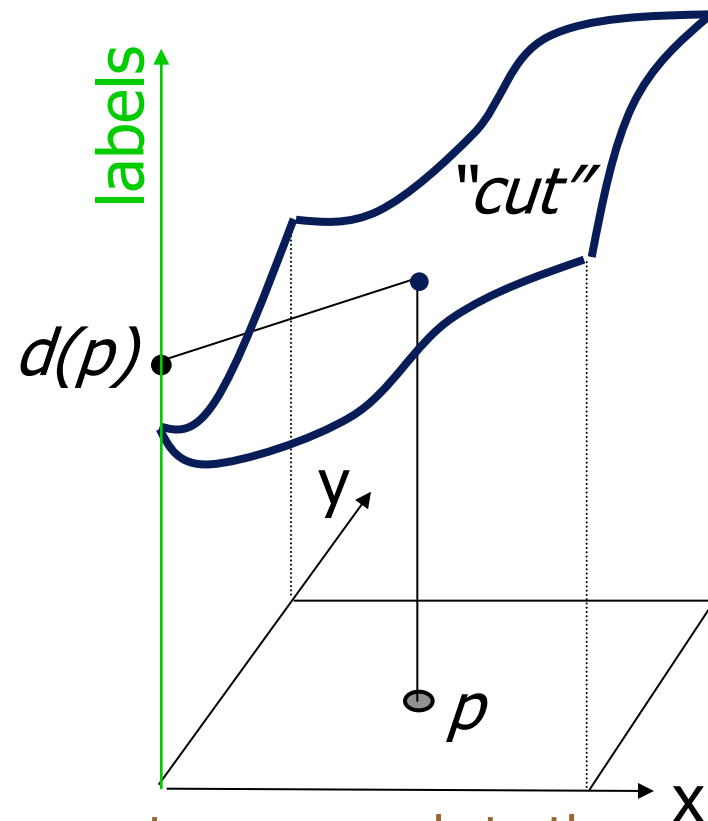
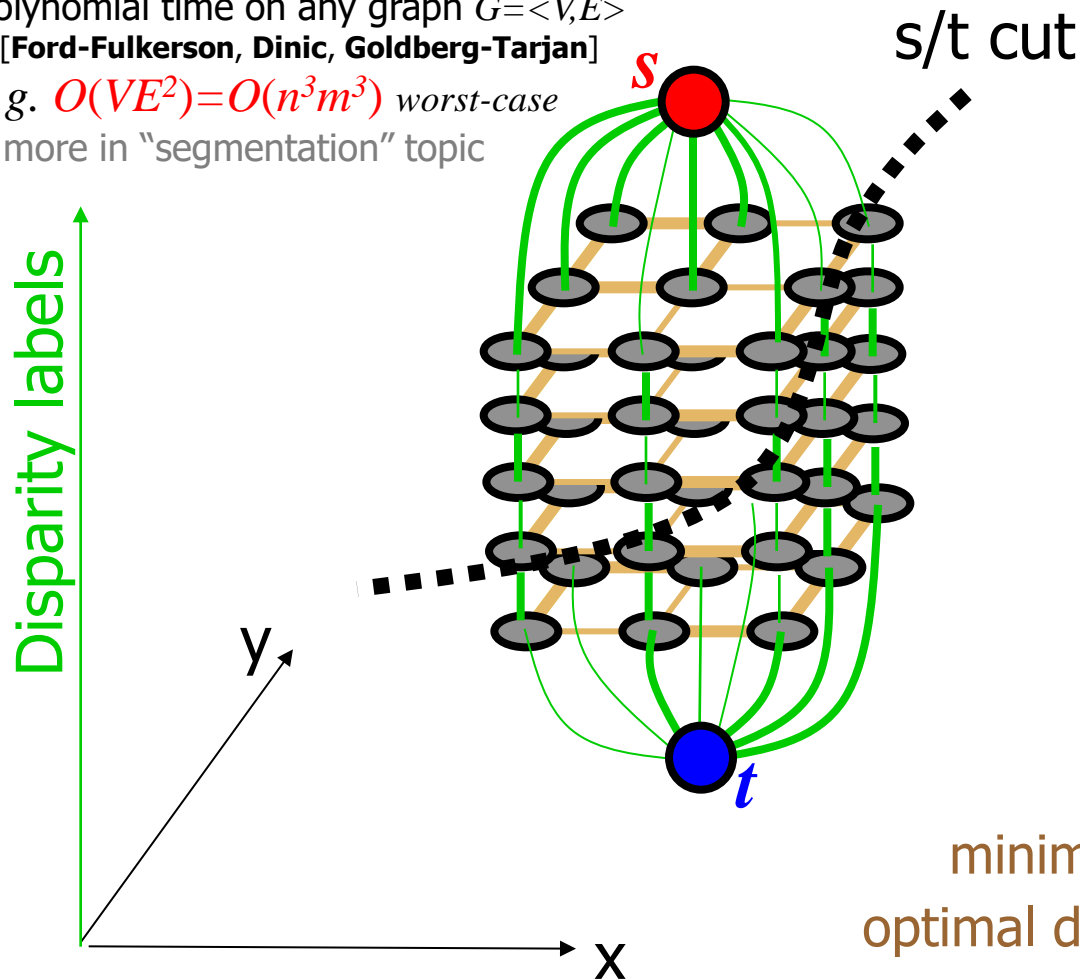
Multi-scan-line stereo with s - t graph cuts [Roy&Cox '98, Ishikawa '98]

optional slides
(grad students)

Minimum s / t cuts can be found in low-order polynomial time on any graph $G = \langle V, E \rangle$

[Ford-Fulkerson, Dinic, Goldberg-Tarjan]

e.g. $O(VE^2) = O(n^3m^3)$ worst-case
more in "segmentation" topic



minimum cut corresponds to the
optimal disparity map $\mathbf{d} = \{d_p\}$ for loss

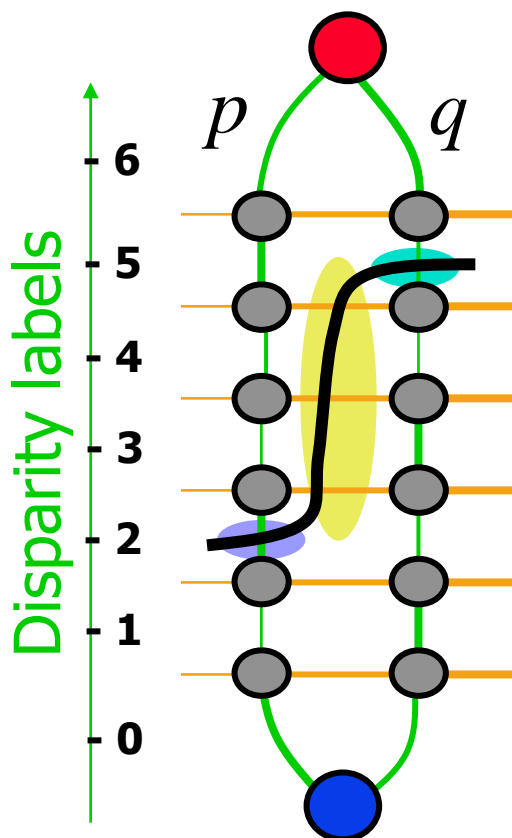
$$E(\mathbf{d}) = \sum_{p \in G} D_p(d_p) + \sum_{\{pq\} \in N} w_{pq} |d_p - d_q|$$

Ishikawa&Geiger 98

What loss function do we minimize this way?

optional slides
(grad students)

Concentrate on one pair of neighboring pixels $\{p, q\} \in N$



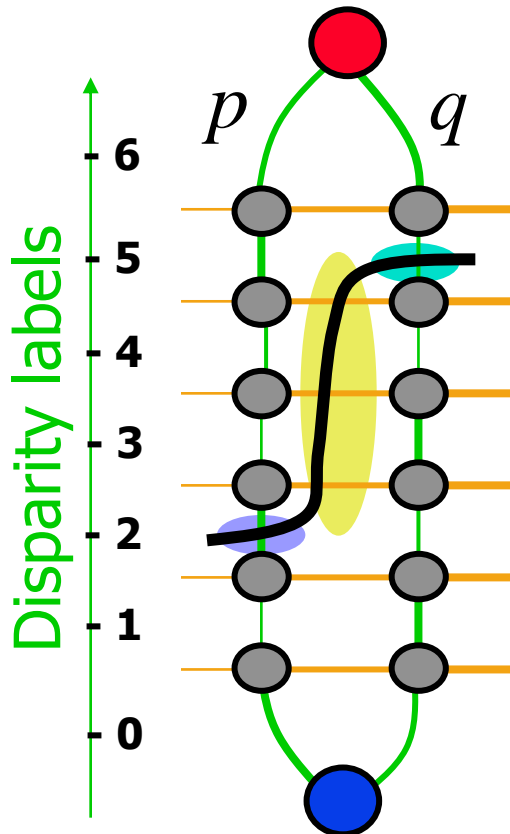
$$E(d_p, d_q) = \begin{array}{l} \text{cost of vertical edges} \\ D_p(2) + D_q(5) + \dots \\ + \\ w_{pq} \cdot |3| + \dots \\ \text{cost of horizontal edges} \end{array}$$

assuming cut has no folds (optional slides later shows how to make sure)

optional slides
(grad students)

What loss function do we minimize this way?

Concentrate on one pair of neighboring pixels $\{p, q\} \in N$

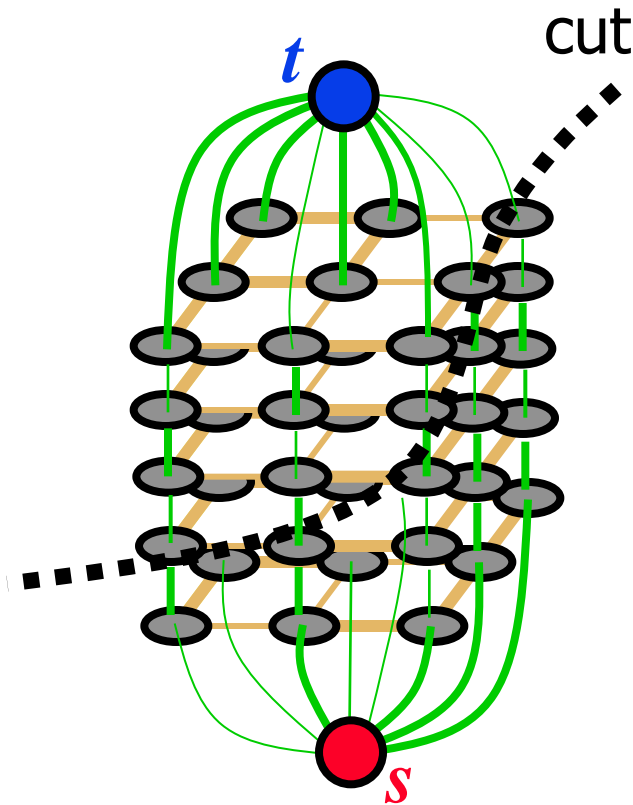


$$E(d_p, d_q) = \begin{array}{l} \text{cost of vertical edges} \\ D_p(d_p) + D_q(d_q) + \dots \\ \\ + \\ w_{pq} \cdot |d_p - d_q| + \dots \\ \text{cost of horizontal edges} \end{array}$$

What loss function do we minimize this way?

optional slides
(grad students)

The combined loss over the entire grid G is



(**photo consistency**, e.g. SSD)
cost of vertical edges

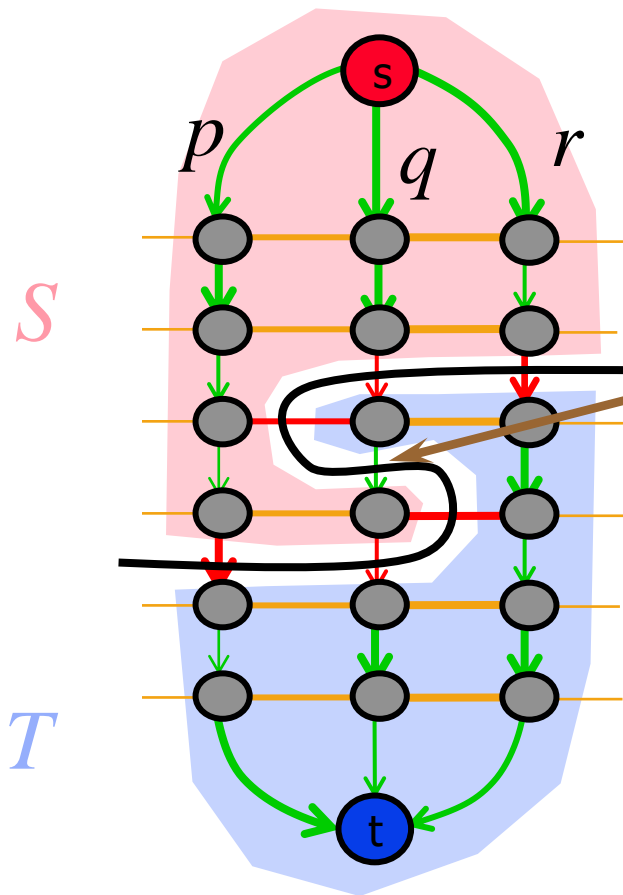
$$E(\mathbf{d}) = \sum_{p \in G} D_p(d_p)$$

$$+ \sum_{\{p, q\} \in N} w_{pq} \cdot |d_p - d_q|$$

cost of horizontal edges
(**spatial consistency**)

How to avoid folding?

consider three pixels $\{p, q, r\}$



introduce directed *t-links*

NOTE: this directed *t-link* is not “severed”

WHY?

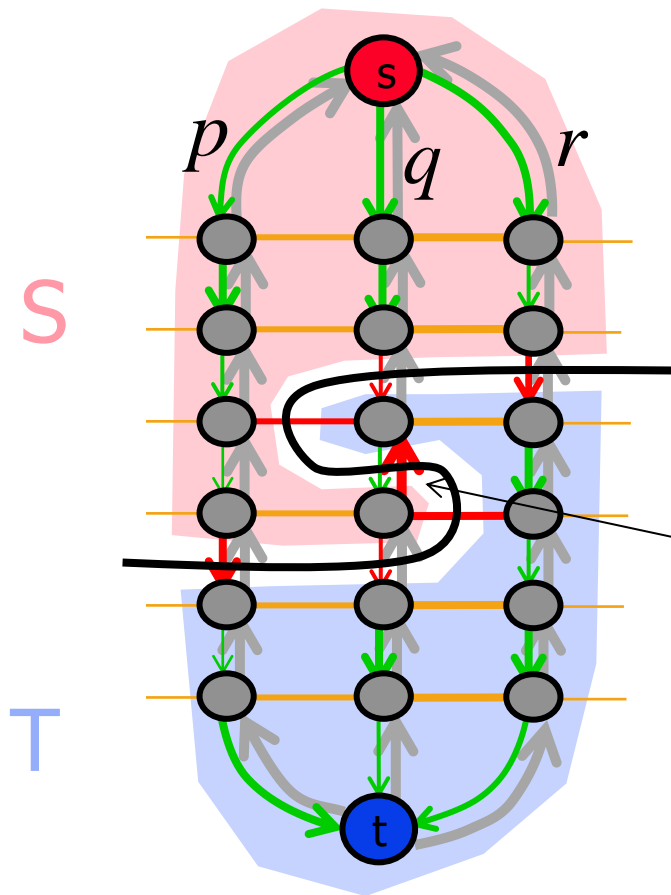
Formally, s/t cut is a partitioning of graph nodes

$$C = \{S, T\} \quad \text{and its cost is} \quad \|C\| = \sum_{\substack{(pq) \in E \\ p \in S \\ q \in T}} c_{pq}$$

only edges from *S* to *T* matter

How to avoid folding?

consider three pixels $\{p, q, r\}$



Solution prohibiting **folds**:

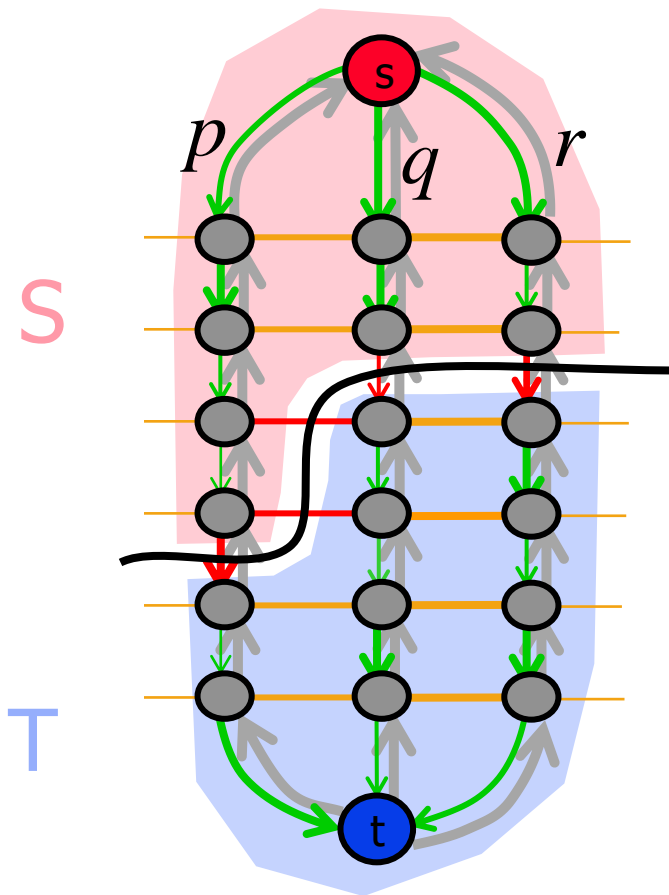
add infinity cost t-links
in the "up" direction



NOTE: **folding cuts** $C = \{S, T\}$
sever at least one of such t-links
making such cuts **infeasible**

How to avoid folding?

consider three pixels $\{p, q, r\}$



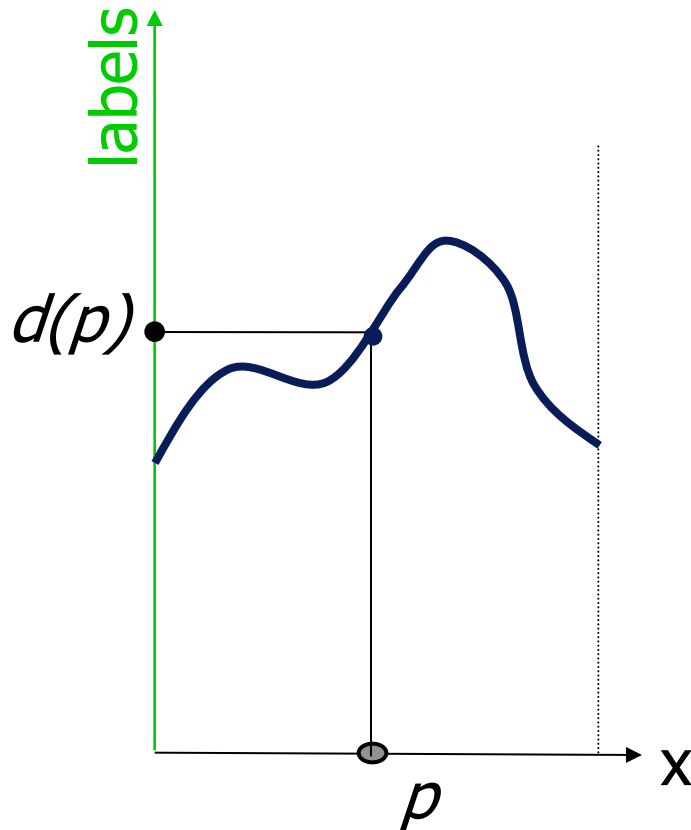
Solution prohibiting **folds**:

add infinity cost t-links
in the “up” direction

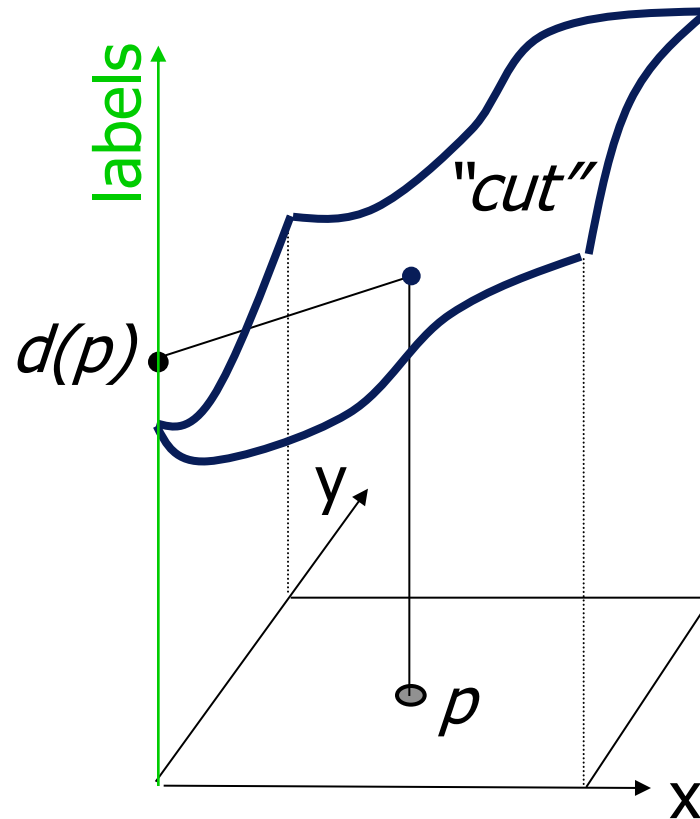


NOTE: **non-folding cuts** $C = \{S, T\}$
do not sever such t-links

Scan-line stereo vs. Multi-scan-line stereo (on whole grid)

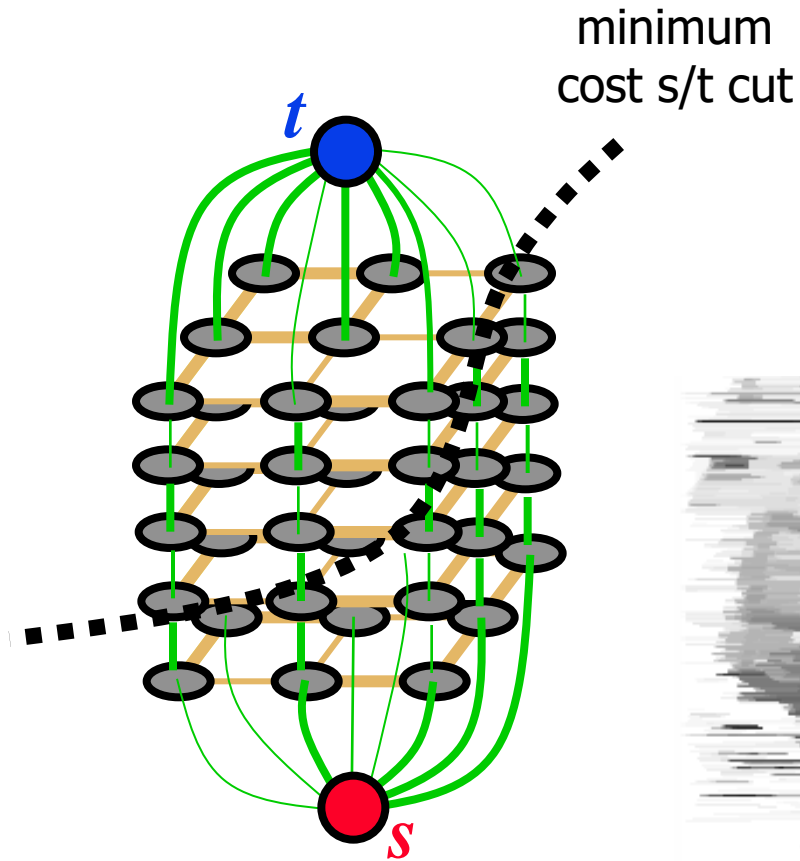


Dynamic Programming
(single scan line optimization)

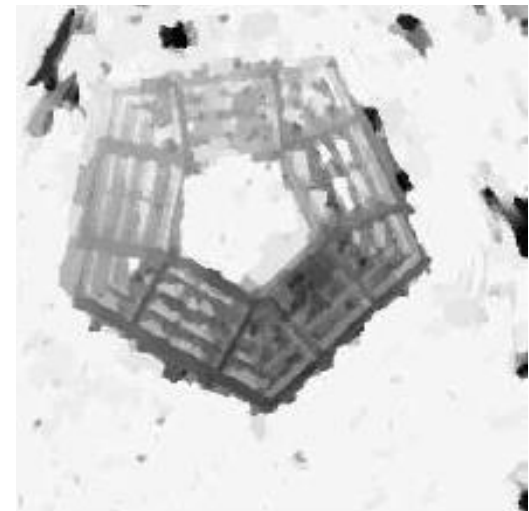


s-t Graph Cuts
(grid optimization)

Some results from Roy&Cox

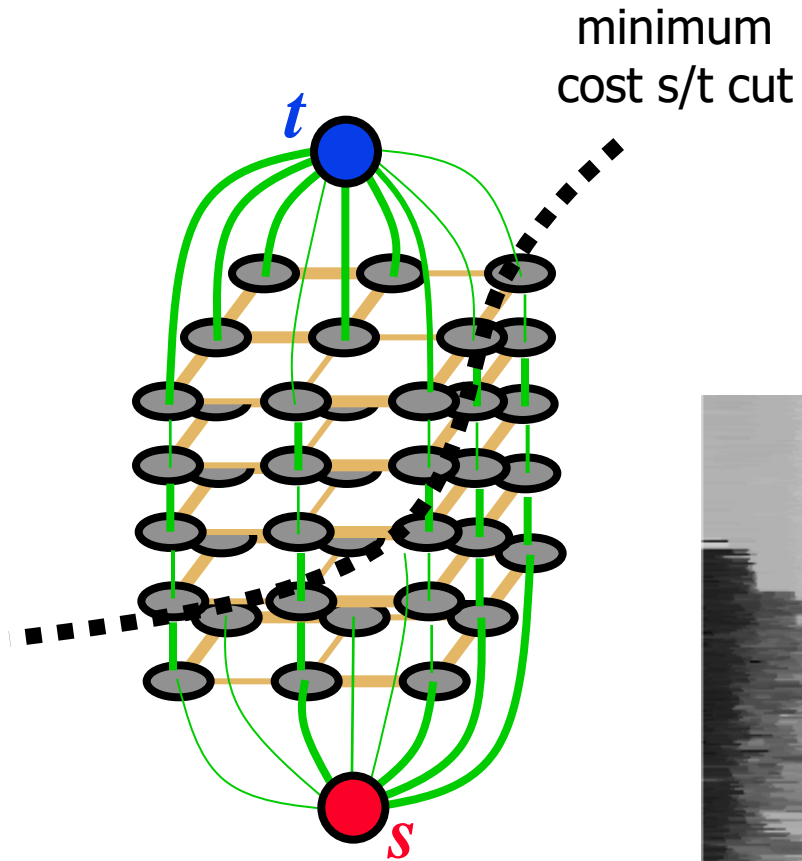


single scan-line stereo
(DP)



multi scan line stereo
(graph cuts)

Some results from Roy&Cox



single scan-line stereo
(DP)



multi scan line stereo
(graph cuts)

Simple Examples: Stereo with only 2 depth layers



binary stereo



essentially,
depth-based binary
segmentation

Simple Examples: Stereo with only 2 depth layers



*background
substitution*



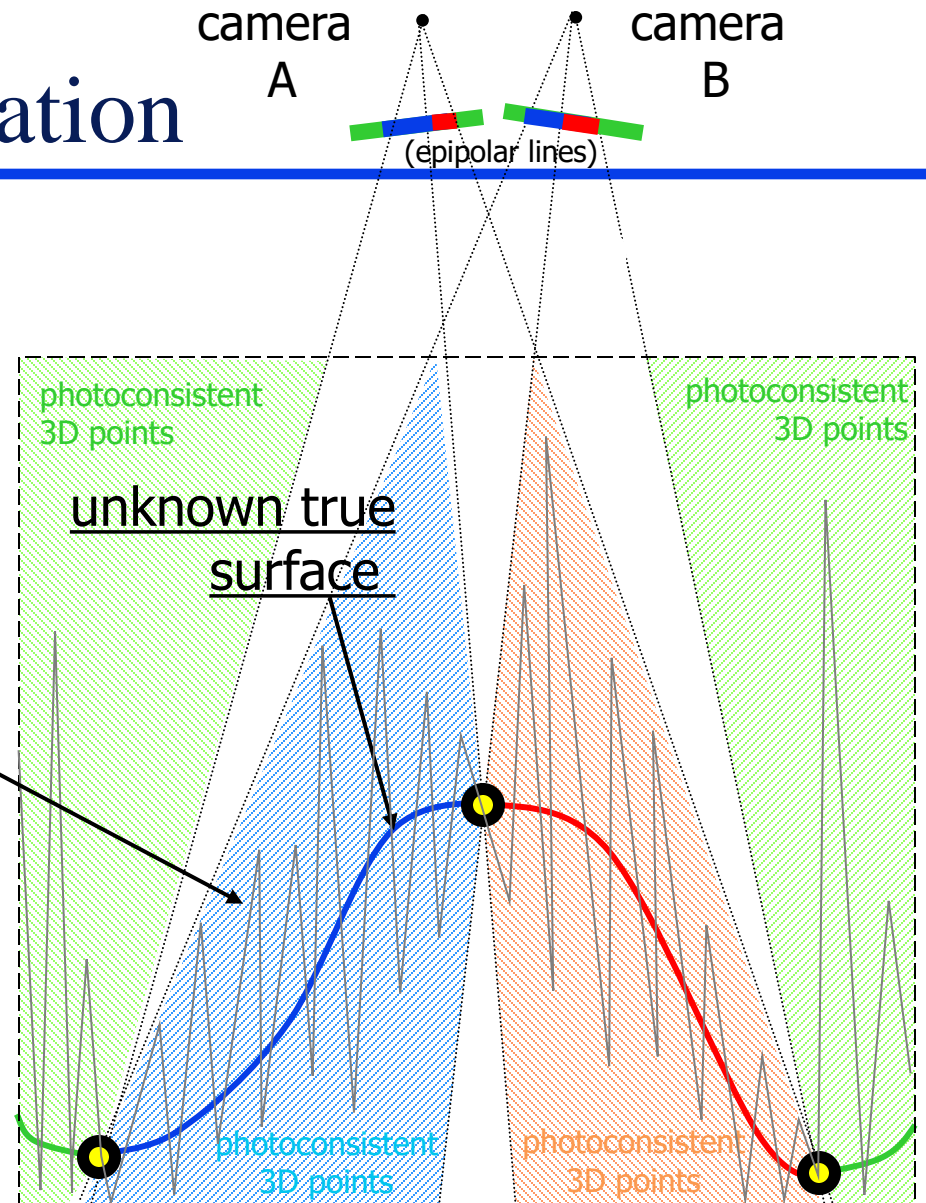
essentially,
depth-based binary
segmentation

Features and Spatial Regularization

$$E(\mathbf{d}) = \sum_{p \in G} |I_p - I'_{p+d_p}|$$

photo-consistency term

photoconsistent depth map



Features and Spatial Regularization


$$E(\mathbf{d}) = \sum_{p \in G} |I_p - I'_{p+d_p}|$$

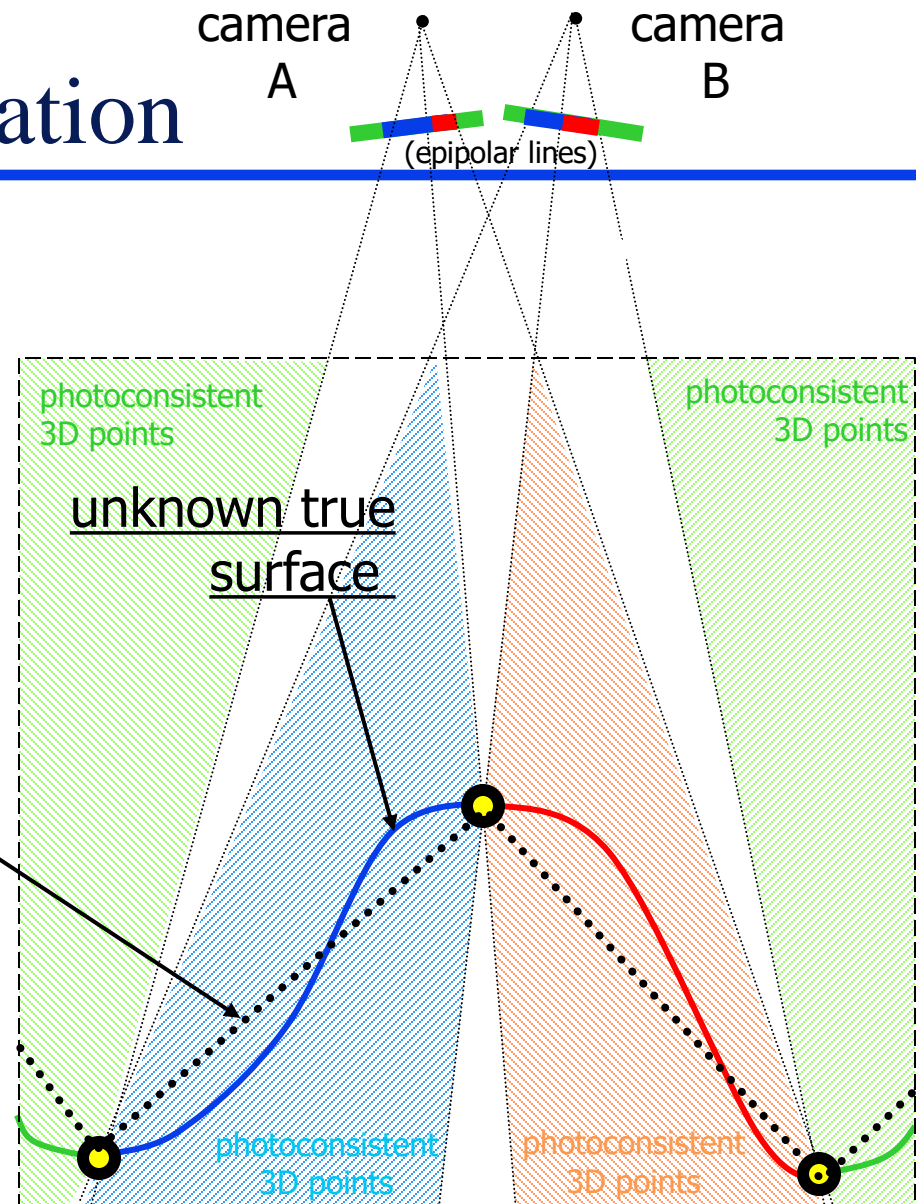
photo-consistency term

$$+ \sum_{pq \in N} w |d_p - d_q|$$

regularization term

regularized depth map

- regularization helps to find **smooth** depth map consistent with points  uniquely matched by photoconsistency
- regularization propagates information from textured regions (features) to ambiguous textureless regions



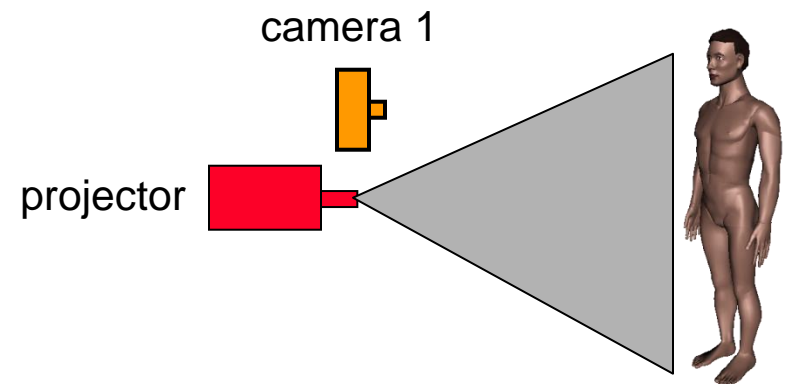
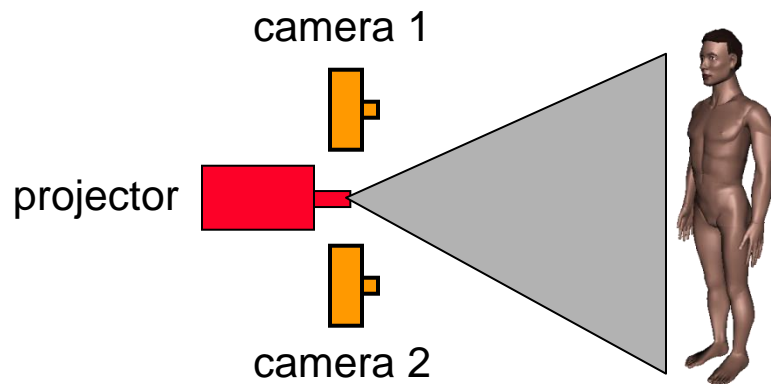
3D volume where surface is being reconstructed
(epipolar plane)

More features/texture always helps!

Active Stereo (with structured light)

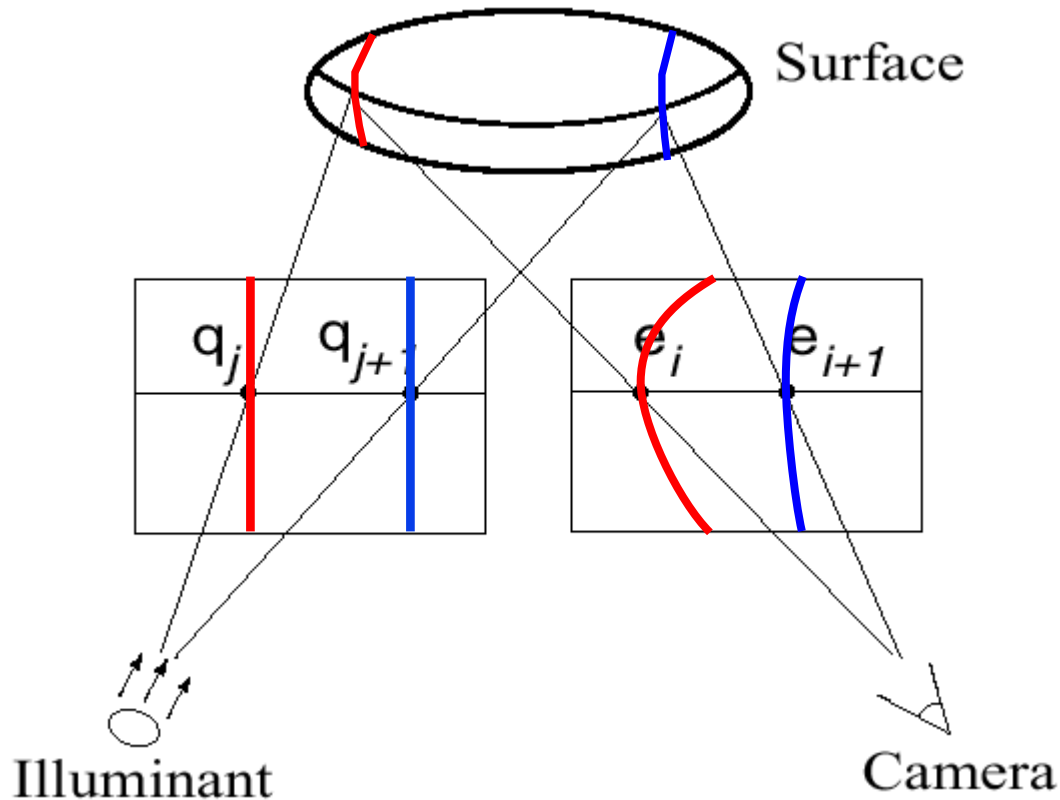


Li Zhang's one-shot stereo

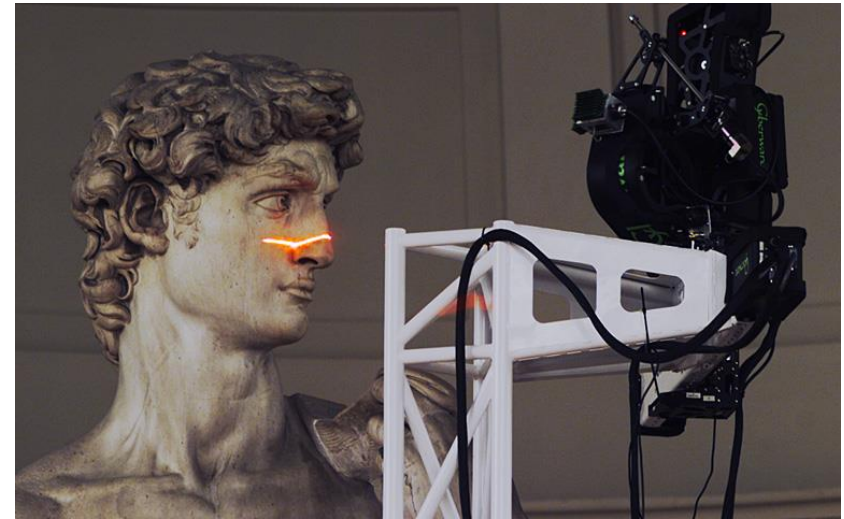
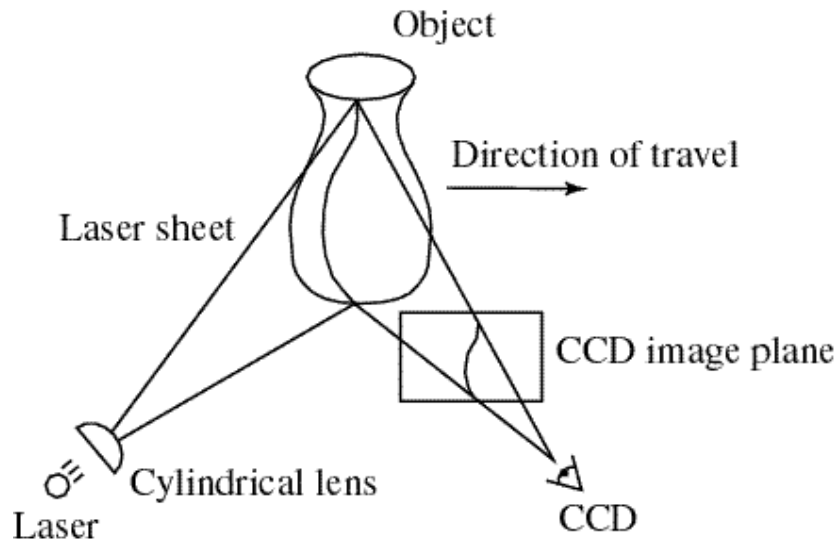


- Project “structured” light patterns onto the object
 - simplifies the correspondence problem

Active Stereo (with structured light)



Laser scanning

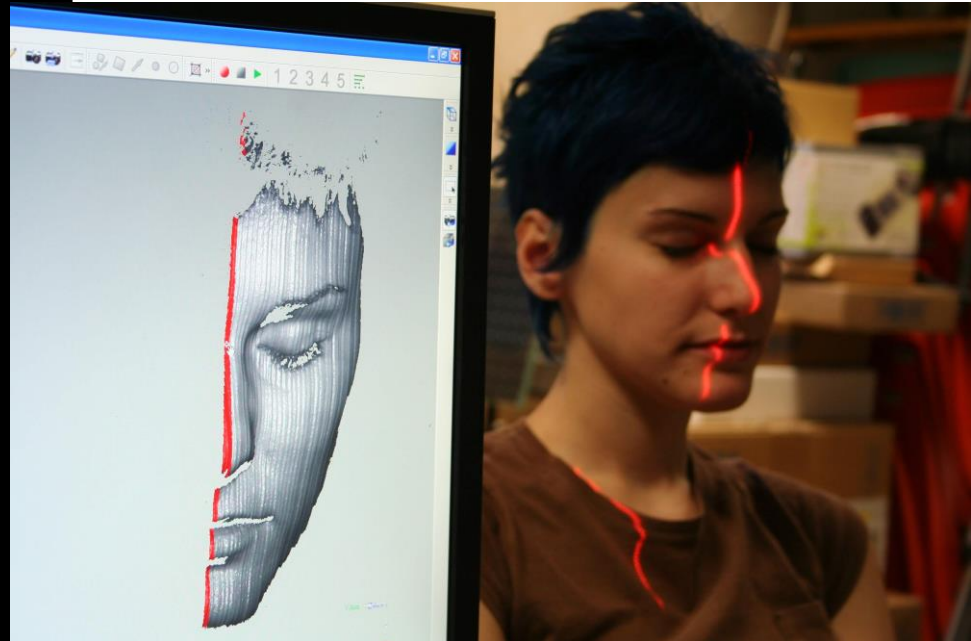
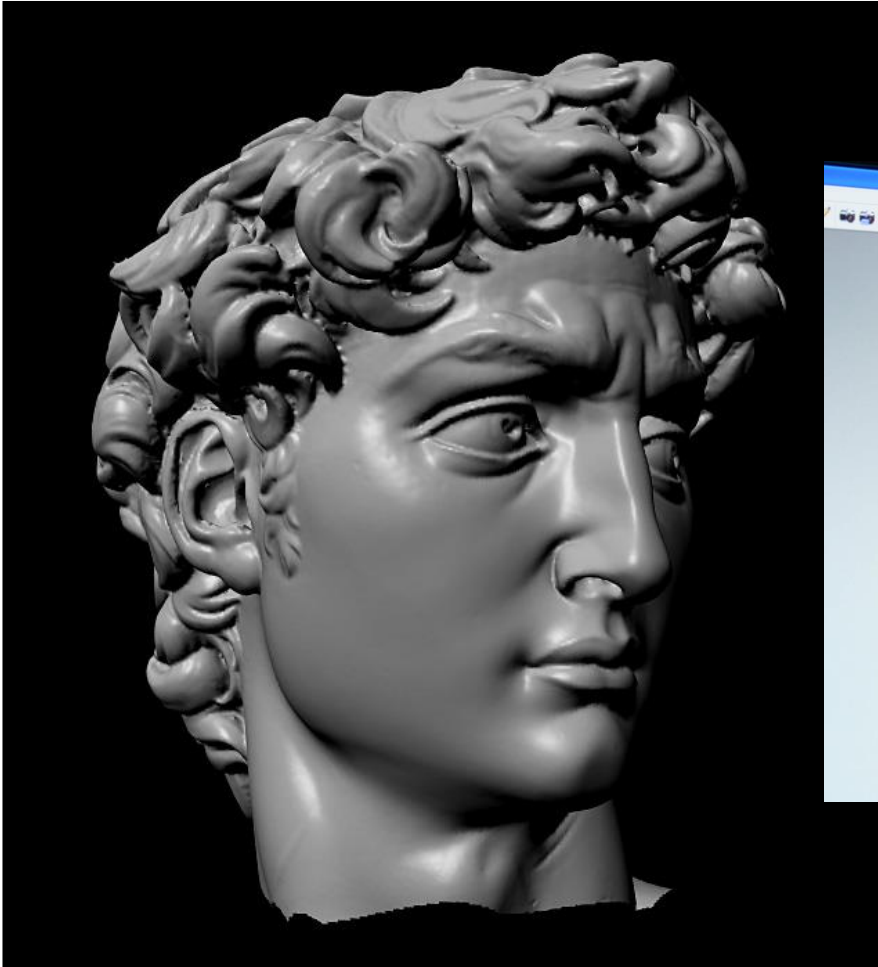


Digital Michelangelo Project [Levoy et al.]
<http://graphics.stanford.edu/projects/mich/>

□ Optical triangulation

- Project a single stripe of laser light
- Scan it across the surface of the object
- This is a very precise version of structured light scanning

Laser scanning

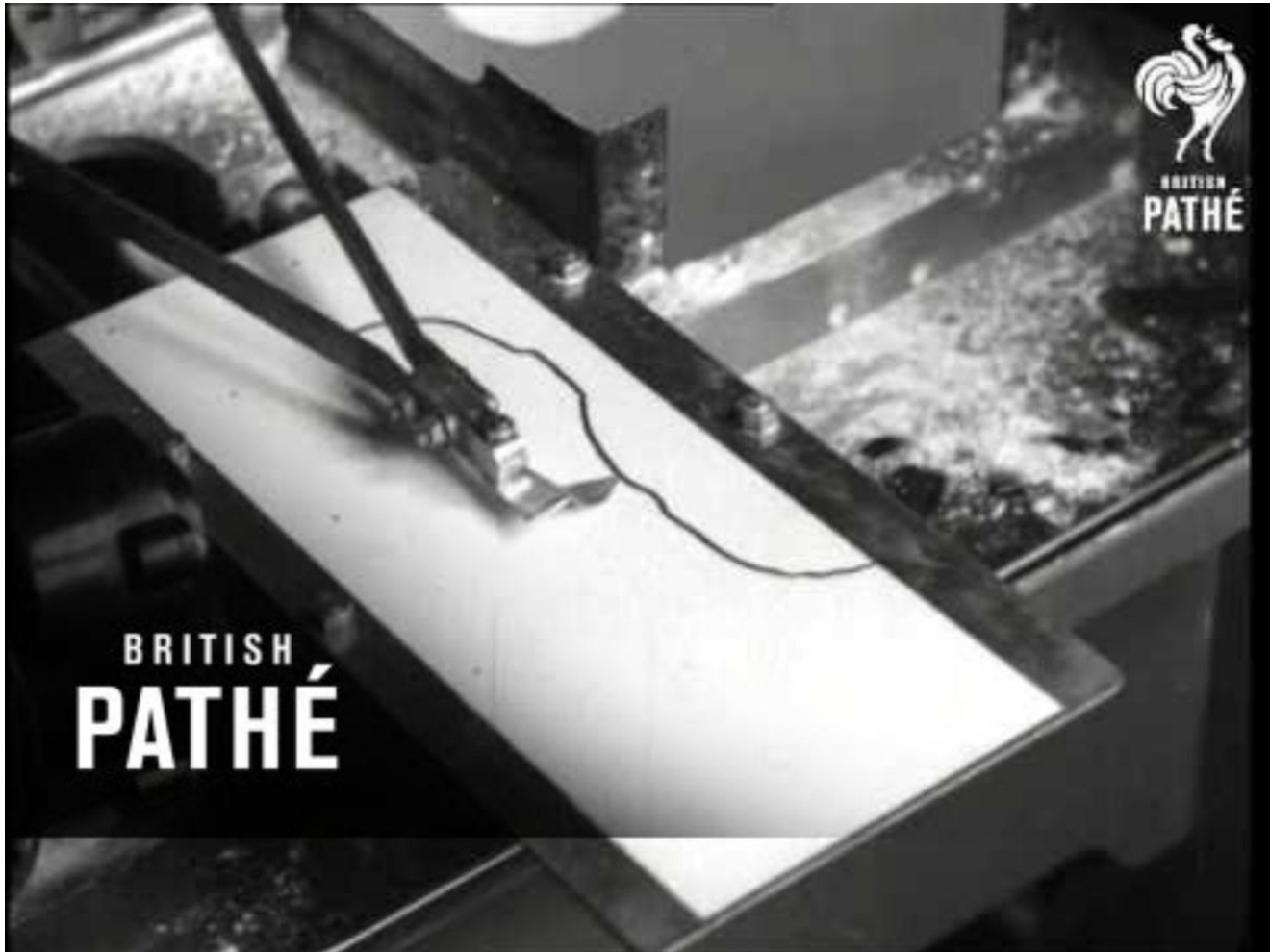


basic laser scanner

Digital Michelangelo Project [Levoy et al.]

<http://graphics.stanford.edu/projects/mich/>

Photo Sculpture (1939)



https://youtu.be/jS_rcwG9mxU?si=JcrzZs2VSZb6yvuS

3D scanning

- Stereo
 - uses photo sensors
 - requires ambient light, does not work in dark environment
 - ambiguous in textureless regions, noisy on specular surfaces
- Active Stereo
 - uses photo sensors and active light sources (e.g. laser)
 - problems with specular or non-reflective surfaces
 - problems with bright ambient light
- Lidar
 - uses *time-of-flight* sensors and active light (laser)
 - good range, no baseline required
 - problems with specular or non-reflective surfaces
 - problems with bright ambient light
 - relatively sparse output (cloud of points)



To produce **dense scene reconstruction**, all methods should address noise and ambiguities by fitting various dense surface models, typically using **surface regularization techniques**

Disparity map $d(p)$ is an example of (regularised) surface model, more in Topic 9B

further considerations:

Robust error/penalty functions

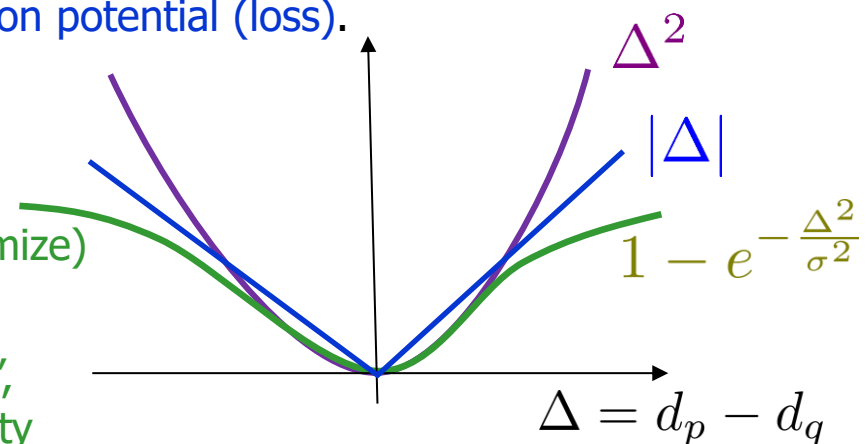
$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\text{photo-consistency}} + \sum_{\{p, q\} \in N} \underbrace{V(d_p, d_q)}_{\text{spatial coherence}}$$

$\underbrace{\quad}_{\text{photo-consistency}} \quad \underbrace{\quad}_{\text{spatial coherence}}$

The last term is an example of **convex** regularization potential (loss).

- easier to optimize, but
- tend to over-smooth

practically preferred
robust regularization
(non convex – harder to optimize)



Note: once deviation/error Δ is "large enough", there is no reason to keep increasing the penalty

further considerations:

Robust error/penalty functions

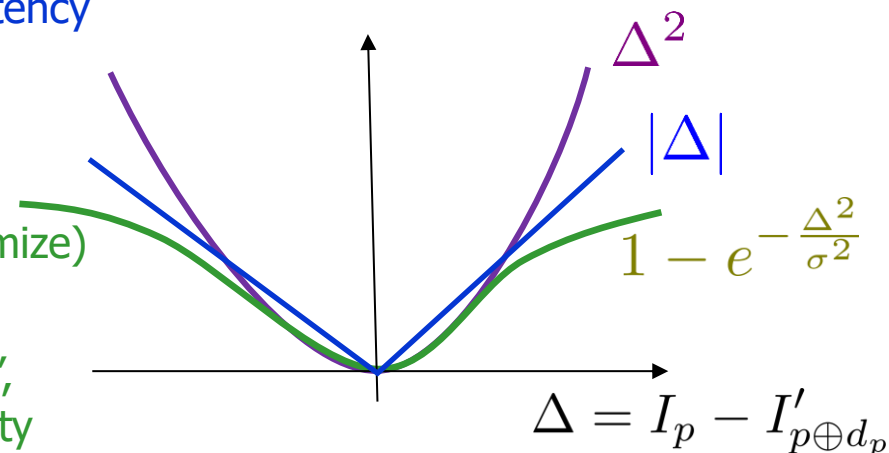
$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\text{photo-consistency}} + \sum_{\{p, q\} \in N} \underbrace{V(d_p, d_q)}_{\text{spatial coherence}}$$

$|I_p - I'_{p \oplus d_p}|$
 $w_{pq} \cdot |d_p - d_q|$

Similarly, robust losses are needed for **photo-consistency** to handle occlusions & “specularities”

practically preferred
robust regularization
(non convex – harder to optimize)

Note: once deviation/error Δ is “large enough”, there is no reason to keep increasing the penalty



further considerations:

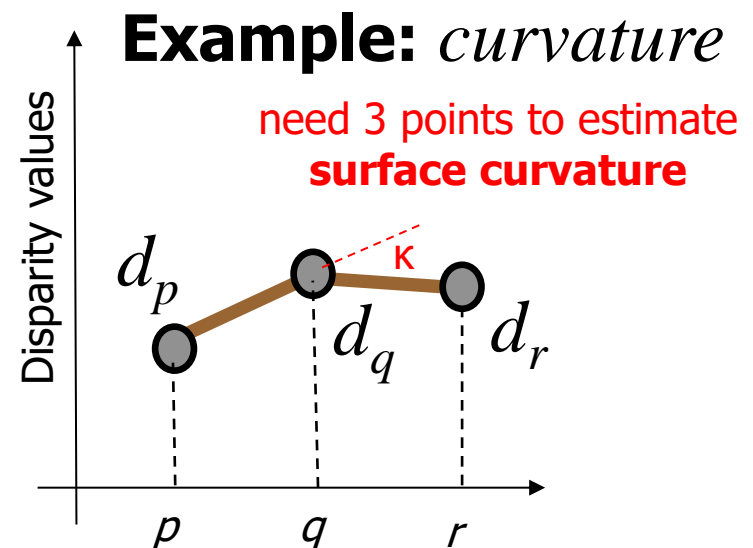
Higher-order regularization

$$E(\mathbf{d}) = \sum_{p \in G} \underbrace{D_p(d_p)}_{\text{photo-consistency } |I_p - I'_{p \oplus d_p}|} + \sum_{\{p, q\} \in N} \cancel{V(d_p, d_q)} V(d_p, d_q, d_r)$$

$p, q, r \in N$ higher-order "coherence"

Many state-of-the-art methods use higher-order regularizers

Q: why penalizing depth curvature instead of depth change?

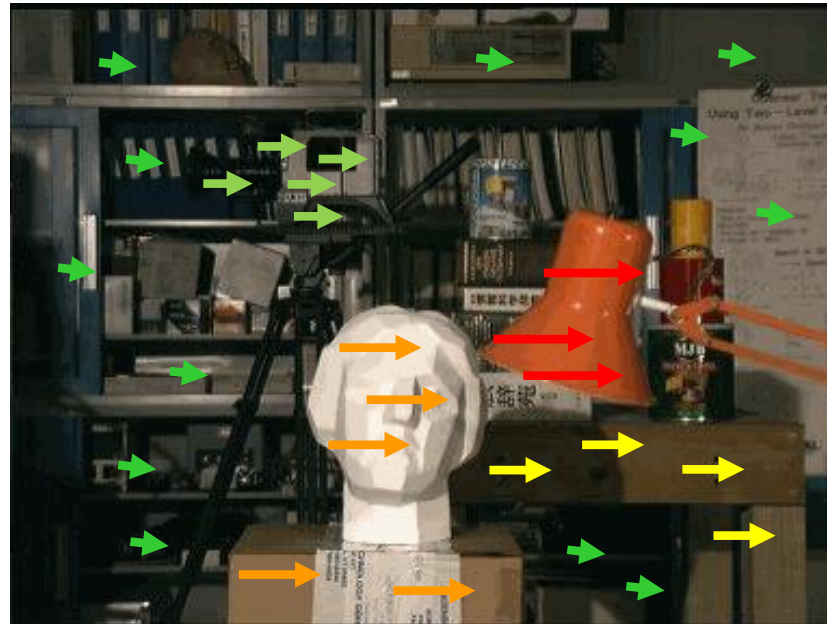


From 1D correspondence (stereo) to 2D correspondence problems (motion)

1D shifts along **epipolar lines**.

Assumption for stereo:

only camera moves,
3D scene is stationary



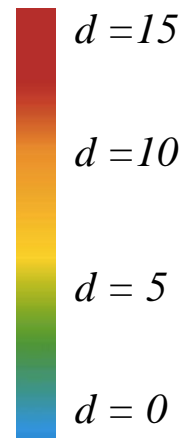
vector field (motion) with a priori known direction

From 1D correspondence (stereo) to 2D correspondence problems (motion)

1D shifts along **epipolar lines**.

Assumption for stereo:

only camera moves,
3D scene is stationary



vector field (motion) with a priori *known direction*

⇒ We estimate only *magnitude* represented by a **scalar field** (disparity map)

From 1D correspondence (stereo) to 2D correspondence problems (motion)

In general, correspondences between two images
may not be described by global models (like *homography*) or
by **shifts along known epipolar lines**.

if 3D scene
is NOT stationary
motion is
vector field
with **arbitrary**
directions
(no epipolar line constraints)



From 1D correspondence (stereo) to 2D correspondence problems (motion)

In general, correspondences between two images **may not be** described by global models (like *homography*) or by **shifts along known epipolar lines**.

For (non-rigid) motion the correspondences between two video frames are described by a general ***optical flow***

if 3D scene
is NOT stationary
motion is
vector field
with **arbitrary**
directions
(no epipolar line constraints)



From 1D correspondence (stereo) to 2D correspondence problems (motion)

$$E(\mathbf{v}) = \sum_{p \in G} \underbrace{D_p(v_p)}_{\text{color-consistency}} + \sum_{\{p, q\} \in N} \underbrace{V(v_p, v_q)}_{\text{regularity}}$$

$$(I_p^t - I_{p+v_p}^{t+1})^2 \quad w \cdot \|v_p - v_q\|^2$$

Horn-Schunck 1981

optical flow regularization

- 2nd order optimization

(pseudo Newton)

- Rox/Cox/Ishikawa's method only works for scalar-valued variables

if 3D scene
is NOT stationary

motion is

vector field

with **arbitrary directions**

(no epipolar line constraints)



SOCIETY OF ROBOTS

optical flow

$$\mathbf{V} = \{v_p\}$$

more difficult problem

need 2D shift vectors v_p

(no epipolar line constraint)

over-smoothed vector field

(robust regularization losses can preserve sharp changes in motion between objects)

From 1D correspondence (stereo) to 2D correspondence problems (motion)

State-of-the-art methods **segment**
independently moving objects

We will discuss
segmentation
problem
next

if 3D scene
is NOT stationary
motion is
vector field
with **arbitrary**
directions
(no epipolar line constraints)



SOCIETY OF ROBOTS

optical flow

$$\mathbf{V} = \{\mathbf{v}_p\}$$

more difficult problem
need 2D shift vectors \mathbf{v}_p
(no epipolar line constraint)

over-smoothed vector field
(robust regularization losses
can preserve sharp changes
in motion between objects)