

# Brief Introduction to Geometry and Vision

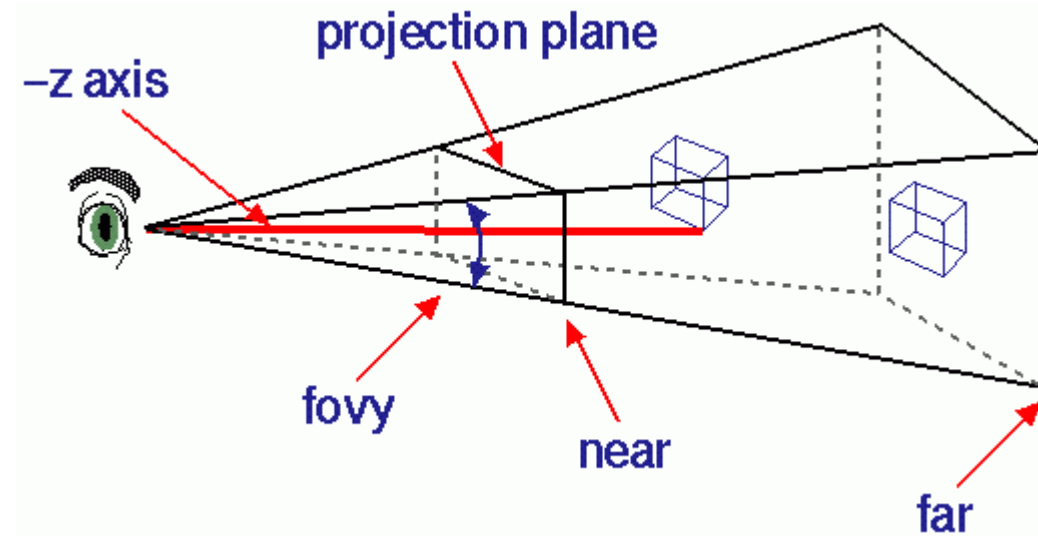
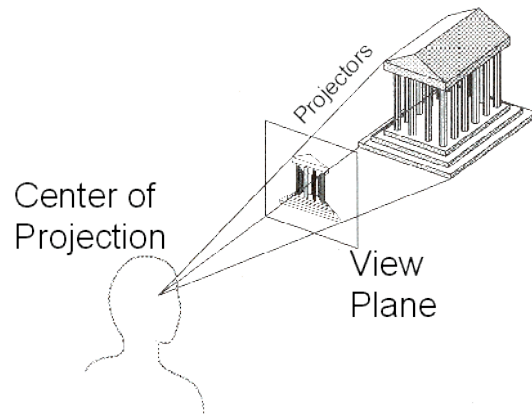
A.L. Yuille (UCLA)

# Plan of Talk

- Four Topics:
- (I) Basic Projection. Perspective. Vanishing Points.
- (II) Camera Calibration. Stereopsis. Essential Matrix. Fundamental Matrix.
- (III) Structure from Motion. Rigid. Extension to Non-Rigid.
- (IV) Geometric Priors. Manhattan World.

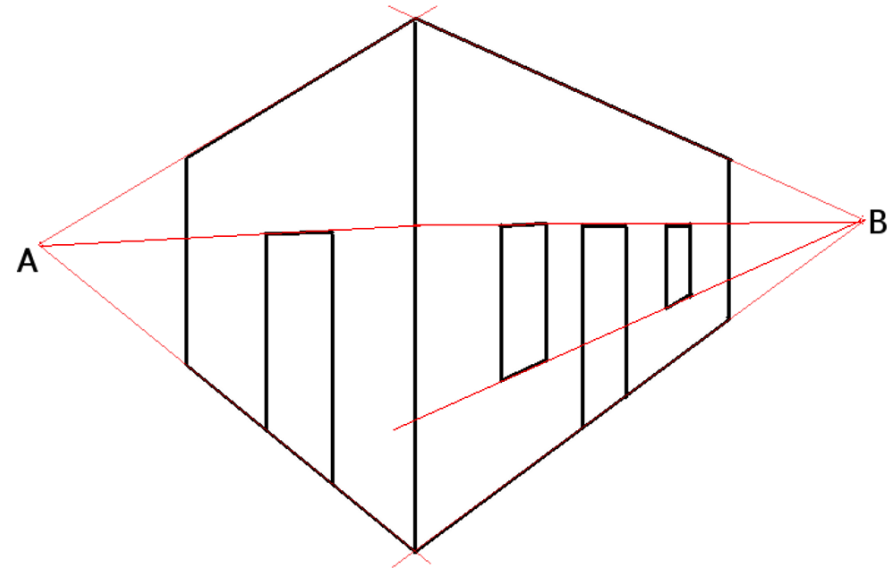
# Geometry of Projection.

- Most analysis is based on the Pinhole camera model.
- Real cameras have lens (W. Freeman's lectures). See Szeliski's book for corrections to the pinhole camera model.



# Properties of Perspective Projection

- Straight lines project to straight lines.
- Parallel lines in space project to lines which converge at a vanishing point.



# Perspective 1

## 1) Perspective

$$u = f \frac{X}{Z}$$

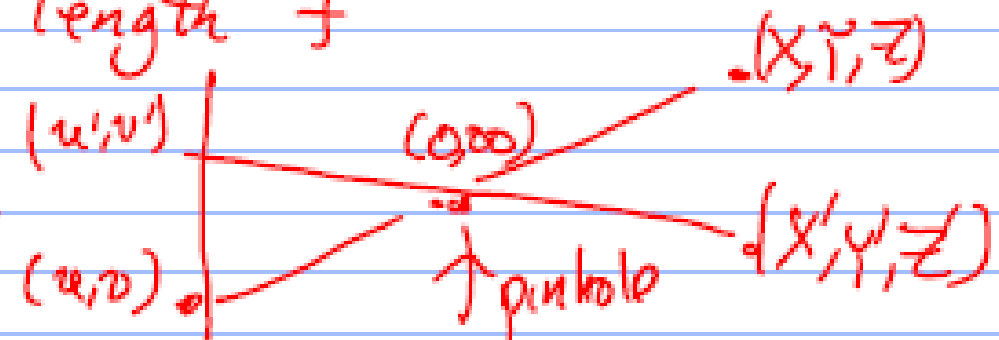
$$v = f \frac{Y}{Z}$$

3D position  $(X, Y, Z)$

image coordinates  $(u, v)$

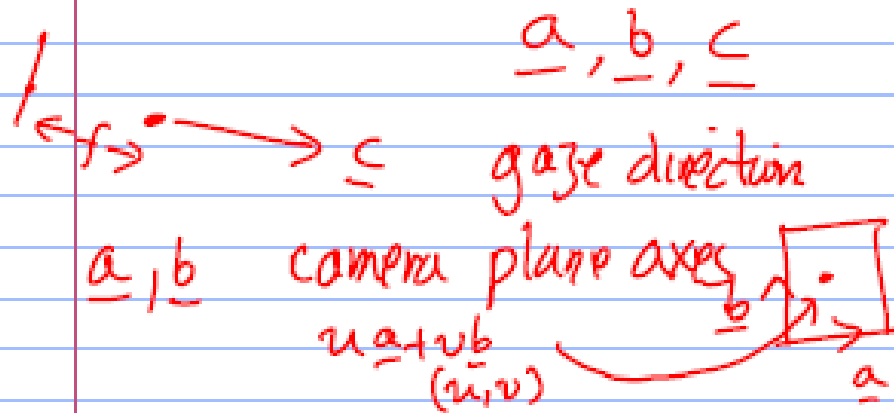
focal length  $f$

- Coordinate system based on camera.
- Origin  $(0, 0, 0)$  at pinhole (lens)
- Camera is pointing in the  $Z$  direction



# Perspective 2

## 2) Camera Parameters

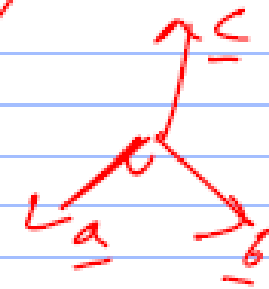


right-angle triad.

$$\underline{a} \cdot \underline{a} = \underline{b} \cdot \underline{b} = \underline{c} \cdot \underline{c} = 1$$

$$\underline{a} \cdot \underline{b} = \underline{b} \cdot \underline{c} = \underline{c} \cdot \underline{a} = 0$$

dot product.



$$\underline{a} \wedge \underline{b} = \underline{c}, \quad \underline{b} \wedge \underline{c} = \underline{a}, \quad \underline{c} \wedge \underline{a} = \underline{b}$$

$\wedge$  cross product

Camera Parameters:  $\underline{a}, \underline{b}, \underline{c}, f, \underline{o}$

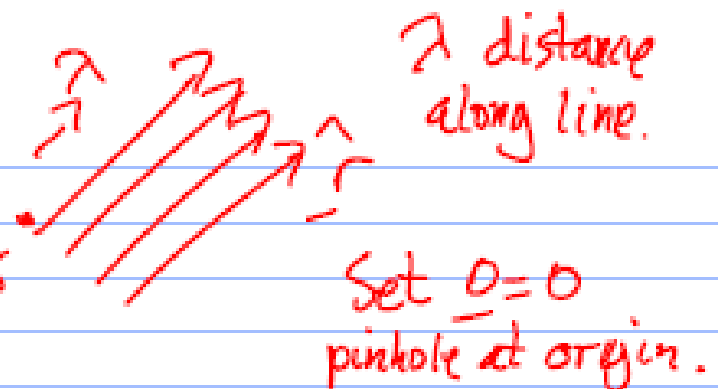
$$u = -f \frac{(\underline{r} - \underline{o}) \cdot \underline{a}}{\underline{r} \cdot \underline{c}}, \quad v = -f \frac{(\underline{r} - \underline{o}) \cdot \underline{b}}{\underline{r} \cdot \underline{c}}$$

More Complex Camera Models — e.g.  $u_0, v_0$  center of image

# Vanishing Points 1.

(3) (i)

Vanishing Points Parallel Lines in Space



$\underline{r}_0 + \lambda \underline{\hat{r}}$   $\lambda$  - length

$\underline{\hat{r}}$  - direction of lines

projects to line in image plane

$$u(\lambda) = \frac{-f(\underline{r}_0 + \lambda \underline{\hat{r}}) \cdot \underline{a}}{(\underline{r}_0 + \lambda \underline{\hat{r}}) \cdot \underline{c}}, \quad v(\lambda) = \frac{-f(\underline{r}_0 + \lambda \underline{\hat{r}}) \cdot \underline{b}}{(\underline{r}_0 + \lambda \underline{\hat{r}}) \cdot \underline{c}}$$

u.s.g. set  $\underline{r}_0 \cdot \underline{c} = 0$ .

$$\underline{r}_0 \rightarrow \underline{r}_0 - \frac{(\underline{r}_0 \cdot \underline{c})}{(\underline{\hat{r}} \cdot \underline{c})} \underline{\hat{r}}$$

$\underline{\hat{r}} \cdot \underline{c} = 0$ , only for lines perp. to direction of gaze  $\underline{c}$ .

$$u(\lambda) = -f \frac{\underline{\hat{r}} \cdot \underline{a}}{\underline{\hat{r}} \cdot \underline{c}} - f \frac{\underline{r}_0 \cdot \underline{a}}{\underline{\hat{r}} \cdot \underline{c}} \frac{1}{\lambda}$$

$$v(\lambda) = -f \frac{\underline{\hat{r}} \cdot \underline{b}}{\underline{\hat{r}} \cdot \underline{c}} - f \frac{\underline{r}_0 \cdot \underline{b}}{\underline{\hat{r}} \cdot \underline{c}} \frac{1}{\lambda}$$

$\lambda$  distance along line in 3D.  $V_\lambda$  is inverse distance

# Vanishing points 2.

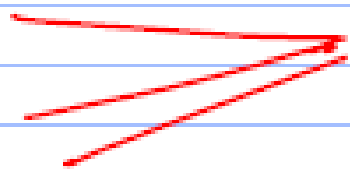
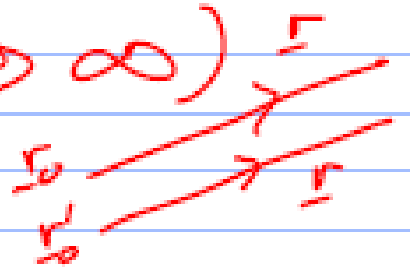
## (4) Vanishing Point (ii)

As  $\frac{1}{\lambda} \rightarrow 0$  (distance in space  $\rightarrow \infty$ )

$$u(\lambda) \rightarrow -f \frac{\hat{r} \cdot \underline{a}}{\hat{r} \cdot \underline{c}}$$

$$v(\lambda) \rightarrow -f \frac{\hat{r} \cdot \underline{b}}{\hat{r} \cdot \underline{c}}$$

independent of  $\underline{r}_0$   
 So all lines in direction  $\hat{r}$  tend to the same vanishing point.



vanishing points may be at infinity.

e.g.



$$\hat{r} \cdot \underline{c} = 0$$

How many vanishing points?  
 see later.



# Linear approximations: Weak Orthographic

- Perspective projection can often be approximated by scaled orthographic projection (e.g., if  $Z$  is constant).
- This is a linear operation.
- Parallel lines project to parallel lines (vanishing points at infinity).
- This is often a good approximation which is easy to use.
- Maths of weak orthographic projection.

# Linear Projection 1

## (5) Linear Projection Models

Can often approximate perspective projection by linear models - simplifies analysis.  
(eg. if relative depth change is small)

$$(u, v) = \begin{pmatrix} K_1 & K_2 & K_3 \\ H_1 & H_2 & H_3 \end{pmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix}$$

affine cameras,  
many variants.

Weak orthographic

$$\underline{K} \cdot \underline{H} = 0 \quad |\underline{K}| = |\underline{H}|$$

$$\underline{K} = K_1, K_2, K_3 \\ \underline{H} = H_1, H_2, H_3$$

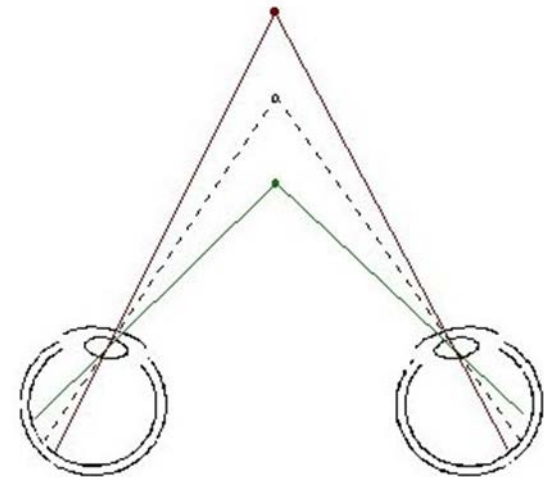
Simplified case

$$u = sX \\ v = sY$$

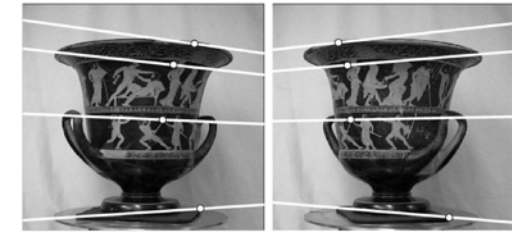
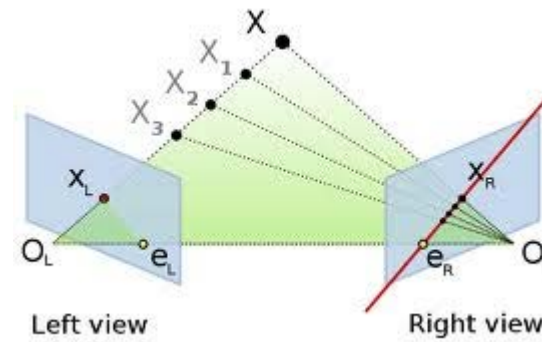
s - scale factor.

# Two Cameras. Binocular Stereo.

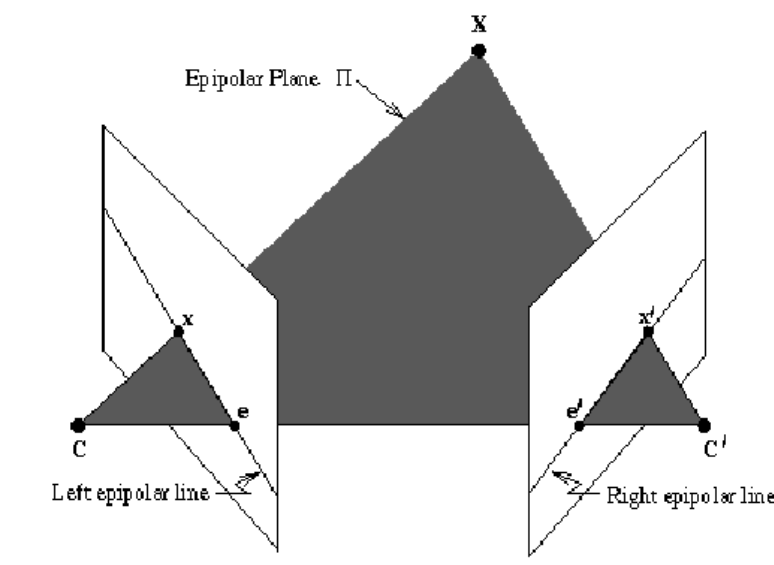
- Binocular stereo.
- Estimate depth from two eyes/cameras by triangulation.
- Requires solving the correspondence problem between points in the two images.
- Correspondence problem is helped by the epipolar line constraint.
- Camera calibration needed.



# Epipolar Lines:



- Points on one epipolar line can only be matched to corresponding epipolar line.
- Epipolar lines depend on the camera parameters.
- If both cameras are parallel, then epipolar lines are horizontal.
- Geometric demonstration of epipolar line constraint.



# Stereo Algorithms can exploit epipolar line constraints

Simplest model: estimate the disparity  $d$  at each point (convert to depth by geometry).

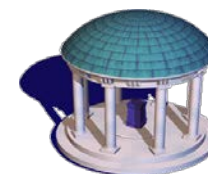
Matching unambiguous, despite epipolar line constraints.

Regularize by smoothness (ordering),

E.g., Marr and Poggio 1978. Arbib and Dev. 1977.

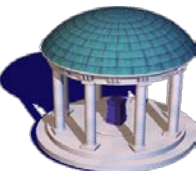
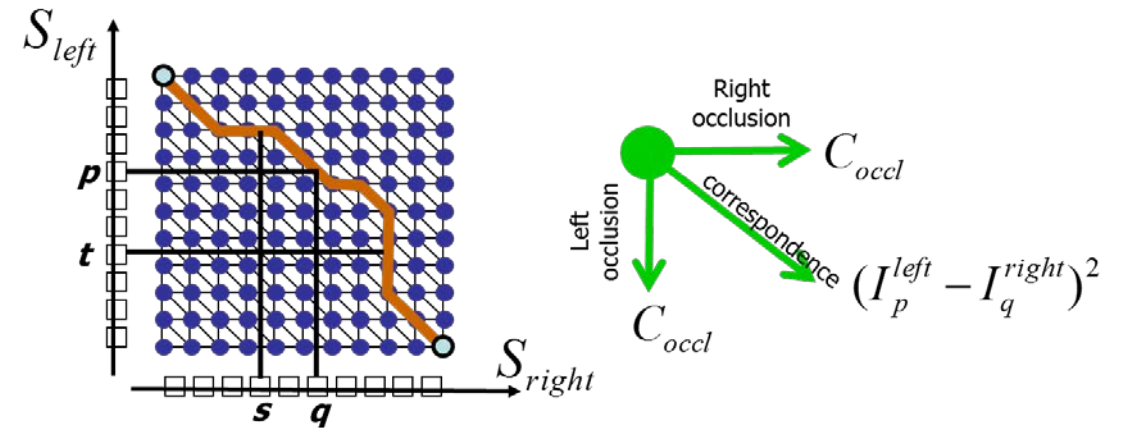
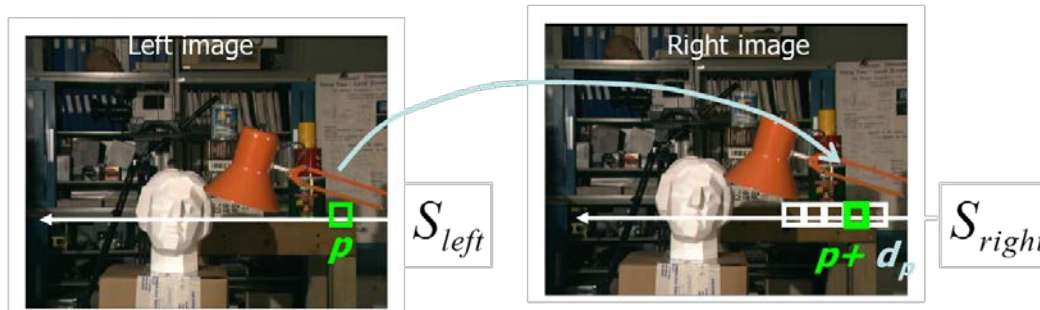
Simple Energy function: (Boykov)

$$E(d_1, d_2, \dots, d_n) = \sum_{p \in S_{\text{left}}} (I_p^{\text{left}} - I_{p+d_p}^{\text{right}})^2 + \sum_{p \in S_{\text{left}}} (d_p - d_{p+1})^2$$

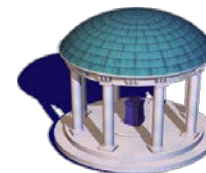


# Exploit Epipolar Line constraint

- The epipolar line constraint reduces correspondence to a one-dimensional problem.
- Dynamic programming can be applied.

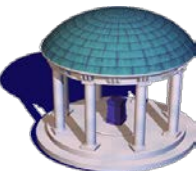
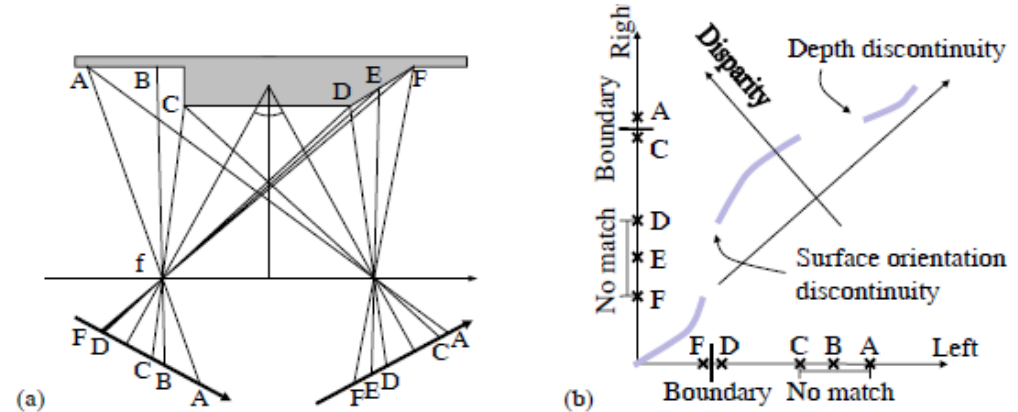


# Results using Dynamic Programming



# Half-Occlusions.

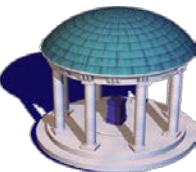
- Da Vinci's stereopsis.
- Points are half-occluded: visible to one eye/camera but not to the other.
- This gives cues for the detection of boundaries.
- Geiger, Ladendorf, Yuille, 1992, Belhumeur and Mumford 1992.
- H. Ishikawa and D. Geiger. 1998 (across epipolar lines).





# Camera Calibration

- Essential Matrix (Longuet-Higgins 1981). Fundamental Matrix (Q.T. Luong and O. D. Faugeras 1992, Hartley 1992).
- More calibration (Z. Zhang 2000).
- More reading on geometry:
- R. Hartley and A. Zisserman. Multiple View Geometry in computer vision. 2003.
- Y. Ma, S. Soatto, J. Kosecka, and S. Sastry. An Invitation to 3-D Vision. 2004.



# Essential Matrix 1

(6) Essential Matrix

Fundamental Matrix (World Coordinates)

$$\underline{X} = (X_1, X_2, X_3) \quad \underline{X}' = (X'_1, X'_2, X'_3)$$

same 3D point  
different coordinate system.

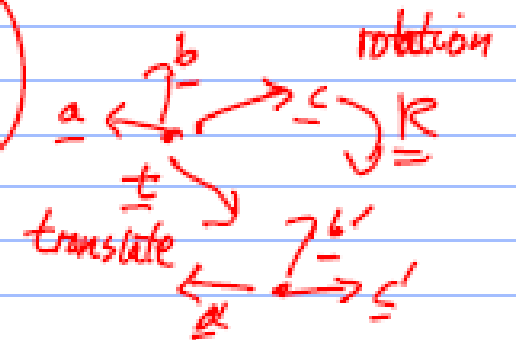
Analysis based on camera coordinates.

$$\begin{pmatrix} u \\ v \end{pmatrix} = \frac{f}{X_3} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix},$$

$$\begin{pmatrix} u' \\ v' \end{pmatrix} = \frac{f}{X'_3} \begin{pmatrix} X'_1 \\ X'_2 \end{pmatrix}$$

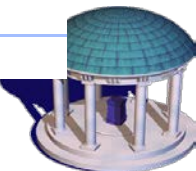
$$\underline{u} \triangleq \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \frac{f}{X_3} \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}, \quad \underline{u}' \triangleq \begin{pmatrix} u' \\ v' \\ 1 \end{pmatrix} = \frac{f}{X'_3} \begin{pmatrix} X'_1 \\ X'_2 \\ X'_3 \end{pmatrix}$$

$$\underline{u}' \triangleq \begin{pmatrix} u' \\ v' \\ 1 \end{pmatrix} = \frac{f}{X'_3} \begin{pmatrix} X'_1 \\ X'_2 \\ X'_3 \end{pmatrix}$$



$$\underline{X} = \underline{R} \underline{X}' + \underline{t}$$

$\underline{R}$  - relative rotation of camera  
 $\underline{t}$  - translation



# Essential Matrix 2

(7)

Essential Matrix

$$\underline{E} = \underline{R} [\underline{t}]_x$$

← rotation → matrix representation of cross product with  $\underline{t}$

$[\underline{t}]_x$  is matrix with components  $\sum_k \epsilon_{ijk} t_k$

Claim  $\underline{u}'^T \underline{E} \underline{u} = 0$

$\epsilon_{ijk}$  anti-symmetric tensor  
imposes constraint on correspond. | Epipolar line  
T - transpose

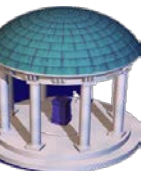
Proof  $\underline{u}'^T \underline{E} \underline{u} \propto \underline{x}'^T \underline{E} \underline{x}$

$$\underline{x}'^T \underline{E} \underline{x} = (\underline{x} - \underline{t})^T \underline{R}^T \underline{R} [\underline{t}]_x \underline{x} = (\underline{x} - \underline{t})^T [\underline{t}]_x \underline{x} = 0$$

$$\sum_{ijk} (x_i - t_i) \epsilon_{ijk} t_k x_j = \sum_{ijk} \epsilon_{ijk} x_i x_j t_k - \sum_{ijk} \epsilon_{ijk} t_i x_j t_k$$

0                      0

Note:  $\epsilon_{ijk} = -\epsilon_{jik}$ ,  $\epsilon_{ijk} = -\epsilon_{ikj}$ ,  $\epsilon_{123} = 1$



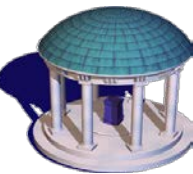
# Structure from Motion: Rigid

- Linear Projection. 3D structure can be estimated by linear algebra (Singular Value Decomposition).
- Camera parameters can also be estimated.
- This estimation is up to an ambiguity.
- Main paper:
- C. Tomasi and T. Kanade. 1991.
- But see also: L.L. Kontsevich, M.L. Kontsevich, A. Kh. Shen.
- "Two Algorithms for Reconstruction Shapes". Avtometriya. 1987



# Structure from Motion: Rigid.

- Linear projection.
- Set of images is rank 3.
- R. Basri and S. Ullman. Recognition by Linear Combinations of Models. 1991.
- Maths of SVD.



# Structure from Motion: Rigid 1

⑧ Structure from Motion: Rigid.

Linear Projection

Set of 3D points

$$\underline{X}^{\mu} = (X_1^{\mu}, X_2^{\mu}, X_3^{\mu})$$

$(X_1, X_2, X_3)$   
position in 3D.

$\mu = 1 \text{ to } N.$

M views

$$\begin{pmatrix} \underline{K}^t \\ \underline{H}^t \end{pmatrix} = \begin{pmatrix} K_1^t & K_2^t & K_3^t \\ H_1^t & H_2^t & H_3^t \end{pmatrix}$$

Measurements in Image

$$u^{\mu,t} = \underline{K}^t \cdot \underline{X}^{\mu} = \sum_{i=1}^3 K_i^t X_i^{\mu}$$

$$v^{\mu,t} = \underline{H}^t \cdot \underline{X}^{\mu} = \sum_{i=1}^3 H_i^t X_i^{\mu}$$

Note:

$u^{\mu,t}$  &  $v^{\mu,t}$  rank 3 (set of image points of an object lie in a 3D space (Basri & Ullman)).



## Structure from Motion: Rigid 2

7) Estimate shape ( $\underline{X}^M$ ) and viewpoint ( $\underline{K}^t, \underline{H}^t$ ).

Assume noise: Gaussian (enables solution by linear algebra.).

$$u^{M,t} = \sum_{i=1}^3 K_i^t X_i^M + n_1^{M,t}$$

$$v^{M,t} = \sum_{i=1}^3 H_i^t X_i^M + n_2^{M,t}$$

$n_1, n_2$  zero mean  
additive Gaussian noise.

Task: Minimize

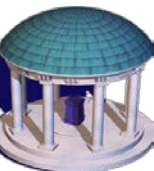
$$E[\underline{K}, \underline{H}, \underline{X}] = \sum_{M,t} \left( u^{M,t} - \sum_{i=1}^3 K_i^t X_i^M \right)^2 + \sum_{M,t} \left( v^{M,t} - \sum_{i=1}^3 H_i^t X_i^M \right)^2$$

w.r.t.  $\underline{K}, \underline{H}, \underline{X}$

simplify

consider

$$\sum_{M,t} \left( u^{M,t} - \sum_{i=1}^3 K_i^t X_i^M \right)^2$$



## Structure from Motion: 3

(10) Solution by Linear Algebra (SVD) up to ambiguity.

$$E[K, X] = \sum_{m,t} \left( u^{m,t} - \sum_{i=1}^3 K_i^t x^m \right)^2$$

Bilinear Problem.

If  $K$  known, solution of  $X$  is linear.

If  $X$  known, solution of  $K$  is linear.

Global minimum can be found by  
Singular Value Decomposition (SVD)

Matrix  $\underline{U}$  - components  $u^{m,t}$

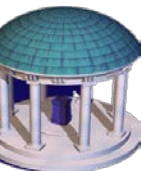
Express  $\underline{U} = \underline{E} \underline{D} \underline{F}^T$

diagonal  $\underline{D} = \begin{pmatrix} d_1 & 0 & 0 \\ 0 & d_2 & \\ & & d_3 \\ & & & \ddots \end{pmatrix}$

CHECK  $\underline{E} \underline{E}^T = \underline{I} = \underline{F} \underline{F}^T$

column's of  $\underline{E}$  and  $\underline{F}$  correspond to eigenvectors

of  $\underline{U} \underline{U}^T$  and  $\underline{U}^T \underline{U}$  respectively.





# Structure from Motion: Rigid 4

(10) Solution by Linear Algebra (SVD) up to ambiguity.

$$E[K, X] = \sum_{i=1}^n \left( u^{M,t} - \sum_{i=1}^3 K_i^t x_i^M \right)^2$$

Bilinear Problem.

If  $K$  known, solution of  $X$  is linear.

If  $X$  known, solution of  $K$  is linear.

Global minimum can be found by

Singular Value Decomposition (SVD)

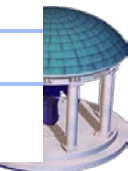
Matrix  $\underline{U}$  - components  $u^{M,t}$

Express  $\underline{U} = \underline{E} \underline{D} \underline{F}^T$

diagonal  $\underline{D} = \begin{pmatrix} d_1 & 0 & 0 \\ 0 & d_2 & \\ & & d_3 \end{pmatrix}$

CHECK  $\underline{E} \underline{E}^T = \underline{I} = \underline{F} \underline{F}^T$

column's of  $\underline{E}$  and  $\underline{F}$  correspond to eigenvectors of  $\underline{U} \underline{U}^T$  and  $\underline{U}^T \underline{U}$  respectively.



## Structure from Motion: 5

(12) Solution

Let  $e_k(t)$   $f_k(\mu)$  be first three  
 $k=1,2,3$  columns of  $\underline{U}$  and  $\underline{V}$

Then solutions are of form:

$$K_i^t = \sum_{k=1}^3 P_{ik} e_k(t)$$

$$X_i^u = \sum_{k=1}^3 Q_{ik} f_k(\mu)$$

where  $\underline{P} \underline{Q}^T = \underline{I}_3$

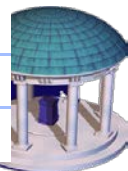
$$\underline{U} = (e_1(t) \ e_2(t) \ \dots)$$

$$\underline{V} = (f_1(\mu) \ f_2(\mu) \ \dots)$$

Same ambiguity

$$\underline{P} \rightarrow \underline{P} \underline{A}$$

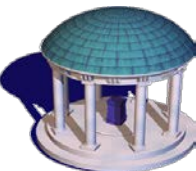
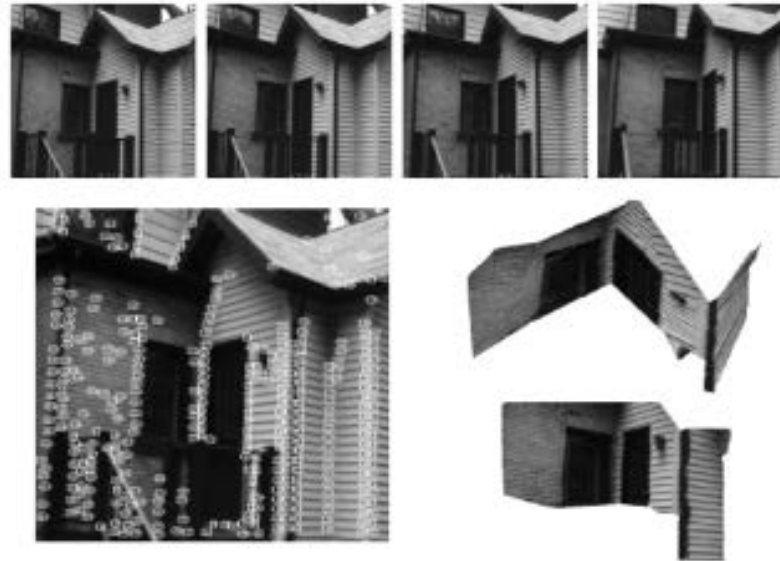
$$\underline{Q}^T \rightarrow \underline{A}^{-1} \underline{Q}$$



# Structure From Motion: Rigid.

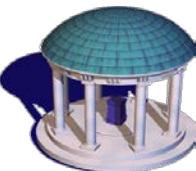
SVD Results: From Tomasi and Kanade. 1991

But see: Kontsevich et al.



## Extension to Non-Rigid Motion

- This approach can be extended to a special class of non-rigid motion.
- The object can be expressed as a linear sum of basis functions. The sum varies over time.
- C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3D shape from image streams. CVPR. 2000.
- Theory clarified by:
- Y. Dai, H. Li, and M. He. A Simple Prior-free Method for Non Rigid Structure from Motion Factorization, in CVPR 2012 (ORAL). IEEE CVPR Best Paper Award-2012. (Code available).
- <http://users.cecs.anu.edu.au/~hongdong/>



# Non-rigid motion

## [12] Relax Rigidity: Basis Function Models (Bregler et al.)

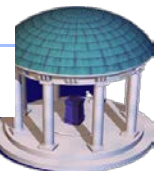
$$\begin{aligned} \underline{u}^{M,t} &= \underline{K}^t \cdot \underline{X}^{M,t} = \sum_{i=1}^3 \underline{K}_i^t X_i^{M,t} && \text{relax rigidity} && \text{Too relaxed} \\ \underline{v}^{M,t} &= \underline{H}^t \cdot \underline{X}^{M,t} = \sum_{i=1}^3 \underline{H}_i^t X_i^{M,t} && X_i^M \rightarrow X_i^{M,t} && \text{Problem is ill-posed.} \end{aligned}$$

Now suppose that the object  $\underline{X}^{M,t}$  can be expressed as a linear combination of bases:

$$\underline{X}^{M,t} = \sum_{a=1}^A \beta_a^t \underline{B}_a^M \quad \leftarrow \text{bases for object}$$

$$\begin{aligned} \underline{u}^{M,t} &= \sum_{i=1}^3 \sum_{a=1}^A \underline{K}_i^t \beta_a^t \underline{B}_{i,a}^M \\ \underline{v}^{M,t} &= \sum_{i=1}^3 \sum_{a=1}^A \underline{H}_i^t \beta_a^t \underline{B}_{i,a}^M \end{aligned}$$

coefficients, depend on  $t$ .  
SVD - twice  
More Complex - ambiguity  
Best Paper CVPR 2012.



# Manhattan World

- Many scenes, particularly man-made scenes, have a natural three-dimensional coordinate systems caused by the structure of the world.

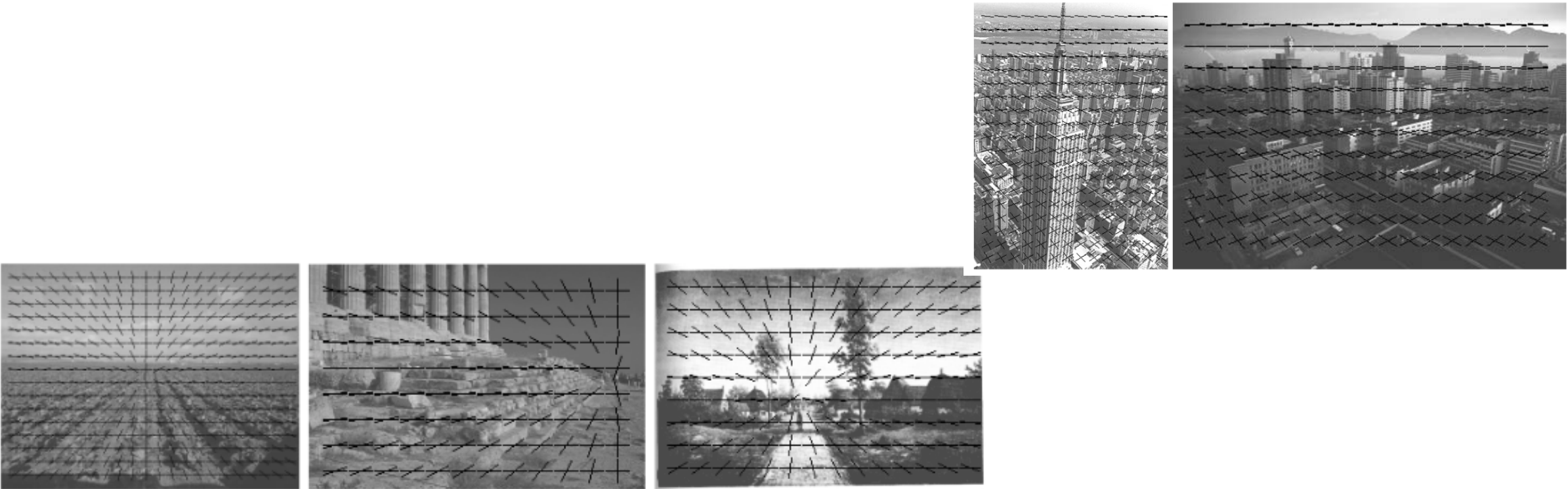
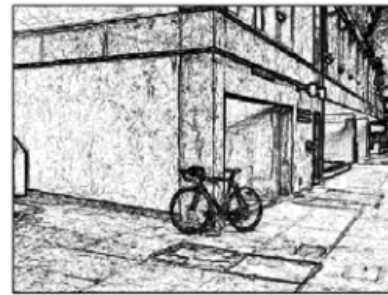


Figure 5: Results on an American mid-west broccoli field, the ruins of the

# Manhattan World: 2

- Back to Ames Room:
- Non-Manhattan Edges:



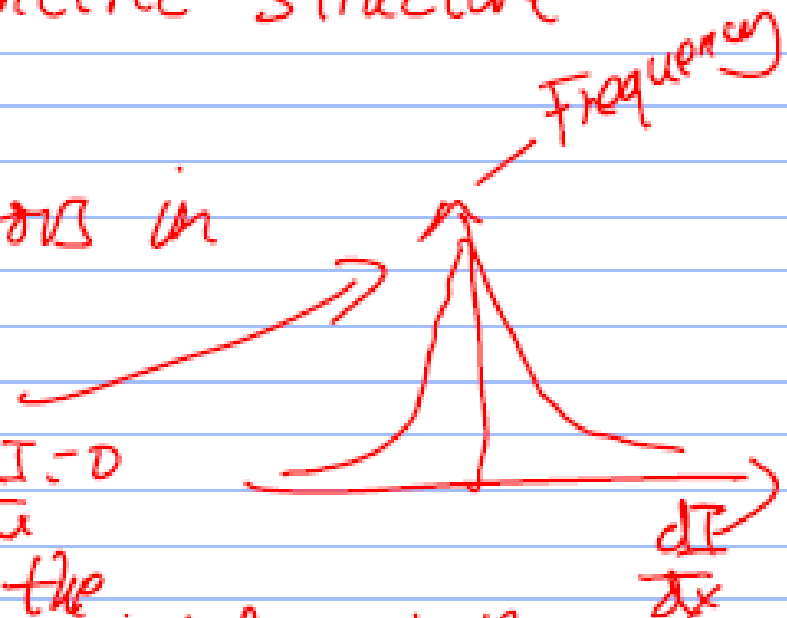
# Manhattan World 1

(13) Manhattan World.

Exploit regularities in the geometric structure of scenes.

Statistics of derivative operators in images (Simoncelli, Zhu, ...)

Statistical justification of piecewise smoothness. Histogram of  $\frac{dI}{dx}$  peaks at  $\frac{dI}{dx} = 0$



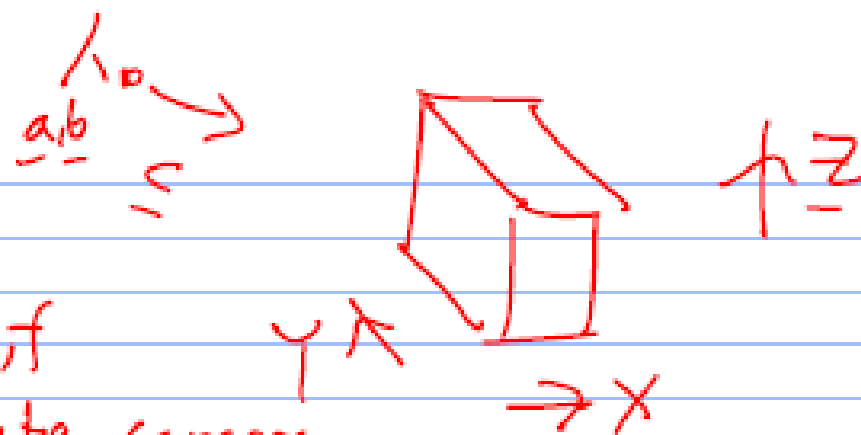
But there are also regularities in the directions of edges, depending on the orientation of the camera relative to the geometric structure of the scene.



# Manhattan World 2

## (24) Manhattan World

- distribution of edges in the image depends on camera parameters  $(a, b, c), f$
- use structure of scene to calibrate camera.
  - estimate  $(a, b, c), f$ .



Assume : A pixel  $\underline{u}$  in image is

- |                                |                         |
|--------------------------------|-------------------------|
| (i) An edge in the X direction | $m_{\underline{u}} = 1$ |
| (ii) " " " " Y " "             | $m_{\underline{u}} = 2$ |
| (iii) " " " " Z " "            | $m_{\underline{u}} = 3$ |
| (iv) An unaligned edge         | $m_{\underline{u}} = 4$ |
| (v) Not an edge                | $m_{\underline{u}} = 5$ |

The  $m$ 's are  
hidden  
variables,

# Manhattan World 3

15)

Generative Model for the image gradient  $\underline{\nabla}I$  at  $\underline{u}$

$$P(\underline{\nabla}I(\underline{u}) \mid m_{\underline{u}}, \underline{\psi})$$

$$\underline{\psi} = (\underline{a}, \underline{b}, \underline{c}, \underline{f})$$

camera parameters.

if  $m_{\underline{u}} = 5$  (no edge)

$P(\underline{\nabla}I(\underline{u}) \mid m_{\underline{u}} = 5, \underline{\psi}) \approx 0$ , unless  $|\underline{\nabla}I(\underline{u})|$  is small

if  $m_{\underline{u}} = 1$  (edge in X direction)

$P(\underline{\nabla}I(\underline{u}) \mid m_{\underline{u}} = 1, \underline{\psi}) \approx 0$ , unless  $|\underline{\nabla}I(\underline{u})|$  is large  
and  $\underline{\nabla}I(\underline{u})$  points in direction  
of projection of X-line.

# Manhattan World 4

(16)

Full Model.

Specify prior  $P(m_u)$

$$P(\underline{\nabla I(u)} | \underline{\psi}) = \sum_{m=1}^5 P(\underline{\nabla I(u)} | \underline{\psi}, m_u) P(m_u)$$

↗ don't care about  $m_u$ ,  
so sum it out.  
(Freeman's talk).

Assume Independence → for Full Image.

$$P(\langle \underline{\nabla I(u)} \rangle | \underline{\psi}) = \prod_u P(\underline{\nabla I(u)} | \underline{\psi}) //$$

Estimate  $\hat{\underline{\psi}} = \underset{\underline{\psi}}{\text{ARG MAX}} P(\langle \underline{\nabla I(u)} \rangle | \underline{\psi})$ .

Very Simple → Works fairly well  
More Sophisticated Models

# Beyond Manhattan:

- There are parallel lines in the scene, but they are not orthogonal.



Edge Segments

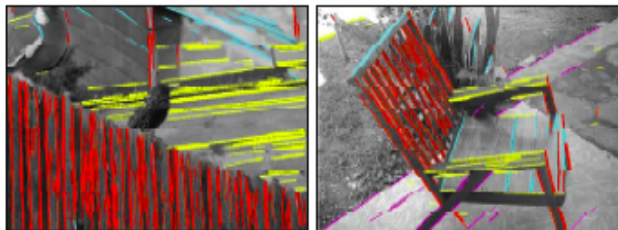
Family-1



Family-2

Family-3

(a) Orthogonal Families

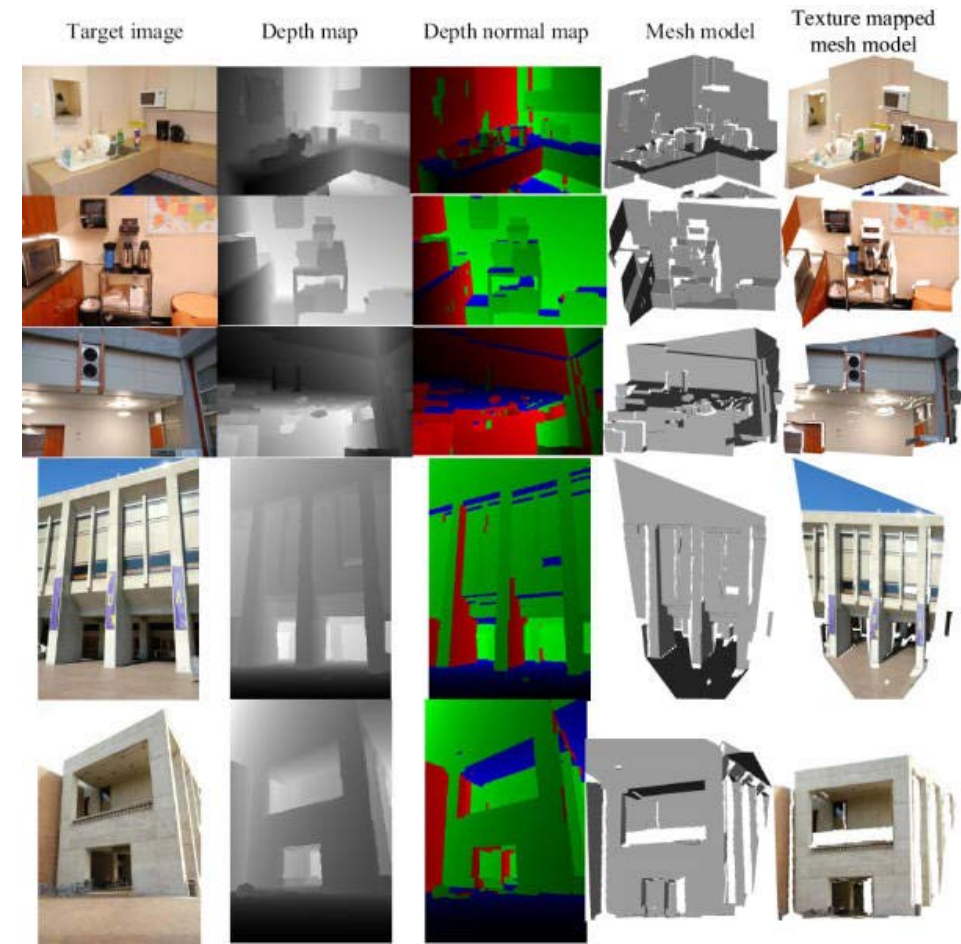


(b) Non-orthogonal Families

(c) Mixed Families

# Manhattan World

- Manhattan World stereo:
- Piecewise planar surfaces with dominant
- Directions.
- Instead of assuming
- Piecewise smoothness.
- Video available:
- <http://grail.cs.washington.edu/projects/manhattan/>
- Y. Fuukawa, B. Curless, S. Seitz, and R. Szeliski. 2009.



# Manhattan World Grammar

- Website: <http://www.youtube.com/watch?v=s0mhpKFv36g>

## Building Reconstruction using Manhattan-World Grammars

Carlos A. Vanegas

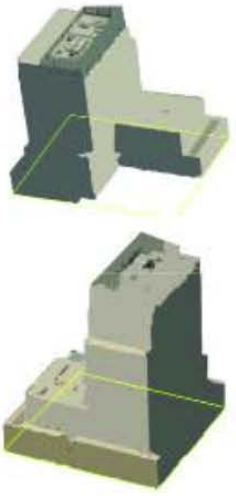
Daniel G. Aliaga

Bedřich Beneš

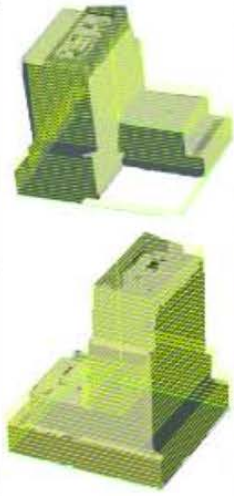
Purdue University



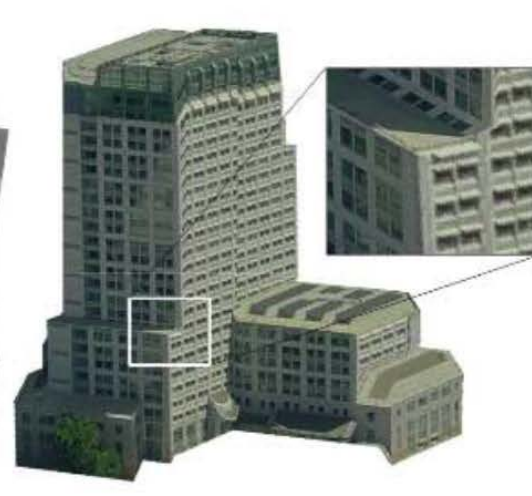
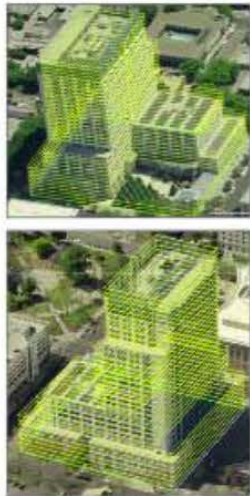
Photos and Footprint



Segmented Photos



Adapted Floors



Reconstructed 3D Model