CMSC 472 / 672

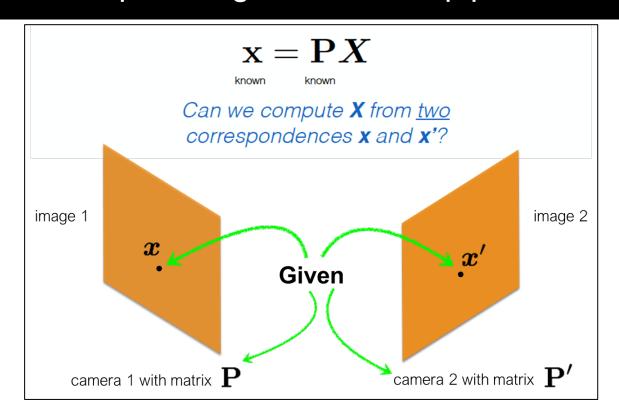
Lecture 16

Stereo Vision

When people ask where you see yourself in 10 years



Recap: Triangulation and Epipolar Geometry

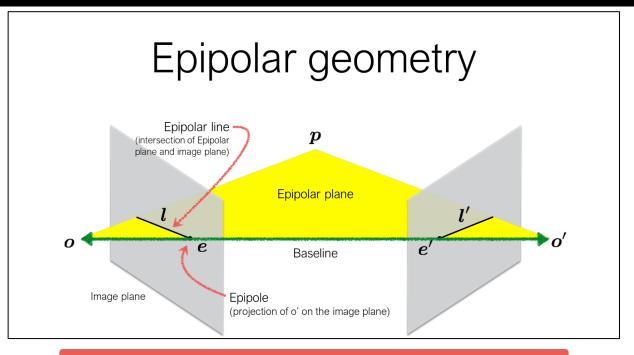


Essential Matrix vs Homography

What's the difference between the essential matrix and a homography?

They are both 3 x 3 matrices but ...

$$m{l}' = \mathbf{E}m{x}$$
 $m{x}' = \mathbf{H}m{x}$ Essential matrix maps a point to a line Homography maps a point to a point



Longuet-Higgins equation $oldsymbol{x'}^{ op}\mathbf{E}oldsymbol{x}=0$

Epipolar lines $egin{aligned} m{x}^ op m{l} &= 0 & m{x'}^ op m{l}' = 0 \ m{l}' &= m{\mathbf{E}}m{x} & m{l} &= m{\mathbf{E}}^Tm{x}' \end{aligned}$

Epipoles $e'^{ op}\mathbf{E}=\mathbf{0}$ $\mathbf{E}e=\mathbf{0}$

(points in normalized <u>camera</u> coordinates)

The fundamental matrix

The Fundamental matrix
is a generalization
of the Essential matrix,
where the assumption of calibrated cameras
is removed

Same equation works in image coordinates!

$$\boldsymbol{x}'^{\top}\mathbf{F}\boldsymbol{x} = 0$$

it maps pixels to epipolar lines

The 8-point algorithm

Assume you have *M* matched *image* points

$$\{\boldsymbol{x_m}, \boldsymbol{x_m'}\}$$
 $m = 1, \dots, M$

Each correspondence should satisfy

$$\boldsymbol{x}_m^{\prime \top} \mathbf{F} \boldsymbol{x}_m = 0$$

How would you solve for the 3 x 3 **F** matrix?

$$\boldsymbol{x}_m'^{\top} \mathbf{F} \boldsymbol{x}_m = 0$$

How many equation do you get from one correspondence?

ONE correspondence gives you ONE equation

$$x_m x'_m f_1 + x_m y'_m f_2 + x_m f_3 + y_m x'_m f_4 + y_m y'_m f_5 + y_m f_6 + x'_m f_7 + y'_m f_8 + f_9 = 0$$

Set up a homogeneous linear system with 9 unknowns

$$\begin{bmatrix} x_1x'_1 & x_1y'_1 & x_1 & y_1x'_1 & y_1y'_1 & y_1 & x'_1 & y'_1 & 1 \\ \vdots & \vdots \\ x_Mx'_M & x_My'_M & x_M & y_Mx'_M & y_My'_M & y_M & x'_M & y'_M & 1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \end{bmatrix} = \mathbf{0}$$

Each point pair (according to epipolar constraint) contributes only one <u>scalar</u> equation

$$\boldsymbol{x}_m^{\prime \top} \mathbf{F} \boldsymbol{x}_m = 0$$

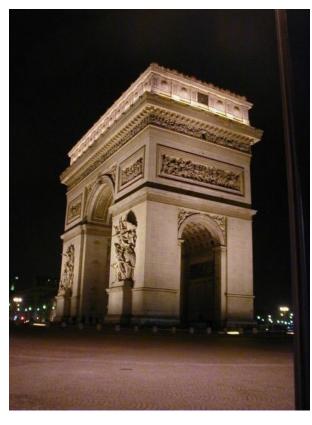
Note: This is different from the Homography estimation where each point pair contributes 2 equations.

We need at least 8 points

Hence, the 8 point algorithm!

Example





epipolar lines



$$\mathbf{F} = \begin{bmatrix} -0.00310695 & -0.0025646 & 2.96584 \\ -0.028094 & -0.00771621 & 56.3813 \\ 13.1905 & -29.2007 & -9999.79 \end{bmatrix}$$

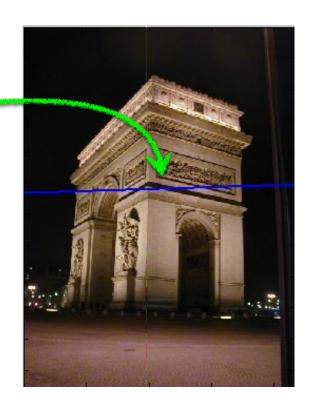
$$x = \begin{bmatrix} 343.53 \\ 221.70 \\ 1.0 \end{bmatrix}$$

$$m{l}' = \mathbf{F} m{x}$$
 $= egin{bmatrix} 0.0295 \\ 0.9996 \\ -265.1531 \end{bmatrix}$

$$m{l}' = \mathbf{F} m{x}$$

$$= \left[egin{array}{c} 0.0295 \\ 0.9996 \\ -265.1531 \end{array} \right]$$





Stereo Imaging



Left image



Right image



Left image



Right image

1. Select point in one image



Left image



Right image

- 1. Select point in one image
- 2. Form the epipolar line for that point in second image



Left image

Right image

- 1. Select point in one image
- 2. Form the epipolar line for that point in second image
- 3. Find matching point along line

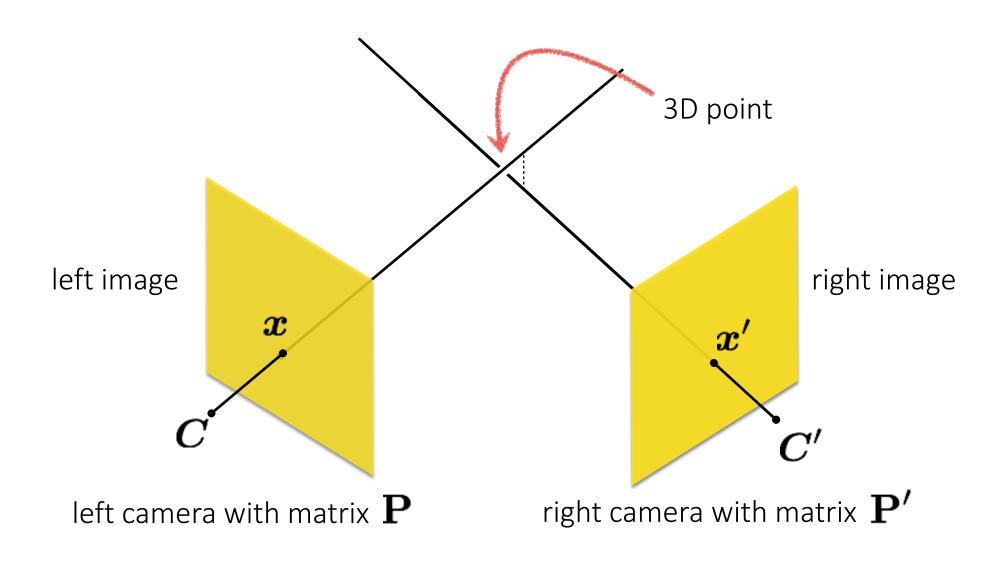


Left image

Right image

- 1. Select point in one image
- 2. Form the epipolar line for that point in second image
- 3. Find matching point along line
- 4. Perform triangulation

Triangulation



Stereo rectification





What's different between these two images?





The amount of horizontal movement is inversely proportional to ...







The amount of horizontal movement is inversely proportional to ...

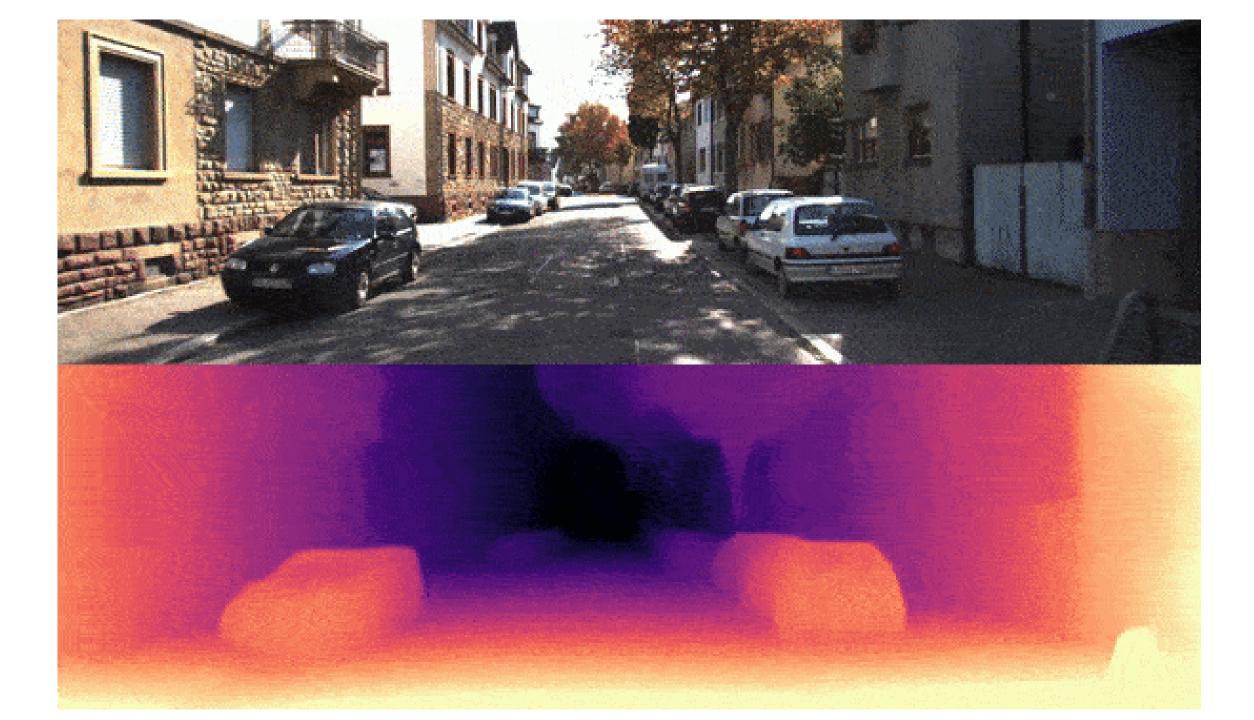


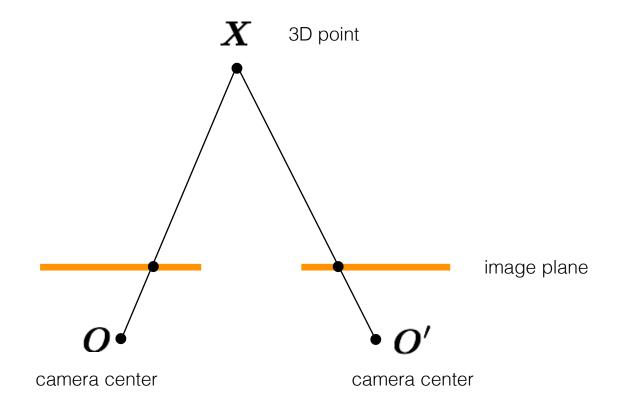


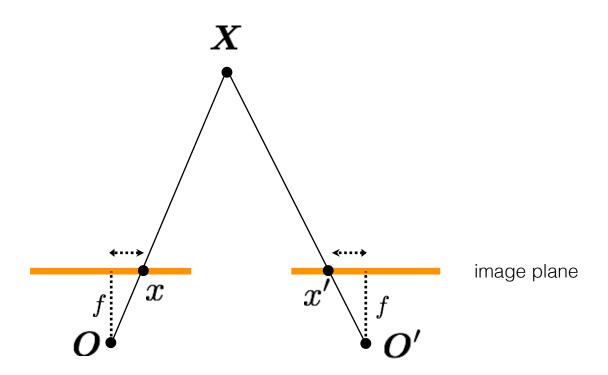


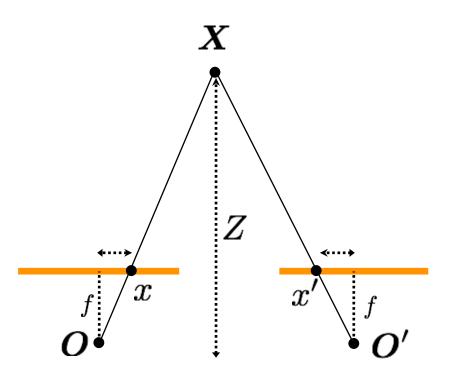
... the distance from the camera.

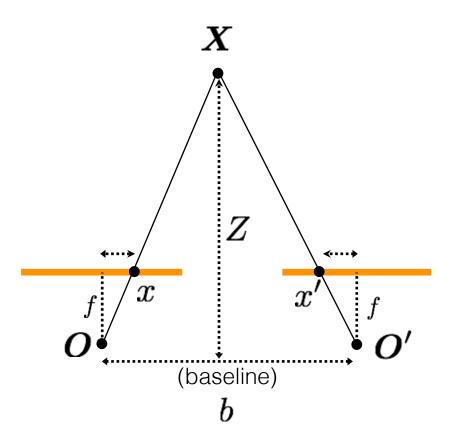
... aka ... *depth*

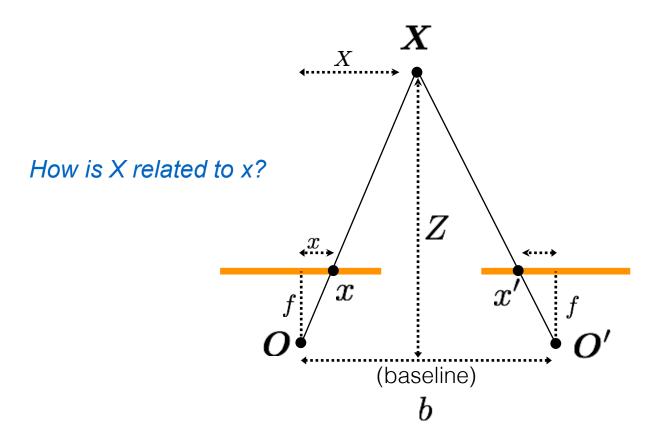


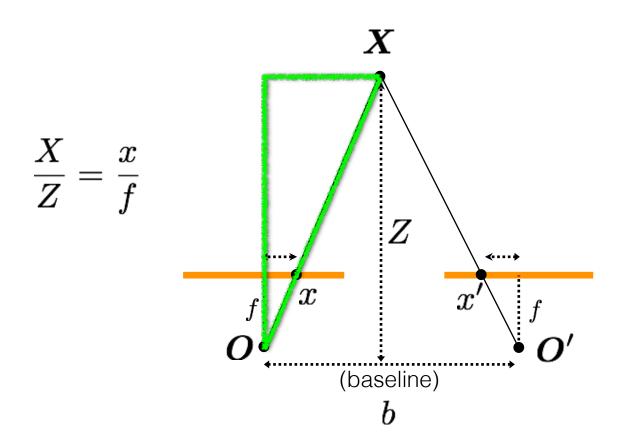


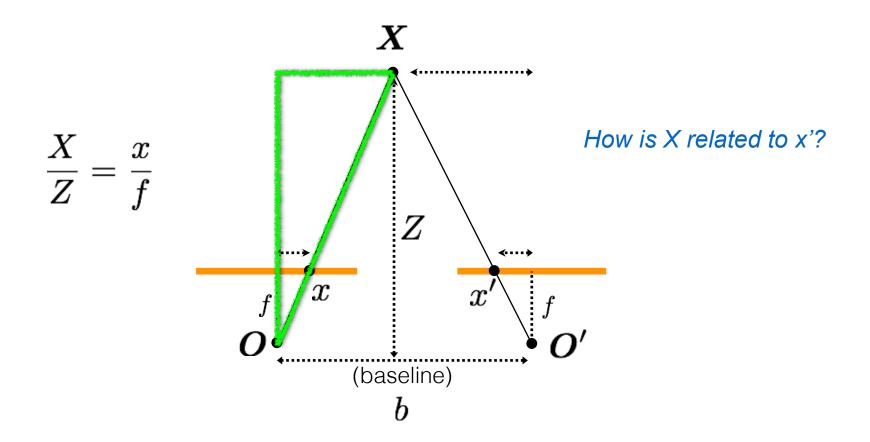


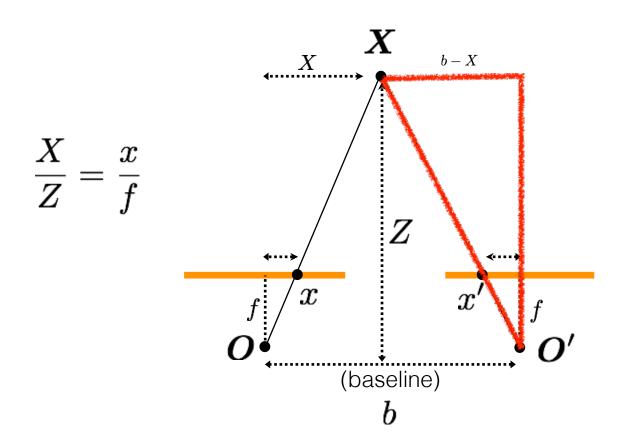




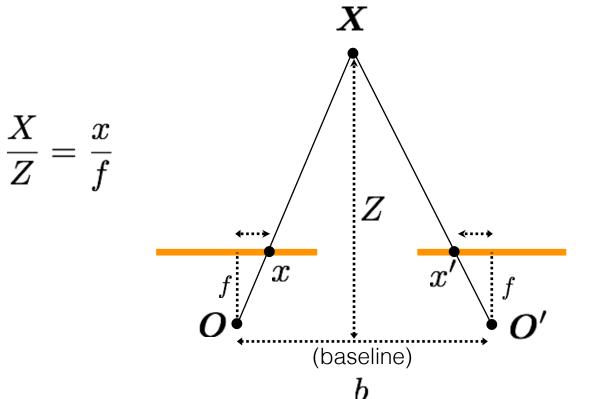






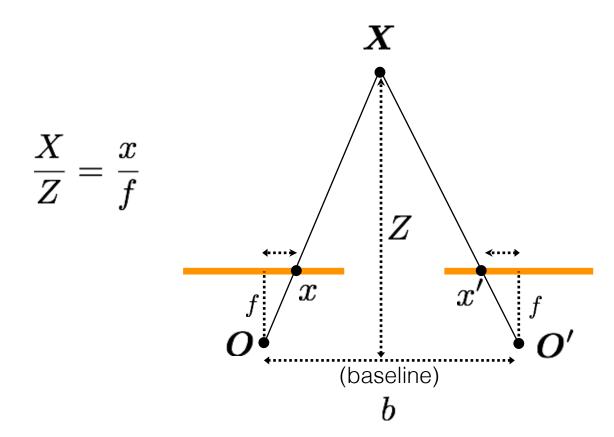


$$\frac{b-X}{Z} = \frac{x'}{f}$$



Disparity

$$d=x-x'$$
 (wrt to camera origin of image plane) $=rac{bf}{z}$



$$\frac{b-X}{Z} = \frac{x'}{f}$$

Disparity

$$d = x - x'$$

$$= \frac{bf}{Z}$$

inversely proportional to depth

Stereoscopes: A 19th Century Pastime







Old **Zeiss** pocket stereoscope with original test image

A **stereoscope** is a device for viewing a <u>stereoscopic pair</u> of separate images, depicting left-eye and right-eye views of the same scene, as a single three-dimensional image.

A typical stereoscope provides each eye with a lens that makes the image seen through it appear larger and more distant and usually also shifts its apparent horizontal position, so that for a person with normal binocular depth perception the edges of the two images seemingly fuse into one "stereo window".

Google Cardboard





Second-generation Google Cardboard viewer

Developer Google

Manufacturer Google, third-party companies

Type Virtual reality platform

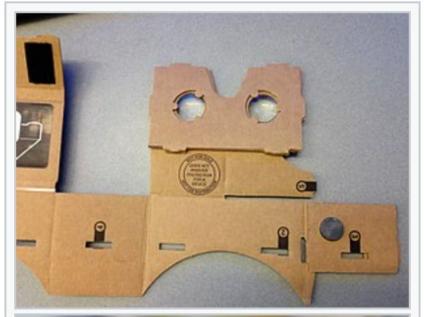
Release date June 25, 2014; 9 years ago

Discontinued March 3, 2021; 2 years ago

(Official viewer, Google Store)

Units 15 million

shipped





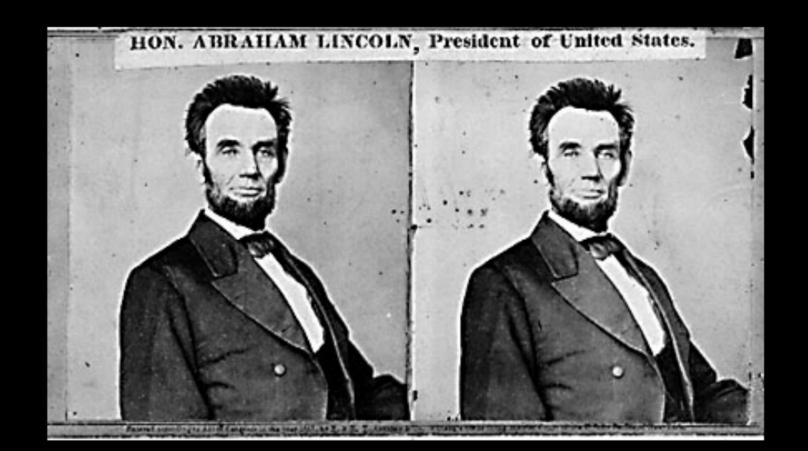
A Cardboard viewer unassembled (top) and assembled (bottom)

Once the kit is assembled, a smartphone is inserted in the back of the device and held in place by the selected fastening device. A Google Cardboard—compatible app splits the smartphone display image into two, one for each eye,

Apps on the mobile phone substitute for stereo cards; these apps can also sense rotation and expand the stereoscope's capacity into that of a full-fledged <u>virtual reality</u> device.

The underlying technology is otherwise unchanged from earlier stereoscopes.

https://en.wikipedia.org/wiki/Google_Cardboard





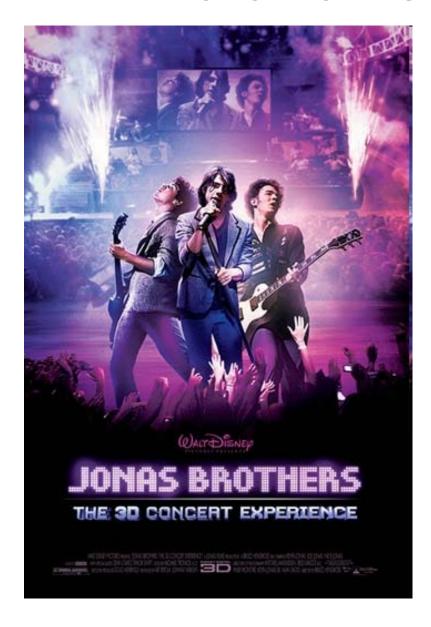
Public Library, Stereoscopic Looking Room, Chicago, by Phillips, 1923

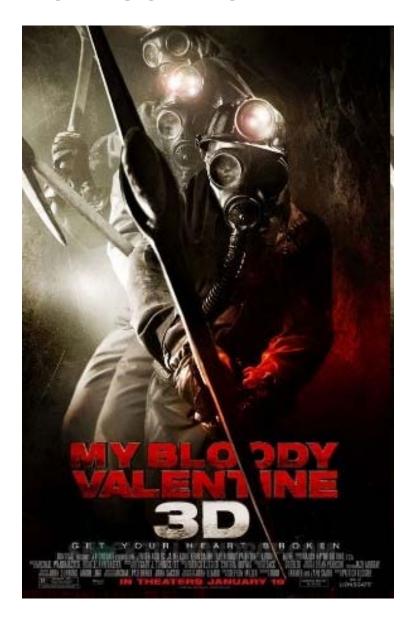




Mark Twain at Pool Table", no date, UCR Museum of Photography

This is how 3D movies work



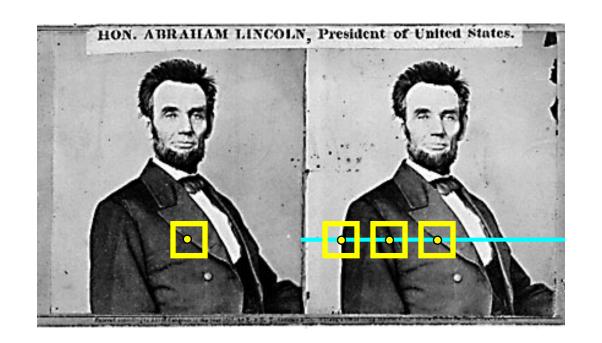


So can I compute depth from any two images of the same object?



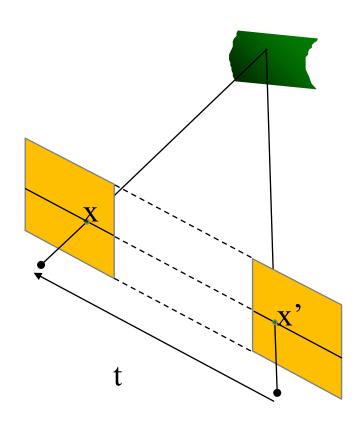


Yes if you can "rectify" them i.e. make epipolar lines horizontal



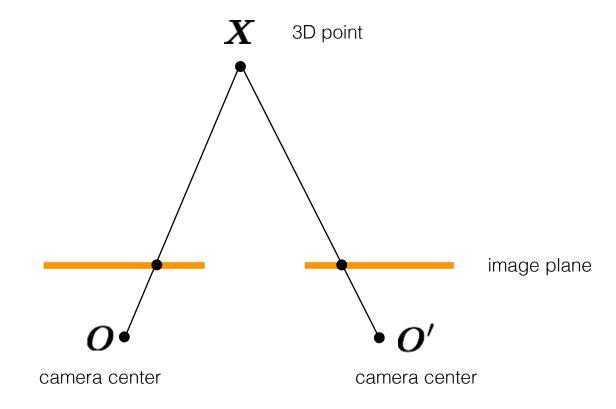
- 1. Rectify images
 (make epipolar lines horizontal)
- 2. For each pixel
 - a. Find epipolar line
 - b. Scan line for best match
 - c. Compute depth from disparity

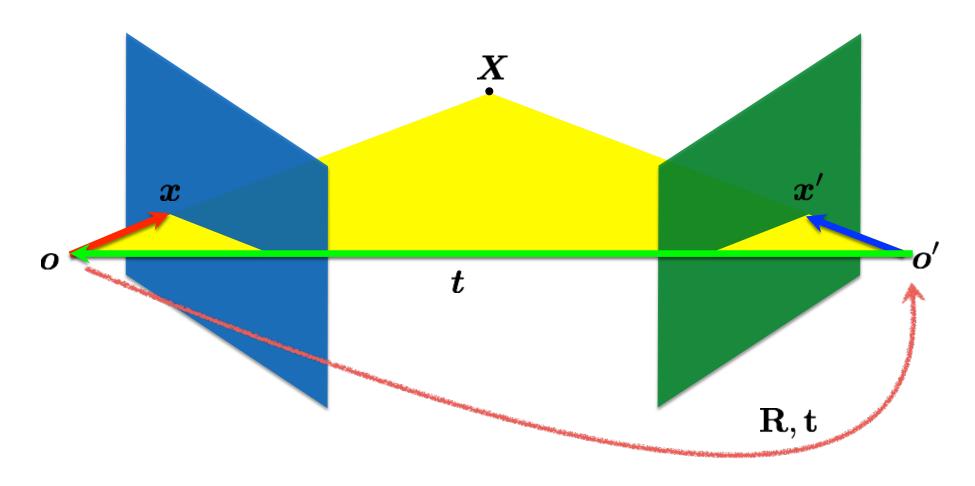
$$Z = \frac{bf}{d}$$



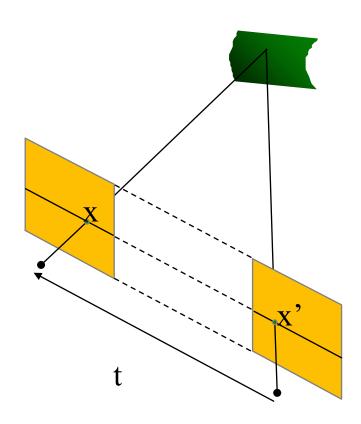
When this relationship holds:

$$R = I \qquad t = (T, 0, 0)$$





$$\boldsymbol{x}' = \mathbf{R}(\boldsymbol{x} - \boldsymbol{t})$$



When this relationship holds:

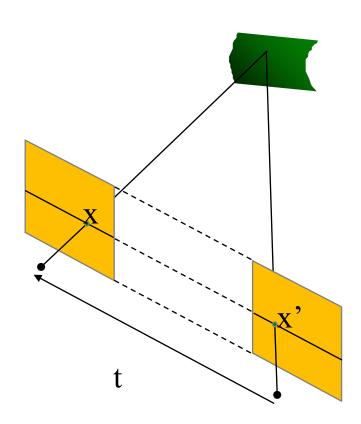
$$R = I \qquad t = (T, 0, 0)$$

Let's try this out...

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

This always has to hold for rectified images

$$x^T E x' = 0$$



Write out the constraint

When this relationship holds:

$$R = I \qquad t = (T, 0, 0)$$

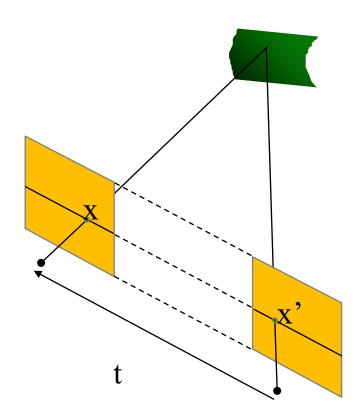
Let's try this out...

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

This always has to hold for rectified images

$$x^T E x' = 0$$

$$\begin{pmatrix} u & v & 1 \\ -T \\ Tv' \end{pmatrix} = 0$$



Write out the constraint

When this relationship holds:

$$R = I \qquad t = (T, 0, 0)$$

Let's try this out...

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

This always has to hold

$$x^T E x' = 0$$

The image of a 3D point will always be on the same horizontal line

$$\begin{pmatrix} u & v & 1 \\ -T \\ Tv' \end{pmatrix} = 0$$

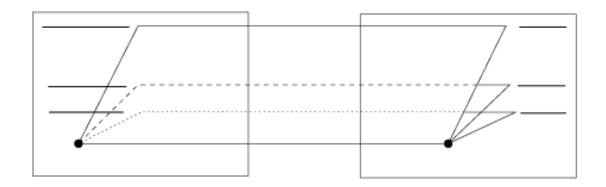
y coordinate is always the same!

Stereo Rectification

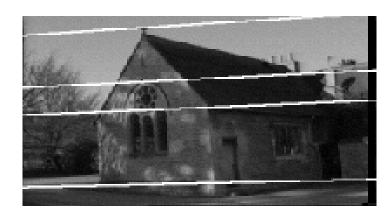
- Rotate the right camera by R

 (aligns camera coordinate system orientation only)
- 2. Rotate (**rectify**) the left camera so that the epipole is at infinity
- 3. Rotate (**rectify**) the right camera so that the epipole is at infinity
- 4. Adjust the scale

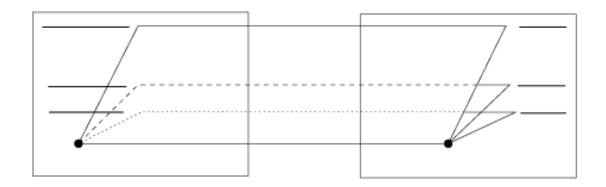
Parallel cameras

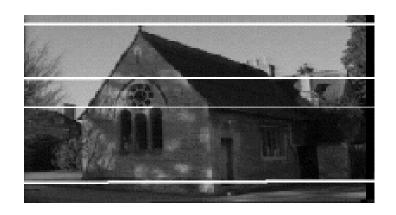


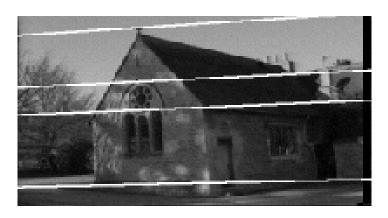




Parallel cameras







epipole at infinity

Setting the epipole to infinity

(Building **R**_{rect} from **e**)

Let
$$R_{
m rect}=\left[egin{array}{c} m{r}_1^{ op} \ m{r}_2^{ op} \ m{r}_3^{ op} \end{array}
ight]$$
 Given: epipole $m{e}$ (using SVD on E) (translation from $m{E}$)

$$oldsymbol{r}_1 = oldsymbol{e}_1 = rac{T}{||T||}$$
 epipole coincides with translation vector

$$m{r_2} = rac{1}{\sqrt{T_x^2 + T_y^2}} \left[egin{array}{c} -T_y & T_x & 0 \end{array}
ight]
ight. egin{array}{c} ext{cross product of e and the direction vector of the optical axis} \end{array}$$

$$\boldsymbol{r}_3 = \boldsymbol{r}_1 \times \boldsymbol{r}_2$$

orthogonal vector

If
$$m{r}_1 = m{e}_1 = rac{T}{||T||}$$
 and $m{r}_2$ $m{r}_3$ orthogonal

then
$$R_{ ext{rect}}oldsymbol{e}_1=\left[egin{array}{c} oldsymbol{r}_1^ opoldsymbol{e}_1\ oldsymbol{r}_2^ opoldsymbol{e}_1\ oldsymbol{r}_3^ opoldsymbol{e}_1 \end{array}
ight]=\left[egin{array}{c} 1\ 0\ 0 \end{array}
ight]$$

At x-infinity

Stereo Rectification Algorithm

- 1. Estimate **E** using the 8 point algorithm (SVD)
- 2. Estimate the epipole **e** (SVD of **E**)
- 3. Build Rrect from e
- 4. Decompose **E** into **R** and **T**
- 5. Set $\mathbf{R}_1 = \mathbf{R}_{\text{rect}}$ and $\mathbf{R}_2 = \mathbf{R}\mathbf{R}_{\text{rect}}$
- 6. Rotate each left camera point (warp image) $[x' y' z'] = \mathbf{R}_1 [x y z]$
- 7. Rectified points as $\mathbf{p} = f/z'[x' y' z']$
- 8. Repeat 6 and 7 for right camera points using \mathbf{R}_2

Use built-in OpenCV functions for this

stereoRectifyUncalibrated()

Computes a rectification transform for an uncalibrated stereo camera.

Parameters

points1 Array of feature points in the first image.

points2 The corresponding points in the second image. The same formats as in findFundamentalMat are supported.

F Input fundamental matrix. It can be computed from the same set of point pairs using findFundamentalMat.

imgSize Size of the image.

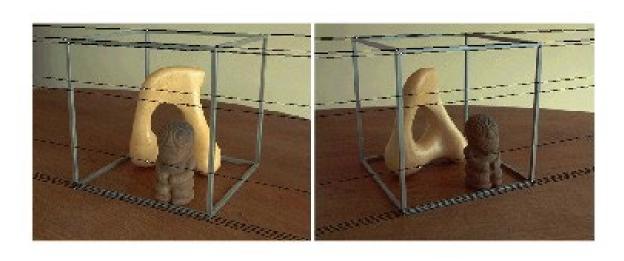
H1 Output rectification homography matrix for the first image.

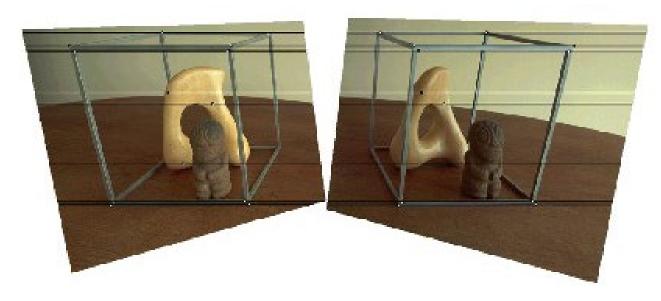
H2 Output rectification homography matrix for the second image.

threshold Optional threshold used to filter out the outliers. If the parameter is greater than zero, all the point pairs that do not comply with the epipolar geometry (that is, the points for which $|points2[i]^T *F *points1[i]| > threshold$) are rejected prior to computing the homographies. Otherwise, all the points are considered inliers.

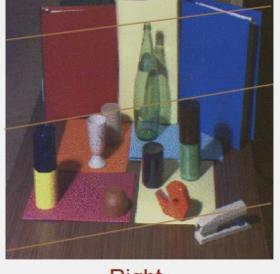
The function computes the rectification transformations without knowing intrinsic parameters of the cameras and their relative position in the space, which explains the suffix "uncalibrated". Another related difference from stereoRectify is that the function outputs not the rectification transformations in the object (3D) space, but the planar perspective transformations encoded by the homography matrices H1 and H2. The function implements the algorithm [88].

Rectification example





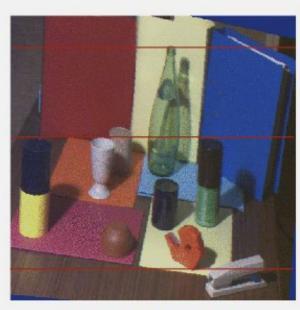




Left Right



Rectified Left



Rectified Right



What can we do after rectification?



Depth Estimation



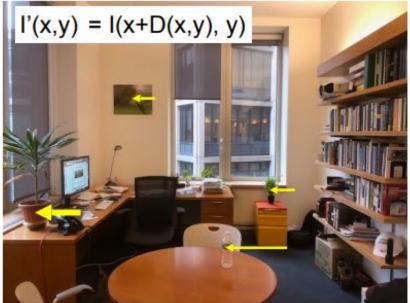


Depth Estimation via Stereo Matching



Disparity map



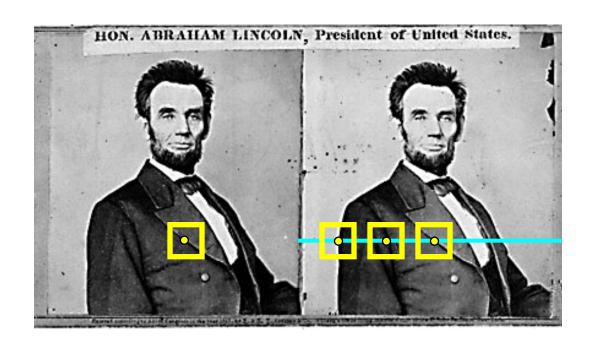


$$Z(x,y) = \frac{f}{D(x,y)}$$

Finding correspondences



We only need to search for matches along horizontal lines.



- 1. Rectify images
 (make epipolar lines horizontal)
- 2. For each pixel
 - a. Find epipolar line
 - b. Scan line for best match

c. Compute depth from disparity

How would you do this?

$$Z=rac{bf}{d}$$

Computing disparity

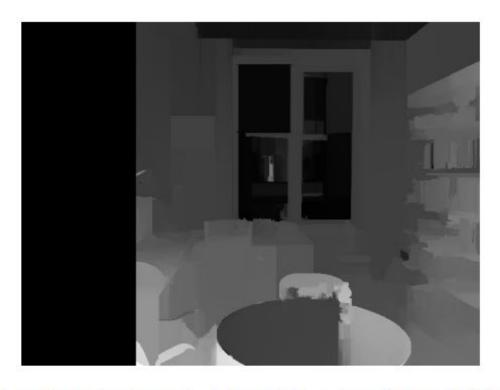




Computing disparity

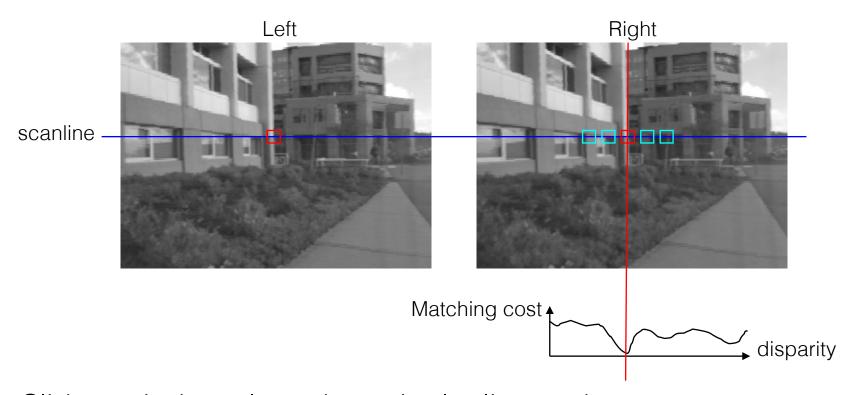






Semi-global matching [Hirschmüller 2008]

Stereo Block Matching

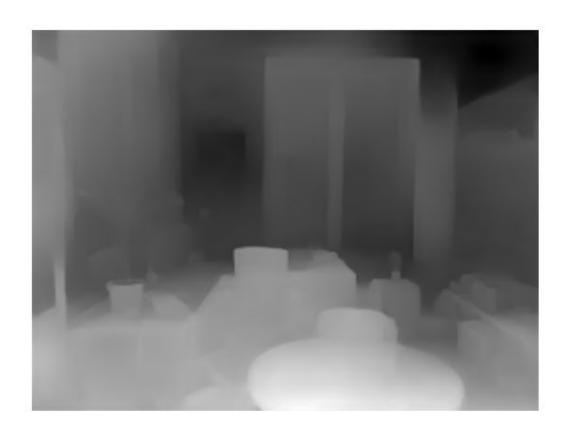


- Slide a window along the epipolar line and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

Depth from Single Image

Can also learn depth from a single image





MegaDepth: Learning Single-View Depth Prediction from Internet Photos

Depth from Single Image

Use inference power of deep learning to regress depth directly from single image

Not as accurate as stereo methods, but still solves ambiguity issues through semantic cues

Depth Map Prediction from a Single Image using a Multi-Scale Deep Network

David Eigen deigen@cs.nyu.edu

Christian Puhrsch

Rob Fergus

cpuhrsch@nyu.edu

fergus@cs.nyu.edu

Dept. of Computer Science, Courant Institute, New York University

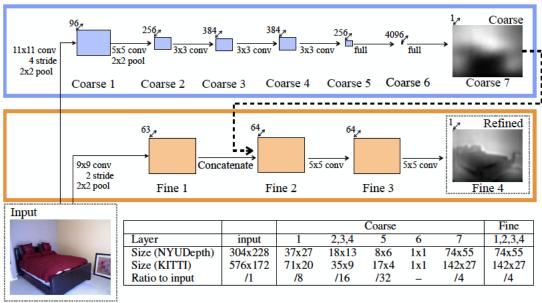


Figure 1: Model architecture.

My Research: ICCV 2021: Answering Questions about Images using Depth Information

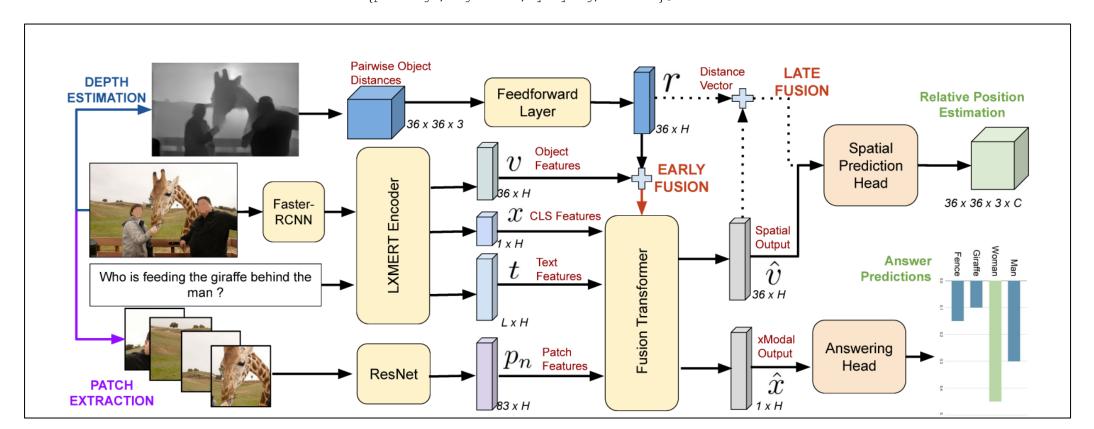


This ICCV paper is the Open Access version, provided by the Computer Vision Foundation.

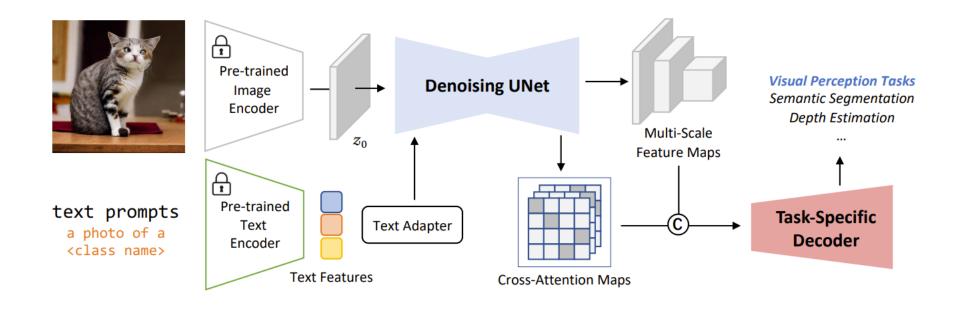
Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

Weakly Supervised Relative Spatial Reasoning for Visual Question Answering

Pratyay Banerjee Tejas Gokhale Yezhou Yang Chitta Baral Arizona State University
{pbanerj6, tgokhale, yz.yang, chitta}@asu.edu



VPD: Language-Guided Depth Estimation



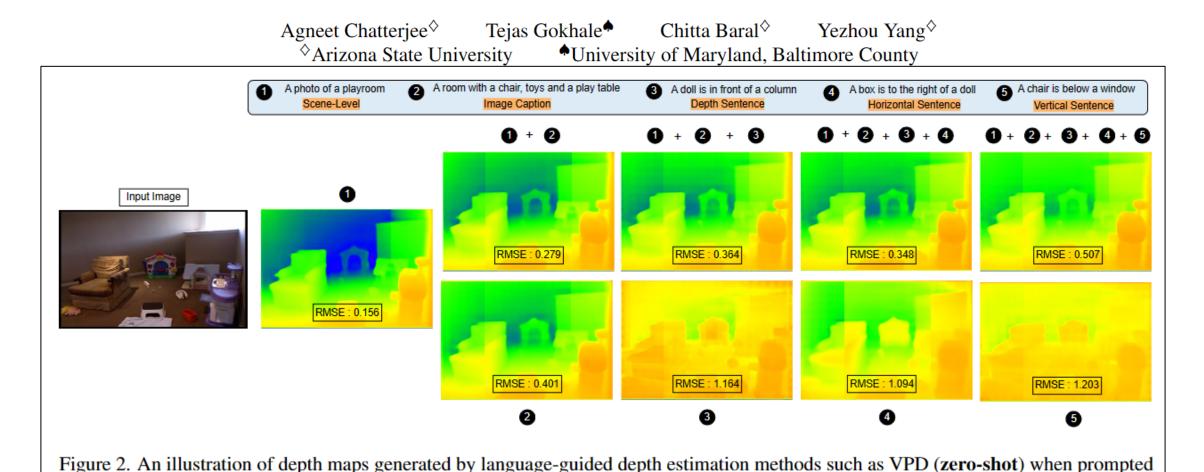
Unleashing Text-to-Image Diffusion Models for Visual Perception

Wenliang Zhao^{1*} Yongming Rao^{1*} Zuyan Liu^{1*} Benlin Liu² Jie Zhou¹ Jiwen Lu^{1†}

¹Tsinghua University ²University of Washington

My Research: CVPR 2024: Quantifying Efficacy of Language-Guided Depth Estimation

On the Robustness of Language Guidance for Low-Level Vision Tasks: Findings from Depth Estimation



with various sentence inputs that we use as part of our study. The first row shows the effect of progressively adding descriptions as input,

while the second row shows depth maps generated by single sentence inputs.