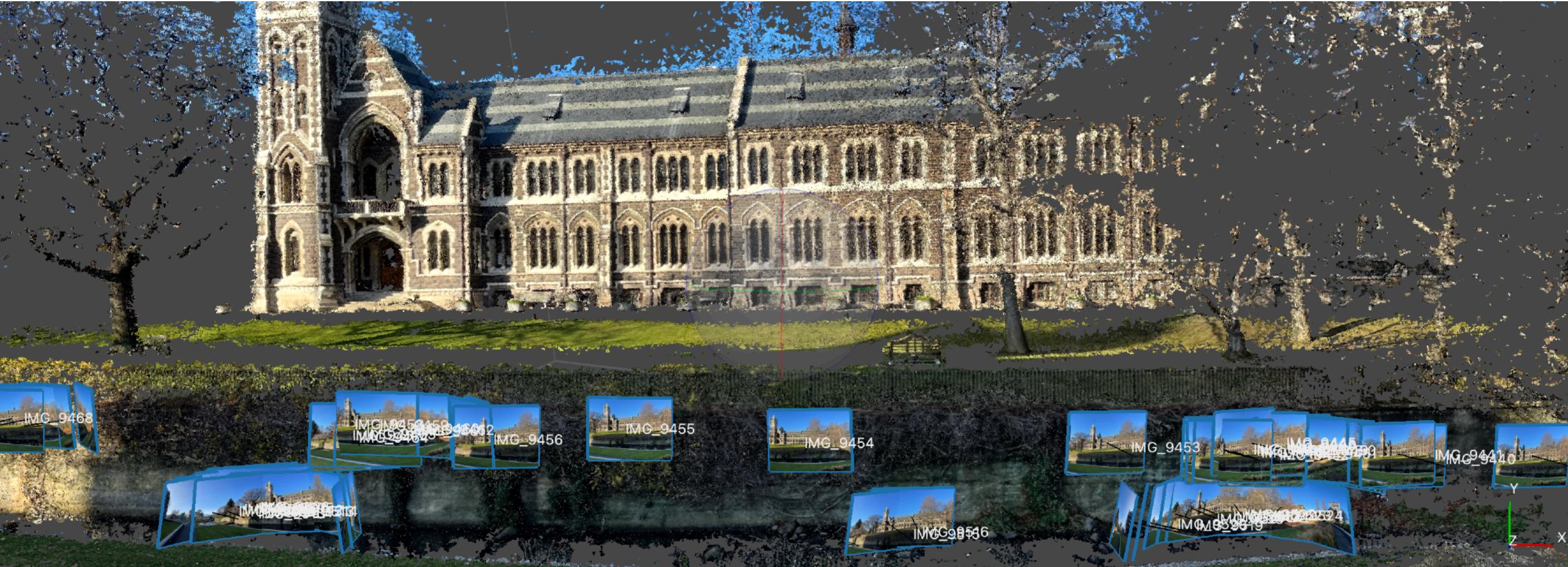


# Visual Computing I:

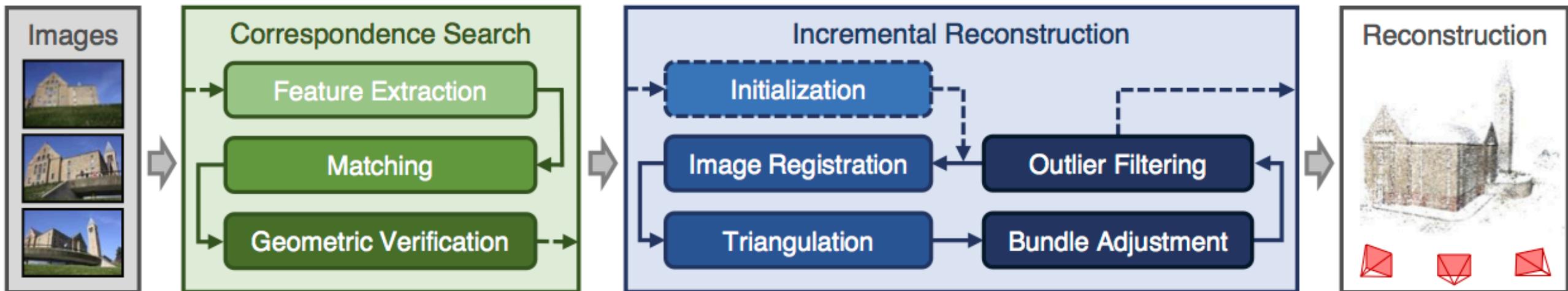
Interactive Computer Graphics and Vision



3D Reconstruction/Recap Lecture

Stefanie Zollmann and Tobias Langlotz

# Colmap Workflow

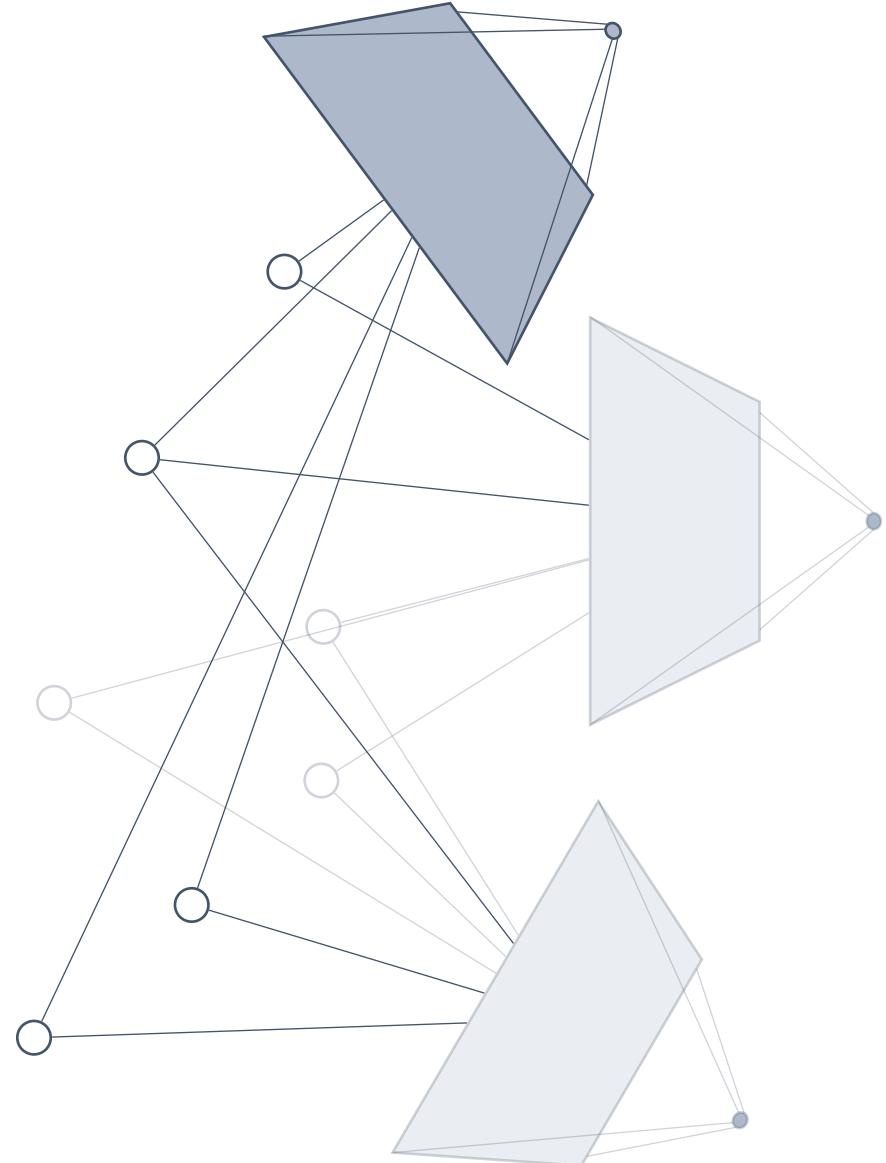


<https://colmap.github.io/tutorial.html>

# **Absolute Pose**

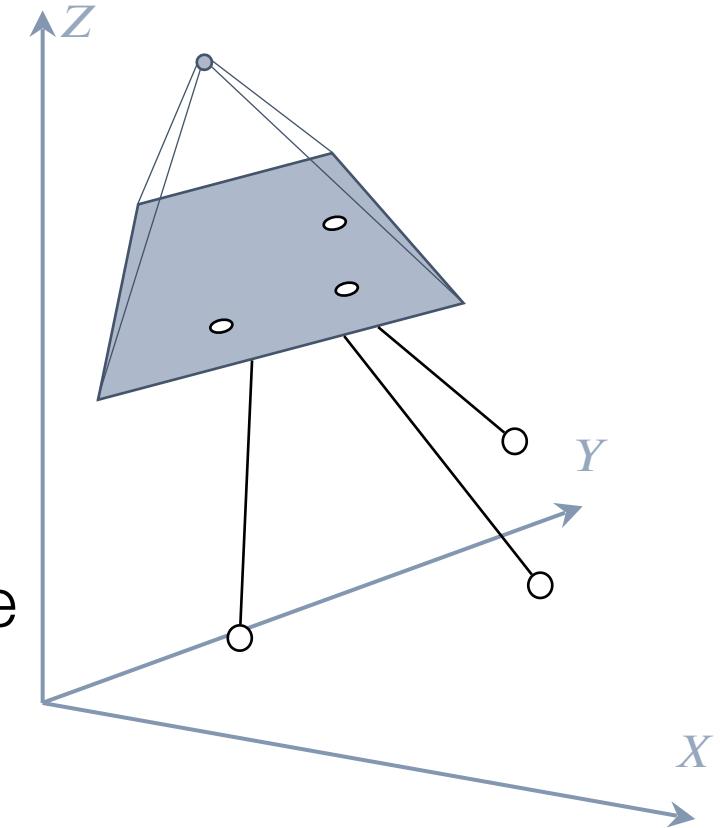
# Multi-View Stereo Pipeline

- Identify best matching pair
  - Two-view stereo for relative pose
  - Reconstruct 3D points
  - Minimise reprojection error
- While still images to add
  - Select next best view
  - *Determine absolute pose*
  - Reconstruct more 3D points
  - Minimise reprojection error



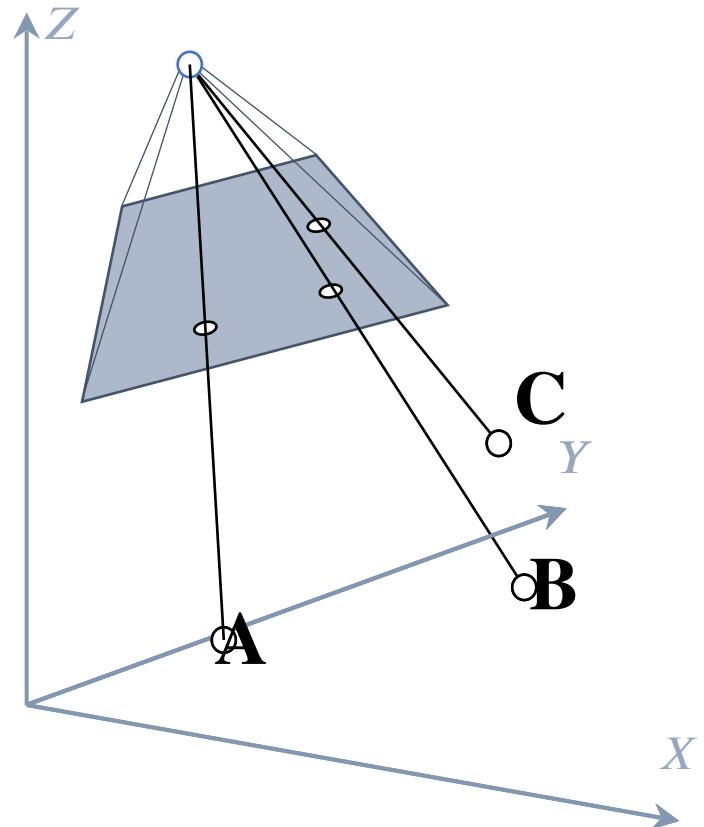
# Absolute Pose Problem

- We have a new image taken with a calibrated camera
- We have already established a coordinate frame
- We have 3D points and their corresponding 2D points from the first image pair
- We found matches from new image to first image pair
- Need to find pose (rotation and translation) of the camera when the new image was taken



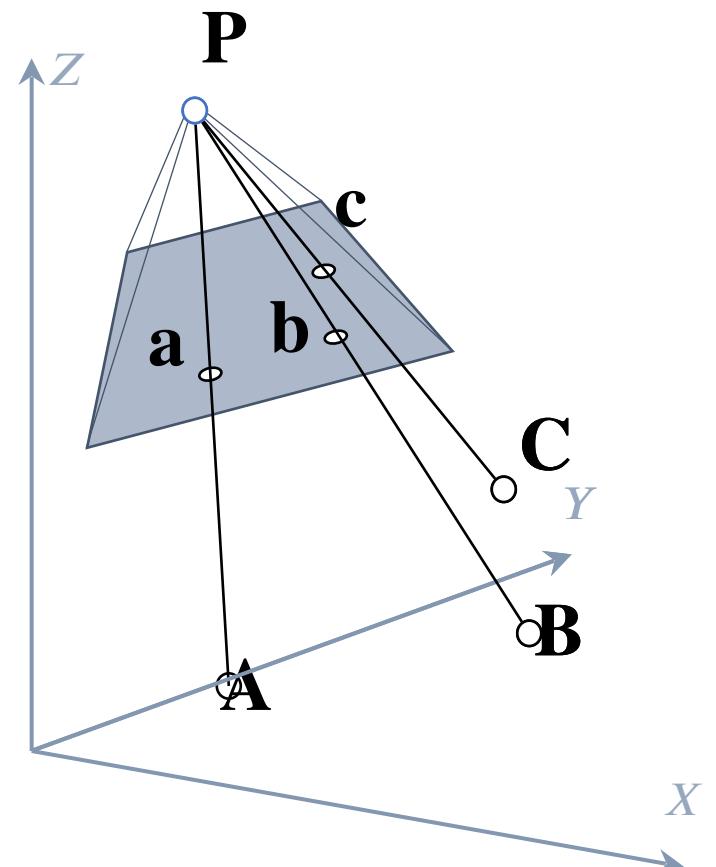
# Perspective- $n$ -Point Pose (PnP)

- Pose from 2D-3D matches
  - Have six unknowns ( $\mathbf{R}$ ,  $\mathbf{t}$ )
  - Each 2D-3D match  $\mathbf{u}_i \leftrightarrow \mathbf{x}_i$  gives
$$k_i \mathbf{u}_i = \mathbf{K}[\mathbf{R} \quad \mathbf{t}] \mathbf{x}_i$$
- How many matches?
  - Each point gives 3 equations
  - Adds 1 unknown ( $k_i$ )
  - $6 + n$  unknowns,  $3n$  equations,  
 $n \geq 3$



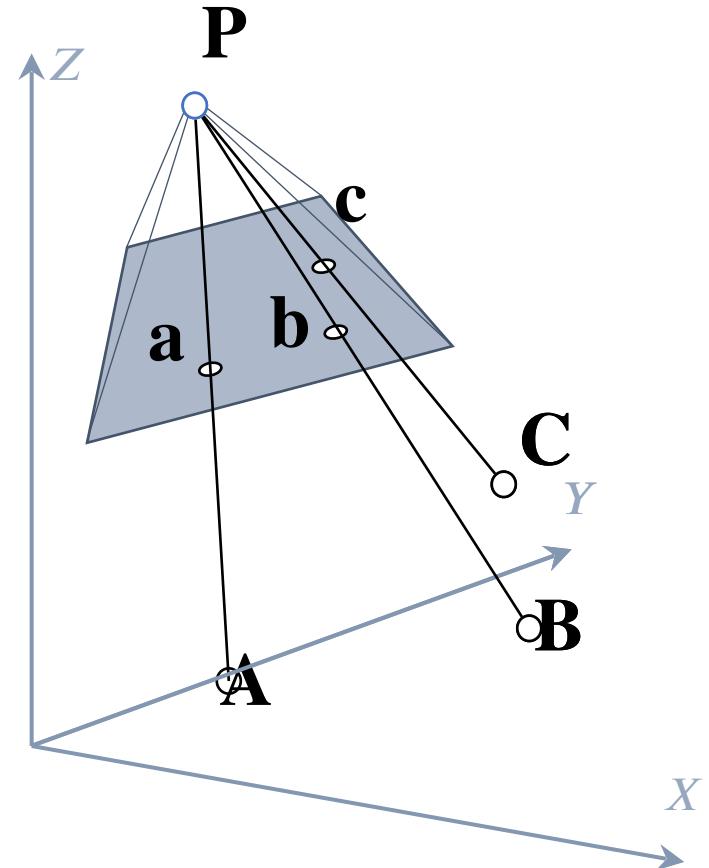
# Perspective-3-Point Pose (P3P)

- Input:
  - The 3D coordinates of three known points in the world (e.g., landmarks)
  - The 2D coordinates of their projections in the camera image
  - The camera's intrinsic calibration (focal length, principal point, etc.)
- Output:
  - The rotation (camera orientation)
  - The translation (camera position)



# Perspective-3-Point Pose (P3P)

- We have three 2D-3D matches
  - Call the 3D points **A**, **B**, and **C**
  - The matching image points are **a**, **b**, **c**
  - The (unknown) camera location is **P**
- Overview of solution
  - Use angles between  $\vec{PA}$ ,  $\vec{PB}$ , and  $\vec{PC}$
  - Find the distances  $|\vec{PA}|$ ,  $|\vec{PB}|$ ,  $|\vec{PC}|$
  - This gives us equations for **P**



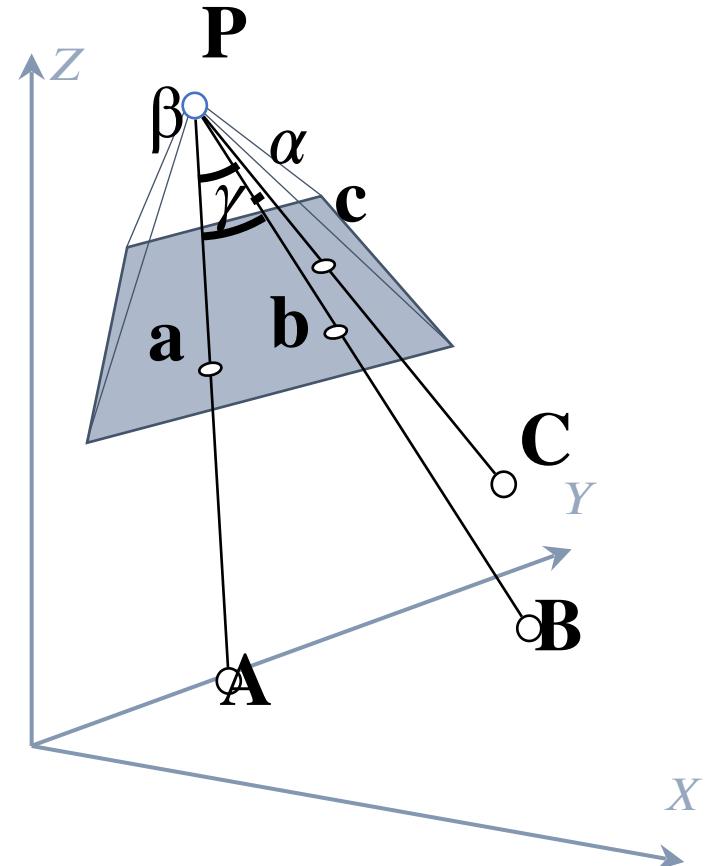
# Perspective-3-Point Pose (P3P)

- Define the angles between rays
- $$\alpha = \angle BPC \quad \beta = \angle APC \quad \gamma = \angle APB$$

- Using dot products:

$$\cos\alpha = \frac{\mathbf{b} \cdot \mathbf{c}}{|\mathbf{b}| |\mathbf{c}|}$$

$$\cos\beta = \frac{\mathbf{a} \cdot \mathbf{c}}{|\mathbf{a}| |\mathbf{c}|} \quad \cos\gamma = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|}$$



# Perspective-3-Point Pose (P3P)

- The law of cosines tells us that

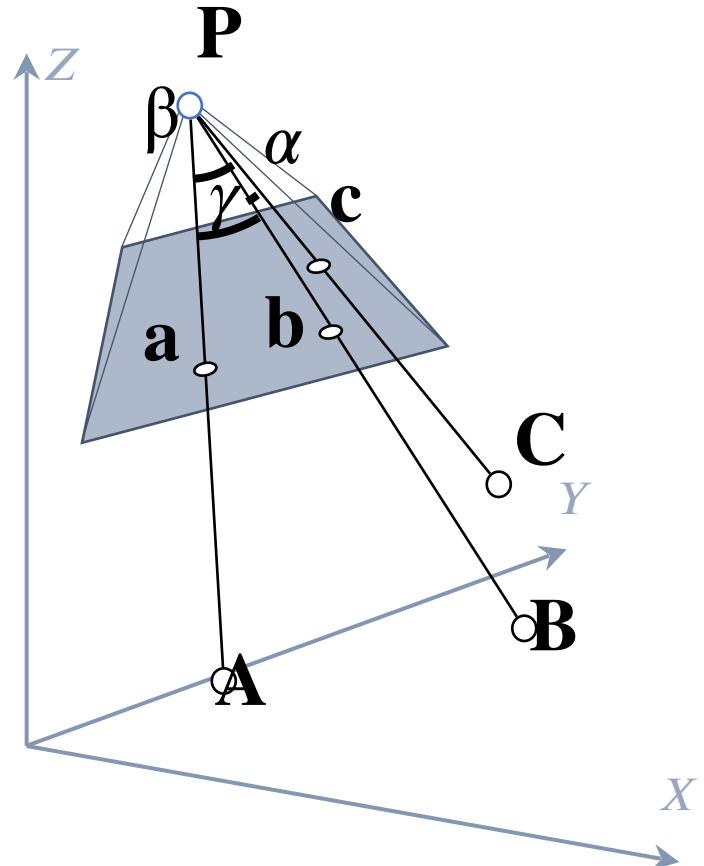
$$|\vec{AB}|^2 = |\vec{PA}|^2 + |\vec{PB}|^2 - 2|\vec{PA}||\vec{PB}|\cos\gamma$$

$$|\vec{AC}|^2 = |\vec{PA}|^2 + |\vec{PC}|^2 - 2|\vec{PA}||\vec{PC}|\cos\beta$$

$$|\vec{BC}|^2 = |\vec{PB}|^2 + |\vec{PC}|^2 - 2|\vec{PB}||\vec{PC}|\cos\alpha$$

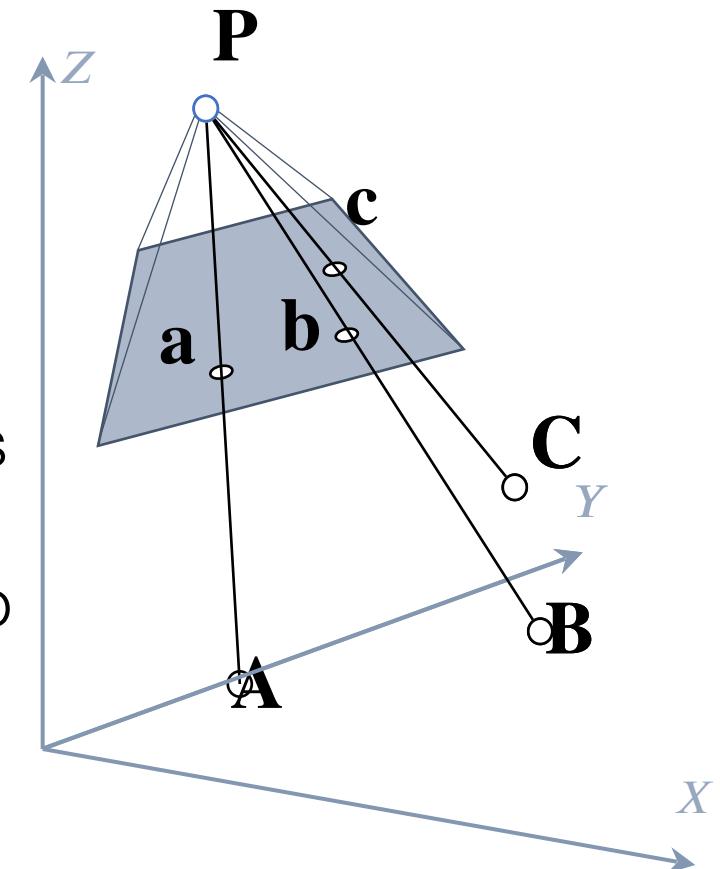
- Solve for

$$|\vec{PA}|, |\vec{PB}|, \text{ and } |\vec{PC}|$$



# Perspective-3-Point Pose (P3P) - Summary

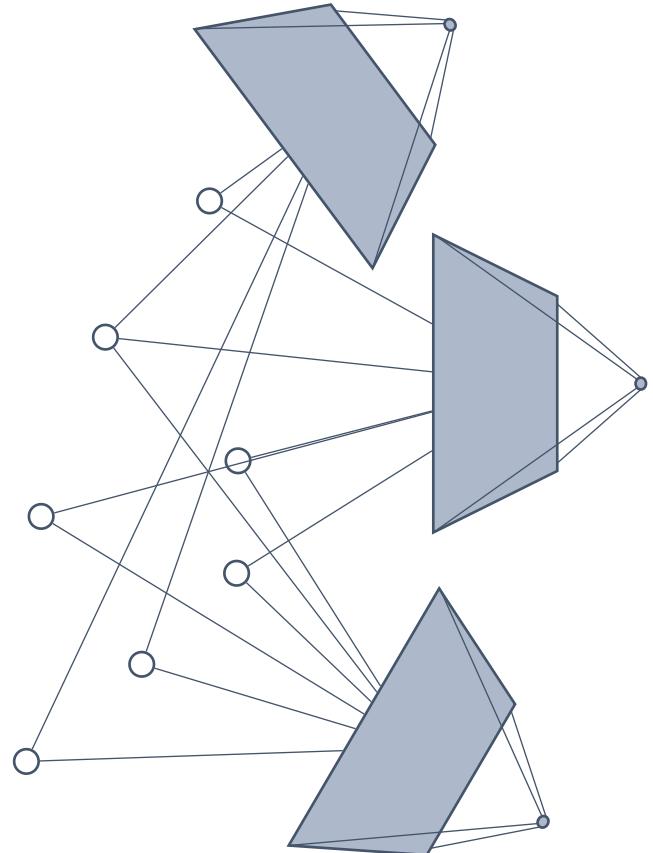
- Input: 3 known 3D world points and their corresponding 2D image points
- Form camera rays: Convert the 2D points into unit direction vectors in the camera frame
- Use the law of cosines: Relate angles between rays (from image) to distances between 3D points (from world) to solve for the distances from the camera to each point
- From these distances, compute the camera center (position) and the rotation that aligns the 3D world points with their camera rays



# **Bundle Adjustment**

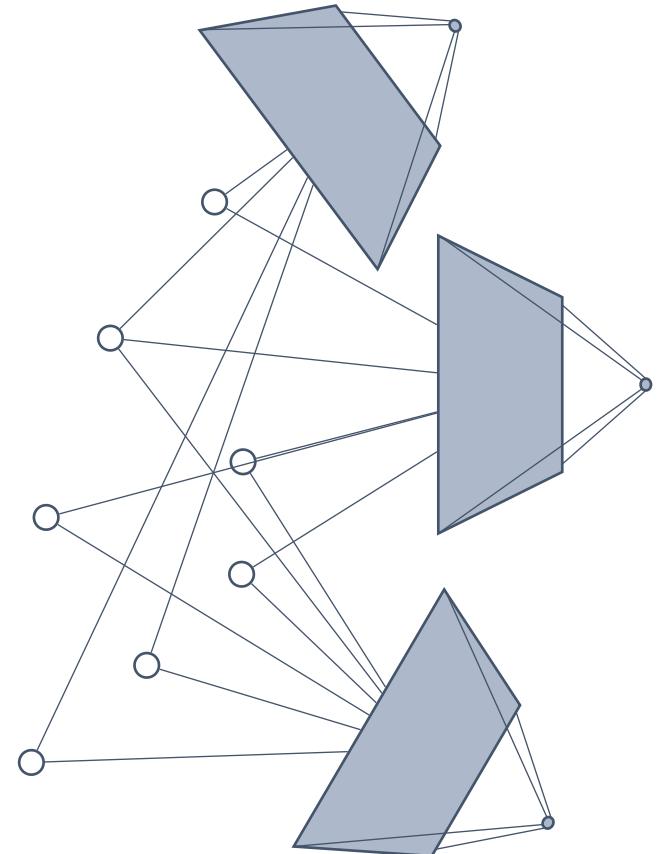
# Multi-View Stereo Pipeline

- Identify best matching pair
  - Two-view stereo for relative pose
  - Reconstruct 3D points
  - Minimise reprojection error
- While still images to add
  - Select next best view
  - Determine absolute pose
  - Reconstruct more 3D points
  - ***Minimise reprojection error***



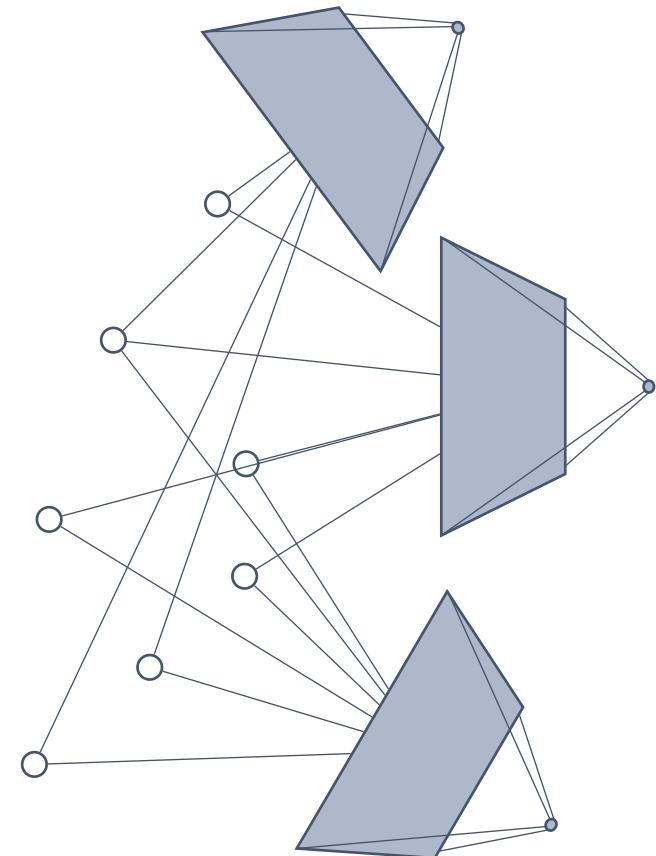
# Bundle Adjustment

- Minimising reprojection error for multi-view stereo
- Jointly refines:
  - Camera parameters (poses + intrinsics)
  - 3D point positions (scene structure)
- Goal: minimize reprojection error across all images
- Geometry forms **bundles** of rays for each camera
- **Adjusting** these gives an optimal reconstruction solution
- Large sparse non-linear least squares
- Levenberg-Marquardt is commonly used for this



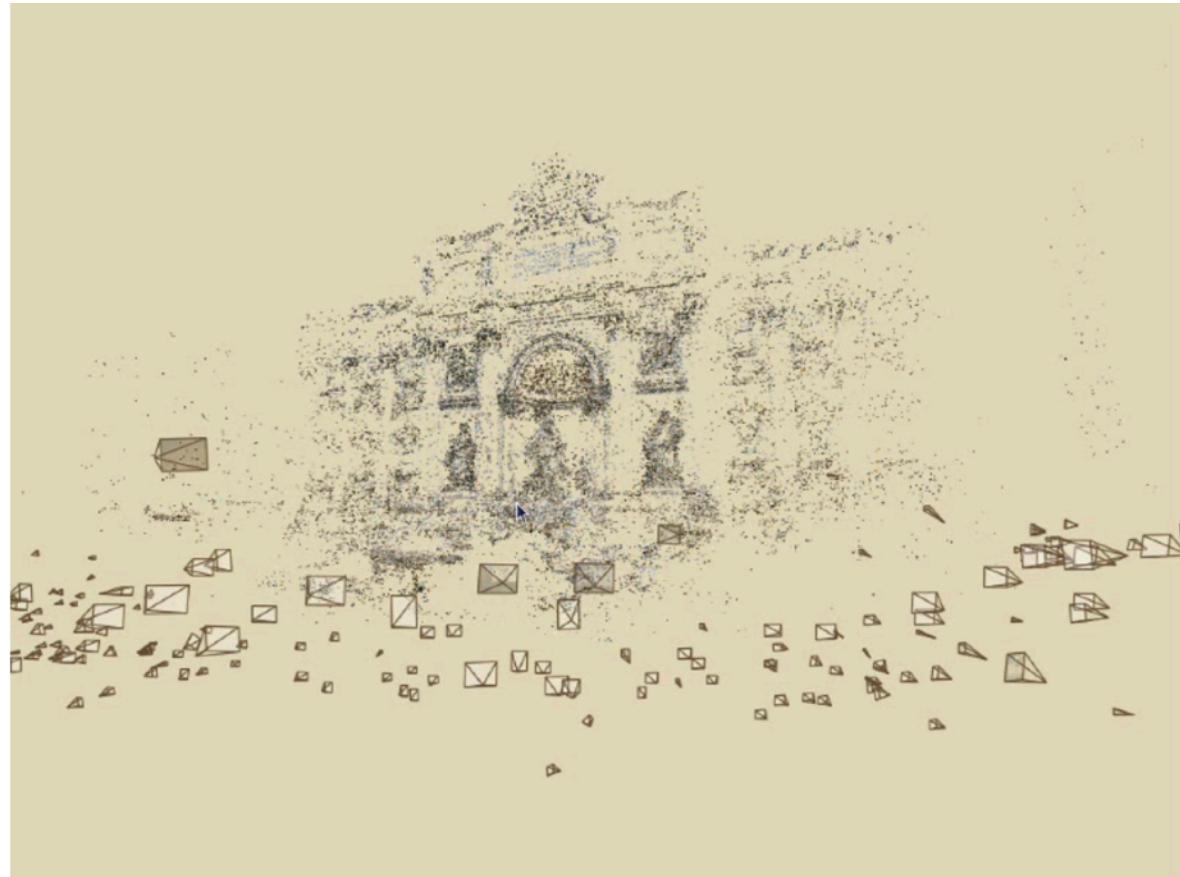
# Reprojection Error

- If we have  $f$  features in  $n$  cameras, reprojection error is:



# Bundle Adjustment

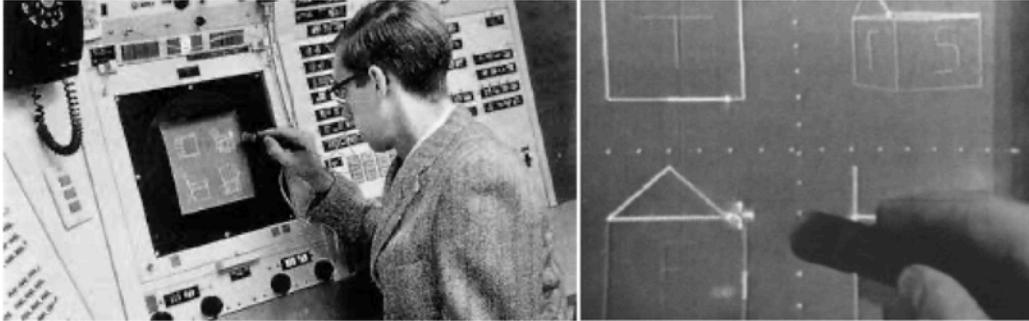
- Why is it Needed?
  - Individual pose or triangulation estimates are noisy
  - Errors accumulate across many views
- Bundle adjustment enforces global consistency



# **Revision**

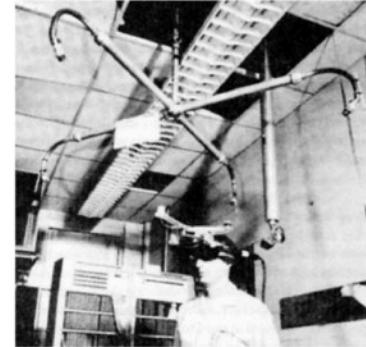
# Introduction

## A Short History of Visual Computing



Ivan Sutherland, *Sketchpad, a Man-Machine Graphical Communication System*, 1963

## A Short History of Visual Computing



Ivan Sutherland, 1968 "A head-mounted three dimensional display"

65

## A Short History of Visual Computing



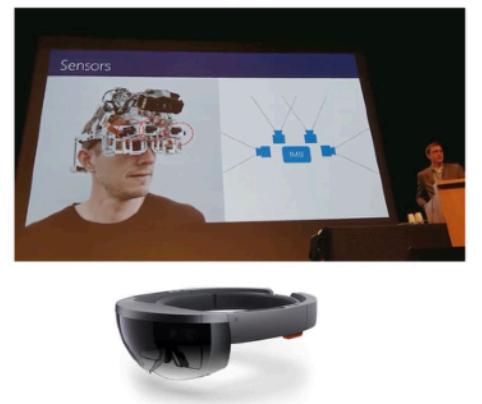
SIFT by David G. Lowe "Object Recognition from Local Scale-Invariant Features" (1999)

## A Short History of Visual Computing



George Klein "Parallel tracking and mapping for small AR workspaces" (PTAM, 2007)

77



79

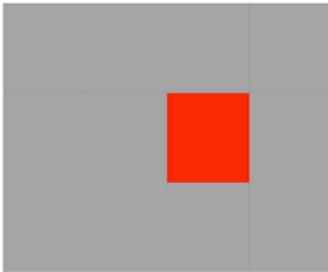
# Digital Images

## Digital Images and 2D Image Representation

- 2D array of pixels – a raster
  - Pixels have a value, or list of values
    - value: usually brightness
    - values: brightness, colour, transparency
- ```
# Create a 2D array
image2D = create_array(height, width)

# Fill with example grayscale values
for y from 0 to height-1:
    for x from 0 to width-1:
        image2D[y][x] = 125

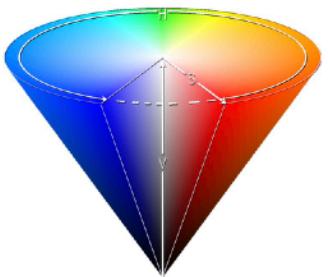
# Access pixel at (row=2, col=3)
image2D[1][2] = ?
```



16

## HSV Colour Space

- How do we describe colour?
  - The sky is a bright blue
  - The trees are a dark green
  - That shirt is a bold colour
- HSV is closer to this:
  - Hue – what general colour is it
  - Saturation – how strong the colour is (distance from grey)
  - Value – how light or dark it is
  - Not perceptually uniform
- Also HSL, L = lightness/luminance
  - from dark to light



17

## Representing Colour

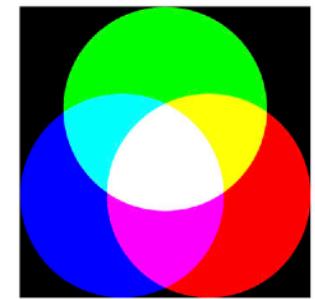
- We perceive 3 colour dimensions
  - Short, medium, long cones
  - Intensity, Red-Green, Blue-Yellow
- Colour representations are usually 3 dimensional
  - Red-Green-Blue -> RGB
  - Cyan-Magenta-Yellow -> CMY (+K)
  - Hue-Saturation-Value ->HSV or HSL
  - And others (e.g. YUV, CIE)



34

## RGB Model

- Most common in computing
  - Most monitors and projectors
  - Additive colour model
  - Not perceptually uniform
- Record/transmit RGB channels
  - Aligns to the types of cone in the eye
  - Mixes three primary light colours
  - Can represent most colours (not all?)
  - Typically within the range [0,255] or [0,1]
    - E.g. Red as (255,0,0) or (1,0,0)
  - Can be abbreviated as hex values:
    - (183, 24, 92) -> #B7185C



35

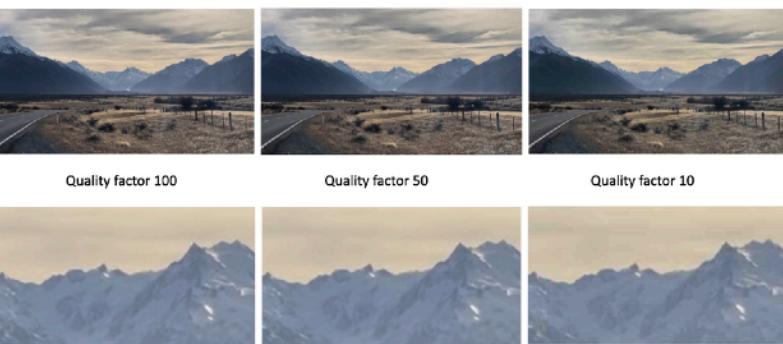
## CMY Colour Space

- Most common for printing
  - Subtractive Colour Model
  - Not perceptually uniform
- CMY is the inverse of RGB
  - Cyan Magenta Yellow (Black)
  - Cyan pigment absorbs red light
  - Magenta pigment absorbs green
  - Yellow pigment absorbs blue
- In theory C+M+Y = Black
  - In practice a 'true Black' ink is used
- Black is represented by K (Key)



43

## JPEG Compression Artifacts



44

# 2D Transforms

## Transformations - Scaling

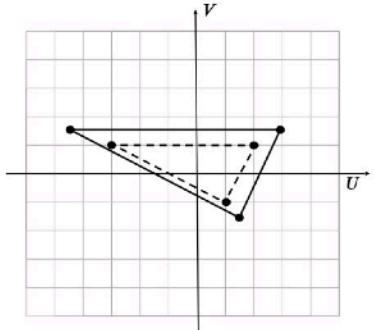
- Scaling points by factor,  $s$   
 $(u, v) \rightarrow (su, sv)$

In vector terms

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = s \begin{bmatrix} u \\ v \end{bmatrix}$$

Example:

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = 1.5 \begin{bmatrix} u \\ v \end{bmatrix}$$



13

## Transformations - Translation

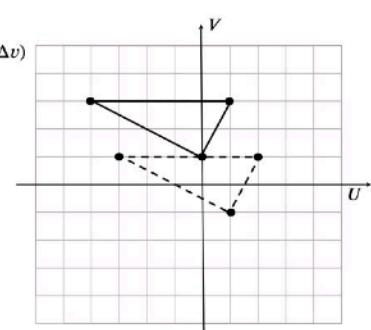
- Shifting a point,  $(u, v)$ , by some offset,  $(\Delta u, \Delta v)$   
 $(u, v) \rightarrow (u + \Delta u, v + \Delta v)$

In vector form:

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix}$$

Example:

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$



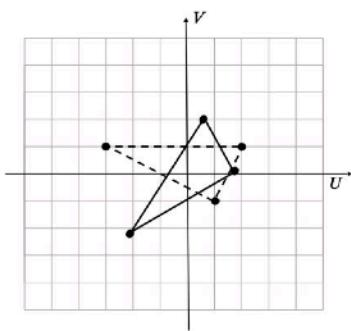
14

## Transformations - Rotation

- Example rotation around 45 degrees

$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$$R\left(\frac{\pi}{4}\right) = \begin{bmatrix} \cos \frac{\pi}{4} & -\sin \frac{\pi}{4} \\ \sin \frac{\pi}{4} & \cos \frac{\pi}{4} \end{bmatrix} R\left(\frac{\pi}{4}\right) = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \approx \begin{bmatrix} 0.7071 \\ 0.7071 \end{bmatrix}$$



15

## Transformations - Combining Transforms

- Rotate  $-45^\circ$  about  $(2,1)$

1. Shift by  $(-2, -1)$

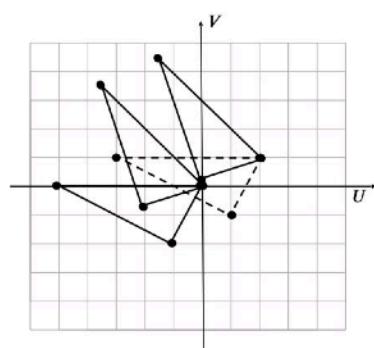
2. Rotate by  $-45^\circ$

3. Shift by  $(2, 1)$

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} -2 \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = R_{-45^\circ} \left( \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} -2 \\ -1 \end{bmatrix} \right)$$

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = R_{-45^\circ} \left( \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} -2 \\ -1 \end{bmatrix} \right) + \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



16

## Homogenous Coordinates

- The solution is counter-intuitive
- 2D points become sets of 3-vectors  
 $\begin{bmatrix} u \\ v \end{bmatrix} \rightarrow k \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}, \forall k \neq 0$
- The vector  $[a \ b \ c]^T$  corresponds to the point  $\left(\frac{a}{c}, \frac{b}{c}\right)$
- These are homogeneous coordinates
- All (linear) transformations become  $3 \times 3$  matrices

## Projective Transforms

- To capture perspective, we generalize to Projective Transforms

Affine = special case where:  $h_{31} = 0, h_{32} = 0, h_{33} = 0$

Properties:

- Lines stay straight
- Parallelism not preserved
- Can map any quadrilateral to any other quadrilateral
- Models perspective projection

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix}$$

$$H = \begin{bmatrix} 2 & 0 & 0 \\ 0.2 & 1 & 0 \\ 0.5 & 0.3 & 1 \end{bmatrix}$$

17

# **Quiz**

# Image Manipulation/Image Filters

## Greyscale



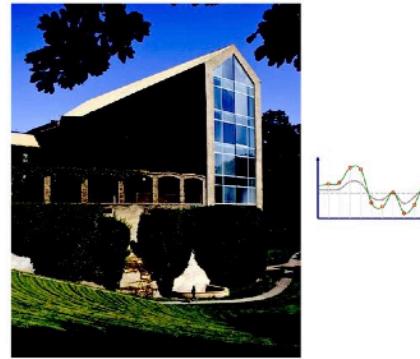
Original



Greyscale  
("black&white" photo)

## Contrast

- Contrast is the difference in luminance or colour that makes an object (or its representation in an image or display) stand out from its background.
- Simpler: Contrast is the difference in luminance or colour for each pixel from the overall mean
- Compute mean luminance  $L^*$  over whole image
- Scale deviation from  $L^*$  for each pixel



15

## Quantization and Dithering

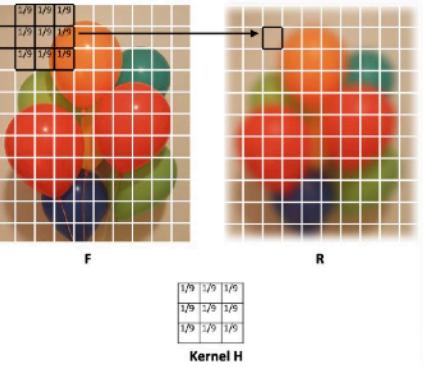
- Quantization:
  - Reduce intensity/colour resolution
    - Greyscale 8bit to 1bit
    - RGB 12 bit to 8 bit
    - RGB 8bit to 5/6 bits (RGB565)
  - Frame buffers have a limited number of bits per pixel
  - Physical devices have a limited dynamic range
  - Reduce image size



33

## Linear Filters and Image Convolutions

- Image filtering with linear filters
- Construct array  $R$  of same size as original image  $F$
- Fill each location in  $R$  with a weighted sum of the pixel values of corresponding location's neighbourhood in  $F$
- Use the same weights each time
- Different sets of weights represent different filters
- **Convolution** with a (filter) *kernel*
- **Shift-invariant:** filtered result depends on image neighbourhood but not on position of neighbourhood
- **Linear:** filtering the sum of two images is equivalent as summing the two filtered images
- **Shift-Invariant Linear Systems (SILS)**



## Image Features & Feature Detection

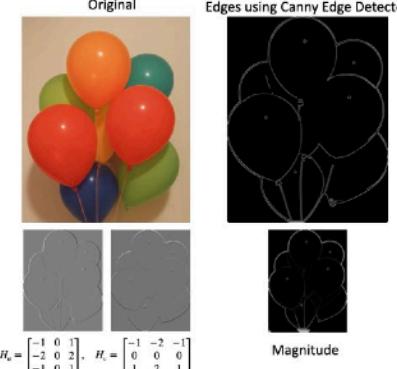
- Often required:
  - Features in an image that can be reliably detected in different images of the same scene
- A 'feature' or 'key point' is
  - A 2D location in an image
  - Accurately/precisely located
  - Can be found in different views
- Feature descriptors
  - Used to tell features apart
  - Ideally robust to changes in lighting, viewpoint, contrast, etc.
- Three common types of features
  - Edges – gradient-based
  - Corners – gradient-based
  - Blobs – bright/dark centre/surround



What are good image features?

## Edge Detection - Canny Edge Detector

- Canny Edge Detector
  - Step 1: Reduce noise (-> Gaussian filter  $G_\sigma$ )
  - Step 2: Compute gradients in x and y direction (-> Sobel operators  $H_x, H_y$ ), gradient is a vector  $\mathbf{g} = [H_x \ H_y]^T$
  - Step 3: compute gradient magnitudes (edge strength) and gradient directions (edge directions)
  - Step 4: select points with maximum magnitudes (above threshold  $T_1$ )
  - Step 5: non-maximum suppression with hysteresis:
    - For all points with magnitudes above high-threshold  $T_1$
    - Trace along the edge (in edge direction) suppress all pixels that are not on traced line (pixel by pixel) and mark as visited
    - Stop tracing when point's magnitude is below low-threshold  $T_2$  or already visited

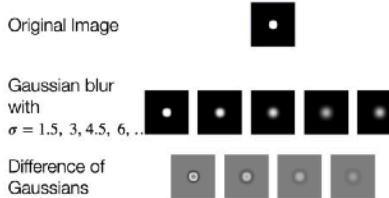


86

# Features

## Blob Detection - Difference of Gaussians (DoG)

- 'Blobs' are alternative features
- A blob is a dark region surrounded by a bright region, or vice-versa
- Blobs have a scale as well as a location

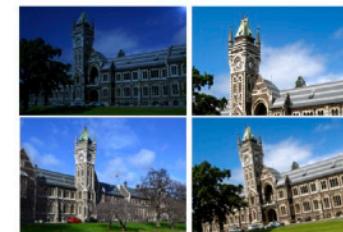


38

## Scale Invariant Features Transform (SIFT)

### Properties of SIFT

- **Robustness:** Invariance to changes in illumination, scale, rotation, affine, perspective.
- **Locality:** robustness to occlusion and clutter.
- **Distinctiveness:** Easy to match to a large database of objects (e.g. different images).
- **Quantity:** Many features (points) can be generated for even small objects.
- **Efficiency:** Computationally "cheap", real-time performance.



Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoints. International journal of computer vision, 60(2), 91-110.

46

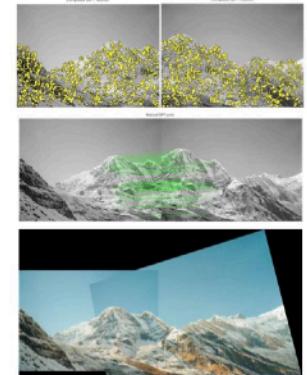
## Scale Invariant Features Transform (SIFT)

### Common problems:

- Panorama stitching
- Image Search
- Object recognition
- Tracking (covered a different time)

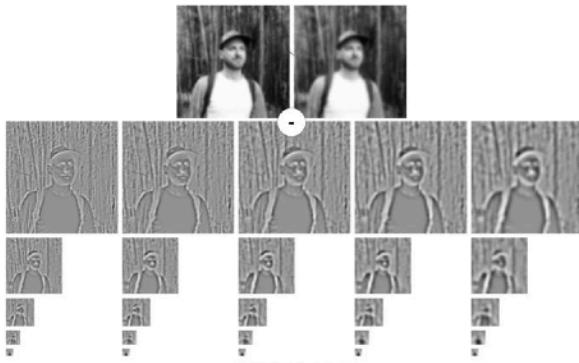
### Scale Invariant Features Transform (SIFT):

- Computer vision algorithm to
  - Detect,
  - Describe,
- And match local features in images



41

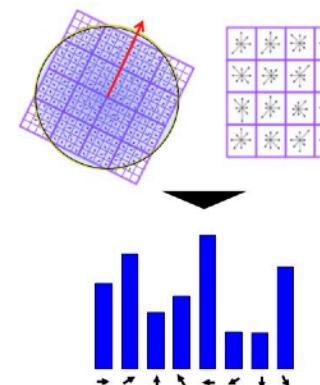
## SIFT - Difference of Gaussians



Otero, I. R. (2015). Anatomy of the SIFT Method (Doctoral dissertation, École normale supérieure de Cachan-ENS Cachan).

## SIFT - Feature description from gradients

- Feature description (orientation and gradient)
- Reorient image patch using main orientation
- Create 16 gradient histograms (8 bins) weighted by magnitude and Gaussian window
- Final keypoint Descriptor - 128 (4x4x8) element vector

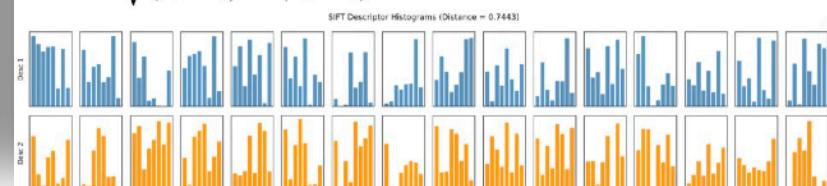


58

## SIFT - Matching

- SIFT descriptor 128 values
  - Usually these are bytes
- Computing the distance between features:

$$d = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots}$$



68

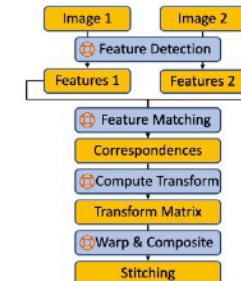
# Images Stitching and Homographies

## How to Create Panoramic Images?



## Image Stitching Process

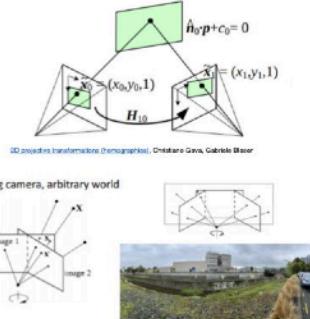
- Four main stages:
  - Detect features in each image
  - Match features between images
  - Estimate a transform
  - Warp one image to the other
- Stitching only works if:
  - The scene is (almost) planar, or
  - The camera only (mostly) rotates



## Homographies in the Stitching Process

- A homography\* is:
  - A linear map between two planes (or views of a plane)
  - A  $3 \times 3$  matrix, up to a scale
- Two images of a scene are related by a homography
  - If the scene is planar, or
  - If the camera only rotates

\* For our purposes – this is one case of a more general concept.



## Homography

$$\begin{bmatrix} u'_i \\ v'_i \\ 1 \end{bmatrix} \equiv \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

The feature location in the second image      Is equivalent to (equal up to a scale)      Some homography (a  $3 \times 3$  matrix)      Applied to...      The feature location in the first image

## RANSAC for Homographies

- Threshold for acceptance
  - Can often be interpreted in terms of the task
  - E.g. for how well does  $H$  align the corresponding points?
  - Set too low -> smaller consensus sets
  - Set too high -> outliers give less accurate models
- Number of trials
  - Resource limits
    - E.g. tracking in 25 fps video, do as many trials as you can in 0.04s
  - Can also be interpreted in a probabilistic manner
    - Replace number of trials with probability of success

## Feathering (Alpha Blending)

- Weighted averaging across overlap
- Pros: simple, fast
- Cons: ghosting if misalignment or parallax



# **Quiz**

# 3D Geometry

## Translation and Scaling in 3D

- Translation by  $(\Delta x, \Delta y, \Delta z)$

$$\mathbf{x}' = \mathbf{T}\mathbf{x}$$

$$\mathbf{T} = \begin{bmatrix} 1 & 0 & 0 & \Delta x \\ 0 & 1 & 0 & \Delta y \\ 0 & 0 & 1 & \Delta z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- Scaling by Scaling by a factor  $s$

$$\mathbf{x}' = \mathbf{S}\mathbf{x}$$

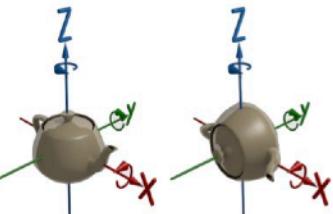
$$\mathbf{S} = \begin{bmatrix} s & 0 & 0 & 0 \\ 0 & s & 0 & 0 \\ 0 & 0 & s & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

## Rotation about the X Axis

- From  $Y$  to  $Z$

- $X$  co-ordinate is unchanged

$$R_X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos(\theta) & -\sin(\theta) & 0 \\ 0 & \sin(\theta) & \cos(\theta) & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$



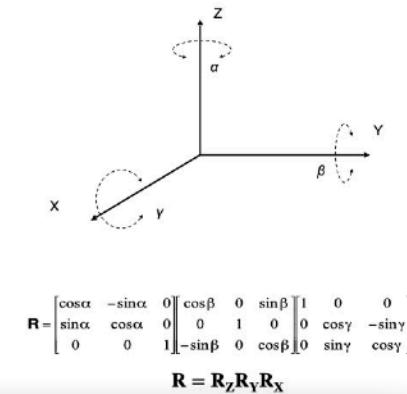
Rotation by 90° about the  $Z$  axis – from  $X$  to  $Y$

## Euler Angles

- Rotation defined by angles of rotation around the  $X$ -,  $Y$ -, and  $Z$ - axes

- Different conventions

- Example Blender Rotation Mode:



$$\mathbf{R} = \mathbf{R}_Z \mathbf{R}_Y \mathbf{R}_X$$

21

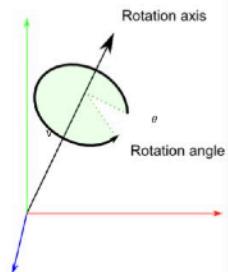
## Quaternion

- Another way of presenting rotations (avoiding Gimbal Lock)
- Can be expressed as a scalar and a 3- vector:  $(a, v)$
- A rotation about the unit vector  $v$  by an angle  $\theta$  can be represented by the unit quaternion

angle in unit-length  
radians axis

$$q(\theta, v) = \left( \cos \frac{\theta}{2}, v_x \sin \frac{\theta}{2}, v_y \sin \frac{\theta}{2}, v_z \sin \frac{\theta}{2} \right)$$

scalar = angle  
vector = axis

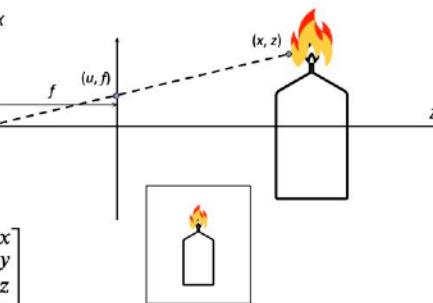


## The Pinhole Camera Model and Projective Geometry

- Removing the sign change

- We can put the image plane in front of the pinhole
- Removes the sign change
- Not practical for real cameras
- The maths works out just fine

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

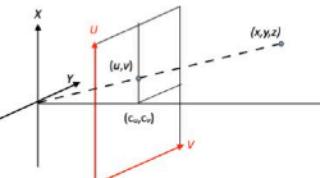


## Camera Parameters - Intrinsic and Extrinsic

- Often break this down into

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

Intrinsic                      Extrinsic



- Most simple case:  $\mathbf{u} = \mathbf{K}[\mathbf{I} \mid \mathbf{0}]\mathbf{x}$

- $\mathbf{K}$ : camera calibration or Intrinsic parameters

- $[\mathbf{I} \mid \mathbf{0}]$ : camera pose or Extrinsic parameters

68

# Stereo Vision/ Epipolar Geometry

## A Simple Stereo System

- Assume a set of 2 pinhole cameras

$$\mathbf{u} = \mathbf{K}[\mathbf{R} - \mathbf{t}]x$$

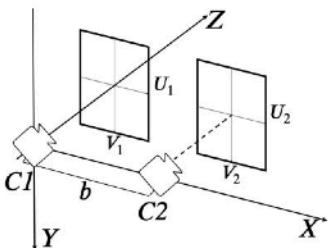
- Simplify for now:  $\mathbf{K} = \mathbf{I}$

- Camera 1 (C1):  $\mathbf{R} = \mathbf{I}$ ,  $\mathbf{t} = \mathbf{0}$

- No rotation, looks along  $Z$

- Camera 2 (C2): shifted along  $X$

$$\mathbf{R} = \mathbf{I}, \mathbf{t} = \begin{bmatrix} b \\ 0 \\ 0 \end{bmatrix}$$



## Project a Point into the Cameras

- Camera 1:

$$\begin{bmatrix} u_1 \\ v_1 \\ 1 \end{bmatrix} \equiv \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \equiv \begin{bmatrix} x/z \\ y/z \\ 1 \end{bmatrix}$$

- Camera 2:

$$\begin{bmatrix} u_2 \\ v_2 \\ 1 \end{bmatrix} \equiv \begin{bmatrix} 1 & 0 & 0 & -b \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x - b \\ y \\ z \\ 1 \end{bmatrix} \equiv \begin{bmatrix} (x - b)/z \\ y/z \\ 1 \end{bmatrix}$$

## Depth and Disparity

- Difference between the views is the disparity

- In this case it is in  $u$ :

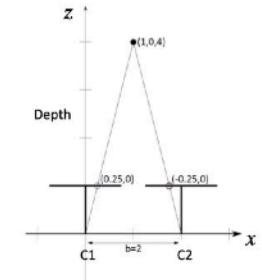
$$u_1 - u_2 = \frac{x}{z} - \frac{x - b}{z} = \frac{b}{z}$$

- The  $v$  values are the same

- If we know  $b$ , can find  $z$ :

$$z = \frac{b}{u_1 - u_2}$$

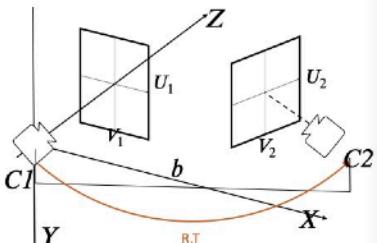
Top Down View



36

## More General Stereo

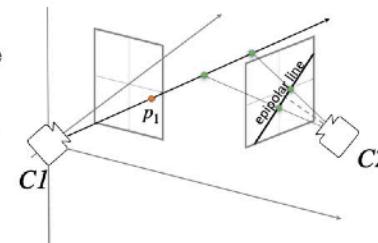
- In the more general case:
  - The cameras have different calibration parameters
  - There is a rotation between the cameras
  - There is translation along all three axes
- We can still align the XYZ coordinate system with camera C1



## Epipolar Geometry

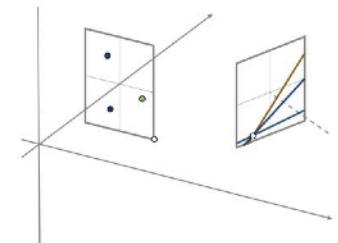
- Consider a point  $p_1$  in the first camera

- This gives a 3D line, along which the point must lie
- The 3D line projects to a 2D line in the second image – this is called an epipolar line



## Epipoles and Epipolar line

- Each point in one image makes an epipolar line
- These all meet at the epipole in the other image
- For simple (parallel) case
  - Line through camera centres is the X axis
  - Epipoles are at infinity
  - Epipolar lines are the rows

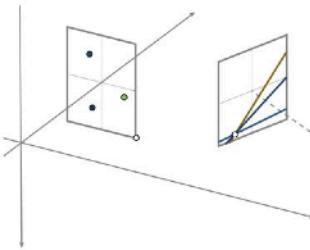


49

# Stereo Vision: Fundamental Matrix

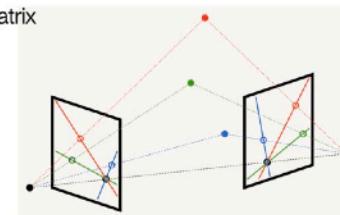
## The Fundamental Matrix

- In general: a point,  $\mathbf{p}$ , lies on a line,  $\mathbf{l}$ , if  $\mathbf{p}^T \mathbf{l} = 0$
- In our case:  $\mathbf{p}_2^T \mathbf{l}_2 = 0$
- If we now use our epipolar constraint:
  - $\mathbf{l}_2 = \mathbf{F} \mathbf{p}_1$
  - Then  $\mathbf{p}_2^T \mathbf{F} \mathbf{p}_1 = 0$



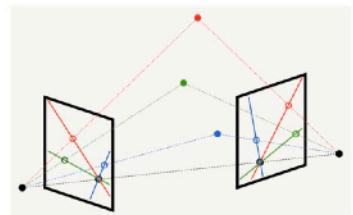
## The Essential Matrix

- For calibrated cameras
  - We can factor out  $\mathbf{K}_1$  and  $\mathbf{K}_2$  in  $\mathbf{F} = \mathbf{K}_2^{-T} [\mathbf{t}]_x \mathbf{R} \mathbf{K}_1^{-1}$ 
    - This gives us the essential matrix  $\mathbf{E} = \mathbf{K}_2^T \mathbf{F} \mathbf{K}_1 = [\mathbf{t}]_x \mathbf{R}$
- $\mathbf{E}$  has only five DoF:
  - 3 for 3D rotation
  - 3 for 3D translation
  - Less 1 for unknown scale



## The 8-Point Algorithm

- Estimate  $\mathbf{F}$  from matches
  - $\mathbf{F}$  is a  $3 \times 3$  matrix – 9 values
  - Defined up to a scale – 8 DoF
  - Need 8 constraints
- $\mathbf{F}$  turns out to be rank 2
  - Only 7 degrees of freedom
  - There is a 7-point algorithm



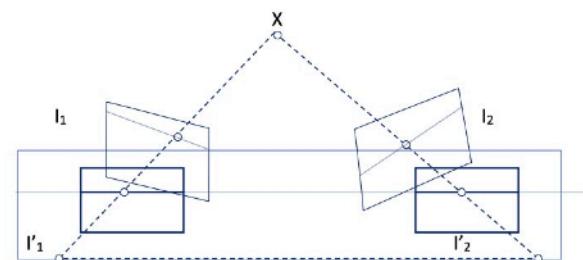
## The Homogeneous Linear System

$$\begin{bmatrix} u_1u'_1 & v_1u'_1 & u'_1 & u_1v'_1 & v_1v'_1 & v'_1 & u_1 & v_1 & 1 \\ u_2u'_2 & v_2u'_2 & u'_2 & u_2v'_2 & v_2v'_2 & v'_2 & u_2 & v_2 & 1 \\ u_3u'_3 & v_3u'_3 & u'_3 & u_3v'_3 & v_3v'_3 & v'_3 & u_3 & v_3 & 1 \\ \vdots & \vdots \\ u_nu'_n & v_nv'_n & u'_n & u_nv'_n & v_nv'_n & v'_n & u_n & v_n & 1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$
$$Af = 0$$

## The 7-Point Algorithm

- Exploiting the geometric constraint that the determinant of  $\mathbf{F}$  must be zero
- The "Two-Solution" Problem: With only 7 points, the system is underdetermined, resulting in two initial 'basis' solutions, let's call them  $\mathbf{F}_1$  and  $\mathbf{F}_2$
- Finding the Right Mix: The true solution  $\mathbf{F}$  is a specific blend of these two basis solutions:  $\mathbf{F} = \alpha \mathbf{F}_1 + (1 - \alpha) \mathbf{F}_2$
- The Cubic Equation: By enforcing the constraint that  $\det(\mathbf{F})=0$ , we get a cubic polynomial for the mixing value  $\alpha$
- Pro: Requires fewer points
- Con: More complex due to solving a cubic polynomial; can have multiple solutions; less robust to noise than the 8-point algorithm with normalisation and RANSAC

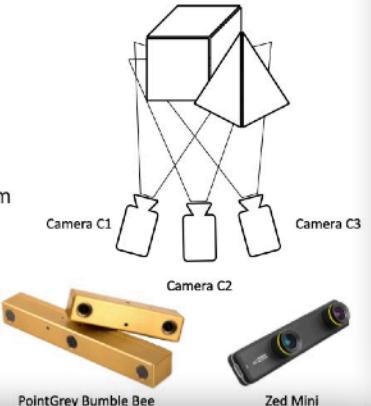
## Stereo Rectification



# Depth Sensing

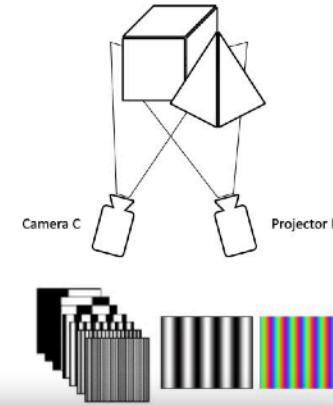
## Passive Range Scanning

- Two or more cameras
  - Stereo camera
  - Two cameras
  - One camera takes several photos from different perspectives
  - Multiple cameras
- Which one is the easiest and why?
- What is the main challenge?

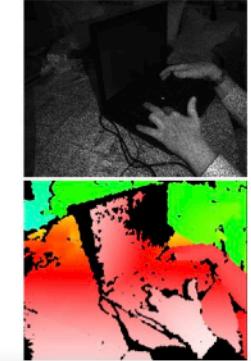


## Structured Light Projection

- Many structured light techniques exist that follow the two objectives:
  - Determine projector-camera correspondences with as few images as possible
  - Be robust against surface modulation and noise
- Codes include:
  - Binary codes, such as the Gray code (most common): <100 images, quite robust
  - Intensity codes that apply cosine patterns and phase shifts: 6 images, but are not robust against strong absorptions
  - Color coding: <3 images, but very fragile (does not work on coloured surfaces)



## Structured Light Projection



MS Kinect 1 (PrimeSense sensor, not MS Kinect 2 uses ToF discussed later)

42

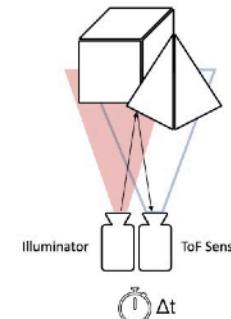
## Structured Light Projection



Raskat et al. "Shaderlamps"

## Optical Time-of-Flight (ToF)

- General idea:
  - Compute depth by round-trip estimation of a light wave emitted and its reflection back to the sensor
- $D = C \cdot \Delta t / 2$
- $\Delta t = 2D/C$
- $D=1m \Delta t = 6.6 \cdot 10^{-9} s$
- $D=0.01m \Delta t = 6.6 \cdot 10^{-11} s$
- $D=0.001m \Delta t = 6.6 \cdot 10^{-12} s$  (pico seconds)
- 3.3 picoseconds for light to travel 1 millimeter



## Depth from Mono / Depth Estimation using ML

- Convolutional neural networks (CNNs) have been applied to monocular depth estimation
- Common architectures include U-Net, ResNet, and DenseNet
- During training:
  - Network is given pairs of images and their corresponding ground-truth depth maps
- Minimize the difference between predicted depth and ground truth depth

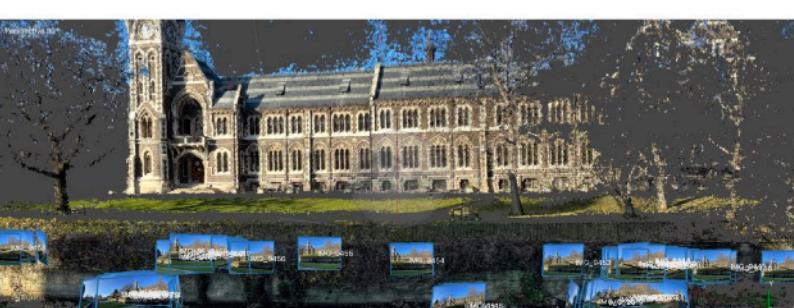


<https://github.com/yuhsuanyeh/BiFuse>

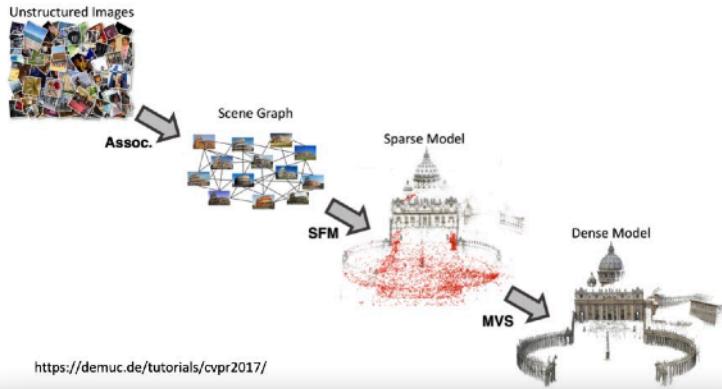
71

# Multiview Geometry

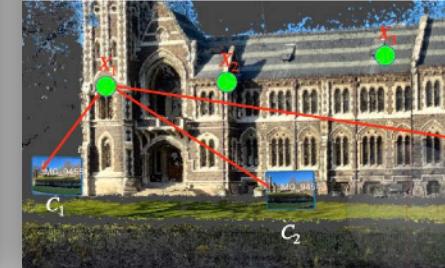
## Goal



## Workflow



## Structure-from-Motion

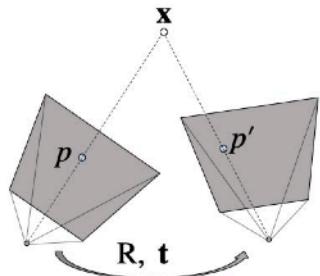


- Estimation of 3D points  $X_i$  and camera poses  $C_j$
- Based on motion (camera views from different viewpoints)

43

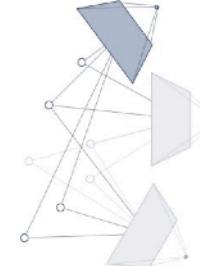
## Decomposing the Essential Matrix

- To recover  $R$  and  $t$ 
  - Compute SVD of  $E$
  - SVD: gives us three matrices  $U$   $S$  and  $V$  transpose:  $E = USV^T$
- From this decomposition, we get:
  - Two possible rotations ( $R$ )
  - Two possible translations ( $t$ )
- That means there are four candidate solutions in total.
- To find the correct one: check which solution places the 3D points in front of both cameras



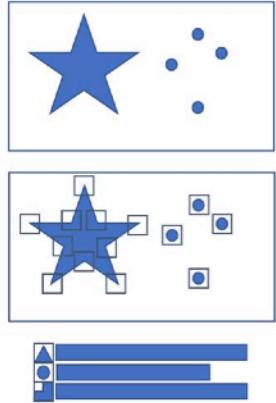
## Multi-View Stereo Pipeline

- Identify best matching pair
  - Two-view stereo for relative pose
  - Reconstruct 3D points
  - Minimise reprojection error
- While still images to add
  - Select next best view
  - Determine absolute pose
  - Reconstruct more 3D points
  - Minimise reprojection error



## Bags of Words (BoW)

- Need equivalent of words
- Cluster features – can pre-train
  - Choose some number,  $k$ , of ‘words’
  - Make  $k$  clusters – e.g.  $k$ -means
  - Cluster centres become words
- Bag of words for images
  - Detect and describe feature
  - Count features closest to each word



63

# **Quiz**

**The end!**