

Answer to Question 4 回答问题 4

(a)

(i) Expressing the cross-ratio between four colinear points  $a, b, c, d$ .

(i)表示四个共线点之间的交比  $a, b, c, d$  。

The cross-ratio of four colinear points  $a, b, c, d$  is defined as:  
四个共线点的交比  $a, b, c, d$  定义为:

$$\text{Cross-ratio} = \frac{(a - c)(b - d)}{(a - d)(b - c)}$$

Alternatively, in terms of distances along the line (assuming points are in order and distances are positive):或者，就沿线的距离而言（假设点有序且距离为正）：

$$\text{Cross-ratio} = \frac{|ac| \cdot |bd|}{|ad| \cdot |bc|}$$

(ii) Proving that the cross-ratio is projective invariant.(ii)证明交比是射影不变的。

**Proof:**

Under a projective transformation, points on a line are transformed via a linear fractional (homography) function:在射影变换下，线上的点通过线性分数（单应性）函数进行变换：

$$x' = \frac{ax + b}{cx + d}$$

For four colinear points  $x_1, x_2, x_3, x_4$ , the cross-ratio is:

对于四个共线点  $x_1, x_2, x_3, x_4$ ，交比为：

$$\text{CR}(x_1, x_2; x_3, x_4) = \frac{(x_1 - x_3)(x_2 - x_4)}{(x_1 - x_4)(x_2 - x_3)}$$

After applying the projective transformation:应用射影变换后：

$$x'_i = \frac{ax_i + b}{cx_i + d}, \quad \text{for } i = 1, 2, 3, 4$$

The cross-ratio of the transformed points is:变换后点的交比为：

$$\text{CR}(x'_1, x'_2; x'_3, x'_4) = \frac{(x'_1 - x'_3)(x'_2 - x'_4)}{(x'_1 - x'_4)(x'_2 - x'_3)}$$

After algebraic manipulation (omitted for brevity), it can be shown that:

经过代数运算（为简洁而省略），可以证明：

$$\text{CR}(x'_1, x'_2; x'_3, x'_4) = \text{CR}(x_1, x_2; x_3, x_4)$$

**Conclusion:**

The cross-ratio is invariant under projective transformations, including perspective projections. Therefore, the cross-ratio of the projected points  $A, B, C, D$  is equal to that of the original points  $a, b, c, d$ .

交比在投影变换（包括透视投影）下保持不变。因此，投影点的交比  $A, B, C, D$  等于原始点的值  $a, b, c, d$ 。

(iii) Calculating  $|bc|$  and discussing advantages and disadvantages of using cross-ratio.

(三)计算  $|bc|$  并讨论使用交叉比率的优点和缺点。

Given:

- Image distances** (in pixels):**图像距离**（以像素为单位）：
  - $|AB| = 20$  px
  - $|BC| = 40$  px
  - $|CD| = 60$  px
- Real-world distances** (in meters):**真实世界距离**（以米为单位）：
  - $|ab| = 6$  m
  - $|cd| = 8$  m
  - $|bc| = s$  m (unknown)

**Step 1: Compute positions along the line in the image.步骤 1：计算图像中沿线的位置。**

Let's assign coordinates along the line in the image:让我们沿着图像中的线指定坐标：

- $x_A = 0$
- $x_B = x_A + |AB| = 0 + 20 = 20$
- $x_C = x_B + |BC| = 20 + 40 = 60$
- $x_D = x_C + |CD| = 60 + 60 = 120$

**Step 2: Compute the cross-ratio from the image points.步骤 2：计算图像点的交比。**

$$\text{CR}_{\text{image}} = \frac{(x_A - x_C)(x_B - x_D)}{(x_A - x_D)(x_B - x_C)} = \frac{(-60)(-100)}{(-120)(-40)} = \frac{6000}{4800} = \frac{5}{4}$$

**Step 3: Assign positions in the real world.第 3 步：分配现实世界中的位置。**

- $x_a = 0$
- $x_b = x_a + |ab| = 0 + 6 = 6$
- $x_c = x_b + s = 6 + s$
- $x_d = x_c + |cd| = 6 + s + 8 = 14 + s$

**Step 4: Compute the cross-ratio from the real-world points.**

**步骤 4：计算现实世界点的交叉比。**

$$\text{CR}_{\text{real}} = \frac{(x_a - x_c)(x_b - x_d)}{(x_a - x_d)(x_b - x_c)} = \frac{(-(6 + s))(- (8 + s))}{(- (14 + s))(-s)} = \frac{(6 + s)(8 + s)}{(14 + s)(s)}$$

**Step 5: Set the cross-ratios equal and solve for  $s$ .第 5 步：设置交叉比率相等并求解  $s$ 。**

$$\frac{(6 + s)(8 + s)}{(14 + s)(s)} = \frac{5}{4}$$

Multiply both sides by  $4(14 + s)(s)$ :两边同时乘以  $4(14 + s)(s)$ ：

$$4(6 + s)(8 + s) = 5(14 + s)(s)$$

Simplify:

$$4(s^2 + 14s + 48) = 5(s^2 + 14s)$$
$$4s^2 + 56s + 192 = 5s^2 + 70s$$

Bring all terms to one side:将所有项移至一侧：

$$5s^2 + 70s - 4s^2 - 56s - 192 = 0$$

Simplify:

$$s^2 + 14s - 192 = 0$$

**Step 6: Solve the quadratic equation.第6步：求解二次方程。**

$$s = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-14 \pm \sqrt{14^2 - 4(1)(-192)}}{2}$$

Compute the discriminant:计算判别式：

$$\sqrt{196 + 768} = \sqrt{964} \approx 31.03$$

Compute  $s$ :

$$s = \frac{-14 \pm 31.03}{2}$$

Positive solution (since distance cannot be negative):正解（因为距离不能为负）：

$$s = \frac{-14 + 31.03}{2} \approx \frac{17.03}{2} \approx 8.52 \text{ m}$$

**Answer:**

The distance between points  $b$  and  $c$  is approximately **8.52 meters**.

点之间的距离  $b$  和  $c$  约为**8.52米**。

Advantages of using cross-ratio:使用交叉比率的优点：

- Projective Invariance:** Allows measurement without knowledge of camera parameters.  
**投影不变性：**无需了解相机参数即可进行测量。
- Flexibility:** Can compute distances using only image measurements and some known real-world distances.  
**灵活性：**可以仅使用图像测量值和一些已知的现实世界距离来计算距离。
- No Calibration Needed:** Useful when camera calibration is not available.  
**无需校准：**当相机校准不可用时很有用。

Disadvantages of using cross-ratio:使用交叉比率的缺点：

- Requires Four Colinear Points:** Must have at least four points in a straight line.  
**需要四个共线点：**必须至少有四个点在一条直线上。
- Sensitivity to Measurement Errors:** Small errors in image measurements can lead to significant errors in calculated distances.  
**对测量误差的敏感性：**图像测量中的小误差可能会导致计算距离中的重大误差。
- Complex Calculations:** Involves solving non-linear equations, which can be computationally intensive.  
**复杂计算：**涉及求解非线性方程，这可能是计算密集型的。
- Assumption of Perfect Geometry:** Assumes points are perfectly colinear and accurately measured.  
**完美几何的假设：**假设点完全共线并且测量准确。

(b)

(i) Formulating the depth  $Z$  using the disparity equation.(i)制定深度  $Z$  使用视差方程。

The depth  $Z$  of an object in stereo vision is given by:

深度  $Z$  立体视觉中物体的形状由下式给出：

$$Z = \frac{f \cdot B}{d}$$

Where:

- $f$  is the focal length. $f$  是焦距。
- $B$  is the baseline distance between the two camera centers.  
 $B$  是两个相机中心之间的基线距离。
- $d$  is the disparity between corresponding points in the stereo images.  
 $d$  是立体图像中对应点之间的视差。

(ii) Deriving the estimated depth error  $\Delta Z$  in terms of  $\Delta d$  and  $Z$ .

(ii)推导估计深度误差  $\Delta Z$  按照  $\Delta d$  和  $Z$  。

**Step 1: Express the actual depth with disparity error.步骤1：用视差误差表达实际深度。**

Let the actual disparity be  $d + \Delta d$ , so the actual depth is:

设实际差距为  $d + \Delta d$ ，所以实际深度为：

$$Z_{\text{actual}} = \frac{f \cdot B}{d + \Delta d}$$

**Step 2: Compute the depth error  $\Delta Z$ .步骤 2：计算深度误差  $\Delta Z$ 。**

$$\Delta Z = Z_{\text{actual}} - Z = \frac{f \cdot B}{d + \Delta d} - \frac{f \cdot B}{d}$$

**Step 3: Simplify using a common denominator.步骤 3：使用公分母进行简化。**

$$\Delta Z = f \cdot B \left( \frac{d - (d + \Delta d)}{d(d + \Delta d)} \right) = -\frac{f \cdot B \cdot \Delta d}{d(d + \Delta d)}$$

**Step 4: Apply Taylor series approximation (assuming  $\Delta d$  is small).**

**步骤 4：应用泰勒级数近似（假设  $\Delta d$  很小）。**

Since  $\Delta d$  is small relative to  $d$ , we approximate  $d + \Delta d \approx d$ :

自从  $\Delta d$  相对于较小  $d$ ，我们近似  $d + \Delta d \approx d$ ：

$$\Delta Z \approx -\frac{f \cdot B \cdot \Delta d}{d^2}$$

**Step 5: Relate  $f \cdot B$  to  $Z$  and  $d$ .第五步：关联  $f \cdot B$  到  $Z$  和  $d$ 。**

From the original depth equation:由原来的深度方程可得：

$$Z = \frac{f \cdot B}{d} \implies f \cdot B = Z \cdot d$$

**Step 6: Substitute back into the error expression.步骤 6：代入错误表达式。**

$$\Delta Z \approx -\frac{(Z \cdot d) \cdot \Delta d}{d^2} = -\frac{Z \cdot \Delta d}{d}$$

**Answer:**

The estimated depth error is:估计的深度误差为：

$$\Delta Z \approx -\frac{Z \cdot \Delta d}{d}$$

**Observation:**

- Depth Error Proportionality:** The depth error  $\Delta Z$  is directly proportional to the estimated depth  $Z$  and the disparity error  $\Delta d$ .  
**深度误差比例：**深度误差  $\Delta Z$  与估计深度成正比  $Z$  和视差误差  $\Delta d$ 。
- Inverse Relation with Disparity:** As disparity  $d$  decreases (object is farther away),  $\Delta Z$  increases, indicating greater sensitivity to disparity errors at greater depths.  
**与视差成反比关系：**作为视差  $d$  减小（物体距离较远）， $\Delta Z$  增加，表明在更深的深度对视差误差更敏感。
- Implication:** Small errors in disparity measurement can lead to significant errors in depth estimation, especially for distant objects.  
**含义：**视差测量中的小误差可能会导致深度估计中的重大误差，尤其是对于远处的物体。

(c)

**Key difference between training previous models (CNNs, ViT) and Large Language-Vision Models (LLVMs):**

**训练以前的模型（CNN、ViT）和大型语言视觉模型（LLVM）之间的主要区别：**

Previous models like Convolutional Neural Networks (CNNs) and Vision Transformers (ViT) are typically trained solely on visual data to learn spatial features from images. In contrast, Large Language-Vision Models (LLVMs) are trained on both visual and textual data simultaneously. LLVMs learn joint representations by aligning images with corresponding textual descriptions, enabling them to understand and generate multimodal content. This multimodal training allows LLVMs to leverage contextual information from language, improving performance in tasks like image recognition, captioning, and cross-modal retrieval.

以前的模型，如卷积神经网络（CNN）和视觉变换器（ViT），通常仅基于视觉数据进行训练，以从图像中学习空间特征。相比之下，大型语言视觉模型（LLVM）同时针对视觉和文本数据进行训练。LLVM 通过将图像与相应的文本描述对齐来学习联合表示，使它们能够理解和生成多模态内容。这种多模态训练使 LLVM 能够利用语言中的上下文信息，提高图像识别、字幕和跨模态检索等任务的性能。