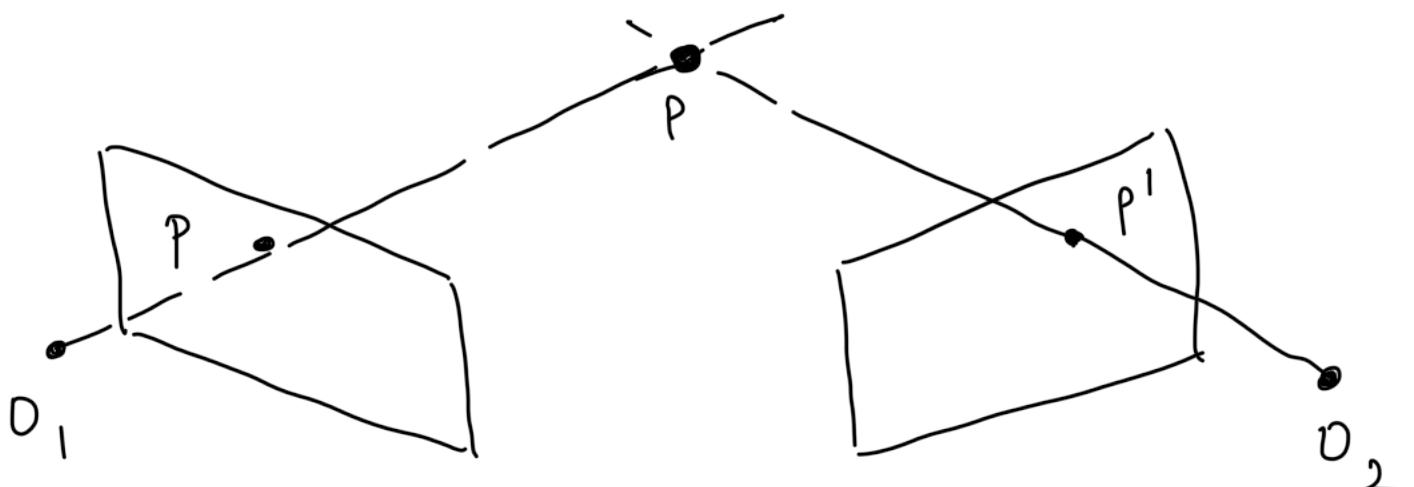


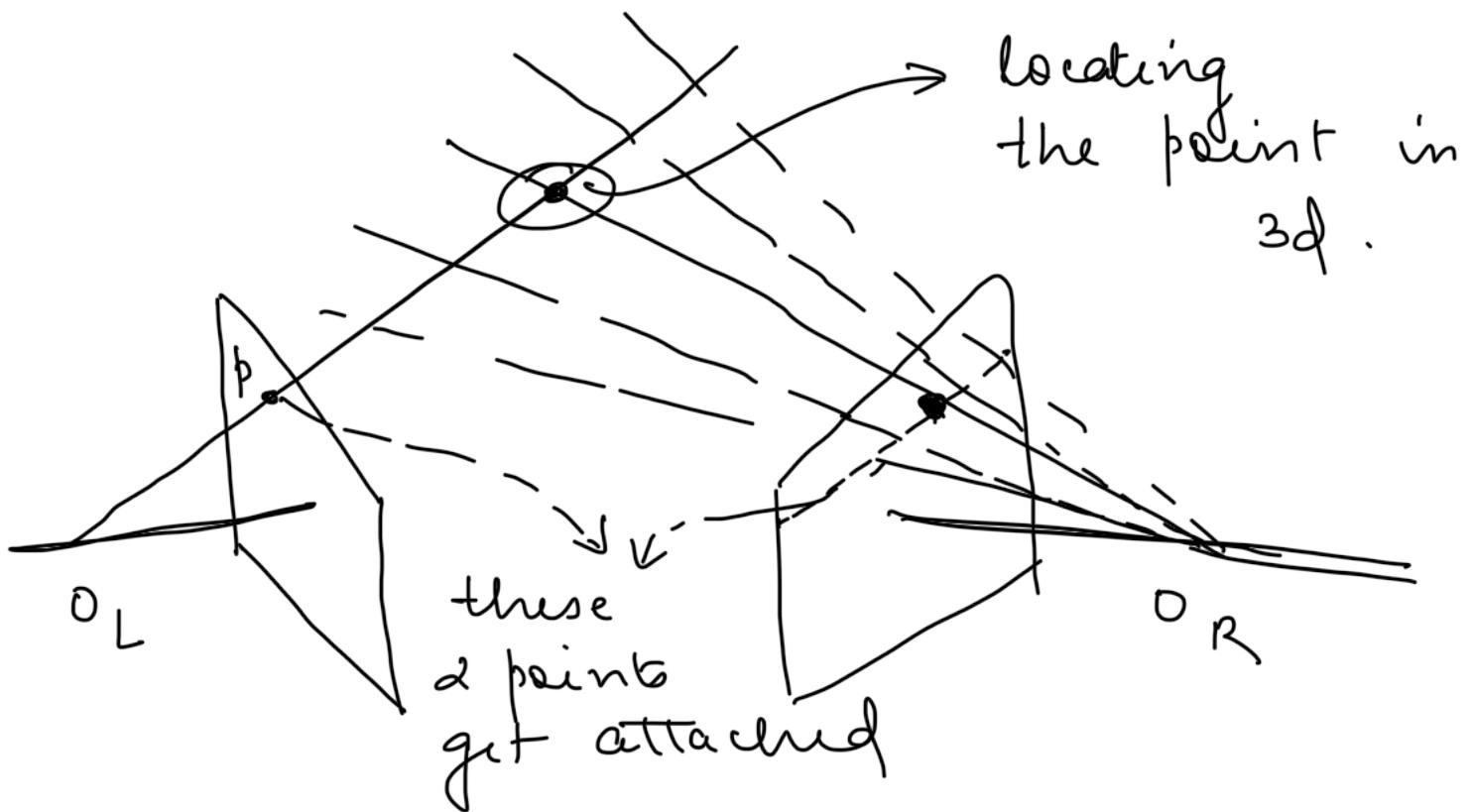
## DLP - Week-11

### L1 Depth Estimation - Interv & Depth From Stereo

- we will be estimating the pixel-level depth relative to the camera position
- applications of the depth estimation
  - 1. autonomous navigation - indoor, roads, aerial
  - 2. VR/AR - virtual meeting rooms, gaming, robotic surgeries
  - 3. scene understanding - semantic segmentation, 3D mapping on a large scale, architectural captures
- understanding depth from stereo



(kind of a parallax error)  
↳ correspondence



## L2 Molecular Depth Estimation - Multi Scale Deep Network

- molecular depth estimation has few challenges -
  1. depth cues like parallax is missing, this time we have only 1 image.
  2. cues available for analysis are vanishing points, decreasing object sizes with depth.
  3. another problem, currently there is no way we can quantify these cues through a rule-based method.
- the ways to handle single image molecular depth estimation methods

1. supervised  $\rightarrow$  multi scale deep network,  
UNet
  2. unsupervised  $\rightarrow$  depth from LR  
consistency.
- Multi Scale deep Network —
    - $\rightarrow$  network has 2 parts — coarse and fine network.
    - $\rightarrow$  coarse network — predicts the overall depth from the scene.
    - $\rightarrow$  refines the global predict locally (fine network)

### ↳ 3 UNet

- the architecture is U shaped
- it has 2 parts — encoder  downsampling + decoder upsampling network.
- encoder  $\rightarrow$  uses conv layers followed by max pooling to finally reduce the spatial dimension while it continues to increase the no. of feature channels. It captures the high-level features & the relevant content from the image.

- decoder → it has the transposed conv. layers to increase the spatial dimensions of feature maps. It goes on to combine low-level features from the downsampling path with high level features to produce detailed output depth maps.
  - skip connec<sup>n</sup> → it directly connects corresponding layers in the down-sampling path to the upsampling path. It preserves the spatial information lost during downsampling. Further, allowing for more precise reconstruction in the output.
  - loss func<sup>n</sup> — MSE → it quantifies the diff. b/w predicted depth values and ground truth values.
 
$$MSE = \frac{1}{n} \sum_{i=1}^N (d_i - \hat{d}_i)^2$$

total no.  
of pixels  
 ↓  
 true value      ↓  
 predicted value
  - loss func<sup>n</sup> — L-depth (MAE)
- $$MAE = \frac{1}{n} \sum_{i=1}^N |d_i - \hat{d}_i|$$

## L4 Left - Right consistency

- problem with supervised learning - is that it requires tons of training data for supervised learning. Also, since we are doing estimation, we need the image depth map pairs. (double the data)
- steps in LR consistency -
  1. use left image of the stereo pair predict the disparity.
  2. true disparity isn't available during testing/training.
  3. next, transform the left to right image using the predicted disparity.
  4. true right image is present in bg during training, can be used for supervision.

