

POLITECNICO DI TORINO

Master's Degree in Mathematical Engineering



Master's Degree Thesis

**Segmenting dynamic points in 3D
scenarios**

Supervisors

Prof. Tatiana TOMMASI
Prof. Chiara PLIZZARI

Candidate

Francesco BORGNA

March 2024

Summary

Ma che dici Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Acknowledgements

ACKNOWLEDGMENTS

*“HI”
Goofy, Google by Google*

Table of Contents

List of Tables	VI
List of Figures	VII
Acronyms	IX
I Related Works	1
1 Datasets	3
2 Neural Rendering	4
3 Photogrammetry	5
3.1 Structure from Motion: SfM	6
3.1.1 Feature Detection and Matching	8
3.1.2 Estimating Fundamental and Essential Matrices	10
3.1.3 Estimating Camera Pose from Essential Matrix	13
3.1.4 Multi-view Structure from Motion	14
3.1.5 COLMAP	14
3.1.6 Monocular Depth Estimation	17
A Galileo	20
B Math Notation	21
Bibliography	22

List of Tables

List of Figures

3.1	Pinhole Camera	5
3.2	Reference systems	7
3.3	Basic SfM Scenario	7
3.4	Feature Detection Example [3]	9
3.5	SIFT features extracted	10
3.6	Features correspondence computed using BruteForce Matcher.	11
3.7	Point correspondence geometry.	11
3.8	12
3.9	Multiple View Scenario.	14
3.10	Building Rome in one day	15
3.11	Monocular datasets	18
3.12	Top: input images. Middle: Inverse depth maps predicted by the presented approach. Bottom: corresponding point clouds rendered from a novel view-point.(Taken from [8])	19

Acronyms

SfM

Structure from Motion

pcd

Pointcloud

Part I

Related Works

-
1. Epic Kitchens
 2. Epic Fields
 3. Photogrammetry
 4. COLMAP
 5. NeRF
 6. NeuralDiff
 7. Monocular Depth Estimation
 8. motion estimation
 9. (N3F)
 10. (Gaussian Splatting)

Chapter 1

Datasets

Chapter 2

Neural Rendering

Chapter 3

Photogrammetry

Photogrammetry is the science and technology of obtaining reliable information about physical objects and the environment through the process of recording, measuring and interpreting photographic images and patterns of electromagnetic radiant imagery and other phenomena[1].

It comprises all techniques concerned with making measurements of real-world objects features from images. Its utility range from the measuring of coordinates, quantification of distances, heights, areas and volumes, preparation of topographic maps, to generation of digital elevation models and orthophotographs. The functioning rely mostly on optics and projective geometry rules.

As first assumption we have the modellization of the camera as a simplified version of itself: the *Pinhole Camera*. As in the first designed cameras(*camera obscura*), in the Pinhole Camera world's light is captured through a pinhole and then projected into the *focal plane*.

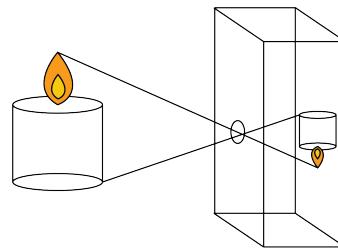


Figure 3.1: Pinhole Camera

A Pinhole camera is characterized by two collection of parameters:

- **Extrinsic** parameters: gives us information on location and rotation in the world.

- **Intrinsic** parameters: gives us internal property such as: focal length, field of view, resolution...

These parameters can be rewritten in their corresponding matrices:

$$Intrinsic = K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

where:

- f_x, f_y are the *focal lengths* of the camera in the x and y directions, they are needed to keep the image aspect ratio.
- c_x, c_y are the coordinates of the *principal point* (the point where the optical axis intersects the image plane).

$$Extrinsic = \begin{bmatrix} \mathbf{R}_{3x3} & \mathbf{t}_{3x1} \\ \mathbf{0}_{1x3} & \mathbf{1}_{1x1} \end{bmatrix}$$

where:

- \mathbf{R}_{3x3} is a rotation matrix ([Aggiungi le tre combinazioni...](#))
- \mathbf{t}_{3x1} is a translation vector

Extrinsic matrix is also known as the 4x4 transformation matrix that converts points from the world coordinate system to the camera coordinate system.

Exploring homogeneous coordinates we can rewrite the image capturing process of a specific camera as the combination of its characteristic matrices:

$$\begin{bmatrix} u \\ v \\ z \end{bmatrix} = \begin{bmatrix} f_x & s & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{R}_{3x3} & \mathbf{t}_{3x1} \\ \mathbf{0}_{1x3} & \mathbf{1}_{1x1} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

3.1 Structure from Motion: SfM

Structure from motion (SfM) is the process of estimating the 3-D structure of a scene from a set of 2-D images. SfM is used in many applications, such as 3-D scanning, augmented reality, and visual simultaneous localization and mapping [2].

We can compute SfM in different ways depending on the data and tools at our disposal. Factors such as the **type** and the **number** of cameras can imply using

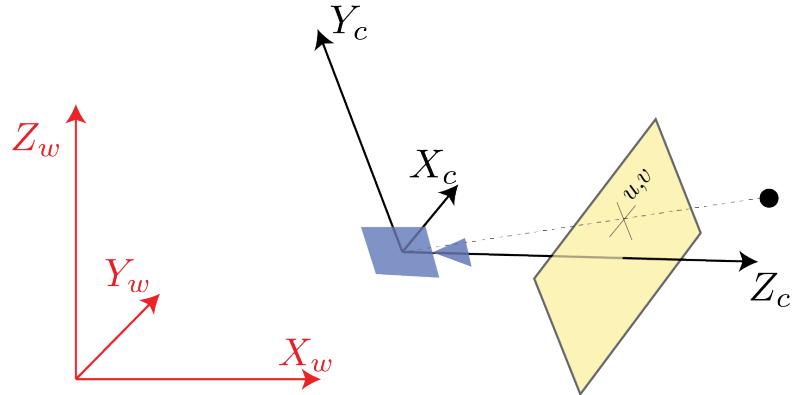


Figure 3.2: Reference systems

different methods. Also the **ordering** of the frames may be relevant for a successful reconstruction.

The most basic scenario consists in two images captured from the same camera from two different point of view:

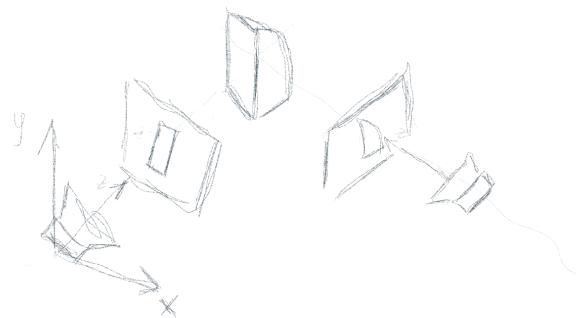


Figure 3.3: Basic SfM Scenario

In this case the 3-D structure can be recovered *up to scale*, meaning that the relations between the obtained points is faithful to the reality, but it can differ from the real size by a *scale* factor. The scale factor could be retrieved if we know the size of an object in the scene or by having informations from other sensors.

In the basic scenario the method consists in the following main steps:

1. Feature Detection and Matching
2. Estimating Fundamental matrix
3. Estimating Essential Matrix from Fundamental matrix
4. Estimating Camera Pose from Essential Matrix

3.1.1 Feature Detection and Matching

In order to find corresponding points in two images we need to detect some points of interest in each image. Instead of trying to match each pixel, which would be computational inefficient and possibly misleading due to color, it has been shown successful to use *feature points*. Feature points are relevant points which encapsulate the local appearance of its surrounding pixels which is invariant under changes in illumination, translation, scale and in-plane rotation. Good features are *unique*, *can be easily tracked* and *can be easily compared*.

A practical example could be to imagine how us humans find correspondances in pictures, looking at image 3.4. If we would be asked to find the exact position of the six patches in the underlying image, the job would be easier for patches E and F, being corners they have less possible misleading correspondances, as happens instead for patches A or B.



Figure 3.4: Feature Detection Example [3]

In the same way feature extractor works for computers. Two of the first algorithms are infact based on corners detection: **Harris Corner Detector** and **Shi-Tomasi Corner Detector**. Depending on the differences between the two images to be compared, different algorithms have been developed. Here I report some of them as reported in [3]:

- **SIFT**: Harris corner detector is not good enough when scale of image changes. Lowe developed a breakthrough method to find scale-invariant features and it is called SIFT
- **SURF** (Speeded-Up Robust Features): faster SIFT version
- **FAST Algorithm for corner detection** All the above feature detection methods are good in some way. But they are not fast enough to work in real-time applications like SLAM. There comes the FAST algorithm, which is really "FAST".
- **BRIEF**(Binary Robust Independent Elementary Features)SIFT uses a feature descriptor with 128 floating point numbers. Consider thousands of such

features. It takes lots of memory and more time for matching. We can compress it to make it faster. But still we have to calculate it first. There comes BRIEF which gives the shortcut to find binary descriptors with less memory, faster matching, still higher recognition rate.

- **ORB(Oriented FAST and Rotated BRIEF)** Open source alternative to SIFT and SURF released by OpenCV.



Figure 3.5: SIFT features extracted

After having found these points of interest, we need to match them from the different pictures. This step is also known as *Feature Matching*. The most basics algorithms are:

- **Brute-Force Matcher:** It takes the descriptor of one feature in first set and is matched with all other features in second set using some distance calculation. And the closest one is returned.
- **FLANN:** which stands for *Fast Library for Approximate Nearest Neighbors*. It contains a collection of algorithms optimized for fast nearest neighbor search in large datasets and for high dimensional features. It works faster than BFMatcher for large datasets.

1

3.1.2 Estimating Fundamental and Essential Matrices

The Fundamental (F) and the Essential (E) matrices allow to relate the projection of a point located in space from one image to the other. These matrices are based

¹This is a footnote with additional information.

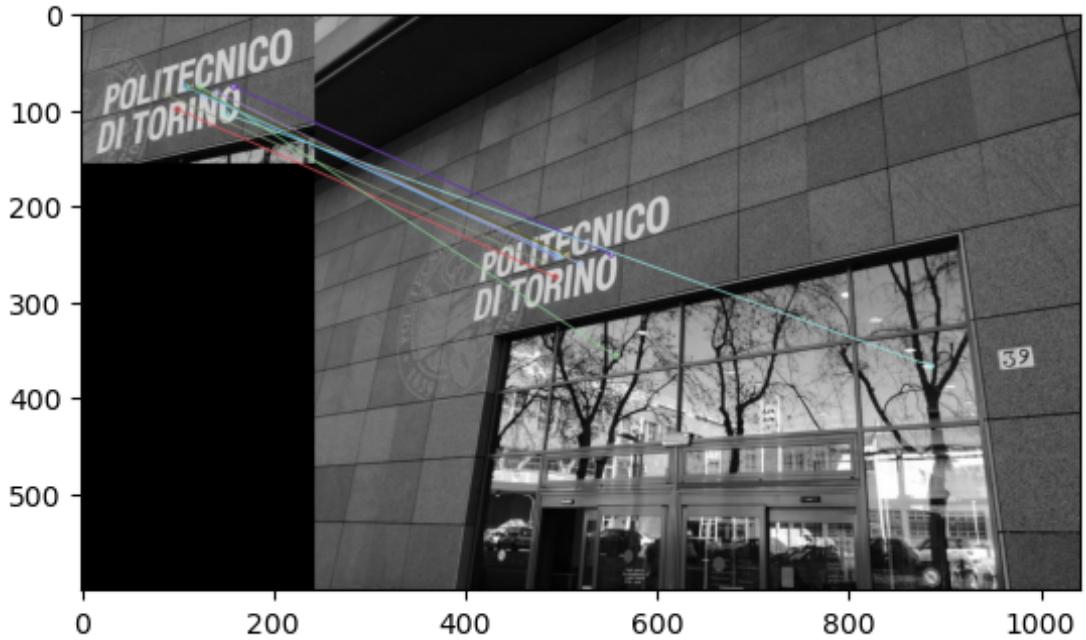


Figure 3.6: Features correspondence computed using BruteForce Matcher.

on the so called *Epipolar Geometry*, which describes the relationship between two images. As we can see from Figure(3.7), given two cameras C_l and C_r the following

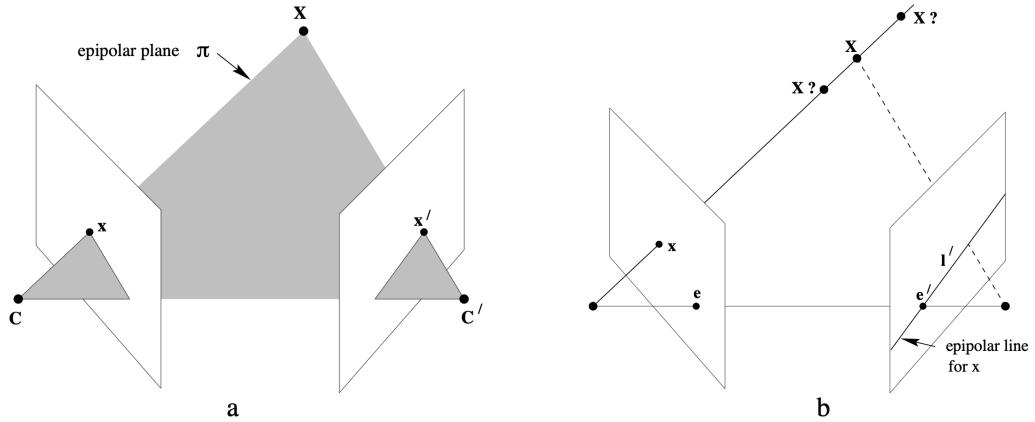


Figure 3.7: Point correspondence geometry.

definitions can be given:

- The **Epipole**: which is the point of intersection of the **baseline**(the line that connects the two camera centers C_l and C_r) with the image plane. In

Figure(3.7) denoted by e_i .

- An **Epipolar plane**, which is any plane containing the baseline.
- An **Epipolar line** which is the intersection of any epipolar plane with the image plane.

Now that we have a clear understanding of the underlying geometry, let's proceed to see the actual derivation of the two matrices.

Let's consider the following scenario, reported in Figure(3.8): we have a point X and two cameras C_l and C_r . The relative positions are respectively x_l and x_r , while the pixel projections are p_l and p_r . If we consider as reference the left camera, the position of the right one is shifted of a vector t .

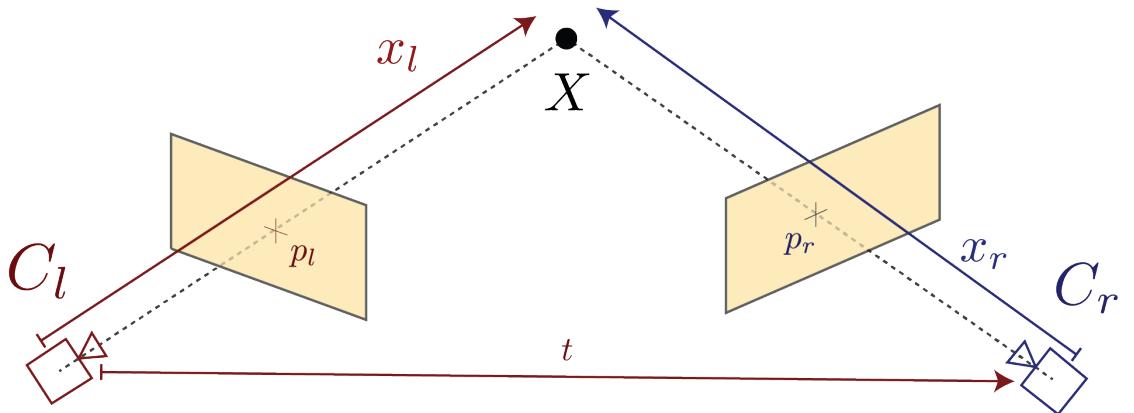


Figure 3.8

We can extract the Essential Matrix by making the following considerations:

$$x_l \cdot (t \times x_l) = 0 \quad (3.1)$$

but x_l can be written as:

$$x_l = Rx_r + t \quad (3.2)$$

So Equation(3.1) becomes:

$$x_l \cdot (t \times (Rx_r + t)) = x_l \cdot (t \times Rx_r) \stackrel{a}{=} x_l^T [t]_{\times} Rx_r = 0 \quad (3.3)$$

where equality (a) is due to the cross-product matrix notation². We call Essential Matrix the term $[t]_{\times} R$, such that:

$$x_l^T [t]_{\times} Rx_r = x_l^T Ex_r = 0 \quad (3.4)$$

In a similar way we can get the Fundamental Matrix, by replacing the relative positions with the pixel positions. Recalling that the pixel positions are linked to the relative positions by:

$$p_l = \frac{1}{z_l} K_l x_l \quad p_r = \frac{1}{z_r} K_r x_r \quad (3.5)$$

where z_i are the focal distances. We can substitute in Eq.(3.4) and obtain:

$$p_l^T z_l K_l^{-1} E K_r^{-1} z_r p_r = 0 \quad z_l, z_r \neq 0 \quad (3.6)$$

Since z_i are constants can be simplified, obtaining:

$$p_l^T K_l^{-1} E K_r^{-1} p_r = p_l^T F p_r = 0 \quad z_l, z_r \neq 0 \quad (3.7)$$

where F is the Fundamental Matrix. We can thus write the relation between the two matrices:

$$E = K_l^T F K_r \quad (3.8)$$

3.1.3 Estimating Camera Pose from Essential Matrix

Since $[t]_{\times}$ is skew symmetric and R is orthonormal(since it is a rotation matrix), if we know the Essential Matrix we can decompose it in its components³using singular value decomposition.

²

$$\mathbf{a} \times \mathbf{b} = [\mathbf{a}]_{\times} \mathbf{b} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

³

$$[t]_{\times} = UW\Sigma U^T \quad (3.9)$$

$$R = UW^T V^T \quad (3.10)$$

$$W = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (3.11)$$

3.1.4 Multi-view Structure from Motion

Now that we have grasped the basics to find images correspondances, let's see how we can relate multiple image to recover the structure of a scene. The last stage is called *bundel adjustment* and it is a iterative algorithm used to adjust structure and motion parameters by minimising a cost function. The possible methods for bundle adjustment are:

- **Sequential:** which work by considering an additional images at each time, extending in this way the initial reconstruction.
- **Factorization:** work by computing camera poses and scene geometry using every image measurement at the same time.

Bundle Adjustment is needed since the image measurements are usually noisy. Minimising an appropriate cost function we can obtain a clean model as in a Linear Regression.

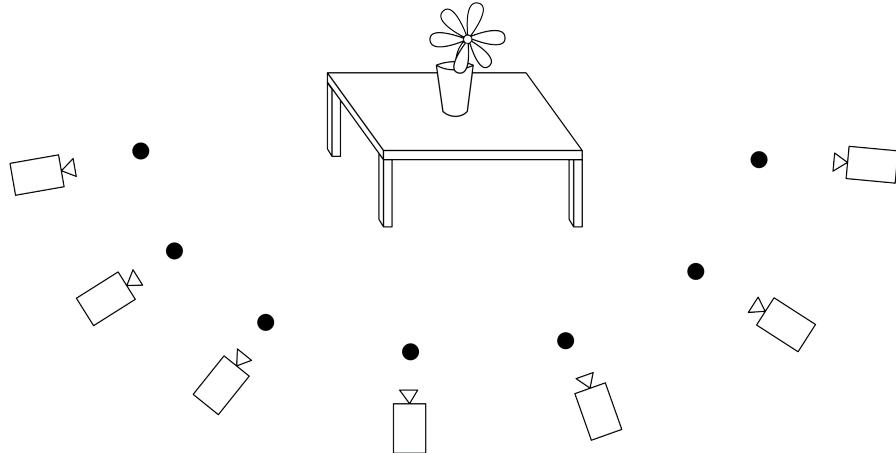


Figure 3.9: Multiple View Scenario.

3.1.5 COLMAP

Abandoning theory and moving on to practice, let's see COLMAP.

"COLMAP is a general-purpose, end-to-end image-based 3D reconstruction pipeline (i.e., Structure-from-Motion (SfM) and Multi-View Stereo (MVS)⁴) with

⁴Even if they seem similar, these two pipelines have different goals. In fact, SfM The primary goal of SfM is to estimate the 3D structure of a scene and camera poses (positions and orientations) simultaneously from a set of 2D images taken from different viewpoints. MVS, on the other hand, is



Figure 3.10: Building Rome in one day

a graphical and command-line interface. It offers a wide range of features for reconstruction of ordered and unordered image collections. The software runs under Windows, Linux and Mac on regular desktop computers or compute servers/clusters. COLMAP is licensed under the BSD License”[4].

The concepts we explained in the theoretical part are not always applicable in real world, or their results are not quite satisfying. In literature a vastity of authors tried and succeeded in obtaining refined algorithms for specific scenarios. Anyway they still lacked a general-purpose method. With COLMAP they managed to compensate for the lack of generalization. The actual implementation and characteristics are explained from the authors in [5, 6].

Usage

The usage of COLMAP is pretty straight forward. After the creation of a project folder with a *images* directory containing the images we want to process we can simply launch the script reported in Listing (3.1).

Listing 3.1: Automatic COLMAP Reconstruction

```

1 #!/bin/bash
2 #The project folder contains a folder "images" with all images.
3
4 DATASET_PATH=/path/to/project
5 colmap automatic_reconstructor \
6   --workspace_path $DATASET_PATH \
7   --image_path $DATASET_PATH/images \
8   --single_camera 1

```

specifically focused on dense 3D reconstruction. It involves estimating the depth or 3D coordinates for every pixel in the images to create a detailed 3D model with meshes and textures.

9

Anyway the program offers ample freedom to the single usage of the various steps need for *SfM* or *MVS*. Here I report the actual pipeline that I used for the extraction of the pointclouds for our experiments.

Listing 3.2: COLMAP Single Commands

```

1  #!/bin/bash
2  DATASET_PATH="/scratch/fborgna/EPIC_Diff/"
3
4  colmap feature_extractor \
5      --database_path $DATASET_PATH/database.db \
6      --image_path $DATASET_PATH/images \
7      --ImageReader.single_camera 1 \
8
9  colmap exhaustive_matcher \
10     --database_path $DATASET_PATH/database.db
11
12 mkdir $DATASET_PATH/sparse
13
14 colmap mapper \
15     --database_path $DATASET_PATH/database.db \
16     --image_path $DATASET_PATH/images \
17     --output_path $DATASET_PATH/sparse \
18     --camera_model SIMPLE_PINHOLE \
19
20 mkdir $DATASET_PATH/dense
21
22 colmap image_undistorter \
23     --image_path $DATASET_PATH/images \
24     --input_path $DATASET_PATH/sparse/0 \
25     --output_path $DATASET_PATH/dense \
26     --output_type COLMAP \
27     --max_image_size 2000
28
29 colmap patch_match_stereo \
30     --workspace_path $DATASET_PATH/dense \
31     --workspace_format COLMAP \
32     --PatchMatchStereo.geom_consistency true
33
34 colmap stereo_fusion \
35     --workspace_path $DATASET_PATH/dense \

```

```

36      --workspace_format COLMAP \
37      --input_type geometric \
38      --output_path $DATASET_PATH/dense/fused.ply
39
40      colmap poisson_mesher \
41          --input_path $DATASET_PATH/dense/fused.ply \
42          --output_path $DATASET_PATH/dense/meshed-poisson.ply
43
44

```

3.1.6 Monocular Depth Estimation

Until now we have always considered scenarios in which multiple images were available from different points of view. What if we have just one image?

This scenario takes the name of *Monocular Depth Estimation*. As we have previously seen, the projection of a scene on the bidimensional image plane of a camera inevitably lose the three-dimensional structure of the scene. In this way we can no more use optics laws, since we do not have enough informations to solve those equations.

Possible solutions include the usage of neural architectures. As a matter of fact, these methods try to learn relations between the possible objects or elements of the image, thus extracting a higher semantic knowledge from the picture. This task is known to be solved particularly good from neural networks, in particular convolutional or transformer based networks (See Section (2)), trained on large amount of data, which are nowadays growing in number and quality.

In particular the advent of transformers has marked a new state of the art, not only in natural language processing, but also in dense prediction. As shown and stated in [7], they "*introduce dense vision transformers, an architecture that leverages vision transformers in place of convolutional networks as a backbone for dense prediction tasks. [...] this architecture yields substantial improvements on dense prediction tasks. [...] we observe an improvement of up to 28% in relative performance when compared to a state-of-the-art fully-convolutional network*". The motivation is due to the different functioning of the two networks. Infact, both methods implies an encoder-decoder structure but the nature of convolutional layer limits its performance. As explained in Sec.?? convolutional networks progressively downsample the input image to extract features at multiple scales. Unfortunately in dense prediction tasks, keeping feature resolution and granularity showed to be essential, since the decoder can't recover the initial informations from the compressed ones. "*Unlike fully-convolutional networks, the vision transformer backbone foregoes explicit downsampling operations after an initial image embedding*

"has been computed and maintains a representation with constant dimensionality throughout all processing stages"[7].

The second motivation that made us choose for the transformer-based architecture is its versatility and generalizability. In [8] the authors introduce a method that enable mixing multiple datasets during training. In this way the model will be effective across a variety of scenarios since each dataset is equally varied and captures the diversity of the visual world. In particular they experimented with eleven datasets which consists of RGB images with corresponding depth annotation of some form:

- Training: DIML Indoor,MegaDepth,ReDWeb,WSVD, 3D Movies
- Testing: DIW,ETH3D,Sintel,KITTI,NYUDv2,TUM-RGBD

Each single dataset has its own characteristic and has its own biases and problems. Mixing them during the train phase allows to mitigate those problems and obtain a good generalization. To evaluate the generalization they used the *Zero-Shot Cross-dataset Transfer* which consists in evaluating on datasets that were not seen during training. This proved to be extremely successful.

Dataset	Indoor	Outdoor	Dynamic	Video	Dense	Accuracy	Diversity	Annotation	Depth	# Images
DIML Indoor	✓			✓	✓	Medium	Medium	RGB-D	Metric	220K
MegaDepth		✓	(✓)		(✓)	Medium	Medium	SfM	No scale	130K
RedWeb	✓	✓	✓		✓	Medium	High	Stereo	No scale & shift	3600
WSVD	✓	✓	✓	✓	✓	Medium	High	Stereo	No scale & shift	1.5M
3D Movies	✓	✓	✓	✓	✓	Medium	High	Stereo	No scale & shift	75K
DIW	✓	✓	✓			Low	High	User clicks	Ordinal pair	496K
ETH3D	✓	✓			✓	High	Low	Laser	Metric	454
Sintel	✓	✓	✓	✓	✓	High	Medium	Synthetic	(Metric)	1064
KITTI		✓	(✓)	✓	(✓)	Medium	Low	Laser/Stereo	Metric	93K
NYUDv2	✓		(✓)	✓	✓	Medium	Low	RGB-D	Metric	407K
TUM-RGBD	✓		(✓)	✓	✓	Medium	Low	RGB-D	Metric	80K

Figure 3.11: Monocular datasets



Figure 3.12: Top: input images. Middle: Inverse depth maps predicted by the presented approach. Bottom: corresponding point clouds rendered from a novel view-point.(Taken from [8])

Appendix A

Galileo

```
1 import os  
2 os.system("echo 1")
```

$\mathcal{O}(n \log n)$
numpy

Appendix B

Math Notation

$$\mathbf{a} \times \mathbf{b} = [\mathbf{a}]_{\times} \mathbf{b} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$
$$\mathbf{a} \times \mathbf{b} = [\mathbf{b}]^T_{\times} \mathbf{a} = \begin{bmatrix} 0 & b_3 & -b_2 \\ -b_3 & 0 & b_1 \\ b_2 & -b_1 & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

Bibliography

- [1] Wolf and Dewitt. 2000; McGlone, 2004 (cit. on p. 5).
- [2] Mathworks 2023. URL: [https://www.mathworks.com/help/vision/ug/structure-from-motion.html#:~:text=Structure%20from%20motion%20\(SfM\)%20is,localization%20and%20mapping%20\(vSLAM\)](https://www.mathworks.com/help/vision/ug/structure-from-motion.html#:~:text=Structure%20from%20motion%20(SfM)%20is,localization%20and%20mapping%20(vSLAM)). (cit. on p. 6).
- [3] OpenCV 2024. URL: https://docs.opencv.org/4.x/d54/tutorial_py_features_meaning.html (cit. on p. 9).
- [4] Colmap 2024. URL: <https://colmap.github.io/index.html> (cit. on p. 15).
- [5] Johannes Lutz Schönberger and Jan-Michael Frahm. «Structure-from-Motion Revisited». In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016 (cit. on p. 15).
- [6] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. «Pixelwise View Selection for Unstructured Multi-View Stereo». In: *European Conference on Computer Vision (ECCV)*. 2016 (cit. on p. 15).
- [7] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. «Vision Transformers for Dense Prediction». In: *ICCV* (2021) (cit. on pp. 17, 18).
- [8] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. «Towards Robust Monocular Depth Estimation: Mixing Datasets for Zero-Shot Cross-Dataset Transfer». In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44.3 (2022) (cit. on pp. 18, 19).