

A comparative study on classical and Deep Learning-based approaches for Stereo Matching

Amresh Kumar (M23CSA004)
Divyaansh Mertia (M23CSE013)
Sujay Kumar Ingle (D23CSA003)
Kulendu Kashyap Chakraborty (C23CS1007)

April 28, 2024

Indian Institute of Technology, Jodhpur

Contents

1	Problem Statement	1
1.1	Introduction	1
1.2	Classical Approaches	1
1.2.1	Winner Takes It All (WTA)	1
1.2.2	Semi-Global Matching (SGM)	1
1.3	Deep Learning-Based Approach	2
1.4	Comparative Analysis	2
1.5	Conclusion	2
2	Methodology	3
2.1	Local matching	3
2.1.1	Sum of absolute differences	3
2.1.2	Sum of squared differences	3
2.2	Cost volume	4
2.3	Matching algorithm	4
2.3.1	Winner-takes-it-all solution	4
2.3.2	Semi-global matching	4
2.4	Evaluation: compare to ground-truth	5
3	Results	6
3.1	Input Images	6
3.2	Classical Approach	7
3.2.1	Winner takes it all	7
3.2.2	Semi Global Matching	8
3.3	DL Approach	9
4	Observations and Conclusion	10
5	Reference	11
Bibliography		12
A	Code	A.1

List of Tables

List of Figures

3.1	Input Images	6
3.2	Winner Takes It All	7
3.3	Semi Global Matching	8
3.4	Practical Stereo Matching via Cascaded Recurrent Network with Adaptive Correlation	9

CHAPTER 1

Problem Statement

1.1 Introduction

This Project presents a comparative study between classical stereo matching techniques and modern deep learning-based approaches. The objective is to analyze the effectiveness of each method in the context of stereo matching.

1.2 Classical Approaches

Classical stereo matching algorithms are based on direct pixel comparisons and optimization techniques. The primary methods evaluated in this study include:

1.2.1 Winner Takes It All (WTA)

The Winner Takes It All approach utilizes cost functions to determine the disparity between stereo images. The key cost functions examined are:

1. Sum of Absolute Differences (SAD) – This method calculates the absolute difference between pixel values in left and right images over a window, which is assumed to be the disparity.
2. Sum of Squared Differences (SSD) – Similar to SAD, but squares the difference to penalize larger disparities more severely.

1.2.2 Semi-Global Matching (SGM)

Semi-Global Matching improves the local methods by considering pixel discontinuities along several lines (paths) across the image, optimizing the global consistency of disparity. The cost functions used include:

1. Sum of Absolute Differences
2. Sum of Squared Differences

1.3 Deep Learning-Based Approach

The deep learning-based approach explored in this study involves the Practical Stereo Matching via Cascaded Recurrent Network with Adaptive Correlation (PSM-CR-AC). This method leverages a neural network architecture that processes stereo images through cascaded stages to refine disparity predictions adaptively. Key features of this approach include:

- **Cascaded Processing:** Uses multiple stages of processing to incrementally refine the disparity estimates.
- **Recurrent Neural Networks:** Employs RNNs to capture spatial dependencies along image rows.
- **Adaptive Correlation:** Implements a learning-based correlation layer that adapts to different image and disparity conditions, improving robustness over traditional fixed correlation methods.

1.4 Comparative Analysis

The comparative analysis aims to evaluate the performance of classical methods against the deep learning-based approach in terms of accuracy, computational efficiency, and robustness across different stereo datasets.

1.5 Conclusion

This study will provide insights into the trade-offs and benefits of classical and deep learning-based stereo matching techniques, highlighting the scenarios where each method excels.

CHAPTER 2

Methodology

Classical Approach

In this code several methods for estimating disparity in a stereo image are implemented. For this purpose a *matching cost volume* is calculated by means of sum of squared differences (SSD), sum of absolute differences (SAD) and normalised cross-correlation (NCC) and then the most appropriate match chosen either by the simple *winner-takes-it-all* approach (WTA) or *semi-global matching* (SGM). For this purpose the given images have to be converted to grayscale.

2.1 Local matching

The following error-measures and correlations will be used for evaluating a corresponding matching cost between two image patches p and q of equal size $W \times H$.

2.1.1 Sum of absolute differences

In case of the sum of absolute differences the matching of two patches p and q is penalised depending on the sum of absolute differences of the two windows according to

$$SAD(p, q) = \sum_{x=1}^W \sum_{y=1}^H |p(x, y) - q(x, y)| \quad (2.1)$$

This means very similar image patches lead to a low SAD while non-matching patches result in a high SAD.

2.1.2 Sum of squared differences

In case of the sum of squared differences the matching process is penalised quadratically instead of linearly making use of the squared difference instead

$$SSD(p, q) = \sum_{x=1}^W \sum_{y=1}^H (p(x, y) - q(x, y))^2 \quad (2.2)$$

2.2 Cost volume

We use these similarity measures to compute a cost-volume CV for a pre-defined range of disparities D

$$CV(x, y, d) = S(I_0(x, y) I_1(x - d, y)) \quad (2.3)$$

where the parameter $d \in \mathcal{D}$ and $\mathcal{D} = \{0, \dots, D - 1\}$ are all valid disparities and S is any of the aforementioned error-measures.

This basically means that we take the left picture and translate the right picture trying to overlap the objects in the two pictures taken from different views. The points at a certain depth have a certain disparity and thus the optimal shift can be used to determine the correct depth. In order to account for a certain deviation we use a certain search window (W, H) rather than trying to match the points directly.

2.3 Matching algorithm

2.3.1 Winner-takes-it-all solution

One fast way of obtaining then the best disparity for each image point would be taking the point with the highest value in the cost volume along the disparity axis according to

$$\bar{d}(x, y) \in \arg \min_d CV(x, y, d) \quad (2.4)$$

This though leads to noisy results as this approach doesn't penalise label changes at all.

2.3.2 Semi-global matching

In semi-global matching a different approach is taken, rather than looking for the best fit on a scanline, a sort of global optimisation is used. Each pixel with a corresponding unary cost given by the cost volume is assigned an additional pairwise cost that depends on wherever the neighbouring pixels have a similar depth value or deviate significantly. This energy can be written as

$$\min_z \left[\sum_{i \in \mathcal{V}} g_i(z_i) + \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i, z_j) \right] \quad (2.5)$$

where \mathcal{V} are the image pixels, \mathcal{E} the edges, the connections between two pixels. The g_i are given by the cost volume and the pairwise cost $f_{i,j}$ defines a penalty for jumps between neighbouring pixels.

$$f_{i,j}(z_i, z_j) = \begin{cases} 0, & \text{if } z_i = z_j \\ L_1, & \text{if } |z_i - z_j| = 1 \\ L_2 & \text{else} \end{cases} \quad (2.6)$$

This is done as following: First messages for all four disparity directions are calculated where the first message in each direction is initialised with $\vec{0}$.

$$m_{i+1}^a(t) = \min_{s \in \mathcal{D}} [m_i^a(s) + f_{i,i+1}(s, t) + g_i(s)] \quad (2.7)$$

This can be done for every direction by a combination of mirroring and transposing the cost volume. Then the beliefs are computed

$$b_i(s) = g_i(s) \sum_{a \in \{L,R,U,D\}} m_i^a(s) \quad (2.8)$$

The correct disparity is then calculated from the believes as follows

$$\hat{d}(x, y) \in \arg \min_d b(x, y, d) \quad (2.9)$$

The last formula contains is intentionally given as \in as the solutions might not be unique.

2.4 Evaluation: compare to ground-truth

The performance of the stereo workflow is evaluated by comparing it with a ground-truth disparity map, in this case with the $accX$ measure

$$accX(z, z*) = \frac{1}{Z} \sum_{x=1}^W \sum_{y=1}^H m(x, y) \cdot \begin{cases} 1, & \text{if } |(z(x, y) - z^*(x, y))| \leq X \\ 0 & \text{else} \end{cases} \quad (2.10)$$

This measure characterises errors less than or equal to X disparities, between the prediction z and the ground truth disparity map z^* with a mask m that contains 1 for the Z valid pixels and 0 for the invalid pixels.

The mask basically excludes pixels that should not be evaluated e.g. because they are occluded in either of the two pictures. The average of the remaining pixels that were estimated correctly is determined. All pixels that guessed the depth correctly (threshold X) are set to 1, all pixels that did not estimate it correctly do not contribute. In this way the $accX$ measures the amount of pixels that were matched correctly to those that could possibly be matched. An $accX$ of 1 would correspond to the ground truth.

CHAPTER 3

Results

3.1 Input Images



(a) Left Image

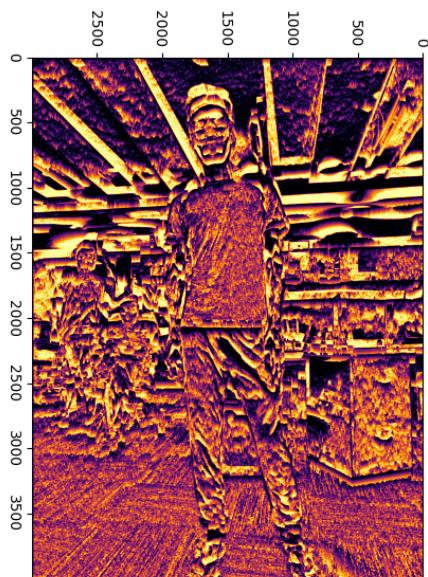


(b) Right Image

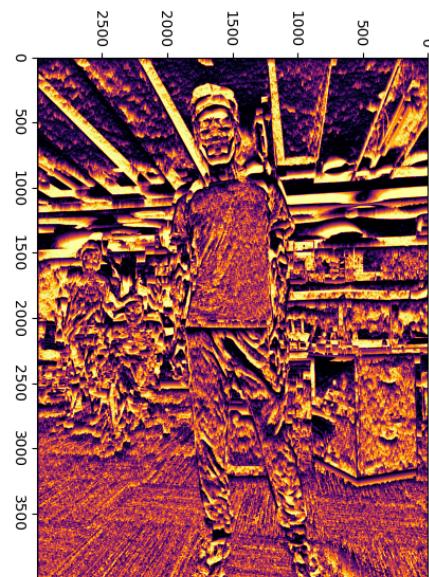
Fig. 3.1: Input Images

3.2 Classical Approach

3.2.1 Winner takes it all



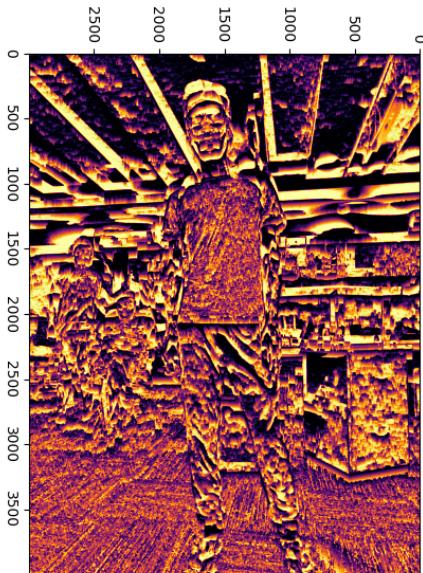
(a) Sum of Absolute Differences



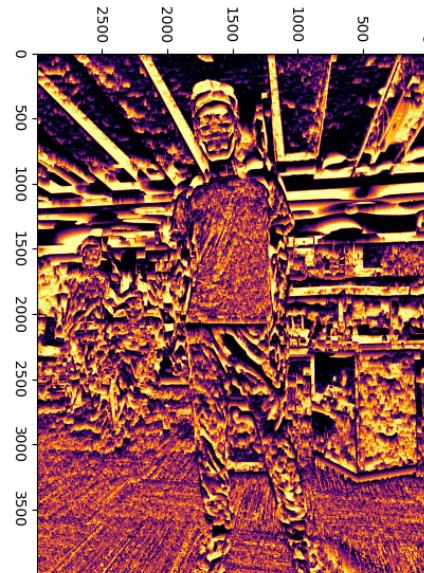
(b) Sum of Squared Differences

Fig. 3.2: Winner Takes It All

3.2.2 Semi Global Matching



(a) Sum of Absolute Differences



(b) Sum of Squared Differences

Fig. 3.3: Semi Global Matching

3.3 DL Approach



Fig. 3.4: Practical Stereo Matching via Cascaded Recurrent Network with Adaptive Correlation

CHAPTER 4

Observations and Conclusion

CHAPTER 5

Reference

Bibliography

APPENDIX A

Code
