# Project 1 - MyAutoPano
# Phase 2 - Deep Learning Approach

Manoj Velmurugan[*], Rishabh Singh[†]

Robotics Engineering

Worcester Polytechnic Institute

Email: [*]v.manoj1996@gmail.com, [†]rsingh8@gmail.com

*Abstract*—**This work presents a comprehensive exploration of homography estimation using both supervised and unsupervised learning approaches. A curated dataset was meticulously generated to train the homography model, ensuring exposure to diverse scenarios. Remarkably, both methodologies, supervised and unsupervised, exhibited good accuracy in homography estimation. The unsupervised approach introduced a Differentiable Pseudo Inverse-based Direct Linear Transform (DLT) algorithm, enhancing the model's robustness.**

**The supervised learning-trained model demonstrated its efficacy in stitching images, showcasing the practical utility of the homography model. This work amalgamates the strengths of supervised and unsupervised learning strategies, offering a versatile framework for homography estimation and real-world image stitching applications.**

**One late day was used!**

## I. KEY CONTRIBUTIONS

- **Input image size choice:** Our network processes (256,256) size images instead of (128,128) to preserve crucial features from the input images.
- **Implementation of Pseudo Inverse-based Differentiable Direct Linear Transform (DLT):** This work incorporates a Pseudo Inverse-based differentiable Direct Linear Transform (DLT) approach.
- **Loss computation focused on the center region in unsupervised learning:** In unsupervised learning, the loss is computed exclusively for the central region of the image to mitigate issues related to black regions.
- **Utilization of HDF5 file format for training efficiency:** The HDF5 file format is employed to store and access data efficiently during training, reducing disk access time.

## II. DATA GENERATION

To train a Convolutional Neural Network (CNN) called HomographyNet for estimating homography between image pairs, synthetic data generation is essential. The MSCOCO dataset, known for its diverse natural scenes, provides a suitable source.

In this work, training and validation - images and labels were obtained by the following,

- **Resize:** Resize the input image to a fixed size of (512, 512).
- **Random Subpatch Selection:** Randomly choose a translation and perturbation for a subpatch of size (256, 256) located at the center of the resized image.
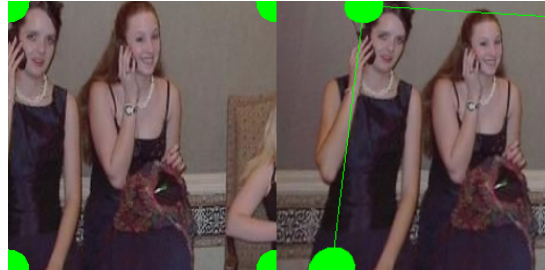


Fig. 1: Sample validation image with corners overlaid

- **Homography Matrix:** Obtain the transformation matrix using `cv2.getPerspectiveTransform` based on the chosen translation and perturbation.
- **Warp Perspective:** Warp the perspective of the resized image using `cv2.warpPerspective` with the obtained homography matrix.
- **Crop Center Portions:** Crop the center portions of both the original (512, 512) size image and the transformed image to obtain patches.
- **Homography Calculation:** Subtract the patch corners to obtain the H4pt (homography matrix represented by four corner points).
- **Data Storage:** Dump the contents, including the original image, transformed image, and H4pt, into an HDF5 file for improved training performance.

A sample generated image along with 4 corners is shown in the figure 1.

Our training dataset contained 60000 samples and testing/validation dataset contained 12000 samples.

## III. HOMOGRAPHY NETWORK

A Convolutional Neural Network (CNN) was formulated to compute $H4$ points from grayscale images of size $256 \times 256$. The loss function used is a simple L2 loss between the predicted 4-point homography $\hat{H}_4$ and the ground truth 4-point homography $H_4$. The loss function $l$ is given by $l = ||\hat{H}_4 - H_4||_2$.

The estimation pipeline is shown in figure 2.

The network architecture is shown in figure 3 and specification is tabulated in table I. Using any simpler model, results
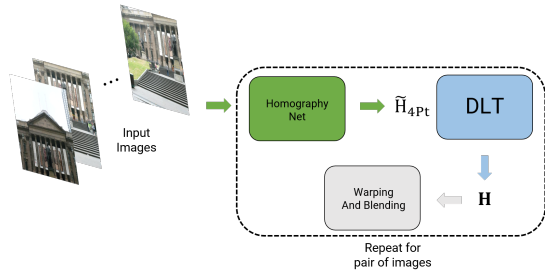
Fig. 2: Overview of supervised homography

| | |
|---|---|
| Parameter count | 1.081M |
| Kernel size | (3,3) |
| Conv layer count | 10 |
| Input size | (batch, 2, 256, 256) |

TABLE I: Homography Network Specification

in higher loss value. Dropout was added to the fully connected layers to improve generalization.

### A. Training

The model was trained with some initial learning rate like 2e-3. It was reduced down to 1e-4 in between training. Reduced training rate did not seem to help and to speed up training, the rate was increased to 2e-3. Complete list of hyper parameters can be found in table II.

Training and validation loss was summed over all the batches and plotted every epoch (figure 4, figure 5).

### B. Homography Result

When there are good features on all 4 corners, the match is pretty close as shown in figure 6. as expected, the detected corner transformations are good (figure 6b).
When the features are sparse (such as sky), the detection is off for a single corner - figure 7.

From the loss curves, the network seems to be still training and this result is expected to improve in future.

## IV. UNSUPERVISED HOMOGRAPHY

In the context of homography estimation, where conventional supervised learning heavily relies on accurately labeled datasets, the demand for large, precisely annotated datasets

| | |
|---|---|
| Learning rate | 1e-4 to 2e-3 (manually adjusted) |
| Batch size | 2048 |
| Optimizer | Adam |
| Loss function | L2 loss |
| Train epochs | 1000 |
| GPU memory usage | 6460 MB |
| Train Batch count | 30 |
| Validation Batch count | 6 |

TABLE II: Supervised Learning Hyper Parameters

| | |
|---|---|
| Epoch | 781 |
| avg training loss | 385.88 |
| avg test loss | 1700.75 |

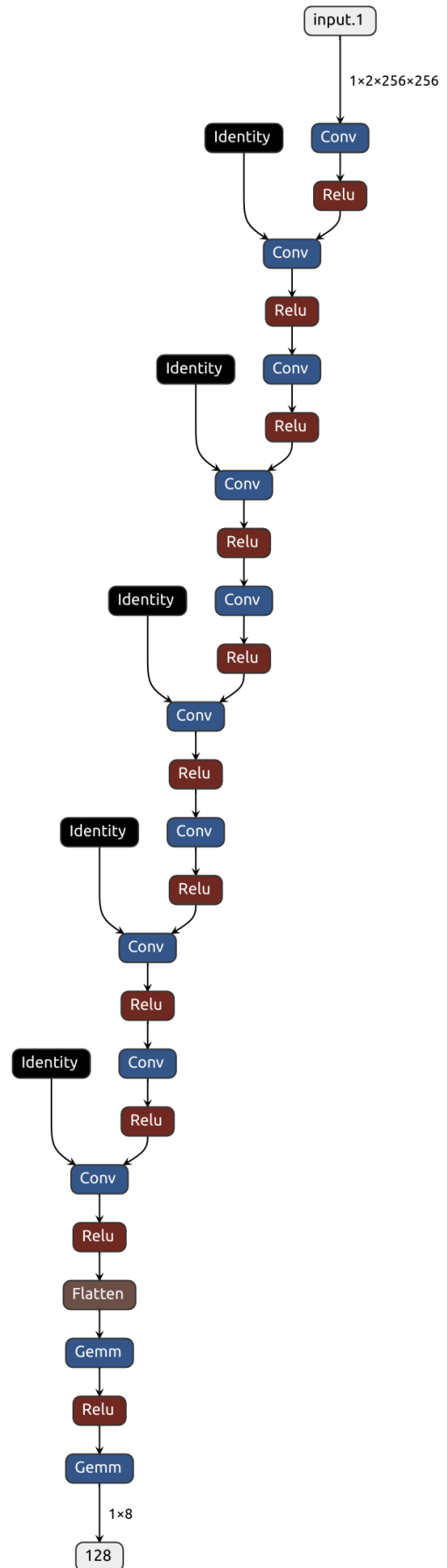TABLE III: Supervised Learning - Loss
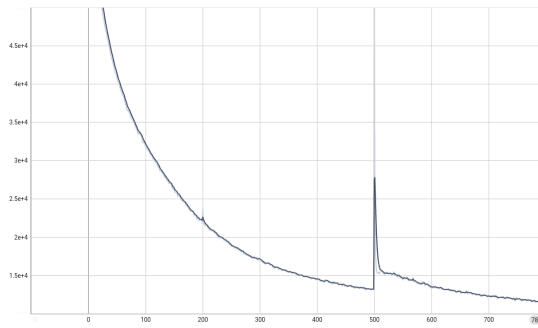


Fig. 3: Homography Network Architecture

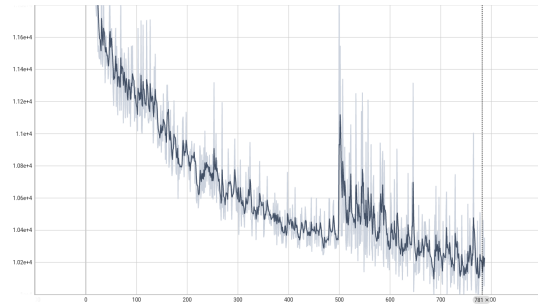Fig. 4: Supervised Learning - Training Loss vs Epoch



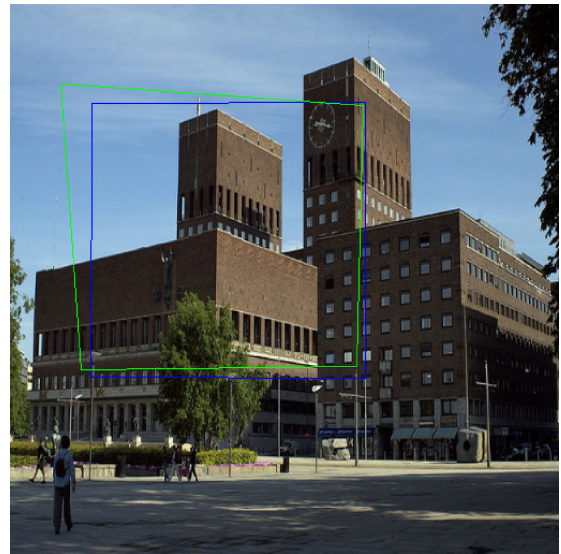Fig. 5: Supervised Learning - Validation Loss vs Epoch



(a) Homography: green - label, blue - prediction



(b) Network inputs overlaid with detected boundary transformation

Fig. 6: Supervised learning - Homography for test image 1



(a) Homography: green - label, blue - prediction



(b) Network inputs overlaid with detected boundary transformation

Fig. 7: Supervised learning - Homography for test image 2

can pose significant challenges. However, an alternative approach emerges with the exploration of unsupervised learning techniques. Rather than depending on explicit labels, we use an implicit evaluation strategy for homography estimation. This novel methodology involves leveraging the predicted corner points from the homography network to warp images. By subsequently comparing these warped images with their originals, we can formulate an unsupervised loss function. This innovative paradigm allows the network to learn the underlying spatial transformations implicitly (figure 8). The ensuing section details the implementation and results of this unsupervised learning approach, shedding light on its potential advantages.

*A. Algorithm*

To facilitate unsupervised learning in homography estimation, the following methodology is used:

1) **H4pt Estimate from Homography Network:** The homography network, utilizing the same architecture as in supervised learning, generates a prediction of the homography matrix ($H_{4pt}$) based on the input images.

2) **Computation of H Matrix using Pseudo-Inverse:** A differentiable and vectorized approach is adopted to compute the homography matrix ($H$) using the pseudo-
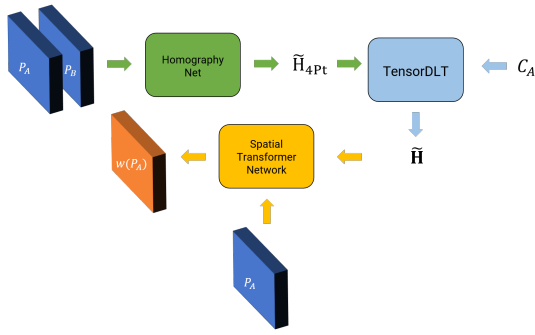
Fig. 8: Overview of unsupervised homography

inverse of least squares problem involving $H_{4pt}$ as shown in listing 1. This enables end-to-end differentiability, crucial for training.

3) **Kornia Package for Image Warping:** The computed homography matrix ($H$) and one of the input images ($A$) are passed to the Kornia package. Leveraging its functionality, the package warps the image ($A$) to generate an estimate of the second image ($B$) based on the homography transformation.

4) **Loss Computation using L1 Function:** To evaluate the dissimilarity between the estimated and original images ($B$), a loss function is computed. Both the estimated and original images are center-cropped to remove undesired black regions, and the L1 loss function is employed for loss computation.

```
A = torch.zeros(batchsize, 8, 6).to("cuda"
    )
# CA is original points, CB is transformed
 corner points
B = predicted + CA
for i in range(0, 4):
    A[:, 2*i, 0] = CA[:, 2*i]
    A[:, 2*i, 1] = CA[:, 2*i+1]
    A[:, 2*i, 2] = 1

    A[:, 2*i+1, 3] = CA[:, 2*i]
    A[:, 2*i+1, 4] = CA[:, 2*i+1]
    A[:, 2*i+1, 5] = 1
A.requires_grad_(True)
B.requires_grad_(True)

# SOLVE A H = B, where H is 6x1 vector for
 each batch
Hvec = torch.linalg.pinv(A) @ B.unsqueeze
    (-1)
```

Listing 1: Differentially solving for the first 6 elements of $H$ matrix

### B. Training

The model was trained with some initial learning rate like 2e-3. It was reduced down to 1e-4 in between training. Reduced training rate did not seem to help and to speed up training, the rate was increased to 2e-3. Complete list of hyper parameters can be found in table IV.

| | |
|---|---|
| Learning rate | 1e-4 to 2e-3 (manually adjusted) |
| Batch size | 2048 |
| Optimizer | Adam |
| Loss function | L1 loss |
| Train epochs | 1000 |
| GPU memory usage | 15336 MB |
| Train Batch count | 30 |
| Validation Batch count | 6 |

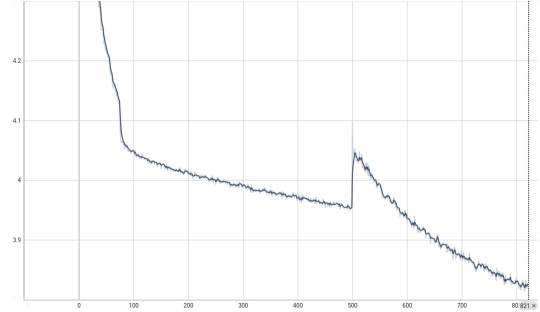TABLE IV: Unsupervised Learning Hyper Parameters



Fig. 9: Unsupervised Learning - Training Loss (summed over batches) vs Epoch

Training and validation loss was summed over all the batches and plotted every epoch (figure 9, figure 10).

### C. Homography Result

When there are good features on all 4 corners, the match is pretty close as shown in figure 11, as expected, the detected corner transformations are good (figure 11b).
When the features are sparse (such as sky), the detection is off for a single corner - figure 12.

From the loss curves, the network seems to be still training and this result is expected to improve in future.

## V. IMAGE STITCHING

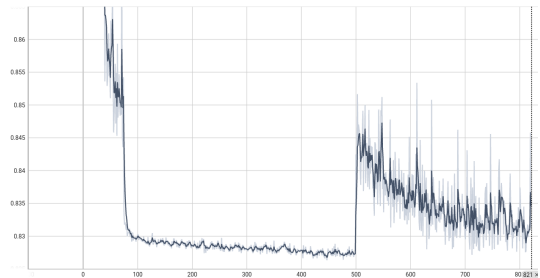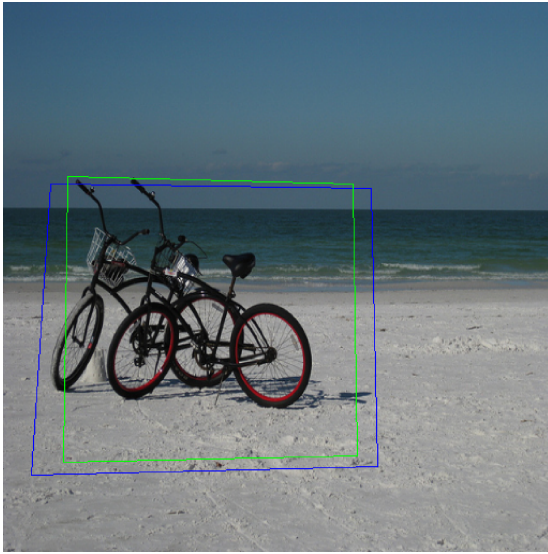The blending algorithm is captured below,



Fig. 10: Unsupervised Learning - Validation Loss (summed over batches) vs Epoch

| | |
|---|---|
| Epoch | 814 |
| avg training loss | 0.127 |
| avg test loss | 0.138 |

TABLE V: Supervised Learning - Loss

(a) Homography: green - label, blue - prediction



(b) Network inputs overlaid with detected boundary transformation

Fig. 11: Unsupervised learning - Homography for test image 1
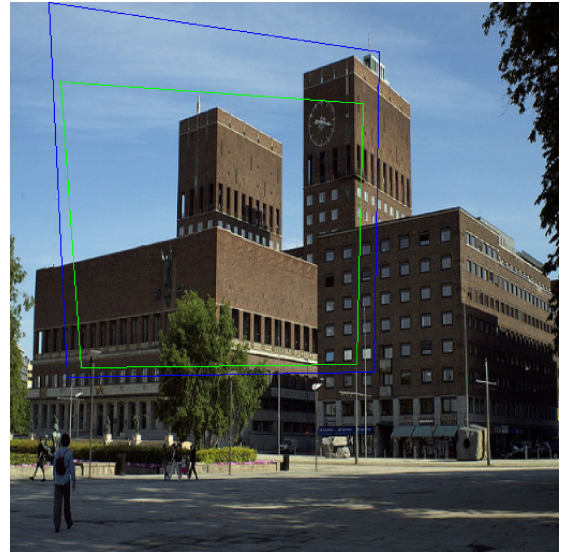


(a) Homography: green - label, blue - prediction



(b) Network inputs overlaid with detected boundary transformation

Fig. 12: Unsupervised learning - Homography for test image 2

1) **Reshape Images:** Reshape images 1 and 2 to $256 \times 256$. Find corresponding transformation matrices $H1$ and $H2$ for this scaling operation.
2) **Convert to Gray Scale and Transform:** Convert the reshaped images to grayscale, transform them to tensors, and feed them to the neural network.
3) **Compute Transformation Matrix H:** From the predicted homography 4-points ($H4pts$), compute the transformation matrix $H$.
4) **Compute Network Transformation Matrix $H_{net}$:** Compute $H_{net} = H1 \cdot H \cdot H2^{-1}$. This matrix can operate on the original image irrespective of its shape.
5) **Apply $H_{net}$ to Image 1:** Apply $H_{net}$ to $img1$, shift the image, and paste it on a large canvas.
6) **Apply $H_{net}$ to Image 2:** Apply $H_{net}$ to $img2$, shift the image, and paste it on another large canvas.
7) **Blend Images:** Blend both images using a weighted average with 0.5 as the weight. Utilize the `cv2.addWeighted` function.

The blending result is shown in figure 13. This result was computed for 998 epoch's model checkpoint.

## VI. KEY FINDINGS

- **The difference between training loss and test loss is not significant in unsupervised learning:** The marginal disparity between training loss and test loss in unsupervised learning scenarios suggests that this approach may mitigate overfitting, leading to improved generalization.
- **Unsupervised learning duration compared to supervised learning:** Unsupervised learning typically requires a more extended training time in comparison to supervised learning. This is attributed to the need for the model to indirectly infer labels through the optimization process of gradient descent.
- **Training epochs for model usability in image stitching:** Our model necessitated training for 998 epochs to achieve usability in the context of stitching images.
- **Sufficiency of Pseudo Inverse-based Direct Linear Transform (DLT):** The utilization of a Pseudo Inverse-based Direct Linear Transform (DLT) approach appears to be adequate for the given task.

## VII. CONCLUSION

As detailed above, we undertook the training of a homography model employing both supervised and unsupervised

Fig. 13: Stitching using supervised Homography network

approaches. A meticulously crafted dataset was generated to fuel the iterative training process, ensuring the model's exposure to diverse scenarios. Remarkably, both the supervised and unsupervised methodologies demonstrated good accuracy in homography estimation. The unsupervised approach leveraged a Pseudo Inverse-based Direct Linear Transform (DLT) algorithm.

In the final stage of our exploration, the model trained through supervised learning was tested in the real-world task of stitching images. This application showcases the versatility and practical utility of the homography model developed in this study. The supervised and unsupervised learning strategies, coupled with the implementation of a specialized DLT algorithm, contributes to a comprehensive and robust framework for homography estimation.