

# **Recognition of Previously Identified Individuals using Affine Moment Invariants as Feature Descriptors in Gait Sequences**

*Internship report submitted in partial fulfilment of the requirements  
for the degree of B.Tech. + M.Tech. in Electronics and Communication  
Engineering*

*by*

Sandesh V Bharadwaj  
(Roll No: ESD15I005)



DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING  
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY,  
DESIGN AND MANUFACTURING, KANCHEEPURAM

December 2019

# Certificate

This is to certify that the Internship report titled "**Recognition of Previously Identified Individuals using Affine Moment Invariants as Feature Descriptors in Gait Sequences**" by **Sandesh V Bharadwaj** (ESD15I005) to the **Indian Institute of Information Technology, Design and Manufacturing, Kancheepuram** for the partial fulfilment of "**B. Tech Electronics and Communication Engineering with Specialization in Design and Manufacturing + M. Tech Signal Processing and Communication System Design**", is a bonafide record of the research work done by him under my supervision from **May to October 2019**.

**Smt. Soma Mitra**

Senior Director  
Advanced Signal Processing Group  
Centre for Development of Advanced Computing (CDAC)  
Kolkata

## *Abstract*

Gait recognition is a biometric technology that identifies individuals in a video sequence by analysing their style of walking or limb movement. However, this identification is generally sensitive to appearance changes and conventional feature descriptors such as Gait Energy Image (GEI) lose some of the dynamic information in the gait sequence. Active Energy Image (AEI) focuses more on dynamic motion changes than GEI and is more suited to deal with appearance changes.

In this paper, we propose a new approach, which allows identifying people by analysing the dynamic motion variations and identifying people in spite of appearance changes, without using a database of predicted appearance changes. In the proposed method, the active energy image is calculated by averaging the difference frames of the silhouette sequence and divided into multiple areas. Affine moment invariants are extracted as gait features for each area. Next, matching weights for all areas are calculated based on the similarity between extracted features and those in the database. Finally, the subject is identified by the weighted combination of similarities in all areas. The CASIA-B Gait Database is used as the principal dataset for the experimental analysis.

## *Acknowledgements*

The research internship at **Centre for Development of Advanced Computing (C-DAC), Kolkata** was a great chance for learning and professional development.

I take this opportunity to express my sincere thanks to **Smt. Soma Mitra** for providing me an opportunity to do the internship work at C-DAC, Kolkata.

I express my sincere thanks and gratitude to **Mr. Kunal Chanda**, who, as my guide and mentor, took a keen interest in my work and guided me all along till the completion of the internship.

I would also like to thank **Mr. Souvik Banik**, **Mr. Washef Ahmed** and **Mr. Sayantan Bhattacharya**, who guided me and supported me during the course of the internship.

This research internship is a big milestone in my career development. The intense and thorough research has exposed me to different aspects of the subject matter and cultivated an interest to pursue it further.

I would also like to thank **IIITDM Kancheepuram** for permitting me to pursue this internship.

# Contents

<b>Certificate</b>	i
<b>Abstract</b>	ii
<b>Acknowledgements</b>	iii
<b>Contents</b>	iv
<b>List of Figures</b>	vi
<b>Abbreviations</b>	vii
<b>1 Introduction</b>	1
1.1 Gait Recognition . . . . .	1
1.2 Background . . . . .	2
1.2.1 Model-based approach . . . . .	2
1.2.2 Motion-based approach . . . . .	2
1.3 Feature Descriptors . . . . .	3
1.4 Literature Survey . . . . .	5
<b>2 Proposed Methodology</b>	32
2.1 Active Energy Image . . . . .	33
2.2 Affine Moment Invariants . . . . .	35
2.2.1 Image moments and moment invariants . . . . .	35
2.3 Estimation of matching weights . . . . .	36
2.3.1 Principal Component Analysis . . . . .	37
<b>3 Experiment and Results</b>	39
3.1 Implementation . . . . .	39
3.2 Experiment and Result . . . . .	40
<b>4 Conclusions and Future Work</b>	42
4.1 Conclusion . . . . .	42
4.2 Future Work . . . . .	43

<b>Bibliography</b>	<b>44</b>
---------------------	-----------

# List of Figures

1.1	Literature Review of Current Gait-Based Person Identification Techniques . . . . .	31
2.1	Gait silhouette sequence of individual . . . . .	33
2.2	Difference Images between consecutive significant frames of gait sequence . . . . .	33
2.3	Active Energy Image . . . . .	34
2.4	Segmented AEI (6 segments) . . . . .	34
3.1	Best Case Results (K=23, M=5) . . . . .	41

# Abbreviations

<b>AEI</b>	Active Energy Image
<b>GEI</b>	Gait Energy Image
<b>AMI</b>	Affine Moment Invariant
<b>PCA</b>	Principal Component Analysis
<b>CCR</b>	Correct Classification Rate

# **Chapter 1**

## **Introduction**

In the modern world, reliable recognition of individuals has become a fundamental requirement in various real-time applications such as forensics, international travel and surveillance. The motivations for using biometrics are diverse and often overlap. They include improving the convenience and efficiency of routine access transactions, reducing fraud, and enhancing public safety and national security.

### **1.1 Gait Recognition**

Gait recognition is a biometric recognition technique that recognizes individuals based on their walk cycle i.e. gait, and has been a topic of continued interest for person identification due to the following reasons:

- First, gait recognition can be performed with low-resolution videos with relatively simple instrumentation.
- Second, gait recognition can work well remotely and perform unobtrusive identification, especially under conditions of low visibility.
- Third, gait biometric overcomes most of the limitations that other biometric identifiers suffer from such as face, fingerprint and iris recognition which have certain hardware requirements that add to the cost of the system.

- Finally, gait features are typically difficult to impersonate or change, making them somewhat robust to appearance changes.

## 1.2 Background

In the last two decades, significant efforts have been made to develop robust algorithms that can enable gait-based person recognition on real-time data. Modern gait recognition methods can be classified into two major groups, namely model-based and motion-based methods.

### 1.2.1 Model-based approach

In model-based methods, the human body structure or motion is described using a mathematical model and the image features are extracted by measuring the structural components of models or by the motion trajectories of the body parts. These models are often simplified based on justifiable assumptions such as the system only pathologically normal gait.

Such systems normally consists of gait capture, a model(s), a feature extraction scheme, a gait signature and a classifier. The model can be 2-dimensional or 3-dimensional structural (or shape) model and motion model that lays the foundation for the extraction and tracking of a moving person. The main advantages of this approach is that it can reliably handle occlusion (especially self-occlusion), noise, scale and rotation well. However, system effects such as viewpoint invariance and effects of physiological, psychological and environmental changes are major limitations in implementing a model-based recognition system.

### 1.2.2 Motion-based approach

Motion-based methods consider the human gait cycle as a sequence of images and extract binary silhouettes from these images. Motion-based approaches are insensitive to the quality of images and have the added advantage of low computational cost compared to model-based approaches.

A baseline algorithm proposed by Sarkar et al. [11] uses silhouettes as features themselves, scaling and aligning them before use. Bobick and Davis [12] proposed the motion-energy image (MEI) and motion-history image (MHI) to convert the temporal silhouette sequence to a signal format. Han and Bhanu [13] used the idea of MEI to propose the Gait Energy Image (GEI) for individual recognition. GEI converts the spatio-temporal information of one walking cycle into a single 2D gait template, avoiding matching features in temporal sequences. GEI is comparatively robust to noise, but loses dynamic variations between successive frames.

Zhang, Zhao and Xiong [14] proposed an active energy image (AEI) method for gait recognition, in which the active regions of a gait sequence are extracted by calculating the difference between two adjacent silhouette images in a gait sequence. AEI focuses more on dynamic regions than GEI, and can alleviate the effect caused by low quality silhouettes. Current research on gait representation include Gait Entropy Image (GEnI), frequency-domain gait entropy (EnDFT), gait energy volume (GEV) etc., which focus more on dynamic areas and reducing view-dependence of traditional appearance-based techniques.

### 1.3 Feature Descriptors

A feature descriptor is a metric or a quantifiable value used to describe an image at a high level perspective. Features related to color, texture, shapes, color blobs, edges and other interest points are contained in an image. Feature descriptors need to be unique but repeatable as well, meaning that the same physical interest points in the image must be detectable under different viewing conditions. Descriptors need to be of proper dimensions as well; large descriptors will have higher computational costs, but small descriptors may discard useful information.

Feature descriptors are of the following types:

1. **Color Descriptors:**

- *Histograms*: It is the distribution of number of pixels of an image. The number of elements in a histogram relates to the number of bits in each pixel of an image.
- *Color Coherent Vector (CCV)* [22] : Each histogram bin is partitioned into two types; coherent type contains pixel value belonging to a large informally colored region.
- *Color Moments and Moment invariants* [23]: Color moments are shift-invariant measures that centralize color distribution in an image. Moment invariants are derived from these color moments, where only central moments are invariants themselves
- *Color Scale Invariant Feature Transform (SIFT) Descriptors* [24]: The local shape of a region is described by edge orientation histograms. SIFT descriptors are shift-invariant and normalized, due to which gradient magnitude changes have no effect. However, these descriptors are not invariant to light color changes, as the R, G and B channels are combined to form the intensity channel.

## 2. Texture Descriptors:

- *Gray level co-occurrence matrix (GLCM)*[25]: Used in motion estimation of images to extract second order statistical texture features.
- *Haralick Texture Feature*[26]: Used for image classification, these 13 features capture information about the patterns emerging in patterns of texture. These features are calculated by using co-occurrence matrix.

## 3. Visual Descriptors:

- *Visual color descriptors*
- *Visual texture descriptors* [28]
- *Visual shape descriptors* [27]

## 4. Frequency Domain Descriptors

## 1.4 Literature Survey

Citation	Date of Publication	Title of Paper	Defn of key terms and concepts	Objectives from Abstract and Conclusion (point-wise)	Research Methodology (include steps of Algorithm point-wise)	Summary of Research Results
[1]	April 23, 2018	<b>Person Identification from Partial Gait Cycle Using Fully Convolutional Neural Network</b>	<ul style="list-style-type: none"> <li>• <b>Gait Recognition</b> - to recognize people by the way they walk</li> <li>• <b>FCN</b> - Fully Convolutional Neural Network</li> <li>• <b>GEI</b> - <u>Gait Energy Image</u>, obtained by averaging binary silhouette over one gait cycle</li> <li>• <b>ReLU</b> - Rectified Linear Unit, a piecewise activation function used in CNNs.</li> <li>• <b>IC-GEI</b> - Incomplete GEI</li> <li>• <b>RC-GEI</b> - Reconstructed Complete GEI</li> <li>• <b>TC-GEI</b> - True Complete GEI</li> <li>• <b>ITCNet</b> - Incomplete -to-</li> </ul>	<ul style="list-style-type: none"> <li>• To identify individuals from gait features when <u>few (or single) frame(s)</u> is available .</li> <li>• To build a fully convolutional neural network for <u>GEI reconstruction</u> from <u>incomplete gait cycle</u>; this is done by training several <u>auto-encoder</u>s independently and then combining them into a model.</li> <li>• Fully reconstruct a <u>true GEI</u> from incomplete</li> </ul>	<p>1. <u>Incomplete GEI computed</u> by averaging silhouette images.</p> $GEI(x, y) = \frac{1}{N} \sum_{i=1}^N f_i(x, y)$ <p>Where <math>f(x, y)</math> denotes <u>binary value</u> of pixel in position <math>(x, y)</math> at time <math>t</math>. <math>N</math> is <u>no. of frames</u> in gait cycle.</p> <p>2. <u>Training for Incremental GEI reconstruction</u> approach using <u>9 FCNs</u> that each single FCN enhances the quality of input GEI.</p> <ul style="list-style-type: none"> <li>• Each FCN is trained on <u>different types</u> of GEI, but has the <u>same architecture</u>.</li> <li>• Structure of FCN – <u>encoder</u> (convolutional) part and <u>decoder</u> (deconvolutional) part. <ul style="list-style-type: none"> <li>◦ <b>Encoder</b> - <u>three convolutional layers</u>. Each layer is followed by a <b>ReLU</b> (activation function), and a <b>pooling layer</b> (down-sampling unit).</li> <li>◦ <u>Batch normalization</u> and <u>dropout</u> is then used to prevent <u>overfitting</u> and increase <u>rate of convergence</u>.</li> <li>◦ <b>Decoder</b> - <u>up-sampling layer</u>, <u>convolutional layer</u>, <b>ReLU</b> layer, batch normalization and dropout.</li> <li>◦ <u>Up-sampling layer</u> is used to remove <u>convolutional artifacts</u> like <u>checkerboard patterns</u>, which occur if</li> </ul> </li> </ul>	<p>1. Reconstruction:</p> <ul style="list-style-type: none"> <li>a. <b>Increasing no. of frames</b> of <b>IC-GEIs</b> generates <b>RC-GEIs</b> closer to <b>TC-GEIs</b>. However, upward trend becomes slower as 18/20f-GEIs are very close to true GEIs.</li> <li>b. <b>OULP</b> – Min. Accuracy - <u>89.35%</u> Max. Accuracy - <u>96.15%</u></li> <li>c. <b>Casia-B</b> – Min. Accuracy - <u>92.27%</u> Max. Accuracy - <u>94.48%</u></li> </ul> <p>2. Recognition:</p> <ul style="list-style-type: none"> <li>a. Recognition between OULP and Caisa -B was compared using <u>Euclidean distance</u> as similarity metric between GEI samples.</li> <li>b. ITCNet <u>improved</u> identification performances for IC-GEIs, especially of those with <u>fewer frames</u>.</li> <li>c. As no. of frames <u>increase</u>, rank-1 and rank-5 identification gets <u>closer</u> to that of TC-GEIs. (CMC curve)</li> </ul>

		<p>complete GEI Network</p> <ul style="list-style-type: none"> <li><b>ROC – Receiver Operating Characteristic Curve</b>, used for measuring performance of classification model, plotting true positive rate vs false positive rate.</li> <li><b>CMC – Cumulative Matching Characteristic</b> curve, metric used to measure performance of identification and recognition algorithms based on precision for each rank.</li> </ul>	<p>te GEI, despite starting frame and number of available frames.</p>	<p><u>transposed convolutional layer</u> is <b>directly used</b> for the encoder output.</p> <ul style="list-style-type: none"> <li>○ <b>Sigmoid function</b> is placed at output of decoder to get a <u>gray-scale output</u> image.</li> <li>• Input image - <math>64 \times 64</math> pixels,</li> <li>• Convolution layer parameters – <ul style="list-style-type: none"> <li>○ Kernel size - <math>(4,4)</math></li> <li>○ Strides - <math>(1,1)</math></li> <li>○ Output Dimensions - <math>128, 64, 32</math></li> </ul> </li> <li>• Pooling Layer: <ul style="list-style-type: none"> <li>○ Kernel size and pooling strides - <math>(2,2)</math> each</li> </ul> </li> </ul> <p>3. <u>ITCNet converter</u>:</p> <ul style="list-style-type: none"> <li>• The trained auto-encoders from Step 2 are <u>stacked</u> together.</li> <li>• Input – any type of incomplete GEI</li> <li>• Output – complete GEI</li> <li>• Auto-encoder in layer <math>i</math> maps <math>mf</math>-GEI to <math>nf</math>-GEI, where <math>n - m = T/10</math>, <math>T</math> = <u>gait cycle length</u>.</li> <li>• <u>Weight initialization</u> using <u>Gaussian distribution</u> with std. deviation <math>\sigma_c</math></li> </ul> $\sigma_c = \sqrt{\frac{2}{f_w * f_h * f_t * f_d}}$ <p>Where <math>f_w</math> – spatial width  <math>f_h</math> – spatial height  <math>f_t</math> – time extent of filter  <math>f_d</math> – no. of filters in convolutional layer</p> <ul style="list-style-type: none"> <li>• <u>Mean Square Error (L)</u> – Loss value between RC-GEI and complete GEI</li> </ul> $L = \sum_{i=1}^N \ x - y\ ^2$ <p>Where <math>x</math> – RC-GEI  <math>y</math> – original complete GEI</p> <ul style="list-style-type: none"> <li>• MSE minimized by using stochastic gradient descent, using Adam optimizer.</li> </ul>	<p>comparison)</p> <p>d. For <u>verification</u>, <u>similarity</u> between two GEIs is computed. If this score is <u>less than or equal to threshold</u>, this is regarded as <u>positive pair</u> (same subject).</p> <p>3 Conclusion – ITCNet successfully reconstructed TC-GEI from IC-GEI regardless of starting frame and no. of available frames. Model improves recognition rate, particularly when <u>only 0.1 part</u> of gait cycle is available.</p>
--	--	---	---	--	--

					<ul style="list-style-type: none"> <li>○ Optimizer parameters – <math>B_1 = 0.8, \beta_2 = 0.99, \epsilon = 10 \times 10^{-8}</math></li> <li>● Learning rate- optimal solution through quantitative testing together with weight decay and momentum parameter.</li> <li>○ Initial value – <math>10^{-3}</math>, decreases every 5 epochs by factor of 10</li> <li>● Reconstruction accuracy – portion of errors <b>&lt;0.08</b> between pixels in TC-GEI and RC-GEI</li> <li>● <b>50 epochs</b>, batches of size 80 for training.</li> <li>● Dropout probability – <b>0.5</b> during training</li> </ul> <p>4. Datasets used –</p> <ol style="list-style-type: none"> <li><i>OULP</i> – OU-ISIR Large Population Dataset, section A used for gait recognition, gait data selected at <math>85^\circ</math>. <b>2254</b> subjects for training set, <b>1000</b> subjects in <b>validation</b> and <b>test</b> sets.</li> <li><i>Casia-B</i> – normal walking section under <math>90^\circ</math> used, with <b>124</b> subjects. Each subject has <b>1 gallery set</b> and <b>5 probe sets</b>.</li> </ol> <p>5. Evaluation Metrics –</p> <ol style="list-style-type: none"> <li><b>Reconstruction</b>:             <ol style="list-style-type: none"> <li>MSE value is computed.</li> <li><b>SSIM</b> – Structural Similarity Index, which is a decimal value between <b>-1</b> and <b>1</b>. If <b>identical</b>, output is <b>1</b>.</li> </ol> </li> <li><b>Recognition</b> –             <ol style="list-style-type: none"> <li>For verification, <b>receiver operating characteristic</b> (ROC) curve and <b>equal error rates</b> (EER)</li> <li>For identification, <b>Rank-1/Rank-5</b> metrics and <b>cumulative matching characteristic</b> (CMC) curve.</li> </ol> </li> </ol>	
--	--	--	--	--	--	--

[2]	Sept. 2018	<b>Appearance and Gait-Based Progressive Person Re-identification for Surveillance Systems</b>	<ul style="list-style-type: none"> <li><b>Progressive search</b> – techniques that can be used</li> <li><b>Gait feature</b></li> <li><b>LSTM</b> – Long short-term memory, special kind of RNN, capable of learning long-term dependencies.</li> <li><b>Siamese network</b> – neural network that contain two or more identical sub-networks.</li> <li><b>GoogLeNet (Inception)</b> – A 22-layer deep CNN model developed by Google, using a CNN inspired by LeNet and implemented an inception module, which used batch normalization, image distortions and several</li> </ul>	<ul style="list-style-type: none"> <li>To implement a progressive person re-identification (<b>PROPRID</b>) approach to simultaneously improve timeline ss and accuracy of identifying the target person</li> <li>To utilize the appearance of the target person as a course filter to reduce computational complexity of the fine search</li> </ul>	<ol style="list-style-type: none"> <li><b>Person image sequences</b> are taken as <u>input</u>, considered as progressive processes</li> <li><b>Appearance-based person coarse filtering</b> – Hair, clothes, skin and other visual information used to filter out dissimilar people in large-scale surveillance videos. <ul style="list-style-type: none"> <li><u>Low-level, mid-level and high-level</u> appearance features are used as coarse filter</li> <li>For low-level and mid-level features – <u>CN, HOG, and LOMO descriptors</u> <ul style="list-style-type: none"> <li>➤ CN is benchmark for person Re-ID on Market-1501 dataset, quantized by BOW model for fast color feature matching.</li> <li>➤ HOG is used to characterize poses and appearances (shape feature)</li> <li>➤ LOMO is used for stable texture representation against changes in the viewpoint.</li> </ul> </li> <li>For high level features, <b>GoogLeNet</b> is used as feature extractor. <ul style="list-style-type: none"> <li>➤ GoogLeNet used ImageNet as training data, cannot apply directly to person re-id.</li> <li>➤ GoogLeNet model is fine-tuned using <u>PRID-2011</u> and <u>iLIDS-VID</u> datasets.</li> <li>All four types of features are integrated by <u>distance-level fusion</u>.</li> </ul> </li> </ul> </li> <li><b>Human silhouette extraction</b> done by <b>Mask R-CNN</b> to generate GEI on PRID-2011 and iLIDS-VID datasets. <ul style="list-style-type: none"> <li>Mask R-CNN is trained (fine-tuned) using <b>Casia-B</b> dataset. It cannot be used</li> </ul> </li> </ol>	<ol style="list-style-type: none"> <li><b>Evaluations of Appearance-Based Person Re-Identification:</b> <ul style="list-style-type: none"> <li>5 approaches evaluated on PRID-2011 and iLIDS-VID: <ul style="list-style-type: none"> <li>➤ CN</li> <li>➤ HOG</li> <li>➤ LOMO</li> <li>➤ GoogLeNet: <u>1024-D feature vector</u> from last pooling layer obtained to represent semantic features.</li> <li>➤ <b>FMLF</b>: <u>Multi-level features combined</u> into appearance features. <i>Fusion weights obtained on validation subset of training data.</i></li> <li>○ GoogLeNet Features are more <u>discriminative</u> and <u>robust</u> than CN, HOG and LOMO descriptors.</li> <li>○ Multi-level fusion features (FMLF) achieve <b>highest accuracy</b>.</li> </ul> </li> </ul> </li> <li><b>Evaluation of Progressive Person Re-Identification:</b> <ul style="list-style-type: none"> <li>Two methods i.e. <b>FMLF</b> and <b>PROPRID</b> were compared on PRID-2011 and iLIDS-VID datasets with the following fusion weights: <ul style="list-style-type: none"> <li>➤ FMLF: <ul style="list-style-type: none"> <li>PRID-2011 – CN = 0.1</li> <li>HOG = 0.2</li> <li>LOMO = 0.4</li> </ul> </li> </ul> </li> </ul> </li> </ol>
-----	------------	--	--	--	--	---

		<p>very small convolutions to reduce no. of parameters.</p> <ul style="list-style-type: none"> <li>• <b>CN</b> – Color Name Descriptor, assigning names to colors for identification</li> <li>• <b>BOW</b> – Bag-of-words model, a vector of occurrence counts of local image features</li> <li>• <b>HOG</b> – Histogram of oriented gradients</li> <li>• <b>LOMO</b> – Local Maximal Occurrence, analyses the horizontal occurrence of local features and maximizes the occurrence for stable representation in viewpoint changes.</li> <li>• <b>Mask R-CNN</b> – Mask regional CNN, consisting</li> </ul>	<p>directly as it is trained on COCO dataset.</p> <p><b>4. Gait Cycle Feature extraction:</b></p> <ul style="list-style-type: none"> <li>○ Gait sequences are processed by <b>ResNet-50</b> model to gain spatial features of each image.</li> <li>○ <u>Spatial features</u> fed into <u>LSTM</u> layers to obtain <u>temporal feature sets</u>.</li> <li>○ Feature sets summarized by <u>mean-pooling</u> to generate sequence-level features.</li> <li>○ <b>Gait recognition</b> performed by computing <u>cosine distance</u> between two gait sequences.</li> </ul> <p><b>Siamese LSTM (SLSTM)</b> Network is used for learning subtle periodic features.</p> <ul style="list-style-type: none"> <li>○ Consists of two LSTM sub-networks, parallel in structure and same parameters.</li> <li>○ SLSTM takes <u>pairs of gait sequences</u> as inputs, maps each gait sequence to a feature vector.</li> <li>○ For all gait sequences, two sequences belonging to same identity selected randomly as <u>positive pairs</u>, whereas two sequences belonging to different identities selected randomly as <u>negative pairs</u>.</li> <li>○ <b>Contrastive loss</b> layer applied to connect the two branches of SLSTM network.</li> </ul> $D(v_i, v_j) = \ v_i - v_j\ _2$ <p>Where <math>v_i</math> and <math>v_j</math> are sequence-level feature vectors</p> <p>The contrastive loss function is:</p> $L(v_i, v_j, Y) = \frac{Y}{2} D^2 + (1 - Y) \frac{1}{2} [\max(m - D, 0)]^2$	<p>GoogLeNet = 0.3</p> <p>iLIDS-VID – CN = 0.1 HOG = 0.2 LOMO = 0.1 GoogLeNet = 0.6</p> <p>➤ PROPRID: PROPRID <u>fuses</u> scores of <u>FMLF</u> and <u>human gait</u> models. Before late fusion, similarity vectors are <u>normalized</u> to (0,1) range. PRID-2011: FMLF – 0.65 Gait – 0.35</p> <p>iLIDS-VID: FMLF – 0.6 Gait – 0.4</p> <p>○ PROPRID achieves <u>81.00%</u> and <u>62.33%</u> for <u>Rank-1</u> recognition on both datasets, which exceeds second-best methods by <u>4.00%</u> and <u>0.33%</u> respectively.</p> <p>3. Comparison with state-of-the-art methods:</p> <ul style="list-style-type: none"> <li>○ <b>ASTPN</b> – Spatial and Temporal Attention Pooling Network, enables feature extractor to be aware of current input video sequence.</li> <li>○ <b>RCN</b> – Recurrent CNN with temporal pooling.</li> <li>○ <b>RFA</b> – Recurrent Feature Aggregation network, based on LSTM.</li> </ul>
--	--	---	---	--

		<p>of a region proposal netowrk and a binary mask classifier.</p> <ul style="list-style-type: none"> <li>• <b>ResNet-50</b> – 50-layer Residual Network, uses skip connection to propagate information over layers.</li> <li>• <b>Contrastive Loss</b> – Distance-based loss function, used to find similarity between image pairs.</li> <li>• <b>Backpropagation</b> – algorithms used to train neural networks by using gradient descent.</li> <li>• <b>FMLF</b>: Fusion of Multi-Level Features, features are combined into comprehensive appearance features. Weights obtained</li> </ul>	<p>where <math>\mathbf{m}</math> is the min. distance margin of different identities' gait sequences. Total loss is sum of all gait sequence pair losses.</p> <ul style="list-style-type: none"> <li>○ In training, both LSTM sub-networks are optimized simultaneously with weight sharing.</li> <li>○ <u>Pairwise gait sequences</u> with similar/dissimilar labels fed into the two LSTM sub-networks <u>separately</u>.</li> <li>○ Outputs are <u>combined by contrastive loss layer</u> to output contrastive loss.</li> <li>○ <u>Backpropagation with contrastive loss</u> is used to <u>train model</u>.</li> <li>○ During <u>testing</u>, Siamese architecture and loss function is <b>discarded</b>.</li> <li>○ Fine-tuned ResNet-50 model and LSTM sub-netowrk <u>jointly used</u> as feature extractor.</li> <li>○ Final decision made by <i>comparing cosine distances</i>.</li> <li>○ Comparison is done using CMC curves.</li> </ul>	<p>PROPRIID achieves the best performance in terms of rank-1 recognition for both datasets.</p>
--	--	---	--	---

			<p>on validation subset of training data.</p> <ul style="list-style-type: none"> <li>• <b>PRID-2011</b> - Person Re-ID Dataset, consists of images extracted from multiple person trajectories recorded from two different, static surveillance cameras. Images contain viewpoint change and differences in illumination, background and camera characteristics.</li> <li>• <b>iLIDS-VID</b> - iLIDS Video re-Identification dataset, 600 image sequences of 300 different pedestrian observed across 2 disjoint cameras</li> </ul>			
[3]	April 2014	<b>Gait-assisted Person Re-identification</b>	<ul style="list-style-type: none"> <li>• <b>Person Re-ID</b></li> <li>• <b>GEI</b> – Gait Energy Image</li> <li>• <b>Frame</b></li> </ul>	<ul style="list-style-type: none"> <li>• Investigate if gait extracted from real-world</li> </ul>	<ol style="list-style-type: none"> <li>1. Silhouette sequence selection is used to ensure reliable gait extraction <ul style="list-style-type: none"> <li>a. Non-zero pixels &gt; 40% of image size for reliable gait extraction.</li> </ul> </li> </ol>	<ol style="list-style-type: none"> <li>1. SAIVT SoftBio Dataset - <ol style="list-style-type: none"> <li>a. Combined Color and Gait Features model for Re-ID: <ul style="list-style-type: none"> <li>i. 150 subjects in</li> </ul> </li> </ol> </li> </ol>

	<b>Wide Area Surveillance</b>	<b>Difference Energy Images (FDEI) -</b> Gait feature used to deal with incomplete silhouettes while retaining shape and motion changes, by constructing multiple thresholded GEIs from single gait cycle, and summing of these GEIs with differences between consecutive silhouettes. Positive values are retained, and negative values discarded. <b>Z normalization -</b> standardization, by converting all indicators to a common scale with average of zero and standard deviation of one i.e	video sequences can be successfully leveraged as an additional feature along with appearance features for person re-ID. • Identify robust gait features that can be extracted from noisy and incomplete silhouettes and retain discriminative capability for Re-ID. • Implement a sparsified representation-based gait recognition method, where probe gait features	<p>b. Aspect ratio of silhouettes over sequence of frames is used to estimate gait period. Assumption - silhouettes extracted and pre-processed.</p> <p>Pre-processing – size normalization and horizontal alignment of silhouettes.</p> <p>2. Gait Period estimation:</p> <p>Gait period estimation is performed using lower half of silhouette region</p> <p>a. Aspect ratio of silhouette over frame sequence represented as 1D temporal signal, which is z-normalized and smoothed using moving average filter.</p> <p>b. Peaks in signal magnified by computing autocorrelation sequence; first derivative of autocorrelation signal used to detect zero-crossings.</p> <p>c. Zero-crossing of positive and negative peaks used to compute distance between prominent peaks</p> <p>d. Average of distances between consecutive peaks =&gt; Gait period in number of frames.</p> <p>3. Gait Features –</p> <p>a. GEI – computed by averaging the silhouettes in spatial domain over gait cycle</p> $GEI[i, j] = \frac{1}{N} \sum_{t=1}^N S_t[i, j]$ <p>Where <math>S_t</math> is silhouette at frame t, (i,j) are spatial image coordinates N is estimated gait cycle period.</p> <p>b. FDEI - deals with incomplete silhouettes while retaining changes in shape and motion.</p> <p>Steps:</p>	<p>8 camera views; only 122 subjects in atleast 2 camera views.</p> <p>ii. Constraint: same subjects' gallery and probe camera views are different.</p> <p>iii. Simple background subtraction followed by low level image processing to extract silhouette sequences for gait feature extraction.</p> <p>iv. 5 randomly selected frames from a sequence used to extract color features</p> <p>v. All frames used to extract gait features.</p> <p>vi. As weight assigned to gait features increases, significant boost to performance.</p> <p>vii. Any weight combination of color and sparsified representation gait similarity-based Re-ID outperforms color and nearest-neighbour gait similarity.</p> <p>viii. Difference in performance between GEI and FDEI is</p>
--	-------------------------------	---	--	---	--

		$x_i = \frac{x_i - \mu}{\sigma}$ <ul style="list-style-type: none"> <li>• <b>SAIVT SoftBio Dataset</b></li> <li>• <b>MCID dataset</b> - multi-camera Re-ID dataset.</li> <li>• <b>Sparsified Representation</b></li> <li>• <b>Bhattacharyya Distance</b> - Measures the similarity of two probability distributions</li> </ul> <p>For probability distributions <math>p</math> and <math>q</math> over the same domain <math>X</math>,</p> $D_B(p, q) = -\ln \sum_{x \in X} p(x)q(x)$ <p>Where</p> $BC(p, q) = \sum_{x \in X} p(x)q(x)$ <p>is the <b>Bhattacharyya Coefficient</b> for discrete probability distributions.</p>	<p>are represented as linear combination of gallery gait features, using</p> <ul style="list-style-type: none"> <li>◦ Gait Energy Images</li> <li>◦ Frame Difference Energy Images</li> </ul>	<p>i. Construct GEIs from sub-cycles -</p> $I_c(i, j) = \frac{1}{N_c} \sum_{t=1}^{N_c} S_t(i, j)$ <p>Where <math>N_c</math> is period of sub-cycle</p> <p>ii. Thresholding GEIs to remove noise, creating dominant GEIs (DGEIs) for each sub-cycle -</p> $DGEI_c(i, j) = \text{threshold}(I_c(i, j))$ <p>iii. Frame difference computed by subtracting consecutive silhouettes and including only positive portion of difference in feature -</p> $DS_t(i, j) = \begin{cases} 0, & \text{if } S_t(i, j) \geq S_{t-1}(i, j) \\ S_{t-1}(i, j) - S_t(i, j), & \text{otherwise} \end{cases}$ <p>iv. FDEI generated by summation of +ve frame difference and corresponding sub-cycle DGEI -</p> $FDEI(i, j) = DS_t(i, j) + DGEI_c(i, j)$ <p>4. Feature Matching for Re-ID -</p> <p>a. Color features: Distribution of colors characterized using weighted HSV algorithm.</p> <p>i. Silhouette divided into three parts - head, torso and legs, by detection of one vertical axis of symmetry and 2 horizontal axes of symmetry.</p> <p>ii. Histogram for each body part weighted by distance from axis of symmetry.</p> <p>iii. Each histogram concatenated channel-wise to generate single color feature descriptor.</p> <p>b. Combined similarity measure: Similarity</p>	<p>negligible.</p> <p>b. Only Gait features for Re-ID:</p> <ul style="list-style-type: none"> <li>i. Re-ID only using gait features; only 23 subjects with usable gait features available across views.</li> <li>ii. Gait features outperform color features significantly</li> <li>iii. FDEI yield better Rank-1 matching accuracy than GEI, yet comparable.</li> </ul> <p>2. MCID dataset - 9-camera network in and around a building; 40 subjects, 19 in multiple cameras. <math>w_{gait} = 0.2</math> and <math>w_{color} = 0.8</math> for best Re-ID accuracy</p> <p>a. Combined Color and Gait features model for Re-ID:</p> <ul style="list-style-type: none"> <li>i. For each camera pair, only subset of probe set i.e closed probe set (intersection b/w probe and gallery set IDs) used to establish Re-ID.</li> <li>ii. Gait features of MCID extracted only from lower body silhouettes (truncated GEIs).</li> <li>iii. Performance consistent w.r.t. SAIVT SoftBio</li> </ul>
--	--	---	---	--	---

				<p>measure between gallery subject G and probe subject P is weighted sum of color feature-based similarity and gait features based similarity i.e</p> $dist(G, P) = w_{color} \cdot d_{color}(G, P) + w_{gait} \cdot d_{gait}(G, P)$ <p>Where <math>d_{color}</math> is color feature-based similarity, <math>d_{gait}</math> is gait feaure-based similarity, <math>w_{color}</math> and <math>w_{gait}</math> are weights given to color and gait similarity, such that:</p> $w_{gait} = 1 - w_{color}$ <ul style="list-style-type: none"> <li>i. Color similarity is calculated using Bhattacharyya distance. Color similarity for given probe-gallery pair is maximum similarity among all probe-gallery frame pairs.</li> <li>ii. Gait similarity is computed as reconstruction error between probe gait features and gallery gait model, utilizing sparsified representation.</li> <li>iii. Given a dictionary matrix built using labeled training features of several subjects, a test feature of i-th subject is linear combination of training images on same subject from dictionary.</li> </ul> <p><math>V = [v_1, \dots, v_n]</math> is a dictionary matrix, where each column is obtained by vectorizing gait features belonging to all gallery subjects. Multiple columns in V may belong to same gallery subject.</p> <ul style="list-style-type: none"> <li>iv. Separate dictionary is built for each gallery subject, ensures that if probe subject is as close</li> </ul>	<p>dataset results.</p> <p>b. Only Gait features model for Re-ID:</p> <ul style="list-style-type: none"> <li>i. FDEI Re-ID within 4% of color accuracy.</li> <li>ii. Combining either GEI or FDEI boots Re-ID accuracy.</li> <li>iii. Pure GEI or FDEI performs poorly on MCID relative to SAIVT SoftBio dataset, due to truncated gait features.</li> </ul> <p><b>Conclusion:</b> Sparsified representation based cross view gait recognition eliminates view angle estimation or gait feature reconstruction for matching.</p>
--	--	--	--	---	--

					<p>as possible to gallery subject in identity and view angle, probe image lies in linear span defined by subset of gait features constituting V.</p> $I_p = V \cdot \alpha$ <p>Goal – to find sparsest <math>\alpha</math> that generated <math>GEI_p</math> in V i.e. solving the <math>l_1</math>- minimization</p> $\hat{\alpha} = \arg \min \ \alpha\ _1 \text{ s.t. } I_p = V \cdot \alpha$ <p>Solved using linear programming, leveraging augmented Lagrange multipliers.</p> <p>v. <math>d_{gait}(G, P) = \ I_p - V \cdot \hat{\alpha}\ _2</math> This is an estimate of how well <math>\hat{\alpha}</math> reproduces <math>GEI_p</math></p> <p>vi. If good quality silhouette sequences available and gait features extractable, <math>w_{gait}</math> is set to non-zero value. In absence of usable gait features, Re-ID established using only color features. Acceptable quality – No. of +ve pixels &gt; 40% of total pixels in image.</p>	
[4]	November 2012	<b>Person Re-identification using View-dependent Score-level Fusion of Gait and Color Features</b>	<ul style="list-style-type: none"> <li><b>View-dependent score-level fusion</b> – adaptive weight control of gait and color features, where greater weight is given to gait features in similar-view</li> </ul>	<ul style="list-style-type: none"> <li>Person re-identification across multiple non-overlapping cameras, using a spatio-temporal histogram as a gradient-based shape and</li> </ul>	<ol style="list-style-type: none"> <li>Color and gait features, observation view extracted from image sequence; color and gait feature matching returns individual distances.</li> <li>Feature Extraction – <ul style="list-style-type: none"> <li>Gait Feature: <ul style="list-style-type: none"> <li>STHOG features used as gait feature containing both shape and motion information.</li> </ul> </li> <li>Pedestrian window size-normalized and divided into multiple spatio-temporal cells.  <math>G = [G_x, G_y, G_t]^T</math> is first order derivative for each of x-, y-, and t-axes from two successive frames.</li> </ul> </li> </ol>	<ol style="list-style-type: none"> <li>Constructed own dataset; <ol style="list-style-type: none"> <li>Walking image sequences of 27 subjects (14 for training, 13 for testing); 7 non-overlapping views ranging from front view to rear-oblique view.</li> <li>Near side view used as query view</li> <li>Size normalization window -&gt; 100x140 pixels</li> <li>Spatio-temporal</li> </ol> </li> </ol>

		<p>matching and to the color features in different-view matching.</p> <ul style="list-style-type: none"> <li>• <b>Assumptions -</b> <ul style="list-style-type: none"> <li>◦ Pedestrian detection and tracking for each camera preprocessed.</li> <li>◦ Observation view for each pedestrian given.</li> </ul> </li> <li>• <b>STHOG -</b> Spatio-temporal histogram of oriented gradient.</li> <li>• <b>Tanimoto distance -</b> Ratio of the intersecting set to the union set as a measure of similarity</li> <li>• <b>Baseline algorithm</b></li> <li>• <b>LLR -</b> Linear Logistic Regression classifier</li> <li>• <b>CMC -</b> Cumulative matching characteristics.</li> </ul>	<p>motion gait feature to discriminate persons with similar color clothing in conjunction with background edge attenuation</p> <ul style="list-style-type: none"> <li>• Implement view-dependent score-level fusion framework to adaptively control weights of gait and color features.</li> </ul>	$\phi = \tan^{-1} \left( \frac{G_y}{G_x} \right) \quad \text{where } \phi \text{ is spatial gradient.}$ $\theta = \tan^{-1} \left( \frac{G_t}{\sqrt{G_x^2 + G_y^2}} \right) \quad \text{where } \theta \text{ is temporal gradient.}$ <p>iii. Orientation of spatial and temporal gradients voted separately according to spatial gradient magnitude into individual 9-bin histograms for each cell, then normalized w.r.t. <math>L_1</math> norm.</p> <p>iv. Both histograms combined into single 18-bin histogram, and histograms from all cells combined into single histogram.</p> <p>v. Effect of background variations must be suppressed; gradient-based background attenuation rather than background subtraction.</p> <p>vi. Attenuates edge at certain position if background edge also exists at same position and edge patterns of background/input images are similar.</p> $G_s = \frac{\sqrt{G_x^2 + G_y^2}}{\left( 1 + K G_x^{B^2} e^{\frac{-z_x^2}{\sigma^2}} \right) \left( 1 + K G_y^{B^2} e^{\frac{-z_y^2}{\sigma^2}} \right)}$ <p><math>K</math> and <math>\sigma</math> are hyper parameters to control background attenuation.  <math>G_x^B</math> and <math>G_y^B</math> are horizontal and vertical gradients in background image.  <math>z_x</math> and <math>z_y</math> denote dissimilarity of horizontal and vertical</p>	<p>cell sizes -&gt; 100 pixels x 20 pixels x 3 frames, given 7 cells per window.</p> <p>e. <math>K = 5.0</math>; <math>\sigma = 10.0</math></p> <p>f. Length of subsequence = 15 frames</p> <p>Repeated experiments for 20 trials, random choice of training and test subjects. Avg. performances of 20 trials evaluated using CMC.</p> <p>2. Results –</p> <p>a. Performance using STHOG feature better than color feature for similar observation views; lower than color features for different observation views</p> <p>b. Adaptive fusion controls weight of STHOG and color features; achieves best performance</p>
--	--	--	--	---	---

				<p>edge patterns in background and input images.</p> $z_x = \max (\ I_{x,y}^B - I_{x,y}\ , \ I_{x+1,y}^B - I_{x,y}\ )$ $z_y = \max (\ I_{x,y}^B - I_{x,y}\ , \ I_{x,y+1}^B - I_{x,y}\ )$ <p>Where <math>I_{x,y}^B</math> and <math>I_{x,y}</math> are pixel intensities at <math>(x,y)</math> in background and input images.</p> <p>b. Color Feature:</p> <ul style="list-style-type: none"> <li>i. Color histogram used with modified HSV system i.e. Histogram for each cell with 10 bins; 7 hue bins and 3 value bins (white, gray, and black).</li> <li>ii. Region-based feature; hence, background region masked out by background subtraction during voting.</li> <li>iii. Color histograms normalized for individual cells and concatenated into single histogram.</li> </ul> <p>3. Matching –</p> <p>a. Features are synchronized when matching using baseline algorithm as follows:</p> <ul style="list-style-type: none"> <li>i. Query sequence segmented into subsequences with certain length (gait period)</li> <li>ii. Each segment phase-synchronized with other sequence by shifting frame to minimize total distance b/w them</li> <li>iii. Baseline algorithm utilizes Tanimoto distance; replaced with <math>L_{0.5}</math> norm to improve robustness wrt outliers.</li> <li>iv. Distance of individual sub-sequences computed and minimum</li> </ul>	
--	--	--	--	--	--

					<p>value chosen as final distance.</p> <p>4. View-dependent Score-level Fusion –</p> <ul style="list-style-type: none"> <li>a. Score-level fusion function maps 2-D distance vectors derived from STHOG and color features into single distance.</li> <li>b. Fusion function trained separately for each discrete observation view difference, due to reliability of STHOG features being highly dependent on observation view difference.</li> <li>c. LLR of likelihood ratio b/w +ve and -ve samples is used for simplicity and high generalization capability.</li> <li>d. Relative distance of each negative sample w.r.t corresponding positive sample is input into fusion framework.</li> <li>e. After training, input 2-D distance vector converted into single fused distance.</li> <li>f. Minimum fused distance between person in one camera image to another camera image used for Re-ID.</li> </ul>	
[5]	June 2013	<b>Gait-Based Person Identification Robust to Changes in Appearance</b>	<ul style="list-style-type: none"> <li>• <b>Affine moment invariants</b> – moment-based descriptors , invariant under general affine transform. Centralized moment of order(<math>p+q</math>) of object O:</li> </ul> $\mu_{pq} = \sum_{(x,y) \in \Omega} x^p y^q$	<ul style="list-style-type: none"> <li>• Implement a robust method of person identification in spite of changes in appearance – without using database of predicted</li> </ul>	<p>1. Average GEI is calculated over one gait cycle and the human body area is divided into multiple areas (5).</p> <ul style="list-style-type: none"> <li>a. Silhouette is extracted by background subtraction</li> <li>b. Human body area is scaled to uniform height and set to 128 pixels</li> </ul> $\bar{I}(x,y) = \frac{1}{T} \sum_{t=1}^T I(x,y,t)$ <p>Where T is no. of frames in one gait cycle  <math>I(x,y,t)</math> represents intensity of pixel <math>(x,y)</math> at time t.</p> <ul style="list-style-type: none"> <li>c. High intensity values – body parts that move little during gait cycle i.e. head,</li> </ul>	<p>1. First Experiment (robust to noise and deficit in silhouette images)</p> <ul style="list-style-type: none"> <li>a. CASIA-B: Each subject has 6 walking sequences</li> <li>i. First four sequences are used for training datasets.</li> <li>ii. CCR calculated by dividing the sequences of each subject into two sets.</li> </ul>

		<p>Where <math>x_g</math> and <math>y_g</math> define center of object.</p> $A_1 = \frac{1}{\mu_{00}} (\mu_{20})$ <p>Where <math>A_1</math> is the first affine moment.</p> <ul style="list-style-type: none"> <li>• <b>CCR</b> – Correct classification rate</li> </ul>	<p>d appearance changes - by dividing the human body image into multiple areas and performing feature extraction.</p>	<p>torso;</p> <p>d. Low intensity values – body parts that move constantly i.e. lower parts of legs, arms;</p> <p>e. Gait cycle estimated by calculating first affine moment invariant at each frame.</p> <p>2. Estimation of matching weight:</p> <ul style="list-style-type: none"> <li>a. Affine moment invariants in database &amp; of subject are whitened at each area</li> <li>b. Distance between features of subject and those of all datasets in database determined (<math>d_{n,s}^k</math>)</li> <li>-</li> </ul> $d_{n,s}^k = \sqrt{\mu_{n,s}}$ <p>Where <math>\mu_{n,s}</math> show whitened affine moment invariants of subject and person in database.</p> <p>c. Whitening is done by –</p> <ul style="list-style-type: none"> <li>i. Applying principal component analysis to calculated affine moment invariants</li> <li>ii. Projecting them to new features space</li> <li>iii. Normalizing projected affine features based on corresponding eigenvalue</li> </ul> $1 \leq n \leq N$ <p>Where N is number of people in database</p> $1 \leq s \leq S$ <p>Where S is no. of sequences of each person</p> $1 \leq k \leq K$ <p>Where K is no. of divided areas</p> <p>d. <math>d_{n,s}^k</math> calculated between features of subject and those of each sequence in database at each area.</p> <p>3. Estimate matching weights based on similarity between features of subject and database.</p>	<p>iii. Testing by 2-fold cross-validation method (124x3 sequences for training, rest for testing)</p> <p>iv. CCR calculated by dividing no. of test datasets correctly classified by that of all test datasets.</p> <p>Proposed method shows highest performance of 97.7% (K=17, M=45). CCR increases with K and M</p> <p>b. CASIA-C:</p> <ul style="list-style-type: none"> <li>i. Sequences divided into two sets.</li> <li>ii. Testing through 4-fold cross validation (153x3 sequences used for training, rest for testing)</li> <li>iii. Highest performance of 94% (K=7, M=65). Although silhouette images of worse quality than CASIA-B, person identification with high performance.</li> </ul> <p>2. Person Identification robust to appearance changes:</p> <ul style="list-style-type: none"> <li>a. CASIA-B-BG and</li> </ul>
--	--	--	---	---	--

					<p>High matching weights – less appearance changes Low matching weights – more appearance changes</p> <p>a. At each area k, select sequences from database if <math>d_{n,s}^k &lt; d_{\min}^k</math> .  <math>\acute{d}_n^k = \min \acute{d}_n^k</math>  <math>\acute{d}_n^k = \frac{1}{S} \sum_{s=1}^S d_{n,s}^k</math></p> <p>b. Consider that similarities of non-selected sequences are low, so distance of these sequences are redefined as <math>d_{\max}</math>.  <math>d_{n,s}^k = d_{\max}</math>  This process allows setting low similarities to areas of each sequence in database.</p> <p>c. Repeat a and b for all areas.</p> <p>d. Sum of distances for all areas calculated by  <math>D_{n,s} = \sum_{k=1}^K d_{n,s}^k</math> and k-nearest neighbor method used to identify subject.  (k=1 in experiment)</p>	<p>CASIA-B-CL dataset used.</p> <p>b. K=17 and M=45 for CCR calculations, similar to first experiment.</p> <p>c. CASIA-B-BG:</p> <ul style="list-style-type: none"> <li>i. Handbag (42 sequences)</li> <li>ii. Shoulder bag (171 sequences)</li> <li>iii. Backpack (30 sequences)</li> <li>iv. Others (3 sequences)</li> </ul> <p>CCR = 91.9% CCR with method without matching weights = 20.2%</p> <p>d. CASIA-B-CL:</p> <ul style="list-style-type: none"> <li>i. Thin coat with hood (30 sequences)</li> <li>ii. Coat (24 sequences)</li> <li>iii. Coat with hood (16 sequences)</li> <li>iv. Jacket (70 sequences)</li> <li>v. Down jacket (62 sequences)</li> <li>vi. Down jacket with hood (28 sequences)</li> <li>vii. Down coat with hood (16 sequences)</li> </ul> <p>CCR = 79% CCR without matching weights method = 22.4%</p>
[6]	September 2013	<b>Person re-identification using height-based gait in</b>	<ul style="list-style-type: none"> <li>• <b>KL-Divergence (relative entropy)</b> - measure of difference between</li> </ul>	<ul style="list-style-type: none"> <li>• To implement person re-identification in color-</li> </ul>	<p>Training dataset of N people</p> <p>a. Height dynamics extraction –</p> <p>i. Multiple people tracking framework (position and appearance-based) to generate trajectories of detected people, using</p>	<p>1. TUM-GAID Dataset –</p> <p>a. Public dataset, containing multiple trials of subjects, captured with fronto-parallel</p>

		<b>colour depth camera</b>	<p>two probability distributions, given by <math>D_{KL}</math></p> <ul style="list-style-type: none"> <li><b>DGEI</b> - Depth gradient GEI</li> </ul> <p>• To formulate a probabilistic matching framework incorporating the selected frequency bins, using the discriminative ability of feature vector obtained</p>	<p>depth camera using height temporal information of people, by using a feature-selection scheme for each person based on KL-divergence to identify a discriminative subset of frequency bins for height-based gait feature.</p> <p>• To formulate a probabilistic matching framework incorporating the selected frequency bins, using the discriminative ability of feature vector obtained</p>	<p>3D color depth cameras.</p> <p>ii. Each 3D bounded box contains only foreground 3D color-depth blobs.</p> <p>iii. Top-down camera view:</p> <ol style="list-style-type: none"> <li>1. Since origin is defined in ground plane, person's height estimated from depth value corresponding to highest point in detected blob.</li> <li>2. To increase robustness, mean of <math>t</math> highest points calculated.</li> </ol> <p>iv. Front-parallel camera view:</p> <ol style="list-style-type: none"> <li>1. Height derived from foreground/depth silhouette height.</li> <li>2. Height/scale variations are adjusted using foreground depth info.</li> </ol> <p>Given apparent person height <math>h</math>, average foreground depth <math>d</math> and pre-defined normalizing depth value <math>s</math>, normalized person height <math>h'</math> given as <math>h' = (h/s)*d</math>.</p> <p>b. Feature vector from frequency response (Frequency domain transformation of height dynamics) –</p> <p>i. k-point DFT used to generate frequency response and obtain k-dimensional feature vector</p> $y_m^n = \left  \text{dft} \left( h_m^n, k \right) \right  v$ <p>Where <math>y_m^n</math> represent absolute real part of k-Fourier coefficients (frequency bins), <math>h_m^n</math> is estimated height dynamics for person <math>n</math> and trajectory <math>m</math>.</p> <p>c. Feature Selection – Technique of identifying subset of feature vectors to enhance overall</p>	<p>color depth camera.</p> <p>b. Trials partitioned into two clusters based on clothing; clothing-1, clothing-2.</p> <p>c. During training, multiple trials of clothing-1 for each subject.</p> <p>d. Test dataset partitioned into two subsets, based on clothing for each subject;</p> <ol style="list-style-type: none"> <li>First test partition – trials of clothing-1.</li> <li>Second partition – trials of clothing-1 and clothing-2.</li> </ol> <p>2. Dataset-1 and Dataset-2 –</p> <ol style="list-style-type: none"> <li>Multiple walking subjects using top-down color depth cameras</li> <li>Dataset-1: 2 sequences @15Hz using 5 ceiling mounted color depth cameras with 5 and 10 subjects, walking in indoor environment with varying illumination.</li> <li>Dataset-2: 3 subjects@15Hz with uniform illumination.</li> </ol> <p><math>k = 256</math>-DFT 10 features selected.</p> <p>Results:</p> <ul style="list-style-type: none"> <li>For first test partition, baseline</li> </ul>
--	--	----------------------------	---	--	---	--

			<p>by integration of height temporal information with a color &amp; height-based appearance model.</p>	<p>identification accuracy We identify subset of frequency bins unique to each person</p> <p>i. For each person <math>n</math>, compute empirical Gaussian mean and standard deviation over set of feature vectors <math>M</math> i.e. <math>\{y_m^n\}_{m=1}^{M_n}</math> and obtain <math>k</math> unimodal Gaussian distributions <math>\{\tilde{N}_n(\mu^k, \sigma^k)\}_{k=1}^K</math> corresponding to <math>k</math> frequency bins.</p> <p>ii. KL-divergence for <math>n</math> and frequency bin <math>k</math>:</p> $v^n(k) = \sum_{\substack{p=1 \\ p \neq n}}^N D_{kl}(\tilde{N}_n(k) \mid \tilde{N}_p(k))$ <p>iii. Large KL divergence values - minimum overlap Small KL divergence values - high degree of overlap</p> <p>iv. KL divergence value tends to zero when two gaussian distributions overlap completely.</p> <p>v. Selected frequency bin indices</p> $\Lambda^n = [\lambda^n(j)]_{j=1}^J$ <p>vi. To avoid overfitting, test trajectories not present in training dataset during feature selection</p> <p>d. Probabilistic Identification</p> <ul style="list-style-type: none"> <li>- Maximum likelihood classification used to identify test person.</li> <li>i. Extract test person's height dynamics and obtain frequency response.</li> <li>ii. Identify label in training dataset using likelihood function.</li> </ul> <p>Frequency response</p> $q = [q(k)]_{k=1}^K$ <p>Label I</p>	<p>algorithm (RGB-HT) performs marginally better than gait-based algorithm.</p> <ul style="list-style-type: none"> <li>• For second test partition, gait-based algorithm performs better than baseline algorithm.</li> <li>• Robust algorithm is limited to batch-mode applications due to supervised feature selection.</li> </ul>
--	--	--	--	--	---

					$l = \arg \max_n \prod_{j=1}^J p(q(\lambda^n(j)); \mu^n(\lambda^n))$ <p>Where <math>q(\lambda^n(j))</math> represents frequency bin <math>j</math> in <math>k</math>-dim test feature vector, <math>\mu^n(\lambda^n(j)) \wedge \sigma^n(\lambda^n(j))</math> are empirical mean and standard deviation.</p> <p>e. Baseline Identification Algorithm –</p> <ul style="list-style-type: none"> <li>i. RGB-height histogram: 4D histogram with discrete bins corresponding to person height and red, green and blue channels.</li> <li>ii. 3D RGB histogram created at each non-overlapping person height bin.</li> <li>iii. Nearest-neighbour classification using min. difference of pairwise assignment distance to mean feature vector (<math>\hat{G}_n</math>) for person <math>n</math> in the training dataset.</li> </ul> $l = \arg \max_n (-d(q, \hat{G}_n))$ <p>f. Combined Feature vector-based algorithm –</p> <ul style="list-style-type: none"> <li>i. Combination of color, person-height and gait as feature.</li> <li>ii. RGB-height histogram and height-dynamics gait features combined i.e. respective maximization-based matching measures are combined.</li> </ul>	
[7]	June 2016	<b>Gait-Based Person Identification Considering Clothing Variation</b>	<ul style="list-style-type: none"> <li>• <b>GEI</b> – Gait Energy Image</li> <li>• <b>GENI</b> – Gait Entropy Image</li> <li>• <b>DFT</b> – Discrete</li> </ul>	<ul style="list-style-type: none"> <li>• To implement the proposed method of person-reidentification</li> </ul>	<ol style="list-style-type: none"> <li>1. Silhouettes extracted from each frame of video sequence using background subtraction</li> <li>2. Silhouettes normalized into 120x88 pixels and registered to form spatio-temporal GSV.</li> <li>3. Gait period detection over complete gait cycle.</li> </ol> <p><b>OU-ISIR treadmill dataset-B</b> was only used for training/testing.</p> <p><b>1. Performance Evaluation of GEI, GENI, DFT and EnDFT Feature</b></p>	

	<b>n</b>	<p>Fourier Transform</p> <ul style="list-style-type: none"> <li>• <b>EnDFT</b> - Frequency-domain Gait Entropy</li> <li>• <b>Sequence combination</b></li> <li>• <b>Dynamic part selection</b></li> <li>• <b>GSV</b> - Gait Silhouette Volume</li> </ul> <p>using a combination of gait sequence and dynamically part weightin g using threshold value and weighted integratio n.</p> <ul style="list-style-type: none"> <li>• To test the proposed method on OU-ISIR Treadmill dataset B, which includes 68 subjects with 32 combinat ions of clothing types.</li> <li>• To calculate GEI, GEnI, DFT and EnDFT over a complete gait cycle and utilize dynamic weightin g on the gait features.</li> </ul>	<p><b>4. Gait Representation -</b></p> <p>a. <b>GEI:</b></p> $GEI(x,y) = \frac{1}{N} \sum_{n=1}^N B(x,y,n)$ <p>b. <b>GEnI:</b></p> <p>Avg. self-information obtained from k number of outputs generated from source that generates k symbols:</p> $-k \sum_{j=1}^k p(a_j) \log p(a_j)$ <p>Where <math>a_j</math> are symbols, <math>p(a_j)</math> are source symbols probability.</p> <p>c. <b>DFT:</b></p> <p><math>G(x,y,k)</math> [DFT] and amplitude <math>A(x,y,k)</math> for temporal axis:</p> $G(x,y,k) = \sum_{n=0}^{N-1} B(x,y,n) e^{-j w_0 k}$ $A(x,y,k) = \frac{1}{N} \sqrt{G(x,y,k)} \sqrt{G(x,y,-k)}$ <p>Where N is no. of frames in gait cycle, <math>w_0</math> is base angular frequency, k is frequency component. DFT has 3 components: (0-2 frequency component) (higher frequency components removed as noise)</p> <ul style="list-style-type: none"> <li>• 1<sup>st</sup> component equivalent to GEI</li> <li>• Middle component shows asymmetry of left and right motion</li> <li>• Last component represents symmetry.</li> </ul> <p>d. <b>EnDFT:</b></p> <p>EnDFT is computed from</p>	<p><b>sequence combination:</b></p> <p>a. CCR calculated for rank-1 - EnDFT has best performance at 70.13% recognition rate.</p> <p>2. <b>Comparison with existing methods:</b></p> <p>Proposed sequence combination +EnDFT has the highest recognition rate relative to existing methods such as <b>Baseline+GEI</b> (60.75%), <b>Part-based (without weighting)</b> (65.99%).</p>
--	----------	--	--	---

					<p>DFT image using the Gait Entropy equation, where <math>p_1(x,y)</math> is the intensity value of DFT.</p> <p>EnDFT gives more weights into dynamic areas; less (near to zero) for static areas; uses all three components of DFT</p> <p><b>5. Dynamic Part Weighting and Sequence Combination -</b></p> <ul style="list-style-type: none"> <li>a. Divide human body structure into small fragments; i.e each small fragment is single row</li> <li>b. Evaluation of each row from bottom of gait image, merging with next immediate upper row and finds recognition rate.</li> <li>c. Process continued till top most row is evaluated.</li> <li>d. Median value b/w local maxima and local minima calculated to divide the body image into multiple areas.</li> <li>e. Median value is calculated for next local maxima and so on to bottom of image.</li> <li>f. Estimate matching weight for each area based on similarity b/w extracted probe features and gallery features for standard clothes</li> <li>g. Probe identified by weighted integration of similarities in all areas</li> <li>h. Sequence combination used to calculate final weight.</li> <li>i. K-nearest neighbor classifier used to classify probe.</li> </ul>	
[8]	February 2019	<b>Gait-based person re-identification</b>	<ul style="list-style-type: none"> <li>• <b>LBP</b> – local binary pattern</li> <li>• <b>HOG</b> – Histogram of oriented</li> </ul>	<ul style="list-style-type: none"> <li>• To implement a new method for gait-based</li> </ul>	<p><b>Signature Extraction -</b></p> <ol style="list-style-type: none"> <li>1. Gait Cycle Estimation</li> <li>2. GEI generation</li> <li>3. Semantic classification: Attributes such as carried objects (backpack,</li> </ol>	CASIA-B database used. Results compared between LBP, HOG and P-LBP: 1. Descriptor Choice

	<b>under covariate factors</b>	gradients • Gait features	person re-identification relying on dynamic selection of human parts by computing a new person descriptor from relevant human parts, selected depending on presence of semantic information.	shoulder bag, and handbag) and clothes are semantic.  a. Offline stage: i. Data preparation of training database related to semantic attribute: Construct 2D table from database; each row represents bounding box's image; each column represents feature. Last column contains semantic attribute class. ii. Predictive model for each attribute class constructed using a. Support vector machines (SVM) b. Tree bagger-based decision tree (FT) c. Neural Network (NN) iii. Online Stage: 1. Extracting features to represent global information of each image. 2. Labels generated depending on semantic info. of image. 4. Feature extraction: Dynamic selection of parts from divided GEI image. (7 Parts) Set of relevant parts V: $V = \{c_1 \cdot G_{H1}(x, y), \dots, c_7 \cdot G_{H7}(x, y)\}$ $\begin{cases} c_i = 0, & \text{if semantic attribute } \exists \\ & c_i = 1, \text{ otherwise} \end{cases}$ <p>where i is index of part <math>G_{Hi}(x, y)</math>.</p> a. Extraction of salient and suitable features to capture gait characteristics using partial-LBP (local binary pattern). b. Takes pixels of image by thresholding 3x3 neighborhood of each pixel with center value, and mapping result as binary number (0/1).	<b>experiments:</b> 204 images used. a. P-LBP has advantages in terms of rank 1; 23.92% better than HOG and 47.05% better than LBP. 2. <i>Semantic Attributes classification:</i> a. 124 person images, 10 sequences per person. b. 744 carrying nothing, 248 with bags, 248 wearing coats. c. Classification done using SVM, NN and FT(tree bagger). d. Neural network gives best accuracy (>89%) for detecting semantic attributes. e. For normal semantic attributes (nothing), SVM shows better performance. 3. <i>Dynamic selection of human parts:</i> a. Comparison between P-LBP 7parts, invariant parts, and dynamic selection; P-LBP dynamic selection gives best results, with rank-1 performance of 52.352%, rank-3 of 69.8% and rank-5 of 78.23%.
--	--------------------------------	------------------------------	--	--	--

					$P - LBP(x_c, y_c) = \sum_{n=0}^7 s(i_n - i_c)$ <p>Where <math>i_c</math> corresponds to value of center pixel, <math>i_n</math> to value of 8 surrounding pixels and <math>s(x)</math> defined as</p> $s(x) = \begin{cases} 1, & \Lambda x \geq 0 \\ 0, & \Lambda x < 0 \end{cases}$ <p>c. P-LBP immediately records bit sequences to maintain more local information.</p> <p>d. Binary sequence is further used as gait signature for each relevant part.</p> <p>5. Signature Matching:</p> <ol style="list-style-type: none"> <li>Measures distance between probe person signature and each gallery image signature</li> <li>Common parts not affected by semantic attribute for gallery and probe image are used.</li> <li>Gallery set arranged according to similarity to generate ranked list.</li> <li>Top rank person's identity assigned to probe image.</li> <li>Euclidean distance used as metric.</li> </ol>	
[9]	May 2015	<b>Enhancing person re-identification by integrating gait biometric</b>	<ul style="list-style-type: none"> <li>• <b>PCA</b> – Principal component analysis</li> <li>• <b>Gabor Features</b></li> <li>• <b>Score level fusion</b></li> </ul>	<ul style="list-style-type: none"> <li>• To implement a method to enhance person-reidentification by integrating gait biometric</li> <li>• Using hierarchical feature extraction and descriptor</li> </ul>	<p><b>1. Hierarchical feature extraction:</b></p> <p>a. Gait feature: Original gait feature is 2048-D vector for each sequence; dimensionality reduction done using PCA to avoid overfitting.</p> <p>b. Appearance features:</p> <ol style="list-style-type: none"> <li>HSV space color histogram: Each dimension of HSV (transformed from RGB and equalized to reinforce illumination robustness) is divided into 128 bins to count number of pixels in corresponding bins</li> </ol>	<p>CASIA Gait Database B used, contains 11 views of 124 individuals with 3 conditions: Bag, clothes, normal. View angles: 0, 18, 36, 54, 72, 90, 108, 126, 144, 162 and 180 degrees.</p> <p>1. Comparison with baselines: Feature-level + MLR and score-level + MLR both perform much better compared to other</p>

			r matching with learned metric matrices.	<p>ii. Gabor features: 2D Gabor filter is Gaussian kernel function, modulated by complex sinusoidal plane wave. Used to generate texture feature in face and fingerprint recognition.</p> $\psi_{u,v}(z) = \frac{\ k_{u,v}\ ^2}{\sigma^2} e^{-\frac{\ k_{u,v}\ ^2 \ z\ ^2}{2\sigma^2}} [e^{i\phi}]$ <p>Where u is the orientation, v is scale of Gabor filters and z = (x,y)</p> $k_{u,v} = k_v e^{i\phi}$ <p>Gabor filter of 5 different scale and eight orientations, with:</p> $\sigma = \pi, k_{max} = \frac{\pi}{2}, f = \sqrt{2}$ <p>c. Metric learning to rank:</p> <ul style="list-style-type: none"> <li>i. General metric learning algorithm, based on structural SVM &amp; view metric learning.</li> <li>ii. Used to build view-angle independent mapping.</li> <li>iii. Data from different views of same person clustered in metric space, so that similarities between people measured as distances in metric space, regardless of view angles.</li> </ul> <p><b>2. Fusion:</b></p> <ul style="list-style-type: none"> <li>a. Score-level fusion: <ul style="list-style-type: none"> <li>i. View-independent global function, formulated as</li> </ul> <math display="block">D_{FIN}(S_1, S_2) = D_{APP}(S_1, S_2) + \theta D_{GEI}</math> <p>Where <math>\theta</math> is the weighting factor. Value of <math>\theta</math> range from 0.001 to 10; best result at <math>\theta=1</math></p> </li> <li>b. Feature-level fusion: <ul style="list-style-type: none"> <li>i. Fuse two type features before calculating distance.</li> </ul> <math display="block">\vec{F} = [\vec{A}, \vec{G}]</math> <p>Where <math>\vec{A} \wedge \vec{G}</math> is</p> </li> </ul>	<p>methods.</p> <p>2. Closed/Open set person-reidentification: Feature-level +MLR performs better than score-level+MLr as GEI has weaker discrimination than color and textures in short period Re-ID. However, GEI is appropriate supplement to appearance features.</p> <p>Integrating gait feature alleviates the negative effects of different appearance of people and make algorithm more robust.</p>
--	--	--	--	---	---

					appearance feature vector and GEI feature vector.	
[10]	Dece mber 2016	<b>Person Re-identification in frontal gait sequences via Histogram of Optic flow Energy Image</b>	<ul style="list-style-type: none"> <li>• <b>Optic flow</b></li> <li>• <b>Histogram of Flow Energy Image</b></li> <li>• <b>GEI</b></li> </ul>	<p>• Implement a novel method of re-identifying people in frontal video sequences, based spatio-temporal representation of gait based on optic flow features, called HOFEI.</p>	<p>1. Histogram of Flow:</p> <ol style="list-style-type: none"> <li>HOG encoding scheme on human detection bounding boxes, by using 'ground truth' annotation.</li> <li>Polar cells used to better represent spatial location of limbs and head w.r.t time.</li> <li>Optic flow image divided into cells using polar sampling strategy; compute histogram of flow orientation weighed by magnitude.  <math display="block">HOF^t = [HOF_1^t \dots HOF_n^t \dots HOF_m^t]</math> Where nR is no. of angular regions (cells) and nB is no. of bins defining each cell.  <math display="block">HOF^t</math> is of dimension 64 with nR=8 and nB=8.</li> <li><math>HOF^t</math> computed for each frame through video sequence S</li> </ol> <p>2. Gait Cycle Estimation:</p> <ol style="list-style-type: none"> <li>To estimate gait period, subset of histogram bins <math>HOF_2^t \wedge HOF_3^t</math> are used, which represent lower limb motions. Amplitude has good SNR ratio for detection.</li> <li>The periodic sinusoidal curve generated by plotting HOF peaks of single leg against frame (as function of time). Frames between two consecutive peaks represents gait cycle.</li> </ol> <p>3. Histogram of Flow Energy Image:</p> <ol style="list-style-type: none"> <li><math>HOF^t</math> representations averaged over full gait cycle to obtain HOF energy image  <math display="block">HOFEI = \frac{1}{t_2 - t_1} \sum_{t=t_1}^{t_2} HOF^t</math> <math>t_1 \wedge t_2</math> are the beginning</li> </ol>	<p>Experiments conducted in two scenarios. Re-ID in controlled scenario vs Re-ID in uncontrolled (busy office) scenario.</p> <p>For first scenario, CASIA Dataset B used:</p> <ol style="list-style-type: none"> <li>contains multiple videos of subjects including normal and apparel change conditions.</li> <li>105 subjects considered, 3 gait cycles extracted for each subject.</li> <li>Stefans implementation, using the Lucas-Kanade method was applied.</li> </ol> <p>For second scenario, HDA Person Dataset is used:</p> <ol style="list-style-type: none"> <li>13 indoor cameras, recording simultaneously for 30 min during busy noon hour inside university building.</li> <li>Only camera 19 was used, containing frontal gait sequences.</li> <li>Resolution of 640x480, frame rate of 5fps.</li> <li>12 people, min. of 3 gait cycles extracted per person for training and testing.</li> </ol> <p>Due to limitations of large video sequences and</p>

					and end frame indices of gait cycle	varying appearance conditions per person, Experiment 1 and 2 excluded in HDA Dataset.  Experiments: 1. Recognition in regular conditions: a. First four sequences used for training, last two used for probe set. b. CMC Curve shows CCR of 74.29%. c. Compared to silhouette based approaches (GEI, FDEI, GHI, GMI), CCR in between others. d. GEI and FDEI have better CCR due to use of segmented binary silhouettes and better classification method. 2. Re-identification under change in appearance: a. Apparel change CCR for bag is 66.67%. b. CCR for coat is 59.05%. c. Lower CCR of coat is due to global change in flow features. 3. Variable Distance to camera: a. CASIA: i. For 1 <sup>st</sup> case study, 6 $D_{far}$ probe and 12 training set ( $D_{middle}$ and $D_{near}$ ) per person. ii. For 2 <sup>nd</sup> case
--	--	--	--	--	-------------------------------------	--

						study, $D_{middle}$ is considered as probe and $D_{near}$ and $D_{far}$ are training sets. iii. For 3 <sup>rd</sup> case study, $D_{near}$ is probe set, $D_{middle}$ and $D_{far}$ are training sets. iv. $D_{middle}$ case outperforms other two cases with <b>33.81%</b> CCR; as trained on extreme ranges, classifier performs interpolation when predicting values for middle range. b. HDA Dataset: i. 24 training descriptors, 12 test probes. ii. $D_{middle}$ has highest CCR of <b>75%</b> , with <b>50%</b> and <b>58.33%</b> for $D_{far}$ and $D_{near}$ cases.
--	--	--	--	--	--	---

FIGURE 1.1: Literature Review of Current Gait-Based Person Identification Techniques

## Chapter 2

# Proposed Methodology

In this project, we proposed a motion-based approach to person identification, using *affine moment invariants* (AMIs) as feature descriptors. *Active Energy Image* (AEI) are generated from the subject silhouette sequence. The AEI is divided into multiple areas and features are extracted for each area. Matching weights are estimated for each area based on similarity between extracted features and those in the database sequences, using Euclidean distance as a metric. Finally, the subject is identified using the weighted sum of similarities of all areas. A nearest-neighbor classifier is implemented for identification of the subject. Experiments are performed on the CASIA-B Gait Database, using the normal walking sequences.

The approach can be divided into the following steps:

1. Extract active energy image (AEI) from the gait sequence and divide the image into multiple areas.
2. Extract affine moment invariants (AMI) at each area to use as gait features. The database consists of a set of affine moment invariants of multiple people who wear standard clothing with no accessories.
3. AEI of subject person is also divided in the same way as that of the dataset and gait features (AMIs) are extracted.

4. Matching weights are estimated at each area based on similarity between features of the subject and those in the database.
5. The subject is identified by the weighted combination of the similarities of all areas. Nearest Neighbor classifier is used to classify the test subject.

## 2.1 Active Energy Image

In gait recognition, two types of information - static and dynamic - are extracted from the gait sequence. According to [13], GEI is efficient in extracting static and dynamic information, both implicitly and explicitly. However, since GEI represents gait as a single image, there is a loss of dynamic information such as the fore-and-aft frame relations. Active Energy Image (AEI) feature representation can be used to solve the above problems.



FIGURE 2.1: Gait silhouette sequence of individual

Given a pre-processed binary gait silhouette[14]  $i = i_0, i_1, \dots, i_N - 1$ , where  $f_j$  represents the  $j^{\text{th}}$  silhouette (Figure 2.1), N is the total number of frames in the sequence, the difference image between frames is calculated as follows (Figure 2.2):

$$I_j = \begin{cases} i_j(p, q), & j = 0 \\ ||i_j(p, q) - i_{j-1}(p, q)||, & j > 0 \end{cases} \quad (2.1)$$



FIGURE 2.2: Difference Images between consecutive significant frames of gait sequence

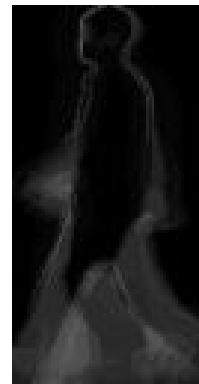


FIGURE 2.3: Active Energy Image

The AEI is defined as (Figure 2.3):

$$A(p, q) = \frac{1}{N} \sum_{j=0}^{N-1} I_j(p, q) \quad (2.2)$$

The active energy image representation is divided into 'K' segments as shown in Figure 2.4:

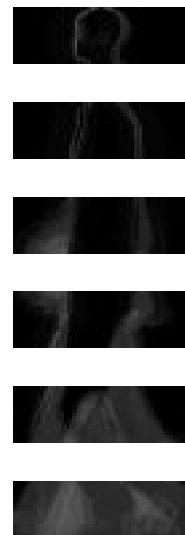


FIGURE 2.4: Segmented AEI (6 segments)

## 2.2 Affine Moment Invariants

### 2.2.1 Image moments and moment invariants

Image moments are certain particular weighted average (moment) of the image pixels' intensities, or a function of such moments, usually chosen to have some attractive property or interpretation. Some properties of images which are found via image moments include area (or total intensity), its centroid, and information about its orientation.

Image moments are well-known in their applications for image analysis, since they are used to derive **moment invariants**, which are invariants to transformations such as translation and scaling. The well-known ***Hu moment invariants*** were shown to be invariant to translation, scale and rotation. However, Flusser[29] showed that the traditional set of Hu moment invariants is neither independent nor complete.

Affine moment invariants, proposed by Flusser and Suk [17] are moment-based descriptors, which are invariant under a general affine transformation. The invariants are generally derived by means of the classical theory of algebraic invariants [18], graph theory [19,20] or the method of normalization [21]. The most common method is the use of graph theory.

The moments describe shape properties of an object as it appears. For an image, the centralized moment of order  $(p+q)$  of an object  $O$  is given by

$$\mu_{pq} = \sum_{(x,y) \in O} (x - x_g)^p (y - y_g)^q A(x, y) \quad (2.3)$$

Here,  $x_g$  and  $y_g$  define the center of the object. More specifically, they are calculated from the geometric moments  $m_{pq}$ , given by  $x_g = \frac{m_{10}}{m_{00}}$  and  $y_g = \frac{m_{01}}{m_{00}}$ .

The affine transformation can be expressed as

$$u = a_0 + a_1x + a_2y \quad (2.4)$$

$$v = b_0 + b_1x + b_2y \quad (2.5)$$

[20] In our method, the number of affine moment invariants ( $\mathbf{A} = (A_1, A_2, \dots, A_M)^T$ ) M is 10. We show 5 such invariants:

$$\begin{aligned}
A_1 &= \frac{1}{\mu_{00}^4} (\mu_{20}\mu_{02} - \mu_{11}^2) \\
A_2 &= \frac{1}{\mu_{00}^{10}} (\mu_{30}^2\mu_{03}^2 - 6\mu_{30}\mu_{21}\mu_{12}\mu_{03} + 4\mu_{30}\mu_{12}^3 + 4\mu_{03}\mu_{21}^3 \\
&\quad - 3\mu_{21}^2\mu_{12}^2) \\
A_3 &= \frac{1}{\mu_{00}^7} (\mu_{20}(\mu_{21}\mu_{03} - \mu_{12}^2) - \mu_{11}(\mu_{30}\mu_{03} - \mu_{21}\mu_{12}) \\
&\quad + \mu_{02}(\mu_{30}\mu_{12} - \mu_{21}^2)) \\
A_4 &= \frac{1}{\mu_{00}^{11}} (\mu_{20}^3\mu_{03}^2 - 6\mu_{20}^2\mu_{11}\mu_{12}\mu_{03} - 6\mu_{20}^2\mu_{02}\mu_{21}\mu_{03} \\
&\quad + 9\mu_{20}^2\mu_{02}\mu_{12}^2 + 12\mu_{20}\mu_{11}^2\mu_{21}\mu_{03} \\
&\quad + 6\mu_{20}\mu_{11}\mu_{02}\mu_{30}\mu_{03} - 18\mu_{20}\mu_{11}\mu_{02}\mu_{21}\mu_{12} \\
&\quad - 8\mu_{11}^3\mu_{30}\mu_{03} - 6\mu_{20}\mu_{02}^2\mu_{30}\mu_{12} + 9\mu_{20}\mu_{02}^2\mu_{21}^2 \\
&\quad + 12\mu_{11}^2\mu_{02}\mu_{30}\mu_{12} - 6\mu_{11}\mu_{02}^2\mu_{30}\mu_{21} + \mu_{02}^3\mu_{30}^2) \\
A_5 &= \frac{1}{\mu_{00}^6} (\mu_{40}\mu_{04} - 4\mu_{31}\mu_{13} + 3\mu_{22}^2)
\end{aligned}$$

### 2.3 Estimation of matching weights

[5] Matching weights are calculated by first whitening the AMI in the database and of the subject at each area. The distance  $d_{n,s}^k$  between the features of the subject and those of all datasets in the database is computed as follows:

$$d_{n,s}^k = ||{}^w A_{SUB}^k - {}^w A_{DB_{n,s}}^k|| \quad (2.6)$$

where  ${}^w A_{SUB}^k$  and  ${}^w A_{DB_{n,s}}^k$  are the whitened AMI of the subject and those of a person in the database.

The whitening of the affine moment invariants is done by:

1. Applying (principal component analysis (PCA)) to the calculated affine invariants and projecting them to a new a feature space

2. Normalizing the projected features based on the corresponding eigenvalue.

(n,s and k) are  $1 \leq n \leq N$  (N is number of persons in database),  $1 \leq s \leq S$  (S is number of sequences of each person), and  $1 \leq k \leq K$  (K is number of divided areas).  $d_{n,s}^k$  is calculated in the Euclidean norm.

To estimate matching weights, we use the similarity between subject features and database features; high matching weights are set to areas with less appearance changes, and low matching weights set to those with more appearance changes.

Steps to estimate matching weights:

1. At each area k, select sequences from the database where  $d_{n,s}^k < \bar{d}_{min}$ . These selected sequences are considered to have high similarity with the subject.  $\bar{d}_{min}$  is defined as:

$$\bar{d}_{min} = \min_n \bar{d}_n^k \quad (2.7)$$

$$\bar{d}_n^k = \frac{1}{S} \sum_{s=1}^S d_{n,s}^k \quad (2.8)$$

We consider that in each area, if at least one sequence of a person in the database is selected, then the matching scores of all sequences of that person are also high.

2. The non-selected sequences are all considered to be low similarities, so the distances of these sequences are redefined as  $d_{max}$  ( $d_{max} = \max_{n,s,k} d_{n,s}^k$ )
3. The above procedure is applied for all areas and the sum of distances for all areas is calculated by  $D_{n,s} = \sum_{k=1}^K d_{n,s}^k$ . Subject is identified by nearest-neighbor method.

### 2.3.1 Principal Component Analysis

In machine learning classification problems, there are often too many features on the basis of which the final classification is done. The higher the number of features, the harder it gets to visualize the training set and then work on it. Sometimes, most of these features are correlated, and hence redundant. This leads to a need for dimensionality reduction, which reduces the complexity of a model and avoids overfitting.

*Principal Component Analysis (PCA)* is one of the most popular algorithms used for dimensionality reduction. Proposed by Karl Pearson, PCA is an unsupervised linear transformation technique used in identifying patterns in data based on the correlation in features. It aims to find the directions of maximum variance in high-dimensional data and projects it onto a new subspace with equal or fewer dimensions than the original one.

The algorithm can be summarized in the following steps:

1. Standardize the d-dimensional dataset: This is usually done by standardizing the data such that the d-dimensional dataset will have a mean of 0 and a standard deviation of 1.
2. Construct the covariance matrix.
3. Decompose the covariance matrix into its eigenvectors and eigenvalues.
4. Sort the eigenvalues by decreasing order to rank the corresponding eigenvectors.
5. Select k eigenvectors which correspond to the k largest eigenvalues, where k is the dimensionality of the new feature subspace ( $k \leq d$ ).
6. Construct a projection matrix W from the “top” k eigenvectors.
7. Transform the d-dimensional input dataset X using the projection matrix W to obtain the new k-dimensional feature subspace.

## Chapter 3

# Experiment and Results

### 3.1 Implementation

The implementation was done in Python 3.6.8, on a Windows 10 laptop, with an Intel i7 8th Generation processor and 8GB of RAM. The following modules with the mentioned versions were used:

- Numpy - 1.17.2
- Scikit-learn - 0.21.3
- SciPy - 1.3.1
- OpenCV - 4.1.1.26
- Pillow - 6.2.0
- Scikit-image - 0.15.0
- ImageIO - 2.6.1
- PyWavelets - 1.0.3

Two custom modules were also written to provide the following functionality for the main program:

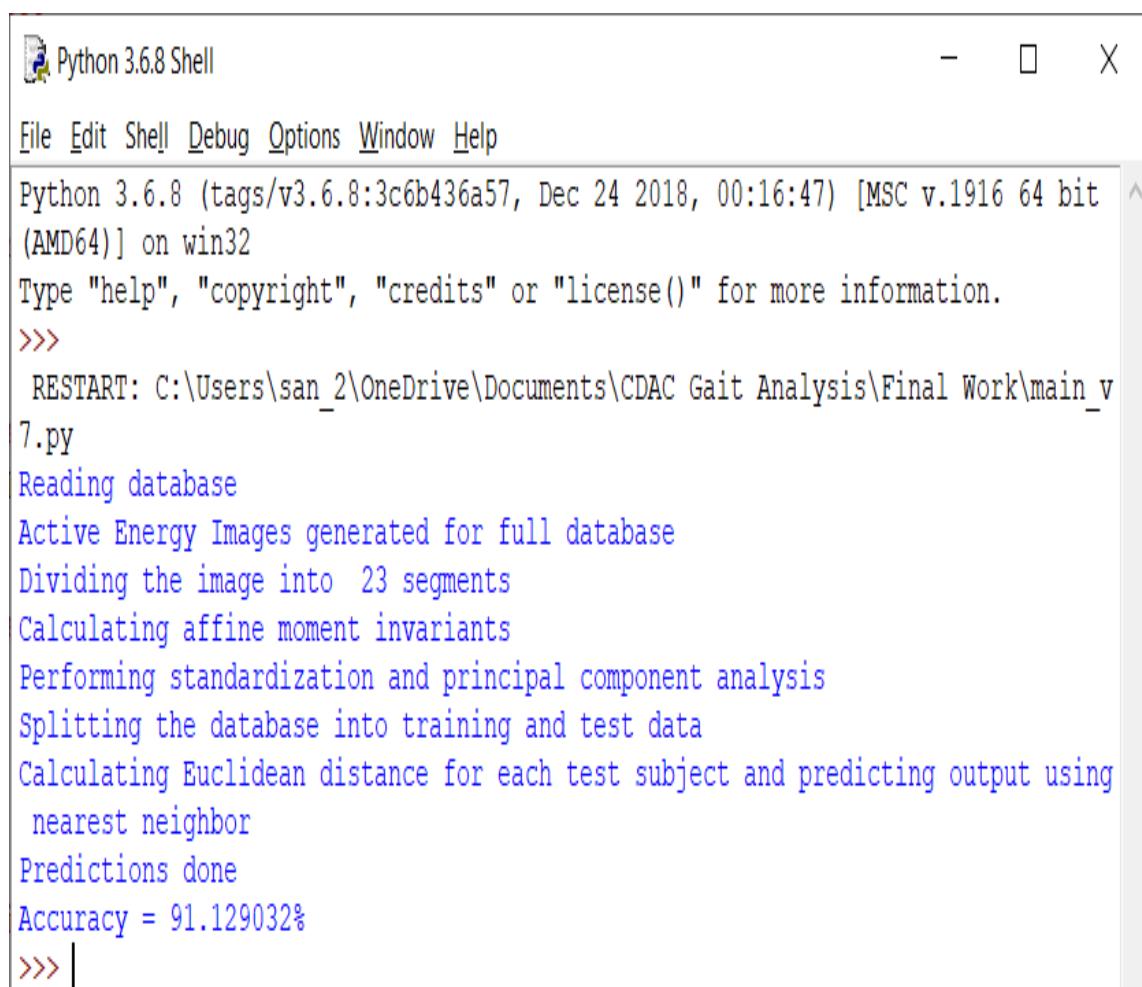
1. *functions\_module.py*: Contains functions to generate the active energy image (AEI), divide the image into areas and compute the accuracy.(Nearest Neighbor Implementation)
2. *ami.py*: Contains functions to extract affine moment invariants from the AEI images.

## 3.2 Experiment and Result

The proposed method was applied to the CASIA-B Gait Dataset. The CASIA-B Dataset is a large multi-view gait database, created in January 2005. There are 124 subjects, and the gait data was captured from 11 views. Three variations, namely view angle, clothing and carrying condition changes, are separately considered.

In the experiment, I used the lateral view (90 degree) standard walking sequences to analyze the prediction rate of the algorithm. The CASIA-B Dataset consists of six standard walking sequences for each person. I calculated CCRs by dividing the six sequences of each subject into two sets;  $124 \times 5$  sequences were used for training and the rest were used for testing.

The total number of affine moment invariants  $M$  were varied from 1 to 10 and the parameter  $K$  (number of divided areas) from 5 to 30. In case of  $K = 23$  and  $M = 5$ , the proposed method shows the highest accuracy of **91.13%** (Figure 3.1).



The screenshot shows a Python 3.6.8 Shell window. The title bar reads "Python 3.6.8 Shell". The menu bar includes File, Edit, Shell, Debug, Options, Window, and Help. The main window displays the following text:

```
Python 3.6.8 (tags/v3.6.8:3c6b436a57, Dec 24 2018, 00:16:47) [MSC v.1916 64 bit  
(AMD64)] on win32  
Type "help", "copyright", "credits" or "license()" for more information.  
>>>  
RESTART: C:\Users\san_2\OneDrive\Documents\CDAC Gait Analysis\Final Work\main_v  
7.py  
Reading database  
Active Energy Images generated for full database  
Dividing the image into 23 segments  
Calculating affine moment invariants  
Performing standardization and principal component analysis  
Splitting the database into training and test data  
Calculating Euclidean distance for each test subject and predicting output using  
nearest neighbor  
Predictions done  
Accuracy = 91.129032%  
>>> |
```

FIGURE 3.1: Best Case Results (K=23, M=5)

# **Chapter 4**

## **Conclusions and Future Work**

### **4.1 Conclusion**

A person re-identification method was proposed that utilizes active energy images, focusing on the dynamic information in the gait sequence, and is robust to appearance changes. In this method, the active energy image of the gait sequence was calculated and the human body area was divided into multiple areas. Affine moment invariants were extracted at each area as gait features and a matching weight was estimated for each area based on the similarity between features of the subject and those in the database. Then, the subject was identified by weighted sum of similarities in all areas using the nearest neighbor classifier.

In this research, we focused on developing a method that utilizes gait features using the dynamic information embedded in the walking speed and style. There are other potential factors that may affect the performance of the algorithm, such as different viewing angle, walking speed, etc. The specific immediate improvement to be done is to increase the current recognition rate without compromising on the dynamic features.

## 4.2 Future Work

Future work on the algorithm will address the following:

1. Improvement of the recognition rate by
  - Extracting features specific to the human body segments such as head, torso, arms and legs, and assigning matching weights based on these body segments.
  - Increasing the number of feature descriptors (AMIs) to get a better measure of variance in dynamic information of the image.
2. Testing and fine-tuning the algorithm on other viewing angles.
3. Testing the algorithm on other gait sequences that include additional clothing and accessories on the human body area.
4. Implementing a neural network to improve estimation of matching weights for similarity of subject and database features.
5. Testing of the algorithm on real-time wide area surveillance to identify and track subjects over multiple non-overlapping cameras.

# Bibliography

- [1] M. Babaee, L. Li, and G. Rigoll, “Person identification from partial gait cycle using fully convolutional neural network,” *Neurocomputing*, 04 2018.
- [2] S. Li, M. Zhang, W. Liu, H. Ma, and Z. Meng, “Appearance and gait-based progressive person re-identification for surveillance systems,” 09 2018, pp. 1–7.
- [3] A. Gala and S. Shah, “Gait-assisted person re-identification in wide area surveillance,” 04 2015, pp. 633–649.
- [4] R. Kawai, Y. Makihara, C. Hua, H. Iwama, and Y. Yagi, “Person re-identification using view-dependent score-level fusion of gait and color features,” 01 2012, pp. 2694–2697.
- [5] Y. Iwashita, K. Uchino, and R. Kurazume, “Gait-based person identification robust to changes in appearance,” *Sensors (Basel, Switzerland)*, vol. 13, pp. 7884–901, 06 2013.
- [6] V. John, G. Englebienne, and B. Krose, “Person re-identification using height-based gait in colour depth camera,” 09 2013.
- [7] I. Ahmed and M. Rokanujjaman, *International Journal of Innovations & Advancement in Computer Science (IJIACS)*, vol. 5, pp. 28–42, 06 2016.
- [8] E. Fendri, I. Chtourou, and M. Hammami, “Gait-based person re-identification under covariate factors,” *Pattern Analysis and Applications*, 02 2019.
- [9] Z. Liu, Z. Zhang, and Q. Wu, “Enhancing person re-identification by integrating gait biometric,” *Neurocomputing*, vol. 168, 05 2015.

- [10] A. Nambiar, J. Nascimento, A. Bernardino, and J. Santos-Victor, “Person re-identification in frontal gait sequences via histogram of optic flow energy image,” vol. 10016, 10 2016, pp. 250–262.
- [11] S. Sarkar, P. J. Phillips, Z. Liu, I. Robledo Vega, P. Grother, and K. Bowyer, “The humanoid gait challenge problem: Data sets, performance, and analysis,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 27, pp. 162–77, 03 2005.
- [12] A. Bobick and A. Johnson, “Gait recognition using static, activity-specific parameters,” vol. 1, 02 2001, pp. I–423.
- [13] J. Han and B. Bhanu, “individual recognition using gait energy image”,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, pp. 316–22, 03 2006.
- [14] E. Zhang, Y. Zhao, and W. Xiong, “Active energy image plus 2dlpp for gait recognition,” *Signal Processing*, vol. 90, pp. 2295–2302, 07 2010.
- [15] K. Bashir, T. Xiang, and S. Gong, “Gait recognition using gait entropy image,” vol. 2009, 01 2010, pp. 1 – 6.
- [16] M. Rokanujjaman, M. Islam, M. A. Hossain, M. Islam, Y. Makihara, and Y. Yagi, “Effective part-based gait identification using frequency-domain gait entropy features,” *Multimedia Tools and Applications*, vol. 74, 11 2013.
- [17] J. Flusser and T. Suk, “Pattern recognition by affine moment invariants,” *Pattern Recognition*, vol. 26, pp. 167–174, 1993.
- [18] S. Marxsen, D. Hilbert, R. Laubenbacher, B. Sturmfels, H. David, and R. Laubenbacher, *Theory of Algebraic Invariants*, ser. Cambridge Mathematical Library. Cambridge University Press, 1993. [Online]. Available: <https://books.google.co.in/books?id=xCneO62aHiIC>
- [19] T. Suk and J. Flusser, “Graph method for generating affine moment invariants,” vol. 2, 09 2004, pp. 192 – 195 Vol.2.
- [20] T. Suk, “Tables of affine moment invariants generated by the graph method,” 01 2005.

- [21] ——, “Affine normalization of symmetric objects,” 10 2005, pp. 100–107.
- [22] G. Pass, R. Zabih, and J. Miller, “Comparing images using color coherence vectors,” in *ACM Multimedia*, 1996.
- [23] J. Flusser, B. Zitova, and T. Suk, *Moments and Moment Invariants in Pattern Recognition*. Wiley Publishing, 2009.
- [24] A. Abdel-Hakim and A. Farag, “Csift: A sift descriptor with color invariant characteristics,” vol. 2, 02 2006, pp. 1978 – 1983.
- [25] B. Pathak and D. Barooah, “Texture analysis based on the gray-level co-occurrence matrix considering possible orientations,” 2013.
- [26] E. Miyamoto and T. E. Merryman, “Fast calculation of haralick texture features,” 2005.
- [27] S. Zhang, Y. Huang, Y. Yu, H. Li, and D. Metaxas, “Automatic image annotation using group sparsity,” 06 2010, pp. 3312–3319.
- [28] B. Manjunath, J. Ohm, V. Vasudevan, and A. Yamada, “Color and texture descriptors,” *Circuits and Systems for Video Technology, IEEE Transactions on*, vol. 11, pp. 703 – 715, 07 2001.
- [29] J. Flusser, “On the independence of rotation moment invariants,” *Pattern Recognition*, vol. 33, pp. 1405–1410, 2000.