

# Study of Person Re-Identification using Adaptive Spatial Temporal Attention Aware Learning with Gaussian Filter

Richa Jain<sup>1</sup>, Prof. Seema Shukla<sup>2</sup>

<sup>1</sup>PG Scholar, <sup>2</sup>Assistant Professor

<sup>1,2</sup>Dept. of ECE, MITS, Bhopal

*Abstract- Individual re-ID is a difficult assignment mostly because of variables, for example, foundation mess, posture, brightening and camera perspective varieties. These components obstruct the way toward extricating hearty and discriminative portrayals, consequently keeping various personalities from being effectively recognized. To improve the portrayal learning, generally neighborhood highlights from human body parts are separated. Be that as it may, the normal practice for such a procedure has been founded on bouncing box part identification. Right now, propose to embrace Adaptive Spatial-Temporal Attention - Aware Learning and Gaussian Filter which, because of its pixel-level exactness and capacity of demonstrating subjective shapes, is normally a superior other option. Our proposed (ASTAL-GF) incorporates learning face to face re-distinguishing proof and impressively outflanks its counter pattern, however accomplishes cutting edge execution. Our proposed techniques improve best in class individual re-recognizable proof on exactness of worldwide and nearby highlights according to combination rate is improve by 1.02%, thus acknowledgment rate may improve. The exactness of nearby and worldwide implanting size is increment by 0.736% and 1.14% separately. The precision relying upon the edges for neighborhood separation and worldwide separation is increment by 0.735% and 1% individually. The exhibition of re-recognizable proof procedure is improved by 2.73%.*

**Index words - Person Re-Identification, Recognition Rate, Adaptive Spatial-Temporal Attention - Aware Learning, Gaussian Filter, Accuracy, Local Distance and Global Distance.**

## I. INTRODUCTION

Individual re-recognizable proof is considered as a huge piece of numerous cameras person on foot following. It is characterized as the way toward recognizing whether to walker pictures of disjoint and non-covering cameras at various time interims is the equivalent or not. A run of the mill individual re-distinguishing proof framework has three stages, to be specific, identifying an individual, following the individual and recovering the individual. This work centers around an audit of individual recovery utilizing generative ill-disposed systems. A few analysts have given exceptional consideration to the issue of individual re-distinguishing proof in the field of PC vision. Disregarding much research right now, stays a difficult issue for the specialists in the field. It has been seen that a large portion of the individual pictures are clicked by disjoint and non-covering cameras introduced in an uncontrolled situation, having a low nature of the pictures. The low-quality pictures it hard to use by the regular frameworks for

separating their highlights to be utilized for distinguishing people's countenances precisely. Along these lines, the highlights speaking to the presence of the individual must be removed based on garments hues or some article with the individual to distinguish them in various stances of numerous cameras. Subsequently, these appearance highlights are progressively reasonable for an individual re-recognizable proof procedure. Be that as it may, the appearance highlights are powerless if there should arise an occurrence of comparable fabric hues. What's more, the presence of people likewise changes in numerous camera sees because of lightning change in various foundation presents. It might likewise be conceivable that various people have comparable appearances or similarity.

## II. BACKGROUND

**Guangyi Chen et. al,** Right now, present a spatial-worldly consideration mindful learning (STAL) strategy for video-based individual re-distinguishing

proof. Most existing individual re-recognizable proof techniques total picture includes indistinguishably from speak to people, which are extricated from the equivalent responsive field across video outlines. In any case, the picture quality might be fluctuating for various spatial areas and changing after some time, which will add to individual portrayal and coordinating adaptively. Our STAL strategy expects to take care of the striking pieces of people in recordings mutually in both spatial and transient spaces. To accomplish this, we cut the video into various spatial-worldly units which safeguard the body structure of an individual and build up a joint spatial fleeting consideration model to get familiar with the quality scores of these units. We assess the proposed technique on three testing datasets including iLIDS-VID, PRID-2011 and the huge scope MARS dataset, and reliably improve the rank-1 precision by an enormous edge of 5.7%, 0.9%, and 6.6% separately, in correlation with the cutting edge strategies.[1]

**O. S. Reshma et. al,** In examination, individual re-ID is additionally an extreme undertaking of coordinating people decided from absolutely totally unique camera sees. It is important applications in AI, risk discovery, human trailing and action examination. Individual re-distinguishing proof is likewise an extreme examination subject because of fractional impediments, low goals pictures and gigantic enlightenment changes. Additionally, individual decided from absolutely totally extraordinary camera sees has significant minor departure from postures and perspectives. This paper outlines the moves identified with the individual re-distinguishing proof together talk about changed methods used individual re-ID.[2]

**Xiaoke Zhu et. al,** Video-based individual re-recognizable proof (re-id) is a significant application by and by. Since enormous varieties exist between various person on foot recordings, just as inside every video, it's trying to lead re-recognizable proof between walker recordings. Right now, propose a concurrent intra-video and between video separations learning (SI2DL) approach for video-based individual re-id. In particular, SI2DL at the same time learns an intra video separation metric and a between video separation metric from the preparation

recordings. The intra-video separation metric is utilized to make every video increasingly minimized, and the between video one is utilized to guarantee that the separation between really coordinating recordings is littler than that between wrong coordinating recordings. Taking into account that the objective of separation learning is to make really coordinating video sets from various people be very much isolated with one another, we additionally propose a couple partition based SI2DL (P-SI2DL). P-SI2DL expects to gain proficiency with a couple of separation measurements, under which any two genuinely coordinating video sets can be very much isolated. Trials on four open person on foot picture succession datasets show that our methodologies accomplish the best in class execution.[3]

**Yueqi Duan et. al,** Right now, propose a setting mindful neighborhood double component learning (CA-LBFL) strategy for face acknowledgment. Not at all like existing learning-based neighborhood face descriptors, for example, separate face descriptor (DFD) and reduced paired face descriptor (CBFD) which gain proficiency with each component code independently, our CA-LBFL misuses the relevant data of contiguous bits by obliging the quantity of movements from various double bits, so increasingly powerful data can be abused for face portrayal. Given a face picture, we first concentrate pixel contrast vectors (PDV) in neighborhood fixes, and become familiar with a discriminative mapping in an unaided way to extend every pixel distinction vector (PDV) into a setting mindful parallel vector. At that point, we perform bunching on the educated twofold codes to develop a codebook, and concentrate a histogram include for each face picture with the scholarly codebook as the last portrayal. So as to misuse nearby data from various scales, we propose a setting mindful neighborhood twofold multi-scale highlight learning (CA-LBMFL) technique to together get familiar with different projection grids for face portrayal. To make the proposed techniques relevant for heterogeneous face acknowledgment, we present a coupled CA-LBFL (C-CA-LBFL) strategy and a coupled CA-LBMFL (C-CA-LBMFL) strategy to lessen the methodology hole of relating heterogeneous faces in the component level, separately. Broad trial results on four generally

utilized face datasets unmistakably show that our strategies beat most best in class face descriptors.[4]

**Weihua Chen et. al**, Individual re-recognizable proof (ReID) is a significant undertaking in wide territory video reconnaissance which centers around distinguishing individuals across various cameras. As of late, profound learning systems with a triplet misfortune become a typical structure for individual ReID. Be that as it may, the triplet misfortune pays primary considerations on getting right requests on the preparation set. It despite everything experiences a more fragile speculation ability from the preparation set to the testing set, in this manner bringing about mediocre execution. Right now, plan a quadruplet misfortune, which can prompt the model yield with a bigger between class variety and a littler intra-class variety contrasted with the triplet misfortune. Therefore, our model has a superior speculation capacity and can accomplish a better on the testing set. Specifically, a quadruplet profound system utilizing an edge based online hard pessimistic mining is proposed dependent on the quadruplet misfortune for the individual ReID. In broad examinations, the proposed arrange beats the vast majority of the cutting edge calculations on agent datasets which unmistakably exhibits the adequacy of our proposed technique. [5]

### III. PROBLEM IDENTIFICATION

The basic objections of my hypothesis work are according to the accompanying:

1. The exactness of worldwide and nearby highlights are restricted according to combination rate, subsequently acknowledgment rate may diminish.
2. The implanting size of highlight portrayal model is additionally critical to the ReID issue. According to increment of worldwide and nearby implanting size, precision may diminish.
3. Re-distinguishing proof precision differs relying upon the edges of worldwide separation and nearby separation. In low light pictures and recordings create low precision for individual re-ID.
4. Execution of re-distinguishing proof procedure is restricted when the lengths of both test and display arrangements are fluctuated.

### IV. PROPOSED METHODOLOGY

The basic procedure of propose methodology Adaptive Spatial-Temporal Attention - Aware Learning with Gaussian Filter (ASTAL-GF) can be explain through following point.

(1) Given a pedestrian video  $X = \{x^t\}_{t=1:T}$ , where

$T$  is the number of video frames and  $x^t$  denotes the  $t^{\text{th}}$  frame. The proposed network is divided into three branches: a global representation branch, local representation branch and gaussian filter representation.

They are fused in an end-to-end framework to learn discriminative person representation in different granularities.

(2) The global representation branch is designed to learn the full body representation of pedestrians. In this branch, at the beginning of the process, image sequence  $X$  is fed into a low level CNN to generate the low-level representations, after that, we apply a residual attention network (RAN) in a high-level CNN to extract global features  $g^t = G(x^t)$ .

(3) The local representation branch is used to address the local variance of person video, e.g., local mismatching due to pose variance. In this branch, we first apply a human pose estimation algorithm to locate the body joints and generate the body part coordinates  $p^{r,t} = \{p_1, p_2, p_3, p_4\}^{r,t}$

based on the estimated joints. The local part generator is pre-trained with the MPII human pose dataset. Then we apply an ROI pooling layer with body part coordinates on the feature maps by the low-level CNN and design a part specific network to generate the local part representation. The part-specific network has the same structure for different body parts but learns the different parameters. The whole processing is formulated as:

$$\{f^{r,t}\}_{r=1:R,t=1:T} = F(x^t, p^{r,t})$$

where  $\{f^{r,t}\}$  represents the local feature of  $r^{\text{th}}$  spatial body part in  $X^t$ .

In addition, we also develop a separate attention branch to learn a joint spatial-temporal attention score map  $a^{r,t} = A(X)$ , which is used to evaluate the qualities of different spatial-temporal units. The

temporal attention focuses on the key frames with rich discriminative information, while the spatial attention identifies the body parts which are not corrupted by occlusions and cluttered background. Finally, we define an aggregation function of local features  $f^{r,t}$  with the attention scores  $a$ ;  $t$  to calculate the distance between two pedestrian video clips, which is formulated as:

$$d_l(i, j) = \psi(f_i^{r,t}, f_j^{r,t}; a_i^{r,t}, a_j^{r,t})$$

where  $i; j$  denote two person videos captured in different cameras. We aggregate the global features with different frames in a temporal pooling layer and calculate the global distance as  $d_g(i, j) = \|g_i - g_j\|_2$ , where  $g_i$  denotes the global feature of  $i$ th person. In the training procedure, we optimize the objective function with both global and local distance, while in the testing procedure, we add them for the final similarity measure.

The objective function of our method is formulated as follows:

$$\min_{G, F, A} = L_{tri}(G, F, A) + L_{cls}(G) + C_{cons}(F)$$

which contains three parts: triplet loss, softmax loss, and consistency constraint.

(a) Triplet Loss: We design the triplet loss to preserve the rank relationship among a triplet of pedestrian videos. In the triplet loss, the distances between feature pairs from the same class are minimized while the ones from different classes are maximized. We calculate the triplet loss with both global and local features as follows:

$$L_{tri} = \sum_{i, j, k \in \Omega} [d_g(i, j) - d_g(i, k) + m_g] + \sum_{i, j, k \in \Omega} \lambda [d_l(i, j) - d_l(i, k) + m_l]$$

where  $m_g$  and  $m_l$  are margin thresholds to limit the gap between the distances from positive and negative samples, and  $[x]^+$  denotes the max function  $\max(0, x)$ .

(b) Softmax Loss: We apply the softmax loss function to learn the identity-specific global representation. Different from triplet loss, the softmax loss focuses on the robustness of person

video representations for identification. The softmax loss is defined as:

$$L_{cls} = \sum_{i \in \Omega} \frac{\exp(W_{y_i}^g g_i)}{\sum_k \exp(W_k^r g_i)}$$

where  $y_i$  is the identity of  $i$ th person and  $W_{y_i}^g, W_k^g$  indicate  $y_i$ th and  $k$ th columns of the softmax matrix.

(c) Consistency Constraint: To preserve the consistency between local and global features, we develop a consistency constraint in our framework, which requests that the identifies of single local feature and global feature are identical:

$$C_{consis} = \sum_{i \in \Omega} \sum_{r, t} \frac{\exp(W_{y_i}^r f_i^{r,t})}{\sum_k \exp(W_k^r f_i^{r,t})}$$

Note that the same parts in different frames share the same softmax matrix.

(4) Gaussian filter representation

Input  $M \times N$  image  $I$ , odd neighborhood size  $NH$ , noise variance  $NV$

$b = \text{floor}(NH/2)$  (half-width of kernel)

(filter interior of image)

for( $r = b \dots M - 2b + 1$ ) (1st  $b$  and last  $b$  rows not filtered)

for( $c = b \dots N - 2b + 1$ ) (1st  $b$  and last  $b$  cols not filtered)

(compute neighborhood statistics)

sum = 0

sumsq = 0 (sum of the squares)

for( $ih = -b \dots +b$ )

for( $jh = -b \dots +b$ )

sum +=  $I[r+ih][c+jh]$

sum +=  $(I[r+ih][c+jh])^2$

end

end (for each pixel in current neighborhood)

## VI. RESULTS AND ANALYSIS

The trial works start with MATLAB R2013a adaptation. Right off the bat, expect unique picture. For this reason `imread()` work has been utilized, at that point show the picture in figure window with figure order. For demonstrating picture in figure

window, we need to utilize imshow() work, at that point following yield has been created.

Table 1: Rank 1 accuracy on iLIDS-VID dataset with different fusion rates.

Fusion Rate	STAL[1]	ASTAL-GF (Proposed)
0	0.788	0.796
0.1	0.804	0.813
0.2	0.824	0.865
0.3	0.828	0.876
0.4	0.81	0.821
0.5	0.822	0.837
0.6	0.801	0.814
0.7	0.805	0.826
0.8	0.798	0.804
0.9	0.796	0.802
1.0	0.782	0.791

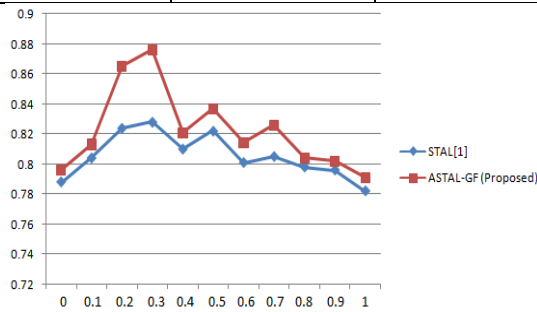


Figure 1: Rank 1 accuracy on iLIDS-VID dataset with different fusion rates.

In our model, we remove both worldwide highlights and neighbourhood body-part highlights to speak to individual video. When testing, we figure the separation frameworks of worldwide and neighbourhood includes autonomously and intertwine them with an equalization rate as the last comparability metric. The rate is set exactly to improve the exhibition. Subsequently, we examined a couple of various combination rates, which run from 0 to 1 with an interim of 0.1. The outcome exhibitions on the iLIDS-VID dataset with various combination rates are obviously, the combination methodology improves the presentation altogether, since the worldwide semantic portrayal and neighbourhood body-part descriptor are correlative. Our model accomplishes the ideal position 1 exactness when the combination rate is set to 0.3.

Table 2: Rank 1 accuracy on iLIDS-VID dataset with different embedding sizes.

Embedding Size	Local		Global	
	STAL[1]	ASTAL-GF (Proposed)	STAL[1]	ASTAL-GF (Proposed)
128	0.815	0.821	0.792	0.801
256	0.812	0.819	0.801	0.812
512	0.824	0.836	0.817	0.824
1024	0.828	0.842	0.828	0.846
2048	0.82	0.831	0.826	0.841

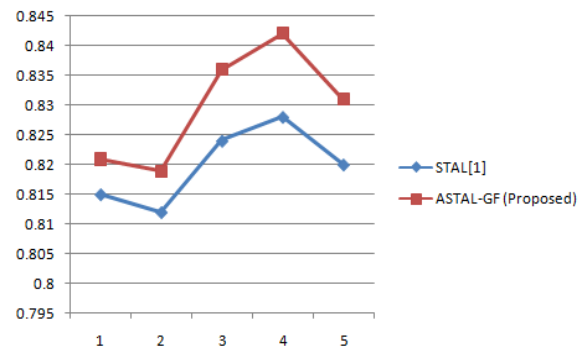


Figure 2: Rank 1 accuracy on iLIDS-VID dataset with different embedding sizes in Local region.

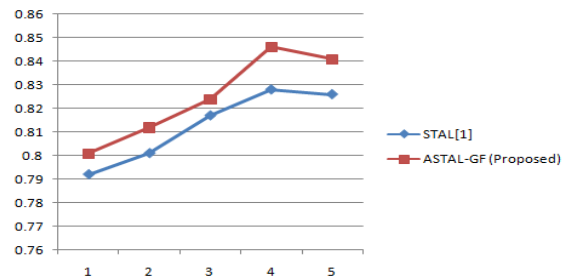


Figure 3: Rank 1 accuracy on iLIDS-VID dataset with different embedding sizes in Global region.

The installing size of highlight portrayal model is likewise vital to the ReID issue. As appeared in Fig. 1, 2, 3 we assessed impacts of various inserting sizes: {128,256,512,1024,2048} on the iLIDS-VID dataset. For the accommodation of examination, we freely explored the worldwide inserting size and nearby installing size. For instance, we fixed the nearby inserting size as 1024 while assessing the impacts of worldwide implanting sizes, and the other way around. The presentation continually increments

as implanting size expanding to 1024. While the development is dormant when the installing size increments from 1024 to 2048. The explanation might be that include parameters are soaked when implanting size is more than 1024. The impacts of both worldwide installing size and nearby inserting size are generally consistency. The worldwide implanting size is somewhat delicate than the nearby one. At last, we pick the 1024 installing size of both worldwide and neighborhood highlight in our analyses.

Table 3: Rank 1 accuracy on iLIDS-VID dataset with different margins (triplet loss).

Margin	Local		Global	
	STAL[1]	ASTAL-GF (Proposed)	STAL[1]	ASTAL-GF (Proposed)
0.4	0.816	0.822	0.818	0.826
0.6	0.819	0.827	0.82	0.834
0.8	0.828	0.831	0.823	0.841
1	0.826	0.83	0.826	0.837
1.2	0.824	0.829	0.828	0.836
1.4	0.822	0.827	0.824	0.831
1.6	0.82	0.825	0.817	0.828
1.8	0.814	0.822	0.809	0.825
2	0.809	0.814	0.802	0.818

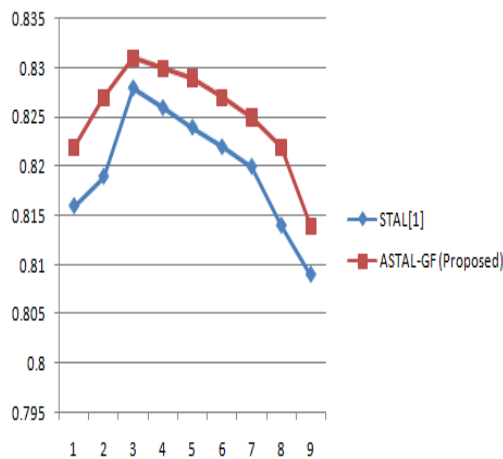


Figure 4: Rank 1 accuracy on iLIDS-VID dataset with different margins (triplet loss) in Local region.

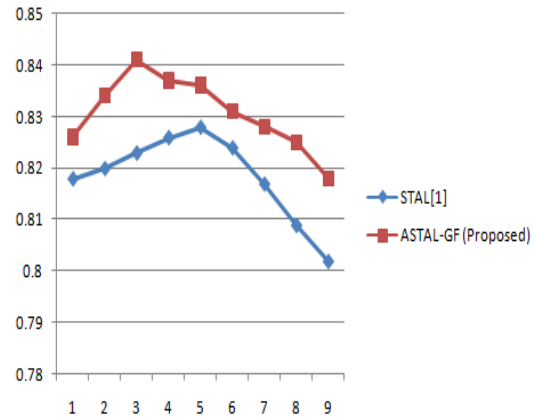


Figure 5: Rank 1 accuracy on iLIDS-VID dataset with different margins (triplet loss) in Global region.

We prepared proposed ASTAL-GF with triplet misfortune, softmax misfortune, and consistency requirement. The edges  $m_l$ ;  $m_g$  in the triplet misfortune are essential hyper parameters which influence the speculation of the model by regularization. Right now, researched how reidentification precision shifts relying upon the edges of worldwide separation and nearby separation. Assessments were performed on the iLIDS-VID dataset with the edges shifted from 0:4 to 2 with a 0.4 interim. When testing the worldwide edge, the nearby edge was fixed at 0.8; while the worldwide edge was set to 1.2 for assessing the neighborhood edge. Fig. 4, 5 represents the rank 1 exactnesses on the iLIDS-VID dataset with various edges. We see that the exhibition first builds comparing to the edge and reaches the ideal around 1. At that point, with the development of edge, the presentation drops step by step.

		Gallery Sequence Length									
		1	2	4	8	16	32	64	128	all	
Probe Sequence Length	1	0.55	0.62	0.64	0.67	0.67	0.69	0.69	0.69	0.69	
	2	0.62	0.68	0.72	0.72	0.74	0.73	0.74	0.74	0.76	
	4	0.65	0.71	0.76	0.77	0.77	0.77	0.76	0.77	0.77	
	8	0.66	0.72	0.76	0.78	0.78	0.79	0.8	0.8	0.79	
	16	0.68	0.73	0.77	0.78	0.79	0.79	0.8	0.8	0.8	
	32	0.68	0.73	0.77	0.78	0.79	0.8	0.81	0.81	0.81	
	64	0.68	0.73	0.77	0.78	0.8	0.81	0.81	0.81	0.81	
	128	0.68	0.74	0.77	0.79	0.81	0.81	0.81	0.82	0.82	
	all	0.68	0.74	0.77	0.79	0.8	0.81	0.81	0.82	0.82	

Figure 6: Rank 1 CMC performance on iLIDS-VID dataset with different lengths of both probe and gallery sequences using STAL[1].

		Gallery Sequence Length						
		1	2	4	8	16	32	64
Probe Sequence Length	1	0.61	0.68	0.69	0.69	0.69	0.7	0.7
	2	0.68	0.72	0.73	0.75	0.77	0.76	0.76
	4	0.7	0.73	0.79	0.79	0.79	0.79	0.78
	8	0.66	0.75	0.79	0.79	0.78	0.79	0.82
	16	0.68	0.76	0.8	0.81	0.82	0.8	0.82
	32	0.68	0.76	0.8	0.81	0.82	0.81	0.84
	64	0.68	0.76	0.8	0.81	0.82	0.84	0.84
	128	0.68	0.77	0.8	0.81	0.83	0.84	0.84
all	0.68	0.77	0.8	0.81	0.81	0.84	0.84	

Figure 7: Rank 1 CMC performance on iLIDS-VID dataset with different lengths of both probe and gallery sequences using ASTAL-GF (Proposed).

We explored how the presentation of the proposed strategy changes when the lengths of both test and exhibition arrangements are differed. We prepared the model on iLIDS-VID dataset as the first setting which chooses 8 casings arbitrarily in a video. During testing, we changed the lengths of both test and display groupings from 1 to 128 in steps relating with the forces of-two and give the reference which utilizes all casings. For example, we fixed the arrangement length in L. On the off chance that the genuine length of a video is more noteworthy than L, we haphazardly chose outlines as the testing grouping; else, we utilized the entire arrangement and arbitrarily examined different edges to supplement the L length. Not the same as the RNN-based technique, we chose irregular casings rather than the first or keep going L edges of the sequential grouping, since edges of a video with various stances and foundation have free accepting in our model.

### VII. CONCLUSION

This proposes a novel separation learning approach, i.e., Adaptive Spatial-Temporal Attention - Aware Learning with Gaussian Filter (ASTAL-GF), for video-based individual re-distinguishing proof, which all the while learns a couple of intra-video and bury video separation measurements. The educated intra-video separation metric can make every video increasingly smaller, with the end goal that the

extricated first-request insights highlight can all the more likely speak to every video. The scholarly between video separation metric can make the separation between genuinely coordinating recordings littler than that between wrong coordinating recordings. The result of this exposition work is as per the following:

1. The precision of worldwide and nearby highlights according to combination rate is improve by 1.02%, subsequently acknowledgment rate may improve.
2. The implanting size of highlight portrayal model is basic for individual ReID, The exactness of neighborhood and worldwide installing size is increment by 0.736% and 1.14% separately.
3. The exactness relying upon the edges of neighborhood separation and worldwide separation is increment by 0.735% and 1% separately .
4. The exhibition of re-recognizable proof procedure is improved by 2.73%, when the lengths of both test and display successions are shifted.

### VIII. FUTURE SCOPE

To additionally improve the factual of scholarly separation measurements, we propose a couple detachment based ASTAL-GF (P-ASTAL-GF), which necessitates that any two really coordinating video sets from various people ought to be all around isolated with one another under the scholarly separation measurements.

Our separation learning structure all the while learns two direct projection frameworks to manage the intra-video and between video varieties. Practically speaking, there may exist nonlinear relationship in the person on foot video information. Subsequently, if the nonlinear relationship can be all around dealt with, better execution might be accomplished by our methodologies. Considering the way that profound learning has shown amazing ability to demonstrate the nonlinearity of tests, we are keen on consolidating the profound learning procedure with our separation learning system in our future work.

## REFERENCES

- [1] Guangyi Chen, Jiwen Lu, Ming Yang and Jie Zhou, "Spatial-Temporal Attention-aware Learning for Video-based Person Re-identification", IEEE Transactions on Image Processing, 2019.
- [2] O. S. Reshma and Reshma Sheik, "Different Techniques for Person Re-Identification: A Review", Asian Journal of Computer Science and Technology ISSN: 2249-0701 Vol.8 No.1, 2019.
- [3] Xiaoke Zhu, Xiao-Yuan Jing, Xinge You, Xinyu Zhang and Taiping Zhang, "Video-based Person Re-identification by Simultaneously Learning Intra-video and Inter-video Distance Metrics", IEEE Transactions on Image Processing, 2017.
- [4] Yueqi Duan, Jiwen Lu, Jianjiang Feng and Jie Zhou, "Context-Aware Local Binary Feature Learning for Face Recognition", IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017.
- [5] Weihua Chen, Xiaotang Chen, Jianguo Zhang, Kaiqi Huang, "Beyond triplet loss: a deep quadruplet network for person re-identification", IEEE Transaction on Pattern Matching, 2017.
- [6] Chi Su, Jianing Li, Shiliang Zhang, Junliang Xing, Wen Gao, Qi Tian, "Pose-driven Deep Convolutional Model for Person Re-identification", IEEE Transaction on Pattern Matching, 2017.
- [7] Haiyu Zhao, Maoqing Tian, Shuyang Sun, Jing Shao, Junjie Yan, Shuai Yi, Xiaogang Wang, Xiaoou Tang, "Spindle Net: Person Re-identification with Human Body Region Guided Feature Decomposition and Fusion", IEEE Conference on Computer Vision and Pattern Recognition, 2017.
- [8] Dangwei Li, Xiaotang Chen, Zhang Zhang, Kaiqi Huang, "Learning Deep Context-aware Features over Body and Latent Parts for Person Re-identification", IEEE Conference on Computer Vision and Pattern Recognition, 2017.
- [9] Ji Lin, Liangliang Ren, Jiwen Lu, Jianjiang Feng, Jie Zhou, "Consistent-Aware Deep Learning for Person Re-identification in a Camera Network", IEEE Conference on Computer Vision and Pattern Recognition, 2017.
- [10] Zhen Zhou, Yan Huang, Wei Wang, Liang Wang and Tieniu Tan, "See the Forest for the Trees: Joint Spatial and Temporal Recurrent Neural Networks for Video-based Person Re-identification", IEEE Conference on Computer Vision and Pattern Recognition, 2017.
- [11] Tetsu Matsukawa<sup>1</sup>, Takahiro Okabe<sup>2</sup>, Einoshin Suzuki<sup>1</sup>, Yoichi Sato<sup>3</sup> "Hierarchical Gaussian Descriptor for Person Re-Identification", IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [12] Mengran Gou, Xikang Zhang, Angels Rates-Borras, Sadjad Asghari-Esfeden, Mario Sznaiar and Octavia Camps, "Person Re-identification in Appearance Impaired Scenarios", IEEE Conference on Computer Vision and Pattern Recognition, 2016.