

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# View-Invariant Intersection Recognition from Videos using Deep Network Ensembles

Anonymous WACV submission

Paper ID 432

## Abstract

This paper strives to answer the following question: Is it possible to recognize an intersection when seen from different road segments that constitute the intersection? An intersection or a junction typically is a meeting point of three or four road segments. Its recognition from a road segment that is transverse to or 180 degrees apart from its previous sighting is an extremely challenging and yet a very relevant problem to be addressed from the point of view of both autonomous driving as well as loop detection. This paper formulates this as a problem of video recognition and proposes a novel LSTM based Siamese style deep network for video recognition. For what is indeed a challenging problem and the limited annotated dataset available we show competitive results of recognizing intersections when approached from diverse viewpoints or road segments. Specifically, we tabulate effective recognition accuracy even as the approaches to the intersection being compared are disparate both in terms of viewpoints and weather/illumination conditions. We show competitive results on both synthetic yet highly realistic data mined from the gaming platform GTA as well as on real world data made available through Mapillary.

## 1. INTRODUCTION

Recognizing an intersection from a different approach sequence or a sequence of viewpoints different from one seen before can be pivotal in various applications that include autonomous driving, driver assistance systems, loop detection for SLAM including multi-robot and multi-session SLAM frameworks. An immediate benefit could be that data can be exchanged between robots only around areas of intersection than across entire runs to detect and close loops, an idea that was put forth in an indoor robotic context earlier [20]. Retrieving previously seen intersections can also be advantageous from the point of view large-scale outdoor topological mapping frameworks.

Intersection recognition from disparate video streams or image sequences is an extremely challenging problem that stems due to large variation in viewpoint, weather and appearance. Complexity also emanates from varying levels of traffic and chaos at a junction<sup>1</sup> between two sequences, as well as due to changing levels of occlusion, illumination between two video sequences of a specific intersection. Lack of annotated datasets on intersections also poses a challenge.

Within the robotics community, loop detection methods have used diverse image recognition and retrieval techniques [16, 15, 6] that have not attempted to detect loops from very diverse viewpoints while they perform admirably in the presence of weather changes [23] or under the duress of varying traffic [17]. In the vision community, CNN [31] features have been used to efficiently compare scenes or structures which are viewed from varying depths, which give a zoom-in vis-a-vis zoom-out effect. Whereas in this paper, we argue that recognizing intersections has many potential benefits and therefore needs an attention and treatment of its own.

We propose a novel stacked deep network ensemble architecture that combines state-of-the-art CNN, Bidirectional LSTM and Siamese style distance function learning for the task of view-invariant intersection recognition in videos. While the CNN component of the ensemble primarily deals with image-level feature representation, the bidirectional LSTM is the key recurrent network component that enables learning of temporal evolution of visual features in videos followed by the Siamese network-based distance metric learning for comparing two input video streams. Our proposed network is conceptualized in Figure 1.

We contribute in the following ways:

1. Firstly we propose a novel and pertinent problem of recognizing intersections when approached from highly disparate viewpoints

<sup>1</sup>In this work, we use the words junctions and intersections interchangeably.

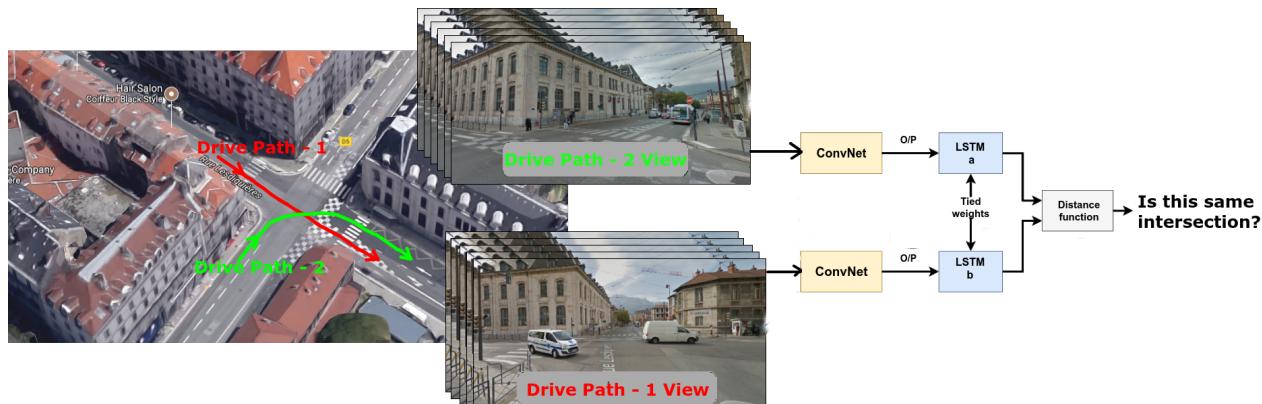
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120

Figure 1. Trajectories shown with red and green are separate drive sessions taken at the same intersection. Corresponding videos are given to our proposed network to determine whether they are belonging to same intersection or not.

121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139

2. We propose an original deep network ensemble. The proposed architecture can handle videos of varying length and compare videos capturing reverse trajectories. To the best of our knowledge, this kind of Siamese network with recurrent LSTM component has been largely unexplored in the video domain. Furthermore, we use the hidden state of LSTM cells, instead of the traditionally used output state, for video representation, which has largely been unexplored too.
3. We showcase the efficacy of the proposed architecture through competitive results on challenging sequences. We use synthetic yet highly realistic data from the GTA [2] gaming platform and actual real world data of Mapillary [4] for this purpose.

140  
141  
142  
143  
144  
145  
146  
147  
148  
149

To the best of our knowledge, this is the first attempt in this direction using the ensemble of state-of-the-art CNN, bidirectional LSTM and Siamese network architecture for recognizing intersection in input video pairs. Highly unpredictable nature of road intersections makes them vulnerable to accidents. Nevertheless, proposing robust solutions for autonomous vehicles aware of the intersections and the being able to recognize them based on a previously explored pathways is much needed.

150  
151

## 2. Literature Review

152  
153  
154  
155

In existing literature around robotic perception, there are limited efforts towards intersection detection and recognition. Some of these use sensors other than cameras such as laser and virtual cylindrical scanners [35, 22].

156  
157  
158  
159  
160  
161

Recently, a framework for detecting intersections from videos was proposed in [11]. They used LSTM based architecture to observe the relative change in outdoor features to detect intersection. Nevertheless, their work did not focus on intersections recognition task. On the other hand, visual loop closure detection techniques in the literature [7, 16, 15]

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193

mainly focused on image level comparison of scenes. These techniques try to find re-occurring scene during the driving session based on fast KD tree based comparison of image features. Although these methods achieve admirable accuracy for loop detection they cannot be improvised for view-invariant recognition over videos.

In recent years, gaming environments have been created and/or used for generating and annotating datasets. [27] used Grand Theft Auto V (GTA) gaming-engine to create large-scale pixel-accurate ground truth data for training semantic segmentation systems. SYNTHIA [28] is a virtual world environment which was created to automatically generate realistic synthetic images with pixel-level annotations for the task of semantic segmentation. [27] and [28] show the added improvement in the real-world setting by using synthetic-datasets for training deep-network models. CARLA [13] is another open-source simulator developed to help in development, training and validation of Autonomous Driving. Visually, GTA environment is more realistic than SYNTHIA and CARLA environments.

Real-world datasets such as [25, 1] are designed with the idea of visual place recognition using images. These datasets are limited in the pathways covered and the number of intersections. Another real world street-level imagery platform [4] consist of images and image-sequences with variability in the pathways and intersections and thus enables sequence learning for intersection recognition. Concurrent works have used this platform to create Mapillary Vistas Dataset [26] which is a semantic segmentation dataset consisting of 25k images. We have used Mapillary dataset to showcase our result on real world scenarios.

Specifically there has been no prior effort towards recognizing intersections when approached from different road segments that constitute the intersection.

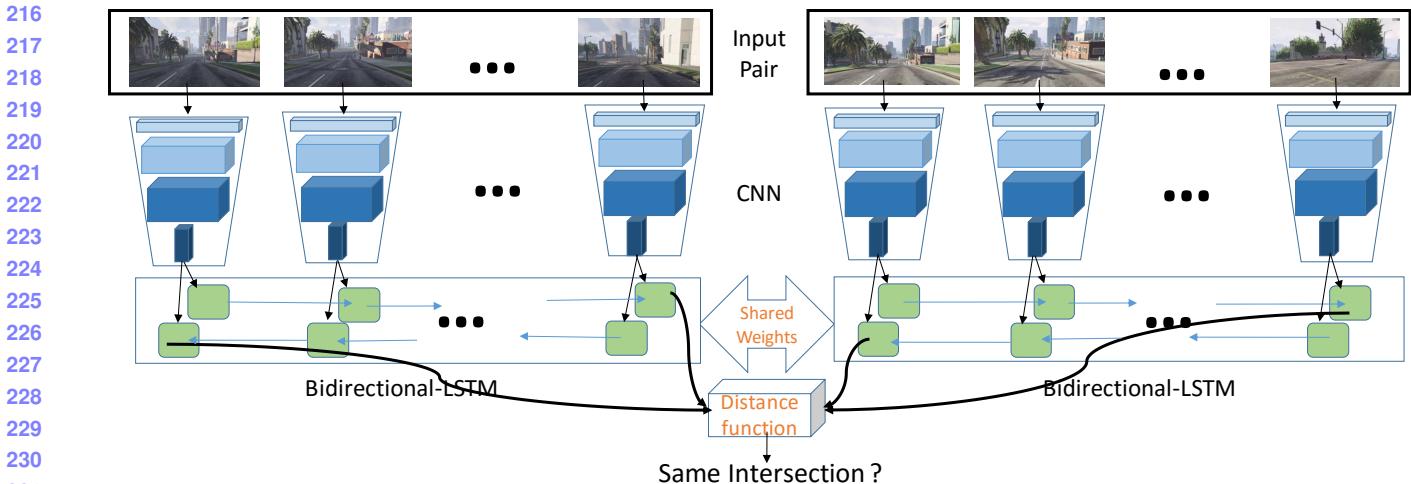


Figure 2. Proposed model is shown with video pairs as input and binary-classification as the output. Features from pretrained CNN network are fed into bidirectional LSTM with shared-weights as shown with hidden units in green (unfolded in time).

### 3. Method

In order to achieve view-invariant intersection recognition from video sequences, our overall approach follows the recent successes of deep network models in learning useful representations of visual data [9]. Our overall strategy comprises of three specific sub-objectives: (i) capture representations of individual video frames that can help identify the intersection in a view-invariant manner; (ii) leverage the temporal correlations between video frames to better capture the unique identity of an intersection across different views; and (iii) use an appropriate distance metric to compare the thus learned spatiotemporal representations of intersection video sequences. We use a stacked deep network ensemble architecture to realize the above strategy.

Our stacked ensemble consists of a state-of-the-art Convolutional Neural Network (CNN) to learn video frame representations, a Bidirectional Long Short-Term Memory (LSTM) Network (a variant of Recurrent Neural Networks, that addresses the vanishing gradient problem) network that captures temporal correlations across video frames, and a Siamese Network that learns a suitable distance metric that can uniquely identify an intersection from these video sequences across views. Figure 2 summarizes the proposed ensemble strategy.

#### 3.1. Convolutional Neural Networks

Convolutional Neural Networks (CNNs) based models have recently enjoyed tremendous success in a variety of tasks including image classification [21], image-segmentation [8], activity-recognition [33], etc. CNN is a biologically-inspired model which sequentially stacks a series of convolution operations on an image followed by non-linear activations and pooling operations, and are ca-

pable of capturing low-level and high-level features from the image. It has been shown that codes learned by the CNN can be used as off-the-shelf features for a variety of tasks [30], [32].

From a plethora of CNN pretrained models, we choose the model most relevant to our tasks which provide invariance in lighting and pose. In particular, we employ AmosNet [12] as our CNN architecture. AmosNet is a variation of CaffeNet [21] trained for Visual Places Recognition under varying weather conditions. We plug in AmosNet with pretrained weights and use the *pool6* outputs as input features for the LSTM.

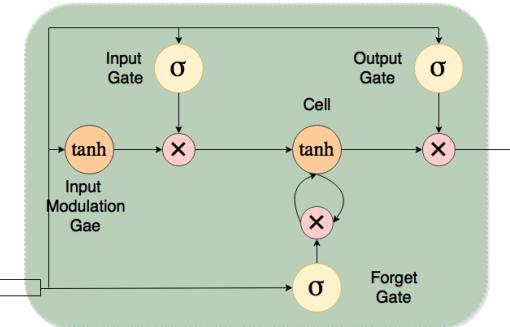


Figure 3. LSTM structure. Expansion of green block shown in Figure 2 (reproduced from [14])

#### 3.2. Long Short-Term Memory Network

Recurrent Neural Networks (RNNs) are neural networks adapted for time-varying sequential data by incorporating a feedback-loop in the architecture (Figure 3). In theory, RNNs should be able to model temporal correlations in long-term sequences; however, they often fail in practice

324 due to the exploding/vanishing gradient problem [10].  
 325

326 Long Short-term Memory (LSTM), introduced by  
 327 Hochreiter and Schmidhuber [18], is a variant of RNN with  
 328 the ability to capture long-term dependencies [19]. Similar  
 329 to RNN, it has a hidden representation which is updated  
 330 at each time step, but it also consists of a memory state  
 331 ( $\mathbf{c}_t$ ), which holds the relevant information at each time step.  
 332 The information entering and exiting out of cell states is  
 333 controlled by gating mechanisms (sigmoid, and tanh functions).  
 334 An input gate  $\mathbf{i}_t$  decides what part of the information  
 335 entering the LSTM through  $\hat{\mathbf{g}}_t$ , should be held in the current  
 336 cell state,  $\mathbf{c}_t$ . The forget gate  $\mathbf{f}_t$  decides what to forget from  
 337 the previous cell state,  $\mathbf{c}_{t-1}$ . The output gate  $\mathbf{o}_t$  determines  
 338 how much of  $\mathbf{c}_t$  should be visible as output of the LSTM  
 339 at time  $t$ . This system is captured by the set of equations  
 340 below:

$$\begin{aligned} \mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i) \\ \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f) \\ \mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o) \\ \hat{\mathbf{g}}_t &= \tanh(\mathbf{W}_g \mathbf{x}_t + \mathbf{U}_g \mathbf{h}_{t-1} + \mathbf{b}_g) \\ \mathbf{c}_t &= \mathbf{i}_t \odot \hat{\mathbf{g}}_t + \mathbf{f}_t \odot \mathbf{c}_{t-1} \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh \mathbf{c}_t \end{aligned}$$

341 Here  $\sigma(\mathbf{x})$  is a sigmoid function which operates on each element  
 342 of the vector  $\mathbf{x}$  separately and outputs a value between  
 343 0 and 1, and  $\odot$  represents the Hadamard product. LSTMs  
 344 are capable of handling variable-length inputs and unfolding  
 345 each input dynamically.

346 Bidirectional RNNs [29] incorporate future and past context  
 347 by running the inverse of the input through a separate  
 348 RNN. A similar extension for LSTMs is called Bidirectional  
 349 LSTMs which have shown good performance in many applications  
 350 that involve sequential data such as text [24]. The proposed model  
 351 employed Bidirectional LSTMs to capture the temporal correlations  
 352 between video frames. This enables our model to identify the intersection  
 353 uniquely whether a vehicle approaches it in the forward or opposite  
 354 direction.

355 Instead of using averaged/fused output of unfolded  
 356 LSTM units or output from last unit we use the hidden-  
 357 representation from the last unfolded time step as a feature  
 358 vector for the videos as shown in Figure 2. This is because,  
 359 we have empirically observed that the final results while using  
 360 LSTM's hidden state outperform results with traditionally  
 361 used regular output state.

### 362 3.3. Siamese Network

363 The third component of our deep network ensemble is a  
 364 Siamese network, which is used to learn data-driven  
 365 distance metrics. A Siamese Network [34] consists of two  
 366 independent networks with same architecture and shared  
 367 weights. Both these networks are merged to learn a distance

368 metric,  $d$ , between the two inputs provided to the respective  
 369 networks.

370 During training, Siamese networks work on triplets  $(\mathbf{x}_1, \mathbf{x}_2, y)$  where  $y$  is the ground truth similarity between  $\mathbf{x}_1$  and  
 371  $\mathbf{x}_2$ , i.e.,  $y = 1$  if the videos denote the same intersection,  
 372 else  $y = 0$ . The networks' weights are then updated by  
 373 minimizing the loss function described below.

374 **Contrastive Loss:** The total loss function over a dataset  
 375  $X = \{(\mathbf{x}_1^i, \mathbf{x}_2^i, y^i), \forall i = 1, \dots, n\}$  is given by:

$$L_{\mathbf{W}}(X) = \sum_{i=1}^n L_{\mathbf{W}}^i((\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)), y^i) \quad (1)$$

376 where  $\mathbf{W}$  are the weights/parameters of the network,  $\phi(\mathbf{x})$   
 377 denotes the output of the last layer of the shared network  
 378 architecture for each individual input  $\mathbf{x}$ , and  $L_{\mathbf{W}}^i$  is  
 379

$$L_{\mathbf{W}}^i = y^i L_{pos}(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)) + (1 - y^i) L_{neg}(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)) \quad (2)$$

380 where

$$L_{pos}(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)) = d(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i))^2 \quad (3)$$

$$L_{neg}(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)) = \max(1 - d(\phi(\mathbf{x}_1^i), \phi(\mathbf{x}_2^i)), 0)^2 \quad (4)$$

381 At test time, the learned distance function is used to  
 382 predict the similarity in the videos using the expression in  
 383 Eq. 5.

$$\hat{y} = \begin{cases} 0 : d(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2)) > \theta \\ 1 : d(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2)) \leq \theta \end{cases} \quad (5)$$

384 We set the value of  $\theta$  to be 0.5 in our experiments. Since  
 385 the number of possible negative pairs can be very high as  
 386 compared to positive pairs, we scale the loss function for  
 387 each positive by a constant factor ( $> 1$ ), given by the ratio  
 388 of negative pairs to positive pairs in the dataset.

### 389 3.4. Training and Network Parameters

390 Although our model can be trained in an end to end manner  
 391 we train it in a greedy way based on empirical observation.  
 392 We use pretrained weights from previous state-of-the-art  
 393 CNN models which provide a meaningful weight-  
 394 initialization to our model and reduces training time. We  
 395 use the contrastive loss described in Eq. 1 to train the Bi-  
 396 directional LSTM, as shown in Figure 2. The proposed model  
 397 is trained using the ADAM optimizer which is a variant of  
 398 Stochastic Gradient Descent (SGD).

399 In our initial experiments, we varied the number of Bi-  
 400 directional LSTM-layers (upto 3), but the performance gain  
 401 was negligible. Similarly, we experimented with different  
 402 values (5, 10, 50, 80, 150, 250) for hidden-unit dimensions  
 403 in the LSTM cell. Empirical results found 80 to be the best  
 404 performing after which the performance saturates. Thus,  
 405 we fix the number of Bidirectional layer to one and hidden  
 406 layer dimension to be 80 in all further experiments.

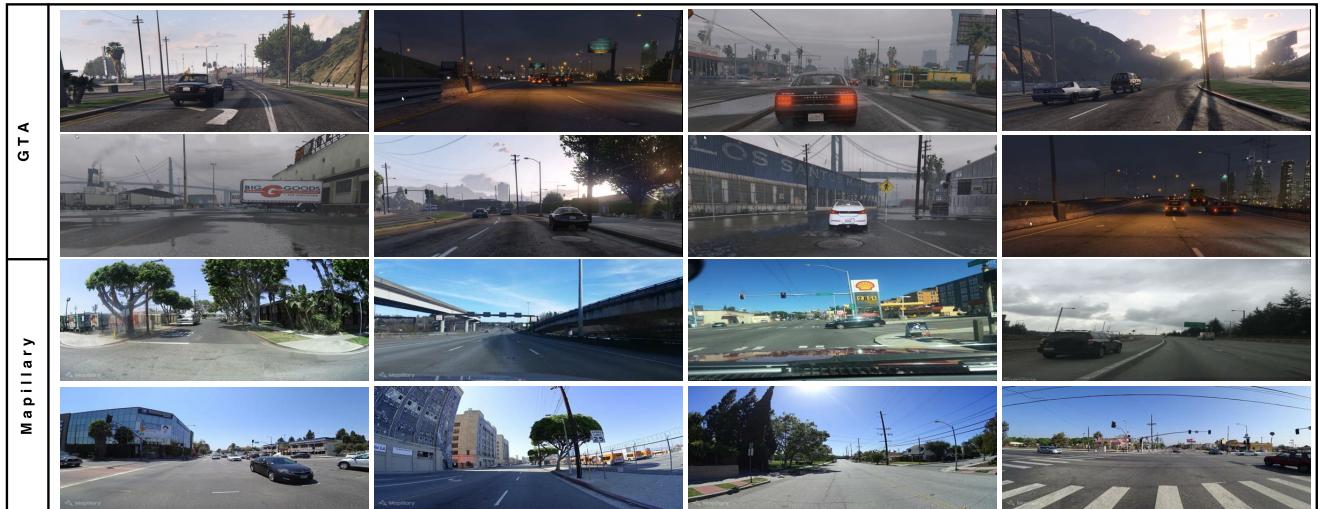


Figure 4. Data Visualization: Random snapshots from GTA environment [2] (Row 1,2) and Mapillary [4] (Row 3,4). Images show different weather and day/night conditions as well as various outdoor scenes for urban scenarios. The game-environment snapshots look visibly similar to that of real world images and have significantly higher variations in weather and lighting. (Best seen in color)

A pair of videos  $x_1, x_2$  are passed as a set of contiguous frames to the model. Each frame in the video is an RGB image with a fixed size ( $227 \times 227 \times 3$ ). For each frame, the CNN feature ( $6 \times 6 \times 256$  which is resized to  $1 \times 9216$ ) is computed. These CNN features from each video are fed into two separate Bidirectional LSTM with tied-weights. The hidden unit representation ( $\phi(x_1), \phi(x_2)$ ) for respective input videos from the last cell in the Bidirectional LSTM network is then used as a feature for the video. Further, we compute the distance  $d$  between  $\phi(x_1)$  and  $\phi(x_2)$  using the Siamese network approach, and use this to predict  $\hat{y}$  if the videos are of same intersection or not as mentioned before in Section 3.3.

## 4. Experiments & Results

In this section, we provide detailed description of datasets synthetically generated and real datasets for various experiments. Subsequently we describe the various experimental scenarios which is followed by the explanation of the quantitative and qualitative results.

### 4.1. Dataset Description

Synthetic data for this task is collected from GTA [2] gaming environment while real world data is mined from publicly available street images/sequences from Mapillary [4]. Figure 4 shows random snapshots captured from GTA environment (Row 1,2) and Mapillary (Row 3,4). One can infer from the figure that the snapshots from GTA environment are visually realistic as well as similar to real-world data from Mapillary and contain the key challenges associated with real environments (as listed in Section 1). However, the GTA environment offers more complexity in

terms of weather and lighting variations as compared to Mapillary.

Each video/sequence consist of a series of frames collected around junction points. Each junction can be traversed in multiple pathways. We refer to the set of all such possible traversals as *trajectories*.

#### 4.1.1 GTA

Here, we collected videos in two different scenarios. Firstly, we choose a set of intersections in the game-environment and then sample them in a *dense* manner, i.e., each possible trajectory in the junction is traversed. Secondly, we traversed a car from one point to another via means of an AI-car driving mode [3]. From such traversal, we break the captured video into small chunks involving intersections and non-intersections. We discard video chunks which have no intersection.

In the first scenario, we selected 27 unique junctions where 9 of them had arbitrary lighting conditions and the rest were in normal daylight condition. In the second scenario, we captured 12 traversals in two arbitrary different lighting variations.

#### 4.1.2 Mapillary

We download images from [4], which is a community-led service for people who collaboratively want to visualize the world with street-level photos. It has more than 200 million images [5]. These are captured in various modes including walking, riding (either a bike or car), panorama and photo-spheres. The data is available under a CC-BY-SA license agreement.

540 Table 1. Number of video pairs in training, testing and the validation set for the different experimental scenarios  
541

| Dataset   | Lighting Setting | Trajectory Relation            | Training |          | Validation |          | Testing  |          |
|-----------|------------------|--------------------------------|----------|----------|------------|----------|----------|----------|
|           |                  |                                | Positive | Negative | Positive   | Negative | Positive | Negative |
| GTA       | day and night    | All combinations               | 1509     | 1440     | 459        | 612      | 106      | 104      |
|           |                  | Overlap in view and trajectory | 640      | 1440     | 168        | 212      | 62       | 104      |
|           | day              | All combinations               | 864      | 1166     | 300        | 653      | 72       | 82       |
| Mapillary | day              | Overlap in view and trajectory | 3080     | 3328     | 421        | 428      | 400      | 400      |
|           |                  | All combinations               | 6409     | 6976     | 350        | 318      | 425      | 420      |

551 We use mapillary’s API to download images and construct *trajectories*. We first query for all the images in a  
552 bounding box defined by the longitudes and latitudes. For  
553 every image in the bounding box, the API provides a tuple  
554 consisting of image-key, latitude, longitude, orientation and  
555 the (video-)sequence it belongs to along with other meta-  
556 data information. These images can be then downloaded  
557 using the image-key. Using two orthogonally overlapping  
558 sequences we get a location in the map which is used to  
559 identify the intersection. Subsequently, this is used to down-  
560 load a subset of continuous images from all the sequences  
561 passing through the intersection. We generate sequences  
562 representing multiple trajectories at an intersection from the  
563 provided tuples.  
564

565 From the visual analysis of the extracted data, we found  
566 that the dataset lacks different lighting settings (predomi-  
567 nantly capturing morning/day hour images). This can be  
568 attributed to the limitations of data collection in real world  
569 settings.  
570

571 We mined around 300,000 images which consisted of  
572 around 1700 usable sequences from around 500 junctions.  
573 For each junction, we get 3 to 6 sequences which implies  
574 that the junctions are not sampled in a *dense* manner.  
575

## 576 4.2. Evaluation Metric

577 We report percentage accuracy as the metric in our  
578 experiments. Accuracy is defined as the ratio of the number  
579 of correctly classified samples (both positive and negative  
580 pairs combined) to the total predictions made.  
581

## 582 4.3. Experimental Setup

583 From the extracted dataset, we first annotate positive and  
584 negative video pairs. Positive (videos involving the same  
585 intersection) pairs are mined from the above datasets ex-  
586 haustively, i.e., we select all possible positive pairs. Neg-  
587 ative pairs generated from the datasets can be very high as  
588 any two trajectories from different junction are treated as a  
589 negative pair. Hence we limit the number of negative pairs  
590 by randomly fixing a subset.  
591

592 To make sure that test, train and validation sets are mu-  
593 tually exclusive, we sampled these from different junctions

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

551 Table 2. Accuracy of our method on different datasets for the  
552 different experimental scenarios

| Dataset   | Lighting setting | Trajectory Relation            | Accuracy |
|-----------|------------------|--------------------------------|----------|
| GTA       | day and night    | All combinations               | 70.95    |
|           |                  | Overlap in view and trajectory | 78.2     |
|           | day              | All combinations               | 76.72    |
| Mapillary | day              | Overlap in view and trajectory | 81.0     |
|           |                  | All combinations               | 72.1     |

553 and/or non-overlapping traversals. Specifically, in GTA for  
554 the test set, we maintain an exclusive set of 3 intersections  
555 which are only sampled from junctions while for training  
556 and validation, we use both scenarios as detailed in Sec-  
557 tion 4.1. Similarly, in Mapillary we maintain the exclusiv-  
558 ity in train, validation and test data by extracting sequences  
559 from mutually exclusive bounding boxes.  
560

561 For a pair of videos involving the same junction there  
562 exists a trajectory-relation between them, which can be  
563 categorized in terms of overlap in view and/or trajectory.  
564 Overlap in trajectory refers to the proximal relation in the  
565 two trajectories while overlap in view refers to trajectories,  
566 viewing the larger part of the same area. For example tra-  
567 jectories moving in opposite direction can be proximal in  
568 distance but can have minimal view overlap.  
569

570 For the purpose of our experiments we categorize the  
571 trajectory-relations into two primary setups: *overlap in view*  
572 and *trajectory* and *All combinations*. The former refers to  
573 parallel and overlapping trajectory-relations while the latter  
574 refers to all the possible trajectory-relations. Similarly, we  
575 categorize the lighting setting into two categories: *day* vs  
576 *day and night*. The names are self-explanatory.  
577

578 We test the robustness and generalizing capability of  
579 our model under various trajectory-relation and lighting  
580 (day/night) settings. The details of the experiment-wise dis-  
581 tribution of samples can be found in Table 1. For GTA,  
582 we sampled the video-pairs in three categories based on  
583 trajectory-relations and lighting conditions. We discarded  
584 the *day* and *Overlap in view and trajectory* setting due to  
585

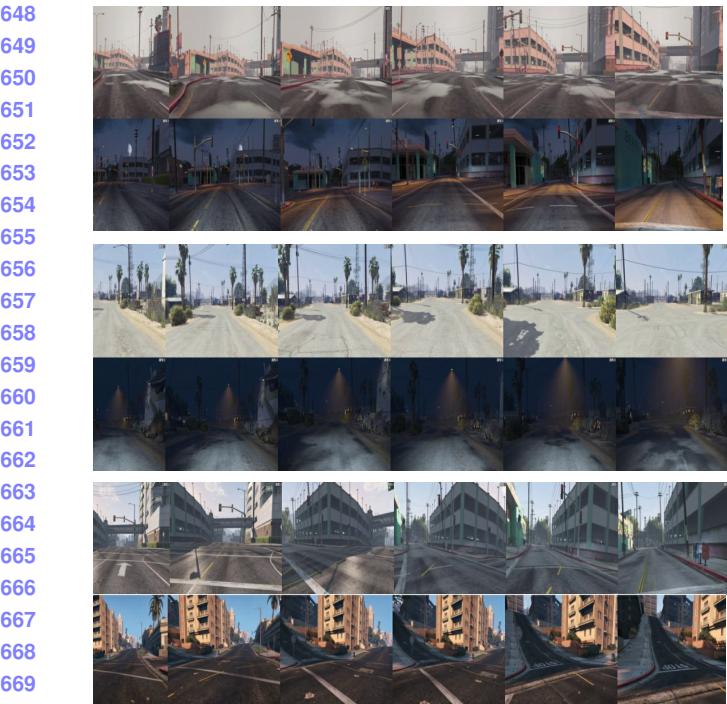
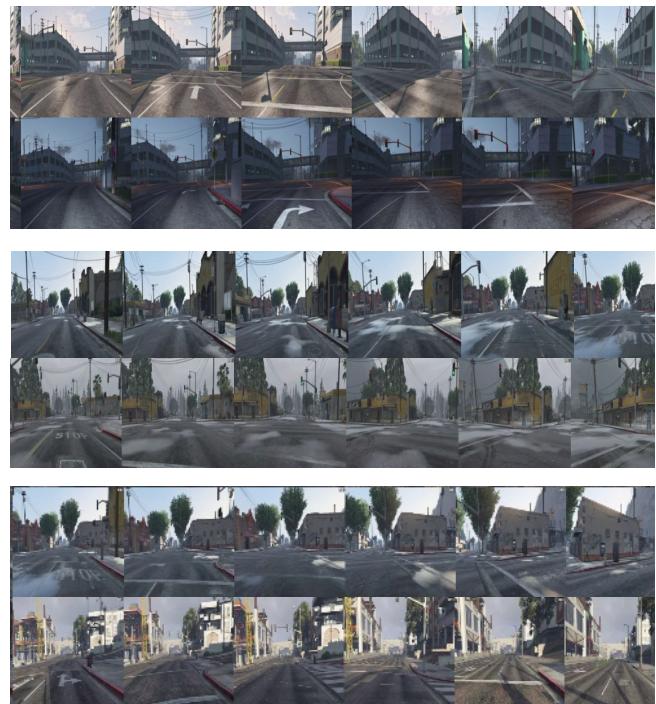


Figure 5. Qualitative results on GTA: Six randomly selected video-pairs (each video is represented as a set of contiguous frames). Row-1,2 show 4 correctly classified positive video-pairs. Row-3 shows 2 correctly classified negative video-pairs. These examples show that our model is able to capture view and lighting variations. Here, we just show 6 similar continuous frames per video, corresponding frames in the two videos need not be temporally aligned (Section 4.6). (Best seen in color)



insufficient training data. For Mapillary, we decided on the two categories as the dataset predominantly consisted of day/morning images.

#### 4.4. Quantitative Results

Table 2 depicts results on GTA and Mapillary dataset under varying scenarios using our proposed model. For GTA, under all possible lighting conditions (*day and night*) the model performs better by 7% when video-pairs have an *overlap in view and trajectory* as compared to *all combinations*. This points to the fact that the generic setting involving all trajectory relations is harder to model as the view variations increase significantly. When the lighting condition is set to *day* only, for all trajectory-relations the model gives an accuracy of 76.72% which is 6.5% greater than *day and night* scenario. This can be attributed to the reduced complexity in lighting variations.

Similarly, for mapillary the model performs better when the trajectory-relations are limited to *overlap in view and trajectory* as compared to *All combinations* by 8.9%.

#### 4.5. Qualitative Results

In this subsection, we fix our model to the most generic scenario (all possible lighting variations and trajectory relations). Figure 5 shows the qualitative result on a subset of 4 positive (Row-1, 2) and 2 negative video pairs (Row-3) from

the test-set of GTA using this model. All these video pairs are correctly classified by the model. It is easy to see that our learned model is able to generalize to varying weather, lighting and view conditions. For example: video-pair in Row-2, Column-2 capture the different view of the same structure in rainy weather. Similarly, in Row-1, Column-1 the model is able to predict similar videos in presence of view and lighting variations. The shown examples are in coherence with the challenges in intersection-detection as mentioned in Section 1. In Figure 5 we show 6 contiguous frames from each sequence, where we manually crop the relevant contiguous frames for illustration purposes (see Supplementary material for full videos).

Similarly, Figure 7 show the qualitative result on 3 correctly classified positive video pairs. Each video-pair has a different trajectory-relation. Row-1 video-pair are similar trajectories, Row-2 video-pair are partially-overlapping trajectories and Row-3 video-pair are orthogonal trajectories.

#### 4.6. Discussion

In addition to qualitative and quantitative results it is important to understand the challenges associated with the problem. Here, we explicitly discuss this in the context of positively labeled video-pair from GTA data shown in Figure 6. On careful observation, it becomes evident that only a few corresponding frames in the two videos are similar

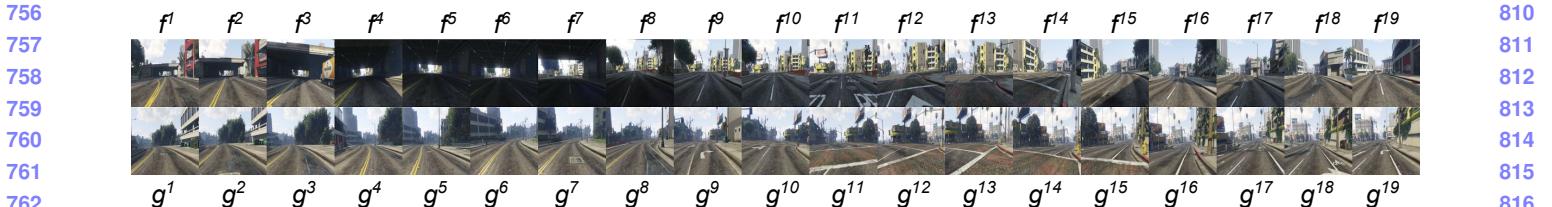


Figure 6. A positively annotated video-pair consisting of 19 frames each. Here,  $f^i$  and  $g^i$  represent the frames in the temporal order for the two videos respectively. The videos have an overlapping trajectory in opposite directions. It is non trivial to observe that a partial structural overlap exists between a few frames (Observe  $f^{19}$  and  $g^1$ ,  $f^{10}$  and  $g^{11}$ ,  $f^{11}$  and  $g^{12}$ ). (Best seen in color)



Figure 7. Qualitative results on Mapillary: . (Best seen in color)

due to the presence of partial overlapping structure in them. Furthermore these similar frames are not in the same order in the two videos. This reflects one should not interpret the quantum of the results in this context.

To show the limitation of our model, we present a failure case shown in Figure 8, where a negative pair is classified as a positive pair by our network. However visual inspection reveals that in this case, due to presence of large illumination change and less discriminative structure, it is challenging even for human beings to correctly classify the video pair.

Thus, in our view the proposed architecture not stellar but performs significantly impressive for such complex and ill posed problem of view-invariant intersection.



Figure 8. A failure case for the model. Row-1, 2 are contiguous frames from the two videos. The video-pairs belong to different junctions but is classified as same intersection. Images in Row-1 are structurally similar to that of Row-2. (Best seen in color)

## 5. CONCLUSIONS

This paper proposed a novel stacked deep network ensemble architecture that combines state-of-the-art CNN, bidirectional LSTM and Siamese style distance function learning for the task of view-invariant intersection recognition in videos. The proposed architecture enables recognizing the same intersection across two videos of variable length having large view variations, inverted trajectory, lightning and weather variations. We have collected annotated data (more than 2000 videos) from GTA [2] and Mapillary [4] and have computed results on this data with varying parameter choices and have reported competitive results.

## References

- [1] Gardens Point Dataset. <https://wiki.qut.edu.au/display/cyphy/Day+and+Night+with+Lateral+Pose+Change+Datasets>.
- [2] Grand Theft Auto V. [https://en.wikipedia.org/wiki/Development\\_of\\_Grand\\_Theft\\_Auto\\_V](https://en.wikipedia.org/wiki/Development_of_Grand_Theft_Auto_V).
- [3] Grand Theft Auto V Auto-drive Mod. <https://www.gta5-mods.com/scripts/vautodrive>.
- [4] Mapillary. <https://www.mapillary.com/app>.
- [5] Mapillary Wikipedia. <https://en.wikipedia.org/wiki/Mapillary>.

- |     |   |     |
|-----|---|-----|
| 864 | [6] S. Achar, C. V. Jawahar, and K. M. Krishna. Large scale           | 918 |
| 865 | visual localization in urban environments. In <i>ICRA</i> , 2011.     | 919 |
| 866 |   | 920 |
| 867 | [7] A. Angeli, D. Filliat, S. Doncieux, and J. Meyer. Fast and        | 921 |
| 868 | incremental method for loop-closure detection using bags of           | 922 |
| 869 | visual words. <i>TRO</i> , 2008.                                      | 923 |
| 870 | [8] V. Badrinarayanan, A. Kendall, and R. Cipolla. SegNet: A          | 924 |
| 871 | deep convolutional encoder-decoder architecture for image             | 925 |
| 872 | segmentation. <i>arXiv preprint arXiv:1511.00561</i> , 2015.          | 926 |
| 873 | [9] Y. Bengio, A. Courville, and P. Vincent. Representation           | 927 |
| 874 | learning: A review and new perspectives. <i>TPAMI</i> , 2013.         | 928 |
| 875 | [10] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term        | 929 |
| 876 | dependencies with gradient descent is difficult. <i>Transactions</i>  | 930 |
| 877 | <i>on neural networks</i> , 1994.                                     | 931 |
| 878 | [11] D. Bhatt, D. Sodhi, A. Pal, V. Balasubramanian, and K. M.        | 932 |
| 879 | Krishna. Have I reached the intersection: A deep learning-            | 933 |
| 880 | based approach for intersection detection from monocular              | 934 |
| 881 | cameras. In <i>IROS</i> , 2017.                                       | 935 |
| 882 | [12] Z. Chen, A. Jacobson, N. Sunderhauf, B. Upcroft, L. Liu,         | 936 |
| 883 | C. Shen, I. Reid, and M. Milford. Deep learning fea-                  | 937 |
| 884 | tures at scale for visual place recognition. <i>arXiv preprint</i>    | 938 |
| 885 | <i>arXiv:1701.05105</i> , 2017.                                       | 939 |
| 886 | [13] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and              | 940 |
| 887 | V. Koltun. CARLA: An open urban driving simulator. In                 | 941 |
| 888 | <i>CoRL</i> , 2017.   | 942 |
| 889 | [14] M. Fayyaz, M. H. Saffar, M. Sabokrou, M. Fathy, R. Klette,       | 943 |
| 890 | and F. Huang. STFCN: Spatio-temporal FCN for semantic                 | 944 |
| 891 | video segmentation. <i>arXiv preprint arXiv:1608.05971</i> , 2016.    | 945 |
| 892 | [15] D. Gálvez-López and J. D. Tardos. Bags of binary words for       | 946 |
| 893 | fast place recognition in image sequences. <i>TRO</i> , 2012.         | 947 |
| 894 | [16] A. Glover, W. Maddern, M. Warren, S. Reid, M. Milford,           | 948 |
| 895 | and G. Wyeth. OpenFABMAP: An open source toolbox for                  | 949 |
| 896 | appearance-based loop closure detection. In <i>ICRA</i> , 2012.       | 950 |
| 897 | [17] A. A. Hafez, M. Singh, K. M. Krishna, and C. V. Jawahar.         | 951 |
| 898 | Visual localization in highly crowded urban environments.             | 952 |
| 899 | In <i>IROS</i> , 2013.  | 953 |
| 900 | [18] S. Hochreiter and J. Schmidhuber. Long short-term memory.        | 954 |
| 901 | <i>Neural computation</i> , 1997.                                     | 955 |
| 902 | [19] A. Karpathy, J. Johnson, and L. Fei-Fei. Visualizing             | 956 |
| 903 | and understanding recurrent networks. <i>arXiv preprint</i>           | 957 |
| 904 | <i>arXiv:1506.02078</i> , 2015.                                       | 958 |
| 905 | [20] A. K. Krishnan and K. M. Krishna. A visual exploration           | 959 |
| 906 | algorithm using semantic cues that constructs image based             | 960 |
| 907 | hybrid maps. In <i>IROS</i> , 2010.                                   | 961 |
| 908 | [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet          | 962 |
| 909 | classification with deep convolutional neural networks. In            | 963 |
| 910 | <i>NIPS</i> , 2012.   | 964 |
| 911 | [22] Q. Li, L. Chen, Q. Zhu, M. Li, Q. Zhang, and S. S. Ge. Intersec- | 965 |
| 912 | tion detection and recognition for autonomous urban                   | 966 |
| 913 | driving using a virtual cylindrical scanner. <i>IET Intelligent</i>   | 967 |
| 914 | <i>Transport Systems</i> , 2013.                                      | 968 |
| 915 | [23] C. Linegar, W. Churchill, and P. Newman. Made to measure:        | 969 |
| 916 | Bespoke landmarks for 24-hour, all-weather localisation               | 970 |
| 917 | with a camera. In <i>ICRA</i> , 2016.                                 | 971 |
|     | [25] M. J. Milford and G. F. Wyeth. Seqslam: Visual route-based       |     |
|     | navigation for sunny summer days and stormy winter nights.            |     |
|     | In <i>ICRA</i> , 2012.  |     |
|     | [26] G. Neuhold, T. Ollmann, S. R. Bulò, and P. Kontschieder.         |     |
|     | The mapillary vistas dataset for semantic understanding of            |     |
|     | street scenes. In <i>ICCV</i> , 2017.                                 |     |
|     | [27] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for    |     |
|     | data: Ground truth from computer games. In <i>ECCV</i> , 2016.        |     |
|     | [28] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and              |     |
|     | A. Lopez. The SYNTHIA Dataset: A large collection of                  |     |
|     | synthetic images for semantic segmentation of urban scenes.           |     |
|     | In <i>CVPR</i> , 2016.  |     |
|     | [29] M. Schuster and K. K. Paliwal. Bidirectional recurrent neural    |     |
|     | networks. <i>Transactions on Signal Processing</i> , 1997.            |     |
|     | [30] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson.   |     |
|     | CNN Features off-the-shelf: an astounding baseline for                |     |
|     | recognition. In <i>CVPR workshops</i> , 2014.                         |     |
|     | [31] N. Sunderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, |     |
|     | B. Upcroft, and M. Milford. Place recognition with                    |     |
|     | convnet landmarks: Viewpoint-robust, condition-robust,                |     |
|     | training-free. In <i>RSS</i> , 2015.                                  |     |
|     | [32] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed,                |     |
|     | D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich.               |     |
|     | Going deeper with convolutions. In <i>CVPR</i> , 2015.                |     |
|     | [33] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri.     |     |
|     | Learning spatiotemporal features with 3d convolutional net-           |     |
|     | works. In <i>ICCV</i> , 2015.   |     |
|     | [34] W. Yih, K. Toutanova, J. C. Platt, and C. Meek. Learning         |     |
|     | discriminative projections for text similarity measures. In           |     |
|     | <i>CoNLL</i> , 2011.  |     |
|     | [35] Q. Zhu, Q. Mao, L. Chen, M. Li, and Q. Li. Veloregistra-         |     |
|     | tion based intersection detection for autonomous driving in           |     |
|     | challenging urban scenarios. In <i>ITSC</i> , 2012.                   |     |

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# Supplementary Material: View-Invariant Intersection Recognition from Videos using Deep Network Ensembles

Anonymous WACV submission

Paper ID 432

## Abstract

This supplement document is organized as follows:

- Section 1 has additional qualitative results on Mapillary [2] dataset where we discuss performance of the proposed network in the context of success and failure cases.
- We show additional quantitative results for the five different experimental scenarios (as discussed in the main paper) in Section 2.
- Additional qualitative results are included in the video-file submitted along with this document. In addition to GTA [1] and Mapillary [2] data, we also show results on publicly available UQ St Lucia Dataset [3]. This data consists of calibrated stereo data from road traversals on a 9.5km circuit around the University of Queensland’s St Lucia campus. This dataset is generally used for testing visual SLAM algorithms.

## 1. Qualitative Results

In this section, we show additional qualitative results on Mapillary dataset. We present and discuss success as well as failure cases and attempt to understand the labyrinthine challenges associated with problem/data and how proposed model deals with them.

Figure 1 shows a correctly predicted positive video pair<sup>1</sup> from Mapillary where both trajectories pass through same intersection. The pair of trajectories are depicted using two sequences of frames. We zoom into regions from both the trajectories to show corresponding common structures. We also show an *intersection-map* which represent the *top-view* of the trajectories on the intersection. As depicted by the *Intersection-map* in Figure 1, the trajectories are perpendicular and hence minimal view overlap exists between them.

<sup>1</sup> We use the terms sample and video pair interchangeably

Figure 3 shows 5 positive samples from Mapillary that are successfully predicted by the network. Each sample depict a unique trajectory-relation and view-overlap as shown by the *intersection-map* (bottom-most row). We color encode (green, cyan, red, brown) the visually perceptible common structures (using ellipses) in the trajectories. We follow the same encoding while plotting the respective *intersection-map* in the last row. In video pair 1, 5, we observe that the overlapping part of trajectories are in the opposite direction but our network managed to exploit the common structures seen from different views. In video pair 2, though the partially overlapping trajectories are in the same direction, we can observe that *shadows* pose a challenge that is successfully handled by the network. This indicate the ability of proposed network to generalize on varying lighting conditions. Video pair 3 shows a trivial scenario where both trajectories are overlapping and network can easily predict them to same. On the other hand, video pair 4 shows the tougher scenario with perpendicular trajectories where view overlap is minimal among all the depicted cases and our network successfully predicted this to be the same intersection. Moreover, we also observe occlusion in a few trajectories (video pair 4 trajectory 2, video pair 5 trajectory 2) due to vehicles.

Figure 2 shows some wrong predictions made by the network. Video pair 1 consist of trajectories from different intersections (negative sample) but is predicted as the same intersection. We believe that the model got mislead here due to lack of unique structures in the scene as most of it has trees. Additionally, the variance in lighting conditions (as the network was trained on *day* conditions in Mapillary) can also pose challenge if not trained on other conditions.

In contrast, video pair 2, 3 are trajectories from the same-intersection but are predicted as different intersection. In video pair 2, 3 we observe a complete lack of visual features in the overlapping part of trajectories. In case of video pair 2, overlapping structure is completely occluded in the first trajectory due the presence of *truck*. In video pair 3, trajectory 1 has the camera capturing *side-view* instead of the *front-view* while also being partially occluded by the mir-



Figure 1. A correctly classified sample from Mapillary. First row show frames (arranged in the temporal order) from the two trajectories respectively where we highlight the overlapping structural content (green and cyan). Second row zooms in for overlapping structures in the sequences. Rightmost figure in the row shows trajectory of both the videos in an *intersection-map* which acts as a *top-view* of the intersection. We color encode the trajectories and overlapping content in the *intersection-map*. (Best seen in color)



Figure 2. Failure Cases: 3 wrong predictions by the network. Video pair 1 consist of trajectories from different intersections but is predicted as the same intersection. In contrast, video pair 2, 3 are trajectories from the same-intersection but are predicted as different intersection.

rror in the *car*. Hence, using purely the visual content for all these faiuar cases, even human might easily fail without some extra context (ground truth labels were generated using GPS coordinates).

## 2. Quantitative Results

We report precision (P), recall (R), F1-score (F1) and true negative rates (TNR) for all the experimental scenar-

ios in Table 1.

For Mapillary, we observe that the model performs better on all metrics when the trajectory-relations are limited to *Overlap in view and trajectory* as compared to *All combinations*. The decrease in precision (by 0.249) and recall (by 0.115) is much higher then the decline in true negative rates (by 0.040). These results are in coherence with the increasing complexity added in trajectory-relations in *All combi-*

216 Table 1. Precision (P), recall (R), F1 score (F1) and true negative rates (TNR) of our method on different datasets for the different experi- 270  
 217 mental scenarios. 271

| Dataset   | Lighting Setting | Trajectory Relation            | P     | R     | F1    | TNR   |
|-----------|------------------|--------------------------------|-------|-------|-------|-------|
| GTA       | day and night    | All combinations               | 0.49  | 0.858 | 0.623 | 0.528 |
|           |                  | Overlap in view and trajectory | 0.588 | 0.613 | 0.60  | 0.83  |
|           | day              | All combinations               | 0.69  | 0.77  | 0.727 | 0.736 |
| Mapillary | day              | Overlap in view and trajectory | 0.759 | 0.903 | 0.824 | 0.716 |
|           |                  | All combinations               | 0.51  | 0.788 | 0.619 | 0.676 |

224

225 *nations* scenario. Figure 3 shows examples with large view 278  
 226 variations and minimal overlap, indicating difficultly level 279  
 227 of the *All combinations* setup. 280

228 Similarly, for GTA, in *day and night* scenario, we observe 283  
 229 that the model performs better on precision (by 0.098) 284  
 230 and true negative rates (by 0.30) when the trajectory- 285  
 231 relations are limited to *Overlap in view and trajectory* as 286  
 232 compared to *All combinations*, whereas it performs poorer 287  
 233 on recall (by 0.245). 288

234 In the common scenario of *day and All combinations*, we 289  
 235 notice that the performance on GTA is better than Mapillary 290  
 236 on all metrics. This can be explained due to complexity 291  
 237 of data annotation in real-world where visually dis-similar 292  
 238 videos can be annotated as same based on their GPS coor- 293  
 239 dinates as shown by an example (Video Pair 2) in Figure 2. 294  
 240 Such data-complexity is slightly lower in GTA due to the 295  
 241 controlled environment settings. 296

## 243 References 297

- 244 [1] Grand Theft Auto V. [https://en.wikipedia.org/wiki/Development\\_of\\_Grand\\_Theft\\_Auto\\_V](https://en.wikipedia.org/wiki/Development_of_Grand_Theft_Auto_V). 1 299
- 245 [2] Mapillary. <https://www.mapillary.com/app>. 1 300
- 246 [3] M. Warren, D. McKinnon, H. He, and B. Upcroft. Unaided 301
- 247 stereo vision based pose estimation. In ACRA, Brisbane, 2010. 302

248 1 303

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269



Figure 3. 5 correctly classified samples from Mapillary. Each sample depicts unique trajectory-relation and view-overlap between trajectories as shown in the *intersection-map* (bottom-most row) which represents the *top-view* of the intersection. We color encode the visually perceptible common structures (cyan, green, red, brown) in the Video pairs. Video pair 2, 3 show examples where trajectory direction and view overlap. Video pair 1, 4, 5 depict cases when the overlap in trajectories and/or the view is minimal. (Best seen in color)