

Deep Learning for Robotics

Mustafa Ozdeser

Submitted to the
Institute of Graduate Studies and Research
in partial fulfillment of the requirements for the degree of

Master of Science
in
Computer Engineering

Eastern Mediterranean University
October 2020
Gazimağusa, North Cyprus

Approval of the Institute of Graduate Studies and Research

Prof. Dr. Ali Hakan Ulusoy
Director

I certify that this thesis satisfies all the requirements as a thesis for the degree of Master of Science in Computer Engineering.

Prof. Dr. Hadi Isik Aybay
Chair, Department of Computer
Engineering

We certify that we have read this thesis and that in our opinion it is fully adequate in scope and quality as a thesis for the degree of Master of Science in Computer Engineering.

Prof. Dr. Marifi Guler
Supervisor

Examining Committee

1. Prof. Dr. Rashad Aliyev

2. Prof. Dr. Marifi Guler

3. Assoc. Prof. Dr. Tolgay Karanfiller

ABSTRACT

The pathway to other computer-vision implementations was opened by sophisticated machine learning methods and simultaneous computing. In particular, the use of neural networks to control temporal information and to use the interaction of human robots for incremental learning. The world is interpreted across time, and time-indexed trajectories execute functions. The deep learning group typically ignored this valuable property. Rather, the emphasis is on developing metrics on single picture tasks or reviewing batch images. Real-time video processing got less coverage. Yet that's just what machines need. Processing single photographs does not have adequate details to track the world and process a batch of pictures.

In order to presume the last segmentation of the file, this network format requires a sequence of images that begin with the current image. We learned how to build and train these networks end-to - end. An detailed series of studies was produced on different systems and benchmarks. We found significant progress over non-recurring equivalents using RFCNN. While not restricted to robots, their influence is most evident. Mostly because robotics need to practice complex logic using minimal train details. This mixture contributes to extreme overfitting in a significantly different area during the study. Simulated results and specific output checks verify the device. We noticed that teaching the robot new things is simple for us, and later the robot would understand and use this knowledge.

Keywords : Deep Learning, Robotic

ÖZ

Gelişmiş makine öğrenimi teknikleri ve eşzamanlı hesaplama, diğer bilgisayarla görme uygulamalarının kapısını açtı. Sonuç olarak bilgisayar görüşü hızla ilerlemesine rağmen, robotik biliminin gerçek dünyadaki uygulamaları kadar etkili değildir. Bu makalede, bu konunun iki makul tetikleyicisini tartışıyor ve alternatif çözümler sunuyoruz. Spesifik olarak, zamansal bilgiyi yönetmek için sinir ağlarını kullanmak ve kademeli öğrenme için insan robot etkileşimini kullanmak gerekir. Dünya zaman içinde yorumlanır ve zaman indeksli yörüngeler işlevleri yerine getirir. Derin öğrenme grubu tipik olarak bu değerli yapıyı görmezden geldi. Bunun yerine, tek resimli görevler için metrikler geliştirmek veya toplu görüntüleri gözden geçirmek vurgulanmaktadır. Gerçek zamanlı video işleme daha az kapsam kazandı. Yine de makinelerin ihtiyacı olan şey bu. Tek bir fotoğrafın işlenmesi, dünyayı izlemek ve bir grup fotoğrafı işlemek için yeterli ayrıntıya sahip değildir.

Gerçek zamanlı strateji ve karar verme ile ertelendi. Bu sorunu çözmek için, çeşitli robotik senaryolarda çok yararlı olan, segmentasyon için tekrar eden tamamen evrişimli bir sinir ağı (RFCNN) öneriyoruz. Bu ağ biçimi, son görüntü segmentasyonunu varsaymak için mevcut resimden başlayarak bir dizi görüntüyü kapsar. Bu ağları uçtan uca nasıl inşa edeceğimizi ve eğiteceğimizi öğrendik. Farklı sistemler ve kıyaslamalar üzerine detaylı bir dizi çalışma üretildi. RFCNN kullanarak tekrar etmeyen eşdeğerlere göre önemli ilerleme bulduk. Derin öğrenme yaklaşımları, erişilebilir en popüler makine öğrenimi çözümleri olsa da, eğitim ve test arasındaki veri dağıtımında yaşanan değişimden muzdariptir. Robotlarla sınırlı olmamakla birlikte, etkileri en belirgindir. Bu karışım, çalışma sırasında önemli ölçüde farklı bir

alanda aşırı uyuma katkıda bulunur. Bu sorunu hafifletmek için, robotun İnsan-Robot Etkileşimi (HRI) aracılığıyla yeni algı bilgileri hakkında düşünebileceği yeni bir model öneriyoruz. Sesi kullanarak insan dostu iletişim için eksiksiz bir HRI programı ve Geste kullanılmaktadır. İnsan geribildirimini kullanarak, bir nesne algılama ağını geliştirmek için aşamalı bir öğrenme yaklaşımı oluşturulur. Simüle sonuçlar ve gerçek performans testleri cihazı kontrol eder. İnsanların robota kolayca yeni şeyler öğretebileceğini ve daha sonra robotun bu bilgiyi bilip kullanacağını gösterdik.

Anahtar Kelimeler: Derin Öğrenme, Robotik

ACKNOWLEDGEMENT

I would like to thank my supervisor Prof. Dr. Marifi Guler. Through his constructive suggestions, he really helped me out. Our meetings were positive, but it also meant that I had critical thought and a new viewpoint on a topic. Finally, I want to thank my parents, whenever I need them, for being there for me. Without them, this feat would not have been practicable.

TABLE OF CONTENTS

ABSTRACT.....	iii
ÖZ	iv
ACKNOWLEDGEMENT	vi
LIST OF TABLES	ix
LIST OF FIGURES	x
1 INTRODUCTION	1
2 LITERATURE REVIEW AND BACKGROUND	4
2.1 The Learning of Deep.	4
2.1.1 The Networks of Recurrent.....	7
2.1.2 Recognition Networks and Object Detection	9
2.2 The Interaction of Human Robot.....	10
3 THE VIDEO SEGMENTATION OF RECURRENT FULLY CONVOLUTION NETWORKS	15
3.1 The Methodology	17
3.1.1 The Segmentation of the Recurrent Unit of Conventional.....	17
3.1.2 The Segmentation of Convolutional Gated Recurrent Unit (Conv-GRU)	18
3.2 Results.....	19
3.2.1 Databases.....	20
3.2.2 The Recurrent Fully Convolutional Network of Results.....	23
3.2.3 Further Analysis.....	26
3.3 Discussion.....	26
3.4 Results of Uncertainty Estimation.....	28

4 THE HUMAN ROBOT INTERACTION IN INCREMENTAL LEARNING	37
4.1 Humans for Incremental Learning	40
4.1.1 The Robotics of Deep Learning.....	40
4.1.2 The Open Set Recognition Multi Class.....	42
4.2 The Structure.....	43
4.3 Investigations	45
4.3.1 The HRI of Incremental Learning.....	46
4.3.2 Recognition Baseline and Object Detection.....	48
4.3.3 The Approach of Incremental Learning	49
4.3.4 The Performance Evaluation of Mock Human-Human Interaction.....	50
4.4 Analysis.....	53
5 CONCLUSION AND FUTURE WORK.....	54
REFERENCES.....	56

LIST OF TABLES

Table 1: Description of the planned network. $F(n)$ shows the thickness of the filter $n \times n$. $P(n)$ Denotes the complete map padding of n zero function.	22
Table 2: Tiramisu MC dropout.....	28
Table 3: Rankinig IOU of variational ratio	29
Table 4: Tiramisu TA-MC accuracy.....	30
Table 5: Tiramisu TA-MC ranking IOU of variational ratio.....	30
Table 6: Tiramisu RTA-MC performance.....	31
Table 7: Tiramisu RTA-MC ranking IOU of variational ratio.....	31

LIST OF FIGURES

Figure 1: Modern classifier deep network.....	5
Figure 2: A completely linked classification network (top) vs a fully integrated segmentation network (bottom). The main change being the full absence of the FCN base.....	6
Figure 3: The GRU architecture.....	9
Figure 4: The pizza maker's proposed human-robot interaction.....	12
Figure 5: Types of interactive robots.....	13
Figure 6: Overview of the proposed recurring method FCN recurrent part for better viewing.....	16
Figure 7: The RFC-VGG architecture.....	18
Figure 8: RFCN and FCN on SegtrackV2.....	20
Figure 9: Qualitative tests for SegtrackV2 and Davis with FC-VGG overlay top picture and RFC-VGG lower segmentation.....	21
Figure 10: Qualitative checks for Synthia experiments in which data is superimposed on the network forecasts.....	24
Figure 11: Qualitative cityscapes tests contain studies that exceed the network feedback forecast. The strange lines are FCN-8s and also their accompanying numbers: the production of RFCN-8s.....	25
Figure 12: PR curve of mean IOU	29
Figure 13: Tiramisu TA-MC PR-curve of mean IO.....	30
Figure 14: Tiramisu RTA-MC PR curve of mean IOU.....	31
Figure 15: Precision-Recall curves of pixel level of Segnet backbone.....	35

Figure 16: Precision-Recall curves of pixel level of Tiramisu backbone.....	36
Figure 17: The HRI of incrementing robot knowledge.....	38
Figure 18: System block diagram.....	43
Figure 19: Visualization of RGB. Objects are detected and located in the scene.	45
Figure 20: Incremental learning pattern.....	46
Figure 21: Multimeter	47
Figure 22: Recognition success in an gradual learning environment in the first step. The robot takes pictures and adds new objects one by one. The displayed values are the highest accuracy.	48
Figure 23: Incremental learning scenario recognition performance with the second approach. The user introduces new objects one by one and the robot collects their images. The values shown are the highest accuracies.....	49
Figure 24: Recognition performance after the first approach is gradually added to new classes. Data are taken from imagenet for new classes. The values shown are the highest accuracies.	50
Figure 25: Identification of success by slowly introducing different second method levels. Data was taken from imagenet for new classes. The values displayed are the most correct. Box plots show the difference in precision, as the relation between new and old class study samples ranges between 0.05 and 0.5. The green line is the accuracy of the 0.1 ratio.....	51
Figure 26: Evaluation of NASA task load index for 4 interfaces tested.	52

Chapter 1

INTRODUCTION

It was still is the aim of the community to move robotics increasingly and benefit them for multiple tasks in place of specific applications. This was not only a conventional approach to engineering but a robot was built and prepared for a very specific mission. Machine learning is used to overcome this problem, because it can carry out a wider range of positions. This sample is gaining reputation, particularly after successful deep learning applications for image analysis.

The essential desire behind deep robotic learning is that they are much more general than any other algorithm for learning. Deep networks have been shown to be capable of high-level thinking and abstraction. Because of this, in unstructured environment, there occurs an idea for robotics. In addition, there are rather efficiency in aspect of advanced numerical libraries and parallel processing and networks. High frequency response module is required for time-critical robotics tasks to control motion. This can be delivered by deep networks on GPUs. These say that the practical use of deep image learning still differs significantly from its use in robotics.

The application of robotics is dependent on time by essence in terms of time information. Robots are active in dynamic environments, time is indexed and the function of time is what they perceive. For instance, a robot tennis player must monitor the ball over time and create a time-indexed path for a successful hit based on its speed

and angle. Note that data can not be collected from a single picture for such an activity, since speed can not be observed. Most computer vision groups don't address these situations that implies that most deep architectures for single image processing are expected.

Training Data Deep networks are highly capable of mapping functions and can learn from input to output any non-linear function. But that comes at a cost. The training information needed to learn a general function for a job through a deep network is proportionate to the difficulty of the problem. For instance, the differentiation of black and white may necessitate only a handful of darkened samples, but thousands are needed to differentiate cats from dogs. Considering the network accuracy can be affected by increasing numbers of data samples with a different appearance. It's an severe case of robotics. Analyzing a household robot assistant and in case that we have gathered sufficient data in the laboratory in order to train them for the fundamental task of handling objects. Then the robot has been sent to a person's home and unfortunate to us, house items look very different from particles in the laboratory. The existing solution to this issue is to boost training samples to include similar examples of possible conducting tests in the computer display community.

First, details from the context and associated works are provided in Chapter 2 to explain the others. Our potential approach for the use of time data for video processing is demonstrated in chapter 3. We define a complete convolutional recurrent network of video segmentation, in particular. We then implement in Chapter 4 an evolutionary method of learning, through which robots explore persons by natural interactions.

Subventions are documented as a new approach for video segmentation centered on repeated networks that are appropriate for recurrent online robotics. There are systematic studies that demonstrate that the method is superior to a single image segmentation by offering an gradual learning paradigm to improve a deep network by normal human experience and implementation , training and proof of comprehension.

Chapter 2

LITERATURE REVIEW AND BACKGROUND

2.1 The Learning of Deep

Over the last decade, deep learning (DL) has had a big influence on data science. This chapter presents the fundamental concepts in this field. It contains both the basic architectures used for designing deep neural networks and a brief overview of some common cases. Neural networks have revolutionized the daily lives of today. Their important impact is also present in the most basic actions such as ordering product online via Amazon's Alexa or spending time online video games with computer agents. The 21st century began with some advances in the field of language speech and processing in neural networks. For example, NNs in imagery are used in lesion detection and segmentation, and with this technology tasks such as text to speech and text to image have improved remarkably. It has a powerful influence and continues to grow. The NN journey started in the mid 1960s when the Perceptron was published. Its development was driven by human neuron activity formulation and human visual perception research. But there was a very rapid deceleration in the field, which lasted nearly three decades.

While some other important developments in the next decade have taken place, such as the development of the long-short memory machine (LSTM), the field has been further deteriorated. There were questions without adequate answers particularly regarding the non-convex nature of the optimisation goals used, the overfitting of

training data and the challenge of disappearing gradients. These difficulties led to two decades of NN. Meanwhile, classic machine learning techniques were developed and attracted a lot of attention from academia and industry. One of the main algorithms was the newly proposed Support Vector Machine (SVM), which had a clear mathematical interpretation for a convex optimization problem. These features enhance its popularity and use in different applications [1].

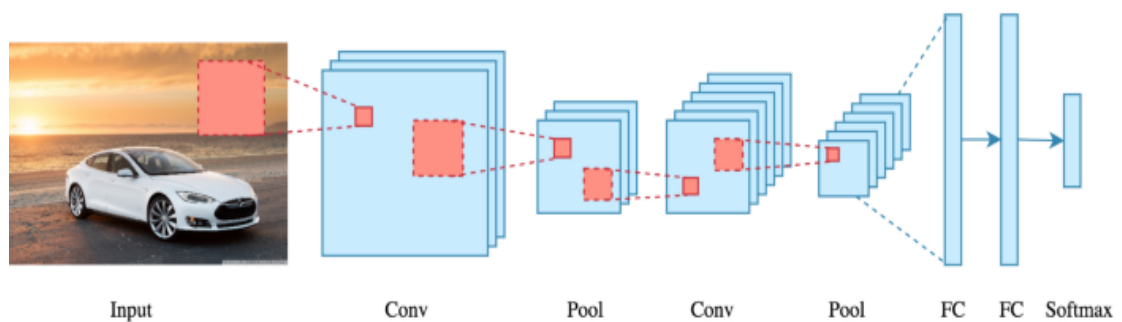


Figure 1: Modern classifier deep network

CNN offers a highly effective way of processing images that produce its reputation with computer vision and machine learning. Utilizing CNN together through fully connected layers, the modern neural network classification architecture was produced.

The Segmentation of Fully Convolutional Networks

For the classification component in the CNN used for classification, the least completely connected classification layers are relevant. However, detailed predictions are possible for all pixels with pixels marking. In [2] the concept of using a completely convolutional neural network educated in the segmentation of semantic pixels is introduced. The FCN architecture [3] is focused on VGG due to its performance in classifying tasks. The completely linked layers of these networks require only fixed size inputs and classification labels to be produced. To solve this issue, a completely

linked layer can be transformed to a convolutional layer. Convolution filters should be used for all spatial input scales, irrespective of the sample dimension. In order for this gross map to be very dense, the initial image size must be sampled. The up-sample process can be conveniently translated between the line section. It was introduced in a new layer with a network upsample. It makes it possible to know the weight by back propagation across the network. The filters for the deconvolution layer are the foundation for the input image reconstruction. Another suggestion is to compile performance maps from the changed representation of the data. However, there has been talk of the usage of upsampling with deconvolution [2].

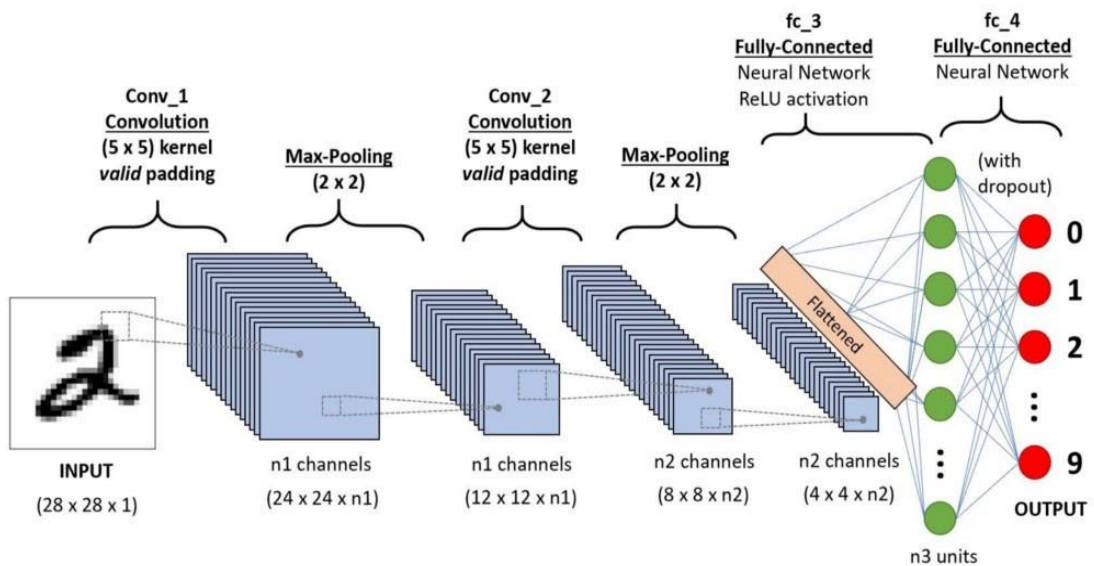


Figure 2: A completely linked classification network (top) vs a fully integrated segmentation network (bottom). The main change being the full absence of the FCN base

Different applications have tried the FCN architecture. In [4] it is used for the location of objects. A individual FCN network has been used in order to predict boundary box sites from an input pyramid. The network has been shown to be able be trained and perform better for multiple tasks either. A modified architecture for visual object

tracking was used in [5]. Characteristic maps from various layers crossed and joined for improved monitoring, two independent FCN divisions. Finally, a complete deconvolution network with stacked deconvolution layers is presented for semantic segmentation in [6]. Multiple deconvolution layers showed positive effects on segmentation accuracy.

2.1.1 The Networks of Recurrent

The RNN [76] are structured to combine neural network architectures with sequential knowledge. By using a hidden unit in each repeating cell, these networks are capable of studying complex dynamics. This unit functions like a dynamic memory that, depending on the state in which the unit is, can be changed. As outlined below, you can model the simplest recurrent device.

$$h_t = \theta\varphi(h_{t-1}) + \theta_x xt \quad (1)$$

$$y_t = \theta_y\varphi(h_t) \quad (2)$$

Recurrent networks have been active in many speech recognition and speech synthesis functions, comprehension of text[69] but their difficulties come with them. The uncontrolled flow of data between units creates problems with gradients disappearing and exploding[7]. The derivative of each node is dependent on all of the preceding nodes during back propagation by recurrent units.

$$\frac{\partial E}{\partial \theta} = \sum_{t=1}^{t=S} \frac{\partial E_t}{\partial \theta} \quad (3)$$

$$\frac{\partial E_t}{\partial \theta} = \sum_{k=1}^{k=t} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial h_\theta} \quad (4)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^t \theta^T \text{diag}[\varphi'(h_{i-1})] \quad (5)$$

Using gated buildings is a response to this problem. Between each node, the gates will regulate back propagation flow.

Long Short Term Memory (LSTM)

There are three gates in each LSTM node, each with learnable weights, which are input, output, and forget gates. The perfect way to recall valuable knowledge from previous states and determine the new state can be discovered through these walls. The element-wise result is the operator.

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad (6)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \quad (7)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \quad (8)$$

$$g_t = \sigma(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (9)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (10)$$

$$h_t = o_t \odot \varphi(c_t) . \quad (11)$$

Gated Recurrent Unit (GRU)

Similar to LSTM, the Gated Recurrent Device utilizes a gated flow system-Controlling. It has a simpler architecture, however, which makes the memory use much quicker and less usage.

$$z_t = \sigma(W_{hz}x_{t-1} + W_{xz}x_t + b_z) \quad (12)$$

$$r_t = \sigma(W_{hr}h_{t-1} + W_{xr}x_t + b_r) \quad (13)$$

$$h_t = \Phi(W_h(r_t \odot h_{t-1}) + W_x x_t + b) \quad (14)$$

$$h_t = (1 - z_t) \odot h_{t-1} + h_{t-1} + z \odot h_t. \quad (15)$$

GRU has no direct control over the sensitivity to memory information, whereas By having an output bolt, LSTM got it. In the way that the memory nodes are modified, these two are both distinct. Using summation after input gate over flow and ignoring gate, LSTM updates its secret state.

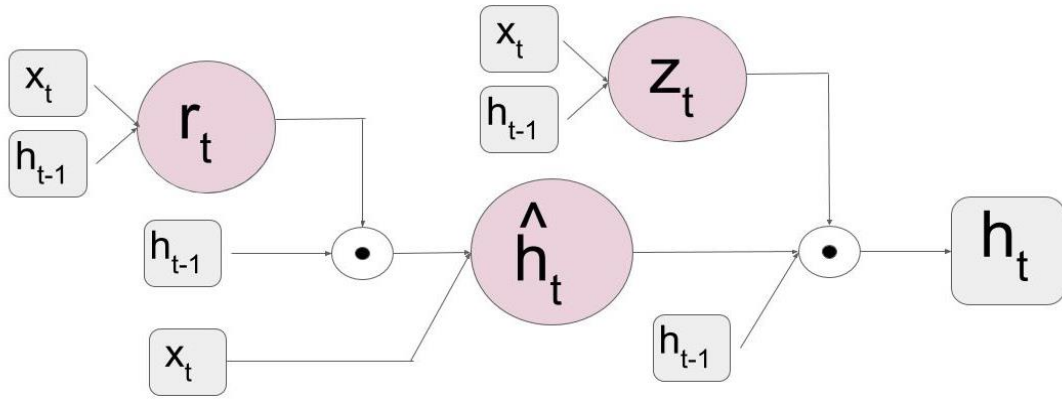


Figure 3: The GRU architecture

2.1.2 Recognition Networks and Object Detection

In a disordered world, object recognition tries to decide the right match boxes. For still photographs, object identification has historically been done by means of an background gradient and strong limits were identified between the target and the context or solid colored patching. The usage of images would track moving items with temporary filters even better. More recently, the apps SIFT[7] and HOG[8] have been used extensively. Most of the researchers in the area have acknowledged this time primarily by combining various low-level ensemble models with high-level filters. The area strategy was always a general concept issue for various detection techniques. This relies on the classification of superpixels or sliding glass. However, their precision relies on the sample fineness, which directly impacts measurement costs. The selection of groups to be observed for professionals is usually limited to extracting very large or small items.

The next move was the development of the first Convolutional Neural Networks [9]. Profound networks have shown that improved object classification capabilities can be derived than conventional HOG and SIFT. The question was how to define the utility of deep classification networks. Multiple work was conducted simultaneously to solve

this issue. At the time, they were state-of-the-art, but no pain. First, as mentioned, sliding window caused calculation difficulties. Second, for pooling layers, using convolutional networks with limited input size constraints. Consequently, either low spatial resolution or high network size. These conditions will reduce precision.

Girshik et al. [13] planned a network object proposal technique to overcome problems with a sliding window approach. This method does not rely on the sliding window method for external proposals. It develops the ideas internally and assess them via a deep network. The regional offer could therefore be accomplished and enhanced. Also, it dramatically reduced the inference time. Item identification mechanisms are typically assisted by a classification mechanism. The observed region is easily but flawed to incorporate into the initial picture and into the classification network. It is faulty because we recalculate very analogous features to the detection network that has already been extracted. It would be better to use intermediate detection network features and the classifier above to note the object. In addition, detection and identification network programming will work together by improved overfitting.

2.2 The Interaction of Human Robot

Autonomous robots became an integral component of major companies' production chain. Nonetheless, due to their mission limitations, their use is less feasible in smaller companies or households. Although these scenarios do not yet have robots with full autonomy, there is a great opportunity for other solutions. The design of an HRI system has a number of factors. In human proximity, the HRI system should be modified based on the device and the user's physical location. If the robot and the customer are in the same place , people can protect the environment by themselves. And if the robot is positioned in a distant location, the robot only sends customer

information. In Robot 's architecture, robot capability restricts the HRI system. For example, robot arms suited for physical combat but have a very limited range. On the other side, UAVs with extensive range are available, but absolutely no physical contact is necessary. The HRI communication system medium and form can be designed to support one or more types of communication devices. It can also define its own communication protocols. The specification explicitly influences the option of the HRI program. For instance, the design of an arm robot interaction software operation and another space operating arm robot is very disparate. The architecture is very identical, however the form of applications built to dramatically alter the design. Latency and acceleration will be recognized for the space robot. Exact haptic feedback and exact working space limits are important for a surgical robot. The predicted autonomy of the robot will differ in autonomy HRI systems. In one step, fully autonomous robots work in close proximity to humans. A good example of this is the roomba vacuum robot systems in which it operates on the basis of sensory inputs only. There are teleoperated robots at the other end of the spectrum which are completely controlled by the human user. However, it is much easier to mix flexibility and collaboration than these two extremes. Semi-autonomy has the benefits all mechanisms provide to solve their inconveniences. In this model, the human being does not regulate the actuator entirely and encourages the robot to assume responsibility for a laborious regulate portion. The robot may use human instruction concurrently to avoid the robot 's complicated high level thinking. However a reasonable balance of control and contact may be much greater than all extremes. Semi-autonomy ensures that all mechanisms can transcend their drawbacks. Human beings in this model do not directly monitor any actuator and require the robot to take care of the laborious monitor portion. In parallel, the robot may use human guidance which bypasses high-level robot thinking. One excellent

example is [14]. Knowing the idea of pizza making is incredibly complicated for robotics or how they are mixed. The physical aspect of the task is at the same time too repeated and time consuming for humans.

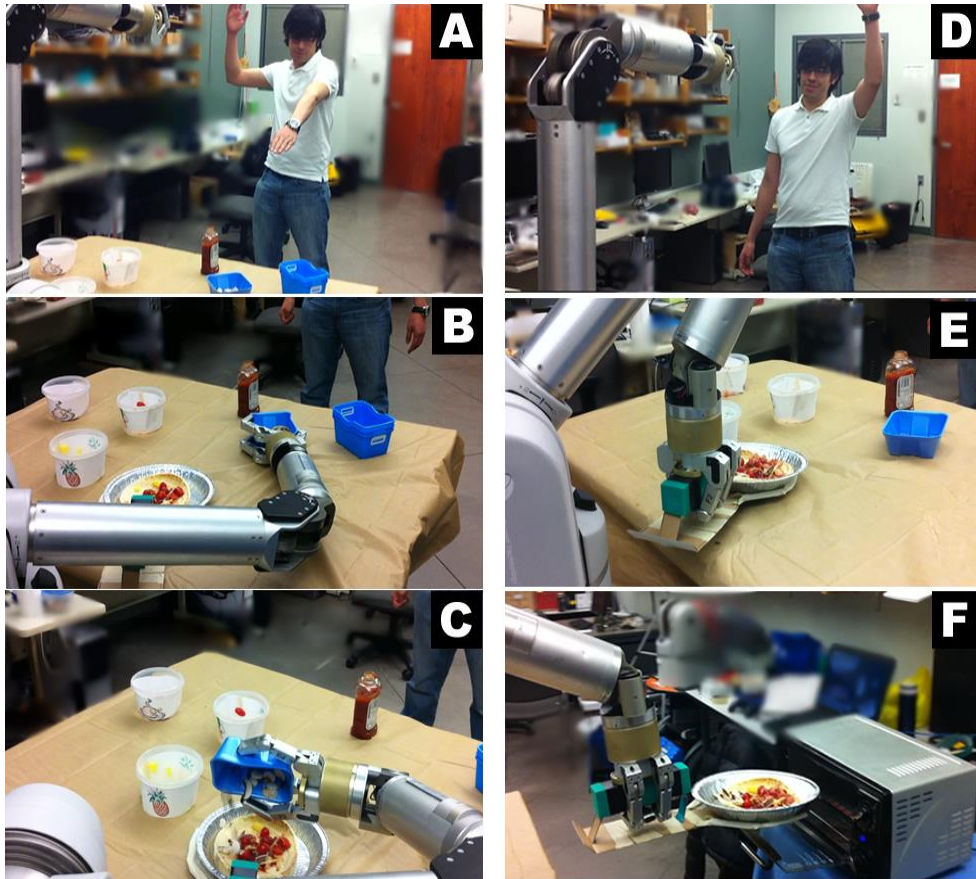


Figure 4: The pizza maker's proposed human-robot interaction [14]

HRI 's past focus was on finding more practical ways to employ human robots [15]. A change in awareness has drawn attention to places where human interaction strategies find the robot's regular communication more humanly. Human [16], voice, body language , facial expressions and physical activity are main contact types. The need stresses the robot's reasoning module because of its intelligible human influence. Such computers are also called virtual robots. Below, I'll analyze prominent plays on social HRI awareness issues. Yet let's look at three key aspects of the social HRI system first.

A social HRI robot's most important feature can be vision. It turns raw signals from environmental sensors into something that makes robot sense. Any definitions include navigation, image detection / tracking, identification of voices, etc. In the mean time, the appearance of a world robot is encapsulated in this section. You may name it the cortex as well. It is liable for utilizing the knowledge on interpretation in decision-making. For eg, a nurse robot will have sufficient medications for the patient. The vision gives a cue to recognise the object. In the intermediate, patients should identify by their expressions and prescribe the correct drug for both of them. The robot has now seen the universe and determined and now needs to take steps or transmit the findings to humans. Thus through different media including voice synthesis, motion, multimedia (text, images), picture, etc.

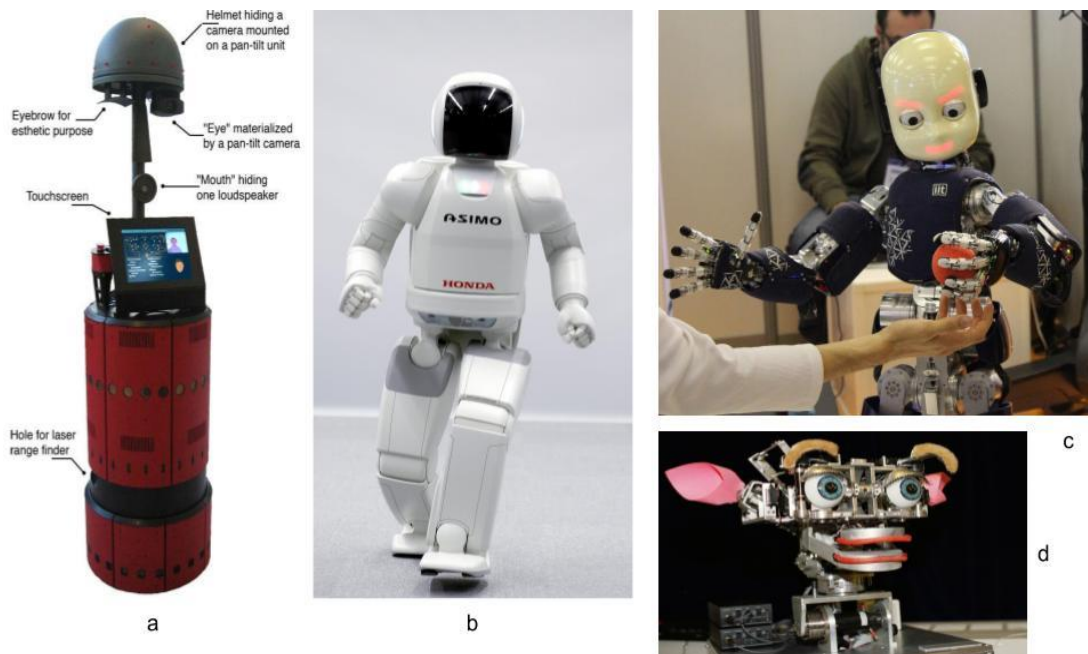


Figure 5: Types of interactive robots

MIT's Kismet [17] and Honda's Asimov [18] are two of the early contributions. The facial characteristics of Kismet and may interact to some degree. It was used extensively in HRI research. Another example of social robots that are primarily

planned for potential growth is ICub [19]. This can identify and recognise and monitor images. Personal machines have often been used as tour guides or leaders. The Swiss National Science Museum tour guides [20] and Rackham [21] were tour guides for the BioSpace Show. They interacted with the visitors and led them to their ultimate destination. We provide speech recognition, emulation and modules for navigation. Face recognition, motion tracking and interpretation of movements are other capabilities.

Chapter 3

THE VIDEO SEGMENTATION OF RECURRENT FULLY CONVOLUTIONAL NETWORKS

Recent research into deep neural networks has greatly enhanced the interpretation of the system. This phenomenon was originally exacerbated by naming objects [9][3][22]. Semantic segmentation is a more complex task as implemented in [23][25] and provides pixel identification. A full convergence network was launched in [23]. Such networks have a basic diagram for each picture and analyze complex network forecasts. This method permitted end-to - end seminal segmentation training. Nevertheless, one aspect of this new phenomenon is that the entire world is not a collection of static pictures. The program creates a great deal of environmental consciousness. The current CNN networks are not easily modified. The best way to use time information in CNN is to patch and distribute several frames as a single source. Minimal variations of this approach are used to describe one million Youtube videos in context[26]. In [28] a convolutional Boltzman restricted system was introduced, which acquires properties such as optical flows from the image sequence.

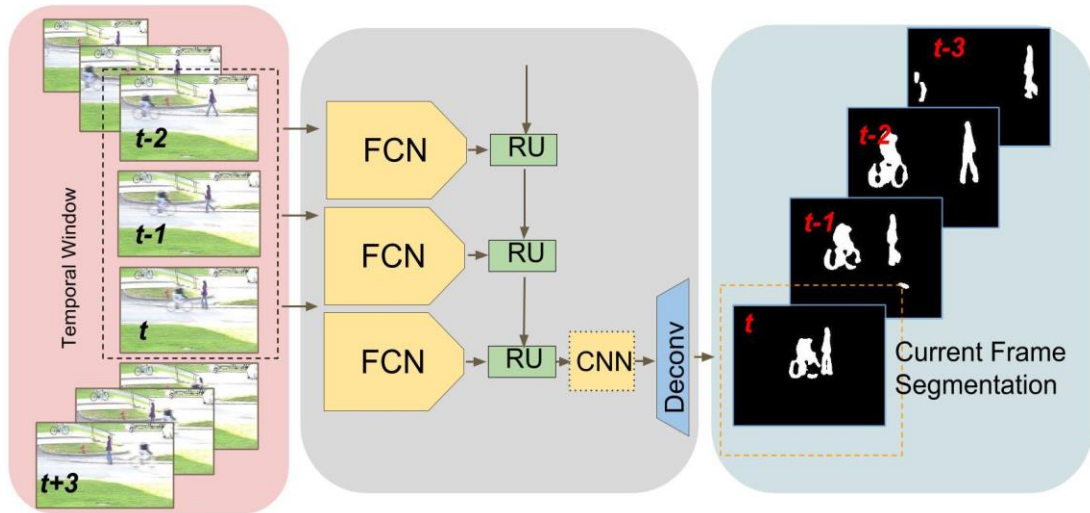


Figure 6: Overview of the proposed recurring method FCN recurrent part for better viewing

Several architectures are proposed for the resolution of the key limitation of recurrent networks, namely the vanishing gradients. Another recently introduced design is the Gated Recurrent Unit (GRU) [34]. LSTM and GRU have been seen to outweigh other practices in [35]. The repeating structures and GRU efficiency were identical to LSTM, but the amount of parameters was the one issue in the previous archives is that tectures only function as vector sequence data. They can not manage data where spatial knowledge, such as photographs or charts, is important. Another work [36] uses GRU to fix spatio-temporal video functionality. Experiments were performed on video subtitling and understanding of human behavior. Using FCNs along with repeating managed systems will overcome much of the inconveniences of prior strategies. Our layout is focused on the recurrent neural network because the learning of temporal dynamics has proven to be efficient an end-to - end video segmentation training system not required for the offline processing of data. This is the first research in our knowledge to pose a persistent, completely convolutional pixel mark network.

3.1 The Methodology

Abstractly, we use both the information of time and space segmentation from a persistent totally convergent network (RFCN). In general, the architecture principle consists of the use of repeated nodes integrating totally convolutional Current unit(RU) operations. The repeating category literally corresponds to LSTM, GRU or Conv-GRU. In comparison to the batch / offline edition that requires the whole video as content, we strive for on-line segmentation through all our networks. The frames are stuck over a glass. Every window is then stretched across the RFCN and gives the last moving window frame a segmentation.

The FCN network's forward propagation. The entire network was eventually taught and the lack of pixel awareness was logarithmic. We have also developed numerous applications around network architectures to use conventional and convolutional recurring unit.

3.1.1 The Segmentation of the Recurrent Unit of Conventional

The Lenet network that was transformed into a truly translated network is our first architecture. Lenet is a well developed, shallow, and popular network for early experiments. This architecture is defined by RFC-Lenet. A 2D map with complex projections is the devolution performance of the FCN and is flattened into a 1D vector as an entry into the repeating array. The repeating device extracts a vector in the sliding window from each frame and leaves the last section of the frame.

Note, RU layer must be used for initial large matrix architecture as it operates on compressed full-size file vectors. We can attach a deconvolution layer after the recurrent node to this problem. This leads us to our second design. The RU gives the

input and output of the rough last frame projections a flattened vector coarse diagram. This gross chart is then translated to detailed forecasts. This is helpful for wider photos that suit certain RU parameters. Reduced size, allows a smaller state and a common local minimum for scanning optimizers in a shorter training time. An example of this technique are the RFC-12s. This is the slightly modified Lenet FCN network version. The only adjustment now is that the replay takes place at the conclusion of the last move.

3.1.2 The Segmentation of Convolutional Gated Recurrent Unit (Conv-GRU)

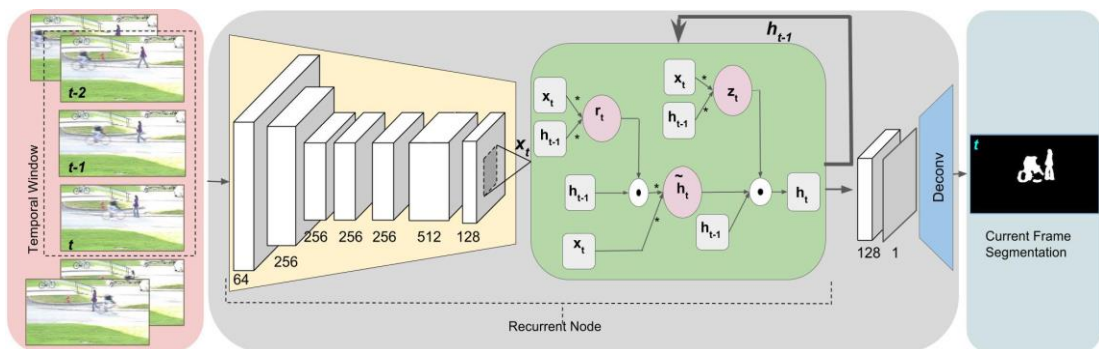


Figure 7: The RFC-VGG architecture

Assume that a repeated unit is placed according to a spatial dimension in a situation where $h \times w$ and has a channel c . Once flattened, it will become a long c -matrix. w . There will previously be constant unit weights $c \times (h \cdot w)^2$ That is spatial area power four. Such matrices can only be preserved for the smallest characteristic maps. Even if there is no computing problem, such a design also introduces a large network variance which prevents generalization. Weights in groundbreaking units are three-dimensional, close to the ordinary convolutional layer and correspond to the input instead of the dot product. In the design, weight patterns are sized $k_h \times k_w \times c \times c \times f$ where k_h are the height of the kernel, kernel width, input numbers and filter numbers, respectively and we should conclude that the maps' spatial bandwidth in kernel scale can be very limited in contrast

with the map's space scale. This method is much more efficient and is simpler to learn weights because of the smaller search regions. For a vertically sequential segment by sector network this approach is used. This layer can be seen on heat maps or charts. This development is directed to the deconvolution layer in the first example and the likelihood map is produced in pixels. In this case, after the recurrent layer a CNN layer would be used to turn its output characteristics into a heat map. The second scenario is summarized of the RFC-VGG. The emphasis is on the VGG-F network [37]. Since VGG-F weight affects the weight of our screens, overcrowding issues are reduced thanks to the detailed imaging. Built into a fully convergent network, the specifically connected layers are substituted by convolutional layers. The last two layers of pooling are that such that VGG-F is properly segmented. A convolution is then used to evaluate replicated units followed by a convolutional chain, along with a breakdown.

$$z_t = \sigma(W_{hz} * h_{t-1} + W_{xz} * x_t + b_z) \quad (16)$$

$$t_t = \sigma(W_{hr} * h_{t-1} + W_{xr} * x_t + b_r) \quad (17)$$

$$h_t = (1 - z_t)xh_{t-1} + zx\hat{h}_t \quad (18)$$

3.2 Results

Each section summarizes our research findings. We define, first, the datasets we used, our training methods, and the hyperparameter sense. Finally, there are objective and qualitative tests. The open source code for RFCNN software can not be released. On top of Theano, we created our own library for arbitrary output [38]. Networks of FCN as a recurrence server. (1) Promotes dynamic representation networks. The main characteristics of this program are: any single CNN and any random number of replicated layers can be the architecture. Both data lengths are authorized for networks. (2) In the recurrent framework are three gated architectures, LSTM, GRU and Conv-GRU. (3) Deconvolution layer and Segmentation FCN interface support.

3.2.1 Databases

There are four datasets used in this article: 1) Moving MNIST. 2) Transfer Recognition [39]. 3) [40]. Segtrack 2. 4) Video segmentation, densely annotated (Davis) [41].

By dynamically changing the characters from the original MNIST, the Moving MNIST Data Set is synthesized. After translation, the segmentation labels are generated using input threshold images. A new structure is considered to be the picture. We may, thus, have an arbitrary image set.

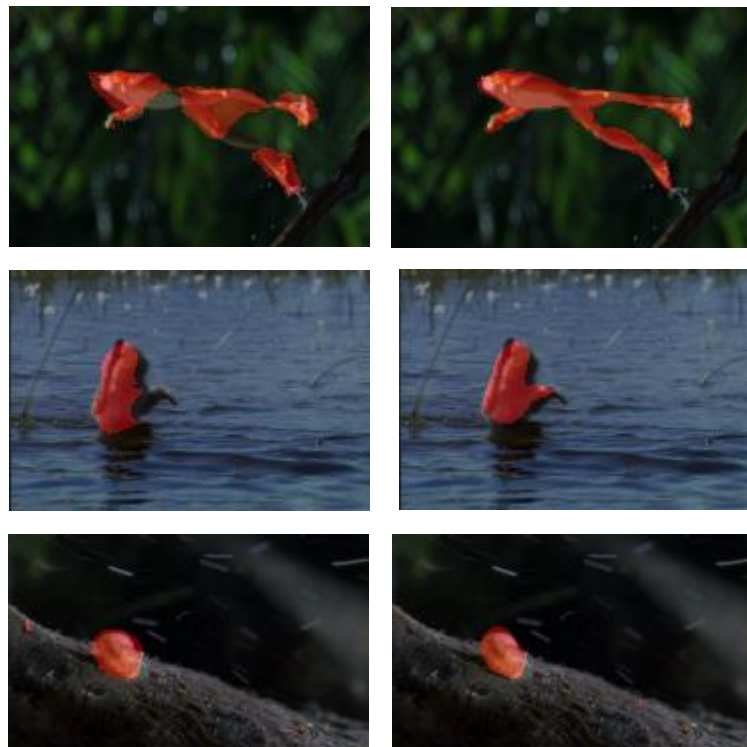


Figure 8: RFCN and FCN SegtrackV2

Shift Detection Dataset [39] A realistic, dynamic range of videos with pixel marking of moving objects is given in this data collection. Indoor and outdoor scenes are also included in the dataset. This depends on segmentation for moving objects. We searched for video clips of identical moving objects, for example vehicles or individuals, to semanticize sequences. Therefore, six videos were chosen:

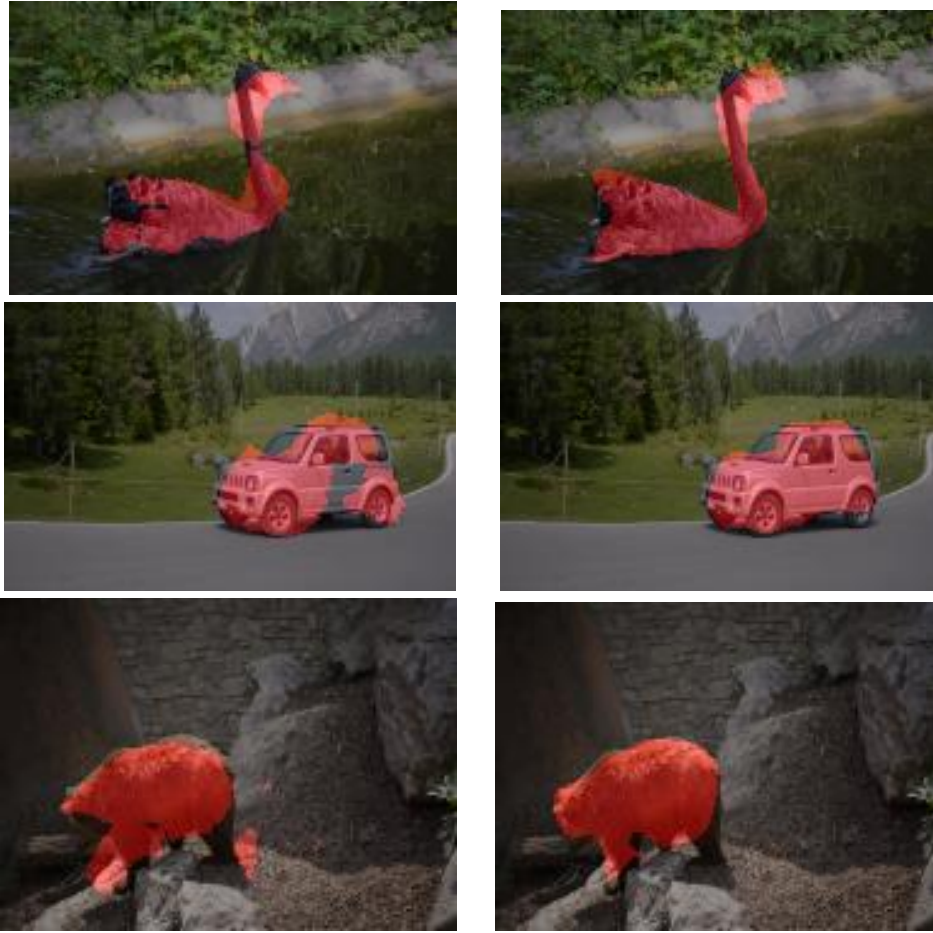


Figure 9: Qualitative tests for SegtrackV2 and Davis with FC-VGG overlay top picture and RFC-VGG lower segmentation.

Table 1: Description of the planned network. F(n) shows the thickness of the filter $n \times n$. P(n) Denotes the complete map padding of n zero function.

Network Architectures					
RFC-Lenet		RFC-12s		RFC-VGG	
input: 28×28		input: 120×180		input: 240×360	
Recurrent Node	Conv: F(5), P(10), D(20)	Recurrent Node	Conv: F(5), S(3), P(10), D(20)	Recurrent Node	Conv: F(11), S(4), P(40), D(64)
	Relu		Relu		Relu
	Pool 2×2		Pool 2×2		Pool 3×3
	Conv: F(5), D(50)		Conv: F(5), D(50)		Conv: F(5), P(2) D(256)
	Relu		Relu		Relu
	Pool(2×2)		Pool(2×2)		Pool(3×3)
	Conv: F(3), D(500)		Conv: F(3), D(500)		Conv: F(3), P(1) D(256)
	Relu		Relu		Relu
	Conv: F(1), D(1)		Conv: F(1), D(1)		Conv: F(3), P(1) D(256)
	-		-		Relu
	-		-		Conv: F(3), P(1) D(256)
	-		-		Relu
DeConv: F(10), S(4)	Flatten	Conv: F(3), D(512)			
Flatten	GRU: W(100×100)	Conv: F(3), D(128)			
GRU: W(784×784)	DeConv: F(10), S(4)	ConvGRU: F(3), D(128)			
		Conv: F(1), D(1)			
		DeConv: F(20), S(8)			

Davis [41] dataset contains 50 high resolution and densely annotated pixel-exact videos. The videos contain several challenges, such as occlusions, rapid movement, nonlinear deformation and motion blur.

Synthia [42] is an urban machine textual segmentation dataset. This contains 13 level pixel rating annotations. Since only half of the data set is required for our road series experiments in the summer.

CityScapes [43] is an actual dataset that focuses on urban scenes captured during driving videos in different cities. There are 5,000 beautifully annotated 20,000 30-class images.

3.2.2 The Recurrent Fully Convolutional Network of Result

The source for this series of studies known as FC-VGG is a fully convolutional VGG. This is compared with the regular RFC-VGG version. To stop overcrowding of VGG tests, the initial five convolution layers are not finalized and pretrained. For these experiments, the data is split into two sections for each sequence, half as workout and half as test outcomes. Statistics reveal that RFC-VGG is 3 to 5 percent greater than the DAVIS to SegTrack architectural data sets. The qualitative RFC-VGG analysis against FC-VGG.

$$p = tp / fp + tp, r = true p / true p + false n \quad (16)$$

$$- \text{measure} + F = 2 * p * r / p + r \quad (17)$$

$$IoU = true p / true p + false p + false n \quad (18)$$

The values are listed as FC-VGG and RFC-VGG categorized as SegTrack V2 and Davis. The SegTrack V2 listed respectively as Precision, Recall, Fmeasure and IoU as 0.7759, 0.6810, 0.7254 and 0.7646. In RFC-VGG, it listed as 0.8325, 0.7280, 0.7767 and 0.8012. In Davis, FC-VGG listed as 0.6834, 0.5454, 0.6066 and 0.6836. In RFC-VGG, the values listed as 0.7233, 0.5586, 0.6304 and 0.6984.

It shows that the use of time data in a regular unit increases the segmentation of objects. It can be explained as the motion of segmented objects in repeated systems has been implicitly recognized. This can also be used for gathering time data from the maps as the number of parameters is minimized by the repetitive conv-GRU unit. The recurrent unit may then establish a movement pattern for segmented systems by building on

detailed information from these maps. The repeating version can then be built using the skip architecture with the aid of an upgraded, fully configured network for more efficient segmentation.

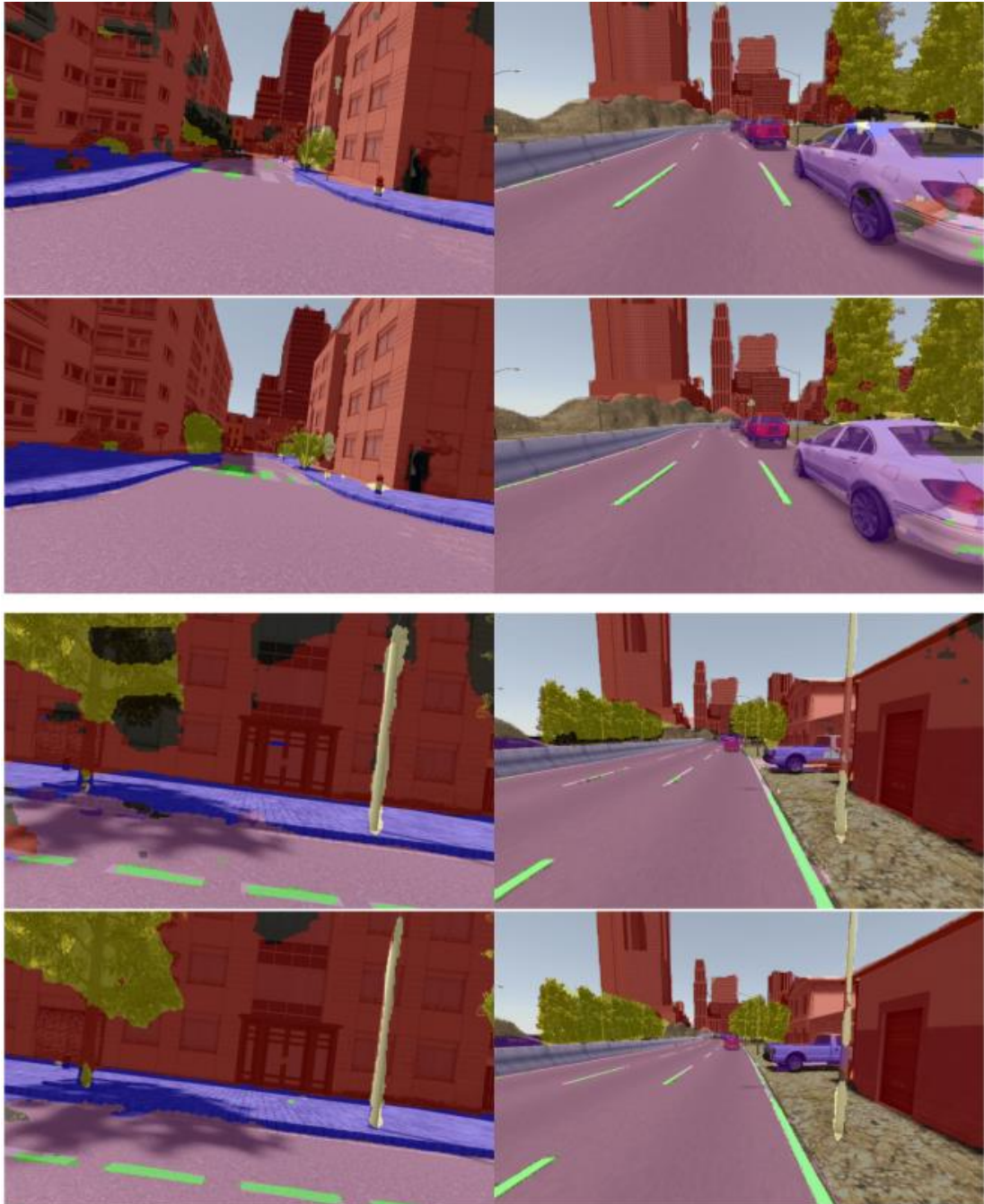


Figure 10: Qualitative checks for Synthia experiments in which data is superimposed on the network forecasts



Figure 11: Qualitative cityscapes tests contain studies that exceed the network feedback forecast. The strange lines are FCN-8s and also their accompanying numbers: the production of RFCN-8s

The values are considered as Semantic Segmentation of RFC-VGG compared to FC-VGG, Synthia Highway Summer Sequence. In FC-VGG the values are listed as 0.755, 0.504, 0.275, 0.946, 0.958, 0.840, 0.957, 0.762, 0.883 and 0.718. In RFC-VGG, the values are listed as 0.812, 0.566, 0.487, 0.964, 0.961, 0.907, 0.968, 0.865, 0.909 and 0.742 under the category named as Mean Class IoU and Per-Class IoU which are Car, Pedestrian, Sky, Building, Road, Side walk, Fence, Vegetation and Pole. The leading value rate illustrates that highest value is RFC-VGG.

The values considered as CityScapes Semantic Segmentation Results for RFCN-8s as against FCN-8s. In FCN-8s, the values are listed as 0.53, 0.917, 0.710, 0.792, 0.683 and 0.585. In RFCN-8s, the values are listed as 0.565, 0.928, 0.739, 0.814, 0.719 and 0.652. The RFCN-8s is the leading one.

3.2.3 Further Analysis

The values are considered as FC-Lenet, LSTM, GRU and RFC-Lenet measurements tested in synthesized MNIST datasets. The respectively categorized names are FC-Lenet, LSTM and GRU which are listed as Precision, Recall and F-measure. In FC-Lenet, the values are listed as 0.868, 0.922 and 0.894. In LSTM, the values are listed as 0.941, 0.786 and 0.856. In GRU, 0.955, 0.877 and 0.914. Lastly, the values are 0.96, 0.877 and 0.916.

3.3 Discussion

FCN-12s pre-training, recall and F-measurement on the ground map of FCN-12, RFC-12s on six sequences of check set motion detection requirements. Set motion detection norm. (d) and (EE) indicate decoupling and final integration of recurring units into the FCN respectively.

In FC-12s, the values are listed as 0.827, 0.585 and 0.685. In RFC-12s(D), the values are 0.835, 0.5887 and 0.69. In RFC-12s(EE), the values are 0.797, 0.623 and 0.7.

Number of layers: An rising number of layers can help to overfit. Therefore, the network should be intelligent enough to grasp the requisite mission. Because the data set becomes more diverse, the required depth becomes-(more classes, smaller areas, more variance in image). For example, even the low level LeNet may have a high output in the fairly simple SegTrackV2, (binary segmentation, dominant ROI). To

achieve respectable performance, Cityscapes and Synthia required substantially deeper networks. RFCN and FCN behaved in this way similarly.

Recurrent location of layer: The location of the recurrent layer has been modified from the third to the fifth one. This choice has a huge effect on computation and efficiency. Recurring layers are too expensive to practice on early layers with broad charts. However, we believe that the complexities should have been described better. For multiple layers of pooling insensitive to place artifacts in the image, which interfere with dynamic extraction. We have typically shown that the repeating layer earlier leads to success. It is well educated (it needs a lot longer to learn before reaching the first layer). We were restricted to pick this layer by the design of the skip. The repeated layer will be before or after both layers of the skip.

Moving Window Size: we measured 1-5 windows, which would be only an additional convolutional layer with a curtain. We have seen a small improvement in size 1, as anticipated, and success improves with size changes at the expense of more rigorous training that ultimately decreases efficiency.

FCN vs RFCN: RFCN 's networks took almost twice the FCN 's preparation period. Adadelta vs. SGD educated all networks differently. Nor have we shown a lot of batch-size flexibility. Interestingly, once RFCN achieves a acceptable degree of precision, it appears more consistently and seamlessly than FCN.

The architecture of the RFCN is very common and can be applied in many ways. This is usable in non-visual amounts. This can be incredibly useful for mobile robots or self-driving vehicles when exact details from its SLAM framework are usable. The

network will derive from this information a next image indicator. To order to learn objective dynamics, you can add optical flow to RFCN. It function can be taught at the same time as the key segmentation target.

3.4 Results of Uncertainty Estimation

Essentially, uncertainty assessment is measured of frame-level and pixel-level metrics. The pixel-level Metric assessment is inspired by precision-recall curve metrics. It indicates that the remaining pixels are exact as pixels that have greater percentile thresholds for an ambiguity.

Pixel level indicators are useful for assessing the calculation of uncertainty. However, calculating pixel ambiguity in actual implementations is challenging to exploit. For eg, the active learning machine needs to figure out which structure is essential to mark instead of choosing which pixel to mark.

Bayesian neural networks are known to model instability in neural networks. Prediction is hard to achieve for the Bayesian Neural Network. Kendall tau is measured on how the sequence of ranks is near to the sequence of ground truth.

Tiramisu MC dropout N=5

Table 2: Tiramisu MC dropout

	Accuracy
Global Accuracy	89.4
Mean Accuracy	75.4
Mean IOU	62.7

PR-Curve

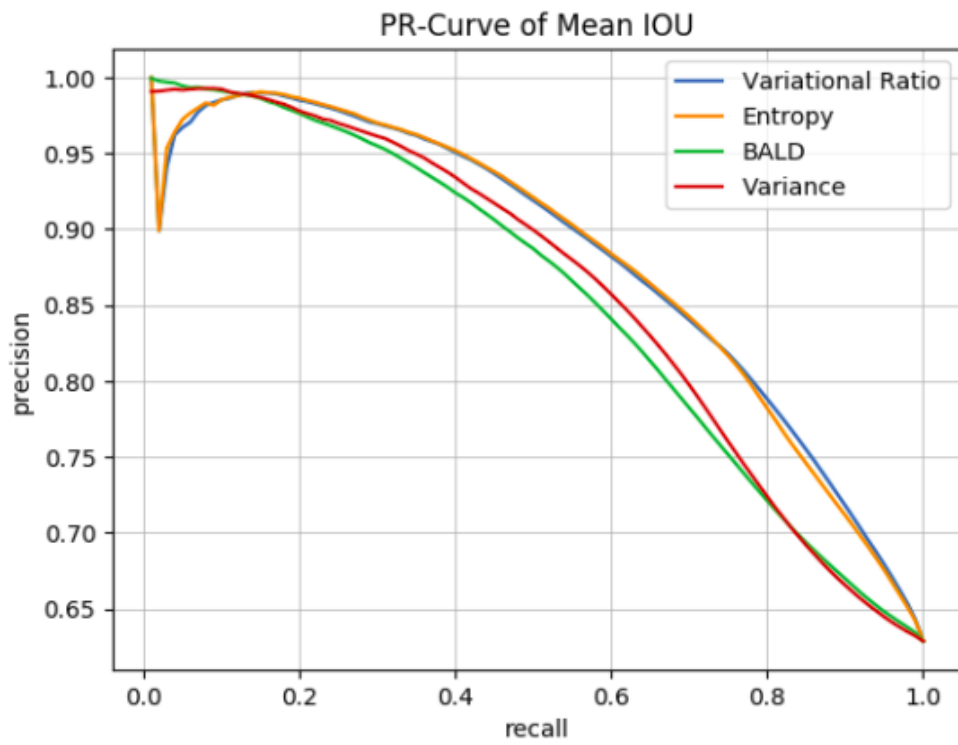


Figure 12: PR curve of mean IOU

Ranking IOU of Variational Ratio

Table 3: Ranking IOU of variational ratio

Percentage	Ranking IOU
10%	43.5
30%	58.1
50%	73.4
70%	85.3

Tiramisu TA-MC

Table 4: Tiramisu TA-MC accuracy

	Accuracy(%)
Global Accuracy	89.7
Mean Accuracy	73.6
Mean IOU	62.3

PR-Curve

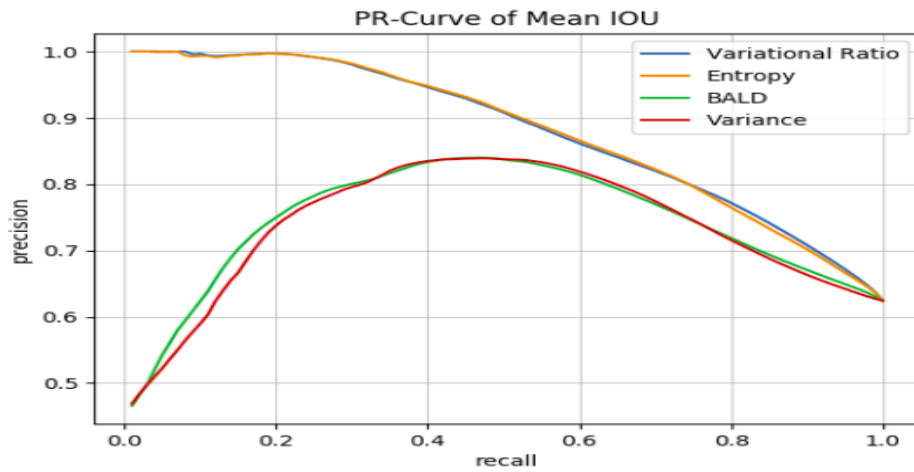


Figure 13: Tiramisu TA-MC PR-curve of mean IO

Ranking IOU of Variational Ratio

Table 5: Tiramisu TA-MC ranking IOU of variational ratio

Percentage	Ranking IOU
10%	34.9
30%	60.8
50%	76.8
70%	87.1

Tiramisu RTA-MC, Performance

Table 6: Tiramisu RTA-MC performance

	Accuracy(%)
Global Accuracy	89.7
Mean Accuracy	74.3
Mean IOU	62.7

PR-Curve

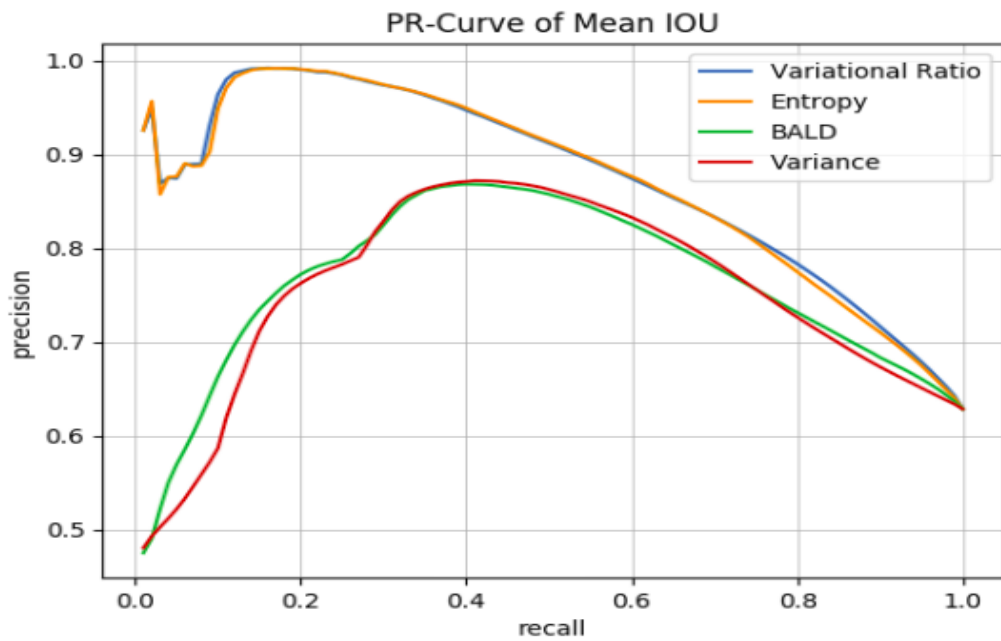


Figure 14 : Tiramisu RTA-MC PR curve of mean IOU

Ranking IOU of Variational Ratio

Table 7: Tiramisu RTA-MC ranking IOU of variational ratio

Percentage	Ranking IOU
10%	43.5
30%	65.3

50%	77.7
70%	86.5

The Tiramisu RTA-MC Ranking IOU of Variational Ratio is evaluated. The Segnet MC method can be respectively listed as Entropy, Variation Ratio, Mean STD and BALD values as 0.648, 0.669, 0.678 and 0.673. In Segnet TA-MC are 0.627, 0.632, 0.541 and 0.528. In Segnet RTA-MC listed as 0.663, 0.675, 0.628 and 0.621. In Tiramisu MC listed as 0.637, 0.654, 0.660 and 0.648. In Tiramisu TA-MC values listed as 0.661, 0.675, 0.664 and 0.627. In Tiramisu RTA-MC values are listed as 0.665, 0.679, 0.636 and 0.613.

The result of Tiramisu backbone in Ranking IoU. Methods respectively listed as MC, TA, RTA and BALD by comparing values in percentages of 10, 30, 50 and 70 classified in Entropy, Variation Ratio and Mean STD. In Entropy metric, MC method 34.9, 61.0, 70.8 and 86.5. In TA, 47.9, 63.9, 71.7, 84.7. In RTA, 47.9, 63.9, 74.2, 86.5. In Variation Ratio, In MC, 34.9, 61.0, 74.2, 86.5. In TA, 47.9, 65.3, 72.5, 85.9. In RTA, 52.2, 65.3, 76.0, 87.8. In Mean STD, the MC listed as 30.5, 63.9, 76.8 and 86.5. In TA, values are listed as 47.9, 75.5, 74.2 and 82.1. In RTA, 43.5, 68.2, 73.4 and 82.2. In BALD, the MC values are listed as 30.5, 62.4, 72.5 and 86.5. In TA, 43.6, 71.1, 71.7 and 80.3. In RTA, values listed as 47.9, 66.8, 71.7 and 81.0.

The SegNet backbone in ranking IOU values are evaluated. Methods respectively listed as MC, TA, RTA and BALD by comparing values in percentages of 10, 30, 50 and 70 classified in Entropy, Variation Ratio and Mean STD. In Entropy metric, MC method 47.9, 59.5, 72.5 and 85.3. In TA, 47.9, 65.3, 69.9, 85.89. In RTA, 47.9, 66.8, 73.4, 89.0. In Variation Ratio, In MC, 47.9, 62.4, 74.2, 85.3. In TA,

52.3,63.9,70.8,85.3.. In RTA, 47.9,65.3,76.8,88.4.In Mean STD, the MC listed as 52.3,66.8,71.7,86.5.In TA , values are listed as 43.6,65.3,65.6,74.2. In RTA, 43.6,69.7,69.9 and 82.8. In BALD, the MC values are listed as 47.9,66.8,71.7,87.8.In TA,43.6,62.4,64.8,74.2 In RTA , values listed as 43.6,69.7,67.3 and 82.2.

In the case of MC dropout, the Neural networks in Bayesian are illustrated that learning a distribution over certain parameters rather than a sequence of deterministic parameters. In the course of training data X and Y , it missions to determine the further distribution of the weight W of the neural network. The Bayesian neural networks with MC dropouts can produce improved estimates of efficiency and uncertainty. However, it does require samples of N times to forecast the images that are N times slower than the initial network. In real-times systems, including self-driving vehicles, which must be predicted and calculated as soon as possible to prevent the MC decline. We suggest the temporary aggregation MC dropout in order to facilitate the operation of the MC dropout.

In Temporal Aggregation MC Dropout, the aggregation of time method MC dropout uses the video templates. As the video has consecutive frames, several separate frames will contain the same artifacts and hence be diverted repeatedly via the Bayesian model. If a video comprises static frames, the average performance of N consecutive frames is the same as MC dropout of N samples.

In Region-Based Temporal Aggregation MC Dropout, Our TA-MC dropout progresses in most situations, but the uncertainty calculation is not correct when the optical flow calculation is incorrect. The optical flow can not be reliable for such environments that involve quickly moving artifacts or occlusion. In order to overcome this issue, we

suggest a geographic time averaging that can dynamically allocate various multiplier factors in each pixel based on its reconstruction error. The effect is a misalignment of the estimated average of motions on the incorrect patch.

In this study , we suggest a regional based temporal aggregation (RTA) approach for simulating the Monte Carlo (MC) video segmentation sampling process. Our RTA approach uses time data by videos beside of this only requires a one-time observation for showing prediction and uncertainty in each frame. RTA will obtain comparable outcomes on the CamVid data set as against the general MC dropout with just drop by 1.2% considering on mean IoU metric and an amazing 10.97 times speeding up the inference process. In addition, by using Entropy and Variance Ratio as the unsecurity approximation parameter, the instability produced by the RTA approach is close to pixel-level metrics and also resulted the MC dropout based on frame-level metrics. In essential applications, it is more necessary to correctly obtain the uncertainty of the instance stage.

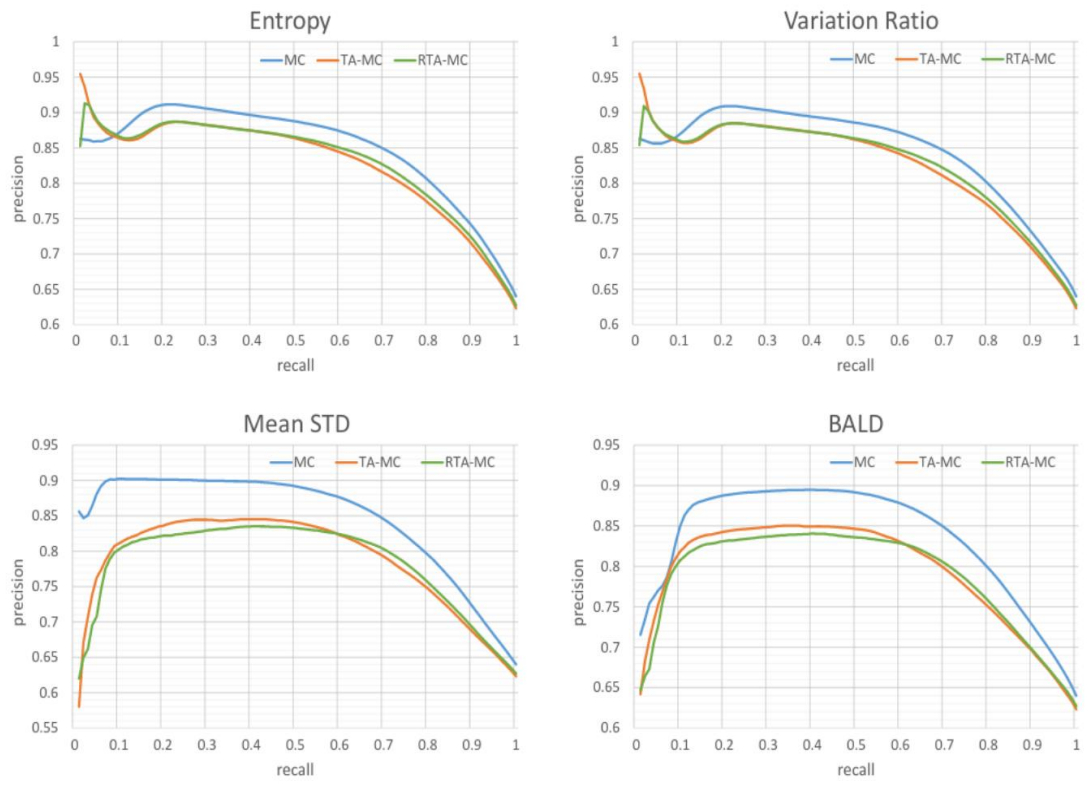


Figure 15: Precision-Recall curves of pixel level of Segnet backbone.

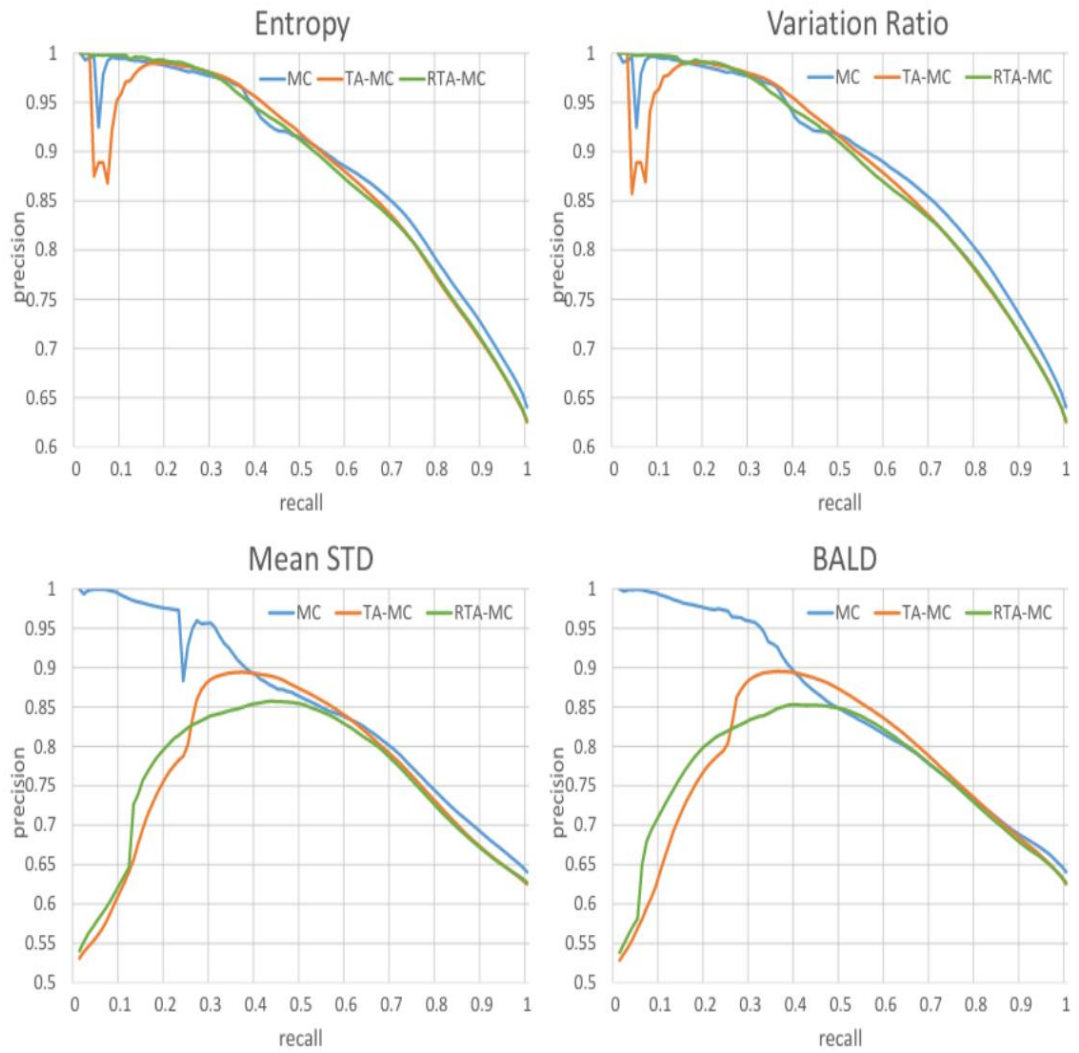


Figure 16: Precision-Recall curves of pixel level of Tiramisu backbone

Chapter 4

THE HUMAN ROBOT INTERACTION IN INCREMENTAL LEARNING

One goal of today's research is to create robot assistants to help people deal with everyday tasks. However, in home environments only a few robot systems were successful. Commercial robots have been designed to do some function, e.g. robot vacuums and lawnmowers. In comparison, several robots are not yet in place that can handle a wide variety of household tasks. One significant explanation for this is that existing robotics are not well-connected. Calling HRI a key component of device architecture, it is simpler to incorporate robotics with daily human activities.

Yanco et al[45] performed a analysis at the latest DARPA Robotics Challenge [45]. The researchers understood this and concentrated on user engagement with the human being in the system model instead of attempting to gain autonomy. In this model, human awareness and advice are used when the computer can not decide itself [47, 48, 49]. In order to provide this guidance it is crucial that, like humans, a common understanding of the earth shared between man and robot is created. For starters, a new apprentice of metalworkers needs to know fast about specific form of materials (stain, bronze, nylon, etc.) and machinery (purple, boiling, frying, dull, design, etc.), etc. Otherwise, the information passed on to it will be difficult to understand.

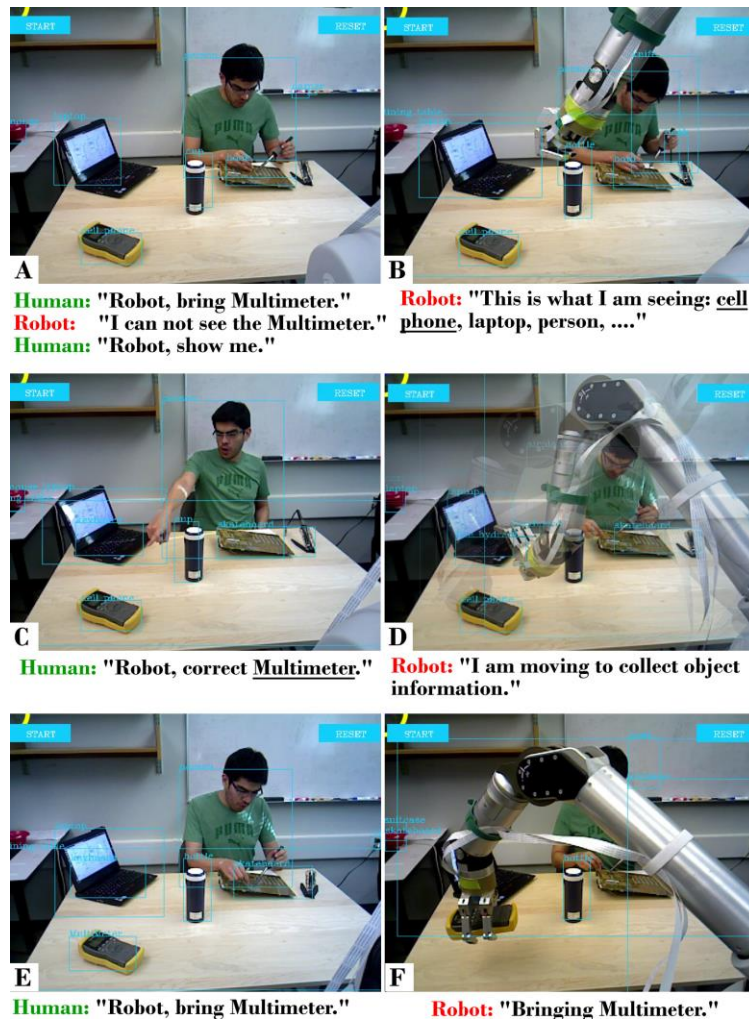


Figure 17: The HRI of incrementing robot knowledge

The human tells the robot to take a circuit board with the multimeter. The system doesn't realize what a mixer is. The human being wants the world's robot image. The robot is iterated over the observed objects by pointing and telling each mark of the objects. The person point and correct the "multimeter" sign, originally identified as "mobile phone." The robot targets the target pointing and takes correcting photographs of the device. The human being is again calling for the multimeter. The work of the robot is good this time. The additional video [50] displays our gui.

Regarding robotics, the same concept can also be recognized. However, robots must grasp the simple universe before they can begin to benefit from guidance. A fundamental concept can be used in target position and entity type recognition. Deep learning shows that certain traditional solutions to computer vision in this area are superior. Most of the first attempts was to kill [12]. The creation and labeling of bounding boxes for each item according to its class was using a convolutional neural network (CNN). This method has attacked growing amounts of national plans both for consistency and, sadly, for computer-based approaches. More recent efforts [13][32] have also allowed use of the RPNs to regress boundaries. RPN will operate on maps rather than photos and thereby bypasses the need to recalculate the feature charts.

CNN applications in a number of datasets of object detection have achieved state-of-the-art. Nonetheless, the tacit premise is that the batch offline testing data collection covers all relevant types of items. Sadly, actual life visibility is special. The algorithm also encounters artifacts not in the training data at the moment of prediction or may appear quite special relative to training instances. The robot knowledge cycle also has to be updated. The IML mainly addresses new instances of established types. Nonetheless, there are two major problems with OSR approaches. One is to identify new categories continuously, and two is to change the system to include the most recent category. Of eg, it is almost impossible to distinguish a sugar box from a detergent box near, often without semanticized labeling or reading the label or having meaning. Nonetheless, we can use a robot 's voice to solve this problem. In particular, robots can communicate with and accept orders or inputs. You will also discover the world with an on-board camera. You propose a separate approach along with this overall plan to

improve the visual perception of the robot gradually. Our task is to detect and place objects. The ability to clarify definitions and to correct HRI 's false interpretations.

4.1 Humans for Incremental Learning

The human-based robot is also referred to in literature as Demonstration Programming (PbD), Demonstration Based (LfD) or Imitation Learning. Automatic learning approaches can be criticized without constructive human involvement, for instance, during the learning process. Throughout this area, the key challenge is to learn the motion direction from user demos for tasks. Although the perceptive efficiency of the robot in this area is overlooked, several primary challenges remain. Interfaces for knowledge sharing and, in particular, strategies for gradual learning. Below is a short summary of work into these problems.

Since pathways for behavior are the key concern in the classroom, most of the interfaces are designed to consider the directions of students. [58][59][60], using individual kinesthetic instruction. The robot guide to the task [67] and teleoperation are the principal approaches for this purpose; the instructor explicitly monitors robotic variables through an interface [61][62]. While these interfaces are suitable for the given direction, they are not standard for human interaction. In addition, more interfaces were created, of course. Additional data are Speech [63][64] and gestures [65]. The robot uses the most sophisticated learning methods to bootstrap and slowly improve its existing abilities [66].

4.1.1 The Robotics of Deep Learning

The popularity in computer vision education was robotic motivation. The lack of reliable and standardized algorithms was one of the key obstacles in the unregulated usage of robotics. In several implementations, methods of deep learning (DL) can now

be implemented with specific and robust efficiency. Classification of photographs and analysis of natural languages in particular. This modern research has gained greatly. DL models can also be used as an image and speech preprocessing tool to transform raw sensor data into a smaller scale unit. Robot grasp recognition [68] utilizes raw pictures to classify grasp points for different items and can be used to catch them. End-to-end DL can also be used when the device is in the network. Input raw to monitor signal production of the robot's effectors. Levine et al.[69] introduced and showed by executing such tasks the viability of such a network. The concept of end-to-end control networks is appealing since comprehensive engineering is required in theory for each node. Nevertheless, several other hyperparameter tunings and specific training treatment in these networks have become evident and essentially quite significant. The key explanation is that all solutions so far suggested need far more knowledge than is technically feasible. It led to the paradigm intensifying instruction, including learning from presentation methods. Deep networks can also practice incredibly complicated computational decision-making, particularly though sufficient knowledge can be used. We accept, thus, that it is generally more realistic to use DL enablers as modules. The sum of the data for a specific mission was trained and finished during these groups.

Detection Networks and Object Localization

The goal is to classify objects and methods of sensing, typically by relation and classification of boxes, to locate all item. Both methods are particularly useful for the comprehension of the scenario and other robots may be envisaged. Overfeat [12] was one of the first CNN to do so. This approach has shown that spatial concepts are both reliable and computationally expensive. Consequently, recent trials [13][32] have employed boundary boxes through National Proposal Networks (RPNs). Other than

input pictures, RPN will work on feature maps, thus eliminating the need to recalculate maps.

4.1.2 The Open Set Recognition Multi Class

The state-of-the-art object detection CNN projects use a number of datasets. They specifically presume that the dataset includes all conceivable types of objects. In the modern world, though, the case is somewhat different. The algorithm faced artifacts which at the time of the estimation were not in the testing datasets. It is particularly important for robotics in families and assistants. Things are specific in-household and work location, and not all of them are accessible in an comprehensive collection of results. The question is not just artifacts. For eg, an assistant robot would handle each patient in the hospital differently, and the patient should be conscious of that. Current recognition techniques may only identify the patient as an person, though.

Literature tackle this issue with gradual machine learning (IML) and transparent range recognition (OSR). The key objective of the IML [52][51] is the treatment of new established class instances. However, two new issues have to be discussed by OSR [56][53]. The first is to create new categories on a continuous basis and the second is to adjust the existing categories.

4.2 The Structure

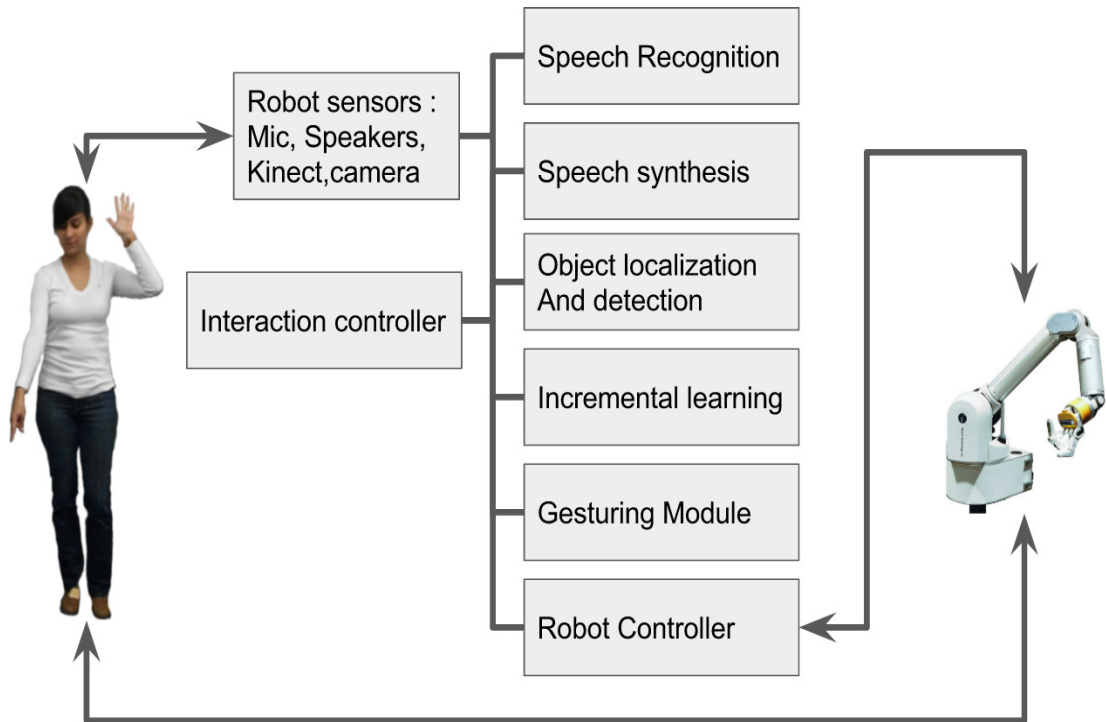


Figure 18: System block diagram

The 7-DOF WAM [70] . It comprises seven modules as shown in the figure. Both devices are fully compatible with ROS [71]. The CMU Sphinx toolkit [72] is included in the speech recognition module. It provides the robot's basic word and sentence recognition for shifting states during interaction. The speech synthesis module is based on the speech synthesis system of the Festival [73] and provides human feedback in a verbal channel.

The location and detection module for objects provides labels and 2D locations of the objects in the scene. The Incremental learning module uses HRI to allow changes in the world of the robot.

The gesture module is based on our previous work [74] [14], which proposes a non-verbal robotic vision system that infer human pointing and performs simple tasks based on commands of human gesture. Our ability to reduce speech description and make interactions more humane is integrated into our system.

Illustrate four major computer contacts. Verbal interaction may establish fundamental verbal communication via the robot for comprehension and integration of human and machine language. The environmental interpretation refers to the relation between speech and action of the robot. The system for identification and positioning includes the 2D border boxes and markings of the artifacts detected. In the correct conversation, participants are utilizing verbal and gestural language to annotate a certain object in the scene that requires clarity. The goal region is hit by a 3D ray with verbal orders from these two stages. The robot is then guided to gather data from this position for the 3D model. The data is collected via a parametric helix curve, which holds the camera in front of the goal. The TLD tracker [75] is used for retrieval to insure that the object is captured during data collection. The tracker 's initial boundary box is presented as the entire picture of robots start a list near the object.

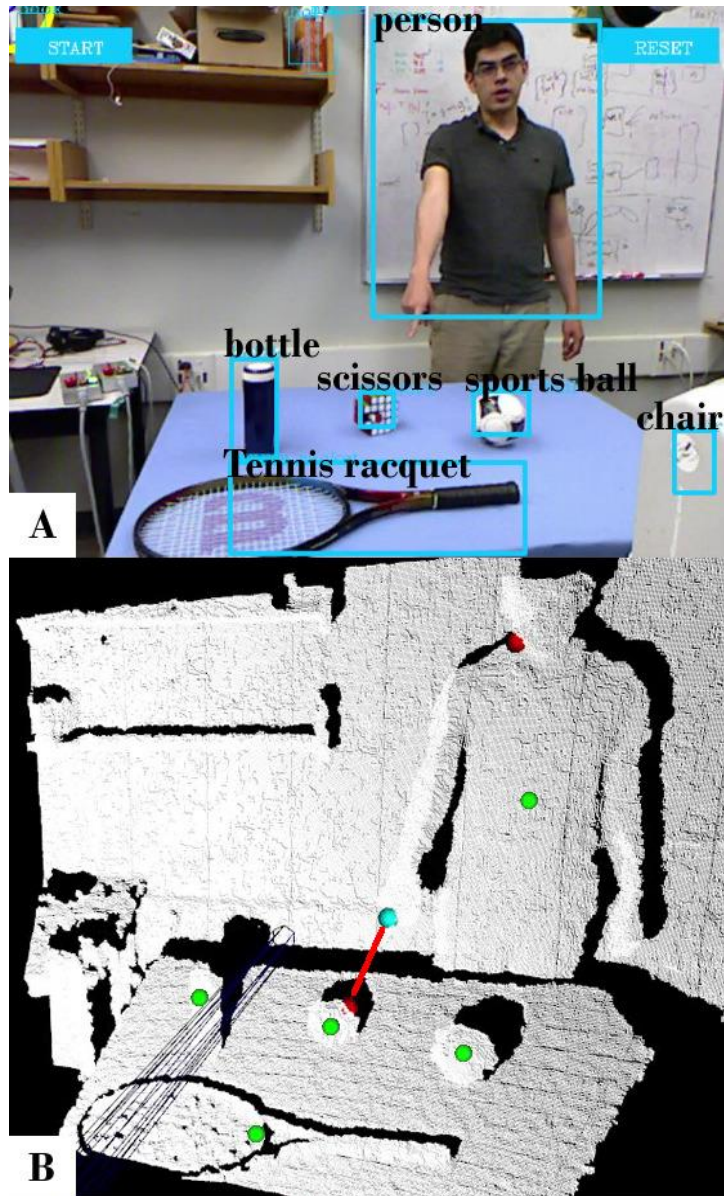


Figure 19: Visualization of RGB. Objects are detected and located in the scene.

4.3 Investigation

To verify our methodology, we evaluated three key components of our process. Target identification and understanding, progressive learning algorithms and, essentially, positive thinking by human input.

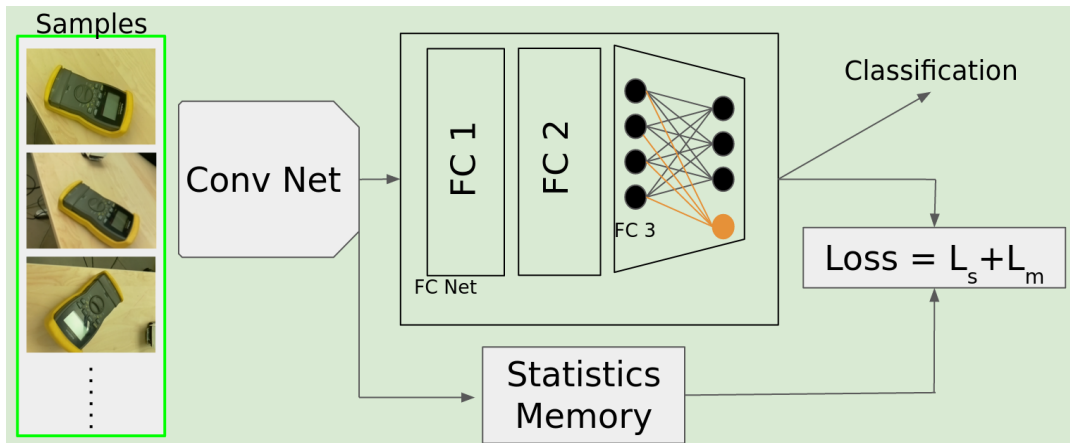


Figure 20: Incremental learning pattern

4.3.1 The HRI of Incremental Learning

The robot begins with the pre-trained recognition and identification of fundamental objects, and at the laboratory we are attempting to teach them to identify new topics. We have thus added new artifacts with our framework one by one to the device. Each new item was retrieved by the robot.

Image of the eye-in-hand. The Fig incorporates quotations from these items. After each package is made, the new component is attached in real time to the robot detection module and then another element is inserted in the same manner. Once and new class is added, the accuracy of the test set MS-COCO plus the new class evaluation section is tested. The comprehensive assessment data set contains the test portion where each new class is evaluated.

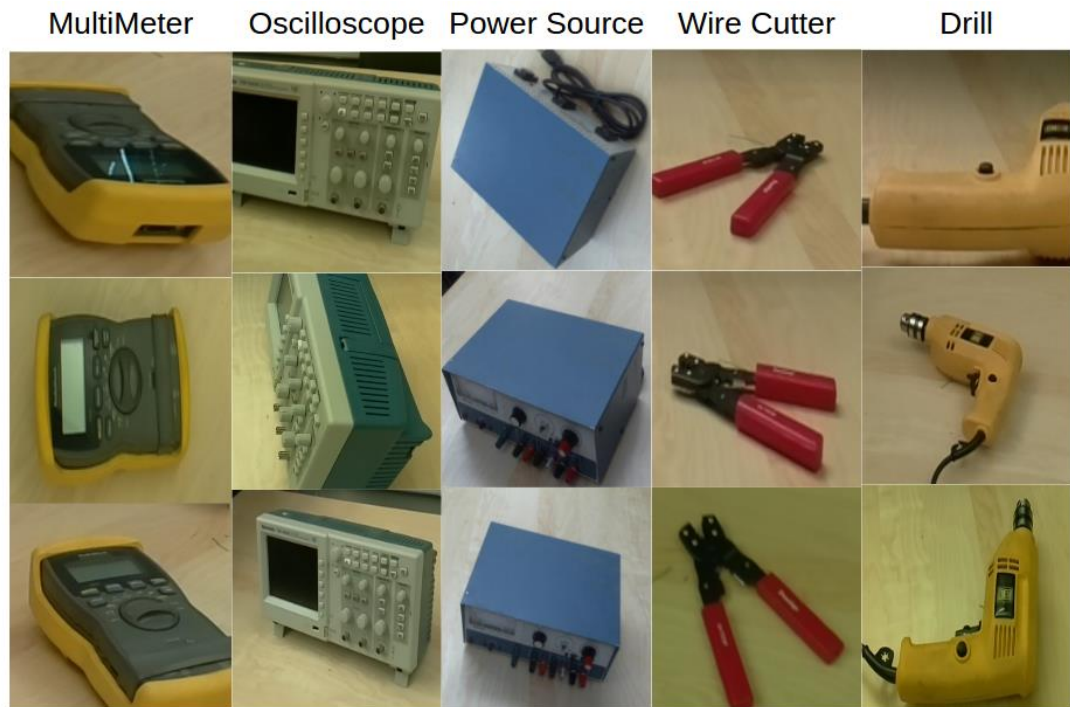


Figure 21: Multimeter

The findings of the first and second strategy studies are presented. All displayed values are of maximum accuracy. However, as we expected, the precision is slightly decreased when new items are added, this drop is not substantial and the pitch is high , particularly during the second method. For a total of 5 percent declining after 11 new objects have been added/ we may therefore conclude that accuracy stays available only though several new objects have been introduced. Remember that in the baseline model there are still 80 popular items, so not much addition is required.

4.3.2 Recognition Baseline and Object Detection

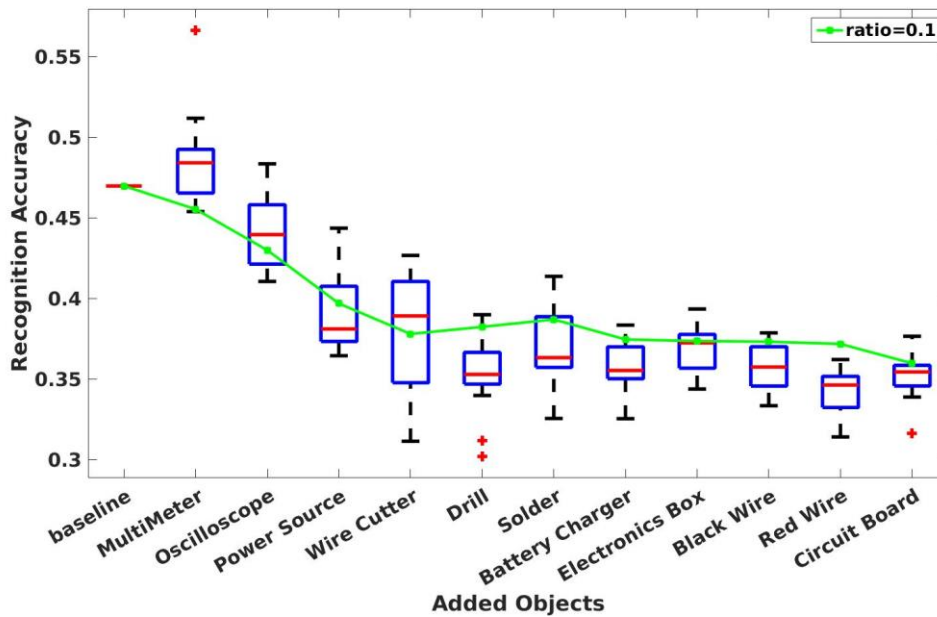


Figure 22: Recognition success in an gradual learning environment in the first step. The robot takes pictures and adds new objects one by one. The displayed values are the highest accuracy.

The Identification and identification time of our computer is 150ms for a GeForce 960 GPU. This is twice as high as 350ms [32] Titan X GPU R-CNN. Average precision (AP) of 0,2026 was achieved in the target detection feature, which is equivalent to the state of the art in .224 [13]. Based on AP metric a calculation is accurate if the land data indicate an IOU of more than 0.5. We have determined the top-1 identity accuracy factor. The assumption is right since the class as far as possible is similar to the simple truth. We have achieved 0.45 accuracy. The success of some MS-COCO identification is unknown to us. However, given the complexities of MS-COCO in comparison with the imagenet, a precision close to this value is necessary.

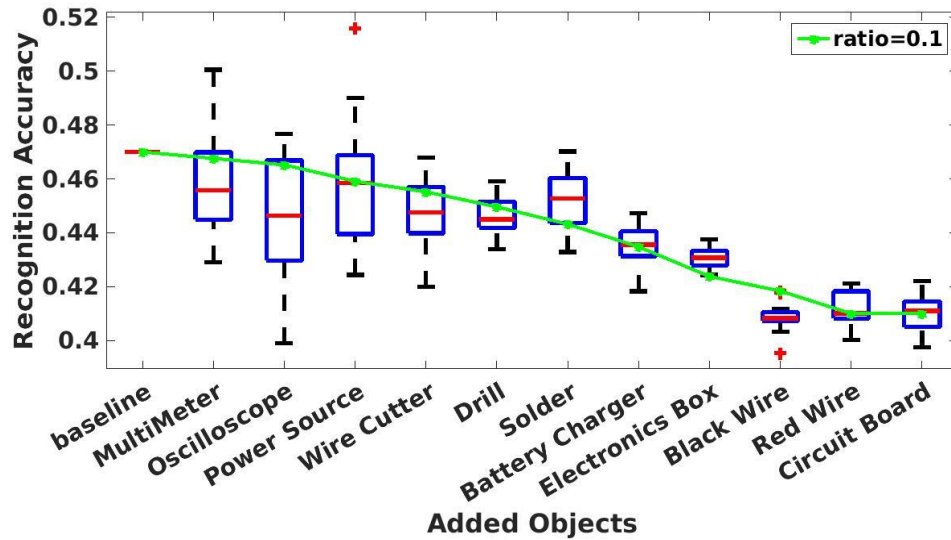


Figure 23: Incremental learning scenario recognition performance with the second approach. The user introduces new objects one by one and the robot collects their images. The values shown are the highest accuracies

4.3.3 The Approach of Incremental Learning

In order to validate our methodology by academics, we checked our incremental learning framework with publicly accessible datasets. We started this experiment. We adopted the same testing procedure as the first study except now imagenet is the named new object data instead of HRI.

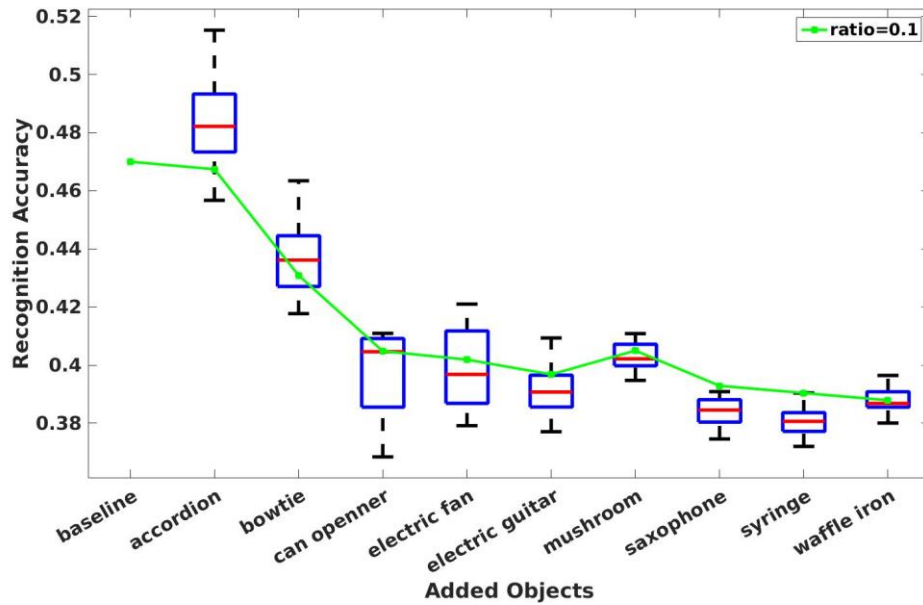


Figure 24: Recognition performance after the first approach is gradually added to new classes. Data are taken from imagenet for new classes. The values shown are the highest accuracies

4.3.4 The Performance Evaluation of Mock Human-Human Interaction

When all items are grouped into two containers, each test session starts with items on the table and ends altered. There have been four mediums tested and listed below. By pointing to them, the instructor displays objects and target bins. The actor follows the hand and maybe the object and bin location in the direction.

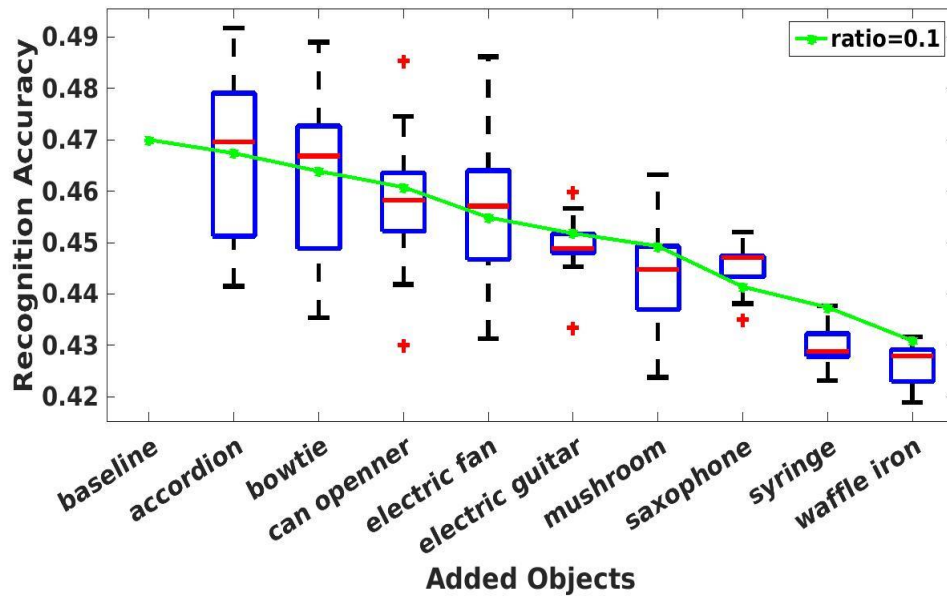


Figure 25: Identification of success by slowly introducing different second method levels. Data was taken from imagenet for new classes. The values displayed are the most correct. Box plots show the difference in precision, as the relation between new and old class study samples ranges between 0.05 and 0.5. The green line is the accuracy of the 0.1 ratio

Speech: The professor explains the item to be presented and displays the containers on the left which are relative locations.

Clicking: both the professor and the participant are used as a guide for this program. When the professor taps on them, he chooses items and selects the appropriate container. The actor watches the console and executes the steps by tapping.

Through this method, the professor will use both his speech and points to express his thoughts. The director listens to the teachers and often tracks their attention and vision to identify the target and the containers.

We evaluated these media in the same way for the same research sample. The actor held the same stuff for both tests. The point and speech mix has a slight advantage on

the others, but there are no definite winners. Time and commitment to accomplish the mission in question. The second argument is attributed to the direct edge of the intellectual demand axis. The next is pressing. It was the most reliable, but it took the most time to finish. The least favorite gui was voice. And there was a strong performance benefit in terms of time and mental demand for physical demand (which is essential for physically handicapped people). Naturally, an object definition appeared to be tougher than in other situations. One lesson from this experience is to turn physical movements into a far more intuitive and comfortable design for humans. Speech often helps particularly if contact between the teacher and the actor became more difficult and required at the higher stage.

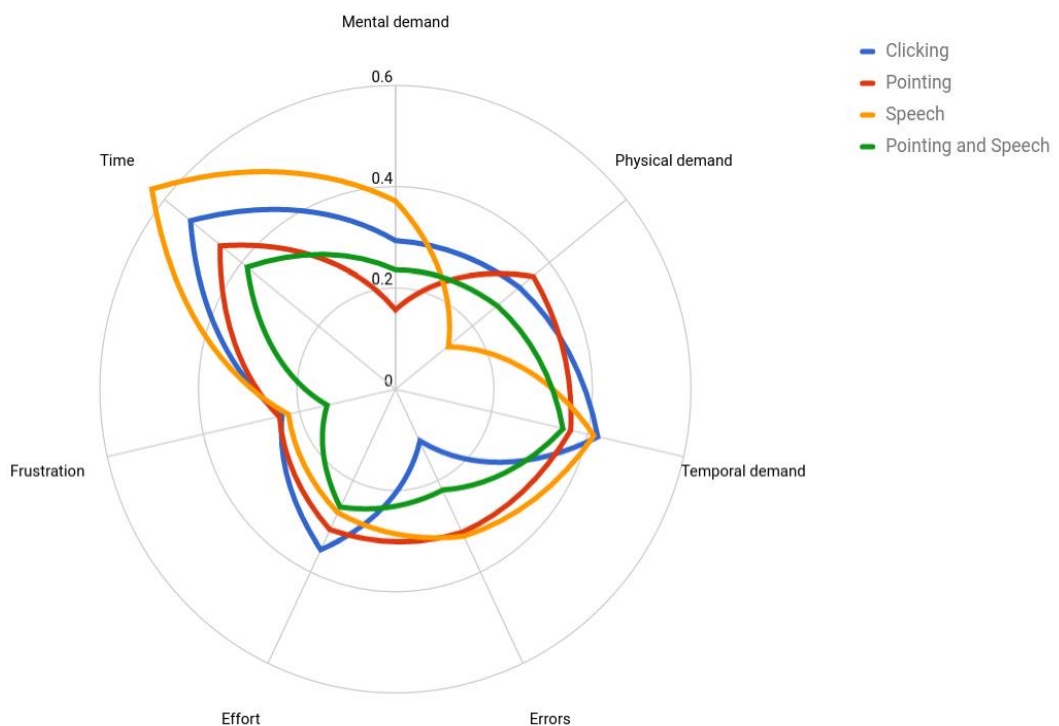


Figure 26: Evaluation of NASA task load index for 4 interfaces tested

4.4 Analysis

We have introduced and built a full HRI program that uses human intelligence to enhance robotic perception. Our program may understand the existence of a new entity and incorporate it to its knowledge base through means of human guidance. The human interface; it makes it easier for people to see and fix the understanding void in movement and expression of robots.

We have proposed two suitable incremental learning models for our system. The strategies were tried and improved. The second gradual learning method became more successful when new artifacts were introduced, indicating fewer efficiency.

While we have seen progressive learning only for visual activities, the concept is not restricted to it and can be applied to teaching. Shift of human to robot trajectory development skills is currently well known. We have learned and used it on our robots from prototype programs. It must also be included in the inclusive program of schooling.

Chapter 5

CONCLUSION AND FUTURE WORK

We addressed some issues in the use of fundamental robotic learning in this article. Namely the lack of time and dependence only on train details that were originally available. Instead we explored how robots can somehow communicate with humans and the world by utilizing their sensors.

We also developed recurrent, completely convergent video segmentation networks that suit well with many robotic activities, such as self-driving vehicles. We demonstrated how to build and train a network like this and then performed experiments with common video segmentation criteria. We have seen gradual changes compared with standard completely convolutional networks.

In a robotic world, we addressed object recognition and perception issues and presented a potential solution for the comprehension of robotics through human robot interactions. We have created a comprehensive interface which helps people to speak, talk, gesture and see. We have seen how this contact will provide the robot with additional knowledge. In the deep object detection-recognition module, two different methods were merged and validated. We also shown how our program can boost the vision of robots with picture evidence in virtual environments and actual everyday artifacts.

An significant thing is the ability of robots to control their surroundings. But these capabilities have proven challenging to learn. When developers tackle such functions of engineering, an ultimate answer is out of control. Nevertheless, as with imagined objects, the robot should not learn how to execute all deceptive functions, as long as certain abilities can be built easily. With this method and trajectory development, you can envision an evolutionary learning cycle for practice. We play with a prototype learning method , which uses a deep network for trajectory learning (exact RNN). It allows different signals and potentially a more general structure to be implemented.

The simulator route planner manages a 7 DoF robot for a mission. The paths are documented and used for training a GRU with several layers, which in the current states provides the next control signal. Promising early tests demonstrate that the analysis can be absorbed by the network. We can see, however, that a deep network needs generalization. The sensory feedback is immediately accessible in the device. The RFCNN is appropriate as it may combine sets of photographs or other sensory signals. The function of separation should be seen as a service. The purpose of directing the training is to generate the correct control signal for the main objective. We must follow this direction in future research and investigate alternative strategies that can direct and develop the network. Several solutions provide memory networks to independently save and activate numerous skills in an reasonable period. It may also be a strong mix of visual servoing and LfD.

REFERENCES

- [1] Lihi Shiloh-Perl and Raja Giryes, Introduction to Deep Learning ,2020.

- [2] Guosheng Lin, Chunhua Shen, Ian Reid, et al. Efficient piecewise training of deep structured models for semantic segmentation. arXiv preprint arXiv:1504.01013, 2015.

- [3] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.

- [4] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. Densebox: Unifying landmark localization with end to end object detection. arXiv preprint arXiv:1509.04874, 2015.

- [5] Lijun Wang, Wanli Ouyang, Xiaogang Wang, and Huchuan Lu. Visual tracking with fully convolutional networks. *In Proceedings of the IEEE International Conference on Computer Vision*, pages 3119–3127, 2015.

- [6] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. *In Proceedings of the IEEE International Conference on Computer Vision*, pages 1520–1528, 2015.

- [7] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.

- [8] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. *In Computer Vision and Pattern Recognition*, 2005. CVPR 2005. IEEE Computer Society Conference on, volume 1, pages 886–893. IEEE, 2005.
- [9] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *In Advances in neural information processing systems*, pages 1097–1105, 2012.
- [10] Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. *In Advances in neural information processing systems*, pages 2553–2561, 2013.
- [11] Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala, and Yann LeCun. Pedestrian detection with unsupervised multi-stage feature learning. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3626–3633, 2013.
- [12] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, *localization and detection using convolutional networks*. arXiv preprint arXiv:1312.6229, 2013.
- [13] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *In Advances in neural information processing systems*, pages 91–99, 2015.

- [14] Camilo Perez Quintero, Romeo Tatsambon, Mona Gridseth, and Martin Jägersand. Visual pointing gestures for bi-directional human robot interaction in a pick-and-place task. In *Robot and Human Interactive Communication (RO-MAN)*, 2015 24th IEEE International Symposium on, pages 349–354. IEEE, 2015.
- [15] Michael A Goodrich and Alan C Schultz. Human-robot interaction: a survey. *Foundations and trends in human-computer interaction*, 1(3):203–275, 2007.
- [16] Cynthia Breazeal, Guy Hoffman, and Andrea Lockerd. Teaching and working with robots as a collaboration. In *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 3*, pages 1030–1037. IEEE Computer Society, 2004.
- [17] Cynthia Breazeal. Toward sociable robots. *Robotics and autonomous systems*, 42(3):167–175, 2003.
- [18] Masato Hirose and Kenichi Ogawa. Honda humanoid robots development. *Philosophical Transactions of the Royal Society of London. A Mathematical Physical and Engineering Sciences*, 365(1850):11–19, 2007.
- [19] Giulia Pasquale, Carlo Ciliberto, Francesca Odone, Lorenzo Rosasco, Lorenzo Natale, and Ingegneria dei Sistemi. Teaching icub to recognize objects using deep convolutional neural networks. *Proc. Work. Mach. Learning Interactive Syst*, pages 21–25, 2015.

- [20] Björn Jensen, Nicola Tomatis, Laetitia Mayor, Andrzej Drygajlo, and Roland Siegwart. Robots meet humans-interaction in public spaces. *IEEE Transactions on Industrial Electronics*, 52(6):1530–1546, 2005.
- [21] Aurélie Clodic, Sara Fleury, Rachid Alami, Matthieu Herrb, and Raja Chatila. Supervision and interaction. In *Advanced Robotics, 2005. ICAR'05. Proceedings., 12th International Conference on*, pages 725–732. IEEE, 2005.
- [22] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [23] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [24] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1529–1537, 2015
- [25] Francesco Visin, Kyle Kastner, Aaron Courville, Yoshua Bengio, Matteo Matteucci, and Kyunghyun Cho. Reseg: A recurrent neural network for object segmentation. *arXiv preprint arXiv:1511.07053*, 2015.

- [26] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. *In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [27] Vincent Michalski, Roland Memisevic, and Kishore Konda. Modeling deep temporal dependencies with recurrent grammar cells. *In Advances in neural information processing systems*, pages 1925–1933, 2014.
- [28] Graham W Taylor, Rob Fergus, Yann LeCun, and Christoph Bregler. Convolutional learning of spatio-temporal features. *In Computer Vision—ECCV 2010*, pages 140–153. Springer, 2010.
- [29] Mircea Serban Pavel, Hannes Schulz, and Sven Behnke. Recurrent convolutional neural networks for object-class segmentation of rgb-d video. *In Neural Networks (IJCNN), 2015 International Joint Conference on*, pages 1–8. IEEE, 2015.
- [30] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [31] Alex Graves. Generating sequences with recurrent neural networks. *ArXiv preprint arXiv:1308.0850*, 2013.

- [32] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Denscap: Fully convolutional localization networks for dense captioning. *arXiv preprint arXiv:1511.07571*, 2015.
- [33] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2625–2634, 2015.
- [34] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- [35] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [36] Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville. Delving deeper into convolutional networks for learning video representations. *arXiv preprint arXiv:1511.06432*, 2015.
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

- [38] Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, James Bergstra, Ian J. Goodfellow, Arnaud Bergeron, Nicolas Bouchard, and Yoshua Bengio. Theano: new features and speed improvements. *Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop*, 2012.
- [39] Nil Goyette, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, and Prakash Ishwar. Changedetection. net: A new change detection benchmark dataset. *In Computer Vision and Pattern Recognition Workshops (CVPRW), 2012 IEEE Computer Society Conference on*, pages 1–8. IEEE, 2012.
- [40] Fuxin Li, Taeyoung Kim, Ahmad Humayun, David Tsai, and James M Rehg. Video segmentation by tracking many figure-ground segments. *In Proceedings of the IEEE International Conference on Computer Vision*, pages 2192–2199, 2013.
- [41] F Perazzi, J Pont-Tuset, B McWilliams, L Van Gool, M Gross, and A Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation.
- [42] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3234–3243, 2016.
- [43] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. *In Proceedings of the*

- IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016.
- [44] Matthew D Zeiler. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- [45] Holly A Yanco, Adam Norton, Willard Ober, David Shane, Anna Skinner, and Jack Vice. Analysis of human-robot interaction at the darpa robotics challenge trials. *Journal of Field Robotics*, 32(3):420–444, 2015.
- [46] Gill Pratt and Justin Manzo. The darpa robotics challenge [competitions]. *Robotics & Automation Magazine*, IEEE, 20(2):10–12, 2013.
- [47] Mona Gridseth, Oscar Ramirez, Camilo Perez Quintero, and Martin Jagersand. Vita: Visual task specification interface for manipulation with uncalibrated visual servoing. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3434–3440. *IEEE*, 2016.
- [48] Adam Eric Leeper, Kaijen Hsiao, Matei Ciocarlie, Leila Takayama, and David Gossow. Strategies for human-in-the-loop robotic grasping. In *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, pages 1–8. ACM, 2012.
- [49] Camilo Perez Quintero, Oscar Ramirez, and Martin Jägersand. Vibi: Assistive vision-based interface for robot manipulation. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4458–4463. *IEEE*, 2015.

- [50] AIXploratoriu. <https://webdocs.cs.ualberta.ca/~vis/HRI/HR> (accessed Feb 2016)
- [51] Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7(Mar):551–585, 2006.
- [52] T Poggio and G Cauwenberghs. Incremental and decremental support vector machine learning. *Advances in neural information processing systems*, 13:409, 2001.
- [53] Abhijit Bendale and Terrance Boulton. Towards open world recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1893–1902, 2015.
- [54] Lorenzo Bruzzone and D Fernandez Prieto. An incremental-learning neural network for the classification of remote-sensing images. *Pattern Recognition Letters*, 20(11):1241–1248, 199
- [55] Walter J Scheirer, Anderson de Rezende Rocha, Archana Sapkota, and Terrance E Boulton. Toward open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(7):1757–1772, 2013.
- [56] Walter J Scheirer, Lalit P Jain, and Terrance E Boulton. Probability models for open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2317–2324, 2014.

- [57] Dana Kulić, Wataru Takano, and Yoshihiko Nakamura. Incremental learning, clustering and hierarchy formation of whole body motion patterns using adaptive hidden markov chains. *The International Journal of Robotics Research*, 27(7):761–784, 2008.
- [58] Aleš Ude, Christopher G Atkeson, and Marcia Riley. Programming full-body movements for humanoid robots by observation. *Robotics and autonomous systems*, 47(2):93–108, 2004.
- [59] Seungsu Kim, ChangHwan Kim, Bumjae You, and Sangrok Oh. Stable wholebody motion generation for humanoid robots to imitate human motions. *In 2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2518–2524. IEEE, 2009.
- [60] Shin’ichiro Nakaoka, Atsushi Nakazawa, Fumio Kanehiro, Kenji Kaneko, Mitsuharu Morisawa, Hirohisa Hirukawa, and Katsushi Ikeuchi. Learning from observation paradigm: Leg task models for enabling a biped humanoid robot to imitate human dances. *The International Journal of Robotics Research*, 26(8):829–844, 2007.
- [61] Adam Coates, Pieter Abbeel, and Andrew Y Ng. Learning for control from multiple demonstrations. *In Proceedings of the 25th international conference on Machine learning*, pages 144–151. ACM, 2008.
- [62] Daniel H Grollman and Odest Chadwicke Jenkins. Incremental learning of subtasks from unsegmented demonstration. *In Intelligent Robots and Systems*

- (IROS), 2010 IEEE/RSJ International Conference on, pages 261–266. IEEE, 2010.
- [63] Maya Cakmak, Crystal Chao, and Andrea L Thomaz. Designing interactions for robot active learners. *IEEE Transactions on Autonomous Mental Development*, 2(2):108–118, 2010.
- [64] Peter Ford Dominey, Manuel Alvarez, Bin Gao, Marc Jeambrun, Anne Cheylus, Alfredo Weitzenfeld, Adrian Martinez, and Antonio Medrano. Robot command, interrogation and teaching via social interaction. *In 5th IEEE-RAS International Conference on Humanoid Robots, 2005.*, pages 475–480. IEEE, 2005.
- [65] Verena V Hafner and Frédéric Kaplan. Learning to interpret pointing gestures: experiments with four-legged autonomous robots. *In Biomimetic neural learning for intelligent robots*, pages 225–234. Springer, 2005.
- [66] R Zoliner, Michael Pardowitz, Steffen Knoop, and Rüdiger Dillmann. Towards cognitive robots: Building hierarchical task representations of manipulations from human demonstration. *In Proceedings of the 2005 IEEE International Conference On Robotics and Automation*, pages 1535–1540. IEEE, 2005.
- [67] Eric L Sauser, Brenna D Argall, Giorgio Metta, and Aude G Billard. Iterative learning of grasp adaptation through human corrections. *Robotics and Autonomous Systems*, 60(1):55–71, 2012.

- [68] Ian Lenz, Honglak Lee, and Ashutosh Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34(4-5):705–724, 2015.
- [69] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.
- [70] Maria Makarov Wam robotic arm. <http://www.barrett.com/products-arm-articles.htm>. Accessed: 2016-09-11.
- [71] Morgan Quigley, Ken Conley, Brian Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, and Andrew Y Ng. Ros: an open-source robot operating system. In ICRA workshop on open source software, volume 3, page 5, 2009.
- [72] CMU Sphinx Open Source Toolkit for Speech Recognition project by carnegie mellon university. <http://cmusphinx.sourceforge.net/>. Accessed: 2016-06-25.
- [73] Alan W. Black and Paul A. Taylor. The Festival Speech Synthesis System: System documentation. Technical Report HCRC/TR-83, Human Communication Research Centre, University of Edinburgh, Scotland, UK, 1997. Available at <http://www.cstr.ed.ac.uk/projects/festival.html>.
- [74] Camilo Perez Quintero, Romeo Tatsambon Fomena, Azad Shademan, Nina Wolleb, Travis Dick, and Martin Jagersand. Sepo: Selecting by pointing as an intuitive human-robot command interface. In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pages 1166–1171. IEEE, 2013.

- [75] Zdenek Kalal, Jiri Matas, and Krystian Mikolajczyk. Pn learning: Bootstrapping binary classifiers by structural constraints. *In Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 49–56. IEEE, 2010.