

RWTH Aachen University
Software Engineering Group

Deep Learning in Autonomous Driving Direct Perception Approach

Seminar Paper

presented by

Bergerbusch, Timo

1st Examiner: Prof. Dr. B. Rumpe

2nd Examiner:

Advisor: Evgeny Kusmenko

The present work was submitted to the Chair of Software Engineering

Aachen, June 21, 2018

The present translation is for your convenience only.
Only the German version is legally binding.

Statutory Declaration in Lieu of an Oath

Last Name, First Name

Matriculation No. (optional)

I hereby declare in lieu of an oath that I have completed the present Bachelor's
thesis/Master's thesis* entitled

independently and without illegitimate assistance from third parties. I have use no other than
the specified sources and aids. In case that the thesis is additionally submitted in an
electronic format, I declare that the written and electronic versions are fully identical. The
thesis has not been submitted to any examination body in this, or similar, form.

Location/City, Date

The German version has to be signed.

Signature

*Please delete as appropriate

Official Notification:

Para. 156 StGB (German Criminal Code): False Statutory Declarations

Whosoever before a public authority competent to administer statutory declarations falsely makes such a
declaration or falsely testifies while referring to such a declaration shall be liable to imprisonment not exceeding
three years or a fine.

Para. 161 StGB (German Criminal Code): False Statutory Declarations Due to Negligence

(1) If a person commits one of the offences listed in sections 154 to 156 negligently the penalty shall be
imprisonment not exceeding one year or a fine.

(2) The offender shall be exempt from liability if he or she corrects their false testimony in time. The provisions of
section 158 (2) and (3) shall apply accordingly.

I have read and understood the above official notification:

City, Date

The German version has to be signed.

Signature

Abstract

Deep learning (DL) came into the foreground of research after some stunning breakthroughs. Especially in image-recognition with novel Convolutional Neural Net (CNN) architectures, which made training on millions of images feasible on off-the-shelf GPUs. These advances are of high relevance to autonomous driving since much of the information for driving is visual.

In this paper we will give a general introduction to the topic of neural networks (NNs) and then state the specialties defining CNNs, which is a subclass of NNs designed for the task of image analysis.

A multitude of frameworks and libraries has been created for DL using all kinds of languages. In case of neural network architectures, it is beneficial to take advantage of domain-specific languages (DSLs), to facilitate the development of new NNs.

Because of that, we will compare the different DSLs CNNArchLang, Caffe, Caffe2, and MxNet, based on the factors of usability, scope of functionality, also regarding training possibilities, and the integration on a subject.

Regarding the task of autonomous driving, we distinguish between the three main paradigms currently used and researched for autonomous driving agents: mediated perception, behavior reflex, and direct perception.

We will analyze a CNN using the direct perception approach and compare it to state-of-the-art implementations of the other paradigms, training on a simulator TORCS or the KITTI database. Furthermore, we state scenarios probably causing problems for the direct perception approach.

Finally, we create an overview of the mentioned languages with a table, which states the functionalities and properties in a nutshell.

Contents

1	Preliminaries	1
1.1	Neural Networks	1
1.2	Convolutional Neural Network (CNN)	2
1.3	AlexNet	4
2	Deep Learning Languages	5
2.1	CNNArchLang	5
2.1.1	Predefined Layers and Functions	6
2.2	Caffe	6
2.3	Caffe2	7
2.4	MxNet	8
3	Comparison: CNNArchLang & Caffe	11
3.1	Implementation	11
3.1.1	Caffe	11
3.1.2	CNNArchLang	12
3.1.3	Comparison	12
3.2	Training	12
3.2.1	KITTI Dataset	13
3.2.2	TORCS	13
4	Available Deep Learning Approaches	15
4.1	Mediated Perception	15
4.2	Behavior Reflex	16
4.3	Direct Perception	17
5	Evaluation of Direct Perception	21

6	Conclusion	23
7	Appendix	25
	Bibliography	28

Chapter 1

Preliminaries

The field of autonomous driving agents has rapidly become a key factor in modern car manufacturing. Current research topic rises the agents from parking or lane keeping assistants to fully autonomous driving agents obeying the traffic rules and having the ability to react to the volatile environment in a reasonable way.

For that, machine learning techniques have proven themselves as an essential part. In order to fulfill the security standards and create a sophisticated agent, it has to be trained on hundred-thousands of scenarios, each having a large set of data attached, for example sensor and image data of multiple cameras.

The approach of *Convolutional Neural Networks* (CNNs) have been proven to be powerful enough to handle this amount of training iterations with a huge number of input variables, while maintaining a large learning capacity. [KSH12]

1.1 Neural Networks

The neural networks are a construct, adapted from biological processes. Its general structure is uncomplicated, but has immense expressiveness. The complete spectrum of what can be modeled by neural networks is currently not known.

A neural net is made out of neurons and has one very basic function:

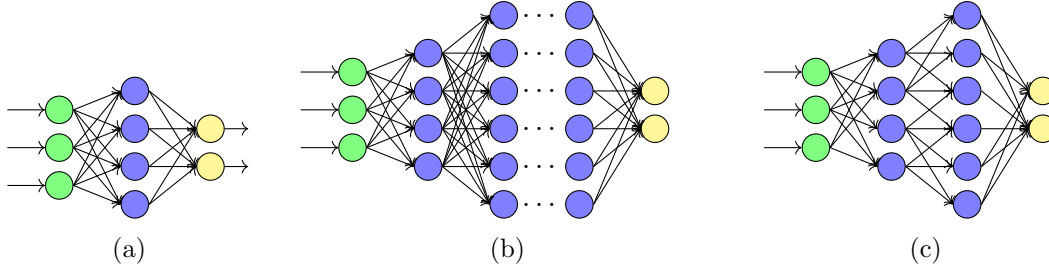
It takes a fixed number $n \in \mathbb{N}$ of the incoming values x_i , where the vector $(x_0, \dots, x_n)^T$ is called a tensor, and multiplies them each with a specific weight w_i , where $0 \leq i \leq n$. Also every neuron contains a bias b , which is a general value subtracted from the sum, so $\sum_{i=0}^n (x_i \cdot w_i) - b$.

Often one applies an activation function to fix the value between 0 and 1. Such a function would be for example the sigmoid function [VS16], which is a non-linear functions. Without using a non-linear function one gets restricted to linear regression and therefore reduces their ability to model more complex functions. [GB10] This non-linear normalized value is forwarded to the neurons of the next layer.

Those neurons are ordered in different layers, visualized in Figure 1.1a.

1. input layer (green):

This layer gets fed with the input values of the problem, which can be for example sensor data or pixel color values.



2. hidden layer (blue):

The hidden layer consists of neurons receiving the values from the previous layer, while not being obliged to have the same number of neurons (c.f. Figure 1.1a). Different architectures can be distinguished to be deep (c.f. Figure 1.1b). This means, that there are multiple hidden layers of neurons instead of just one. The number of neurons a layer contains is again independent of the previous layers. Also for the hidden layers, it is possible to be not fully connected (c.f. Figure 1.1c). Thus some neurons do not forward their value to every neuron of the next layer. There is no rule dictating the best architecture, considering number of layers, neurons per layer or the connectivity.

3. output layer (yellow):

The output neurons produce the value of the neural network. Depending on the neural networks purpose it can be for example a confidence value of a classification, like recognizing a stop sign, or the value of changing the steering wheel angle.

In order to train a neural network one has to define the behavior the neural network should have. In an image classification example concerning traffic signs, one should know what the correct class of a given image of a sign is, i.e. a stop sign. Those are called labels. A neural network can then be trained by giving it values for the input layer and comparing the values of the output layer with the solutions it should have resulted in. The difference can then be checked. Such a difference can be simply `true/false` or a value indicating how big the difference is.

In the example of signs: a classification of a “speed limit 70”-sign as a “speed limit 50”-sign is still wrong, but not as bad as the classification as a “stop”-sign.

Using this difference values, the neural network can use backpropagation [Kin94] to adjust the weights w_i and biases b to improve the output iteratively.

Further information about the underlying training algorithms like gradient descent, newtons method or conjugate gradient [Roj13] is not given here in order to keep the paper in a justifiable length.

1.2 Convolutional Neural Network (CNN)

A CNN is a special class of deep feed-forward neural networks. One of the main design goals of a CNN is, that they require a minimal amount of preprocessing. This is an important aspect, because they are often fed with images. Preprocessing high resolution images is very costly in terms of computational time. In the context of autonomous driving the time is even more crucial, since the driving agent needs to be able to react to spontaneous events.

Like most parts of neural networks, also the CNNs are inspired by biological processes. It is mainly based on the connectivity pattern of an animals visual cortex, where special neurons respond only to stimuli of their receptive field, represented as rectangles lying in the image. Partial overlapping guarantees a complete coverage of the field of view. Those rectangles are often called kernel, feature or filters.[MMM03]

The four basic layer types, which are mostly used for CNNs are:

1. fully-connected layers:

equal to the fully connected layers described for neural networks (c.f. Section 1.1)

2. convolutional layers:

A convolutional layer compares the given image with a list of features, the net derived, to indicate a certain item. In this context it could be an other car driving ahead of the host. Those features are not a-priori given, but learned by the CNN. It is simple to understand using an example:

$$\begin{array}{|c|c|c|c|} \hline 7 & 6 & 5 & 5 \\ \hline 6 & 9 & 1 & 3 \\ \hline 4 & 1 & 2 & 8 \\ \hline 2 & 9 & 3 & 1 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline -1 & 2 \\ \hline 2 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 17 & 22 & 7 \\ \hline 20 & -5 & 9 \\ \hline 2 & 21 & 20 \\ \hline \end{array}$$

The left most is the input and the second one is the kernel. The result is computed as a piece-wise multiplication and then adding them up. The convolution can differ in kernel size and stride, which denotes the movement of the kernel. In the example we have a kernel of size 2×2 and a stride of 1.

So for example (red case): $7 \cdot (-1) + 6 \cdot 2 + 6 \cdot 2 + 4 \cdot 0 = 17$

The number now represents the similarity or likelihood of the feature to be in that position. Edge detection can be applied easily using such a layer. [HW68]

3. rectified linear units (ReLU):

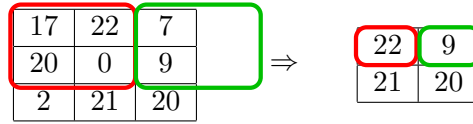
Those ReLUs perform a normalization step, by applying the function $\max(0, x)$ to every element x . So for the example:

$$\begin{array}{|c|c|c|} \hline 17 & 22 & 7 \\ \hline 20 & -5 & 9 \\ \hline 2 & 21 & 20 \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|c|} \hline 17 & 22 & 7 \\ \hline 20 & 0 & 9 \\ \hline 2 & 21 & 20 \\ \hline \end{array}$$

The main purpose of such a ReLU is to create the prerequisite of mathematical functions, only accepting non-negative values. Often researchers tend to use the function $\ln(1 + e^x)$ in order to have a differentiable approximation of ReLU.

4. pooling layers:

The pooling layer reduces the input by having a windows size, which is typically 2 or 3, and a stride, which is typically 2, and map the values within the window to one single value. In our example this is done using the maximum value.



For the red case, the value 22 is the maximum of the window containing $\{17, 22, 20, 0\}$. For the green case, the 9 calculated is the maximum of the window $\{7, 9\}$ since through a stride of 2 the window is partially not contained in the input.

The application of convolution and pooling has the advantage that it reduces the effort to train a CNN. The weights and biases of neurons of each receptive field are equal. This is reasonable since for example a speed limit road sign should be identified independent, whether it is located next to the road, like on a normal road, or above the road, like on an highway. [L⁺15]

1.3 AlexNet

The *AlexNet* was one of the best performing CNN architectures in 2012 and is still performing good enough to be used for various problems. It was originally trained on the ImageNet subsets of ILSVRC-2010 and ILSVRC-2012 ¹ and became famous, because of its result being way ahead of all other competitors.

A highly optimized GPU implementation of this architecture combined with innovative features is publicly available. Those features lead to performance improvement and shorter training time.[KSH12]

An important note is that the original test is dated back to 2012 and therefore was used with an overall GPU memory of 6 GB, with which training took about six days. With modern hardware like new GPUs, SLI usage (having multiple GPUs) or even clusters, the training can be done faster, or the model can be trained with more data to improve performance. The improvement caused only by hardware can be roughly grasped through [SCE⁺17].

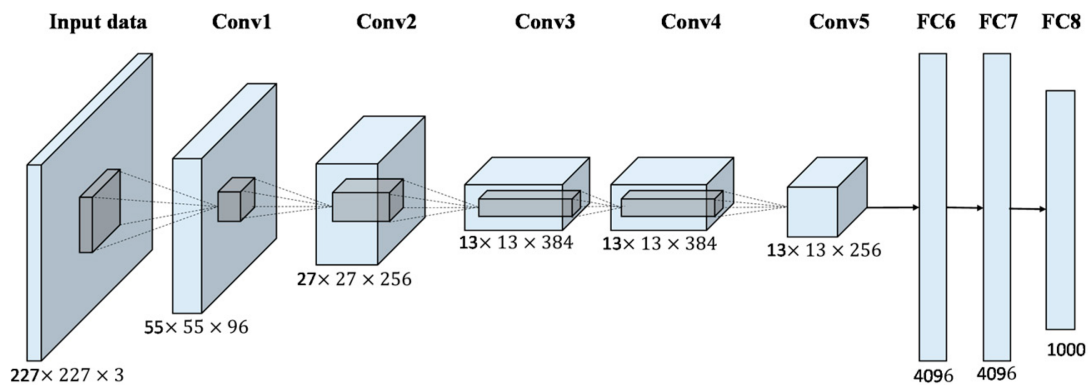


Figure 1.2: The AlexNet-architecture. It consists of 5 convolutional layers (at the beginning), and three fully-connected layers (at the end). [HZCZ17] [KSH12]

¹Further information: <http://www.image-net.org/challenges/LSVRC/>

Chapter 2

Deep Learning Languages

Constructing a CNN from scratch in any typical language like Java, C++, or Python is very elaborately and has a high error potential. Even libraries in any such language often encounter the problem of over-complication due to their own style, syntactical and semantic architecture. Therefore there is a need for specialized languages.

The Deep Learning Languages (DLLs) are part of the Domain Specific Languages (DSL). Their main goal is to provide an easy to understand and less verbose way of describing a CNN with its different layers and connections, without losing expressiveness. Also one wants to have simple build presets, maybe even already trained.

For that we consider four deep learning languages and analyse them on the previously mentioned properties.

2.1 CNNArchLang

The whole description is based on [TvWH17] and especially [Tim18]

One language for modeling CNNs is CNNArchLang. This language is developed at the Chair of Software Engineering, especially Thomas Michael Timmermanns, at the RWTH Aachen University and part of the MontiCAR language family. The main purpose of its creation is the necessity of special properties not given by other CNN-languages: *C&C integration* and *type-safe component interface*. Its basic structure is very similar to Python to improve understanding, based on familiarity, and have an equal non-typed syntax.

One advantage of CNNArchLang is that it is designed to be very simplistic and less verbose than most other CNN-languages. It does so by moving, from defining a CNN by every single neuron, to the definition via layers only used a specialized notation, explained in the next paragraph. For that specific purpose many layers are already defined (c.f. Section 2.1.1). New layers can be constructed by combining predefined layers.

This slightly reduces the expressiveness, since the possibility of performing computations on single tensors is lost. Such low-level operations are used extremely rarely and are not a drastic disadvantage.

In contrast to other languages for deep learning, CNNArchLang does not denote the connections of layers directly, but tries to model the data flow through the network. For that specific task it contains two main operators:

->: Serial Connection:

This orders two elements sequentially. This means it denotes the first elements output as the second elements input.

|: Parallelization:

This allows the split of the network into separate data streams, which can be computed in parallel.

Since serial connection has the higher precedence one has to use brackets. Also to merge the splitted streams, created by |, one can use the operators: `Concatenate()`, `Add()` and `Get(index)`.

2.1.1 Predefined Layers and Functions

Different CNNs often use a similar basic set of layers, but arrange them differently. For that purpose there are some layers already defined by `CNNArchLang` to simplify the usage. There is for example the `FullyConnected`-layer with parameters for the number of units within and whether they should use a bias value (c.f. Section 1.1), the `Convolutional-2D`-layer with parameters for the kernel size, number of filters, the stride, padding and the usage of biases. Further information on any of these parameters or other predefined layers can be found in Section 1.2 and [TvWH17].

Also there are already defined functions like the `Sigmoid`, `Softmax`, `Tanh` or `ReLU`. One important aspect is that every argument has to be named.

The distinction, that `CNNArchLang` is not a framework itself, is very important. It is used to create the code to function in the `MxNet` (see Section 2.4).

2.2 Caffe

Caffe is a full deep learning framework, created by Yangqing Jia during his PhD at UC Berkeley. It is a framework specially build to deal with multimedia input formats using state-of-the-art deep learning techniques. It comes as a BSD-licensed C++ library offering Python and MATLAB bindings. One of the reasons, why it is so well known and frequently used, is because of its design based on expressiveness, speed and modularity.

Using Nvidias Deep Neural Network library `cuDNN` as a wrapper of the `CUDA` functionality, Caffe can use the GPU in order to process even faster and learn in a rate of 40 million pictures per day. The possibility of using multiple GPUs in SLI is not stated and therefore not taken into account. [She14]

However, the possibility of using a cloud system is mentioned. But whether it is a simple decentralization, or the possibility to train the network using the combined power of multiple computers to train, is not mentioned. [JSD⁺14]

Caffe tries to improve its readability by abstracting from the actual implementation using a graph-oriented way of network definition. For that, Caffe uses two elements to represent the network:

- **Blob:**

A 4-dimensional array storing data, like images, parameters or parameter updates. These blobs are the communication between layers.

```

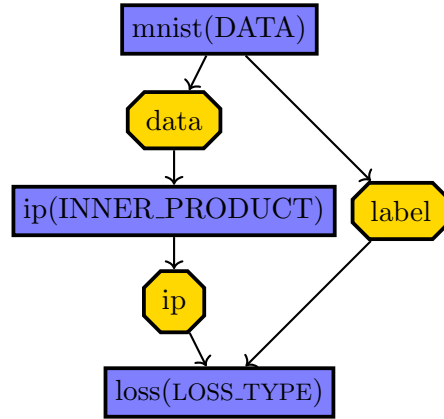
name: "loss-net"
layers {
  name: "mnist"
  type: DATA
  top: "data"
  top: "label"
  data_param { ### }
}

layers {
  name: "ip"
  type: CONVOLUTIONAL
  bottom: "data"
  top: "conv"
  convolutional_param { ### }
}

layers {
  name: "loss"
  type: SOFTMAX_LOSS
  bottom: "ip"
  bottom: "label"
  top: "loss"
}

```

(a)



(b)

Figure 2.1: An example of a Softmax loss network. Blue boxes are the different layers and the yellow boxes are the blobs. Note that in (a) the parameters are not mentioned since they do not add value to this example. There are things defined like kernel size, image scaling or image origin.

- **Layers:**

A Layer as described in Section 1.1 and Section 1.2.

The whole model gets saved as a Google Protocol Buffer, which is a language-neutral structurization, with important key features, like size-minimal and efficient serialization, and efficient interfaces for C++ and Python. [JSD⁺14] [Var08] An example of a graph and its corresponding Protocol Buffer representation, written in prototxt, are given in Figure 2.1.

The biggest advantage of Caffe is the huge community providing a large set of presets, like layers or pre-trained nets, and also a huge forum to ask other users about problems regarding ones project. This is such an advantage, because for every module it is guaranteed to pass a coverage test. [JSD⁺14]

Caffe is written using as plain text named prototxt and therefore can be run via command line. One problem mentioned with prototxt is, that even though there are many nets already defined, the creation of new nets often is highly verbose and repetitive. There are no shortenings. For example [Tim18] mentioned an example net written in CNNArchLang with 36 lines and in Caffe with 6700 lines. This is also due to the fact that, even if a layer can be constructed as a composition of existing layers, one often has to define the forward-, backward-propagation and gradient-updates.

2.3 Caffe2

The framework Caffe2 is the successor of Caffe. Caffe2 is developed by Facebook and its current main usage is the phrase-wise translation in social networks. Keeping the modularity of Caffe in mind, Caffe2 is also designed so that it can be up-scaled as well as mobile deployed. Also Caffe2 is designed in such a fashion that it can easily adapt to drastic changes like quantized computing. [Caf18]

Since the whole architecture is rewritten from scratch and regarding its roughly one year of existence, the library performs relatively well, but does not have the impact to outperform Caffe [Hei17]. One downside of the rewriting is that huge parts of the framework are not sufficiently documented. Caffe2 tries to improve, but still has large holes within the documentation. [Tim18]

Caffe2 offers programs, which allow the user to convert Caffe and PyTorch models to Caffe2. This makes the switching to Caffe2 much easier, since the users do not have to rewrite their models. [Caf18]

The main difference between Caffe and Caffe2 in terms of designing a neural network is that in Caffe2 the user uses `Operators` as the basic units instead of layers. Even though they are similar to the layers of Caffe, they are more flexible and adaptable. Partly based on the popularity of Caffe, Caffe2 also has a huge community and a large set of preset `Operators`, which can be used. [Caf18]

2.4 MxNet

Another framework often mentioned in research is the MxNet. This deep learning framework is part of the Apache Software Foundation. Also it is said to be “Amazons deep learning framework of choice” [Yeg16] and featured to be a preset on the Amazon Web Services (AWS). [CL]

MxNet tries to combine the advantages of imperative frameworks like numpy or MATLAB with the advantages of declarative frameworks, like Caffe, Caffe2, or Tensorflow.

The advantages of imperative and declarative approaches can best be understood using an example. Let the example be to compute: $a = b + c$

imperative:

Procedure: check the ability of b and c to be added. If so strictly compute the sum and declare a as the same type as b and c

Advantage: very straightforward, works well with typical structures, debugger and third party libraries

Usefull for: natural parameter updates and interactive debugging

declarative:

Procedure: compute the computation graph and store the values of b and c as data bindings

Advantage: perform computation as late as possible, leading to good optimization possibilities

Usefull for: specifying computation structure, optimization

By combining both approaches MxNet is able to provide a superset programming interface compared to Caffe. [CLL⁺15]

Also MxNet is able to reduce memory to a minimum by performing everything they can in-place, and free or reuse as fast as possible. Thus, the memory usage of MxNet is outperforming Caffe. [CLL⁺15] A new benchmark with Caffe2 has not been done yet.

Another very big upside of MxNet is the possibility to use not only multiple GPUs in an SLI connection, but also multiple computers or even servers to train a neural network simultaneously. This results in an outstanding scalability. [CLL⁺15]

Similar to Caffe2, MxNet also allows the deployment of trained models to low-end devices using Amalgamation (c.f. [MxNa]) or the AWS.

Due to the fact that MxNet has found its way into Apache Incubator, and therefore it is an Open-Source project, the creation of additional functions and nets is quite simple and is not bound to a given preset. Thanks to the community, also a variety of nets constructed, from which some are already pre-trained, are free to use. [MxNb]

Chapter 3

Comparison: CNNArchLang & Caffe

In this chapter we want to compare the AlexNet (c.f. Section 1.3) implemented using CNNArchLang and Caffe. This net is used in the direct perception approach (c.f. Section 4.3) and therefore crucial for its performance. For both we take an in depth look at the predefined architecture by their respective language. We do that, since both implementations are done by language experts and build to be as efficient and precise in the language as possible.

Further, we state the currently most used methods to train an autonomous driving agent. Those training methods are also used for the approaches in Chapter 4. The important properties required are stated.

3.1 Implementation

3.1.1 Caffe

The implementation of the AlexNet using Caffe is given partly in Figure 7.1. The whole net has a total number of 284 lines. Obviously the code was written very verbosely. Every layer has to be explicitly specified, even if they have a very similar structure to a previous layer.

For example comparing the pooling layer “pool1” from line 51 to 61 and the pooling layer “pool2” from line 99 to 109 in Figure 3.1:

<pre> ⋮ 51 layer { 52 name: "pool1" 53 type: "Pooling" 54 bottom: "norm1" 55 top: "pool1" 56 pooling_param { 57 pool: MAX 58 kernel_size: 3 59 stride: 2 60 } 61 } ⋮ </pre>	<pre> ⋮ 99 layer { 100 name: "pool2" 101 type: "Pooling" 102 bottom: "norm2" 103 top: "pool2" 104 pooling_param { 105 pool: MAX 106 kernel_size: 3 107 stride: 2 108 } 109 } ⋮ </pre>
(a) "pool1"	(b) "pool2"

Figure 3.1

The lines first 3 and last 6 lines are completely similar. The other 2 lines are just different regarding the name of the incoming and outgoing connections. Creating a huge and deep net would lead to an enormously large description file.

3.1.2 CNNArchLang

The implementation of the AlexNet using CNNArchLang can be seen in Figure 7.2. The complete script has 43 lines and defines the same net construction as the 284 line definition in Caffe. This shows the efficient language design used in the creation of CNNArchLang (c.f. Section 2.1).

The two pooling operations of Figure 3.1 can be located in line 32 using the Python like syntax of definition and the sequential connection `->` (c.f. Section 2.1).

Using those techniques CNNArchLang is able to write even complex neural network using a few lines of code. This and the syntax itself create an easy to read program.

3.1.3 Comparison

Due to complications regarding the CNNArchLang-inclusion in an existing program infrastructure, as well as Caffe, currently having issues with their build-script there is no possibility to directly compare times taken for training the neural network.

Nevertheless based on the usage of MxNet, by CNNArchLang (c.f. Section 3.1.2 and Chapter 6), it is suspected to be much faster than the Caffe approach.

About the effectiveness and other performance indicators, such as error rate, there can not be any profound reasoning without an actual implementation and testing.

3.2 Training

In order to train a CNN, independent of the underlying approach (c.f. Chapter 2), one has to obtain a huge database of input images and the labels, i.e. actions or values the CNN should have computed. For autonomous driving the training has to be rigorous. Otherwise the car driven by the agent will take damage by just slight changes of the circumstances. Also different scenarios have to be trained. Only training the driving on a road without

other cars and simply following the lane is a very disparate task compared to overtaking a slower car.

For that, the following sources of such databases are currently state-of-the-art.

3.2.1 KITTI Dataset

The KITTI dataset is a 6 hour recoding by the Karlsruhe Institute of Technology. They mounted various cameras and laser scanners on a VW Passat and drove around the german city Karlsruhe.

During those hours they collected a total amount of 180 GB of data. This data includes images, in different channels from the drivers point of view, sensor data of distances, steering angle, acceleration/braking, current speed, GPS coordinates, and others.

While other test sets are often developed using a very specific setup for a corresponding approach, the KITTI dataset has such a high variety of data captured, resulting in many appearances in other scientific papers. It has become one of the default datasets to compare different approaches on. [GLSU13]

3.2.2 TORCS

A game called **The Open Racing Car Simulator** or short TORCS is a racing game specially designed for artificial intelligence research. It is designed to be modular in order to retrieve every kind of data one needs for their approach. It also offers a documented API, to create for example a driving agent.

The possibility to collect any kind of data is what makes this game so popular in current research. The advantage over the KITTI dataset is that, if one needs a very specific value not included within the KITTI dataset, the approach can not be trained with it. Neglecting, whether such a value is meaningful in terms of the ability to collect it during real driving, there would be an other specialized set collecting only the data this specific approach needs. But this specialized set may lack a value, which is required in order to compare it with an another approach. So there is no common training set, on which a comparison could base.

The disadvantage of the training via TORCS compared to training via KITTI is, that TORCS only uses artificial images, other artificial agents, and is completely exempt from any kind of data noise, like rain disturbing sensor data or sunlight blending the camera. Also other cars behave different in a game than in real life. The KITTI includes some of those problems. To which extend is arguable. [WEG⁺00]

Chapter 4

Available Deep Learning Approaches

There are various approaches of autonomous driving agents, making a variety of assumptions and differ in numerous options. But they can be mostly categorized into two major groups of approaches: mediated perception approaches and behavior reflex approaches. [CSKX15]

In this paper we further analyse a suggested third group, called direct perception, which can be traced to [Gib79] in the mid 50's, but was sharply criticized by researchers of the other two groups, i.e. in [Ull80].

All these three groups differ in the way of interpreting the given sensor data and whether or not to create a some what bigger picture based on consequent data.

4.1 Mediated Perception

The mediated perception approach is a multi-component continuous process. Every component recognizes specific aspects for driving. For example traffic signs, lanes and other cars. Those components are then combined into one single world state representing the cars surrounding, based on the sensor data. [GLSU13]

These world states are 3D models of the current world. Cars are identified using a classifier and then often surrounded by a 3D bounding box. An example can be seen in Figure 4.1. By comparing different frames generated, one can estimate the speed and distance to those objects and derive an A.I. based precedence behavior. [GLSU13][CSKX15]

The often stated problems with such approaches are, that computing such a scene is costly in terms of computation time. Some information is irrelevant, redundant or even misleading due to inaccuracy of sensors. To perform a right turn the sensor information of the distance to a car left behind me is irrelevant, but becomes very important when taking a left turn.

Additionally, many of the subtasks are still open research topics themselves. For example a reliable lane detection throughout various weather conditions or even a new road not having any drawn lines yet. [Aly08]

Also mediated perception approaches require very detailed information up front, like up-to-date maps.



Figure 4.1: An example of a scene using 3D bounding boxes and path prediction.

The approach of mediated perception is a reasonable and very sturdy way of handling such a complex task, but has its drawbacks regarding computational time and additional knowledge.

4.2 Behavior Reflex

The behavior reflex approach of constructing a reliable autonomous driving agent can be dated back to 1989, where researchers tried to directly map a single frame to a decision of a steering angle. For such approaches a quite simple neural network was created.

The network ALVINN consisted of a single hidden layer, used back-propagation and is fed by two cameras: a 30×32 pixel video and a 8×32 pixel range finder retina. The input neurons fired depending on the blue color band of its pixel, because it is believed to infer the highest contrast between road and non-road. The difference in color of road and non-road was fed back to the network. The bias (activation level) is proportional to the proximity of the corresponding area, based on the likelihood and importance of having road in particular fields of the image.[Pom89]

For example, having recognized that the road abruptly ends right in front of the car is more important than recognizing that there is a road in the top left corner.

Such systems, even though they are simple compared to the in Section 4.1 mentioned mediated perception approaches, have been proven to have the capability to perform simple tasks, like lane keeping perfectly. It can elegantly be trained by having a human drive a car with the cameras equipped and forward the images to the neural network and adding the current steering angles as a label.[CSKX15]

The problem with behavior reflex approaches is, that they reach their limits very early, when adding more complex scenarios. Having simple alternations to the trained scenarios,

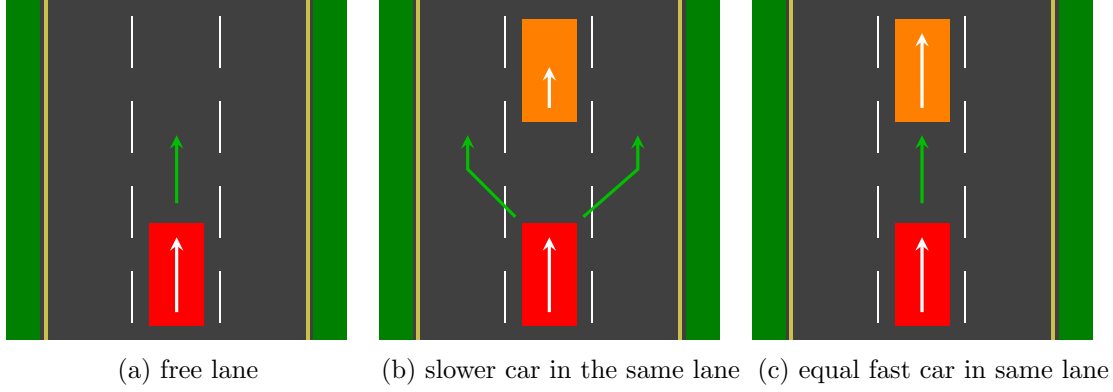


Figure 4.2: The 3 scenarios causing problems with behavior reflex approaches. The red block is the agent, the orange block the other car, the white arrows indicate the velocity and the green arrows the logically deduced behaviors.

which enforce a different behavior, is very hard to train to such a neural network. For example comparing a simple straight 3 lane road with the car in the middle, as sketched in Figure 4.2. The system is confidently able to hold the angle and make small adjustments to stay in the lane (Figure 4.2a). But what if on the same road there is an other car in the middle lane in front of the agent, which is slower? Having quite the same input the system would have to overtake the car left or right (considering an American highway) (Figure 4.2b). Now also considering a car in front, which has the same speed. One can simply stay in the lane (Figure 4.2c). This maneuver is very hard to train to a simple neural network like ALVINN.

4.3 Direct Perception

The direct perception is the third group of approaches, which can be dated back to the 1954 and was initially mainly researched by James J. Gibson. [Gib54] The approach is based on analyzing a picture, not simply deducing a steering angle, or velocity change, like the behavior reflex approaches (cf. Section 4.2), but also performing further computation without parsing it into a 3D world state model like the mediated perception approaches (cf. Section 4.1). [CSKX15]

So it is a third paradigm, which can be interpreted as a hybrid of the two other paradigms. The approach tries to identify only the meaningful affordance indicators and makes a decision based on those parameters.

The CNN gets trained to make a reasonable assumption about the distance of other vehicles, their velocity and lane marks. These assumptions are further called affordance indicators.

We further consider a design based on [CSKX15] and their way of training.

The original paper [CSKX15] stated a total number of 13 indicators separated into two states to be sufficient for their design of a highway only autonomous agent. The states are: in line driving (following the lane) and on line driving (changing lanes). The values themselves can be categorized as: preceding car distances, distances to the lane markers and the steering angle. The indicators and their affiliation to the states are listed in Figure 4.3 and visualized in Figure 4.4.

-
- always:
- 1) angle: angle between the car's heading and the tangent of the road
 - i) "in lane system", when driving in the lane:
 - 2) toMarking LL: distance to the left lane marking of the left lane
 - 3) toMarking ML: distance to the left lane marking of the current lane
 - 4) toMarking MR: distance to the right lane marking of the current lane
 - 5) toMarking RR: distance to the right lane marking of the right lane
 - 6) dist LL: distance to the preceding car in the left lane
 - 7) dist MM: distance to the preceding car in the current lane
 - 8) dist RR: distance to the preceding car in the right lane
 - ii) "on marking system", when driving on the lane marking:
 - 9) toMarking L: distance to the left lane marking
 - 10) toMarking M: distance to the central lane marking
 - 11) toMarking R: distance to the right lane marking
 - 12) dist L: distance to the preceding car in the left lane
 - 13) dist R: distance to the preceding car in the right lane
-

Figure 4.3: The affordance indicators and their affiliation states

Based on the current state, all affordance indicators of the other state are not used, since the other state is defined to be inactive.

In order to identify the current state of the host car, every state has their respective region, where they are active with an overlapping region for smooth transitioning.

Training is done by feeding the CNN with images of a current driving scenario as input and give estimated values to the mentioned affordance indicators active in the current state. Those values get forwarded to a controller dealing with the car. In the original approach the training was done via TORCS (c.f. Section 3.2.2).

The controller is setup to take the affordance values and compute a decision from it in an imperative way. For example the steering can be computed as

$$steerCmd = C \cdot (angle - \frac{dist_center}{road_width})$$

where *dist_center* is the center of the current state (line we are on, or middle of the lane), a coefficient *C* suited for the current conditions, like weather or speed, and *angle* $\in [-\pi, \pi]$ as the affordance indicator.

Also using the controller to compute the accelerating and braking one can adjust the desired speed (*desired_speed*) based on the affordance indicators. The *desired_speed* is set to the speed the agent wants to drive and can drop, if a drastic steering motion is considered or accelerate to overtake a slower car.

In a one lane scenario with a preceding car driving slower than the *desired_speed* there is no space to overtake. For that the controller has the velocity control car-following model:

$$v(t) = v_{max}(1 - e^{-\frac{c \cdot dist(t)}{v_{max}} - d})$$

where v_{max} is the maximal speed the agent is allowed to drive, *c* and *d* are coefficients specific to external conditions, like a wet lane or the cars potential of accelerating and braking, and *dist(t)* is the distance to the preceding car. This distance is given through the trained CNN.

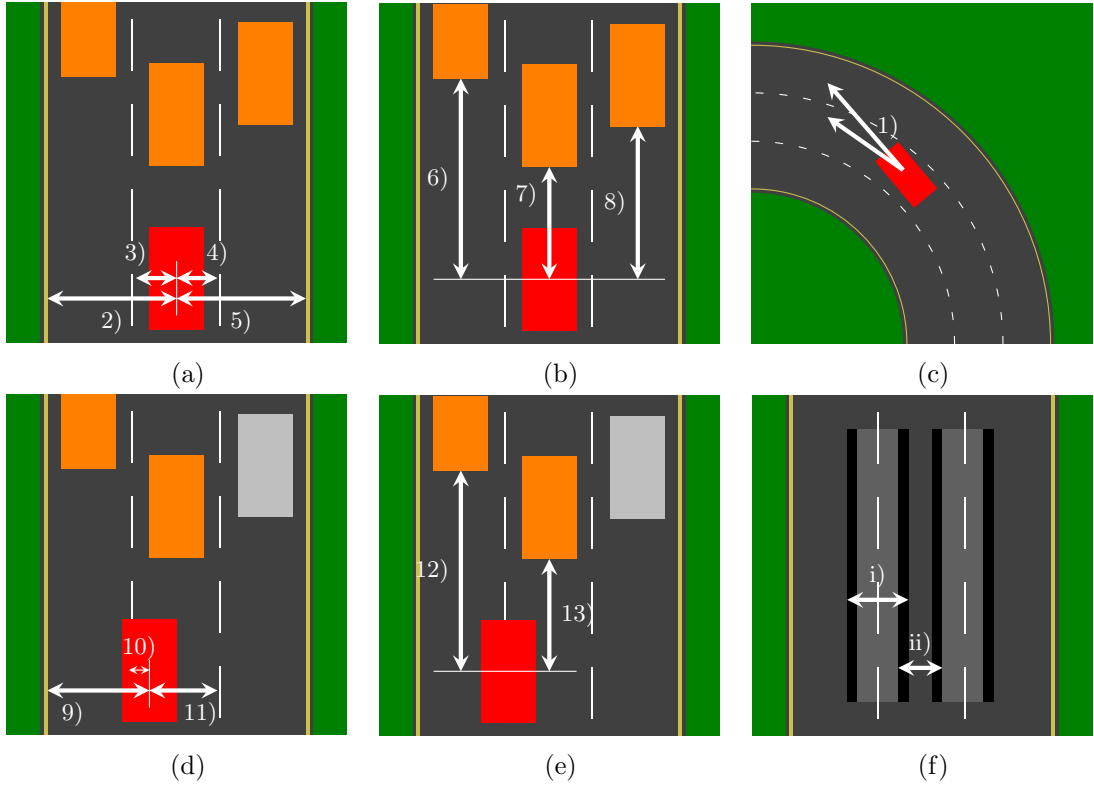


Figure 4.4: Visualization of the affordance indicators listed in Figure 4.3. The red block is the agent, the orange blocks are other cars and the white arrows indicate the velocity. Gray cars are neglected, because of the currently active system.

Chapter 5

Evaluation of Direct Perception

In this chapter we take a look back at the direct perception approach and evaluate its performance. For that we compare the, via TORCS (c.f. Section 3.2.2), trained version to a behavior reflex approach (c.f. Section 4.2). In addition to that, there is also a comparison of the direct perception approach, trained via KITTI (c.f. Section 3.2.1), to a mediated perception approach stated in [GLW⁺14].

We base the comparison on the running examples in [CSKX18], which are written using Caffe.¹

A comparison with the direct perception written in CNNArchLang is highly desired, but currently not possible due to the previously mentioned facts of incompatibility.

Finally, we take a look at possible scenarios causing problems, which may emerge.

The directed perception as stated in [CSKX15], and discussed in Section 4.3, is designed to handle highway driving tasks, such as driving in a lane, overtaking slower cars, detecting the lane configuration, and breaking to avoid a collision.

The behavior reflex approach was able to follow empty lanes perfectly, but was completely unable to have an acceptable behavior considering other cars. Neither a sufficient speed regulation nor the task of staying in lane was observable. The agent left the track various times and had multiple collisions.

On the other hand the direct perception approach manages to change lanes smoothly, avoid collisions and stay in lane. Due to the speed regulation in the controller, stated in Section 4.3, the agent is able to perform an emergency brake, if necessary. So considering those scenarios the direct perception outperforms the behavior reflex approach.

In order to compare the direct perception with the state-of-the-art mediated perception approach, the training is done using the KITTI dataset (c.f. Section 3.2.1) and also combines two CNNs for near and far perception, both using the direct perception approach. It shows, that the direct perception approach is able to perform roughly as good, even though they restrict themselves to the cars closest to the host car. So the direct perception is sufficient for real world examples as well. [CSKX18] [CSKX15]

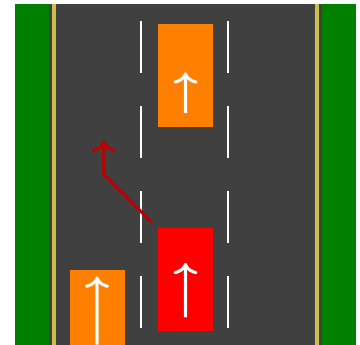


Figure 5.1: A scenario that has to be considered to create a full autonomous driving agent

¹Visualization: <http://deepdriving.cs.princeton.edu/>

Despite the two mentioned comparisons there are still others that need further investigation. Considering more complex tasks, which mediated perception approaches are able to handle, the direct perception needs to prove itself. Considering classification tasks like road signs, pedestrians detection or traffic light detection, including its current light showing, still have to be done in order to create a sufficient agent for real-life cars.

Also more complex scenarios like busy intersections have to be solved. Scenarios as sketched in Figure 5.1, where the overtaking can not take place since a car left and possibly a bit behind of the host is even faster. A number of those scenarios can be managed by a sophisticated controller, but this would again take more time and provides less flexibility.

So concluding, the direct perception definitely is state-of-the-art and has the potential to found the base of a sophisticated autonomous driving agent, if considered as a predictor of distances, which they call affordance indicators. But whether or not the direct perception can handle a complete real-life scenario is still up to prove.

Chapter 6

Conclusion

In this final chapter, we want to recap the languages and frameworks about their abilities and provided functionality.

Property	CNNArchLang	Caffe	Caffe2	MxNet
General Information				
is full framework	X	✓	✓	✓
SLI usage	(✓) ⁵	X ¹	X ¹	✓
mult. computers	(✓) ⁵	✓	✓	✓
Nets Supported				
typical CNNs	✓	✓	✓	✓
arbitrary CNNs	X	✓	✓	✓
Recurrent NNs	X	✓	✓	✓
Constructs				
predefined NNs	✓	✓	✓	✓
pre-trained NNs	X	✓	X ²	✓
arbitrary net creation	✓	X	X	✓
predefined functions	✓	✓	✓	✓
simple function creation	✓	X	X ²	✓
low-level operations	X	X	X	✓
Language bindings				
C++	X	✓ ³	✓	✓
Python	✓	✓ ³	✓	✓
MATLAB	X ²	✓ ³	X	✓
Others	–	Apache Spark ⁴	–	R, Go, Julia, Perl JavaScript, Scala

¹ CUDA does not generally support SLI

⁴ made by a third party[TM16]

² not clearly stated, but also not denied

³ nets still written in prototxt (c.f. Section 2.2)

⁵ through the usage of MxNet

Overall, we can conclude that CNNArchLang has the potential to become a recognized domain specific language for deep learning. The python like syntax and simplistic way of defining a net is superior to the protocol buffer approach of Caffe. Since CNNArchLang is not a full framework, but compiles its code to MxNet code, it benefits from MxNets ability of SLI and even cluster usage. The possibility to define recurrent neural networks

is still up to be implemented and a library of pre-trained models would also increase its potential.

Summarizing, the most important aspect of state-of-the-art deep learning frameworks is the efficiency. Through various methods and researches like KITTI and TORCS the data to train a network is available.

Also to mention is the data estimation in [Grz17]. There, a fleet of 100 cars is analyzed based on how much useful data is produced during one year. After conservative preprocessing this leads to 104 TB of data, which would take an AlexNet roughly 1.2 years to train, if only a single GPU is used. The possibility of using multiple machines is a key element. Using 18 DGX-1 Volta systems, the same effort of training can be done within seven days.

So CNNArchLang is on the right path using MxNet in order to have the cluster computing ability.

Chapter 7

Appendix

```
1 name: "AlexNet"
2 layer {
3   name: "data"
4   type: "Input"
5   top: "data"
6   input_param {
7     shape: {
8       dim: 10
9       dim: 3
10      dim: 227
11      dim: 227
12    }
13  }
14 }
15 layer {
16   name: "conv1"
17   type: "Convolution"
18   bottom: "data"
19   top: "conv1"
20   param {
21     lr_mult: 1
22     decay_mult: 1
23   }
24   param {
25     lr_mult: 2
26     decay_mult: 0
27   }
28   convolution_param {
29     num_output: 96
30     kernel_size: 11
31     stride: 4
32   }
33 }
34 layer {
35   name: "relu1"
36   type: "ReLU"
37   bottom: "conv1"
38   top: "conv1"
39 }
40 layer {
41   name: "norm1"
42   type: "LRN"
43   bottom: "conv1"
44   top: "norm1"
45   lrn_param {
46     local_size: 5
47     alpha: 0.0001
48     beta: 0.75
49   }
50 }
51 layer {
52   name: "pool1"
53   type: "Pooling"
54   bottom: "norm1"
55   top: "pool1"
56   pooling_param {
57     pool: MAX
58     kernel_size: 3
59     stride: 2
60   }
61 }
62 layer {
63   name: "conv2"
64   type: "Convolution"
65   bottom: "pool1"
66   top: "conv2"
67   param {
68     lr_mult: 1
69     decay_mult: 1
70   }
71   param {
72     lr_mult: 2
73     decay_mult: 0
74   }
75   convolution_param {
76     num_output: 256
77     pad: 2
78     kernel_size: 5
79     group: 2
80   }
81 }
82 layer {
83   name: "relu2"
84   type: "ReLU"
85   bottom: "conv2"
86   top: "conv2"
87 }
88 layer {
89   name: "norm2"
90   type: "LRN"
91   bottom: "conv2"
92   top: "norm2"
93   lrn_param {
94     local_size: 5
95     alpha: 0.0001
96     beta: 0.75
97   }
98 }
99 layer {
100  name: "pool2"
101  type: "Pooling"
102  type: "Pooling"
103  bottom: "norm2"
104  top: "pool2"
105  pooling_param {
106    pool: MAX
107    kernel_size: 3
108    stride: 2
109  }
110 }
111 layer {
112   name: "conv3"
113   type: "Convolution"
114   bottom: "pool2"
115   top: "conv3"
116   param {
117     lr_mult: 1
118     decay_mult: 1
119   }
120   param {
121     lr_mult: 2
122     decay_mult: 0
123   }
124   convolution_param {
125     num_output: 384
126     pad: 1
127     kernel_size: 3
128   }
129 }
130 layer {
131   name: "relu3"
132   type: "ReLU"
133   bottom: "conv3"
134   top: "conv3"
135 }
136 layer {
137   name: "conv4"
138   type: "Convolution"
139   bottom: "conv3"
140   top: "conv4"
141   :
142   :
143   :
```

Figure 7.1: The Caffe implementation of the AlexNet. This is only the first 169 lines. The whole net has a size of 284 lines of code. The main reason of this is the highly verbose Protocol Buffer prototxt style. Further explanation in Section 2.2. [Caf]

```

1 architecture Alexnet(img_height=224, img_width=224, img_channels=3, classes=10){
2     def input Z(0:255)^(img_channels, img_height, img_width) data
3     def output Q(0:1)^(classes) predictions
4
5     def split1(i){
6         [i] ->
7         Convolution(kernel=(5,5), channels=128) ->
8         Lrn(nsize=5, alpha=0.0001, beta=0.75) ->
9         Pooling(pool_type="max", kernel=(3,3), stride=(2,2), padding="no_loss") ->
10        Relu()
11    }
12    def split2(i){
13        [i] ->
14        Convolution(kernel=(3,3), channels=192) ->
15        Relu() ->
16        Convolution(kernel=(3,3), channels=128) ->
17        Pooling(pool_type="max", kernel=(3,3), stride=(2,2), padding="no_loss") ->
18        Relu()
19    }
20    def fc(){
21        FullyConnected(units=4096) ->
22        Relu() ->
23        Dropout()
24    }
25
26    data ->
27    Convolution(kernel=(11,11), channels=96,
28        stride=(4,4), padding="no_loss") ->
29    Lrn(nsize=5, alpha=0.0001, beta=0.75) ->
30    Pooling(pool_type="max", kernel=(3,3),
31        stride=(2,2), padding="no_loss") ->
32    Relu() ->
33    Split(n=2) ->
34    split1(i=[0|1]) ->
35    Concatenate() ->
36    Convolution(kernel=(3,3), channels=384) ->
37    Relu() ->
38    Split(n=2) ->
39    split2(i=[0|1]) ->
40    Concatenate() ->
41    fc(->=2) ->
42    FullyConnected(units=10) ->
43    Softmax() ->
44    predictions
45 }

```

Figure 7.2: The AlexNet implemented in CNNArchLang. This is the whole program to describe the whole net. Further explanation in Section 2.1. [TvWH17]

Bibliography

- [Aly08] Mohamed Aly. Real time detection of lane markers in urban streets. In *Intelligent Vehicles Symposium, 2008 IEEE*, pages 7–12. IEEE, 2008.
- [Caf] Caffe model zoo - alexnet. https://github.com/BVLC/caffe/blob/master/models/bvlc_alexnet/deploy.prototxt. [Online; accessed 30-May-2018].
- [Caf18] Caffe2. What is caffe2?, 2018. [Online; accessed 28-May-2018].
- [CL] Tianqi Chen and Mu Li. Mxnet: Flexible and efficient library for deep learning.
- [CLL⁺15] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- [CSKX15] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Computer Vision (ICCV), 2015 IEEE International Conference on*, pages 2722–2730. IEEE, 2015.
- [CSKX18] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. <http://deepdriving.cs.princeton.edu/>, 2018.
- [GB10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [Gib54] James J Gibson. A theory of pictorial perception. *Audiovisual communication review*, 2(1):3–23, 1954.
- [Gib79] James J Gibson. *The ecological approach to visual perception*. Psychology Press, 1979.
- [GLSU13] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [GLW⁺14] Andreas Geiger, Martin Lauer, Christian Wojek, Christoph Stiller, and Raquel Urtasun. 3d traffic scene understanding from movable platforms.

IEEE transactions on pattern analysis and machine intelligence, 36(5):1012–1025, 2014.

- [Grz17] Adam Grzywaczewski. Training ai for self-driving vehicles: the challenge of scale. Technical report, Tech. rep., NVIDIA Corporation. URL <https://devblogs.nvidia.com/parallelforall/training-self-driving-vehicles-challenge-scale>, 2017.
- [Hei17] Heise Verlag. Machine learning: Facebook erweitert einsatz von caffe2, 2017. [Online; accessed 25-May-2018].
- [HW68] David H Hubel and Torsten N Wiesel. Receptive fields and functional architecture of monkey striate cortex. *The Journal of physiology*, 195(1):215–243, 1968.
- [HZCZ17] Xiaobing Han, Yanfei Zhong, Liqin Cao, and Liangpei Zhang. Pre-trained alexnet architecture with pyramid pooling and supervision for high spatial resolution remote sensing image scene classification. *Remote Sensing*, 9(8):848, 2017.
- [JSD⁺14] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [Kin94] Werner Kinnebrock. *Neuronale Netze: Grundlagen, Anwendungen, Beispiele*. Oldenbourg, 1994.
- [KSH12] 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.
- [L⁺15] Yann LeCun et al. Lenet-5, convolutional neural networks. URL: <http://yann.lecun.com/exdb/lenet>, page 20, 2015.
- [MMM03] Masakazu Matsugu, Katsuhiko Mori, Yusuke Mitari, and Yuji Kaneda. Subject independent facial expression recognition with robust face detection using a convolutional neural network. *Neural Networks*, 16(5-6):555–559, 2003.
- [MxNa] Amalgamation: Making the whole system a single file. https://mxnet.incubator.apache.org/faq/smart_device.html. [Online; accessed 25-May-2018].
- [MxNb] Mxnet model gallery. <https://github.com/dmlc/mxnet-model-gallery>. [Online; accessed 25-May-2018].
- [Pom89] Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. In *Advances in neural information processing systems*, pages 305–313, 1989.
- [Roj13] Raul Rojas. *Theorie der neuronalen Netze: eine systematische Einführung*. Springer-Verlag, 2013.

- [SCE⁺17] Vivienne Sze, Yu-Hsin Chen, Joel Einer, Amr Suleiman, and Zhengdong Zhang. Hardware for machine learning: Challenges and opportunities. In *Custom Integrated Circuits Conference (CICC), 2017 IEEE*, pages 1–8. IEEE, 2017.
- [She14] Evan Shelhamer. Deep learning for computer vision with caffe and cudnn. <https://devblogs.nvidia.com/deep-learning-computer-vision-caffe-cudnn/>, 2014.
- [Tim18] Thomas Timmermanns. Modelling languages for deep learning based cyber-physical systems. 2018.
- [TM16] Michael Thomas and Gabriela Motroc. Caffeonspark: Yahoo macht deep-learning-software quelloffen. <https://jaxenter.de/caffeonspark-yahoo-deep-learning-software-quelloffen-35668>, 2016.
- [TvWH17] Thomas Timmermann, Michael von Wenckstern, and Malte Heithoff. Cnnarchlang. <https://github.com/EmbeddedMontiArc/CNNArchLang>, 2017.
- [Ull80] Shimon Ullman. Against direct perception. *Behavioral and Brain Sciences*, 3(3):373–381, 1980.
- [Var08] Kenton Varda. Protocol buffers: Google’s data interchange format. *Google Open Source Blog*, Available at least as early as Jul, 72, 2008.
- [VS16] David H Von Seggern. *CRC standard curves and surfaces with mathematica*. CRC Press, 2016.
- [WEG⁺00] Bernhard Wymann, Eric Espié, Christophe Guionneau, Christos Dimitrakakis, Rémi Coulom, and Andrew Sumner. Torcs, the open racing car simulator. *Software available at http://torcs.sourceforge.net*, 4, 2000.
- [Yeg16] Serdar Yegulalp. Why amazon picked mxnet for deep learning. URL: <https://www.infoworld.com/article/3144025/cloud-computing/why-amazon-picked-mxnet-for-deep-learning.html>, 2016.