

UNIVERSITÀ DEGLI STUDI DI BRESCIA
FACOLTÀ DI INGEGNERIA



CORSO DI LAUREA MAGISTRALE IN INGEGNERIA INFORMATICA

*ESTRAZIONE AUTOMATICA DI INFORMAZIONE
SIMBOLICA DA IMMAGINI*

(Automatic extraction of symbolic information from images)

Relatore:

Ch.mo Prof. Riccardo Cassinis

Correlatore:

Ch.mo Prof. Marco Ragni

Laureando:

Francesco Bonfadelli

Matricola 83174

Anno Accademico 2012-2013

Contents

Sommario	6
Introduzione	6
L'ambiente di lavoro	7
Lo stato dell'arte	8
ACT-R	8
OpenCV	9
Gli obiettivi	11
I Requisiti	13
Requisiti Funzionali	13
Requisiti Non Funzionali	14
Il processo di sviluppo	15
Scrum	15
Il Processo di Sviluppo Adottato	17
Il design del software	18
Implementazione e testing	20
Conclusioni	21
 1. Introduction	 24
 2. State Of The Art	 26
2.1. ACT-R	26
2.1.1. The Theory Behind ACT-R	26
2.1.2. The Architecture	28
2.1.3. ACT-R Models	30
2.2. OpenCV	30
2.2.1. Introduction to Computer Vision	31
2.2.2. OpenCV History	34
2.2.3. Main Features	34
2.2.4. Architecture	35

3. Objective	37
3.1. The Context	37
3.1.1. The Goal of the Team	37
3.2. The Goals of the Software	38
3.3. Requirements	39
3.3.1. Functional Requirements	40
3.3.2. Non-Functional Requirements	41
4. Development Process and Tools	42
4.1. Scrum	42
4.1.1. The Scrum Team	43
4.1.2. Events	45
4.1.3. Artifacts	47
4.2. The Adopted Development Process	49
5. Design	51
5.1. The Architecture	51
5.2. Processes	52
5.3. Communication	53
5.4. Vision Operation Module	54
5.4.1. Architecture	54
5.4.2. Feature Extraction Classes	55
5.5. Class Hierarchy	57
6. Implementation And Testing	59
6.1. Implementation	59
6.1.1. N-sides Shape Detection Algorithm	59
6.2. Testing	61
7. Conclusions	63
A. Appendix A	i
A.1. Details About The Architecture	i
A.1.1. Overview	i
A.1.2. Server Module	ii
A.1.3. Utility Part	iii

A.2. Message Structure	v
A.2.1. Request Format	v
A.2.2. Response Format	v
A.3. Tools Used in The Development	vii
References	ix

List of Figures

0.1. <i>Esempio di un istanza del gioco di Rush Hour</i>	12
0.2. <i>Diagramma delle classi che illustra le classi più importanti del Vision Module</i>	19
2.1. <i>Structure of ACT-R.</i>	29
2.2. <i>Representation of an image in a computer vision system.</i>	32
2.3. <i>How the 2d representation changes with the viewpoint</i>	33
3.1. <i>Example of level of Rush Hour game</i>	39
4.1. <i>Overview of the Scrum Process</i>	42
4.2. <i>The Scrum Team</i>	43
4.3. <i>The Scrum Events</i>	45
4.4. <i>Sprint Backlog and Burn Down Chart</i>	48
5.1. <i>High level scheme of the architecture</i>	51
5.2. <i>Scheme of processes and threads of the architecture</i>	52
5.3. <i>Class diagram illustrating a high level overview of the vision operation module</i>	55
5.4. <i>Class diagram illustrating the main roles in the feature extraction process</i> .	56
5.5. <i>Class diagram illustrating the class hierarchy of the recognized objects</i> . . .	57
A.1. <i>Class diagram which illustrates the whole architecture of the Vision Module</i>	ii
A.2. <i>Class diagram describing the implementation part of the software dedicated to the communication</i>	ii
A.3. <i>Class diagram showing the utility part of the Vision Module</i>	iv

Lo stato dell'arte

Di seguito si descrivono i due principali strumenti utilizzati per questo progetto: la libreria di computer vision *OpenCV* e l'architettura cognitiva *ACT-R*.

ACT-R

ACT-R è un'*architettura cognitiva*, cioè l'implementazione di una teoria riguardante il sistema cognitivo umano. Come tale, esso modella struttura e comportamento del cervello umano cercando di spiegare come le diverse componenti collaborino tra loro formando la mente umana.

La teoria di ACT-R si basa sulla *teoria unificata della cognizione* (Unified theory about cognition), sviluppata da John Robert Anderson, docente presso la Carnegie Mellon University. Il concetto fondamentale di tale teoria è che la mente umana si comporti come un sistema integrato. Secondo essa, infatti, il cervello è separato in diversi moduli indipendenti, ognuno dei quali è dotato delle proprie funzioni ed è incaricato di svolgere compiti specifici. Nel momento di attuare un comportamento, tuttavia, i moduli interagiscono tra loro per raggiungere l'obiettivo comune, proprio come un sistema integrato. La comunicazione tra le diverse parti avverrebbe grazie a specifiche connessioni tra i moduli [ABB⁺04].

Un altro pilastro fondamentale su cui si basa la teoria di ACT-R è la distinzione tra *memoria dichiarativa* e *memoria procedurale*. La prima rappresenta fatti e nozioni che l'essere umano sa ed è conscio di sapere. Per richiamare questo tipo di conoscenza, l'essere umano deve effettuare un processo consapevole. La seconda, invece, si riferisce a tutte quelle abilità e capacità che l'essere umano sa ma che ha imparato in maniera implicita. Esempi di questo tipo di conoscenza sono la lettura e la scrittura [And76].

Gli elementi di base dell'architettura di ACT-R sono *chunk* e *produzioni*. I chunk rappresentano la memoria dichiarativa e sono delle strutture dati caratterizzate da un tipo, chiamato *type*, e da una lista di coppie, ognuna delle quali è costituita da un attributo, chiamato *slot*, e da un valore, chiamato *value*. Le produzioni possono essere paragonate a funzioni e definiscono la sequenza di azioni che possono essere effettuate. Ogni produzione presenta un insieme di precondizioni che determinano le condizioni che devono essere verificate affinché la funzione possa essere eseguita [Bota].

Dal punto di vista dell'architettura, ACT-R è organizzato in *moduli*. Ciascun modulo è costituito da chunk e da produzioni e svolge un insieme determinato di funzioni cognitive. Essi sono l'analogo dei gruppi di neuroni che si attivano nel cervello nel momento in cui viene effettuata una determinata azione da parte dell'essere umano. I moduli sono

tramite matrici di numeri. Tale rappresentazione rende complesse molte delle operazioni di visione, comprese quelle più semplici come, per esempio, il riconoscimento di un semplice oggetto. Purtroppo al momento questa modalità di descrizione è la migliore alternativa tra le possibilità esistenti [BK08].

Oltre alla rappresentazione matriciale, uno dei motivi che rendono così complicata la computer vision è la necessità molto frequente di ricostruire un mondo tridimensionale a partire da un'immagine bidimensionale. Questo è un problema mal posto, che ha come diretta conseguenza l'esistenza di infinite soluzioni, cioè infinite ricostruzioni tridimensionali del mondo a partire dall'immagine bidimensionale. In particolare, la soluzione dipende dall'angolo con cui si analizza l'immagine [BK08].

Un fattore che aumenta ulteriormente la complessità dei compiti è la presenza del rumore, che raggruppa tutte le possibili cause di variabilità presenti nel mondo: tempo atmosferico, luminosità, riflessi, movimenti, deformazioni dovute alle lenti o ai sensori e altri fenomeni. Il rumore può essere classificato in *noto a priori*, quindi prevedibile, e *non noto a priori*, quindi non prevedibile. Il primo, essendo conosciuto, può essere corretto fino a essere completamente rimosso. Questo è il caso, ad esempio, della deformazione delle lenti delle videocamere. Il secondo, invece, essendo casuale, non può essere del tutto eliminato ma può essere solamente ridotto. Tipicamente si utilizzano metodi statistici per gestire questo secondo tipo di rumore [BK08].

Un espediente usato per rendere più facili le operazioni di computer vision è l'introduzione di informazione contestuale. In particolare, più il problema è vincolato, più è facile utilizzare i vincoli per semplificare il problema, più la soluzione ottenuta è affidabile. In questo modo, però, la soluzione perde di generalità e diventa specifica per il problema. Le tecniche di machine learning, che estraggono automaticamente l'informazione contestuale e la utilizzano per risolvere i problemi, possono essere sfruttate nelle operazioni di computer vision [BK08].

Storicamente, il progetto OpenCV viene avviato nel 1999 da Intel Research, divisione del gruppo Intel dedicata alla ricerca. L'iniziativa è finalizzata a migliorare le applicazioni che richiedono un alto carico di lavoro della CPU, in particolare quelle legate alla visione artificiale. Le intenzioni iniziali comprendono la creazione di codice ottimizzato per un'infrastruttura standard di base per la visione, la diffusione di essa tra gli sviluppatori, la portabilità del codice e la possibilità di creare sia applicazioni di natura commerciale che gratuite grazie all'infrastruttura. Nel 2000 viene rilasciata la prima versione alpha, seguita da cinque versioni beta tra il 2001 e il 2005 che portano nel 2006 a OpenCV 1.0. Nel 2008 l'incubatore tecnologico Willow Garage comincia a supportare il progetto e nel 2009 viene rilasciata la versione 2.0. Come stabilito nel piano di rilascio, ogni sei mesi

del compito: infatti è desiderabile che l'architettura cognitiva prenda decisioni riguardanti l'ambiente esterno dopo aver ricevuto ed elaborato tali dati inerenti ad esso.

Il Modulo Vision Operation è diviso logicamente in tre gruppi di classi: quella che si dedica all'elaborazione di immagini e video, quella che rappresenta le forme riconosciute nelle immagini e quella che gestisce le funzioni di utilità e le eccezioni. In figura 0.2, si può vedere un diagramma delle classi, realizzato nel linguaggio UML, che descrive ad alto livello i primi due gruppi di classi del Vision Module.

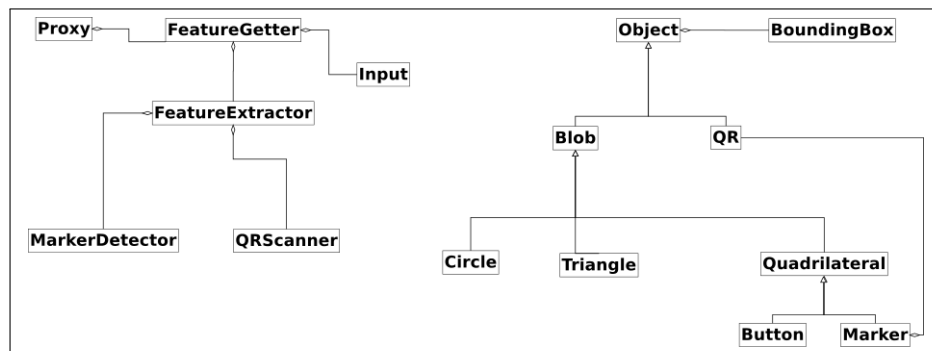


Figure 0.2.: Diagramma delle classi che illustra le classi più importanti del Vision Module

Nel diagramma, sulla sinistra si può vedere l'insieme di classi che elaborano le immagini. In particolare, la classe *Feature Extractor* è la classe fondamentale di questo insieme, in quanto è quella che effettua tutte le elaborazioni sull'immagine, riconoscendo le forme contenute in essa e creando la corrispondente gerarchia di oggetti. Oltre a tali operazioni, essa riconosce i colori, calcola le distanze tra oggetti, stabilisce le loro posizioni reciproche e li confronta sulla base della dimensione. La classe *Input* si occupa di interagire con il filesystem e l'hardware del calcolatore; in particolare, acquisisce immagini e video. *Feature Getter* rappresenta un livello di astrazione e offre un'interfaccia semplice a *Proxy*, classe incaricata di preparare i dati in un formato standard per la comunicazione con ACT-R. Le classi *MarkerDetector* e *QRScanner* svolgono funzioni non presentate tra i requisiti enunciati in questo documento. Esse acquistano senso se si analizzano gli obiettivi globali del software. Le classi sono comunque presenti nel diagramma perché influenzano l'architettura del sistema.

Sulla destra nel diagramma è presentata la gerarchia degli oggetti. *Object* è la super classe e rappresenta qualsiasi oggetto che può essere riconosciuto dal sistema. Esso viene specificato dalle due sotto-classi *Blob* e *QR*. *Blob* rappresenta tutte le possibili forme geometriche e la sua presenza ha senso per distinguerle da oggetti di altra natura, come *QR* o altri che verranno introdotti in futuro. *QR* e *Marker*, esattamente come *MarkerDetector*

1. Introduction

This document describes the development of a software module which implements simple visual operations and communicates their results to the cognitive architecture ACT-R¹. In particular, the software has to analyze images and process them by recognizing the objects they contain, their colors, their dimensions and their positions.

The activity has been carried on in the *Center for Cognitive Science* of *Albert-Ludwigs-Universität Freiburg*. As the name suggests, the field of research of the center is cognitive science². The research, currently, is focused on the spatial reasoning³ and the researchers use ACT-R as a support instrument for their studies. In this context, the work discussed in this document represents a part of the work developed by a team of three people, whose ultimate objective is to make ACT-R perception more similar to human one.

One of the main limitations of the cognitive architecture is that it works in a virtual environment, too easy for representing the reality. The new software represents a starting point for allowing ACT-R to process objects directly in the real world, overcoming this limit.

Stated the impossibility of processing the visual data directly inside the cognitive architecture, the module has been created externally as an independent software. This fact brings two main consequences: the use of OpenCV library for implementing the visual operations and the introduction of a client server architecture in order to implement the communication between ACT-R and the software.

This document is organized in chapters, each of which focuses on specific aspects of the creation process. Chapter 2 presents the context of work by describing the cognitive architecture ACT-R and the computer vision library OpenCV. Chapter 3 explains in detail the objective of the work and defines the requirements. Chapter 4 describes Scrum, which is the chosen development framework, and motivates its adoption. Chapter 5

¹ACT-R is a cognitive architecture, that is a framework that models structure and behavior of the human brain. See act-r.psy.cmu.edu.

²Cognitive science is an interdisciplinary study which has the goal of understanding the nature of the human mind. See www.cognitivesciencesociety.org.

³Spatial reasoning is a sub-field of cognitive science focused on the reasoning regarding object in the space.

presents the main aspects of the architecture of the software and chapter [6](#) deepens the description of computer vision algorithms and describes the choices about the testing activity. After the Conclusions, in Appendix [A](#) the reader can find details about the implementation of the software.

2. State Of The Art

This chapter describes the main working instrument used for this work: the cognitive architecture *ACT-R* and the computer vision library *OpenCV*.

2.1. ACT-R

ACT-R, that stands for *Adaptive Control of Thought-Rational*, is a cognitive architecture, i.e. a computer implementation of a theory about human cognition. As such, it models the structure and behavior of the human brain trying to explain how all the components of the brain work together and form the human mind.

The following sections describe the theory on which ACT-R relies, how the theory is implemented in the framework and the concept of *model*.

2.1.1. The Theory Behind ACT-R

The following section describes two fundamental assumptions on which ACT-R relies, the *Unified Theory About Cognition* and the classification of the human memory in *declarative* and *procedural*. Both of these assumptions are important because they determinate the structure of the framework.

Unified Theory About Cognition

ACT-R implements the homonym theory developed by John Robert Anderson, professor of psychology and computer science at Carnegie Mellon University. Such theory tries to explain the overall behavior of the human mind through connections between its well-defined components that, combined together, form an integrated system. The following quote, that explains the meaning of integrated system, comes from Allen Newell, the man who inspired Anderson in creating ACT-R theory [ABB⁺04].

A single system (mind) produces all aspects of behavior. It is one mind that minds them all. Even if the mind has parts, modules, components, or whatever, they all mesh together to produce behavior. Any bit of behavior has

causal tendrils that extend back through large parts of the total cognitive system before grounding in the environmental situation of some earlier times. If a theory covers only one part or component, it flirts with trouble from the start. It goes without saying that there are dissociations, independencies, impenetrabilities, and modularities. These all help to break the web of each bit of behavior being shaped by an unlimited set of antecedents. So they are important to understand and help to make that theory simple enough to use. But they don't remove the necessity of a theory that provides the total picture and explains the role of the parts and why they exist [New94, p.17-18].

Declarative memory and procedural memory

In psychology, *memory* is defined as the processes by which information is encoded, stored and retrieved [BEAA09].

ACT-R's most important assumption about knowledge is based on Anderson's theory about memory. Anderson divides memory into *declarative* and *procedural* [ABB⁺04].

Declarative memory refers to all the information that can be consciously recalled. This kind of knowledge comprehends facts and notions that human beings explicitly know. To call back this kind of information, there must be a conscious process by the human being. For this reason, this kind of memory is also called *explicit*.

In contrast, procedural memory refers to all that notions or skills that human beings have but which they learnt in an implicit way. Examples of this knowledge are driving, reading and writing. In this case, in order to call back this kind of information, the human being does not need a conscious process. That is why this kind of memory is also called *implicit* [And76].

When a person starts learning typewriting, for example, an attempt he can make in the beginning is trying to memorize the layout of the keyboard. The aware knowledge of all the positions of the keys is the declarative memory. After having become a skilled typewriter, the same person will type faster putting his fingers on the right keys and pushing them in the correct order, without thinking anymore about the positions of the keys on the keyboard. Moreover, if someone asks him where the position of a certain character is on the keyboard, he will probably answer that he can not say it without looking at it. This is because, now, for this task he is using his procedural memory [And93].

2.1.2. The Architecture

This section, after having defined *chunks* and *productions*, the building blocks of ACT-R structure, presents the architecture of the framework.

Chunks and productions

In ACT-R, declarative memory is represented by structures, called *chunks*, and procedural memory by rules, called *productions*. Chunks and productions are the basic building blocks of ACT-R *models*. The concept of model will be described in 2.1.3.

The *chunks* are data structures which are defined by their *type* and their *attribute list*. This is a tuple of pairs, each of which is made up by a fixed part and a variable part. The fixed part is the *name* of the attribute and is called *slot*. The variable part is the *value* of the attribute. Each chunk has also a *name* but it is not considered to be a part of the chunk itself, as it does not exist in ACT-R theory. It is used only for convenience to reference the specific chunk when writing models. The chunk-types can be organized into hierarchies.

The *productions* are ACT-R equivalent of functions. They define sequences of actions and can be fired only if a set of preconditions is satisfied. They can be represented as *if-then* rules, where the *if-part* is a set of conditions that must be true for the production to apply and the *then-part* is the action of the production and consists of the operations the model should perform when the production is selected and used.

In general there could be some conflicts between productions. This happens when preconditions of two or more productions are satisfied at the same time. In these cases the production to be fired is the one with the highest *utility value*, a numeric quantity which gives a priority measure. It can be set a priori by the modeler or learnt while the model is running. The latter, together with the fact that in the calculation of this value the probability of reaching the goal and the time estimate are included, constitutes the basis of the learning mechanism embedded in ACT-R: the more successfully a production is, the more its estimated probability and utility grow, increasing the probability for that production to be selected again [Bota].

Organization of the Modules

All the activities carried out by the human brain, like talking or moving, are performed by neurons located close together in a well defined and limited area of the cortex. Trying to imitate this "architecture", ACT-R's framework is structured in different *modules*, each of which represents one specific function of the human brain.

Figure 2.1 shows the modular structure of ACT-R. In the picture you can see two groups of modules, separated by the *procedural module*.

The first group comprehends *visual*, *aural*, *manual* and *vocal* modules. These let the model interact with the environment. The *visual module* is responsible for recognizing objects in the visual scene and shifting the focus to them. Similarly, the *aural module* identifies sounds and moves the attention to them. The *manual module* can move the virtual hands and perform actions like pressing the key on a keyboard or moving the mouse while the *vocal module* controls the virtual voice.

The other group comprehends *goal module*, *imaginal module* and *declarative module*. These represent the internal information of the model. The *goal module* provides the system with the structure of the goal of the task, defined as a chunk. The *imaginal module* has to contain and update the current context relevant to the current task. The *declarative module* provides the model with a declarative memory, thus it stores the declarative chunks generated by the model and provides a mechanism for retrieving them.

Finally, the *procedural module* is responsible of the communication and the coordination of all the other modules.

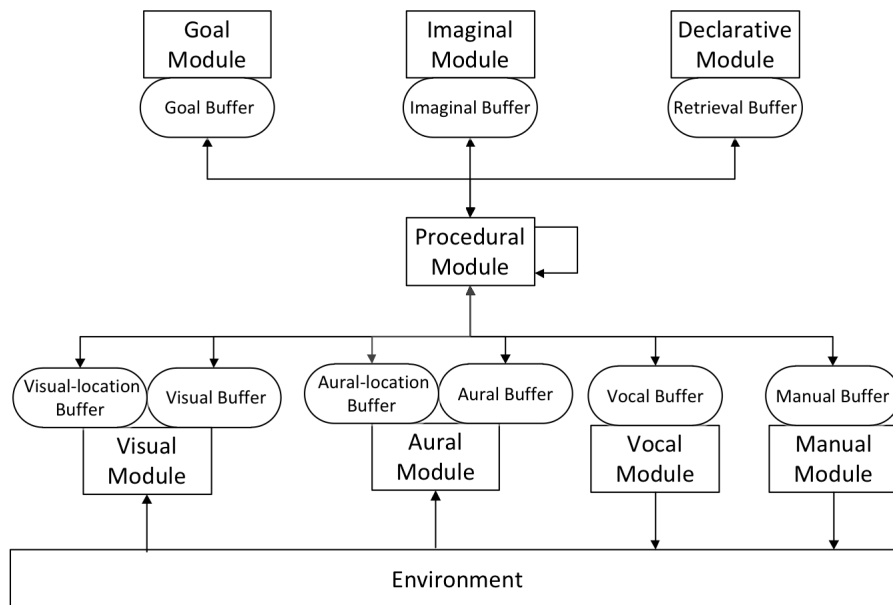


Figure 2.1.: Structure of ACT-R.

All the modules are independent of each other, they do not share variables or information. They can communicate with each other through *buffers*, which represent the interface of a module towards the others. A module can have no buffers as well as one or more than one. The communication consists in exchanging chunks. Each module can

read chunks from every buffer but it can make changes only to the chunks in its own buffers. Moreover each buffer can hold one chunk at a time.

Although modules usually work in a parallel way, their interactions can be only serial. There are two reasons for this limitation: the first one is that the structure of the buffers can hold only one chunk at a time and the second one is that only one production can be fired at a time [Bota].

2.1.3. ACT-R Models

ACT-R is able to perform several task-independent reasoning. To adapt these reasoning to a single task it is necessary an additional layer called *model*.

Each model contains the specific modelers' assumptions about the task within ACT-R view of cognition. That knowledge is expressed through productions that, when the model is running, interact with the modules, querying them and reading their buffers.

Running a model produces a series of atomic cognitive operations which step-by-step leads to the solution of the task. Each operation is associated with quantitative and qualitative measures like the correctness of the goal and the time necessary to complete the operation. This fact gives the model the possibility to predict the sequence of cognitive actions produced by human beings when they try to solve the same task. Comparisons with human performances can be useful to evaluate the quality of the models [Sea02].

ACT-R is written in Lisp and provides its own lisp-like language for the modelers to write the models.

2.2. OpenCV

OpenCV, an abbreviation that stands for *Open Source Computer Vision*, is a computer vision library that was originally developed by Intel and, later on, by Willow Garage. It is a cross-platform library and can run under Linux, Windows and Mac OS X. It is released under a BSD license, thus it is free and open source. In the beginning it was developed in C and C++ and afterwards it was expanded by the addition of interfaces for other languages as, for example, Java, Python, Ruby and Matlab.

OpenCV is designed for computational efficiency and with a strong focus on real-time applications. On Intel architectures there is the possibility to further optimize the library thanks to Intel's *IPP* library. IPP library, that stands for *Integrated Performance Primitives*, contains a series of low-level optimized routines in many different algorithmic areas. If this library is installed, OpenCV uses it automatically at runtime.

OpenCV offers a simple to use infrastructure that helps programmers create quite complicated vision applications in an acceptable time. The version 2.4 has more than 2500 algorithms that cover many areas of vision. The library has been used in many applications as, for example, mine inspection and robotics. It also contains *MLL*, a full and general purpose *Machine Learning Library*, which provides functions for pattern recognition and clustering [BK08, Ope12c].

The following sections contain an introduction to computer vision, a brief history of OpenCV library, a description of its main features and an overview of its architecture. For more information about computer vision see texts by Trucco [TV98] for a simple introduction, Forsyth [FP11] as a comprehensive reference, and Hartley [HZ03] and Faugeras [Fau93] for how 3D vision really works.

2.2.1. Introduction to Computer Vision

Computer vision is the transformation of image or video data into a new representation. Such representation can be either another image, for example a gray-scale version of the original one, or even a decision based on the information extracted from the input data. Often computer vision algorithms use contextual information in order to simplify the problem.

People not expert in computer vision may mislead the complexity of its tasks. Often, in fact, even a very simple operation for the human brain, like recognizing an animal in a picture, represents a very complex task for a machine. Human visual system is very complex. It divides the visual signals into many channels, each of which brings a different kind of information to the brain. The brain uses an attention system which identifies important parts of the image and focuses on them, excluding the analysis of the remaining parts. Moreover, further information coming from all the other senses and the experience acquired during the whole life of the person completes the analysis of the image, in a way still not completely understood. All these factors create the visual perception of human beings.

In a machine vision system, instead, when the computer receives an image or a video from a camera, it receives simply a grid of numbers. In most of cases, there are no mechanisms for pattern recognition, automatic control of aperture and focus of the camera nor cross-associations with other senses or experience. Figure 2.2 explains this concept. The machine does not recognize automatically the side mirror of the car, it sees just a grid of numbers. A goal of computer vision is to turn the grid that describes this particular object into the object itself.

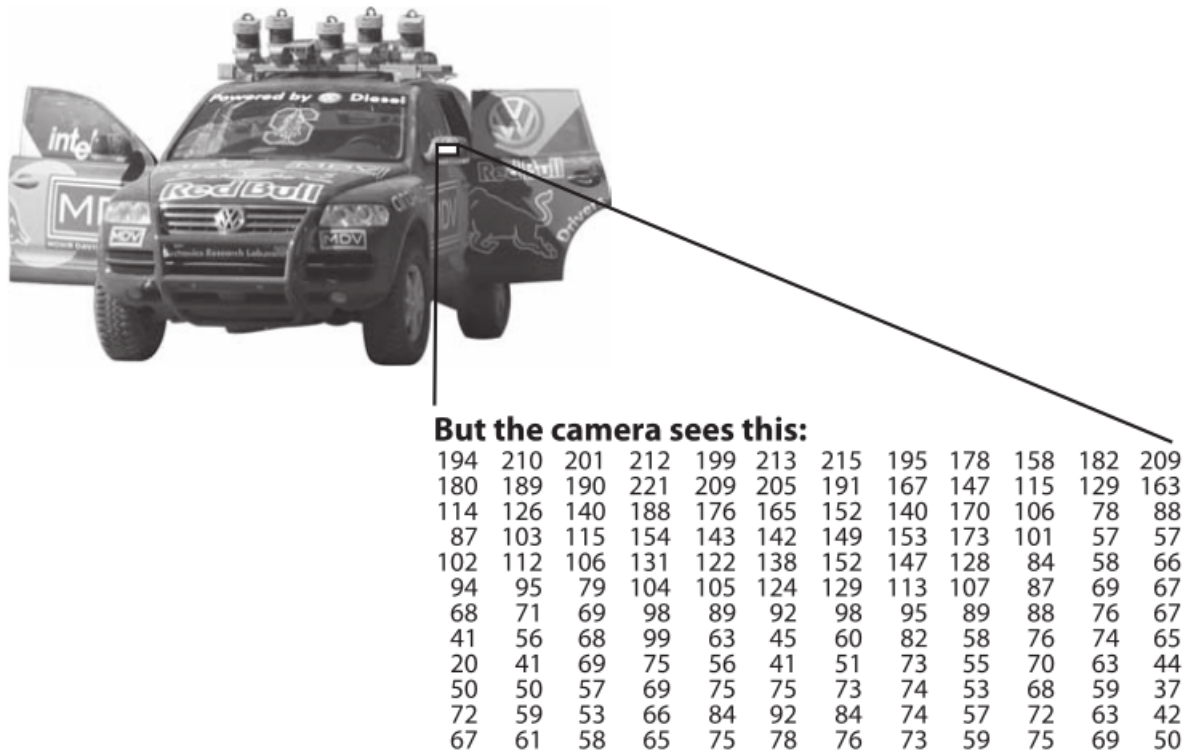


Figure 2.2.: Representation of an image in a computer vision system.

Figure 2.3 shows another aspect that makes computer vision tasks so difficult to solve. An image is a 2D description of a 3D world. The problem in the way it is posed is impossible to solve. There are infinite possible ways to reconstruct a 3D world from a 2D image. As shown in the picture, the solution changes with the point of view of the camera. This is formally an ill-posed problem.

Moreover, generally speaking, data is affected by noise. The word noise comprehends all the possible variations in the world, like weather, lightning, reflections and movements, and all the imperfections affecting the input mechanisms, like errors of the lens or of the sensors.

Some causes of noise are easy to remove. The errors caused by the lens, for example, can be mathematically described by polynomial functions, thus it is quite simple to correct them. The environment variations, instead, are not well known a priori, hence are more difficult to delete. Typically, this category of noise is dealt with statistical methods, that generally do not remove it completely but attenuate it. Say, for example, to have many different pictures of the same visual scene, every one taken in different conditions of weather, lightning, reflections and so on. In this case the noise is considered, statistically

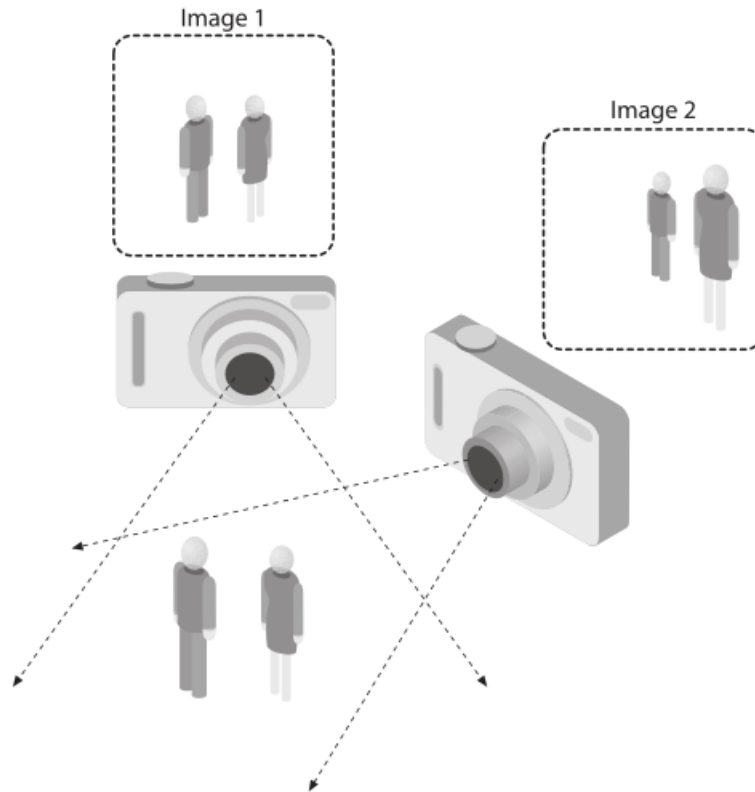


Figure 2.3.: How the 2d representation changes with the viewpoint

speaking, a white noise. A statistical approach suggests to make the average of all the pictures. In this way, the features of the image remain, while the noise is reduced by the average operation.

A common way to simplify the problem of computer vision is to add contextual knowledge to the solution. Consider the example of a mobile robot which has to find and pick up staplers in a building. The robot can use the fact that the staplers are always on desks. This fact gives the robot an implicit information about the size and the position of the object, that brings it to exclude false "recognitions". Moreover, if all the desks have rectangular shapes, are made of wood and have the same color, the computer vision system can use this additional information to exclude other false staplers, increasing the accuracy of the research. As general rule: *"the more constrained a computer vision context is, the more we can rely on those constraints to simplify the problem and the more reliable our final solution will be"* [BK08, 5]. Of course, in this way, the solution is not general purpose any more.

Machine learning techniques allow to extract variables that model contextual information from a limited number of training cases and using them to solve other instances of

the modeled problem. Differently from the previous example, such variables are not set a priori by the modeler but are learnt from a series of trials. For this reason they can be hidden parameters or may even not represent physical quantities and remain mathematical variables [BK08].

2.2.2. OpenCV History

The OpenCV Project starts in 1999 as an Intel Research initiative aimed to improve CPU intensive applications as a part of projects including real-time ray tracing, which in computer graphics is a particular technique for generating images, and 3D display walls. The early goals of the project are: developing optimized code for basic vision infrastructure, spreading this infrastructure to developers and making it portable and available for free, using a license that allows the developers to create both commercial and free applications.

The first alpha version is released to the public in 2000, followed by five beta versions between 2001 and 2005, which lead to version 1.0 in 2006. In 2008, the technology incubator Willow Garage begins supporting the project and, in the same year, version 1.1 is released.

In October 2009, OpenCV 2.0 is released. It includes many improvements, such as a better C++ interface, more programming patterns, new functions and an optimization for multi-core architectures. According to the current OpenCV release plan, a new version of the library is delivered on a six-months basis [Ope12a].

2.2.3. Main Features

OpenCV offers a wide range of possibilities. First of all, it provides an easy way to manage image and video data types. It also offers functions to load, copy, edit, convert and store images and a basic graphical user interface that allows the developers to handle keyboard and mouse and display images and videos. The library permits to manipulate images even with matrix and vector algebra routines. It supports the most common dynamic data structures and offers many different basic image processing functions: filtering, edge and corner detection, color conversion, sampling and interpolation, morphological operations, histograms and image pyramids. Beyond this, it integrates many functions for structural analysis of the image, camera calibration, motion analysis, object recognition and machine learning [Aga06].

2.2.4. Architecture

Since version 2.2, the OpenCV library is divided into several modules, built in library files. Each module serves a specific purpose and includes functions in order to accomplish it.

The *core module* contains the core functionalities of the library, i.e. the basic data structures, both in C and C++ languages, and a variety of functions used by all the other modules. Such functions have many purposes: operating on arrays and matrices, drawing shapes on the images, storing and restoring library data structures and primitive data types in XML and YAML formats and other.

The *imgproc module* contains the main image processing functions. Processing an image means “using higher-level operators that are defined on image structures in order to accomplish tasks whose meaning is naturally defined in the context of graphical, visual images” [BK08, 128]. Such functions are, for example, operations to perform various linear or non-linear filtering on 2D images, geometrical transformations of 2D images (resize, affine and perspective warping, generic table-based remapping), color space conversions, histogram analysis, feature detection and object detection.

The *highgui module* allows the library to interact with the operating system, the file system and hardware such as cameras. Its interface permits to open windows, display images, read and write graphics-related files, both images and video, and handle simple mouse, pointer, and keyboard events. It also allows to add simple user interface elements to the displayed windows.

The *features2d module* allows to use different algorithms to implement the *interest points processing*. This technique relies on the idea that instead of analyzing the complete image it is better to work on a limited number of particular pixels. Such points must be distinguishable features in all the images in which they must be detected. To do this, the module offers functions for feature detection, feature description and descriptor matching.

The *calib3d module* focuses on camera calibration and stereo vision. Camera calibration permits to modify images and videos by correcting the errors derived by the use of lenses. Stereo vision allows to have benefits from all the advantages given by using two equal cameras instead of one, like for example measuring distances with high precision.

The *video module* implements some techniques for video processing, in particular *motion estimation*, *feature tracking*, and *foreground extraction* algorithms. Motion estimation techniques define motion vectors that describe the transformation from consequent frames. Feature tracking algorithms improve motion estimation by working only with some features instead of sets of pixels. Foreground extraction algorithms try to recognize

foreground objects of interest comparing an observed image with an estimate of the same image as if it contained no objects of interest.

The *objdetect module* allows the detection of particular categories of objects such as faces or cars. To do so, it uses functions for training the software in recognizing a particular class of objects and verifying if an input image contains such kind of object once the learning is terminated.

OpenCV also includes other utility modules: the *ml module* contains machine learning functions, the *flann module* computational geometry algorithms, the *contrib module* some contributed code, the *legacy module* some obsolete code and the *gpu module* support for gpu acceleration [Lag11, Ope12b, BK08].

3. Objective

In the first section the chapter introduces the purposes of the working team and motivates the choices of the objectives of the software. Then, in the second one, it presents the goals of the software and, in the third, the requirements.

3.1. The Context

The work presented in the following chapters is a subset of the work of a research group, which aims to provide ACT-R of a visual system which is as similar as possible to the human one. As this objective is very broad, this work focuses only on a particular subset of tasks.

To reach such goal, the researchers continuously monitor the progresses in the computer vision research field in order to find useful innovative techniques to combine with consolidated solutions. Such search, anyway, it is not easy. In fact, although computer vision studies are very active and find every year new algorithms able to solve different problems, often such algorithms introduce contextual information that makes them tasks specific: they are suitable to solve only the tasks for which they are designed. For this reason, it is very difficult to find a general purpose solution which is able to make machine vision system behave like the human one.

Like most of research works, this project has to deal with a certain degree of uncertainty that makes difficult to set any time constraints. The author of this work, however, has been inserted for a limited time interval in a team composed by three people, whose aim is to implement particular solutions in order to solve specific tasks. For this reason the following parts of this document describe only the features developed during this amount of time by the author, explaining them in the context of the objectives of the team.

3.1.1. The Goal of the Team

The goal of the team is to create a general purpose software, which receives as input a video stream or an image and performs an object recognition process on them. ACT-R

uses the information obtained by this tool to orientate inside a building. In particular, the software analyzes a video looking for objects and transmits all the recognized items to a model developed in ACT-R. The model has to distinguish different categories of rooms basing on the information about the objects they contain.

As a starting point for such software, the team has to implement a tool for recognizing geometrical shapes in the input images. Such software has to accomplish two goals: the first one is to extract a subset of all the features in the image, in order to simplify the object recognition process; the second one is to make easier the definition of models for the modelers. This tool, which represents the main topic of this work, is described better in the next sections.

3.2. The Goals of the Software

The software to be designed and developed is a visual module for the cognitive architecture ACT-R, whose aim is to improve ACT-R visual perception in order to make it more similar to the human one.

For developing a model in ACT-R, the modelers need a series of parameters on which basing their assumptions about the model itself and its performances. The collection of such data is accomplished thanks to experiments in which human beings take part. For every experiment a specific program is created in order to measure some variables as, for example, the time needed to complete a task and the accuracy of the result. The parameters obtained by these experiments are then used as a reference for the development of ACT-R model.

Image 3.1 shows a test case of an experiment in which the goal is to solve some levels of the game *Rush Hour*. As shown, the grid contains some colored rectangles, each of which represents a car. The player must free the red one making it go out of the grid through the exit on the right. The cars can be moved only in the horizontal and vertical directions, according to their orientation. The lower the number of moves, the better is such solution.

In order to conduct an experiment, the modelers prepare a set of instances of the problem and every person who is going to participate to the experiment has to face a fixed number of them. For each instance, the ad-hoc program shows the image with the initial configuration (figure 3.1) and the person has to give the solution of the game by clicking on the rectangles in the order he thinks to be correct. As the program never updates the image, the player has to imagine the solution. The software manages to interpret the direction of the shift by analyzing the movement of the eye, thanks to a

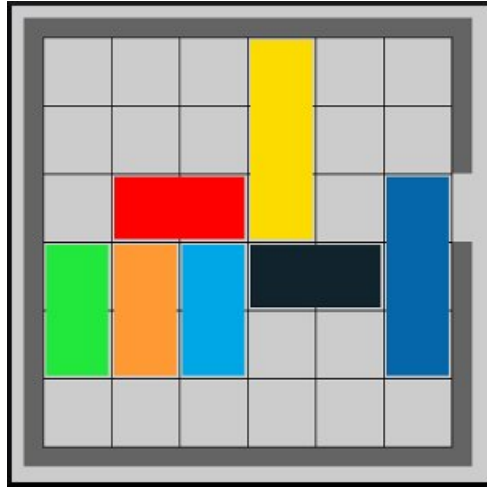


Figure 3.1.: Example of level of Rush Hour game

module of *eye tracking*. Moreover, it records the movements of the mouse, the clicks of the mouse, the time needed to give the solution and the correctness of the solution.

The model written in ACT-R, in order to solve a particular instance of Rush Hour problem, needs to receive as input the initial configuration. At the moment, such configuration is not an image, like the one shown to people during the experiment; rather it is a list of objects written in the language of the model. A desirable improvement of the model should be introducing an object recognition algorithm able to detect the rectangles in the image and generate the correspondent object list in ACT-R language. This approach would make ACT-R visual system more similar to humans'. Moreover, it would make more scalable the process of creating instances for the experiments.

The purposes described above represent the goals of the tool. To achieve them, the software analyzes the image created for the user, processes it extracting the features needed by ACT-R and stores the information in dedicated data structures. Then, when ACT-R requests them, it communicates all the stored data to the cognitive architecture using a format compatible with ACT-R model.

3.3. Requirements

The requirements are grouped in the following sections according to the functional/non-functional classification.

3.3.1. Functional Requirements

The objective of the work is to design and implement a standalone software module which receives as input an image and is able to analyze it and extract some features of it.

The input image is a color image which contains simple shapes. The shapes must not be overlapping and must have the same color hue.

The software must recognize simple shapes, in particular:

- triangles;
- rectangles;
- quadrilaterals;
- circles.

For each shape, it must calculate:

- area;
- perimeter;
- dimensions;
- rotation;
- a rectangular bounding box;
- center;

Moreover, the software has to:

- recognize the color of a single pixel;
- recognize the color of a shape;
- calculate distances between objects;
- make dimensional comparisons between objects;
- calculate the relative position of one object in respect with another one.

The software must be able to communicate with ACT-R. In particular, ACT-R must signal to it which is the image to process and ask for the features to be extracted. The module must return all the information extracted by ACT-R.

3.3.2. Non-Functional Requirements

Product Requirements

The software module must:

- be multi-purpose;
- be portable;
- work in background;
- communicate with ACT-R.

Multi-Purpose The software should be easily adapted to work with all the experiments that will be put in place in the future. Moreover, there should be the possibility to adapt the software in order to use it as a shape recognition tool in the orientation inside a building. For this topic, see [3.1.1](#).

Portability The software has to work on several operating systems.

Communication with ACT-R The software must transmit the extracted information to ACT-R in a standardized format.

Organizational Requirements

- The language of the implementation of the software must be C++;
- the computer vision library to be used must be OpenCV;
- a strict monitoring of the work is required.

4. Development Process and Tools

This chapter, after having introduced the agile development framework *Scrum*, describes the development process adopted by the team in order to complete the work.

4.1. Scrum

Scrum is a framework for the agile management of the project development. As such, it does not define the technical way in which the developers must do the job but it follows the development process. The framework has many dimensions: roles, events, rules and artifacts within the framework serve specific purposes and contribute the Scrum's success [Botb].

The following sections describe one by one these components. The first section describes the main roles of people within the framework, the second focuses on the events and the meetings necessary for the implementation of the methodology and the third one puts in evidence the artifacts needed.

All the information about Scrum are taken from [Botb]. See it for a more detailed description.

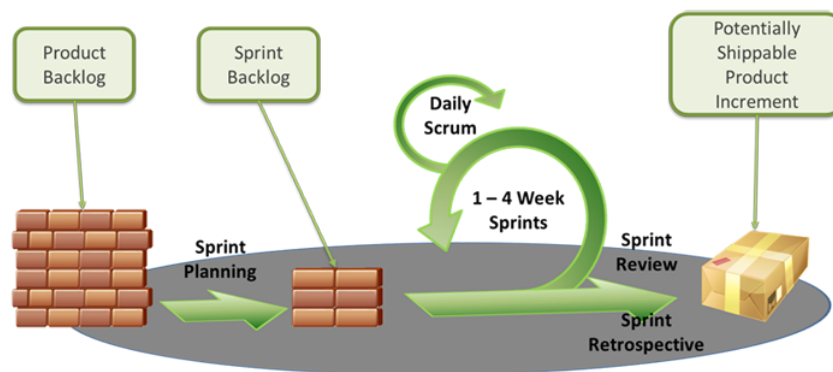


Figure 4.1.: Overview of the Scrum Process

4.1.1. The Scrum Team

The *Scrum Team* is composed of a *Product Owner*, the *Development Team*, and a *Scrum Master*. These figures have different roles but all of them have to reach the same goal, add to the product a finite and usable set of new features at constant time intervals.

Scrum Teams are self-organized and cross-functional. Self-organization gives the members the possibility to choose the best way to accomplish their work, without being directed by others outside the team. This fact not only avoids the team to be directed by someone who has not the technical knowledge to solve specific problems and, consequently, can not find the best ways to achieve the goals, but also allows the Team to optimize its own development process. Moreover, the consequent improvement of the sense of ownership that the Team feels on the product increases the motivation of the members. Cross-functionality gives the team all the competencies needed to accomplish the work without depending on others not part of the team. In this way, all the tasks can be accomplished directly by the team, thus the team is never blocked waiting for the work of someone outside it.

The delivery of the products is iterative and incremental. This fact guarantees a constant feedback on the correctness of the job and ensures that a deliverable version of the product is always available.

The following paragraphs will describe each of the three roles.



Figure 4.2.: The Scrum Team

Product Owner

The *Product Owner* is the person who represents the stakeholders of the product and is responsible for the performance of the team. His role is to define the requirements for

each new version of the product, assign them the priorities and explain them to the team in detail. The requirements are called *Backlog Items* and are included in the *Product Backlog*. This document is described in more details in section 4.1.3.

Scrum rules require a person to assume the role of Product Owner to represent the will of a committee. In this way, he assumes an intermediation role for all the communications between the developers and the stakeholders and becomes a reference point for both these categories. On one side, he has to give the stakeholders all the information about the progress of the development and the performances of the team. On the other hand, he has to decide the priorities of the tasks and deliver the team the messages of the stakeholders.

Development Team

The *Development Team* consists of professionals programmers, whose role is to implement the new functionalities to the product. The team is usually composed of a number of members which varies from three to nine. One team must be self-organizing, this means that only the members can decide step by step the tasks in which the items are split and how to accomplish each of them in order to add functionalities to the product. Every team is also cross-functional, i.e. is composed by people who have different skills. In this way it can be autonomous and it does not have to depend on other people outside the team to accomplish its job. The more synergy the Development Team has, the higher its overall efficiency and effectiveness are.

Scrum Master

The *Scrum Master* has the role to verify that the Scrum process is understood and put in place. He does this by checking that everyone in the team follows Scrum theory, practices, and rules. The Scrum Master is the enforcer of the rules and interacts with all the people inside and outside the Scrum Team in order to avoid useless interactions. On one side, he helps the Product Owner finding techniques for managing the *Product Backlog*, teaching him how to communicate in clear way with the Development Team and understanding and practicing the agile development. On the other hand, he coaches the Development Team in self-organization and cross-functionality, he protects it from unhelpful interruptions and keeps it focused on the tasks. For this, often this role is referred as a "*servant-leader for the Scrum Team*" [Botb, p.6].

4.1.2. Events

Prefixed meetings in Scrum have the purpose of marking time and hence minimizing the need for meetings not defined in Scrum, which can distract the members of the team from their tasks. In this events, different roles in the Scrum Team can interact with regularity, without the need for continuous interruptions. All the events are time-boxed, i.e. every event has a maximum duration. This ensures that only the appropriate amount of time is spent planning, avoiding wasting time.

The following sections describe the *Sprint*, the *Sprint Planning Meeting*, the *Daily Scrum*, the *Sprint Review* and the *Sprint Retrospective*.



Figure 4.3.: The Scrum Events

Sprint

The *Sprint* is the primary unit of development in Scrum. It is a time-box with a prefixed duration and can last from one week to one month. During this period of time the team creates a self-contained portion of product.

Each Sprint starts with the *Sprint Planning Meeting* and ends with the *Sprint Review Meeting* and the *Sprint Retrospective Meeting*. All these events will be described with more details in the next sections.

Sprint Planning Meeting

The *Sprint Planning Meeting* is a meeting at which all the member of the Scrum Team participate. It is the beginning of every Sprint. In this meeting the team decides the new

features the product will include by the end of the Sprint.

At first the Product Owner informs the team about the features that he wants to be added to the product, each one with its priority. Each one is an item in the *Product Backlog*, which is described in the section 4.1.3. Then, the Development Team estimates how long it would take to add every new feature to the product and, consequently, how many of them will be added by the end of the Sprint. Often in this part of the meeting the Product Owner and the Development Team have a talk in order to define better the requirements. This fact guarantees that the developed features adhere strictly to the requirements and leads to a new time estimation for most of the goals.

The goals are split into tasks, each of which takes no more than two days to be accomplished by the Development Team. Every task and the delivery plan are collected in the *Sprint Backlog*, described in the section 4.1.3.

Every Sprint must have a *Sprint Goal*. The Sprint Goal acts like a sort of motivation which reminds the Development Team during the whole Sprint which is in a wider context the goal of their activities. For this reason the Sprint Goal should not be changed during the Sprint.

Daily Scrum

The *Daily Scrum* is a fifteen minutes event for the Development Team and the Scrum Master, that allows the team to synchronize the work and plan the next day activities. This purpose is achieved analyzing the work done in the last day and forecasting the work for the next one. During the meeting, each Development Team member explains what he did since the last meeting, what he is going to do before the next meeting and the problems he had to front. The meeting has the purpose of evaluating the progress towards the Sprint Goal and analyzing the trend of the progress in comparison with the Sprint Backlog. This allows the Scrum Team to measure its speed and, consequently, to make better time estimations in the next Sprints.

Sprint Review

The *Sprint Review* is a time-boxed meeting which closes every Sprint. In this meeting the increment of the work is analyzed and, if needed, the Product Backlog is updated. During the meeting, the Product Owner analyzes what has been completed and what has not. The Development Team describes the problems it encountered, how it solved them and shows the new functionalities added to the product. Then the whole team discusses of new features to be added, old ones to be deleted or how to improve other ones and, more

generally, collaborates on what to do next, updating consequently the Product Backlog.

Sprint Retrospective

The *Sprint Retrospective* is a time-boxed meeting that takes place after every Sprint Review. It is an opportunity for the Scrum Team to analyze its Scrum implementation and create a plan to improve the next Sprint. During the meeting the focus is set on people, relationships, processes and tools. The most important successes are shown and a plan for implementing potential improvements is created. The purpose of this meeting is to make the implementation of the Scrum method more effective by optimizing the development process and introducing techniques that make the method more productive and enjoyable.

4.1.3. Artifacts

Scrum's *artifacts* represent the work in many different ways. The basic artifacts required by the framework are the *Product Backlog* and the *Sprint Backlog*.

Product Backlog

The *Product Backlog* is an artifact which contains the list of the requirements for a product. The elements of the list, called *Backlog Items*, are sorted by priority.

Every Item must have a name, a description, an estimate and a priority. For each Item, the Product Owner, representing the stakeholder, sets the priority and the Development Team defines the estimate.

The Product Backlog is dynamic. Its earliest versions only contain the initially and best understood requirements of the product. Then, as long as the product evolves and new requirements are introduced, it is continuously updated. In this way, it always contains all the features, enhancements and fixes that must be added in the future to the product.

As the higher ordered Items are more important than lower ordered ones, they are clearer and more detailed. The team can make more precise estimates thanks to the greater clarity and increased detail.

The activity of adding detail, estimates, and order to the items in the Product Backlog is called *grooming*. The Scrum Team decides how and when to do it. Anyway the grooming activity must not consume more than the 10% of the capacity of the team.

Sprint Backlog

The *Sprint Backlog* contains the Product Backlog Items selected for the current Sprint and a plan for completing all the modifications to the product and realizing the Sprint Goal. The Sprint Backlog contains all the functionalities that should be added to the product by the end of the Sprint, each of which with its estimates and priorities, and, at the same time, monitors the state of the work during every day of the Sprint.

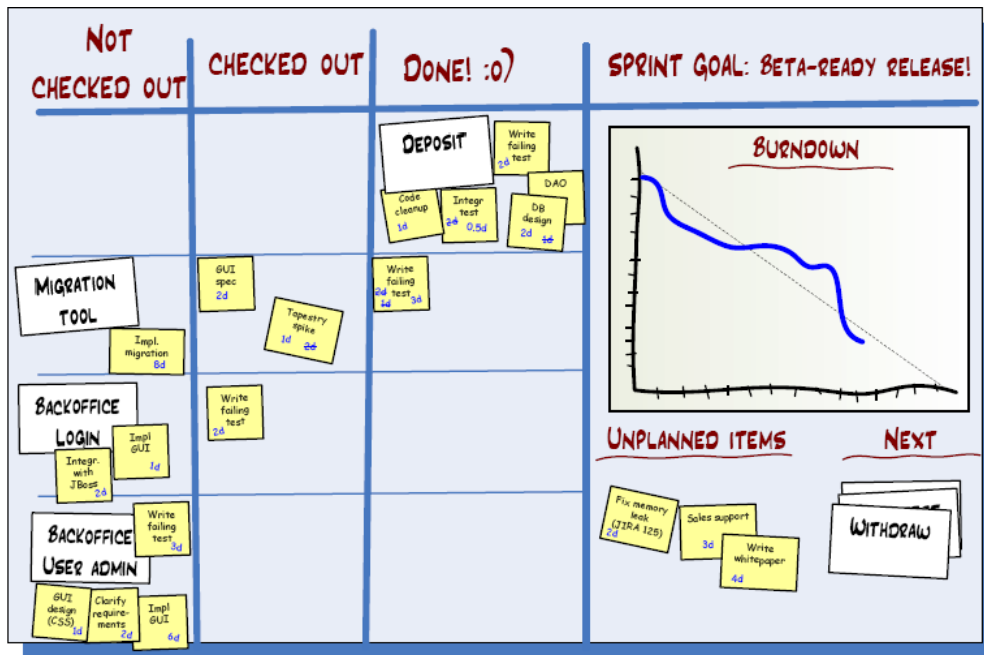


Figure 4.4.: Sprint Backlog and Burn Down Chart

The level of detail of the work is such that it can be measured during the Daily Scrum. In this meeting, in fact, the amount of work done in the day is calculated and the activities for the next day are planned. Usually every Product Backlog Item is split into tasks, each of which must be accomplished in not more than sixteen hours of work. The tasks are never assigned; rather, each member of the team takes charge of some of them during the Daily Scrum, according to the set priority and his own skills.

The Sprint Backlog is updated every day during the Daily Scrum. If some new work emerges to be necessary in order to complete some requirements, the Development Team adds it to the Sprint Backlog. If there are some elements which are deemed unnecessary, they are removed.

The state of the ongoing activities can be monitored by a *Burn Down Chart*, which measures the number of completed tasks per day and compares it with an ideal trend, which represents the estimates of the tasks. This chart allows to analyze the accuracy of

the estimates and to have a visual representation of the ongoing work. The knowledge of the speed of the team allows the team to make better estimations in the future Sprints.

4.2. The Adopted Development Process

The development team decided to adopt an incremental and iterative development process, implementing the *Scrum* framework.

The main motivation of this choice is the capacity of this method to front very well the modifications of the requirements. In most of cases, in fact, even if the requirements at the beginning seem to be well defined and clear, they change during the development. For example, it could happen that new requirements reveal to be more important than other ones or that after some clarifications between the client and the developers some features are redefined, some are deleted and some others are added. The possibility of an incremental and iterative development process to make continuous reviews of the product faces better such changes in the requirements.

Moreover, the introduction of meetings in which is possible to discuss the state of the work and the working method improves considerably the performances of the team itself. In this way, in fact, the supervisors of the development team can postpone all the not essential communications until the next meetings, without the need to interrupt his work nor the developers'. This leads to a more focused and less stressed work both for the developers and the supervisors.

The possibility to check the work at every Sprint increases the compliance between the system and the requirements. During the Sprint Reviews, the new increment of the software is shown. In this meetings, the clients give the Scrum Team a feedback about what they want and if the feature developed satisfy their need. So if there have been some misunderstandings in the comprehension of the requirements, they can be fixed in the next Sprints. This continuous product revision leads to a strict adherence between the system and the requirements.

The adopted development process allows the Product Owner and the Scrum Master to evaluate the performance of the team. As the features to be added and the estimates are decided during the Spring Planning Meeting, the Product Owner and all the members of the team have the plan of the whole Sprint. During the Sprint Review, the members can monitor how the actual work has been done in comparison with the plan. On a shorter period, the Daily Scrum is used to monitor the activity of the days, while the Product Backlog contains the history of all the modifications done on the software since it was born. For all there reasons, with Scrum, the performances of the team are constantly

measured.

In order to facilitate the Scrum method, the team decided to adopt *Git* as *Version Control System*. The features of this instrument allow the team to focus on the design and implementation activities, relieving it from additional "overhead" tasks. This word refers to all the tasks which are different from design and implementation, as for example, the management of concurrent accesses, the administration of different versions, the direction of backup and recovery operations. Moreover, the possibility of switching easily between different versions of code and analyzing the differences between them brings to a series of advantages for the developer. For example, it makes easier locate bugs in the code and allows the programmers to try different solutions and instantly revert the not working ones, all without interfering with the working code. Finally, it allows to monitor the amount of work and who does it.

5. Design

This chapter describes the design of the software, focusing on its main important aspects. The first section describes the architecture from a high-level point of view. The second one shows the organization of processes and threads. The third one explains the choices about the communication mechanism and the last one illustrates the architecture of the visual module.

5.1. The Architecture

This section gives a high level overview of the architecture of the software by introducing its main building blocks and the main modules they interact with.

Image 5.1 on the left side contains the vision module, composed by the sub-modules *Server* and *Visual Operations*, and, on the right, Lisp interpreters that run ACT-R models. The two parts communicate to each other thanks to a client-server communication mechanism based on TCP/IP protocol in which the Lisp interpreter represents the client and the visual module the server.

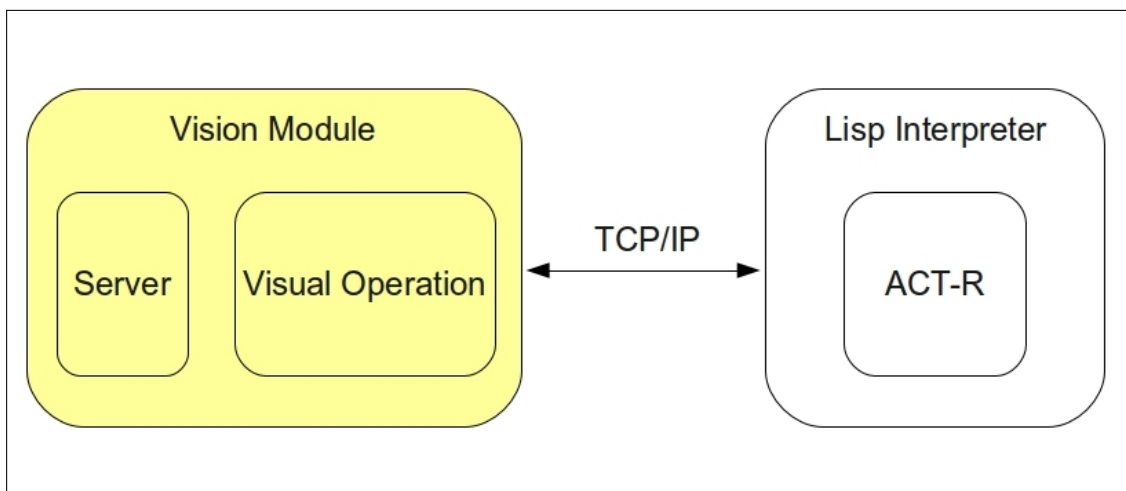


Figure 5.1.: High level scheme of the architecture

ACT-R part is not discussed in this work but is fundamental for the choices about the communication protocol. This document focuses on the *Vision Module*, whose main operations are extracting features from the input data and sending them to the models in a way that can be understood by the cognitive architecture. These functions are realized by the two sub-modules contained in Vision Module.

5.2. Processes

Figure 5.2 represents a scheme of processes and threads that compose the architecture. As shown in the picture, the Lisp interpreter and the Visual Module run on two separated processes. The processes are independent and their management is left to the operative system.

The design on the communication protocol is such that the two processes can be run only on the same computer. The images to be processed, in fact, at the moment are not sent as messages from the client to the server but are stored on the computer which runs both the processes. Anyway, it is enough to modify the message exchange by adding a message sent from the client to the server encapsulating the image data in order to run them on two different computers.

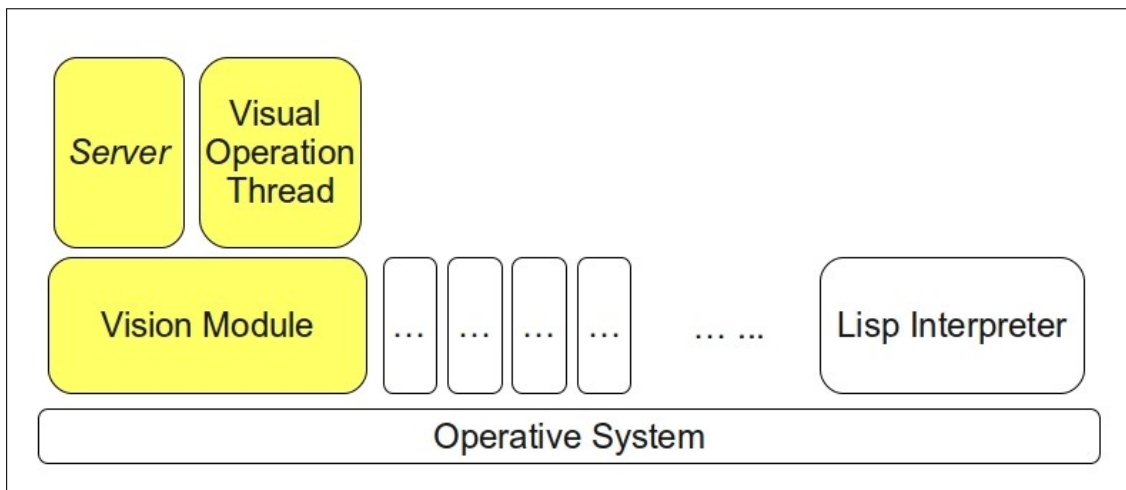


Figure 5.2.: Scheme of processes and threads of the architecture

The Vision Module is composed by two threads: one for the server and the other for the visual processing. The server thread runs in background while the visual operation one is the main thread. The problems regarding race conditions are dealt with *semaphores*; all

the other issues related to threads management are automatically handled by the *Boost*¹ library. The thread organization of ACT-R is not discussed in this work.

5.3. Communication

In order to exchange information between each other, ACT-R and the developed module use a synchronous *client-server* communication mechanism based on TCP. ACT-R models represent the clients, the developed software the server. The format used for the messages is JSON.

The client-server approach has been chosen because it is standardized and is supported by the interpreters that run ACT-R. Moreover, the time spent for coding and decoding the message is negligible if compared to the time needed for the image processing. Finally, the possibility of using shared resources and the transparency offered by the mechanism ensure very high flexibility and scalability.

Another solution analyzed for the communication between ACT-R and the Visual Module is compiling the Visual Module as a static library and then loading it directly in the Lisp interpreter. The latter is able to call the static functions thanks to *foreign functions*. In this way, the communication is more efficient because all the overhead introduced due to the inter process communication is avoided. This method, though, is not maintainable. In fact, every Lisp interpreter has a different method of calling foreign functions [SWI]. As consequence of this, every time the interpreter is changed, the interface of the visual module must be re-written. At design time it has been decided to prefer maintainability to an irrelevant improvement of the performances, thus the client-server has been chosen.

TCP has been preferred to UDP mainly because of the reliability guaranteed by that protocol. Secondly, the bi-directionality of the communication and the small quantity of data exchanged make it more suitable for the purpose.

The communication is synchronous. The client, once having sent the request, stays pending until it receives the response. This choice is coherent with the application domain: as the software module acts like a vision system for the cognitive architecture, it makes sense that the model does not continue the reasoning until the available visual information is received.

The Lisp interpreter which runs ACT-R models is the client, while the developed software is the server. When a model needs to acquire information about an image, it connects to the server. Once connected, it can encode and send a request to the server. After the

¹Boost is a set of libraries for C++ that support many issues including multi-threading. For more information see the official website at www.boost.org

server has received and decoded successfully the request, it processes the input data. When the visual information is ready to be sent, the server encodes it and send it back to the client. After this, the message is decoded and its content is translated in chunks and then delivered to the model, that uses it for the cognitive reasoning.

All the messages exchanged between client and server are encoded in JSON. Appendix [A.2](#) contains an example of the messages.

5.4. Vision Operation Module

This section describes the design of the Vision Operation Module. For reasons of clarity, the description of the architecture follows a top-down approach and makes use of diagrams. The first subsection gives an overview of the module and the following ones study in deep specific aspects of it.

To better explain the contents and the structure of the design, the description is supported by UML diagrams. The following sections describe the module at different levels of detail and the diagrams reflect this approach. In subsection [5.4.1](#), for example, they are used only to present how the classes are organized in the overall structure and that is why they contain only the names of the classes. Further details are added in the diagrams in the next subsections, according with the top-down approach of the chapter.

The reader will notice that some classes that appear in the following diagrams are not described in the chapter. The reason of this is that they do not accomplish to any functions described in the requirements section of chapter [3](#). Their presence makes sense in the context of the whole software, described in [3.2](#). It has been decided not to remove such classes from the diagrams because they influence the architecture and motivate particular design and implementation choices.

This section explains only the most important architectural aspects of the software. Not particularly interesting classes, like, for example, the utility ones are not described here. For a more detailed description about all the aspects of the architecture, see appendix [A.1](#).

5.4.1. Architecture

The vision operation module is composed of three main parts, one for processing images and videos, one for containing the hierarchy of the objects recognized during such processing and one containing the exception classes and the utility functions used by all the other components.

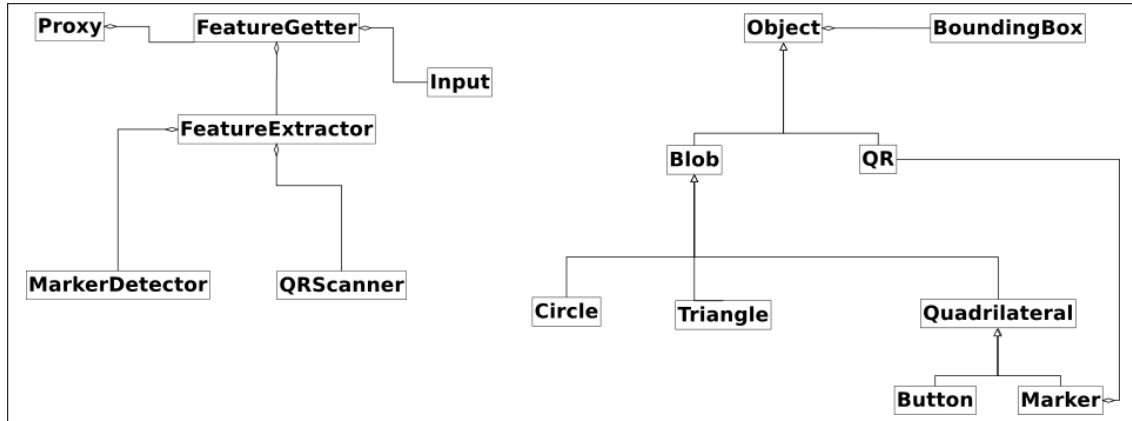


Figure 5.3.: Class diagram illustrating a high level overview of the vision operation module

Image 5.3 shows a class diagram which describes two of the parts that compose the software. The part dedicated to the processing of images and videos is located on the left and is composed by the classes *Proxy*, *FeatureGetter*, *FeatureExtractor*, *Input*, *MarkerDetector* and *QRScanner*. The object hierarchy, which is shown on the right, contains the classes *Object*, *BoundingBox*, *Blob*, *QR*, *Circle*, *Triangle*, *Quadrilateral*, *Button* and *Marker*.

The third part is not included in the section because it contains only utility classes and functions. The reader can find a description of that part in appendix A.1.

5.4.2. Feature Extraction Classes

The classes that play the role of acquiring the images and extracting features from them are *Input*, *FeatureGetter* and *FeatureExtractor*. *Proxy* has the role of intermediation between this part of the software and the one which manages the communication with ACT-R.

Figure 5.4 contains the class diagram that designs the interface of these classes. *FeatureExtractor* represents the core of this part. It receives in input an image and processes it in order to extract the features that contains. For each extracted feature it creates a correspondent element in the *object hierarchy*. Such hierarchy is described in section 5.5. Once the objects are generated, the class creates a list containing them and delivers it to *FeatureGetter*. *FeatureExtractor* is also responsible for accomplishing other goals described in the requirements, like recognizing colors of pixels and objects, calculating the distances of two objects as well as the relative position of one object in respect with another one and making dimensional comparisons between two objects. These features

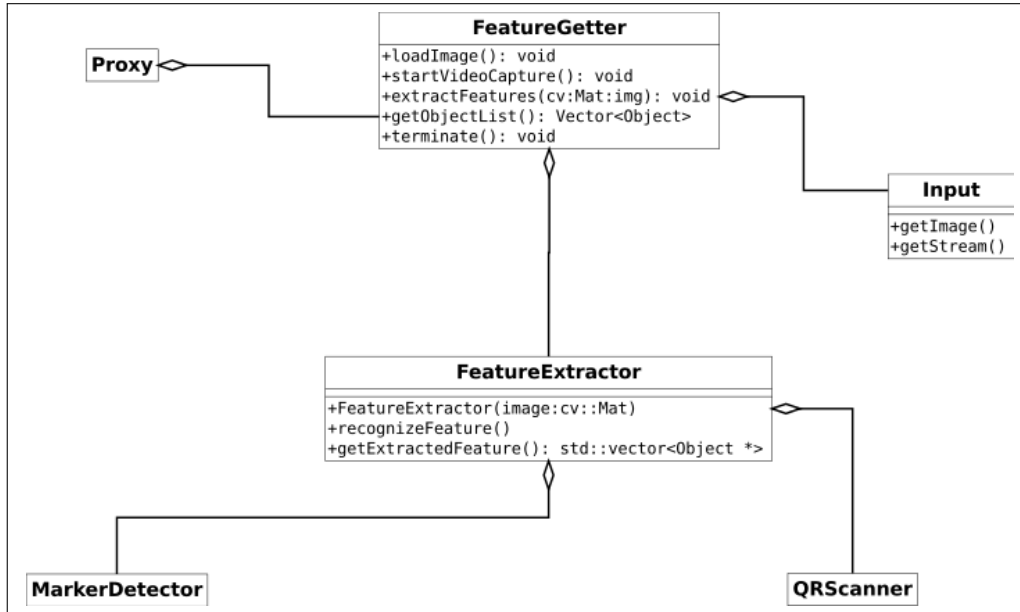


Figure 5.4.: Class diagram illustrating the main roles in the feature extraction process

have not been included in the initial design of the class.

Input is responsible for interacting with the file system and the hardware of the computer. In particular its tasks are loading images from files, saving images to files and loading video streams from the web-cam. The acquired data are delivered to *FeatureGetter*.

FeatureGetter represents a level of abstraction that avoids to call directly the functions of *FeatureExtractor* and *Input*, providing a simpler interface. Its functions create the *Input* object, receive from it the input data, delivers it to *FeatureExtractor* for the extraction of the features and delivers the received object list to *Proxy*.

Proxy represents a further level of abstraction between *FeatureGetter* and the part of the software which manages the communication with ACT-R. Most of the procedures for the communication with the cognitive architecture, in fact, need to use specific data types. The main function of this class is to transform each object of the object list in the type necessary for the communication to ACT-R.

As the classes *MarkerDetector* and *QRScanner* do not accomplish to any functions described in the requirements section of chapter 3, they are not described in this work.

5.5. Class Hierarchy

The class hierarchy contains all the objects that are recognized during the processing of the image. Each feature has a correspondent object that describes it.

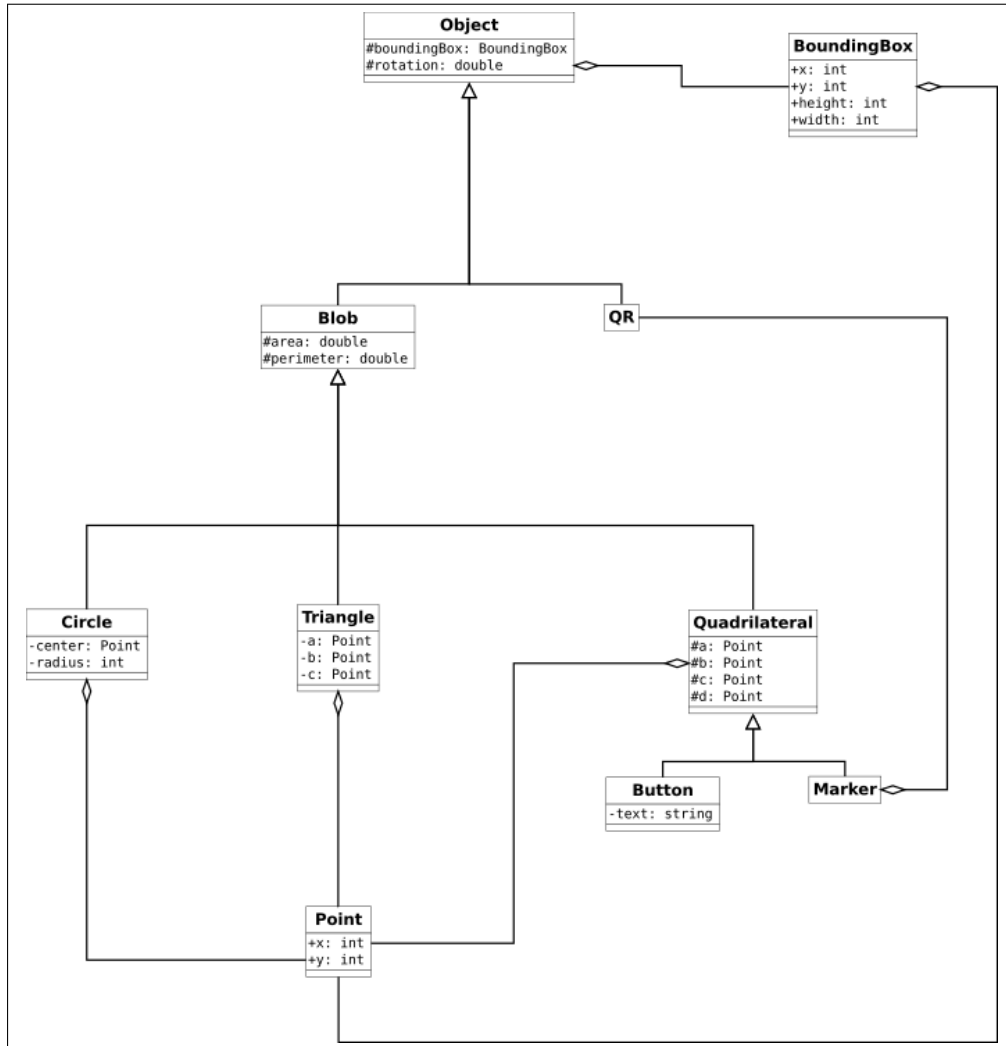


Figure 5.5.: Class diagram illustrating the class hierarchy of the recognized objects

Figure 5.5 contains a class diagrams that specifies such hierarchy. *Object* is the super-class that represents all the possible objects recognized in an image during the feature extraction process. All the elements of this type must have two attributes: a *rotation* value and a *boundingBox*. The latter is defined by the class *BoundingBox*. Such element can have only rectangular shape, in fact its attributes are the coordinates of a point and horizontal and vertical dimensions.

The type *Object* can be specified by two sub-types: *Blob*, that represents a generic

shape, and *QR*, representing a QR element. The presence of QR is not important for the shape recognition task but it justifies the presence of Blob in the hierarchy. Without QR class, in fact, Blob could have been incorporated in Object. Blob is provided of the attributes *area* and *perimeter*, which, as suggested by the names, represent respectively the area and the perimeter of the shape.

The classes *Circle*, *Triangle* and *Quadrilateral* extend *Blob*. They represent the simple shapes that can be recognized during the recognition process. Each item is defined starting by basic geometrical elements: the circle by center and radius, the triangle by three points, the quadrilateral by four points. *Point* class, in particular, has been introduced in this diagram for reasons of completeness.

The class *Button* specifies *Quadrilateral* by adding the attribute *text*, which represents the message of such button. The class *Marker* is not described in this document.

6. Implementation And Testing

This chapter describes the solutions implemented for recognizing shapes, and presents the way in which the testing operations have been conducted.

6.1. Implementation

As regarding the detection of shapes in images, two different algorithms are implemented, one for circles and one for shapes of n sides. In particular, the solution for recognizing circles makes use mainly of a function provided by the OpenCV library, therefore it is not described here. The algorithm for detecting n -sides shapes, instead, is more complicated and not directly implemented in the vision library. Below, the reader can find the description of such method.

6.1.1. N-sides Shape Detection Algorithm

The algorithm used in order to detect n -sides shapes in the images is based on a procedure of corner detection. In order to capture all such objects, the procedure is executed three times for each color level of the image, every time with different parameters. As the corner detection algorithm identifies *false corners*, it requires functions for detecting and removing them. In addition, the repetition of the procedure leads to the detection of many sets which describe the same shape. Consequently the algorithm requires procedures that remove the different copies of the same object.

Listing 6.1 contains the pseudo code of the algorithm. The first operation is a technique for reducing the noise in the image and is followed by two procedures, one for splitting the image in its three channels and the other one for detecting corners. The latter is applied with three different thresholds for each level of the colored image, red, green and blue. This guarantees that all the shapes in the image are recognized. These three functions are offered by OpenCV library.

After the corner detection procedure, an operation of selection is applied, based on the number of detected vertexes. In order to recognize a n -sides shape, only the sets

containing n or more vertexes are kept, the other are discarded. For the recognition of quadrilaterals, for example, all the sets with less than four elements are discarded.

```

function detect_polygon(image, sides_number){

    noise_reduction(image);

    image_levels[3] = split_in_level(image);

    corners = [];

    foreach level of image_levels do{
        foreach threshold of thresholds do{
            detected_corners[] = corner_detection(level,threshold);
            corners.add(detected_corners);
        }
    }

    foreach corner_set of corners do{
        if (corner_set.size() < sides_number)
            corners.remove(corner_set);
    }

    foreach corner_set of corners do{
        remove_in_line_points(corner_set);
    }

    foreach corner_set of corners do{
        if (corner_set.size() != sides_number)
            corners.remove(corner_set);
    }

    remove_duplicates(corners);

    return corners;
}

```

Listing 6.1: Pseudo code of the n -sides shape detection algorithm

The fact that not only the sets with n sides are included in the list is a consequence of the corner detection algorithm. It often happens, in fact, that, due to noise in the image or imperfection of the corner detection procedure, a single segment is detected as a set of

two or more shorter ones. Each couple of such elements defines a *false corner*. For this reason, even a set of points containing more than n items can define an n -vertex shape.

In order to delete the false corners, an algorithm analyzes every point set and calculates the slopes of the straight lines obtained by each couple of adjacent vertexes. If the slopes of two consecutive segments are very close one to each other, the vertex which is common to both the segments is a false corner and therefore is removed.

After this, another selection based on the number of vertexes is applied. This time all the sets not containing n vertexes are discarded.

Another problem of the algorithm is that, due to different thresholds applied to different levels of the image, a shape is often defined by more sets of points. Two points representing the same vertex of the same shape but belonging to different sets in general do not have the same coordinates. The gap in some cases can be even greater than 15 pixels. For example, the same triangle can be specified by $\{(0;0),(200;0),(100;200)\}$ and $\{(2;2),(199;0),(101;197)\}$.

In order to detect and remove these copies, a function has been introduced that checks all the possible couples of sets. It controls the similarity of two elements by analyzing the positions of the pixels and making comparisons between the areas and the perimeters of the two shapes. When two sets specify the same object, the second one is discarded.

Once removed the copies, the algorithm returns a list containing the remaining sets of points, each of which represents a shape of the desired number of sides.

At the moment this algorithm has been implemented and tested only for quadrilaterals and triangles. Anyway, it can be adapted for recognizing shapes with more than four sides. In fact, for each new category of shape, the only function to be modified is the filter for removing copies. A refactoring operation, anyway, can make this function work with general n -sides shapes and remove the need of adapting it for every new category of object. To maintain the coherence with the architecture, each new shape requires its own object in the class hierarchy, described in 5.5. The creation of a new class represents the only necessary implementation detail to be added in order to make the algorithm detect a new n -sides shape.

6.2. Testing

The testing activity depends strictly on the class creation. In particular, a class creation is always followed by the creation of its correspondent test module, which contains all the procedures required for checking the correctness of constructors, destructors and methods of the class to be tested.

Differently from the *test driven development* approach, in this case the test cases are not written before the definition of single methods; rather they are developed when the creation of the whole class is completed. As consequence of this, a test case does not necessarily fail the first time it is run.

Every method is tested ensuring that all the possible flows of the execution are followed. According to the *equivalent partitioning* technique, input data are grouped in different categories by following the taxonomy valid/not-valid and a test case with as input a representative value for each partition is created. Differently from the specification of this methodology, more than one value per category is used for the test cases instead of one.

Moreover, conforming to the *boundary value testing* technique, input values on the boundaries of data categories are introduced in the test cases in order to make the test cover as wide as possible.

The test cases are organized in test suites. All the suites specific of a class are included in the same test module. All the test operations are automatically realized thanks to the CUTE software.

The objective of the testing operation is to evaluate correctness and reliability of the system.

7. Conclusions

The objective of this work is to develop a software module which receives images as input and analyzes them by extracting particular information like the shapes they contain, their colors and their positions and communicates them to ACT-R when the cognitive architecture requests it. The ultimate purpose of this is to create a vision module for the cognitive architecture in order to make its visual perception as similar as possible to human one.

The developed software uses OpenCV library to implement the visual operations and communicates the information to ACT-R thanks to a client-server communication method. This solution is simple, fast and allows to use a library which provides ready and tested functions for computer vision.

The software at the moment recognizes quadrilaterals, circles and rectangles and makes qualitative and quantitative evaluations about colors and relative positions and dimensions of the objects. In particular, the tool is used in the Rush Hour experiment for analyzing the input image, extracting the list of objects it contains and communicating it to ACT-R. In this way the cognitive architecture does not start working from a list of objects but it processes an image, even if indirectly. This fact approaches the ACT-R perception to human one.

The restricted number of features offered by the visual module at the moment represents the biggest limit of this solution. However, it is enough to add new functions to this tool in order to extend more and more the scope of this module and use it increasingly together with ACT-R.

In order to adapt the system to work with other experiments, the first improvements that can be done are adding the recognition of new kinds of shape, for example ellipses, and a tool for *Optical Character Recognition*¹. In fact, at the moment, images used in the experiments conducted with the cognitive architecture, in general are quite simple. They normally contain only shapes and may include some text. The introduction of methods for recognizing new shapes and texts make the system more scalable with the definition

¹An Optical Character Recognition, also known as OCR, is a software which recognizes and decodes text in images.

of new experiments that need to process images.

Adding new functions for object recognition can provide ACT-R with the ability of processing images directly in the real world. This can lead to interesting applications in many research fields, including the spatial reasoning. For example, the team developed a system for robots navigation with landmarks inside buildings. The robot, together with the vision module, is able to identify a landmark, move close to it, decode the information it stores and send it to ACT-R, which can then communicate the robot how to behave.

In conclusion, the innovation of this software module consists in providing ACT-R with the possibility of processing images directly in the real world. This overcomes one of the biggest limits of the cognitive architecture: working in a virtual environment. Shifting the vision tasks to the computer vision library brings an immediate advantage: it allows to use ready and tested solutions to accomplish many visual tasks. Moreover, as the computer vision research field is very active and it produces continuously new solutions, always new possibilities will be available for visual experiments with the cognitive architecture. For all these reasons, the development of this module opens a wide range of possibility for the research in the field of cognitive science.

A. Appendix A

The appendix contains further details about the implementation. The first section adds information about the architecture of the software, the second one defines the format of the messages exchanged between client and server and the third one lists all the tools used during the development.

A.1. Details About The Architecture

This section describes aspects about the design of the architecture which have been omitted in chapter 5. In particular, it deepens the architecture of the Visual Module. At first it gives an overview of the whole classes that compose such module and, then, it focuses on the aspects not described in the design chapter: the architecture of the Server Module and of the Utility Part of Visual Module.

Again the description is helped by class diagrams. This time, such diagrams are obtained from the source code thanks to a tool of reverse engineering.

A.1.1. Overview

Figure A.1 shows the class diagram of the Vision Module. According to the design specifications defined in chapter 5, the module is composed by the sub-modules *Server* and *Visual Operation*. The latter is split in three parts, each one with its own tasks.

Classes *Server* and *Session* represent the Server sub-module. All the other classes compose the Visual Operation one.

Classes *Proxy*, *FeatureGetter*, *Input*, *FeatureExtractor*, *MarkerDetector* and *QRScanner* have the role of extracting features from the input data and preparing the messages to be sent to the cognitive architecture. *Object*, *BoundingBox*, *Blob*, *QR*, *Circle*, *Triangle*, *Quadrilateral*, *Marker* and *Button* form the object hierarchy. *InputException*, *NotOverlappedSegmentException*, *ParallelLinesException*, *VerticalLineException*, *CoincidentLinesException*, *FourPointsSorter*, *StraightLine*, *Segment* and *Point* represent the utility classes.

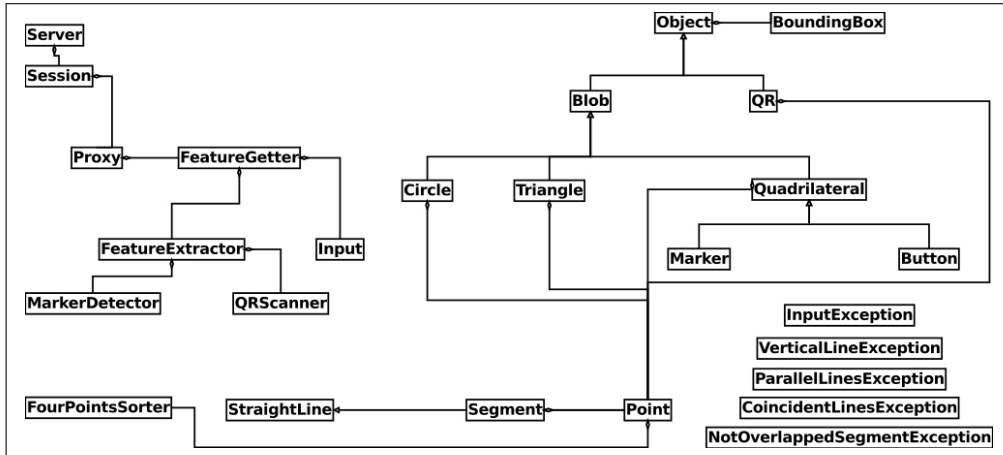


Figure A.1.: Class diagram which illustrates the whole architecture of the Vision Module

A.1.2. Server Module

The sub-module Server, responsible for the communication with the Lisp interpreter, is composed by the classes *Server* and *Session*. Figure A.2 contains the classes that implement this service.

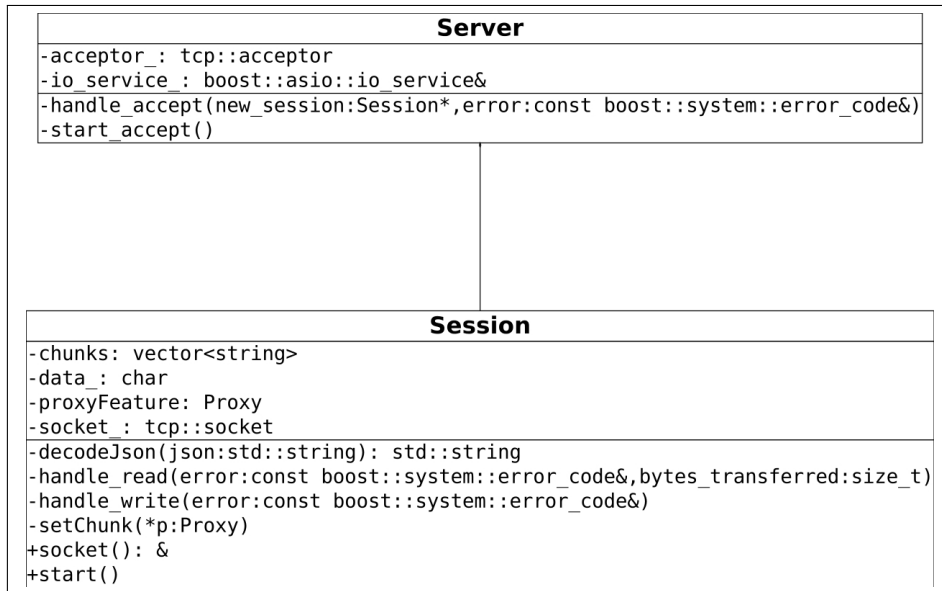


Figure A.2.: Class diagram describing the implementation part of the software dedicated to the communication

An instance of the class *Server* is running in background and listening for new tcp connections. When a client asks for a connection, the *boost::asio::io* service calls the *Server* callback, which creates a new instance of *Session*. The Boost library is responsible

for managing either Server and Session. The two classes have been developed following the Boost library specifications. When the Session object receives a message requesting the extraction of the features, it creates an instance of Proxy and calls the proper methods of this object in order to complete its task.

A.1.3. Utility Part

The classes which form the utility parts of the software are divided in two categories: the exception classes and the ones representing all the other utilities. To these, two modules containing functions are added. Figure A.3 shows a class diagram containing all the classes that constitute this part of the software. The modules, as not classes, can not be shown in the diagram but the prototypes of the functions they contain are listed below.

```

    cv::vector<cv::vector<cv::Point> > deleteOverlapped(cv::vector<cv::vector<cv::Point> >
oldList);
    cv::vector<cv::vector<cv::Point> > squaresSort(cv::vector<cv::vector<cv::Point> > squareList
);
    cv::vector<cv::vector<cv::Point> > deleteFalseSquares(cv::vector<cv::vector<cv::Point> >
oldList);
    cv::vector<cv::Point> sort4PointsClockwise(cv::vector<cv::Point> points);
    cv::vector<cv::vector<cv::Point> > deleteFalseTriangles(cv::vector<cv::vector<cv::Point> >
oldList);

    double angle( cv::Point pt1, cv::Point pt2, cv::Point pt0 );
    double getDistance(cv::Point a, cv::Point b);
    bool similar(cv::vector<cv::Point> a, cv::vector<cv::Point> b);
    bool tooClose(cv::Point a, cv::Point b);
    bool isFalse(cv::vector<cv::Point> square);
    bool isTriangle(cv::vector<cv::Point> terna);

    int getMinMax(const std::vector<Point>& coords, int type);
    std::vector<Point> sort4PointsClockwise(std::vector<Point> points);
    double erone(Point a, Point b, Point c);
    bool inLinePoints(int ax, int ay, int bx, int by, int cx, int cy );
    double getDistance(Point a, Point b);
    bool areSame(double a, double b);
    double getMin(const std::vector<double>& values);
    bool fileExists(char * path);
    string intToString(int a);

```

Listing A.1: Prototypes of the functions included in *extractorUtils.h* and *utils.h*

Classes *InputException*, *NotOverlappedSegmentException*, *ParallelLinesException*, *VerticalLineException* and *CoincidentLinesException* are the exceptions classes. As such, they are used to handle all the anomalous or exceptional events occurring at run-time.

StraightLine, *Segment* and *Point* represent geometrical entities. These elements serve as basic structures for creating the objects in the hierarchy and making evaluations between

couples of objects. Such evaluations can be both qualitative and quantitative, as, for example, the distances between the centers of two objects or the relative positions of one object in respect with another.

FourPointsSorter and the function modules contain the basic methods which are fundamental for all the operations executed by all the other classes of the software. For reasons of completeness the prototypes of the functions are listed in listing A.1.

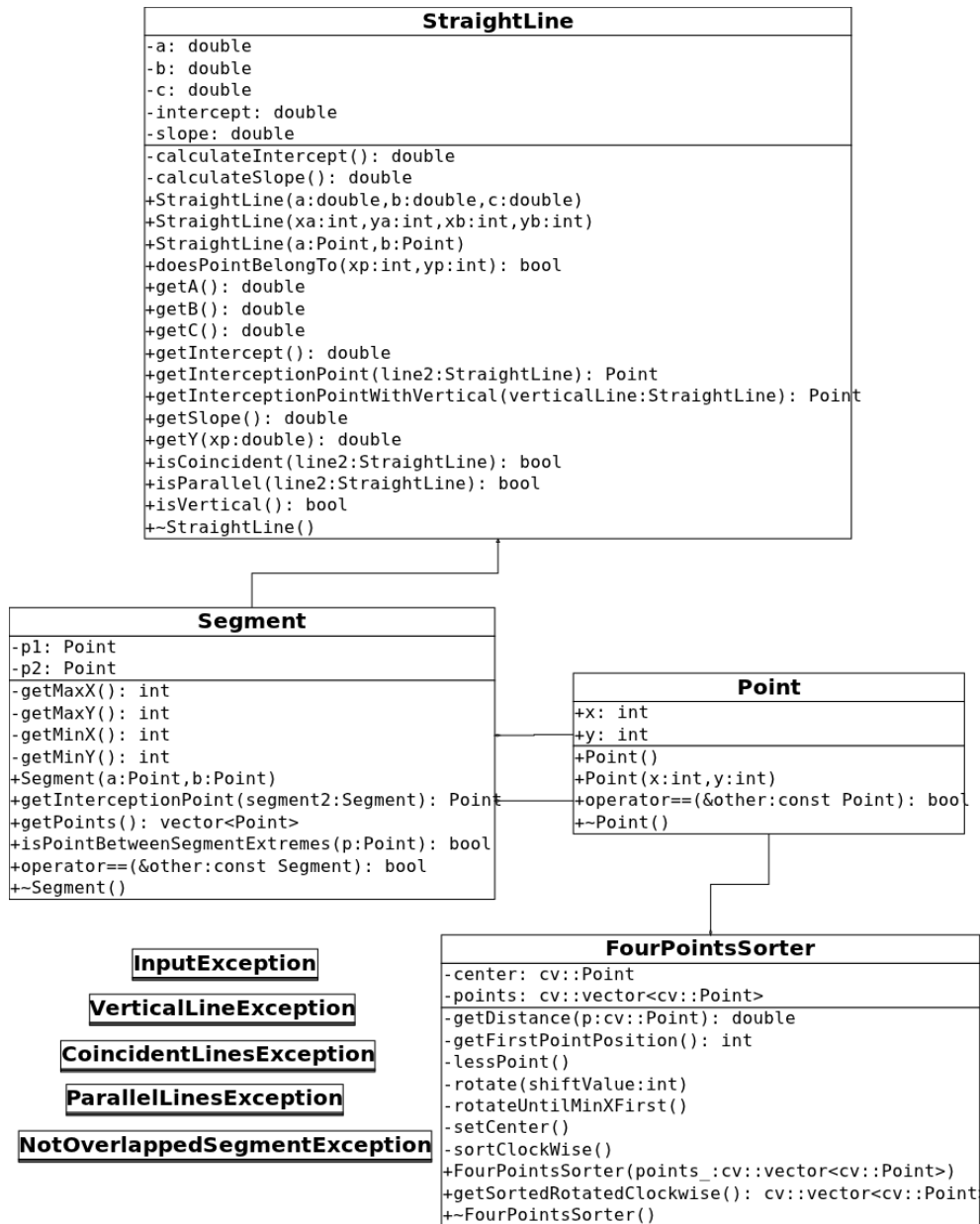


Figure A.3.: Class diagram showing the utility part of the Vision Module

A.2. Message Structure

During the communication between the Lisp interpreter and the Vision Module, two types of message are used:

- the requests, sent from the client to the server;
- the responses, sent from the server to the client.

A.2.1. Request Format

The request message consists of two parts, one fixed and one variable. The fixed part is the attribute `cmd`. The variable part is the value. The only meaningful value for this work is `getFeature`. This command tells the server to extract all the possible features from the input image. Listing A.2 shows the message for this request:

```
{"cmd": "getFeature"}
```

Listing A.2: Request Message

A.2.2. Response Format

The response message contains a list of all the shapes recognized in the input image. Each item of the list is represented by a list of attributes:

- **Bbox**: a rectangular bounding box, defined by two points;
- **Color**: the color of the shape;
- **Type**: the type of the shape;
- **Vertices**: the vertexes of the shape, if present.

Listing A.2.2 shows an example of response, in which a rectangle and a triangle are detected:

```
[
  {
    "Bbox": [
      {
        "x": 30,
        "y": 7
```

```

        },
        {
            "x":44,
            "y":22
        }
    ],
    "Color":"green",
    "Type":"Quadrilateral",
    "Vertices":[
        {
            "x":30,
            "y":7
        },
        {
            "x":30,
            "y":21
        },
        {
            "x":44,
            "y":22
        },
        {
            "x":44,
            "y":7
        }
    ]
},
{
    "Bbox":[
        {
            "x":4,
            "y":7
        },
        {
            "x":21,
            "y":21
        }
    ],

```

```
    "Color": "red",
    "Type": "Triangle",
    "Vertices": [
        {
            "x": 12,
            "y": 7
        },
        {
            "x": 4,
            "y": 21
        },
        {
            "x": 21,
            "y": 21
        }
    ]
}
```

Listing A.3: Example of Response Message

A.3. Tools Used in The Development

In addition to ACT-R and OpenCV, the tools used for the development are:

- **Eclipse**¹ CDT as IDE for writing C and C++ code;
- **CUTE**² as testing tool;
- **Git**³ as distributed version control system;
- **Trac**⁴ as web-based project management and bug tracking system;
- **Make**⁵ as build automation tool;

¹Eclipse: www.eclipse.org

²CUTE: www.cute-test.com

³Git: www.git-scm.com

⁴Trac: trac.edgewall.org

⁵Make: www.gnu.org/software/make

- **Boost**⁶ library for managing all the issues related to multithreading;
- **Dia**⁷ as tool for creating and reading UML diagrams;
- **cpp2dia**⁸ as tool for creating automatically UML diagrams from source code;

⁶Boost: www.boost.org

⁷Dia: www.dia-installer.de

⁸cpp2dia: cpp2dia.sourceforge.net

References

- [ABB⁺04] John R. Anderson, Daniel Bothell, Michael D. Byrne, Scott Douglass, Christian Lebiere, and Yulin Qin. *An integrated theory of the mind*. PSYCHOLOGICAL REVIEW, 111:1036–1060, 2004.
- [Aga06] Gady Agam. *Introduction to programming with OpenCV*. Technical report, 2006. www.cs.iit.edu/~agam/cs512/lect-notes/opencv-intro/opencv-intro.html.
- [And76] J.R. Anderson. *Language Memory Thought*. The Experimental Psychology Series/ Arthur W. Melton consulting ed. Taylor & Francis Group, 1976.
- [And93] J.R. Anderson. *Rules of the Mind with Mac Dis*. Taylor & Francis Group, 1993.
- [BEAA09] A.D. Baddeley, M.W. Eysenck, M.C. Anderson, and M. Anderson. *Memory*. Taylor & Francis Group, 2009.
- [BK08] G. Bradski and A. Kaehler. *Learning OpenCV: Computer Vision with the OpenCV Library*. O’Reilly Media, 2008.
- [Bota] Dan Bothell. *ACT-R 6 Reference Manual*. <http://act-r.psy.cmu.edu/actr6/reference-manual.pdf>.
- [Botb] Dan Bothell. *Scrum Guide 2011*. www.scrum.org/Scrum-Guides.
- [Fau93] O. Faugeras. *Three-Dimensional Computer Vision: A Geometric Viewpoint*. Artificial Intelligence Series. Mit Press, 1993.
- [FP11] D.A. Forsyth and J. Ponce. *Computer Vision: A Modern Approach*. Pearson, 2011.
- [HZ03] R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge books online. Cambridge University Press, 2003.

-
- [Lag11] R. Laganière. *OpenCV 2 Computer Vision Application Programming Cookbook*. Packt Publishing, Limited, 2011.
- [Leg05] P. Legrenzi. *Prima lezione di scienze cognitive*. Universale Laterza. Prime lezioni. Laterza, 2005.
- [New94] A. Newell. *Unified Theories of Cognition*. The William James Lectures. Harvard University Press, 1994.
- [Ope12a] Opencv change logs, 2012. opencv.willowgarage.com/wiki/OpenCV.
- [Ope12b] Opencv documentation page, 2012. docs.opencv.org.
- [Ope12c] Opencv web page, 2012. opencv.willowgarage.com.
- [r8C] Project r8 -[cspace]. <http://www.sfbtr8.spatial-cognition.de/R8>.
- [Riz12] Enrico Rizzardi. *Robotica cognitiva: aspetti percettivi e cognitivi nella navigazione con punti cospicui*, 2011/2012.
- [Sea02] Andrew Sears. *The Human-Computer Interaction Handbook: Fundamentals, Evolving Technologies and Emerging Applications*. Lawrence Erlbaum, 2002.
- [SWI] Swig documentation page. www.swig.org/Doc1.3/Lisp.html.
- [TV98] E. Trucco and A. Verri. *Introductory Techniques for 3-D Computer Vision*. Prentice Hall, 1998.