

COMENIUS UNIVERSITY IN BRATISLAVA
FACULTY OF MATHEMATICS, PHYSICS AND INFORMATICS
DEPARTMENT OF APPLIED INFORMATICS



LEARNING MULTISENSORY INTEGRATION
AND COORDINATE TRANSFORMATION IN A
SIMULATED HUMANOID ROBOT

Master thesis

2019
BC. MARTIN KELLNER

UNIVERZITA KOMENSKÉHO V BRATISLAVE
FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

UČENIE MULTISENZORICKEJ INTEGRÁCIE A
TRANSFORMÁCIE SÚRADNÍC V SIMULÁTORE
HUMANOIDNÉHO ROBOTA

Diplomová práca

Študijný program: Aplikovaná informatika
Študijný odbor: Aplikovaná informatika
Školiace pracovisko: Katedra aplikovanej informatiky
Školiteľ: prof. Ing. Igor Farkaš, Dr.

Bratislava, 2019
Bc. Martin Kellner



Univerzita Komenského v Bratislave
Fakulta matematiky, fyziky a informatiky

ZADANIE ZÁVEREČNEJ PRÁCE

Meno a priezvisko študenta: Bc. Martin Kellner
Študijný program: aplikovaná informatika (Jednoodborové štúdium, magisterský II. st., denná forma)
Študijný odbor: aplikovaná informatika
Typ záverečnej práce: diplomová
Jazyk záverečnej práce: anglický
Sekundárny jazyk: slovenský

Názov: Learning Multisensory Integration and Coordinate Transformation in a Simulated Humanoid Robot
Učenie multisenzorickej integrácie a transformácie súradníc v simulátore humanoidného robota

Anotácia: Spracovanie senzorických podnetov v mozgu zahŕňa multisenzorickú integráciu (kombinovanie vodítko do jedného spoločného podnetu) a transformácie súradníc (referenčných rámcov, napr. z retinotopického na telo-centrický, ovplyvnenej modulačnou premennou, napr. pozíciou očí). Táto schopnosť je kľúčová v kognitívnej robotike, ak chceme vybaviť robota schopnosťou operovať autonómne v 3D priestore.

Anotácia: 1. Naštudujte si problematiku z kognitívnej neurovedy o multisenzorickej integrácii a referenčných rámcoch (súradnicových systémoch).
2. Implementujte a natrénujte model umelej neurónovej siete, ktorá sa naučí integrovať vizuálnu a proprioceptívnu informáciu, vykonávajúc prepočet súradníc v úlohe týkajúcej sa ruky a očí, s využitím simulovaného robota iCub.
3. Vyhodnoťte a analyzujte správanie natrénovaného modelu.

Literatúra: Makin J., Fellows M., & Sabes P. (2013). Learning multisensory integration and coordinate transformation via density estimation. PLOS: Comput. Biol., 9(4).
Švec M., Farkaš I. (2014). Calculation of object position in various reference frames with a robotic simulator. In Proceedings of the 36th Annual Conference of the Cognitive Science Society, Quebec, Canada.
Tikhonoff V., Fitzpatrick P., Nori F., Natale L., Metta G., & Cangelosi A. (2008). The iCub humanoid robot simulator. Advanced Robotics, 1(1), 22-26.

Vedúci: prof. Ing. Igor Farkaš, Dr.
Katedra: FMFI.KAI - Katedra aplikovanej informatiky
Vedúci katedry: prof. Ing. Igor Farkaš, Dr.
Dátum zadania: 16.10.2017

Dátum schválenia: 20.10.2017
prof. RNDr. Roman Ďurikovič, PhD.
garant študijného programu

.....
študent

.....
vedúci práce



Univerzita Komenského v Bratislave
Fakulta matematiky, fyziky a informatiky

ZADANIE ZÁVEREČNEJ PRÁCE

Meno a priezvisko študenta: Bc. Martin Kellner
Študijný program: aplikovaná informatika (Jednoodborové štúdium, magisterský II. st., denná forma)
Študijný odbor: aplikovaná informatika
Typ záverečnej práce: diplomová
Jazyk záverečnej práce: anglický
Sekundárny jazyk: slovenský

Názov: Learning Multisensory Integration and Coordinate Transformation in a Simulated Humanoid Robot
Učenie multisenzorickej integrácie a transformácie súradníc v simulátore humanoidného robota

Anotácia: Spracovanie senzorických podnetov v mozgu zahŕňa multisenzorickú integráciu (kombinovanie vodítok do jedného spoločného podnetu) a transformácie súradníc (referenčných rámcov, napr. z retinotopického na telo-centrický, ovplyvnenej modulačnou premennou, napr. pozíciou očí). Táto schopnosť je kľúčová v kognitívnej robotike, ak chceme vybaviť robota schopnosťou operovať autonómne v 3D priestore.

Anotácia: 1. Naštudujte si problematiku z kognitívnej neurovedy o multisenzorickej integrácii a referenčných rámcoch (súradnicových systémoch).
2. Implementujte a natrénujte model umelej neurónovej siete, ktorá sa naučí integrovať vizuálnu a proprioceptívnu informáciu, vykonávajúc prepočet súradníc v úlohe týkajúcej sa ruky a očí, s využitím simulovaného robota iCub.
3. Vyhodnoťte a analyzujte správanie natrénovaného modelu.

Literatúra: Makin J., Fellows M., & Sabes P. (2013). Learning multisensory integration and coordinate transformation via density estimation. PLOS: Comput. Biol., 9(4).
Švec M., Farkaš I. (2014). Calculation of object position in various reference frames with a robotic simulator. In Proceedings of the 36th Annual Conference of the Cognitive Science Society, Quebec, Canada.
Tikhonoff V., Fitzpatrick P., Nori F., Natale L., Metta G., & Cangelosi A. (2008). The iCub humanoid robot simulator. Advanced Robotics, 1(1), 22-26.

Vedúci: prof. Ing. Igor Farkaš, Dr.
Katedra: FMFI.KAI - Katedra aplikovanej informatiky
Vedúci katedry: prof. Ing. Igor Farkaš, Dr.
Dátum zadania: 16.10.2017

Dátum schválenia: 20.10.2017
prof. RNDr. Roman Ďurikovič, PhD.
garant študijného programu

.....
študent

.....
vedúci práce

DECLARATION: I hereby declare that this thesis is my own work and that all the sources I have used or quoted have been indicated and acknowledged as complete references.

.....

Abstrakt

Slovenský abstrakt v rozsahu 100-500 slov, jeden odstavec. Abstrakt stručne sumarizuje výsledky práce. Mal by byť pochopiteľný pre bežného informatika. Nemal by teda využívať skratky, termíny alebo označenie zavedené v práci, okrem tých, ktoré sú všeobecne známe.

Kľúčové slová: jedno, druhé, tretie (prípadne štvrté, piate)

Abstract

Abstract in the English language (translation of the abstract in the Slovak language).

Keywords:

Contents

1	Introduction	1
2	Related literature and theoretical focus	2
2.1	Reference of Frame	2
2.2	Sensorimotor transformations	4
2.2.1	Spatial Transformations for Eye–Hand Coordination	4
2.3	Gain fields	9
2.3.1	Links between Gain Modulation and Coordinate Transfor- mations	10
2.4	Releated researches	13
2.4.1	Learning using Back-Propagation	13
2.4.2	More Biologically Plausible Learning Rules Than Back- Propagation	15
3	Methods	18
3.1	ICub - Humanoid platform	18
3.2	Artifical neural networks	18
3.2.1	UBAL	18
4	Experiment	19
4.1	Data collecting	19
4.2	Training and Testing	19
4.3	Results	19

List of Figures

2.1	Schematic illustration of the spatial representations of objects in frames of reference (A), vectorially (B) and by coordinate systems (C). On the left, the frame of reference moves with the passenger; on the right, the observer's frame of reference is fixed to the earth.	3
2.2	Gaze-centered encoding of reach space. (a) Drawing depicting the egocentric visual directions of the hand, an orange, and an apple, as shown by the gray arrows. (b) Side view of a human brain showing areas (highlighted in yellow and orange) that encode reach space in gaze-centered coordinates: V1, striate cortex; PPC, posterior parietal cortex. (c) The consequence of an eye movement on the gaze-centered representation of the visual field. The head/eye diagrams depict current gaze position; the circles represent the gaze-centered representation of this visual scene (dotted lines represent visual horizontal and vertical axes and intersect at the fovea). If the person in the upper diagram looks at the orange, the hand and apple are represented in the left visual field (upper circle). In contrast, if the person fixates the apple (lower diagram), the orange and hand are now represented in the right visual field (lower circle). If the orange and the hand were no longer visible when the eye movement occurred, the brain would need to remap their position by taking the intervening eye movement into account (Crawford J., et al. (2004)).	6

- 2.3 Conceptual scheme for spatial transformations in eye– hand coordination. To illustrate the model, consider the following “task”: a subject looks at a briefly flashed target (F) with the arm at resting position (A). Then (B) the subject makes 1) an upward eye movement, followed by 2) a reaching or pointing movement toward the remembered target location (E). We hypothesize that the brain uses the following stages to do this. C: an early representational stage. Target location is stored in eye coordinates, such that this representation (E) must be counterrotated (updated) when the eye rotates. D: comparison stage. Updated target representation (E) is compared with an eye-centered representation of current hand location to generate “hand motor error” in eye coordinates (Buneo et al. 2002). E: visuomotor execution stage. “Hand motor error” signal is rotated by eye orientation and head orientation (or perhaps by gaze orientation) to put it into a body coordinate system appropriate for calculating the detailed inverse kinematics and dynamics of the movement. This last stage would also have to include internal models of the geometry. 8
- 2.4 Working principle of gain fields, based on Zipser and Andersen (1988). The upper part of the panel shows the hypothetical receptive fields of two neurons that are gain modulated (e.g., by eye or hand position) in opposite ways without shifting. For example, the three lines in each graph could represent visual receptive fields mapped relative to gaze at a leftward eye position (red solid line), a central eye position (green dashed line), and a rightward eye position (blue dotted line). Here, eye position modulates the strength of response of two neurons, but does not cause them to shift. However, summation of these two gain-modulated neural responses results in shifting receptive fields in the output, e.g., eye position (or in other cases hand position) has shifted the receptive field. (Blohm G. and Crawford J.D (2009)) 10
- 2.5 A coordinate transformation performed by the visual system. While reading a newspaper, you want to reach for the mug without shifting your gaze. The location of the mug relative to the body is given by the angle between the two dashed lines. For simplicity, assume that initially the hand is close to the body, at the origin of the coordinate system. The reaching movement should be generated in the direction of the mug regardless of where one is looking, that is, regardless of the gaze angle x_{gaze} . The location of the target in retinal coordinates (i.e., relative to the fixation point) is x_{target} , but this varies with gaze. However, the location relative to the body is given by $x_{target} + x_{gaze}$, which does not vary with gaze. Through this addition, a change from retinal, or eye-centered, to body-centered coordinates is performed. (Salinas E. and Sejnowski T. (2001)) 11

- 2.6 Visual responses that are gain-modulated by gaze angle. The response of a parietal neuron as a function of stimulus location was measured in two conditions, with the head turned to the right or to the left, as indicated in the upper diagrams. In these diagrams, the cross corresponds to the location where gaze was directed, called the fixation point; the eight dots indicate the locations where a visual stimulus was presented, one location at a time; and the colored circles show the position of the recorded neuron's receptive field. This was centered down and to the left of the fixation point. In the diagrams, the rightmost stimulus corresponds to 0 degrees, the topmost one to 90 degrees, and so forth. The dashed line indicates the direction straight ahead. The graph below plots the neural responses in the two conditions, indicated by the corresponding colors. The continuous lines are Gaussian fits to the data points. When the head is turned, the response function changes its amplitude, or gain, but not its preferred location or its shape. (Data redrawn from Brothie and others 1995.) 12

List of Tables

Introduction

Related literature and theoretical focus

2.1 Reference of Frame

To better understand the term of *reference frame* or *frame of reference*, consider the following situation: In a moving train is one passenger which lets a book fall. Another person is watching the case and standing on the ground. From the perspective of the passenger, the book is falling straight down, but the observer (the man standing on the ground) sees the book dropping down along a curved path because of the movement of the train. The situation mentioned above has shown that the description of a physical phenomenon depends on the position from where is the case observed. Using reference frames on the previous illustrating example, we can say the book is falling straight down in the train's reference frame but is dropping along a curved path in the earth's frame of reference. About the state of motion of the passenger, we might say he is stationary in the train's reference frame and moving in the earth's frame of reference (Soechting and Flanders, 1992).

Reference frames have often been used by physicists and engineers, but neuroscientists have also adopted the term. In the mathematical sense, the frame of reference is similar to a coordinate system that is defined by a set of axes and an origin. The origin can be anywhere in space, and the orientation of axes can be chosen arbitrarily because a reference frame is only characterized by the state of motion relative to an object. Therefore, the train's reference frame that was mentioned in the illustrated example has the same state of motion as the train, and any point in the reference frame can be defined by its position along each of coordinate axes. Alternatively, the location of a point can be defined by way of vectors. A vector has two properties: a magnitude and a direction; the magnitude is the length of a line segment between the origin and the point, and its direction is from the origin to the point (Soechting and Flanders (1992), Kumar and Barve (2002)).

To switch the reference frame into a term in neuroscience we have to make some changes against the mathematical definition because of neural systems do not report the position of an object as a vector or coordinates, but neurons of the visual system encode information on a restricted area of space. The cells will respond to its *receptive field* or *response field*. That means the firing of the neurons will change depending on stimuli coming from a special area of visual space (Batista 2002). If we want to identify in which frame of reference, for

example, neurons encoding spatial information about the location of an object, then the firing activity of the neurons should stay the same as long as the image of the object falls on the same locus on the retina and stay constant. Once that happened we can say the location is encoded in a *retinocentric* frame of reference (Soechting and Flanders, 1992).

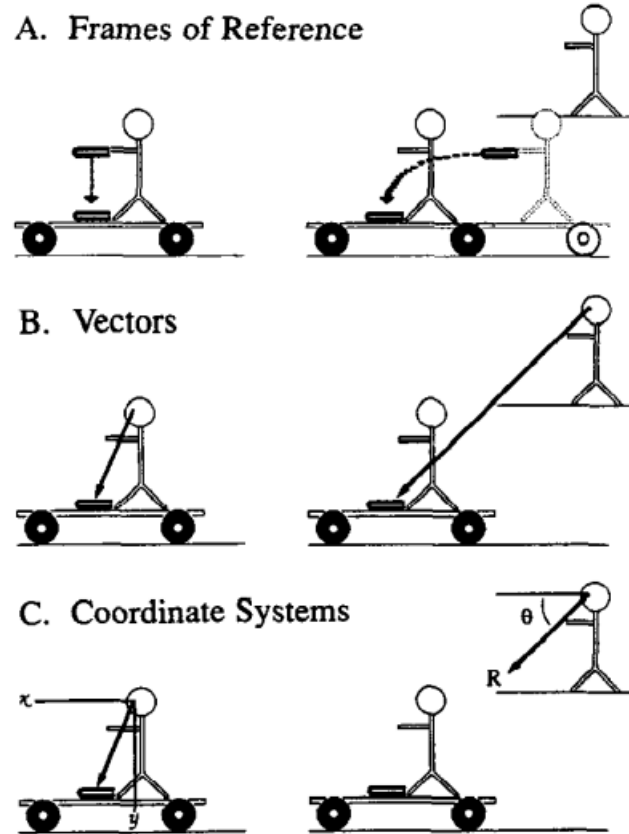


Figure 2.1: Schematic illustration of the spatial representations of objects in frames of reference (A), vectorially (B) and by coordinate systems (C). On the left, the frame of reference moves with the passenger; on the right, the observer's frame of reference is fixed to the earth.

Different brain regions encode spatial locations in different RFs. In neuroscience are RF mostly divided to two main groups: *ego-centric* and *allocentric* reference frames. In an egocentric RF, the location of an object referenced to an observer. At a neural level, an egocentric RF means the receptive field, or the response of the neural population is attached to a reference point (Committeri, G., et al (2004)). For example, an eye-centred RF moves with the eyes (Colby, 1998). Spatial locations are mainly encoded in egocentric reference frames, especially space coding neurons in parieto-frontal cortex are associated with reference frames centred to the eye, the head or the hand (Colby, 1998; Hagler et al., 2007; Sereno and Huang, 2014). Also in the monkey's posterior parietal cortex and in connected regions of the premotor cortex has found neurons that encode space in RFs centred to parts of the body (Cohen & Andersen, 2002; Colby, 1998). Coding of space in allocentric RFs is not so well-studied like egocentric RFs, but it has shown (Fink, Marshall, Shah, et al., 2000; Galati, Lobel, et al., 2000; Honda, 1998; Fink 1997) that in the posterior parietal cortex and the dorsal premotor

cortex (PMd), and early visual areas, neurons report during object-based spatial judgement. Also, it has been proven that during solving of a complex task, such as landmark knowledge, orientation in large-scale space and navigation, some operation referred to both types of reference frame, allocentric and egocentric, which are difficult to untangle (Committeri G., et al (2004)).

2.2 Sensorimotor transformations

When humans, animals, or insect want to carry out basic operations to interact with the world, then they use their bodies and senses. To manipulate with an object, moving in space, or doing other daily completion, the brain must continuously update the internal representation of the world, deal with restrictions that can occur while performing an operation and change body configuration to achieve the goal of the action. For instance, during the reaching of an object, the human brain must take into account the location of the target, the current position of the particular hand to generate the right movement of the hand. One non-neuroscientist might say the operations as above mentioned are not very complicated because we do not need to spend a lot of effort to get the goal, but when we look closely at the neural level, the task is more complicated like it would seem. Information about the location of an object can be produced by stimuli from different modalities, each of these modalities provides neural information that is encoded in a different frame of reference. (Moving in ...) Therefore, neuroscientists focus their research to clarify the following questions: How is the information represented in each reference frame? How is the information from different modalities combine to encode the correct representation of space? How is information in particular reference frame transform to another one? *Coordinate transformation*, in neuroscience, is a term to call transformation from one reference frame to another. We closely discuss eye-hand coordinate transformation, that is central to many human daily activities, in the next section.

2.2.1 Spatial Transformations for Eye–Hand Coordination

Eye-hand coordination is reviewed by Crawford J., et al. (2004) in detail. The authors also bring assumptions about using gaze-centred representations during the process of coordinate transformation that needed to perform a task such as reaching or grabbing an object. In this chapter, we closely discuss their conclusions.

To reach an object the brain needs information about the location of the reaching object, and the position, for the purpose, is encoded relative to a body part. When we realise the fact that arms are anchored to shoulders, and spatial information on object comes mainly from projections on retinas, then we can realise the brain must perform coordinate transformation between at least two representation. First of them is the representation specifying the object's location relative to retinas, and the second representation within which the activation of arm's muscles must be determined to reach the target. Therefore we can say the purpose of eye-hand coordinate transformation (or system that is responsible to do that) is to guide movements of hands using vision. This section discusses how

the brain deals with the transformation of visual stimuli into hand motion for controlling and grabbing a target.

A robot system for controlling the robot's arm to reach or to grab an object is usually in comparison with the human control system is simpler. Such a robot system is usually reduced by using visual feedback, basically, to drive the hand to a point seen by the visual system comparing the current position of the hand with the target position. This can work in such a robotic system because the speed of visual feedback is not so limited as a brain control system. The robotic system visual feedback is limited by computer processing time and speed of electrical flow, whilst the speed of neural conduction and processing time in a real primate brain is not adequate to perform fast hand movement that would be way off the field of vision before new visual feedback would arrive to accurately update the movement. Consequently, a primate brain cannot use only visual feedback to drive hand movement and must use another way to perform correct and fast movements: to use internal models of the physical system and world based on initial conditions and control eye-hand coordinate transformation in feedforward way. Despite this, we cannot say that visual feedback does not contribute to guiding hand movements. Visual feedback helps to achieve the best performance of grabbing a goal, and are essential to deal with unexpected events or a danger that can occur during performing a particular movement. To better understanding of such a feedforward transformation is helpfull to divide it into concentual steps.

Visual Representations of Reach Space

Internal models, for such operations as coordinate transformations, are described by using the concept of reference frames (2.1). Final reference frames for eye-arm coordination associated with points in upper arm and shoulder, and not hand-centrered representation like we could wrongly consider of. The first stimuli are controlled in a retinal frame, and there is a problem with gaze shifting. Every time when eyes (or/and head) shift(s) gaze the visual relationship between sensory apparatus and the external world (Hallett and Lightstone 1976). One solution, how the brain could deal with that problem, is waiting until the gaze shifting gets finished and then update its visual information (O'Regan and Noe 2001), but that could cause that the original target of interest could be moved into less-sensitive peripheral retina or even out of the visual field. It Is good to realise that would also produce long lags in processing time and also would introduce redundant visual computations, and thus, to perform visual guidance of arm performing a sequence of saccades (note) would be impossible. Therefore, there must the necessary representations for future action be stored, and either in an eye-independent form or in such form that is internally updated when the gaze is shifting whether caused by movements of eyes or head. (Duhamel et al. 1992).

Henriques et al. (1998) were asking themselves the question: How the system in the human brain, which appears to be used for early planning of pointing and reaching a target (McIntyre et al. 1997; Vetter et al. 1999), stores these early motor representations during eye movements? In order to answer the question, they performed the experiment described by FIGURE ?? . Human subjects, which were performing the experiment, was asked to point toward the location using hand in total darkness. At the start of the experiment the target location,

that was pointed, was centred to subject's fovea, then subjects shift the gaze by moving of eyes, and in the end, perform mentioned pointing toward the remembered location. Henriques et al. (1998) tried to determine whether the responses

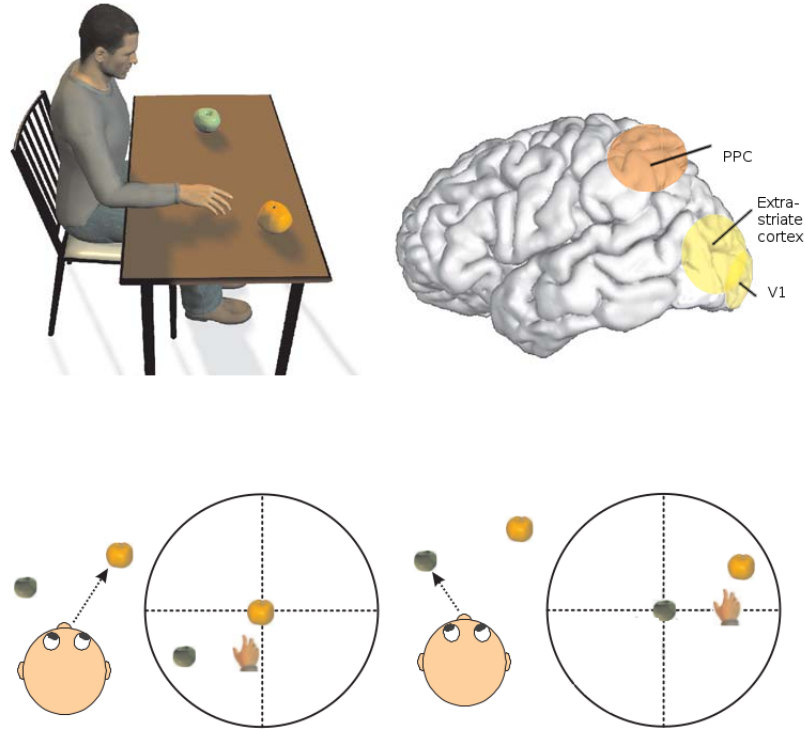


Figure 2.2: Gaze-centered encoding of reach space. (a) Drawing depicting the egocentric visual directions of the hand, an orange, and an apple, as shown by the gray arrows. (b) Side view of a human brain showing areas (highlighted in yellow and orange) that encode reach space in gaze-centered coordinates: V1, striate cortex; PPC, posterior parietal cortex. (c) The consequence of an eye movement on the gaze-centered representation of the visual field. The head/eye diagrams depict current gaze position; the circles represent the gaze-centered representation of this visual scene (dotted lines represent visual horizontal and vertical axes and intersect at the fovea). If the person in the upper diagram looks at the orange, the hand and apple are represented in the left visual field (upper circle). In contrast, if the person fixates the apple (lower diagram), the orange and hand are now represented in the right visual field (lower circle). If the orange and the hand were no longer visible when the eye movement occurred, the brain would need to remap their position by taking the intervening eye movement into account (Crawford J., et al. (2004)).

of such a pointing task was affected by an intervening eye displacement by comparing them to pointing to remembered foveal targets or to retinally peripheral targets. When subjects were pointing toward using an updated gaze-centered representation, pointing behaviour would echo pointing to peripheral targets, whereas pointing using a nonretinal representation would cause no effect. This experiment has cleared up the question about early motor representation and has supported the idea of using an updated gaze-centered frame of reference. Other studies have brought findings that during pointing of auditory and propriocep-

tive targets such a gaze-centered updating has been also recorded (Pouget et al. 2002). Evidence of coding the spatial information of an object in a gaze-centred frame of reference can be also found in both early visual pathway (i.e., retina, lateral geniculate nucleus, striate cortex) and also later in the visual pathway (e.g., extrastriate and parietal cortex).

The gaze-centered representation moves depending on the gaze because the retinal projection of the project depends on the gaze. Therefore, the positions of a reaching target can be defined by the target's direction (horizontal and vertical angular eccentricity) relative to *fovae*, and by the distance between the goal and eyes. Early visual areas carry out the computation of target direction as seen by a virtual eye which we can imagine placed between right and left eye, whereas the required distance is determined by taking into account monocular information (accommodation, relative object size, shading, perspective, etc.), and binocular information (retinal disparity and convergence). Later during early movement planning, signals that contain the information seem to be merged into a single gaze-centred representation of the space in the posterior parietal cortex.

How we has mentioned above gaze-centerered representation must be updated to maintain the stable representation of the world. The brain must update this representation all time while eyes are moving, because the gaze changed. This process which is performing during eyes rotation is named *updating* or *remapping*, and it also occurs during movements of the head and the body. An illustrating example is in FIGURE 1(c), the observer originally looks at the orange, and the apple is placed to left-down relative to the orange. The orange is at the centre of the gaze-centred representation and the apple located on the lower-left quadrant of the visual field. A fast orientation movement approaching the apple cause that the objects are remapped by the same rotation as the movement, but in the opposite direction, what means, the apple is mapped at the centrum, and the orange is located on the upper-right quadrant. It has shown the updating can be achieved even if vision is removed, only by using information about eyes movements.

Developing the reach plan

To correctly generate a reach plan the information about initial hand position must be linked to. (Buneo et al. 2002) have brought suggestions that comparison between initial hand location and gaze-centered representation of the visual target is done earlier and in a gaze-centered reference frame, even when the hand is out of visual field, what means, proprioceptive signals that carry the information about hand location must be transformed into gaze-centered coordinates. (Buneo et al. 2002) found such a kind of signals in the process of a gaze-centered transformation in parietal area 5. Findings of the mentioned earlier comparison do not want to say that a next reference frame transformation is not required to reach a target.

The next important problem, which the brain must solve to accuracy generate and perform a reach plan is how ego-center in translocated during head rotation, because rotation of head cause eye's translation with respect to the shoulder. Ignorance of eye translation component would cause erroneous reach pattern at noncentral head position.

A summary of conceptual and physiological models of visually guided reach is described in the figure 2.3. The first stage of figure corrections to the fact that 3D

representations of target direction are stored and maintained in a retinal frame. The second stage shows an illustration of the transformation according to above-mentioned Buneo et al. (2002) schema to compute the hand displacement in a retinal reference coordinates. There are other reference frames transformations from gaze coordinates into shoulder representations that are not included in the figure and are necessary to reflect the eye-head-shoulder system, but such models like these illustrated can be useful in telling us what signals are required and used for a transformation, but obviously, all of such kinds of explicit intermediate representations are not used by brain directly, and also do not tell us how and where mentioned signals are coded.

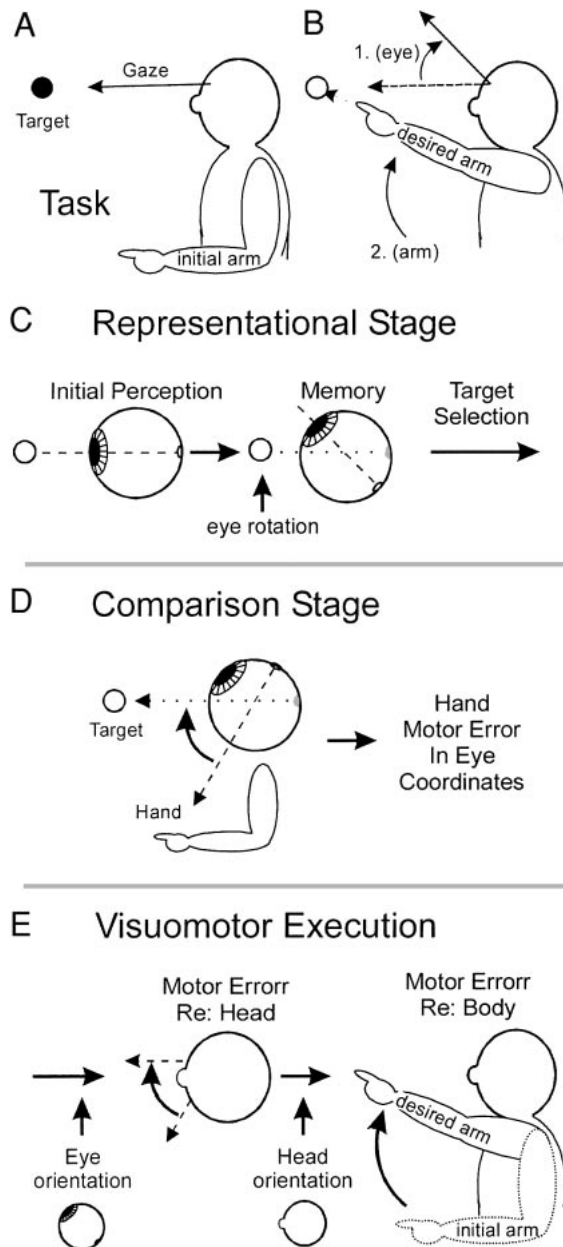


Figure 2.3: Conceptual scheme for spatial transformations in eye-hand coordination. To illustrate the model, consider the following “task”: a subject looks at a briefly flashed target (F) with the arm at resting position (A). Then (B) the subject makes 1) an upward eye movement, followed by 2) a reaching or pointing movement toward the remembered target location (E). We hypothesize that the brain uses the following stages to do this. C: an early representational stage. Target location is stored in eye coordinates, such that this representation (E) must be counterrotated (updated) when the eye rotates. D: comparison stage. Updated target representation (E) is compared with an eye-centered representation of current hand location to generate “hand motor error” in eye coordinates (Buneo et al. 2002). E: visuomotor execution stage. “Hand motor error” signal is rotated by eye orientation and head orientation (or perhaps by gaze orientation) to put it into a body coordinate system appropriate for calculating the detailed inverse kinematics and dynamics of the movement. This last stage would also have to include internal models of the geometry.

To design a complete 3D reach plan, the brain must know the amplitude of the desired movement and the direction. It has shown that in the dorsal premotor cortex the direction and the amplitude are encoded together inside the same neuron, but there is also evidence that direction and amplitude might also be

encoded independently. One explanation of this contradiction might be that, at the neural level, amplitude and direction are encoded together, but are used by different mechanisms to generate muscle activation. In account must be taken the rotation of head and body and also values of correct forces for specific muscles. An amplitude and a direction are properties of vectors, therefore, we can reformulate the problem of designing a reach plan into terminology of vectors as it is often done.

The conclusion of the problem of developing reach plan brain must design a hand movement vector, which depends on the difference between the hand location and the location of the target. The areas in the brain that are responsible for calculating the vector must have spatial information of the hand and the goal position, and the spatial information must be represented in the same reference frame. To determine the movement vector the particular areas could use information from different sources and also perform multisensory integration to get the most likely estimation of both locations, the hand and the target. One of approaches that allow us to examine how populations perform the task of coordinate transformation is using artificial neural networks. During studying of neural encoding of a cognitive function, we very often focus on changing response magnitude of neurons as known as *gain field* what we are discussing in the next section 2.3. (Crawford J., et al. (2004),).

2.3 Gain fields

The term of *gain field* comes from Andersen and Mountcastle (1983). The authors tested the visual receptive field for a specific neuron at different eye positions. The receptive field of a neuron or a neural population can be described as a specific region in sensory space which stimuli can have an influence on the activity of a single neuron or a neural population. For instance, activations of neurons encoding the gaze-centered representation of a target's location should be the same while gaze angles stay constant, contrariwise, firing activity of the neurons change with every eye motion. During the mentioned study the neuron's action potential frequency was changing in a manner of multiplication the frequency by gaze angle, that is scaled by some constant. Both of the shape and the locations of the visual receptive field was unchanged, only the receptive field was scaled by some *gain* factor. This phenomenon was the first time characterized in neuron within LIP and visual area 7A. Another study (Zipser and Andersen (1988)) which has dealt with the task of coordinate transformation from visual target position and signals of gaze position into the position representing by the target location in a space-fixed reference frame, trained a neural network to perform the transformation. Detailed analysis of their neural model has shown the network was able to develop visual receptive fields modulating by eye position in a very close way compared to modulation in parietal cortex. The figure 2.4 illustrates and can help to understand better the result of the study. At the neural population level gain fields modulate responses of single neurons, therefore, the brain can increase or decrease the population output of a neural population, that can allow comparing something like the output of the above mentioned artificial neural network with information about current hand location, and that might

use the process of computation the hand movement to reach a target. It has been shown that also other types of signals produce gain modulation and the modulations, and gain fields were found in many other brain areas. (Blohm G. and Crawford J.D (2009))

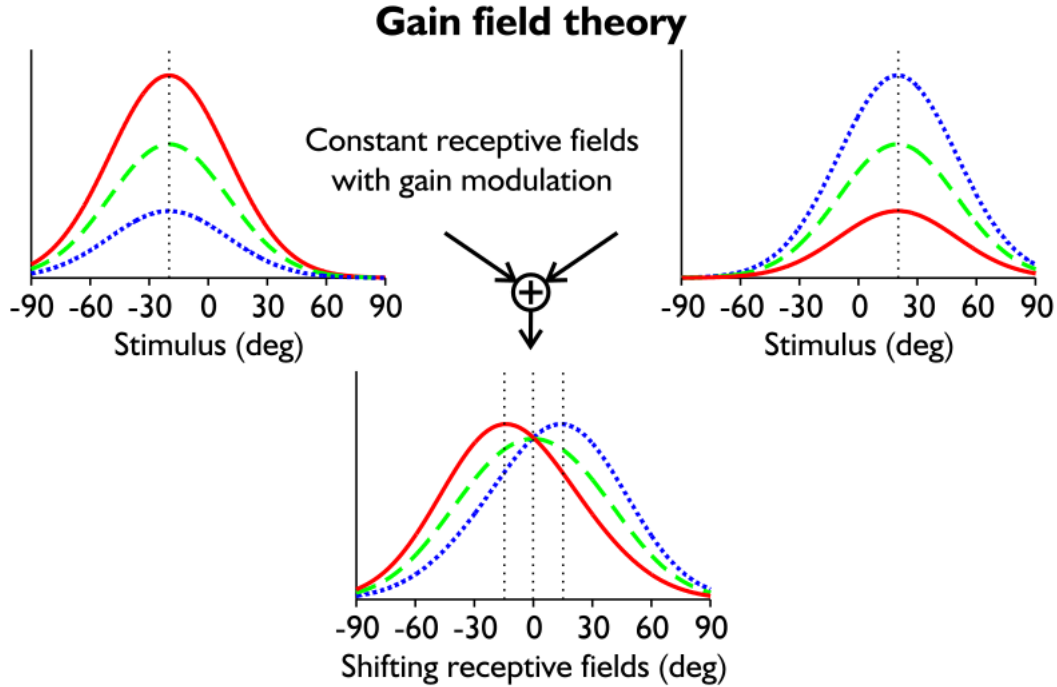


Figure 2.4: Working principle of gain fields, based on Zipser and Andersen (1988). The upper part of the panel shows the hypothetical receptive fields of two neurons that are gain modulated (e.g., by eye or hand position) in opposite ways without shifting. For example, the three lines in each graph could represent visual receptive fields mapped relative to gaze at a leftward eye position (red solid line), a central eye position (green dashed line), and a rightward eye position (blue dotted line). Here, eye position modulates the strength of response of two neurons, but does not cause them to shift. However, summation of these two gain-modulated neural responses results in shifting receptive fields in the output, e.g., eye position (or in other cases hand position) has shifted the receptive field. (Blohm G. and Crawford J.D (2009))

2.3.1 Links between Gain Modulation and Coordinate Transformations

In order to demonstrate the role of gain modulation with performing the task of coordinate transformation in the brain, we will discuss again results that Zipser and Andersen (1988) has brought. They trained a neural network to solve coordinate transformation from gaze signal and retinal location to stimuli in body-centered reference frame. During the artificial neural network, which was training to learn the transformation had as the input information on mentioned sensory stimuli and the output was the body-centered representation of locations. The model was trained using the one of major-used training in machine learning called *backpropagation*. This algorithm allows to *learn* values of *hidden* units (neurons) in order to map inputs to corresponding outputs with the effort to achieve as low

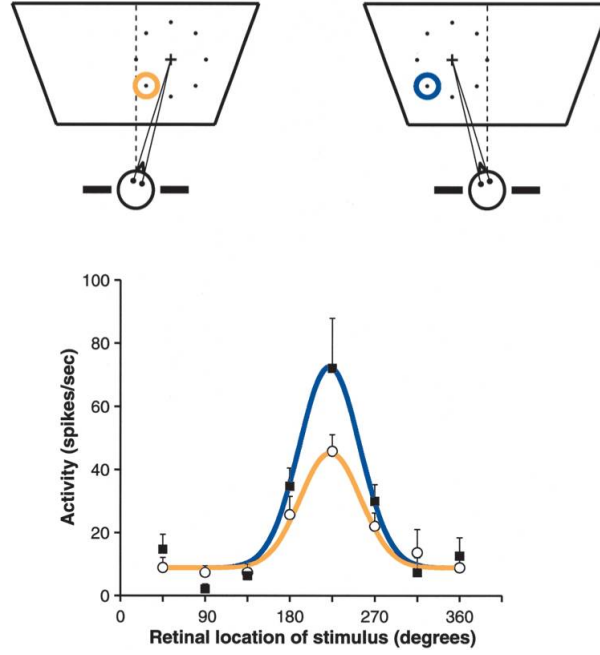


Figure 2.6: Visual responses that are gain-modulated by gaze angle. The response of a parietal neuron as a function of stimulus location was measured in two conditions, with the head turned to the right or to the left, as indicated in the upper diagrams. In these diagrams, the cross corresponds to the location where gaze was directed, called the fixation point; the eight dots indicate the locations where a visual stimulus was presented, one location at a time; and the colored circles show the position of the recorded neuron's receptive field. This was centered down and to the left of the fixation point. In the diagrams, the rightmost stimulus corresponds to 0 degrees, the topmost one to 90 degrees, and so forth. The dashed line indicates the direction straight ahead. The graph below plots the neural responses in the two conditions, indicated by the corresponding colors. The continuous lines are Gaussian fits to the data points. When the head is turned, the response function changes its amplitude, or gain, but not its preferred location or its shape. (Data redrawn from Brotchie and others 1995.)

following equation:

$$r = f(x_{target} - a) g(x_{gaze}), \quad (2.1)$$

where $g(x_{gaze})$ is the function which result is the gain field of particular neuron, and $f(x_{target} - a)$ represent the response of neurons as a curve with a single pick (Figure 2.6). Then, downstream neurons driven by a population such as the above mentioned must have responses R theoretically formed like the following function of x_{target} and x_{gaze} :

$$R = F(c_1 x_{target} + c_2 x_{gaze}), \quad (2.2)$$

where both constants c_1 and c_2 depend on the synaptic weight. The receptive field of the downstream neuron is represented like the peaked function R . This mathematical formulation is only a simplification of the problem, but the authors found conditions under those this can happen and experimentally confirmed that the downstream neurons encode the sum of $x_{target} + x_{gaze}$, what means, the downstream neurons those must encode target in a body-centered reference frame respond as a function of x_{target} swifts when x_{gaze} changes (Salinas and Ab-

bott (1995)). Similar swift was demonstrated in several cortical areas associated with a variety of coordinate transformations.

2.4 Releated researches

Artificial neural networks are mathematical model inspired by real brain neural circuits and used in studying cognitive functions of humans, whether in order to get to understand how are performed particular brain's tasks, or to implement human behaviours into a robotic humanoid. We will not discuss artificial models in detail in this section. This section is focused on significant research works that have been dealt with coordinate transformation, we will discuss their results and methods.

2.4.1 Learning using Back-Propagation

Zipser and Andersen (1988)¹ have been first who trained an artificial neural model and has brought findings that artificial neural networks are able to decode spatial transformation.

Authors have been concerned with the question: How the brain perform such as coordinate transformations which translate sensory inputs into motor commands. They mainly have focused on the brain's area named 7a and neural processes that are performed by this area. Earliel research findings have proposed that the area 7a is most likely used to perform spatial transformations. This area contains cells that to receive a convergence of both retinal signals and eye-position as a non-linear interaction, and visual responses have seemed as modulated by a function of ey position multiplied by the response of retinal receptive fields, what basically means, that visual receptive field remains retinotopic, but the intensity of response is modulated by position of eyes.

To facilitate a comparison of artificial neural cells with cells of area 7a, they analysed experimental data collected in studies with awake, unanaesthetized monkeys. The experimental data was collected by recording neuronal activity during performing various visuospatial tasks. They have focused on three major kinds of neurons in 7a area: the cells responding to eye-position only, the visual cells responding to retinal stimuli, and the cells that respond to both of mentioned cell classes those interactions produce an eye-dependent representation of visual targets. In order to get to test eye-position sensitivity the following experiment was performed: The animal fixated on a small point at different eye positions, their heads were fixated, and the experiment took place in total darkness. A linear increase was recorded in activity for a range of horizontal or vertical eye positions, but some cells showed more complex coding of eye position. To test the receptive fields of the visual cells a flashing spot stimulus was used at different locations in the visual fields while the animal fixated on a target at a single eye position. The neurons, that were specially tuned, were the largest group and showed a convergence of eye position and retinal position, what means, the visual response of the spatially tuned neurons alters as a function of eye position, what was observed by collecting data under the condition in which the visual stimulus

¹All infomation in the section 2.4.1 was received from Zipser and Andersen (1988) article if

appears at the peak locations in the retinal receptive field with animal fixating at different eye positions (možno doplnit).

Model and algorithm

They used a three-layer neural network. The network was trained in order to get the mapping of visual targets to head-centred coordinates. The input data had two parts, an array of units that represent the visual stimulus, and representation of eye position via a set of units. Both of the retinal and eye position inputs was formed using characteristics of the cells in PPC that response to visual stimuli, respectively, to eye stimuli. As the output layer of the network was used two representations: the first one was a gaussian format in which units had gaussian receptive fields coding location in head-centered reference frame. The other one was a monolonic format in which the activity of each neuron is a linear function of the location of the stimulus in head-centered representation. These two formats were chosen because they represent the most common types of coding formats discovered in brain cells. The hidden layer was trained using algorithm named *backpropagation*.

The backpropagation algorithm (Rumelhart et al., 1986a) is used to train hidden layers of a feedforward neural network. The information that is provided by the input x is propagated up to the hidden units at each layer to produce output \hat{y} , this part is called forward propagation. Each artificial neuron consists of two properties, a weight and an activation function, therefore the output \hat{y} depends on these properties of each neuron in hidden layers. The backpropagation is an algorithm that is used to find minimum of a cost function of mapping x to y through updating the weights. The second part consists from backward flow through the network to compute gradient, this is named backpropagation. The gradient which is calculated is the gradient of the cost function that expresses the error between the expected output value and the output value that is the result of the current state of hidden layers during the learning procedure. The gradient is computed by *chain rule* what allows to compute derivatives of function (cost function) formed by composing other functions whose derivatives are known. Suppose that $x \in \mathbf{R}^m, y \in \mathbf{R}^n$, function g maps from \mathbf{R}^m to \mathbf{R}^n , and f maps from \mathbf{R}^n to \mathbf{R} , then if $y = g(x)$ and $z = f(y)$, then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}. \quad (2.3)$$

The same equation only in vector notation:

$$\nabla_x z = \left(\frac{\partial y}{\partial x} \right)^T \nabla_y z, \quad (2.4)$$

where $\frac{\partial y}{\partial x}$ is the $n \times m$ Jacobian matrix of g .

The error E that is the result of the cost function is calculated through the composition of activation functions of neurons, what is a continuous and differentiable function of weights w_1, w_2, \dots, w_l in the network, perforce, the minimalization of E using an iterative process of gradient descent, consist of calculating the gradient:

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_l} \right), \quad (2.5)$$

and following the *chain* rule mentioned above each weight is updated as following:

$$\Delta w_i = -\gamma \frac{\partial E}{\partial w_i} \text{ for } i = 1, \dots, l, \quad (2.6)$$

where γ represent a learning rate variable, which defines the speed of learning (Goodfellow, I. et al (2017), Rojas, R (1996)).

Results

The model was trained by 1,000 trials. After successful training, the authors compared the experimental and trained receptive fields. This comparison has shown that the trained model generate single-peak receptive fields and also produce moderately complex these all are similar to those experimentally observed. The characteristics of generated respective fields of the trained network were obtained by holding the eye-position input into the network constant and simulating visual stimulation at the same retinal position as used in the above-mentioned experiment on area 7a.

By this comparison, Zipper and Andersen have shown that their model trained through backpropagation algorithms was able to learn receptive fields similar to those experimentally observed during performing the task of the coordinate transformation.

2.4.2 More Biologically Plausible Learning Rules Than Back-Propagation

Navarro et al (2018) have studied a coordinate transformation from eye-centered to head-centered coordinates (reference frames) using self-organising artificial neural networks. They have provided a biologically plausible training algorithm to demonstrate how head-centered response in the brain could be developed through an unsupervised process of visually-guided learning.

Coordinate transformation as one of the most studied domains in neuroscience has been examined by using several approaches. A lot of the approaches/works have been inspired by the early work of Zipper and Andersen (1988) (ref), and the models of those works was trained the synaptic weights matrix using a form of supervised. The problem of the usage of supervised learning is that such a training algorithm does not reflect the real learning process in the brain, because supervised learning algorithms oppose the Dale's law that says a given presynaptic neuron cannot be both excitatory (causing an increase of the likelihood of postsynaptic action potential) and inhibitory (causing a decrease of the likelihood of the postsynaptic action potential) across its efferent synapses.

In order to provide a more biologically plausible way, their laboratory has published a neural network of the visually-guided development of head-centered visual neurons those depend on associative learning rules and does not require a supervised approach. Their approach has consisted of an unsupervised competitive form of *trace learning*. Such an approach is similar to that in higher networks layers to bind together inputs that tend to occur close together in time, for instance, if a primate performs eyes movement more often compared to setting the position of the head, then retinal images that occur together in time will correspond to different positions of the eyes but to the identical position of the head.

Trace learning, in this case, will learn neurons in higher layers to respond to the position of visual stimuli in an identical location in the head-centred reference frame across different retinal position.

Model and Algorithm

By authors proposed and used model is a composition of four main components to learn visual representation in the head-centered reference frame. The first part encodes both the retinal position of the visual target and the position of the eyes through a population of input neurons. The second part consists of a population of the output neurons those compete among each other through mutual inhibitory interactions. The third, the local synaptic *trace learning* rule, which is used to update feedforward synaptic connections between neurons of the input population and the output population. The last component of the model relates to experimental findings in which primates prefer to set their gaze by moving their eyes rather by moving the head, what means, there are periods when the eyes are moving and the head is stationary with respect to the visual environment, and thus, the visual target also remains stable in head-centered reference frame.

This model consisted of the four parts allow to train synaptic connections matrix through performing visually gained training, when the eyes move around a virtual environment that contains a stationary visual target, and the head is also stationary. In this case, the visual target falls on different retinal location, but the head-centered position of the target stays unchanged. Series of the positions of eyes and retinal locations of the visual target are represented by retinotopic neurons and their responses are gain modulated by eyes positions. The trace learning rule is able to bind together inputs corresponding to a visual target on the same position in the head-centered reference frame onto the same output neurons, although the locations of the visual target in eye-centered reference frame are different.

The network architecture of the model consists of two layers. The first layer is a population of input neurons that encode both the position of eyes and the retinal locations of visual targets. The second layer is a population of the output neurons that compete among each other to create a representation of patterns in the input layer.

During training, each period of an epoch consists of a visual target fixating in a given location in head-centered reference frame and the eyes perform a series of saccades to different eyes positions.

Trace learning rules use a temporal trace of recent neural activity to get postsynaptic neurons to bind together inputs that occur close together in time. The trace value $d_i(t)$ for the i^{th} neurons in the output layer is given by following equation:

$$\tau \frac{dq_i}{dt} = -q_i + v_i, \quad (2.7)$$

where v_i is the instant firing rate of the neuron, and τ_q is a time constant common for all output neurons. The strength of the synapse from the j^{th} input neuron to the i^{th} output neuron is expressed by:

$$\frac{dw_{ij}}{dt} = \alpha q_i v_j^I, \quad (2.8)$$

where α is the learning rate, v_j^I is the instant firing rate of the j_{th} input neuron, and q_i is the trace value of i^{th} output neuron. After each weight update, the length of the weight vector for each output neuron i , that is $w_i = (w_{i1}, \dots, w_{iN_i})$ where there are N_I input neurons, is renormalised by:

$$w_i := \frac{w_i}{\|w_i\|} \quad (2.9)$$

Competitive learning in self-organizing maps consists of the comparison between a given input pattern and each weight vector w_i of output neurons. The weight vector that is the most similar to the given pattern is declared as the winner, and the winner and its neighbours are moved close to the input pattern by an update weight rules as 2.8.

Methods

3.1 ICub - Humanoid platform

3.2 Artifical neural networks

3.2.1 UBAL

Experiment

4.1 Data collecting

4.2 Training and Testing

4.3 Results

Bibliography

- [1] ANDERSEN, Ra and MOUNTCASTLE, Vb. The influence of the angle of gaze upon the excitability of the light- sensitive neurons of the posterior parietal cortex. *The Journal of Neuroscience*. January 1983. Vol. 3, no. 3p. 532–548. DOI 10.1523/jneurosci.03-03-00532.1983.
- [2] BATISTA, Aaron. Inner space: Reference frames. *Current Biology*. 2002. Vol. 12, no. 11., doi:10.1016/s0960-9822(02)00878-3.
- [3] BLOHM, Gunnar and CRAWFORD, J. Douglas. Fields of Gain in the Brain. *Neuron*. 2009. Vol. 64, no. 5p. 598–600. DOI 10.1016/j.neuron.2009.11.022.
- [4] BLOHM, G., KHAN, A.z. and CRAWFORD, J.d. Spatial Transformations for Eye–Hand Coordination. *Encyclopedia of Neuroscience*. 2009. P. 203–211. DOI 10.1016/b978-008045046-9.01102-5.
- [5] COLBY, Carol L. Action-Oriented Spatial Reference Frames in Cortex. *Neuron*. 1998. Vol. 20, no. 1p. 15–24. DOI 10.1016/s0896-6273(00)80429-8.
- [6] COMMITTERI, G., et al. Reference Frames for Spatial Cognition: Different Brain Areas are Involved in Viewer-, Object-, and Landmark- Centered Judgments About Object Location. *Journal of cognitive neuroscience*. November 2004. Vol. 16, p. 1517–1535. DOI 10.1162/0898929042568550.
- [7] CRAWFORD, J. D., MEDENDORP, W. P. and MAROTTA, J. J. Spatial Transformations for Eye–Hand Coordination. *Journal of Neurophysiology*. 2004. Vol. 92, no. 1p. 10–19. DOI 10.1152/jn.00117.2004.
- [8] FINK, G. Space-based and object-based visual attention: shared and specific neural domains. *Brain*. January 1997. Vol. 120, no. 11p. 2013–2028. DOI 10.1093/brain/120.11.2013.
- [9] FINK, G. R., MARSHALL, J. C., SHAH, N. J., WEISS, P. H., HALLIGAN, P. W., GROSSE-RUYKEN, M., ZIEMONS, K., ZILLES, K. and FREUND, H.-J. Line bisection judgments implicate right parietal cortex and cerebellum as assessed by fMRI. *Neurology*. 2000. Vol. 54, no. 6p. 1324–1331. DOI 10.1212/wnl.54.6.1324.
- [10] GALATI, Gaspare, PELLE, Gina, BERTHOZ, Alain and COMMITTERI, Giorgia. Multiple reference frames used by the human brain for spatial perception and memory. *Experimental Brain Research*. 2010. Vol. 206, no. 2p. 109–120. DOI 10.1007/s00221-010-2168-8.

- [11] GALATI, Gaspare, LOBEL, Elie, VALLAR, Giuseppe, BERTHOZ, Alain, PIZZAMIGLIO, Luigi and BIHAN, Denis Le. The neural basis of egocentric and allocentric coding of space in humans: a functional magnetic resonance study. *Experimental Brain Research*. May 2000. Vol. 133, no. 2p. 156–164. DOI 10.1007/s002210000375.
- [12] GOODFELLOW, Ian, BENGIO, Yoshua, COURVILLE, Aaron and UNDEFINED, undefined undefined. *Deep Feedforward Networks*. In : *Deep learning*. Cambridge, MA : MIT Press, 2017. p. 200–203.
- [13] HONDA, M. Cortical areas with enhanced activation during object-centred spatial information processing. A PET study. *Brain*. January 1998. Vol. 121, no. 11p. 2145–2158. DOI 10.1093/brain/121.11.2145.
- [14] HAGLER, D. J. Jr., RIECKE, L., and SERONO, M. I. Parietal and superior frontal visuospatial maps activated by pointing and saccades. *Neuroimage*. 2007 May 1; 35(4): 1562–1577. DOI 10.1016/j.neuroimage.2007.01.033
- [15] KUMAR, Arvind and BARVE, Shrish. *Newtonian Relativity*. In : *How and why in basic mechanics*. Hyderabad : Universities Press, 2002. p. 114–117.
- [16] ROJAS, R. The Backpropagation Algorithm. In : *Neural Networks*. Berlin : Springer-Verlag, 1996. p. 155–157.
- [17] SALINAS, Emilio and SEJNOWSKI, Terrence J. Book Review: Gain Modulation in the Central Nervous System: Where Behavior, Neurophysiology, and Computation Meet. *The Neuroscientist*. 2001. Vol. 7, no. 5p. 430–440. DOI 10.1177/107385840100700512.
- [18] SOECHTING, J. and FLANDERS, M. Moving In Three-Dimensional Space: Frames Of Reference, Vectors, And Coordinate Systems. *Annual Review of Neuroscience*. January 1992. Vol. 15, no. 1p. 167–191. DOI 10.1146/annurev.neuro.15.1.167.
- [19] SERENO, Martin I and HUANG, Ruey-Song. Multisensory maps in parietal cortex. *Current Opinion in Neurobiology*. 2014. Vol. 24, p. 39–46. DOI 10.1016/j.conb.2013.08.014.
- [20] XING, Jing and ANDERSEN, Richard A. Memory Activity of LIP Neurons for Sequential Eye Movements Simulated With Neural Networks. *Journal of Neurophysiology*. 2000. Vol. 84, no. 2p. 651–665. DOI 10.1152/jn.2000.84.2.651.
- [21] ZIPSER, David and ANDERSEN, Richard A. A back-propagation programmed network that simulates response properties of a subset of posterior parietal neurons. *Nature*. 1988. Vol. 331, no. 6158p. 679–684. DOI 10.1038/331679a0.