




AGI Brain: A Learning and Decision Making Framework for Artificial General Intelligence Systems Based on Modern Control Theory

Mohammadreza Alidoust^(✉) 

Mashhad, Iran

Abstract. In this paper a unified learning and decision making framework for artificial general intelligence (AGI) based on modern control theory is presented. The framework, called AGI Brain, considers intelligence as a form of optimality and tries to duplicate intelligence using a unified strategy. AGI Brain benefits from powerful modelling capability of state-space representation, as well as ultimate learning ability of the neural networks. The model emulates three learning stages of human being for learning its surrounding world. The model was tested on three different continuous and hybrid (continuous and discrete) Action/State/Output/Reward (ASOR) space scenarios in deterministic single-agent/multi-agent worlds. Successful simulation results demonstrate the multi-purpose applicability of AGI Brain in deterministic worlds.

Keywords: Artificial general intelligence · Modern control theory · Optimization · Implicit and explicit memory · Shared memory · Stages of learning · Planning · Policy · Multi-Agent · Emotions · Decision making · Continuous and hybrid ASOR space

1 Introduction

In this paper, AGI Brain, a learning and decision making framework for AGI is proposed which has a unified, simple structure and tries to emulate the stages of human learning. Based on Wang's classification [1], AGI Brain looks at intelligence as a form of optimality and tries to duplicate intelligence using a unified approach by applying state-space representation (e.g. see [2]) and neural networks as its modelling technique. In AGI Brain, intelligence is defined as “optimizing the surrounding world towards common goals”, counting the agent's body as a part of the surrounding world. AGI Brain is a model based algorithm and delays its decision making stage until it built a model upon collected data from interaction with the environment. It estimates the reward value using feedforward neural networks. Like reinforcement learning (RL) [3], AGI Brain works in both continuous and discrete Action/State/Output/Reward (ASOR) spaces. But, unlike RL, AGI Brain works only in deterministic worlds with immediate rewards as yet. AGI Brain benefits from multi-agent capability. In the multi-agent case, the agents can share their experiences easily, and they can also benefit from shared memories.

2 AGI Brain

2.1 The Agent and the World

Consider an intelligent agent ω living in the world γ which also contains the object ψ . The agent ω benefits from an artificial brain Γ which controls the behavior of ω for achieving its goals (Fig. 1).

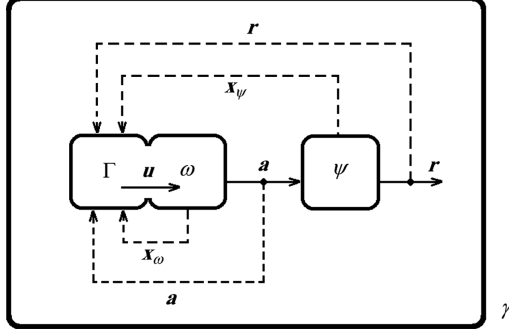


Fig. 1. The world γ consisting of the artificial agent ω and the object ψ .

At every time step n , the artificial brain Γ produces commands u (e.g. hormones or neural signals) which change the states of ω 's body, i.e. x_ω , which then leads to performing action a on the object ψ . This action changes ψ 's states x_ψ , which consequently leads to ψ 's response r . Like a natural brain, the Γ can observe these values by its sensors.

We model the agent ω by its input u , its states x_ω and its output, i.e. action a ;

$$\omega : \begin{cases} x_\omega(n+1) = f_\omega(x_\omega(n), u(n)) \\ a(n) = g_\omega(x_\omega(n), u(n)) \end{cases} \quad (1)$$

And, the object ψ by its input a , its states x_ψ and its output, i.e. response r ;

$$\psi : \begin{cases} x_\psi(n+1) = f_\psi(x_\psi(n), a(n)) \\ r(n) = g_\psi(x_\psi(n), a(n)) \end{cases} \quad (2)$$

The functions f and g are columns and they are, in general, complex nonlinear functions of the states and the inputs.

Please note that bold letters represent vectors or matrices. For simplicity, here we assume that both of ω and ψ change simultaneously.

From the Γ 's viewpoint, the ω 's body is an actuator for performing the Γ 's commands, so, the world γ can be modeled by its input u , its states $x = x_\gamma$ (vector of all states contained in the world γ , i.e. combination of x_ω and x_ψ), as well as its outputs $y = y_\gamma$ (vector of all outputs contained in the world γ , i.e. combination of a and r).

$$\mathbf{x} = \mathbf{x}_\gamma = \begin{bmatrix} \mathbf{x}_\omega \\ \cdots \\ \mathbf{x}_\psi \end{bmatrix}, \mathbf{y} = \mathbf{y}_\gamma = \begin{bmatrix} \mathbf{a} \\ \cdots \\ \mathbf{r} \end{bmatrix} \quad (3)$$

Thus, the discrete time state-space representation of the world γ will be as follows;

$$\gamma : \begin{cases} \mathbf{x}(n+1) = \mathbf{f}(\mathbf{x}(n), \mathbf{u}(n)) \\ \mathbf{y}(n) = \mathbf{g}(\mathbf{x}(n), \mathbf{u}(n)) \end{cases} \quad (4)$$

And the continuous time state-space representation of the world γ will be as follows;

$$\gamma : \begin{cases} \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)) \\ \mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)) \end{cases} \quad (5)$$

2.2 Inside the Artificial Brain

As mentioned in the previous section, in AGI Brain we look at intelligence as a form of optimality. This optimality should be observable and also measurable during the artificial life of an intelligent agent (or a swarm of them). Therefore, an optimality criterion must be defined first. We define the artificial life of an artificial intelligent agent ω living in the world γ as completing a task Q (or a series of tasks) on the object(s) ψ (or on its own body). Completion of task Q requires maximization of reward function R as the task completion criterion, which is a linear combination of objectives \mathbf{J} and their corresponding importance coefficients \mathbf{P} called personality, as well as an artificial brain Γ which is responsible for this optimization procedure.

For the reward function R we have;

$$R = \mathbf{P}^T \mathbf{J} \quad (6)$$

Thus, completion of task Q will be as follows;

$$Q : \text{Max} [R = \mathbf{P}^T \mathbf{J}] \quad (7)$$

The aim of the artificial brain Γ is to complete Q by calculating and then producing the optimal control signal \mathbf{u}^* such that the reward function R is maximized with respect to the state equations that govern the world γ , i.e. equation (4) (or (5)). So, in discrete time case, the brain Γ has to solve the following optimization problem at every time step n ;

$$\begin{aligned}
\mathbf{u}^*(n) &= \underset{\mathbf{u} \in \aleph}{\text{ArgMax}} [R = \mathbf{P}^T \mathbf{J}] \\
\text{s.t.} \quad & \begin{cases} \mathbf{x}(n+1) = \mathbf{f}(\mathbf{x}(n), \mathbf{u}(n)) \\ \mathbf{y}(n) = \mathbf{g}(\mathbf{x}(n), \mathbf{u}(n)) \end{cases}
\end{aligned} \tag{8}$$

Where \aleph is the set of all possible alternatives. The optimal command \mathbf{u}^* is the decision made by the Γ and will be transmitted to the agent's body to produce a desired action. Please note that the same problem applies to the continuous time case.

Remark 1: Objectives. The vector of objectives \mathbf{J} is a vector that contains the objectives of the agent ω on a benefit-cost basis and relates the components of the world γ to the agent's goals. In general, it is a function of time, input, states and outputs as well as the desired states \mathbf{x}_d and desired outputs \mathbf{y}_d of the world γ , i.e. $\mathbf{J} = \mathbf{J}(n, \mathbf{u}, \mathbf{x}, \mathbf{x}_d, \mathbf{y}, \mathbf{y}_d)$.

Remark 2: Personality. Personality vector \mathbf{P} is a dynamical vector which regulates the behavior of the agent and plays as the necessary condition for attaining an objective. Each element of the dynamical vector \mathbf{P} , i.e. p_k , is the coefficient of a corresponding objective j_k at the current time. In general, \mathbf{P} is a function of the current states and outputs of the world γ , i.e. $\mathbf{P} = \mathbf{P}(\mathbf{x}_0, \mathbf{y}_0)$, and based on the current situations alters the importance of each objective to be achieved, in order to regulate the behavior of the agent ω . For instance, \mathbf{P} activates fight-or-flight behavior by deactivating other objectives when ω senses a predator approaching. Also it leads to more diversity in multi-agent environments.

2.3 Learning and Decision Making

Equation (8) has two parts; the reward function and the state equations that govern the world γ . For solving Eq. (8), the artificial brain Γ has to search its available alternative set \aleph . For each alternative $\mathbf{u}^k \in \aleph$, the Γ must solve the state equations first and then evaluate the reward function. The reward function is pre-defined by the designer, but since most environments are unknown, the state equations of such environments are not available. Therefore, the agent does not know what the consequences of performing a command $\mathbf{u}^k \in \aleph$ are, how it changes the world and how the world may seem after performing each command.

In this major case, the agent has to build a model of the world which enables the agent to estimate the consequences of performing each alternative $\mathbf{u}^k \in \aleph$ on the world during the decision making stage. So, by using an estimator, the state equations of Eq. (8) turn into the following estimation problem:

$$\left\langle \begin{matrix} \hat{\mathbf{x}}(n+1) \\ \hat{\mathbf{y}}(n+1) \end{matrix} \right\rangle \xleftarrow{\text{Estimator}} \left\langle \begin{matrix} \mathbf{x}(n) \\ \mathbf{y}(n) \\ \mathbf{u}(n) \end{matrix} \right\rangle \tag{9}$$

Where $\hat{\mathbf{x}}(n+1)$ and $\hat{\mathbf{y}}(n+1)$ are the estimated states and outputs of the world γ after performing command $\mathbf{u}(n)$ on the γ with initial conditions $\mathbf{x}(n)$ and $\mathbf{y}(n)$.

For building such an estimator the agent needs a set of collected data from the world, i.e. observations during its learning stages (which will be described in the following paragraph), as well as a memory to learn those collected data. In AGI Brain, due to the high power of function approximation of neural networks, the agent is equipped with two neural network memories: explicit memory (EM) and implicit memory (IM). The EM is used as the estimator model in Eq. (9) and the IM is used for direct production of the optimal command \mathbf{u}^* without solving Eq. (8).

According to [4], there are three stages of human skill learning; (1) Cognitive: in which movements are slow, inconsistent, and inefficient, (2) Associative: in which movements are more fluid, reliable, and efficient, (3) Autonomous: in which movements are accurate, consistent, and efficient. Regarding the accuracy of the produced signals and with some connivance in the meanings, we implemented the above learning stages in our model as infancy, decision making and expert respectively. During these three learning stages, the artificial brain Γ produces three differently-produced control signals, and stores the consequences of them in its memories, EM and IM.

Infancy Stage. In the infancy stage, the agent tries to collect data from its surrounding world by interacting with it. During this stage, at every time step n , the artificial brain Γ exerts randomly-generated commands \mathbf{u} to the world with initial conditions $\mathbf{x}(n)$ and $\mathbf{y}(n)$, and then observes the results $\mathbf{x}(n+1)$ and $\mathbf{y}(n+1)$.

Each observation vector \mathbf{o} at time step n has the following form;

$$\mathbf{o}^n = [\mathbf{x}(n), \mathbf{y}(n), \mathbf{u}(n), \mathbf{x}(n+1), \mathbf{y}(n+1)]^T \quad (10)$$

The agent stores these data in its memories EM and IM. For the EM, each observation is split into vectors $\mathbf{I}_{EM}^n = [\mathbf{x}(n), \mathbf{y}(n), \mathbf{u}(n)]^T$ and $\mathbf{T}_{EM}^n = [\mathbf{x}(n+1), \mathbf{y}(n+1)]^T$ which are trained to the EM as its inputs and targets respectively. And, for the IM, each observation is split into vectors $\mathbf{I}_{IM}^n = [\mathbf{x}(n), \mathbf{y}(n), \mathbf{x}(n+1), \mathbf{y}(n+1)]^T$ and $\mathbf{T}_{IM}^n = [\mathbf{u}(n)]$ which are trained to the IM as its inputs and targets respectively. During this stage, these observations will not be used for producing the next commands.

Decision Making Stage. Using EM as the estimator in Eq. (9) and substituting Eq. (9) with state equations of Eq. (8), the decision making problem of Eq. (8) turns into the following equation;

$$\begin{aligned} \mathbf{u}^*(n) &= \underset{\mathbf{u} \in \aleph}{\text{ArgMax}} [R = \mathbf{P}^T \mathbf{J}] \\ \text{s.t.} \quad & \left\langle \hat{\mathbf{x}}(n+1) \right\rangle \xleftarrow{EM} \left\langle \begin{matrix} \mathbf{x}(n) \\ \mathbf{y}(n) \\ \mathbf{u}(n) \end{matrix} \right\rangle \end{aligned} \quad (11)$$

For solving Eq. (11), at every time step n , the artificial brain Γ searches \aleph in three stages: (1) Estimation: using its EM, the Γ estimates how the world γ would seem after executing each alternative $\mathbf{u}^k \in \aleph$, i.e. estimating the states $\hat{\mathbf{x}}(n+1)$ and the output

$\hat{\mathbf{y}}(n+1)$ of the world γ for each alternative $\mathbf{u}^k \in \aleph$, given the current states $\mathbf{x}(n)$ and outputs $\mathbf{y}(n)$, (2) Computation: computing the influence of the estimations on reward function $R = \mathbf{P}^T \mathbf{J}$ with respect to $\mathbf{J} = \mathbf{J}(n, \mathbf{u}, \hat{\mathbf{x}}, \mathbf{x}_d, \hat{\mathbf{y}}, \mathbf{y}_d)$ and $\mathbf{P} = \mathbf{P}(\mathbf{x}_0, \mathbf{y}_0)$, (3) Comparison: selecting the alternative which maximizes the reward function R the most as the optimal decision \mathbf{u}^* . The learning process of the agent is also continued in this stage by learning observations with its IM and EM at certain time steps.

Planning. Planning happens when the agent cannot solve a problem at once (because of its limitations) and has to solve it by making decisions in a series of consecutive steps. During this process, the agent estimates the total value of the reward function over different possible series of consecutive inputs, i.e. $\mathbf{U} = \{\mathbf{u}(n_1), \mathbf{u}(n_2), \dots, \mathbf{u}(n_f)\} \in \aleph$, that here we call them policies. Using its EM, the Γ estimates the consequences of the first member of each randomly-generated policy on current states and outputs of the γ , and then the consequence of the second input on the changed γ , and so on. By computing and then comparing the overall reward of all the policies, the Γ selects the optimal policy \mathbf{U}^* which has the maximum reward value over the time horizon n_f . Thus, planning can be considered as an extension of decision making and we have;

$$\mathbf{U}^* = \left\{ \mathbf{u}(n) \mid \underset{\mathbf{u} \in \aleph}{\text{ArgMax}} \sum_{n=n_1}^{n_f} [R = \mathbf{P}^T \mathbf{J}] \right\}$$

s.t.

$$\left\langle \begin{matrix} \hat{\mathbf{x}}(n+1) \\ \hat{\mathbf{y}}(n+1) \end{matrix} \right\rangle \xleftarrow{EM} \left\langle \begin{matrix} \mathbf{x}(n) \\ \mathbf{y}(n) \\ \mathbf{u}(n) \end{matrix} \right\rangle \quad (12)$$

Expert Stage. In the expert stage, the IM takes over the decision making unit. Given the initial conditions \mathbf{x}_0 and \mathbf{y}_0 as well as the desired states \mathbf{x}_d and outputs \mathbf{y}_d , the IM produces the command \mathbf{u} which transmits the current states and outputs to the desired states and outputs. In this stage the commands are produced faster and they are more accurate and efficient. The learning process of the agent is also continued in this stage by learning observations with its IM and EM at certain time steps.

Role of Emotions (Stress). The role of stress as a key factor in decision making has been implemented in our model as the exploration/exploitation ratio regulator. It empowers the agent to a further exploration of the world during its second and third stage of learning. Here, we define stress as the distance between the agent's current state/output and its desired state/output, that produces stress signal s which changes the probability of selecting the optimal decision \mathbf{u}^* (or optimal policy \mathbf{U}^*) as follows;

$$p^{\mathbf{u}=\mathbf{u}^*} = \frac{1}{1 + e^s} \quad (13)$$

Data Sharing and Shared Memories. In the multi-agent scenarios, the agents can share their experience in the form of observation vectors of Eq. (10). This helps the agents to benefit from the experiences of the other agents and make better decisions.

At a higher level, all of the agents can benefit from one shared explicit memory (SEM) and one shared implicit memory (SIM) which are trained with the observations of all of the agents. For training SEM and SIM, a matrix of observations \mathbf{O}^n is formed from the observation vector \mathbf{o}^n of each agent.

3 Simulation

3.1 Continuous ASOR Space

Function Optimization. Assume a world γ_1 which contains a function f as the object ψ_1 which is going to be optimized by a group of agents $\omega_1^1, \omega_1^2, \dots, \omega_1^N$, who have these objectives: (1) j_1 : finding the maxima of the function, (2) j_2 : social behavior by moving close to the group, and (3) j_3 : following the leader, i.e. the agent whose function value is more than the other agents. The agents have different personalities. Coefficient p_1 is positive but p_2 and p_3 are small normally distributed random numbers, so that some agents may try to get far from the group and/or the leader. Agents with negative value of p_3 get nervous when they are far from the leader and may not execute their optimal decision. They select their next moves during their three stages of learning using their SEM and SIM.

$$\gamma_1 : y = f(x_1, x_2) \\ = 20 \exp(-0.01 \sqrt{x_1^2 + x_2^2}) + \sum_{i=1}^{20} 10 \exp(-0.05 \sqrt{(x_1 - x_{1i})^2 + (x_2 - x_{2i})^2}) \quad (14)$$

Where x_{1i} and x_{2i} are random numbers that satisfy the equation: $x_{1i}^2 + x_{2i}^2 = 1000^2$. The simulation results are depicted in Fig. 2.

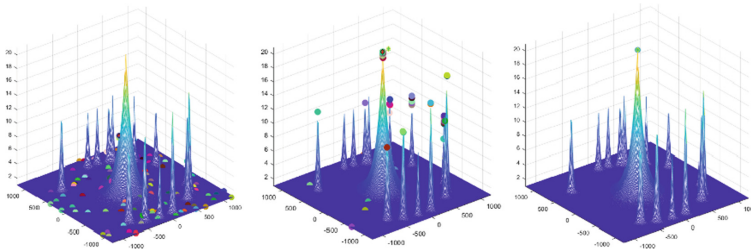


Fig. 2. Simulation results of AGI Brain on γ_1 : small spheres represent agents' positions at: (Left) the infancy stage, (Center) At epoch 35 of the decision making stage some agents reached local and global maxima, (Right) Expert stage: Starting from any initial position, they all reached the global maxima using SIM.

Tracking Control of a MIMO Nonlinear System. Assume a world γ_2 which contains a continuous stirred tank reactor (CSTR) as the object ψ_2 , which is a multi-input multi-output (MIMO) plant described by the following state-space equations¹:

$$\gamma_2 : \begin{cases} \frac{dh(t)}{dt} = w_1(t) + w_2(t) - 0.2\sqrt{h(t)} \\ \frac{dC_b(t)}{dt} = (C_{b1} - C_b(t))\frac{w_1(t)}{h(t)} + (C_{b2} - C_b(t))\frac{w_2(t)}{h(t)} - \frac{k_1 C_b(t)}{(1 + k_2 C_b(t))^2} \end{cases} \quad (15)$$

Where $h(t)$ is the liquid level, $C_b(t)$ is concentration of the output product, $w_1(t)$ and $w_2(t)$ are input flow rates, $C_{b1} = 24.9$ and $C_{b2} = 0.1$ are input concentrations, and $k_1 = k_2 = 1$ are constants associated with the rate of consumption. The single agent ω_2 has to control the liquid level and product concentration on a predefined reference $y_d = [18, 20]^T$ by its available actions $u(n) = [w_1(n), w_2(n)]^T$. Figure 3 illustrates the results of simulation.

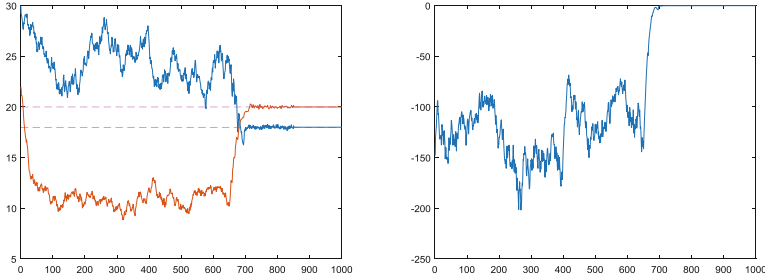


Fig. 3. Simulation results of AGI Brain on γ_2 : (Left) Plant outputs during: Infancy stage ($0 \leq n \leq 650$), decision making stage ($651 \leq n \leq 850$), and expert stage ($851 \leq n \leq 1000$). Blue line: liquid level, red line: liquid concentration, dashed red line: level reference, and dashed magenta line: concentration reference, (Right) Total reward value.

3.2 Hybrid ASOR Space

Animats. The term animats coined by S.W. Wilson refers to artificial animals which interact with real or artificial ecosystems and have the following key properties: (1) autonomy, (2) generality, and (3) adequacy. They exist in a sea of sensory signals, capable of actions, act both externally and internally, and their certain signal or absence of them have special status [5].

Assume a world γ_3 which is an artificial ecosystem that contains a cattle of grazing animats $\omega_3^1, \omega_3^2, \dots, \omega_3^N$, pasturages $\psi_3^{ps1}, \psi_3^{ps2}, \dots, \psi_3^{psM}$, and a predator ψ_3^{Pr} . The animats have the following objectives: (1) Staying alive by (a) j_1 : grazing for increasing

¹ Although the describing equations of the first and the second scenario are available, AGI Brain considers them as unknown environments, and from the observations that are gathered by interaction with these unknown environments, it builds their models in its memories and then completes its required tasks based on the models.

their energy, (b) not getting hunted by the predator by j_2 : moving as far as possible from it, and j_3 : by moving close to the cattle, (2) j_4 : Reproduction by mating, and (3) j_5 : Searching for new pasturages ψ_3^{psl} . The animats have the following discrete actions set: $\aleph = \{\text{up,down,right,left,eat,mate}\}$.

The rules of the world γ_3 are as follows: The animats can move by performing $u_1 = \text{up}, u_2 = \text{down}, u_3 = \text{right}$ and $u_4 = \text{left}$. They can only increase their energy by performing action $u_5 = \text{eat}$ when they are near a pasturage, otherwise they lose one unit of energy. They can only reproduce by performing action $u_6 = \text{mate}$ when they are mature and they have enough energy. The animats get hunted and die by the predator ψ_3^{Pr} when they are very close to it. Also, they die whether their energy level is equal to zero, or they get old. Corresponding elements of the dynamic personality vector $P(x_0, y_0)$ will change when different situations occur, e.g. when the animat ω_3^k observes the predator near to it, $p_2(x_0, y_0)$ and $p_3(x_0, y_0)$ are increased while other elements are decreased to zero. They select their next actions during their three levels of learning using their SEM and SIM. Figure 4 illustrates snap shots of simulation.

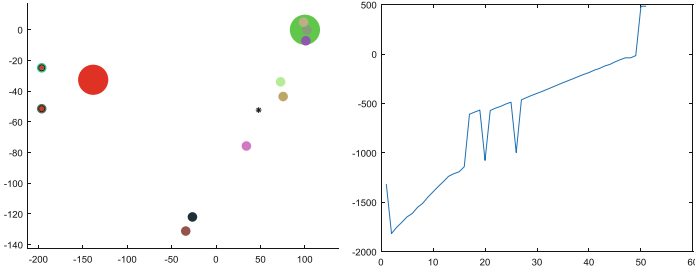


Fig. 4. Simulation snapshots of AGI Brain on γ_3 : (Left) snapshot of the agents' positions during decision making stage. The big green circle represent a pasturage, small circles represent the agents moving towards the pasturage, big red circle represent the predator, and small circles with red star represent hunted agents, (Right) Reward value of agent ω_3^k during decision making stage (color figure online).

4 Discussion and Future Works

In this paper, the AGI Brain as a learning and decision making framework was introduced. The model utilizes optimization as its theory and state-space representation as its technique. It implements the three levels of human learning using its IM and EM. Also it incorporates planning as well as some mental factors of the human intelligence like emotions and different personalities. AGI Brain was tested on two category of problems with continuous and hybrid ASOR deterministic spaces.

As can be seen in the simulation section, the model returned successful results in the above-mentioned scenarios. Utilizing different personalities helped in different behavior as well as exploration/exploitation diversity. For instance, in the first scenario, agents with different personalities found other local maxima of the function. Also,

implementation of the role of stress resulted in more exploration, e.g. in the second scenario. Additionally, the model benefits from ultimate learning power of neural networks. The learning problem of neural networks when dealing with discrete (categorical) ASORs was solved by incorporating pattern recognition neural networks, and high precision learning achieved, e.g. in the third scenario. In multi-agent scenarios, the agents utilize SEM and SIM where they can share their experiences with other agents.

Despite the successful results of AGI Brain on current scenarios, still more developments, implementations and tests on other scenarios must be performed, in order to guarantee the multi-purpose applicability of AGI Brain. Currently, the model is parameter-dependent, e.g. the time horizon of planning. In order to determining the *sufficient* amount of a parameter for successfully accomplishing a policy on a desired task, the model should also learn these parameters during its learning stages, preferably during its infancy stage. The model works only in deterministic environments as yet. For stochastic environments, the IM and EM must be empowered with memories which are able to correctly estimate, for example, the results of an action in a stochastic world. For large alternatives sets \aleph , the performance speed of the model can be improved by searching over smaller random subsets of \aleph (whose elements are randomly selected from \aleph), instead of searching over \aleph . The size of those subsets of \aleph with respect to the size of \aleph determines a trade-off between speed and accuracy.

The ideal realization goal of the AGI Brain is to be implemented as a full-neural-network artificial brain in order to seem like a natural brain. To this end, proper neural networks -that are able to completely mimic the optimization behavior of the AGI Brain- must be developed.

References

1. Wang, P.: Artificial General Intelligence, A gentle introduction, [Online]. <https://sites.google.com/site/narswang/home/agi-introduction>. Accessed Jan 2019
2. Paraskevopoulos, P.: Modern Control Engineering. Taylor & Francis Group, Didcot (2002)
3. Sutton, R.S., Barto, A.G.: Reinforcement Learning: An Introduction, 2nd edn. MIT Press, Cambridge (2018)
4. Wulf, G.: Attention and Motor Skills Learning. Human Kinetics, Champaign (2007)
5. Strannegard, C., Svangard, N., Lindstrom, D., Bach, J., Steunebrick, B.: Learning and decision making in artificial animals. J. Artif. Gen. Intell. **9**(1), 55–82 (2018)