



HAL
open science

The Illusion of Internal Joy

Claude Touzet

► **To cite this version:**

Claude Touzet. The Illusion of Internal Joy. Artificial general intelligence. 4th International Conference, AGI 2011, Mountain View, CA, USA, August 3-6, 2011. Proceedings, 6830, Springer, pp.357-362, 2011, LNAI, 10.1007/978-3-642-22887-2_43 . hal-01338041

HAL Id: hal-01338041

<https://amu.hal.science/hal-01338041v1>

Submitted on 30 Jun 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

The Illusion of Internal Joy

Claude Touzet

Adaptive and Integrative Neurosciences Lab, University of Provence
13331 Marseille, France
claude.touzet@univ-provence.fr

Abstract. J. Schmidhuber proposes a "*theory of fun & intrinsic motivation & creativity*" that he has developed over the last two decades. This theory is precise enough to allow the programming of artificial agents exhibiting the requested behaviors. Schmidhuber's theory relies on an explicit '*internal joy drive*' implemented by an '*information compression indicator*'. In this paper, we show that this indicator is not necessary as soon as the '*brain*' implementation involves associative memories, *i.e.*, hierarchical cortical maps. The '*compression factor*' is replaced by the '*smallest common activation pattern*' in our framework, with the advantage of an immediate and plausible neural implementation. Our conclusion states that the '*internal joy*' is an illusion. This remind us of the eliminative materialism position which claims that '*free-will*' is also an illusion.

Keywords: theory of neural cognition, internal joy drive, motivation, consciousness, cortical maps, unsupervised learning, associative memories.

1 Introduction

J. Schmidhuber build his "*theory of fun & intrinsic motivation & creativity*" [1] on the maximization of an '*internal joy*' that drives a reinforcement learning process. He proposes an operational description of it, a necessary step in order to provide an artifact with fun, motivation and creativity. The intrinsic reward is computed as the compression progress expressed as the number of saved bits [2]. A number of systems have been build following this recommendation, which exhibit the desired creativity behavior [3].

We have to be aware that, as soon as something such as '*intrinsic motivation*' is defined, reinforcement learning becomes *de facto* the unique valid candidate for the learning process. Reinforcement learning is certainly a very efficient way to acquire behaviors [4,5], but it is not the only one. Supervised learning and self-organization do exist. They are not considered as valid candidates for the learning of '*creativity*' behaviors because:

- in the case of supervised learning, its '*supervision*' would limit the freedom that we think is necessarily involved by '*creativity*',
- in the case of self-organization, its (self-)'*organization*' only allows to represent the data, and therefore (again) lacks the ability or freedom involved in '*creativity*'.

Both opinions are misplaced. There is no link between the learning process and the ability to escape the learned knowledge. '*Generalization*' is the ability to process correctly new unknown inputs (following the 'rules' extracted from previously learned knowledge). The generalization quality depends only on the learning samples (not the learning process) AND the implementation.

It is our goal to show here that the '*generalization*' associated to a hierarchical cortical maps implementation is able to create any behavior involving '*intrinsic motivation*'. If successful, our demonstration will also state that any learning (not only reinforcement learning) may contribute to '*creativity*' behaviors, and that '*intrinsic motivation*', also referred sometimes as '*internal joy*', is an illusion.

2 Cortex organization

The cortex is a hierarchy of cortical maps, each cortical map acting as a self-organizing associative memory that preserves the topology and distribution of the input data [6]. Each cortical map acts as a novelty filter, and stops any known situation (or part of situation). Only new unknown situation (or part of it) is allowed to proceed along the hierarchy towards maps of higher level of abstraction (Fig. 1). Behaviors are automatically generated through the cooperation of pairs of cortical maps.

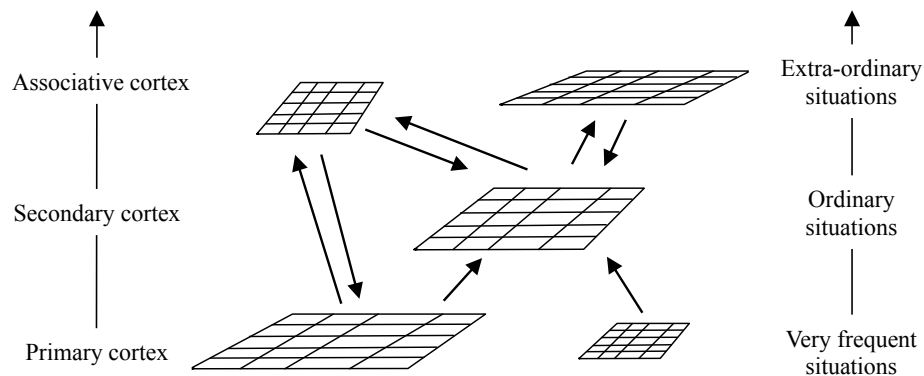


Fig. 1. The cortex is organized as a hierarchy of cortical maps, from the primary cortex receiving sensory inputs, to the secondary cortex allowing for inputs fusion, to the associative cortex where map encoding is sensory modality independent. The cortical maps are novelty filters: as the information progresses in the hierarchy, it is stopped as soon as it is recognized. only uncommon (extra-ordinary) situations (for a given individual) reach the '*abstract*' levels that account for the '*goals*'.

2.1 Behaviors are goal-directed

A behavior is a sequence of actions that can be related to the same goal. The goal is a specific situation, which will end the behavior as soon as it is reached, after what a

new behavior starts (with a new goal). During a given behavior, each action is chosen in order to reduce the distance between the present situation and the goal-situation.

Within our framework, a situation is not a x-y-z vector coding for a given location in the real world, but a vector in the cortical maps space. The cortex is build of several hundreds cortical maps, each one devoted to a specific category of information. The cortical maps pattern of activations at any given time is representative of the input (*i.e.*, experienced) situation combined with the memorized lifelong experience encoded in the synapses. A goal (situation) is therefore a multidimensional vector sharing similarities with the experienced situation.

2.2 Coordination of pairs of cortical maps

The work described in [7] explains how two associative memories (cortical maps) cooperate in order to produce goal-directed behaviors (Fig. 2). To resume (on map 1), a cluster of activity - labeled '*goal*' - is defined by some external (or internal) inputs. On the same map, the experienced current situation is represented by the fact that it activates a particular cluster of cortical columns. If there is a difference between both activities, then both activities (experienced situation and goal) will help activate cortical columns neighbors of the experienced situation. These neighbors columns encode for a situation that is close, but nevertheless different, from the current one. The difference between the two (neighbor) activities selects cortical columns on associated cortical maps (map 2). These cortical columns have memorized the actions associated to any variation of activity (of map 1). This action is done in the real world, and put the agent in a new situation (that should be closer to the goal), and everything starts again.

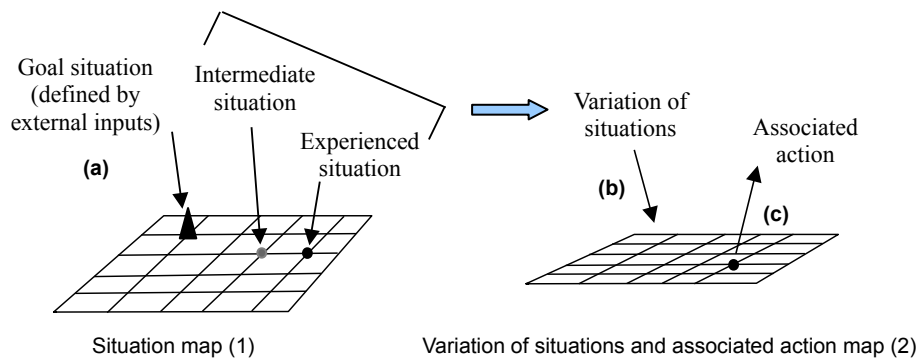


Fig. 2. Two cortical maps cooperate in order to generate a sequence of actions (*i.e.*, a behavior) in response to experienced situations. Each intersection represents a cortical column. The topology of the input situation space is preserved by the first cortical map (1). Therefore, the situation neighbor to the input situation and closer to the goal situation (a) is an intermediate situation in the process devoted to reach this '*goal*'. The variation of activity between experienced and intermediate situations (b) serves as input to the second map (2), where it activates the muscle command (action) associated (c) to this variation of situations.

3 Illusions

The definition of an illusion that is of interest here is "*a misinterpretation of a true sensation*". We claim that the '*internal joy drive*' and '*intrinsic joy*' are both misinterpretations.

3.1 The illusion of '*joy driven*' behaviors

When a behavior's goal is not seen or perceived or understood by the observer (who can be the individual himself), then the behavior is said to be '*driven by joy*'. Why '*joy*'? Just because it is easier to believe that the individual is looking for, or responding to, something pleasant when behaving - instead of the opposite (searching for bad things). Using the information provided in section 2, we claim that any behavior is '*goal driven*' and that the '*internal progress*' is in fact a '*distance-to-goal progress*'. As long as there is a discrepancy between the goal and the experienced situation locations on the cortical map(s), the behavior will occur.

3.2 The illusion of '*internal joy*'

If we assume that '*joy driven*' behaviors are just '*goal driven*' behaviors whose goals are non explicit (to a human observer), then it follows that we must get better acquainted with these goals. Moreover, since they are related to '*joy*', it is of tremendous importance for our well-being to know more about them.

First things first: how are the goals defined and selected? Cortical maps are associative memories (*i.e.*, content addressable memories). Hebbian learning induces the storage of the activity patterns generated by the (lifelong) experienced situations. The number of samples required to (self-)organize a map is several times the number of cortical columns of the map. Therefore, only the most represented situations (among the samples) are going to be memorized. For each subset of situations, the memorized information is the most redundant one, *i.e.*, the most shared cortical column activities: the '*smallest common activation pattern*'.

Next thing on the list: how does a '*joy driven*' behavior start? Let's imagine that the individual is in a situation where there is no explicit '*goal*'. There is nevertheless a situation to experience (even if it is sensory deprivation). This experienced situation will activate a representation of it by activating (some) cortical columns on some cortical maps. This activity acts as a probe of the associative memories and will sooner or later activate an already stored activity pattern: a '*goal*'. Now, the system (*i.e.*, the cortex) is faced with two different patterns of activity: one representing the experienced situation and one the '*goal*'. A behavior will emerge.

It may happen from time to time that the experienced novel situation matches exactly an already memorized - but still never experienced (otherwise it would have been stopped shortly after the primary or secondary cortex) - pattern of activity (Fig. 3). In this case, the neural pattern of activity is minimum, which allows for a much better memorization of the event (in the episodic memory). Minimum activity means that the number of cortical columns involved is minimum - but their electric

activity is maximum! It follows that this particular event will become easily remembered and may often serve as a goal (*i.e.*, an attractor situation) in the following life of the individual. Less neural activity will be required to activate this particular representation. Immediate recognition and better memorization: all the ingredients of a very meaningful experience that could be named 'joy' or 'beauty'.

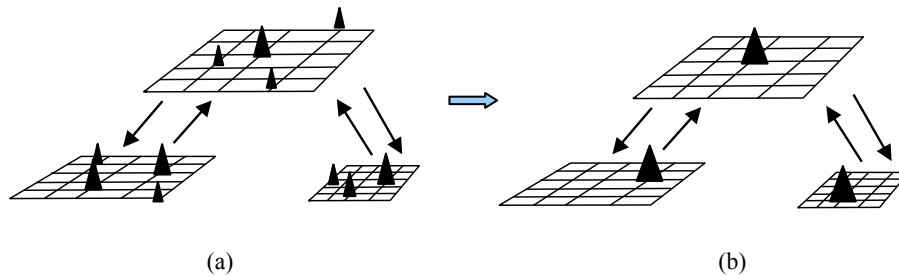


Fig. 3. (a) Columns activity (black triangles) associated to an experienced extra-ordinary situation, that quickly resume in an activity (b) involving less columns (each one exhibiting stronger activity, *i.e.*, bigger triangles). The pattern (a) is dependent on the life-long learning. It contains in essence the pattern (b), which was only waiting for this extra-ordinary experience in order to emerge. The memorization and future activation of pattern (b) is much easier than was pattern (a) activation. The quick transition from (a) to (b) is certainly a unique experience, that an observer may call 'joy'.

4 Discussion & Conclusion

We must point out the similarities between Schmidhuber's compression index and the '*smallest common activation pattern*'. Both consider that the size of the representation of information by the brain is the important factor. Both also acknowledge the fact that as soon as there is no more any variation of this '*factor*', then the behavior stops.

Looking at the biological plausibility of each proposal, a compression index needs to be computed and therefore requires resources. On the opposite, there is no computational resource required in our proposal.

Last but not least, our proposal allows for any learning (supervised, reinforcement and seal-organization) to occur, not only reinforcement learning.

'*Curiosity*' and '*boredom*' are not equivalent to '*internal joy*' even if they can be implemented using a reinforcement signal computed by the difference between the expected situation and the experienced situation [8]. To us, the discrepancy between expected and experienced situations is the root cause of the attentional processes, with the involvement of more and more higher (or abstract) cortical maps as the discrepancy hold [9]. In our view of the attentional processes, there is no need for a reinforcement signal computed elsewhere to achieve '*attention*', and therefore also '*curiosity*' and '*boredom*'.

The goal situations should be considered as attraction basins since they act in twisting the behavior of the individual towards specific representations (the attractive ones). A naive observer will see an individual whose behaviors tend to favor, even research,

some specific situations. It follows that '*internal joy*' is only a side-effect of the brain learning and memory processes. These explanations define '*internal joy*' as an illusion, as it is also the case of '*consciousness*' and '*intelligence*' in the eliminative materialism paradigm.

We may feel as if something (joy) was lost in the argument, but the gain is enormous. With a clear explanation of the '*internal joy*' illusion, it can now be efficiently looked-for by any individual. In particular, since goal situations depend only on the individual experiences, then the illusions of '*joy*' (or '*beauty*') and '*joy drive*' (or '*motivation*') are individual dependent. There is no need to covet what makes somebody else happy: our happiness lies in the (brain-)crystallization of our life, and will evolved as time passes in order to account for any variation in our daily experience. With a careful design of our daily experiences, we can change what makes us happy, and definitely improve our odds in an happiness quest.

Regarding *in silico* implementation of (human) cognition, we feel that it is definitely within our reach. The burden has been moved from the mysterious '*internal joy*' definition to the more accessible sample learning. The condition to respect in order to build a human-like cognition is that the learning base must only involve human-like situations.

References

1. Schmidhuber J.: Formal Theory of Creativity, Fun, and Intrinsic Motivation (1990-2010). IEEE Transactions on Autonomous Mental Development, 2(3), pp. 230--247 (2010)
2. Schmidhuber J.: Driven by compression progress: A simple principle explains essential aspects of subjective beauty, novelty, surprise, interestingness, attention, curiosity, creativity, art, science, music, jokes. In G. Pezzulo, M. V. Butz, O. Sigaud, and G. Baldassarre, editors, Anticipatory Behavior in Adaptive Learning Systems, from Sensorimotor to Higher-level Cognitive Capabilities, LNAI. Springer (2009)
3. Schmidhuber J.: Simple Algorithmic Principles of Discovery, Subjective Beauty, Selective Attention, Curiosity & Creativity. In V. Corruble, M. Takeda, E. Suzuki, eds., Proc. 10th Intl. Conf. on Discovery Science, pp. 26--38, LNAI 4755, Springer (2007).
4. Santos J. M. , Touzet C.: Exploration Tuned Reinforcement Function. Neurocomputing, Vol. 28 , No. 1-3, pp. 93--105 (1999)
5. Touzet C.: Q-learning for robots. In: The Handbook of Brain Theory and Neural Networks (Second Edition), M. Arbib editor, MIT Press, pp. 93--937 (2003)
6. Kohonen T.: Self-Organizing Maps. Springer Series in Information Sciences, Vol. 30, Third Ed., 501 pages. ISBN 3-540-67921-9 (2001)
7. Touzet C.: Modeling and Simulation of Elementary Robot Behaviors using Associative Memories. International Journal of Advanced Robotic Systems, Vol 3, n° 2, pp. 165--170 (2006)
8. Schmidhuber J.: A possibility for implementing curiosity and boredom in model-building neural controllers. In: J. A. Meyer and S. W. Wilson, editors, Proc. of the International Conference on Simulation of Adaptive Behavior: From Animals to Animats, pp. 222--227. MIT Press/Bradford Books (1991)
9. Touzet C.: Conscience, intelligence, libre-arbitre ? Les réponses de la Théorie neuronale de la Cognition (in French). Ed. la Machotte, 156 pages, ISBN: 978-2-919411-00-9 (2010)