

Andreas Holzinger

706.315 Selected Topics on Knowledge Discovery:  
Interactive Machine Learning

2015W, SE, 2.0 h, 3.0 ECTS

Week 45 - 06.11.2015 10:00-11:30

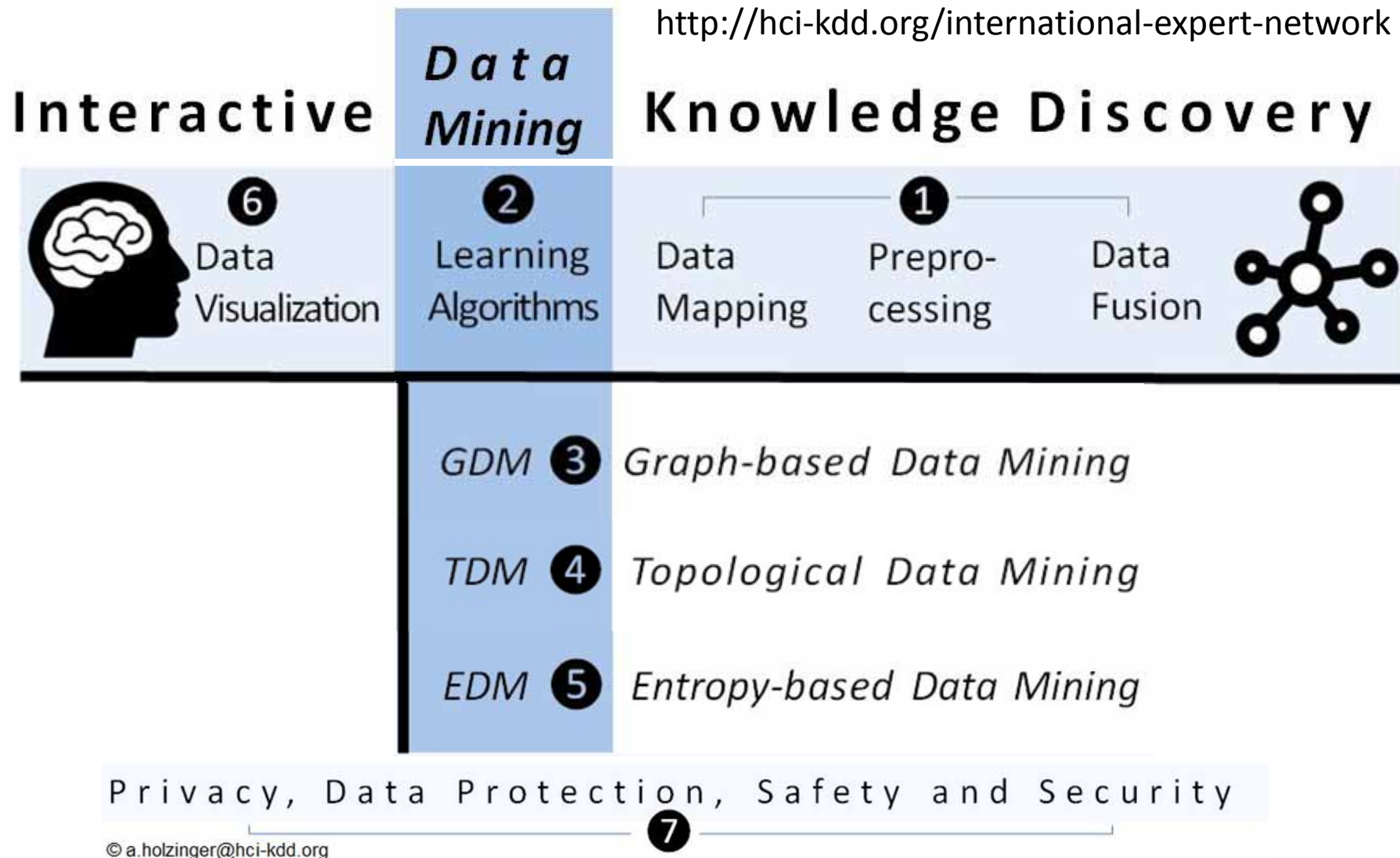
# Reinforcement Learning (RL)

a.holzinger@hci-kdd.org

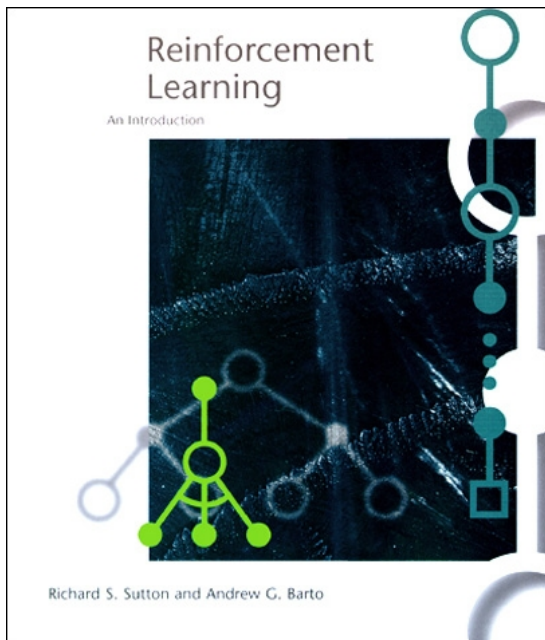
<http://hci-kdd.org/lv-706-315-interactive-machine-learning>



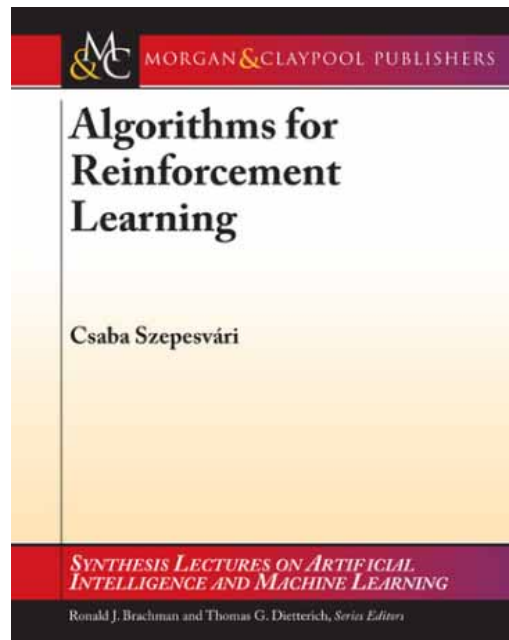
<http://hci-kdd.org/international-expert-network>



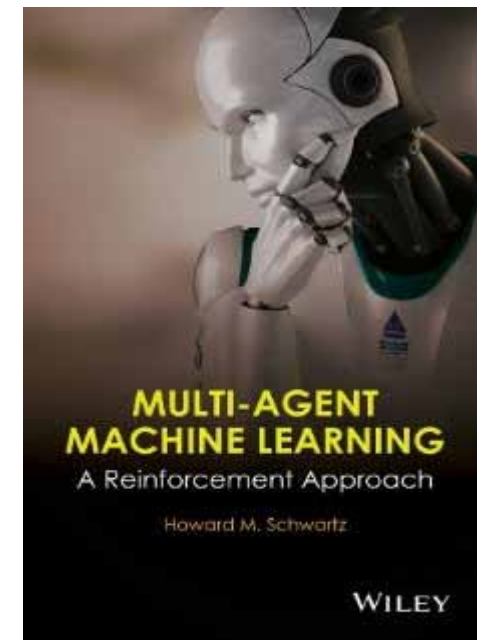
Holzinger, A. 2014. Trends in Interactive Knowledge Discovery for Personalized Medicine: **Cognitive Science meets Machine Learning**. IEEE Intelligent Informatics Bulletin, 15, (1), 6-14.



Sutton, R. S. & Barto, A. G. 1998.  
Reinforcement learning: An introduction, Cambridge, MIT press  
<http://webdocs.cs.ualberta.ca/~sutton/book/the-book.html>  
Second edition in preparation:  
[people.inf.elte.hu/lorincz/Files/RL\\_2006/SuttonBook.pdf](http://people.inf.elte.hu/lorincz/Files/RL_2006/SuttonBook.pdf)



Szepesvári, C. 2010.  
Algorithms for reinforcement learning. Synthesis lectures on artificial intelligence and machine learning, edited by R.J. Brachman and T. G. Dietterich, Morgan & Claypool.  
<http://www.ualberta.ca/~szepesva/RLBook.html>



Schwartz, H. M. (2014). Multi-agent machine learning: A reinforcement approach: John Wiley & Sons.  
<http://www.sce.carleton.ca/faculty/schwartz/index.html>

- **1) What is RL? Why is it interesting?**
- **2) Decision Making under uncertainty**
- **3) Roots of RL**
- **4) Cognitive Science of RL**
- **5) The Anatomy of an RL agent**
- **6) Example: Multi-Armed Bandits**
- **7) RL-Applications in health**
- **8) Future Outlook**

- 1) Given  $x, y$ ; find  $f$  that map a new  $x \mapsto y$   
(S/U/R?)
- 2) Finding similar points in high-dim  $X$  (S/U/R)?
- 3) Learning from interaction to achieve a goal  
(S/U/R)?
- 4) Human expert provides examples (S/U/R)?
- 5) Automatic learning by interaction with  
environment (S/U/R)?
- 6) The agent gets a scalar reward from the  
environment (S/U/R)?

# 1) What is RL?

## Why is it interesting?

*“I want to understand intelligence and how minds work. My tools are computer science, statistics, mathematics, and plenty of thinking”  
Nando de Freitas, Univ. Oxford and Google.”*



1-S; 2-U; 3-R; 4-S; 5-R; 6-R

- I) Supervised learning (classification)

- $y = f(x)$
- Given  $x, y$  pairs; find a  $f$  that map a new  $x$  to a proper  $y$
- Regression, logistic regression, classification
- Expert provides examples e.g. classification of clinical images
- Disadvantage: Supervision can be expensive

- II) Unsupervised learning (clustering)

- $f(x)$
- Given  $x$  (features only), find  $f$  that gives you a description of  $x$
- Find similar points in high-dim  $X$
- E.g. clustering of medical images based on their content
- Disadvantage: Not necessarily task relevant

- III) Reinforcement learning

- $y = f(x)$
- more general than supervised/unsupervised learning
- learn from interaction to achieve a goal
- Learning by direct interaction with environment (automatic ML)
- Disadvantage: broad difficult approach, problem with high-dim data



- Reinforcement Learning is the **oldest approach**, with the longest history, thus can provide insight into understanding human learning [1]
- RL is the **“AI problem in the microcosm”** [2]
- Future opportunities are in Multi-Agent RL (MARL), Multi-Task Learning (MTL), Generalization and **Transfer-Learning** [3], [4].

[1] Turing, A. M. 1950. Computing machinery and intelligence. *Mind*, 59, (236), 433-460.

[2] Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. *Nature*, 521, (7553), 445-451, doi:10.1038/nature14540.

[3] Taylor, M. E. & Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. *The Journal of Machine Learning Research*, 10, 1633-1685.

[4] Pan, S. J. & Yang, Q. A. 2010. A Survey on Transfer Learning. *IEEE Transactions on Knowledge and Data Engineering*, 22, (10), 1345-1359, doi:10.1109/tkde.2009.191.



- **General Purpose Learning Machines**
- Learning and General
- AGI = Artificial General Intelligence
- Remember the General Problem Solver (GPS) by Newell & Shaw

Brooks, R., Hassabis, D., Bray, D. & Shashua, A. 2012. Turing centenary: Is the brain a good model for machine intelligence? Nature, 482, (7386), 462-463, doi:10.1038/482462a.

The video player displays a presentation slide titled "Reinforcement Learning Framework". The slide illustrates the interaction between an Agent and an Environment. The Agent (represented by a person icon) sends ACTIONS to the Environment (represented by a tree icon). The Environment sends OBSERVATIONS back to the Agent. A GOAL is indicated by a speech bubble pointing to the Agent. The slide also includes logos for Google DeepMind and General Artificial Intelligence.

THE ROYAL SOCIETY

Transforming our future  
conference series  
royalsociety.org/events

Watch more videos at:  
royalsociety.org/videos

THE ROYAL SOCIETY

1:34:00 / 1:56:50

CC HD

Future directions of machine learning: Part 2

The Royal Society

Subscribe 16,807

10,704 views

+ Add to Share ... More

93 2

<https://youtu.be/XAbLn66iHcQ?t=1h34m>

# Deep Learning

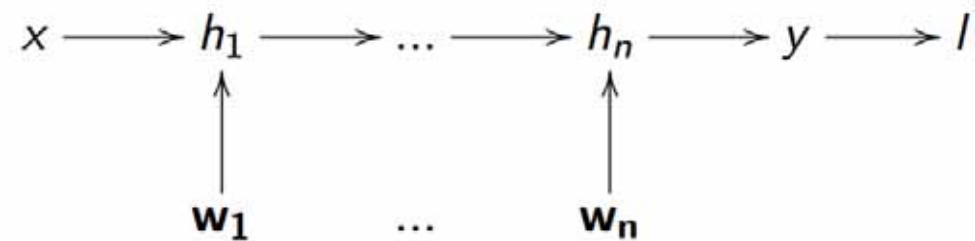
- “Deep Learning” is a buzzword for rebranding neural networks, in particular deep belief networks made of multiple layers

<http://www.forbes.com/sites/kevinmurnane/2016/04/01/what-is-deep-learning-and-how-is-it-useful/#6271f90510f0>

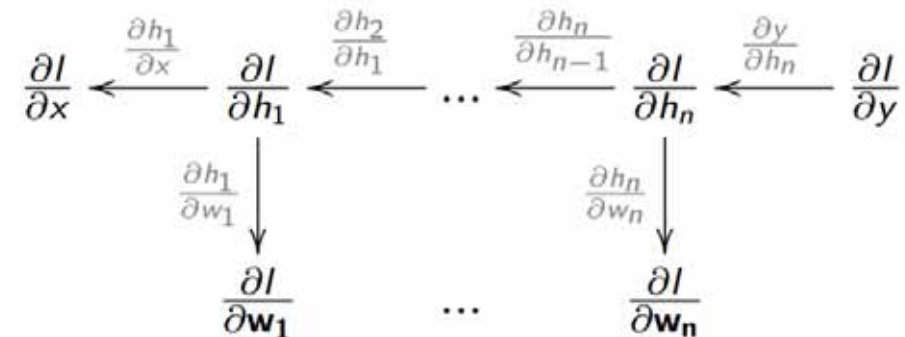
<http://www.forbes.com/sites/kevinmurnane/2016/04/01/thirteen-companies-that-use-deep-learning-to-produce-actionable-results/#7e42ab467967>

<http://www.deeplearningbook.org/>

- ▶ A **deep representation** is a composition of many functions

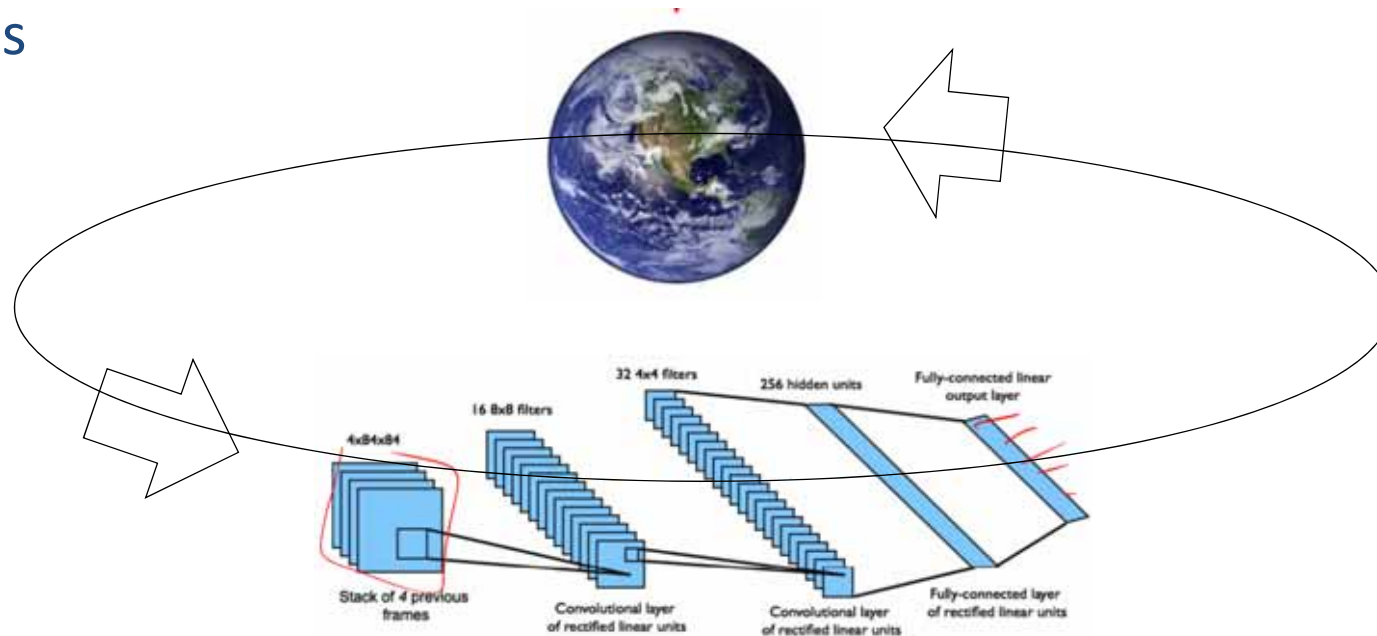


- ▶ Its gradient can be **backpropagated** by the chain rule



# Deep Reinforcement Learning

- Combination of deep neural networks with reinforcement learning = Deep Reinforcement Learning
- Weakness of classical RL is that it is not good with high-dimensional sensory inputs
- Advantage of DRL: Learn to act from high-dimensional sensory inputs



Volodymyr Mnih et al (2015), <https://sites.google.com/a/deepmind.com/dqn/>  
<https://www.youtube.com/watch?v=iqXKQf2BOSE>



- RL = general framework for decision making
- An agent is able to act
- Each action influences the agents future state
- The success of learning is measured by a scalar
- Goal: Select actions to maximize future rewards

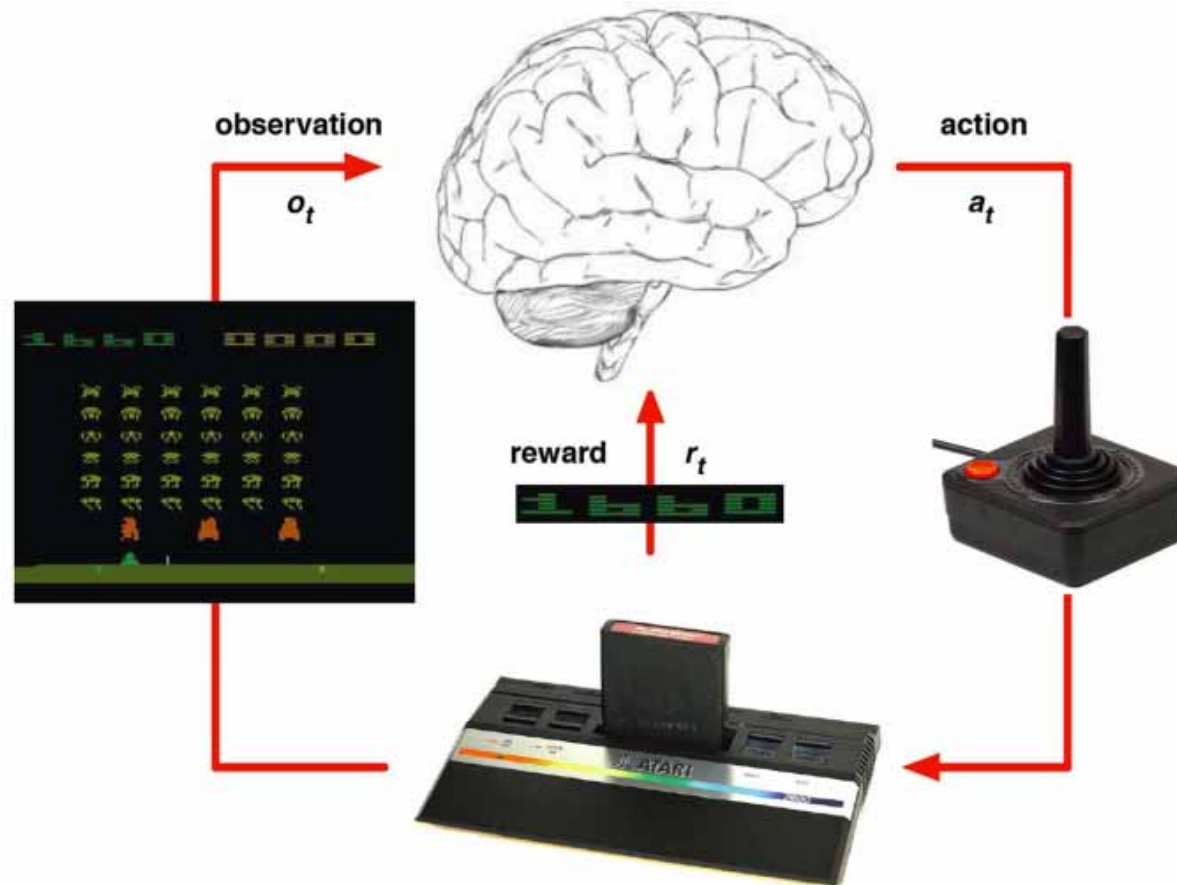
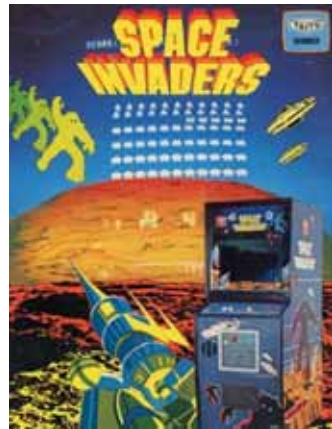
- DL is a general framework for representation learning
- 1) Given an objective
- 2) Learn representation required to achieve 1)
- Directly from raw input signals
- Using minimal domain knowledge

- Single agent shall solve any human-level task
- RL defines the objective
- DL provides the mechanism
- RL + DL → General intelligence
- this is the grand goal of Google Deepmind:  
Solve intelligence ... then solve everything else!

- Playing games: Atari [1], Go [2], ...
- Exploring worlds: Labyrinths, 3D-worlds, ...
- Controlling physical systems: manipulate, ...
- Interacting with humans: recommender, ...

[1] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. & Riedmiller, M. 2013. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602.

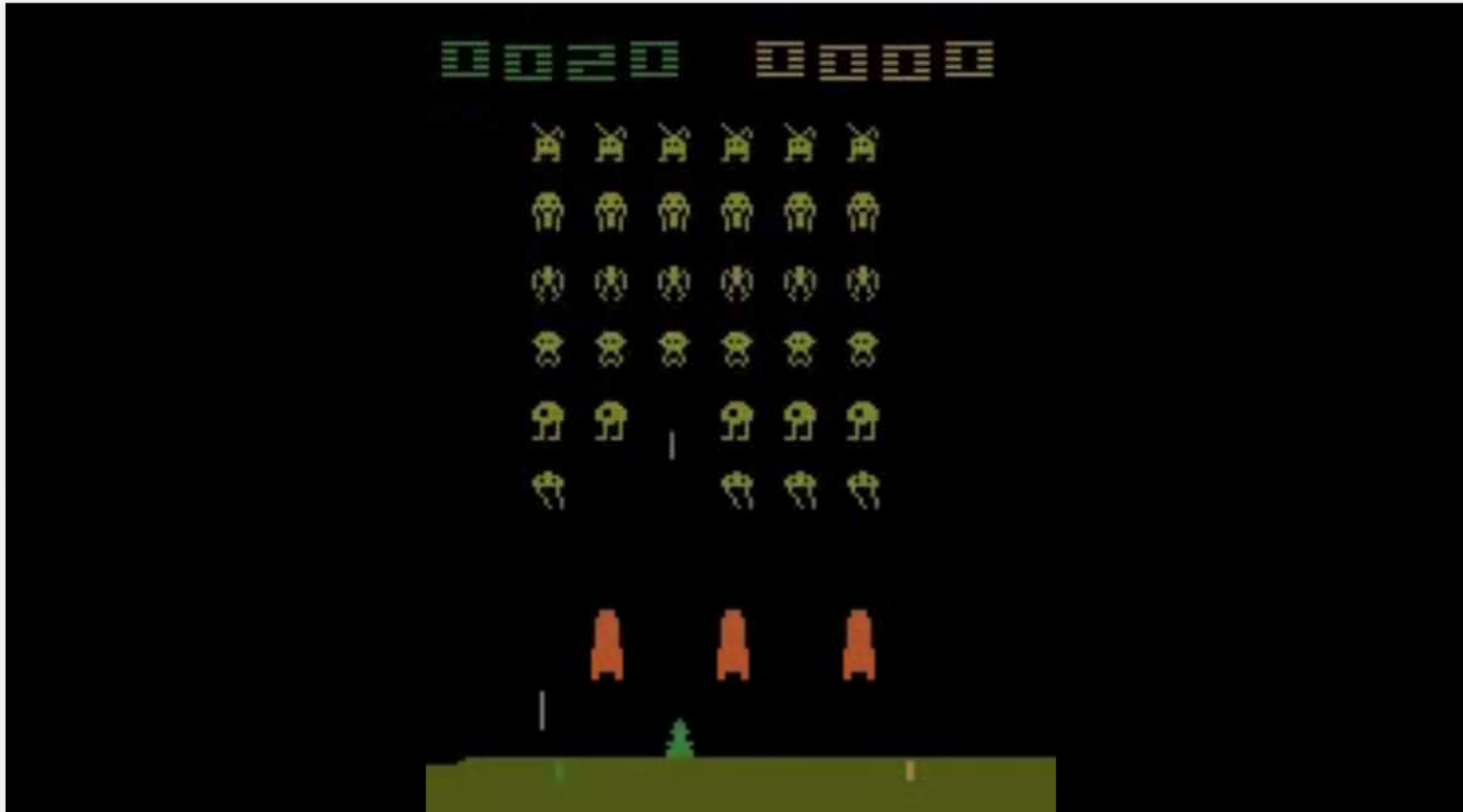
[2] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T. & Hassabis, D. 2016. Mastering the game of Go with deep neural networks and tree search. Nature, 529, (7587), 484-489, doi:10.1038/nature16961.



Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S. & Hassabis, D. 2015. Human-level control through deep reinforcement learning. Nature, 518, (7540), 529-533, doi:10.1038/nature14236



reinforcement learning space invaders



## Deep Q network playing Space Invaders



eldubro

Subscribe 11

1,855

Add to Share More

1 0

Up next



**Richard S. Sutton**

Professor of Computing Science, University of Alberta.  
Bestätigte E-Mail-Adresse bei richsutton.com  
Zitiert von: 46277

artificial intelligence reinforcement learning machine learning cognitive science computer science

**Yu-Jen Chen**

Electrical Engineering, Chung Cheng University  
Zitiert von: 28358

Reinforcement Learning Robotics

**Thomas Dietterich**

Distinguished Professor of Computer Science, Oregon State University  
Bestätigte E-Mail-Adresse bei cs.orst.edu  
Zitiert von: 26014

Machine Learning Computational Sustainability Artificial Intelligence Reinforcement Learning

**Michael L. Littman**

Professor of Computer Science, Brown University  
Bestätigte E-Mail-Adresse bei cs.brown.edu  
Zitiert von: 25879

Artificial Intelligence Reinforcement learning

**Sandeep Singh**

Professor, Computer Science & Engineering, University of Michigan  
Bestätigte E-Mail-Adresse bei umich.edu  
Zitiert von: 20923

Reinforcement Learning Computational Game Theory Artificial Intelligence

**Michael J. Frank**

Professor, Brown University  
Bestätigte E-Mail-Adresse bei brown.edu  
Zitiert von: 11482

Computational Psychiatry Dopamine Cognitive Control Reinforcement Learning Computational Neuroscience

**Robert Babuska**

Professor of Intelligent Control and Robotics, Delft University of Technology  
Bestätigte E-Mail-Adresse bei tudelft.nl  
Zitiert von: 10567

Computational Intelligence Systems and Control Robotics Nonlinear System Identification Reinforcement learning

**Chuck Anderson**

professor of computer science, colorado state university  
Bestätigte E-Mail-Adresse bei cs.colostate.edu  
Zitiert von: 7635

machine learning reinforcement learning brain-computer interface neural networks

**Csaba Szepesvari**

Department of Computing Science, University of Alberta  
Bestätigte E-Mail-Adresse bei cs.ualberta.ca  
Zitiert von: 6719

machine learning learning theory online learning reinforcement learning Markov Decision Processes

**Jan Peters**

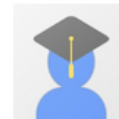
Professor at Technische Universität Darmstadt and Researcher at MPI for Intelligent ...  
Bestätigte E-Mail-Adresse bei ias.tu-darmstadt.de  
Zitiert von: 6668

Robot Learning Reinforcement Learning Machine Learning Robotics Biomimetic Systems

**Thore Graepel**

Research Scientist, Google DeepMind, and Professor of Computer Science, UCL  
Bestätigte E-Mail-Adresse bei ucl.ac.uk  
Zitiert von: 5931

Machine Learning Probabilistic Modelling Reinforcement Learning Deep Learning

**Alan Pickering**

Professor of Psychology  
Bestätigte E-Mail-Adresse bei gold.ac.uk  
Zitiert von: 5462

personality learning reward cognitive control reinforcement learning

**Daeyeol Lee**

Professor of Neurobiology, Yale University School of Medicine  
Bestätigte E-Mail-Adresse bei yale.edu  
Zitiert von: 5110

Neuroscience decision making neuroeconomics reinforcement learning prefrontal cortex

**Lihong Li (李力鸿)**

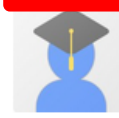
Researcher, Microsoft Research  
Bestätigte E-Mail-Adresse bei microsoft.com  
Zitiert von: 4974

Reinforcement Learning Machine Learning Artificial Intelligence

**Yael Niv**

Professor of Psychology and Neuroscience, Princeton University  
Bestätigte E-Mail-Adresse bei princeton.edu  
Zitiert von: 4865

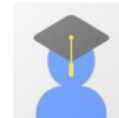
reinforcement learning neuroeconomics fMRI cognitive neuroscience computational neuroscience



University of California at Berkeley

Bestätigte E-Mail-Adresse bei berkeley.edu  
Zitiert von: 4639

Decision-making reinforcement learning

**Doina Precup**

McGill University  
Bestätigte E-Mail-Adresse bei cs.mcgill.ca  
Zitiert von: 4638

Artificial Intelligence machine learning reinforcement learning

**Naoshige Uchida**

Professor of Molecular and Cellular Biology, Harvard University  
Bestätigte E-Mail-Adresse bei mcb.harvard.edu  
Zitiert von: 4409

Neurobiology Decision Making Reinforcement learning Dopamine Olfaction

**Michael Bowling**

University of Alberta  
Bestätigte E-Mail-Adresse bei cs.ualberta.ca  
Zitiert von: 4380

Artificial Intelligence Machine Learning Game Theory Reinforcement Learning Computer Games

Status as of 03.04.2016



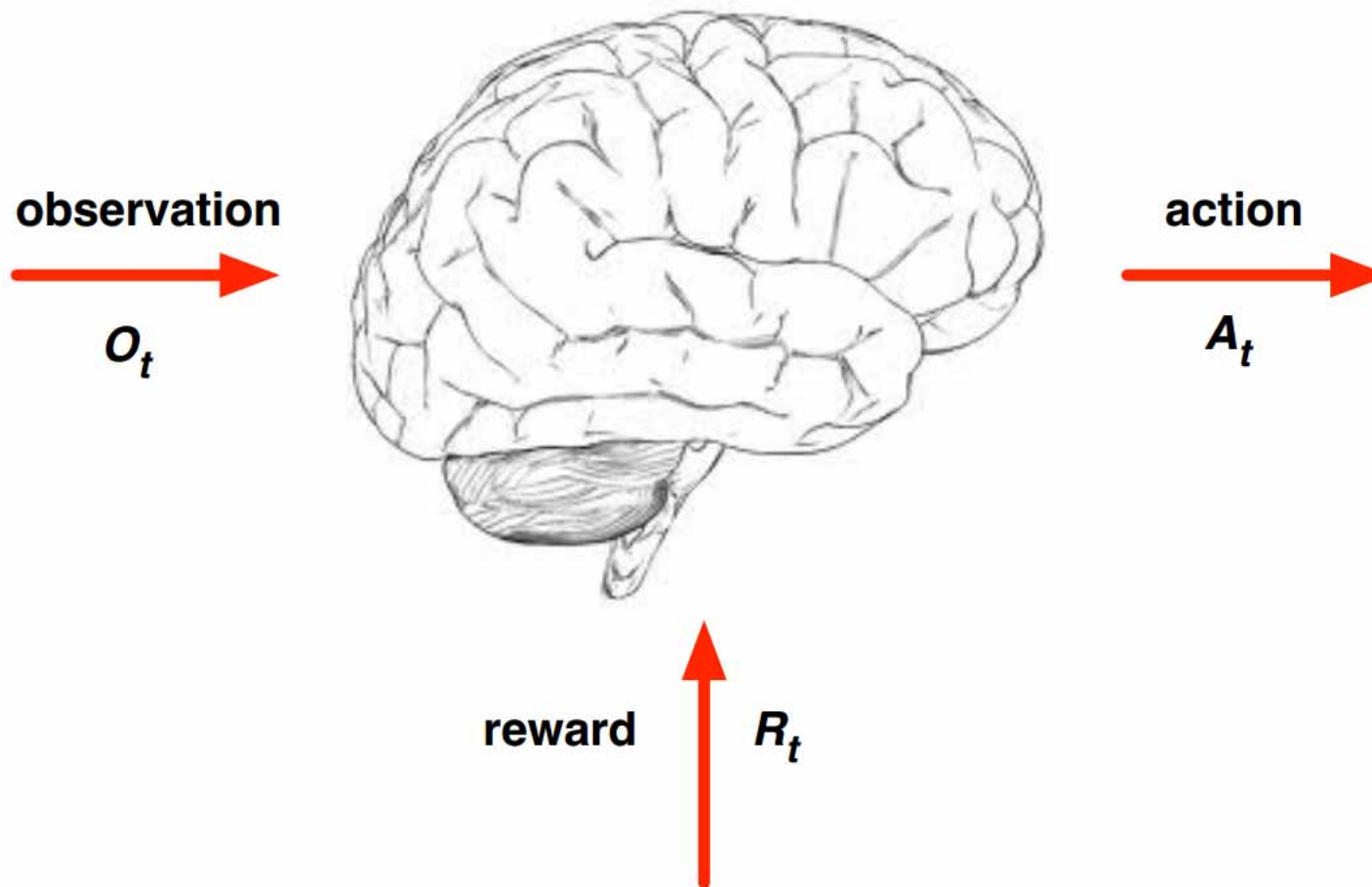
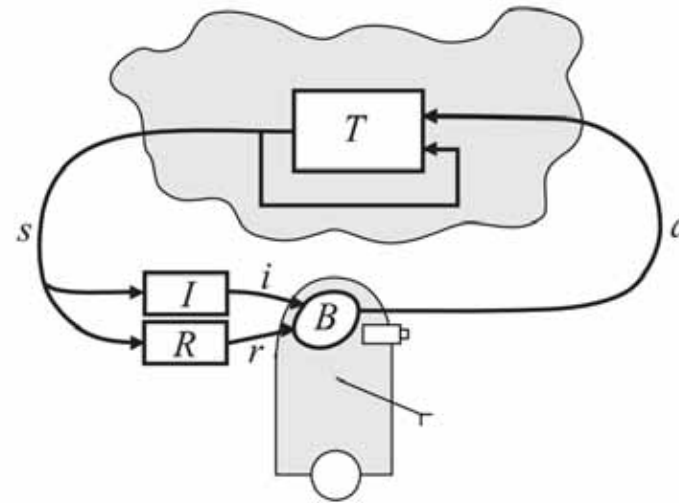


Image credit to David Silver, UCL



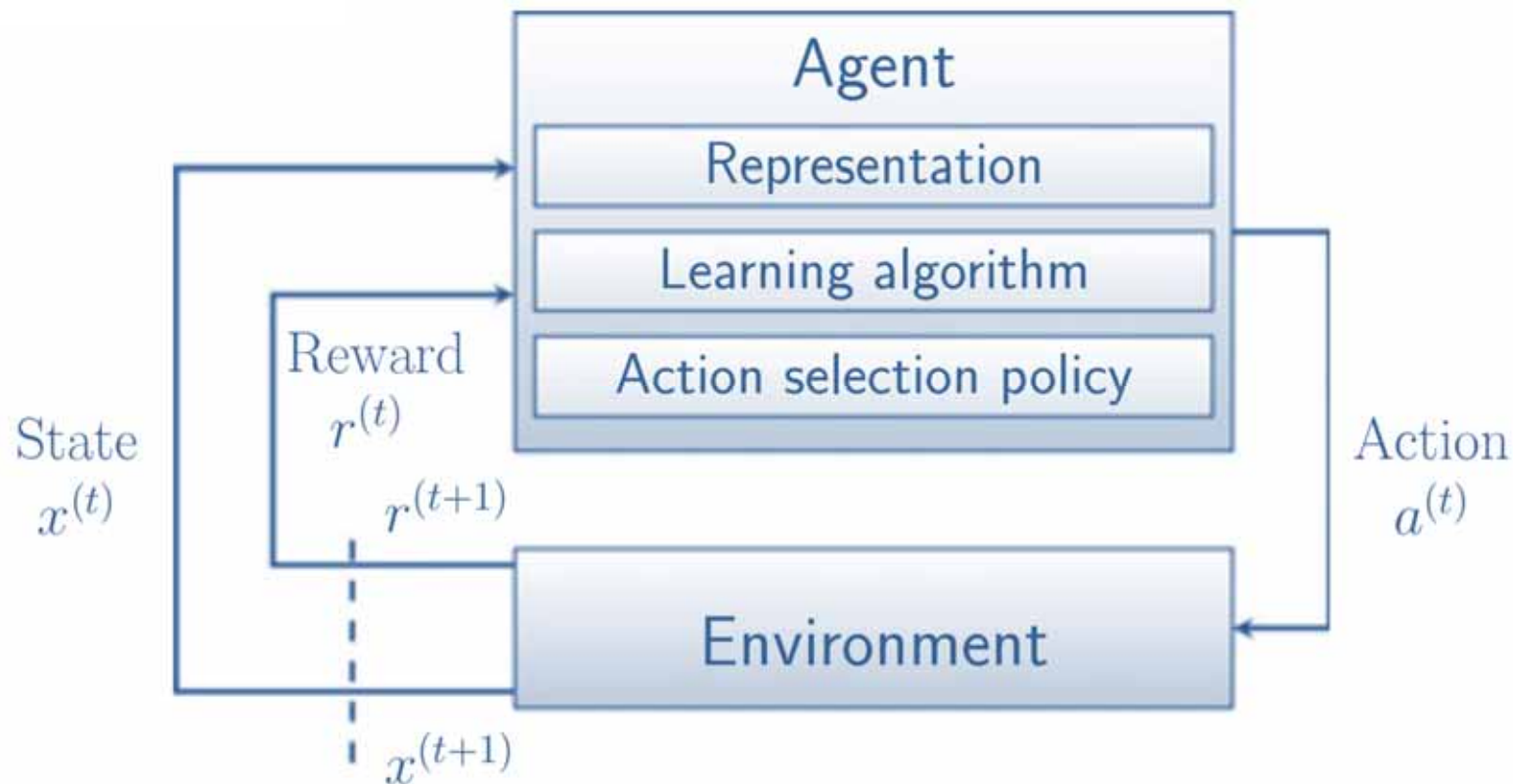
```
initialize  $V(s)$  arbitrarily
loop until policy good enough
  loop for  $s \in \mathcal{S}$ 
    loop for  $a \in \mathcal{A}$ 
       $Q(s, a) := R(s, a) + \gamma \sum_{s' \in \mathcal{S}} T(s, a, s') V(s')$ 
       $V(s) := \max_a Q(s, a)$ 
    end loop
  end loop
end loop
```

Kaelbling, L. P., Littman, M. L. & Moore, A. W. 1996. Reinforcement learning: A survey. Journal of Artificial Intelligence Research, 4, 237-285.

```

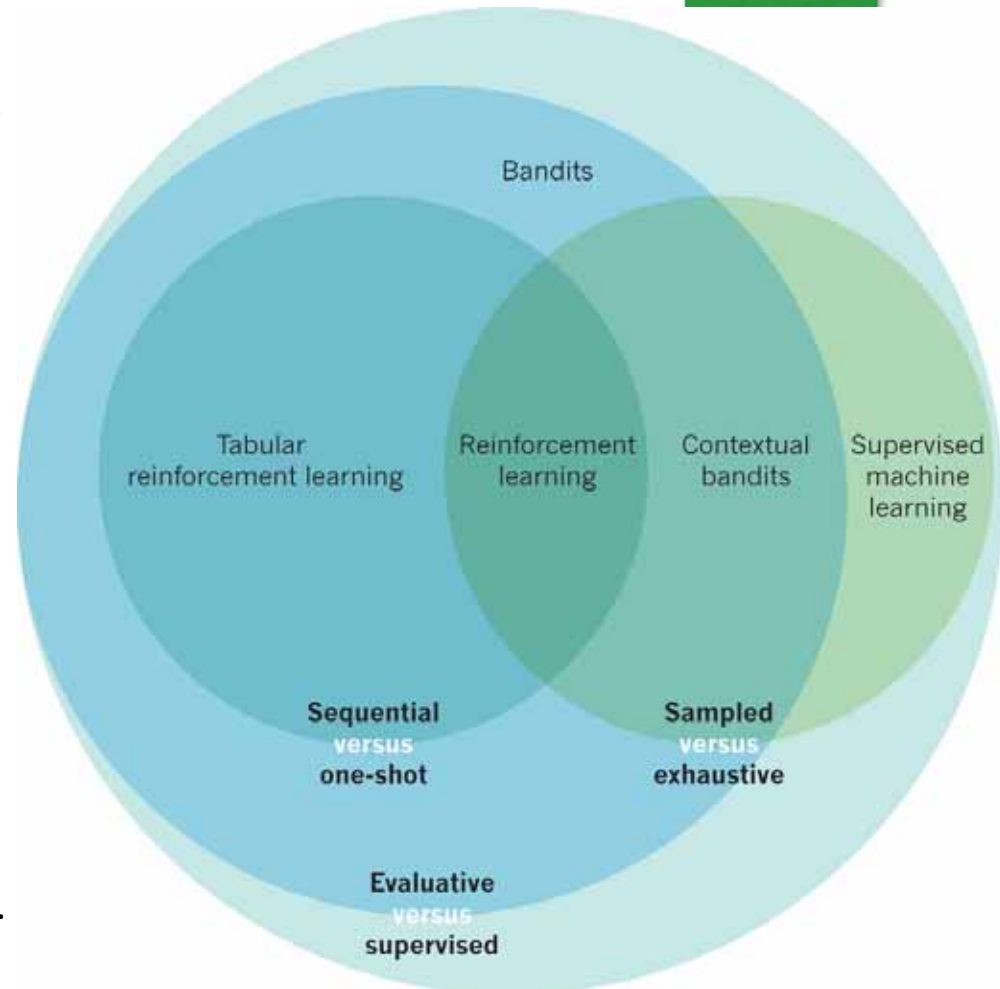
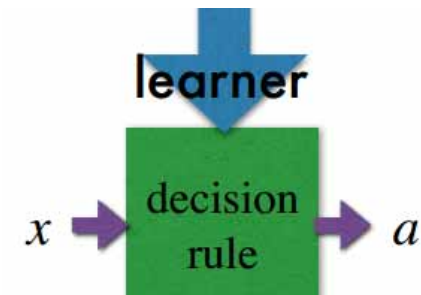
for  $t = 1, \dots, n$  do
  The agent perceives state  $s_t$ 
  The agent performs action  $a_t$ 
  The environment evolves to  $s_{t+1}$ 
  The agent receives reward  $r_t$ 
end for
    
```

**Intelligent behavior** arises from the actions of an individual seeking to **maximize its received reward** signals in a **complex and changing world**



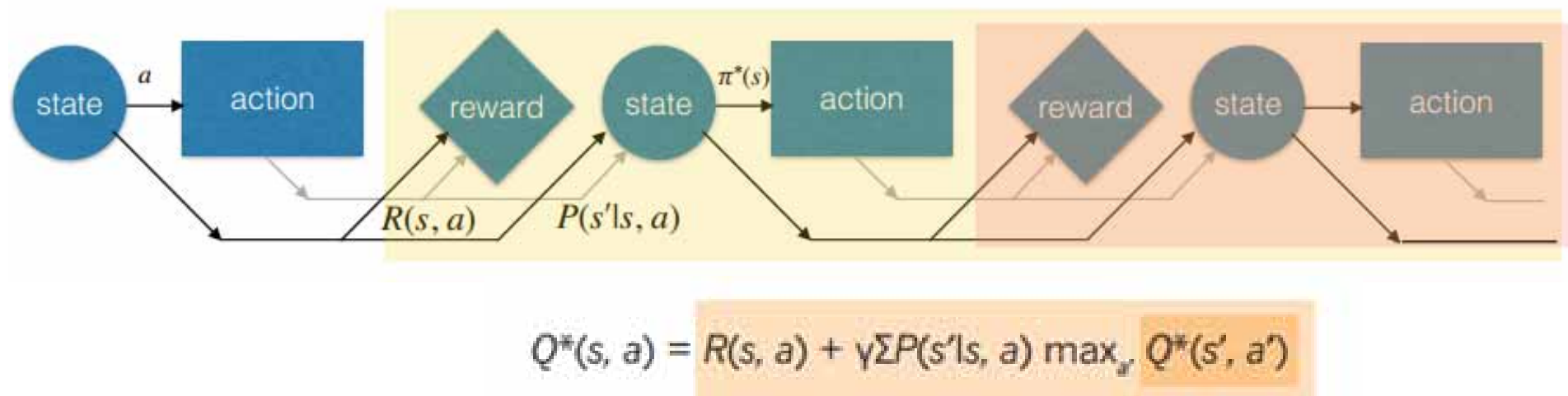
Sutton, R. S. & Barto, A. G. 1998. Reinforcement learning: An introduction, Cambridge MIT press

- Supervised:  
Learner told best  $a$
- Exhaustive:  
Learner shown every possible  $x$
- One-shot: Current  $x$  independent of past  $a$



Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451.

- Markov decision processes specify setting and tasks
- Planning methods use knowledge of  $P$  and  $R$  to compute a good policy  $\pi$
- Markov decision process model captures both sequential feedback and the more specific one-shot feedback (when  $P(s'|s, a)$  is independent of both  $s$  and  $a$ )



Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451.

- 1) Observes
- 2) Executes
- 3) Receives Reward
- Executes action  $A_t$ :
- $O_t = sa_t = se_t$
- Agent state =  
environment state =  
information state
- Markov decision  
process (MDP)

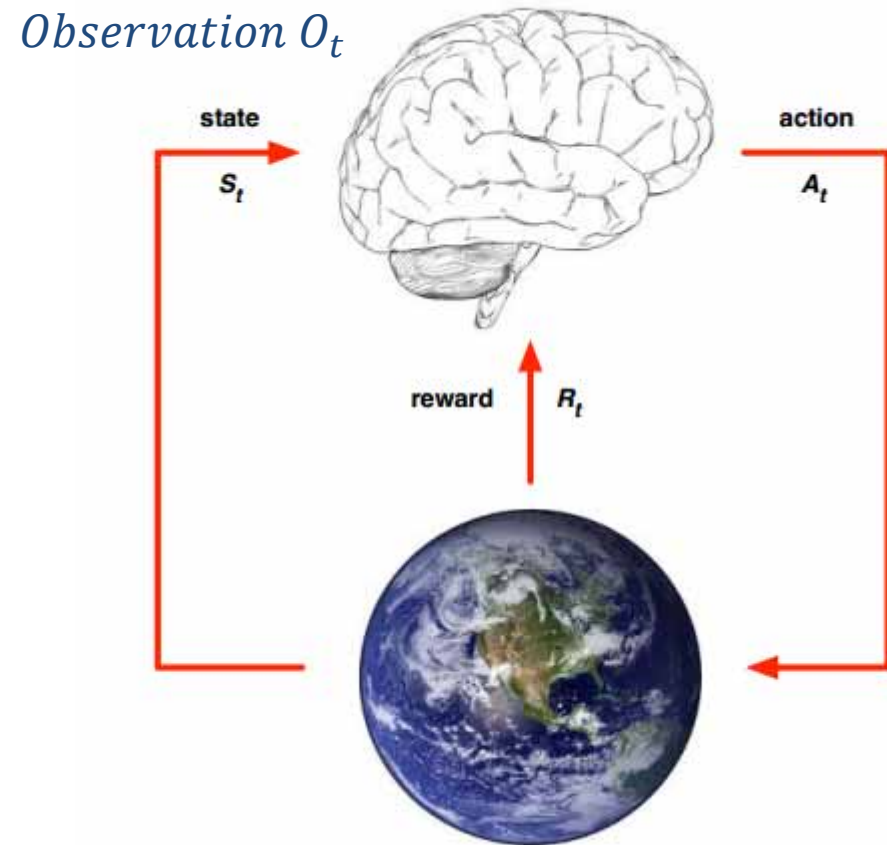
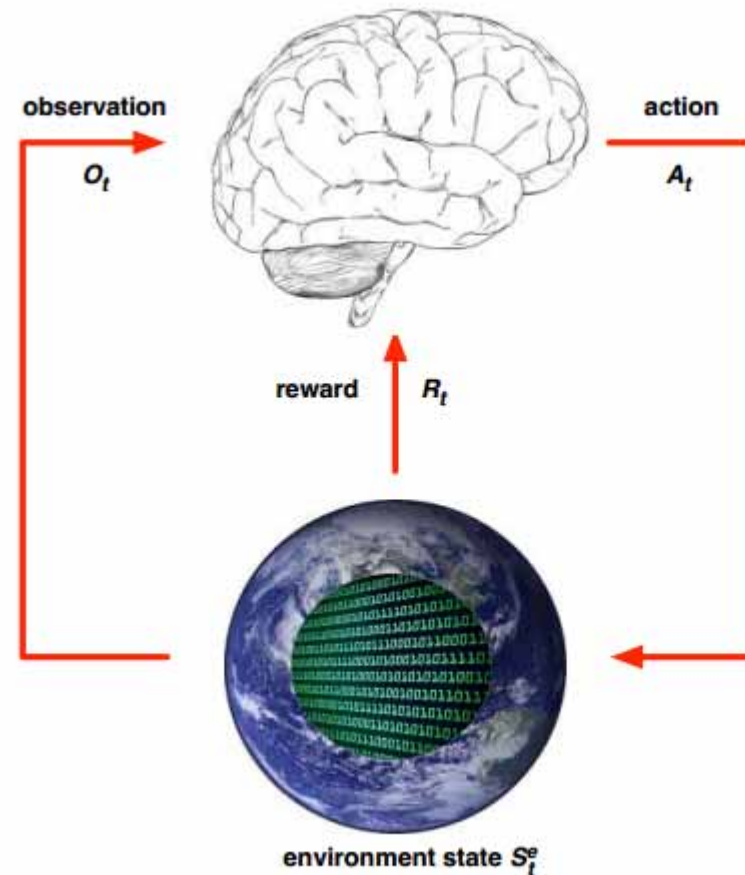


Image credit to David Silver, UCL

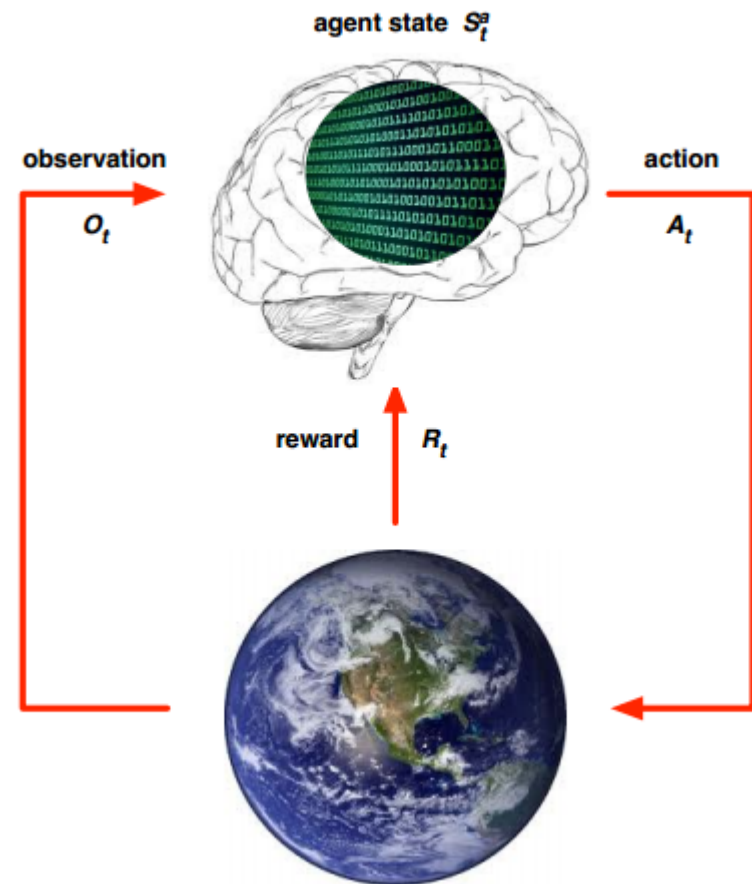
- i.e. whatever data the environment uses to pick the next observation/reward
- The environment state is not usually visible to the agent
- Even if  $S$  is visible, it may contain irrelevant information
- A State  $S_t$  is Markov iff:

$$\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t]$$





- i.e. whatever information the agent uses to pick the next action
- it is the information used by reinforcement learning algorithms
- It can be any function of history:
- $S = f(H)$

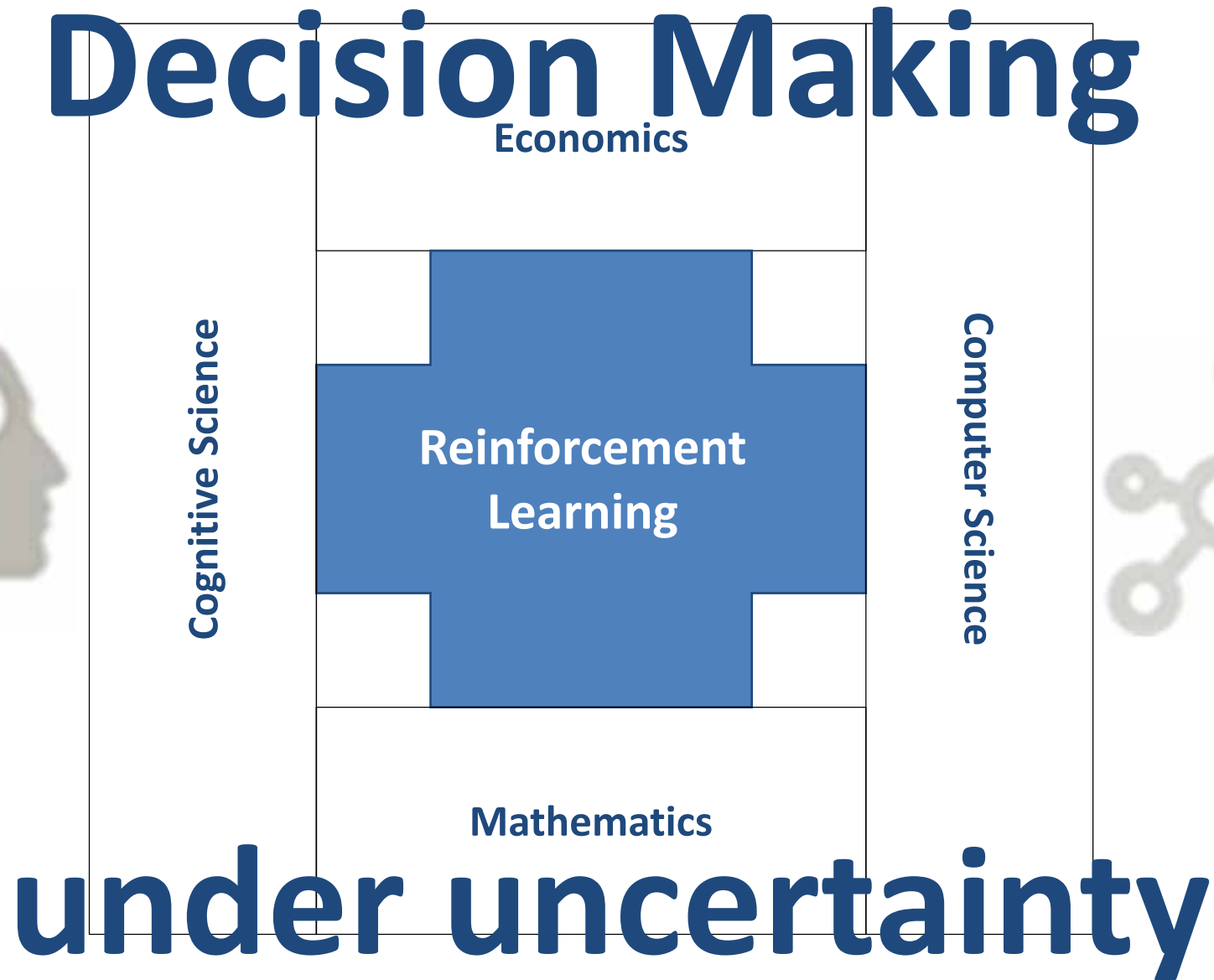


$$H_t = O_1, R_1, A_1, \dots, A_{t-1}, O_t, R_t$$

- RL agent components:
  - Policy: agent's behaviour function
  - Value function: how good is each state and/or action
  - Model: agent's representation of the environment
- Policy as the agent's behaviour
  - is a map from state to action, e.g.
  - Deterministic policy:  $a = (s)$
  - Stochastic policy:  $(a|s) = P[A_t = a | S_t = s]$
- Value function is prediction of future reward:

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

- Partial observability: when agent only indirectly observes environment (robot which is not aware of its current location; good example: Poker play: only public cards are observable for the agent):
- Formally this is a partially observable Markov decision process (POMDP):
  - Agent must construct its own state representation  $S$ , for example:
- Complete history:  $S_t^a = H_t$
- Beliefs of environment state:  $S_t^a = (\mathbb{P}[S_t^e = s^1], \dots, \mathbb{P}[S_t^e = s^n])$
- Recurrent neural network:  $S_t^a = \sigma(S_{t-1}^a W_s + O_t W_o)$



## 2) Decision Making under uncertainty



Source: Cisco (2008).  
Cisco Health Presence  
Trial at Aberdeen Royal  
Infirmary in Scotland

3 July 1959, Volume 130, Number 3366

# SCIENCE

## Reasoning Foundations of Medical Diagnosis

Symbolic logic, probability, and value theory  
aid our understanding of how physicians reason.

Robert S. Ledley and Lee B. Lusted

The purpose of this article is to analyze the complicated reasoning processes inherent in medical diagnosis. The importance of this problem has received recent emphasis by the increasing interest in the use of electronic computers as an aid to medical diagnostic processes

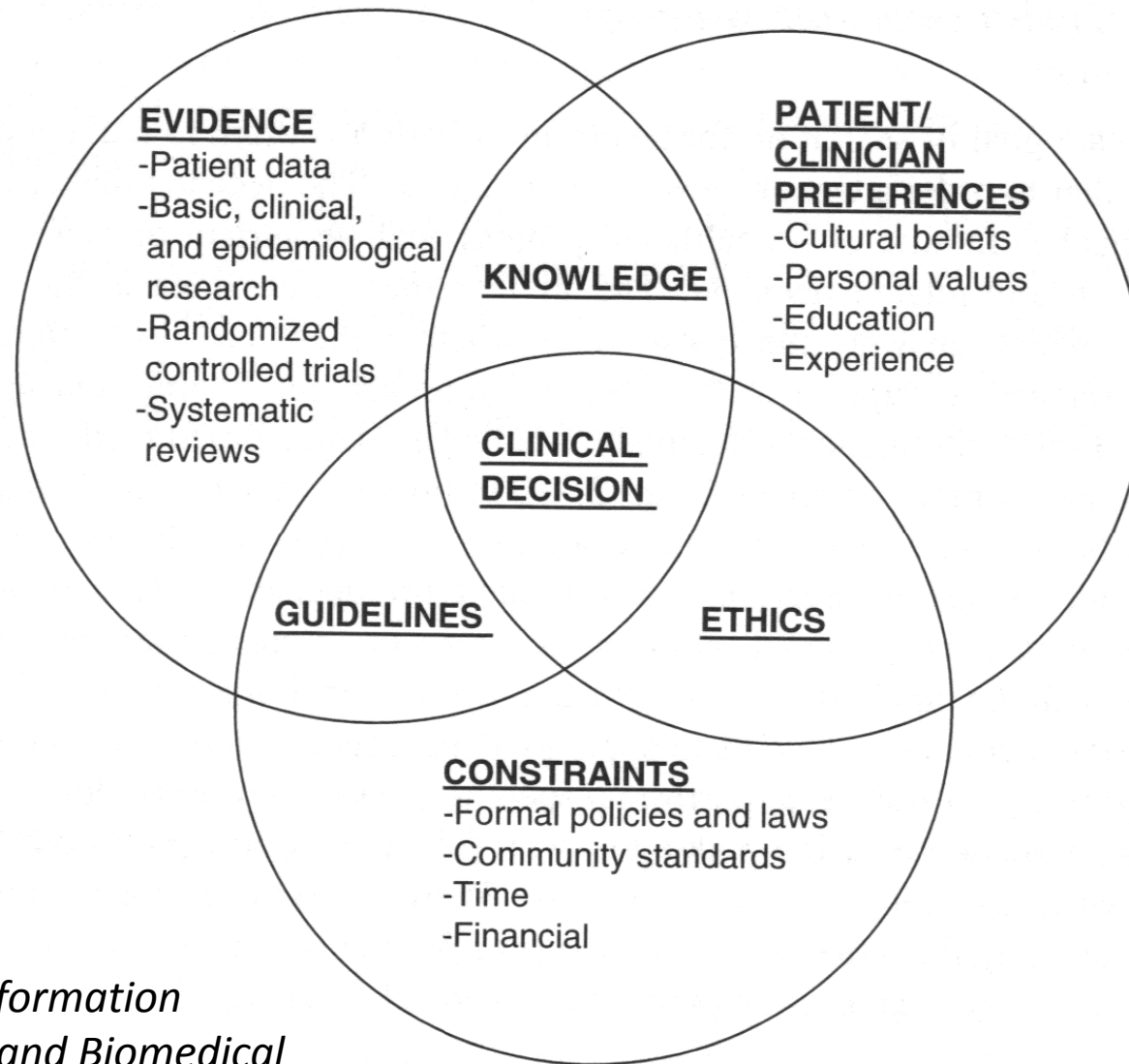
fitted into a definite disease category, or that it may be one of several possible diseases, or else that its exact nature cannot be determined." This, obviously, is a greatly simplified explanation of the process of diagnosis, for the physician might also comment that after seeing a

ance are the ones who do remember and consider the most possibilities."

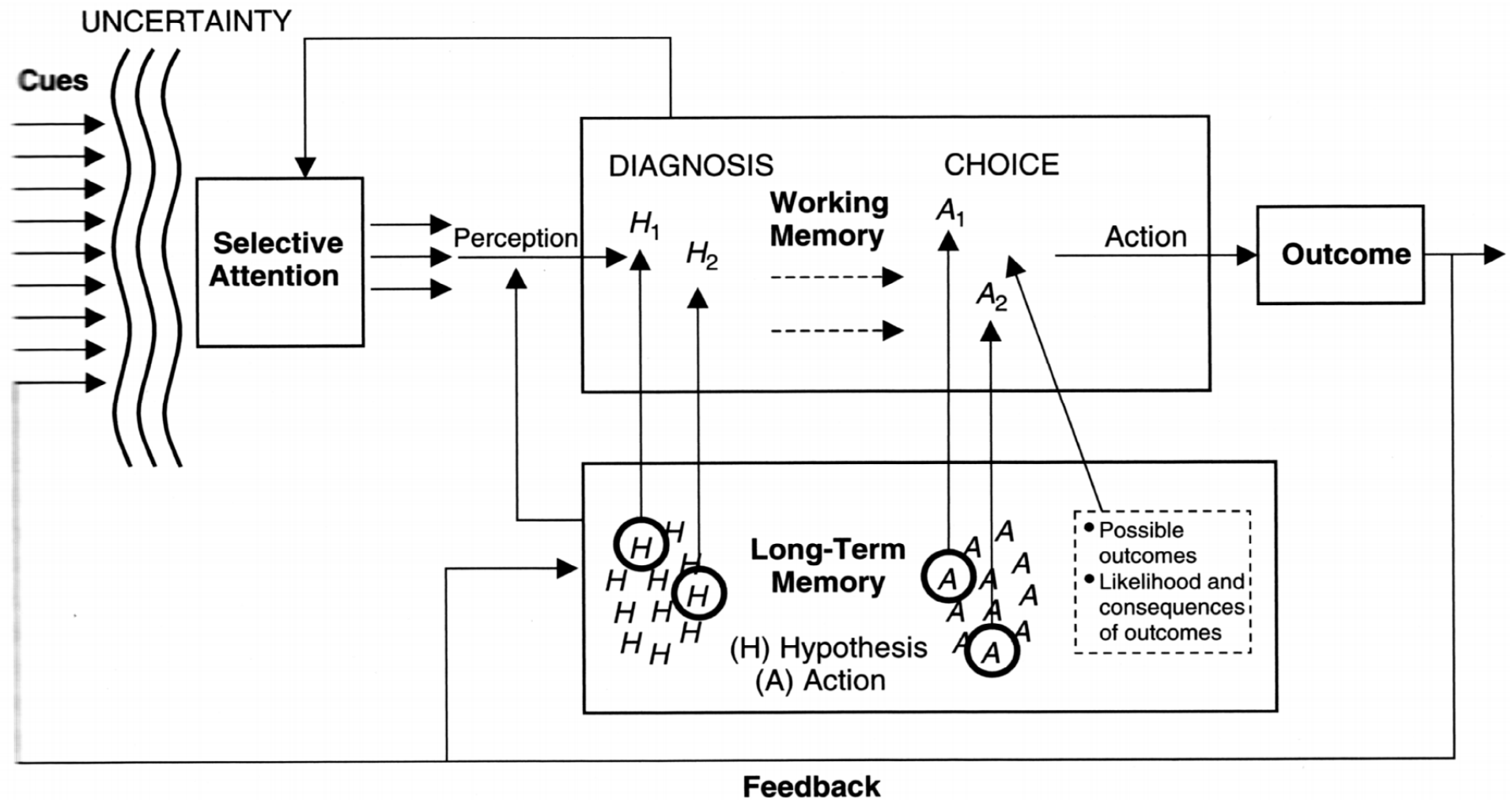
Computers are especially suited to help the physician collect and process clinical information and remind him of diagnoses which he may have overlooked. In many cases computers may be as simple as a set of hand-sorted cards, whereas in other cases the use of a large-scale digital electronic computer may be indicated. There are other ways in which computers may serve the physician, and some of these are suggested in this paper. For example, medical students might find the computer an important aid in learning the methods of differential diagnosis. But to use the computer thus we must understand how the physician makes a medical diagnosis. This, then, brings us to the subject of our investigation: the reasoning foundations of medical diagnosis and treatment.

Medical diagnosis involves processes that can be systematically analyzed, as well as those characterized as "intangible." For instance, the reasoning foundations of medical diagnostic procedures





Hersh, W. (2010) *Information Retrieval: A Health and Biomedical Perspective*. New York, Springer.



Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus (OH), Charles Merrill.

**Medical Action ...**

**is permanent decision making**







Stanford Heuristic Programming Project  
Memo HPP-78-1

February 1978

Computer Science Department  
Report No. STAN-CS-78-649

E. Feigenbaum, J. Lederberg, B. Buchanan, E. Shortliffe

Rheingold, H. (1985) *Tools for thought: the history and future of mind-expanding technology*. New York, Simon & Schuster.



DENDRAL AND META-DENDRAL:  
THEIR APPLICATIONS DIMENSION

by

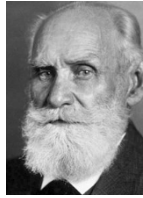
Bruce G. Buchanan and Edward A. Feigenbaum

COMPUTER SCIENCE DEPARTMENT  
School of Humanities and Sciences  
STANFORD UNIVERSITY



Buchanan, B. G. & Feigenbaum, E. A. (1978) DENDRAL and META-DENDRAL: their applications domain. *Artificial Intelligence*, 11, 1978, 5-24.

# 3) Roots of RL



Ivan P. Pavlov (1849-1936)  
1904 Nobel Prize  
Physiology/Medicine

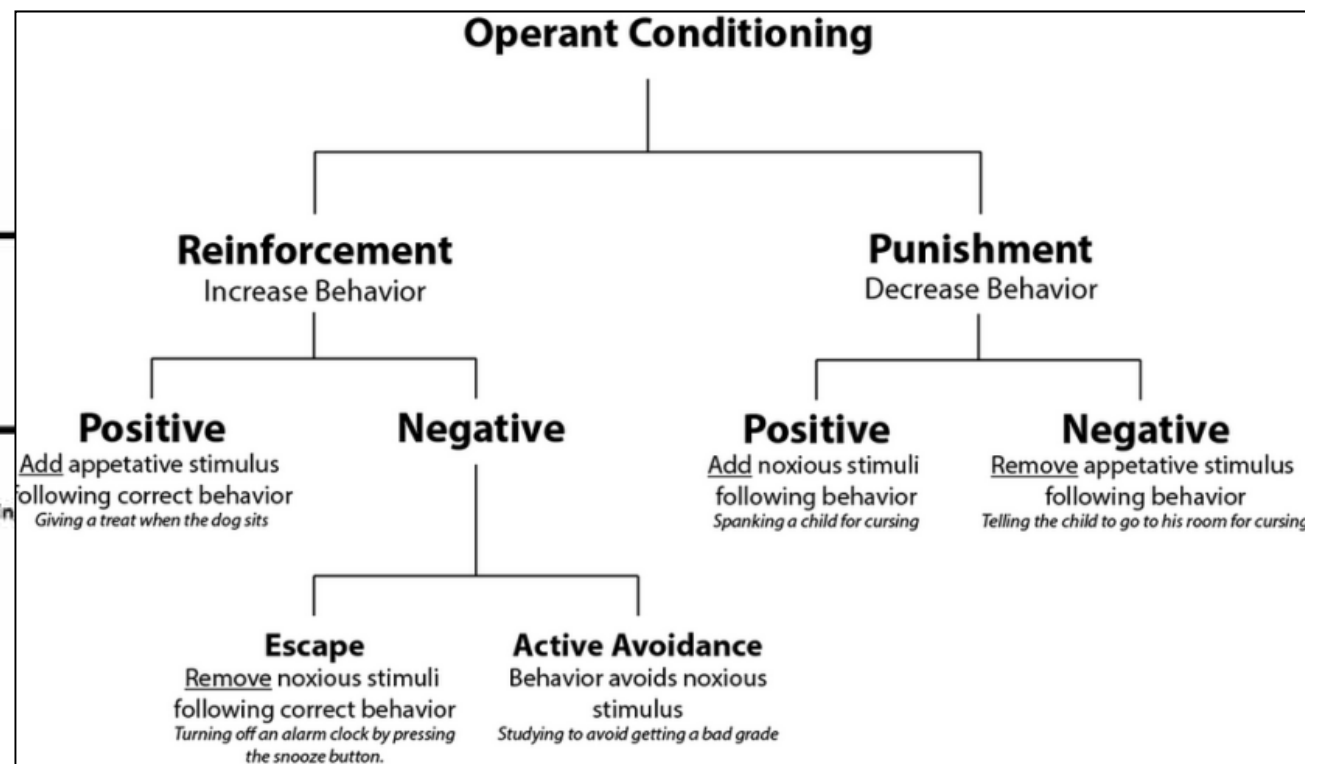


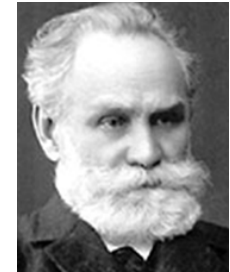
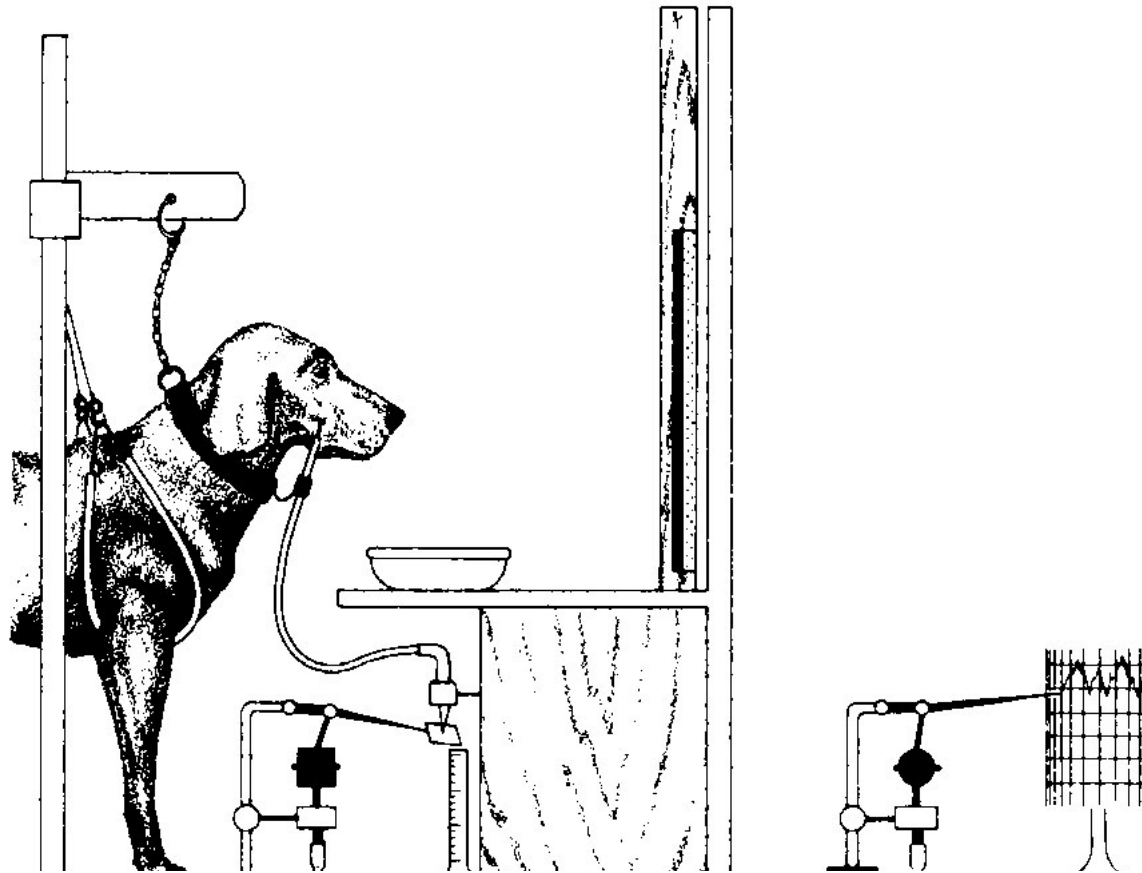
Edward L. Thorndike  
(1874-1949)  
1911 Law of Effect



Burrhus F. Skinner  
(1904-1990)  
1938 Operant Conditioning

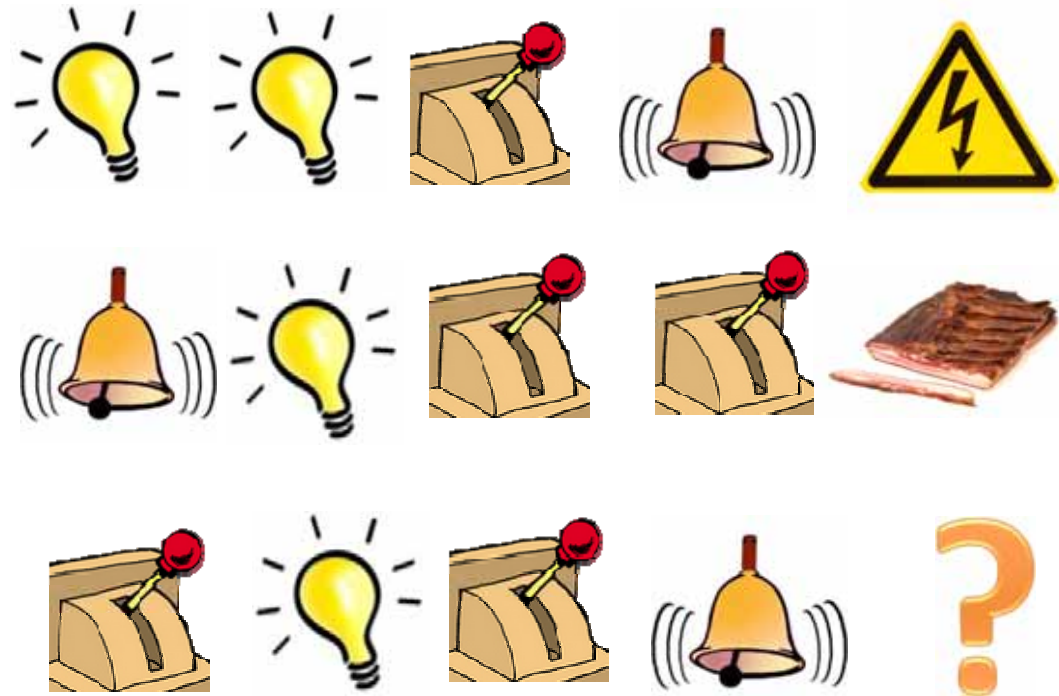
	Punishment (decreasing behavior)	Reinforcement (increasing behavior)
Positive (adding)	adding something to decrease behavior	adding something to increase behavior
Negative (subtracting)	subtracting something to decrease behavior	subtracting something to increase behavior





- *Classical (human and) animal conditioning*: “the magnitude and timing of the conditioned response changes as a result of the contingency between the conditioned stimulus and the unconditioned stimulus” [Pavlov, 1927].





- What if agent state = last 3 items in sequence?
- What if agent state = counts for lights, bells and levers?
- What if agent state = complete sequence?



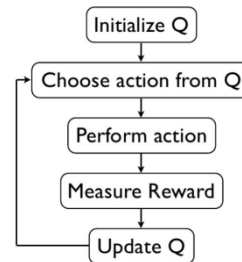
<https://webdocs.cs.ualberta.ca/~sutton/book/the-book.html>



Turing, A. M. 1950. Computing machinery and intelligence. Mind, 59, (236), 433-460.



Richard Bellman 1961. Adaptive control processes: a guided tour. Princeton.



Watkins, C. J. & Dayan, P. 1992. Q-learning. Machine learning, 8, (3-4), 279-292.



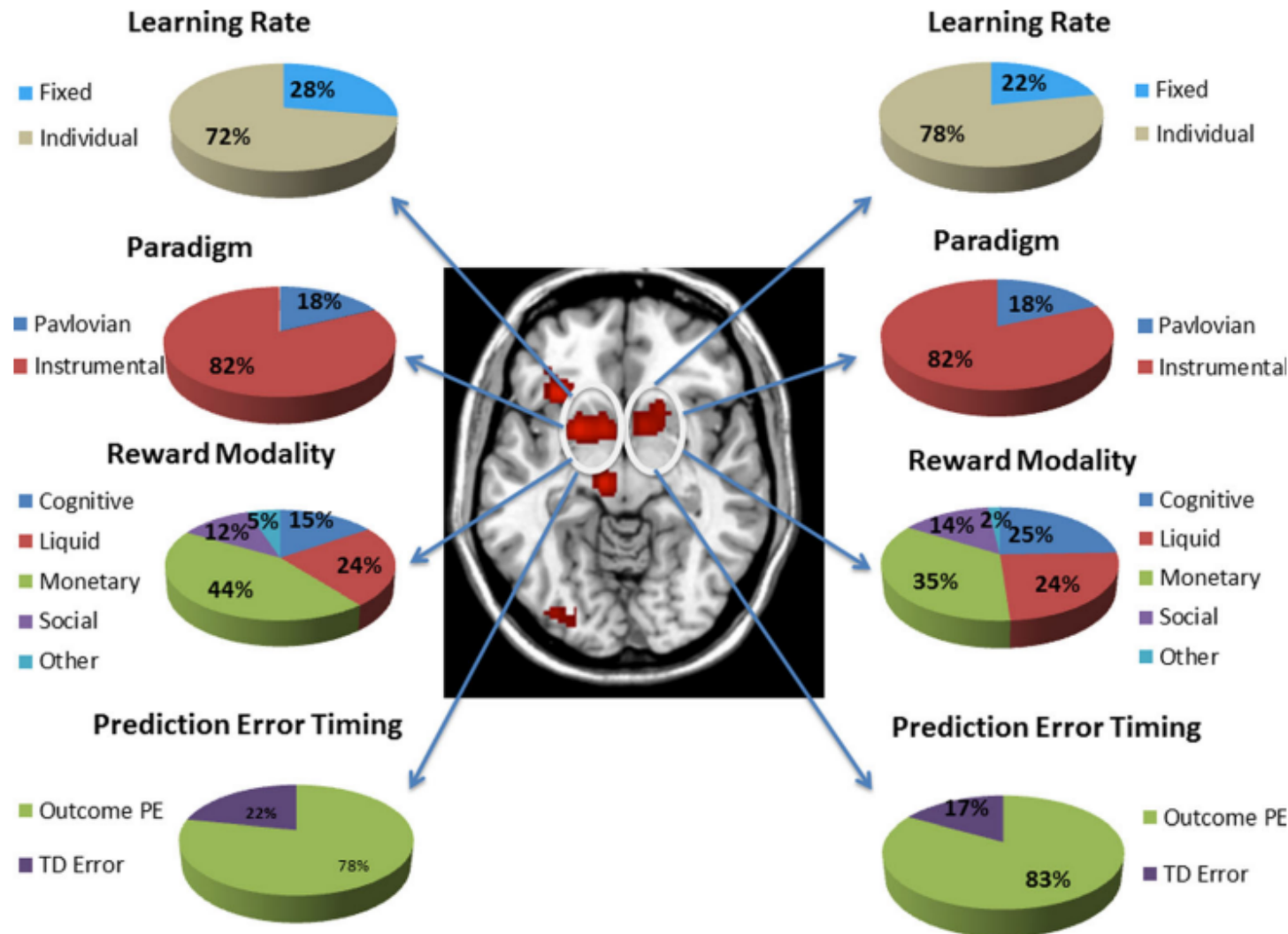
Sutton, R. S. & Barto, A. G. 1998. Reinforcement learning: An introduction, Cambridge, MIT press.



Littman, M. L. 2015. Reinforcement learning improves behaviour from evaluative feedback. Nature, 521, (7553), 445-451.

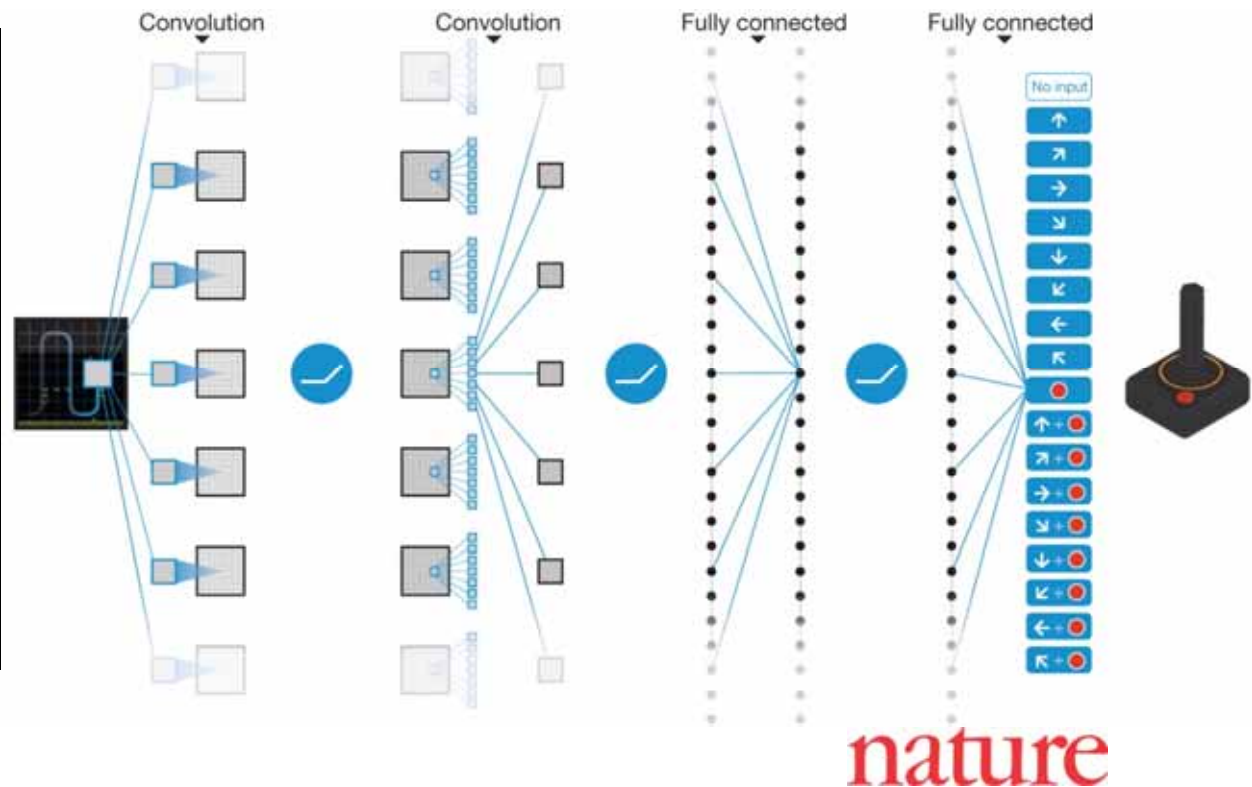
### Excellent Review Paper:

Kaelbling, L. P., Littman, M. L. & Moore, A. W. 1996. Reinforcement learning: A survey. Journal of Artificial Intelligence Research, 4, 237-285



Chase, H. W., Kumar, P., Eickhoff, S. B. & Dombrovski, A. Y. 2015. Reinforcement learning models and their neural correlates: An activation likelihood estimation meta-analysis. *Cognitive, Affective & Behavioral Neuroscience*, 15, (2), 435-459, doi:10.3758/s13415-015-0338-7.

Deep Q-networks (Q-Learning is a model-free RL approach) have successfully played Atari 2600 games at expert human levels



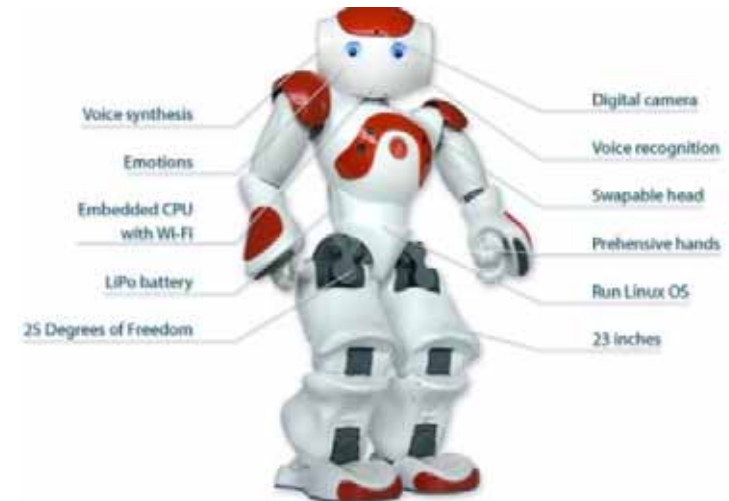
Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S. & Hassabis, D. 2015. Human-level control through deep reinforcement learning. *Nature*, 518, (7540), 529-533, doi:10.1038/nature14236



[http://images.computerhistory.org/timeline/timeline\\_ai\\_robotics\\_1939\\_elektro.jpg](http://images.computerhistory.org/timeline/timeline_ai_robotics_1939_elektro.jpg)



1985



<http://cyberneticzoo.com/robot-time-line/>







<http://www.neurotechnology.com/res/Robot2.jpg>



<https://royalsociety.org/events/2015/05/breakthrough-science-technologies-machine-learning>

Kober, J., Bagnell, J. A. & Peters, J. 2013. Reinforcement Learning in Robotics: A Survey. The International Journal of Robotics Research.



Nogrady, B. 2015. Q&A: Declan Murphy. Nature, 528, (7582), S132-S133, doi:10.1038/528S132a.

- 1943 McCulloch & Pitts – artificial neuron
- 1949 Hebb – learning as synapse modification
- 1957 Rosenblatt – Perceptron (SVM)
- 1969 Minsky & Papert – limits of perceptrons (begin of AI winter)
- 1974 Backpropagation – renaissance of NN
- 1982 Hopfield Network (recurrent NN)
- 1985 Boltzmann machine (stochastic RNN)
- 2010 Backpropagation (GPU's – deep learning is hype)



# 4) Cognitive Science of RL Human Information Processing



Salakhutdinov, R., Tenenbaum, J. & Torralba, A. 2012. One-shot learning with a hierarchical nonparametric Bayesian model. *Journal of Machine Learning Research*, 27, 195-207.

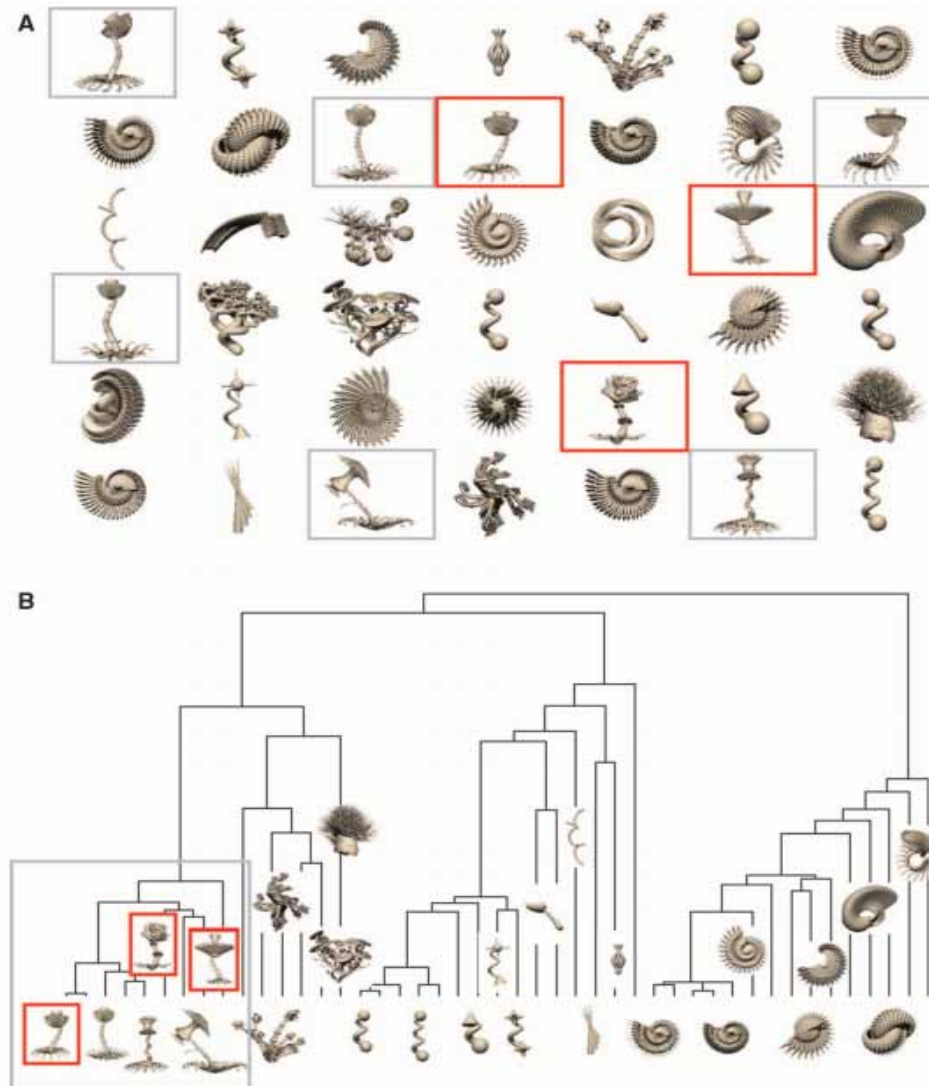
Quaxl

Quaxl

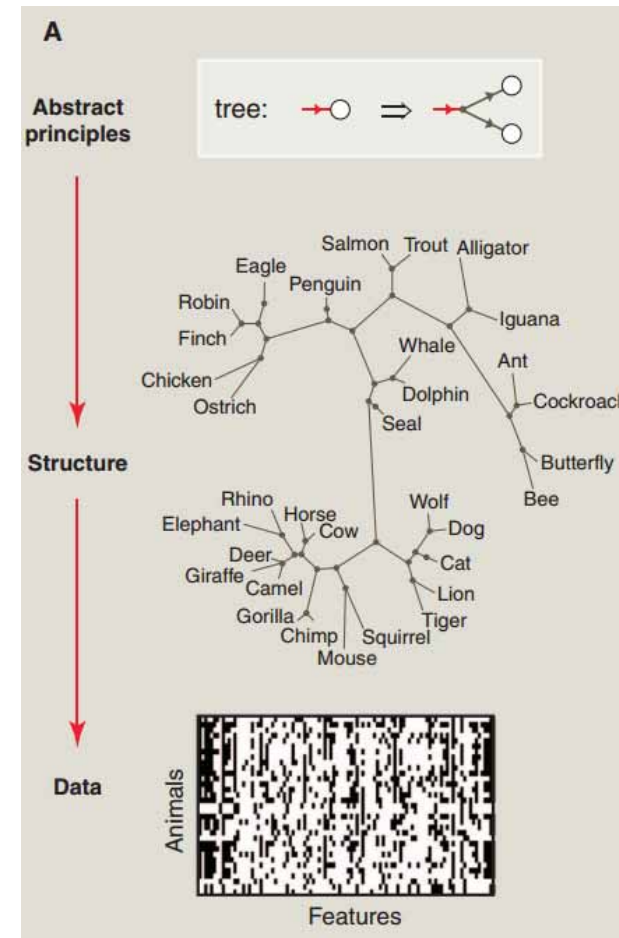


Quaxl

Salakhutdinov, R., Tenenbaum, J. & Torralba, A. 2012. One-shot learning with a hierarchical nonparametric Bayesian model. Journal of Machine Learning Research, 27, 195-207.



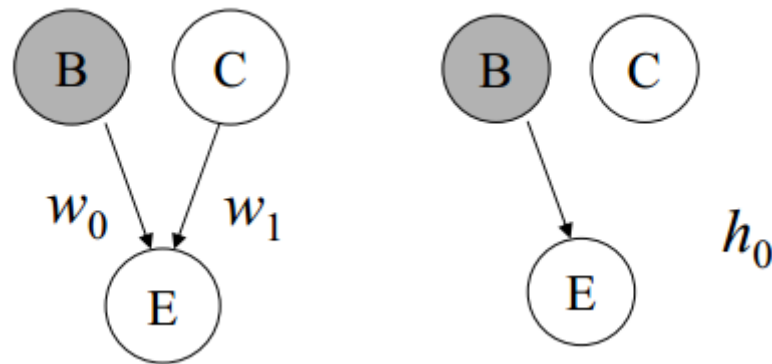
$$P(h|d) = \frac{P(d|h)P(h)}{\sum_{h' \in H} P(d|h')P(h')} \propto P(d|h)P(h)$$



Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. 2011. How to grow a mind: Statistics, structure, and abstraction. *Science*, 331, (6022), 1279-1285.

- which is highly relevant for ML research, concerns the factors that determine the subjective difficulty of concepts:
- Why are some concepts psychologically extremely simple and easy to learn,
- while others seem to be extremely difficult, complex, or even incoherent?
- These questions have been studied since the 1960s but are still unanswered ...

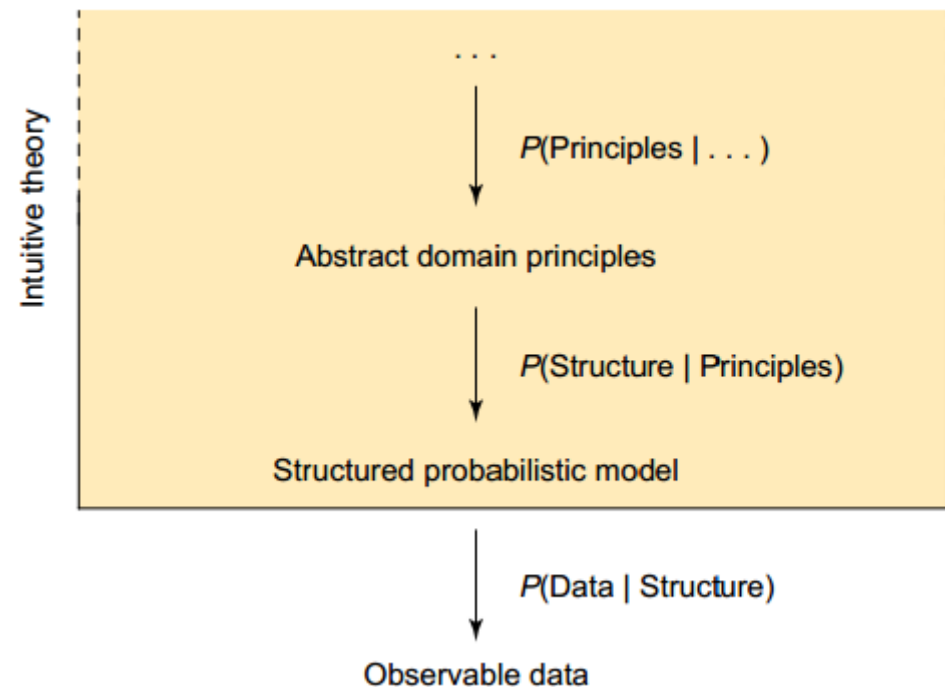
Feldman, J. 2000. Minimization of Boolean complexity in human concept learning. *Nature*, 407, (6804), 630-633, doi:10.1038/35036586.



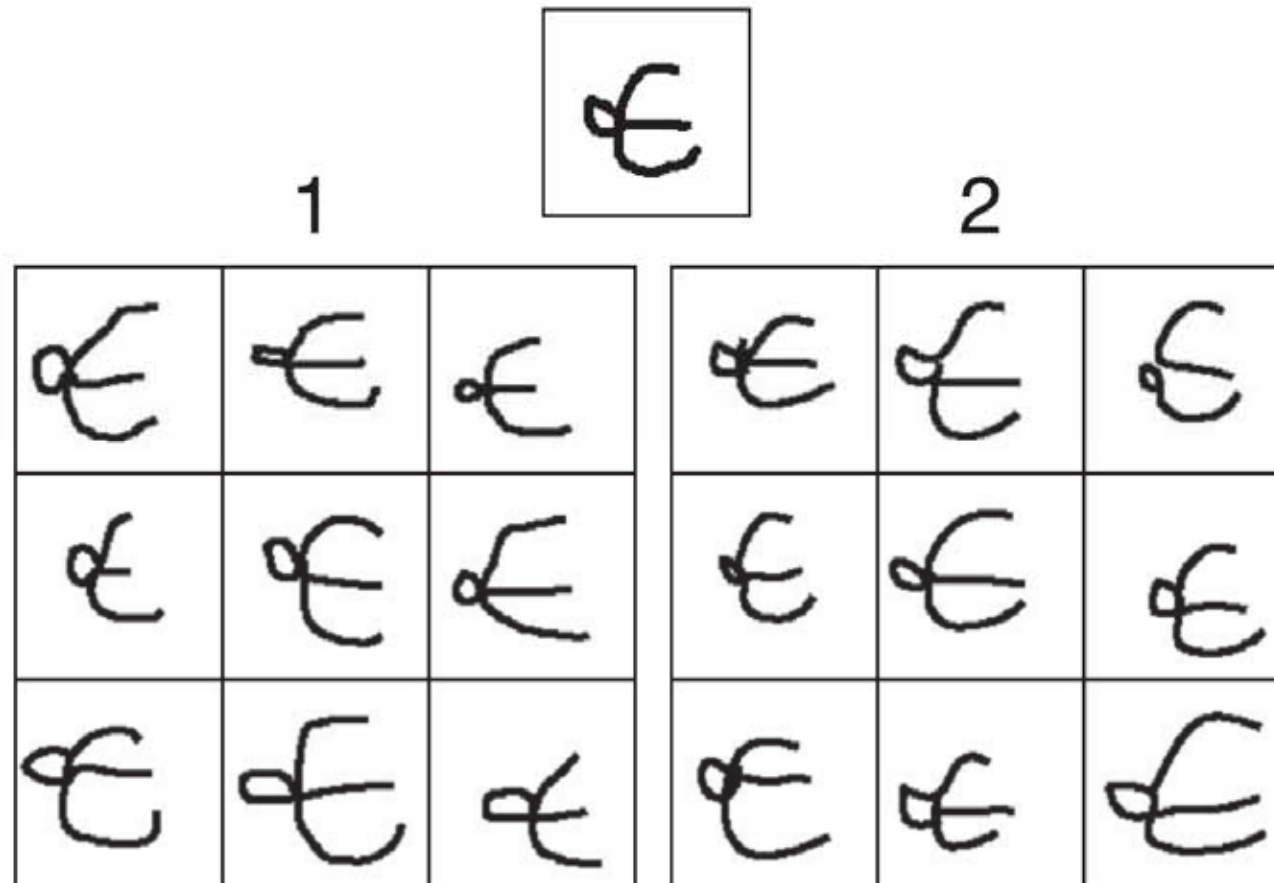
- Cognition as probabilistic inference
  - Visual perception, language acquisition, motor learning, associative learning, memory, attention, categorization, reasoning, causal inference, decision making, theory of mind
- Learning concepts from examples
- Learning and applying intuitive theories (balancing complexity vs. fit)

- Similarity
- Representativeness and evidential support
- Causal judgement
- Coincidences and causal discovery
- Diagnostic inference
- Predicting the future

Tenenbaum, J. B., Griffiths, T. L. & Kemp, C. 2006. Theory-based Bayesian models of inductive learning and reasoning. Trends in cognitive sciences, 10, (7), 309-318.



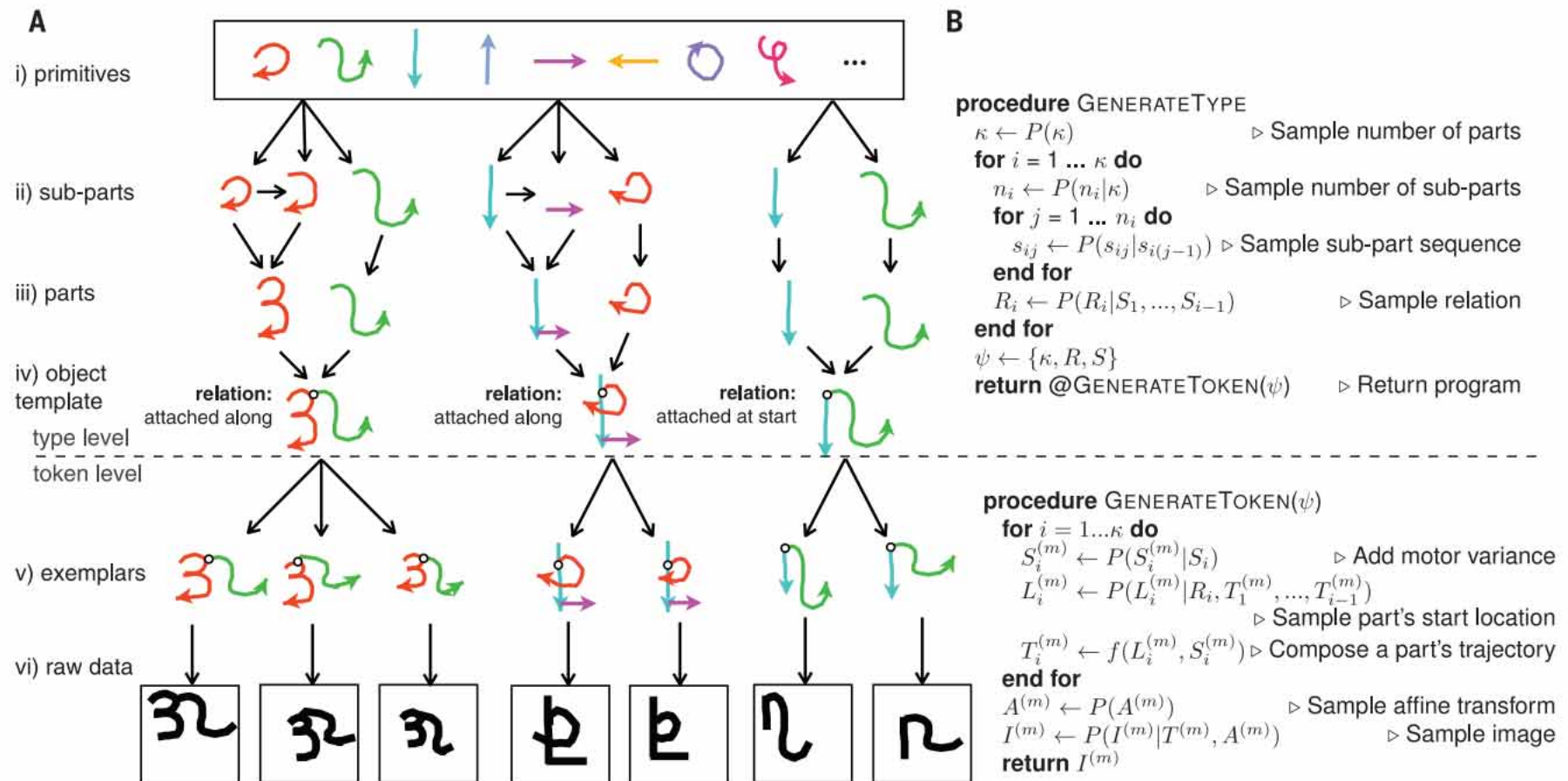




Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. 2015. Human-level concept learning through probabilistic program induction. *Science*, 350, (6266), 1332-1338, doi:10.1126/science.aab3050.

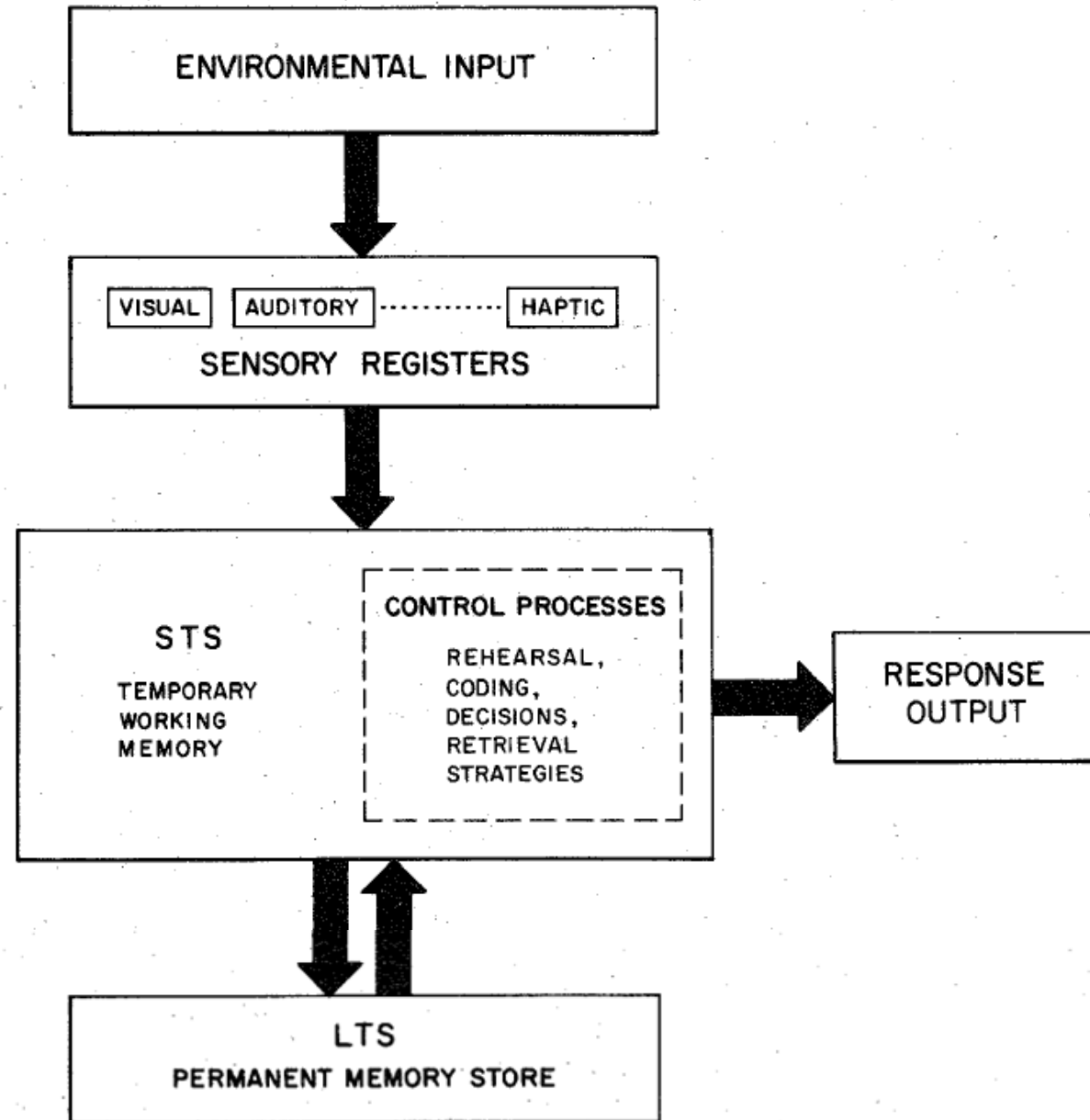


A Bayesian program learning (BPL) framework, capable of learning a large class of visual concepts from just a single example and generalizing in ways that are mostly indistinguishable from people

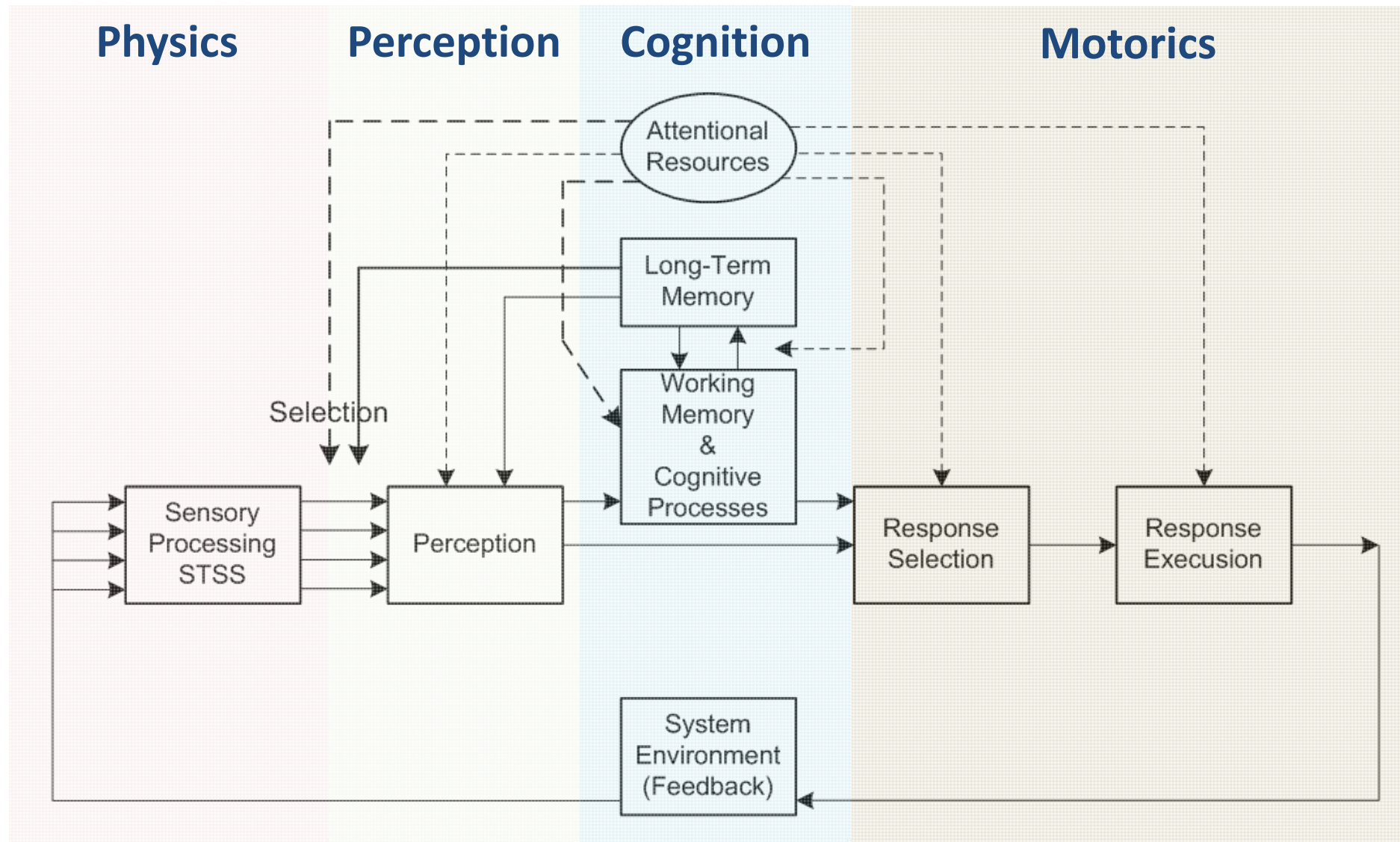


Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. 2015. Human-level concept learning through probabilistic program induction. Science, 350, (6266), 1332-1338, doi:10.1126/science.aab3050.

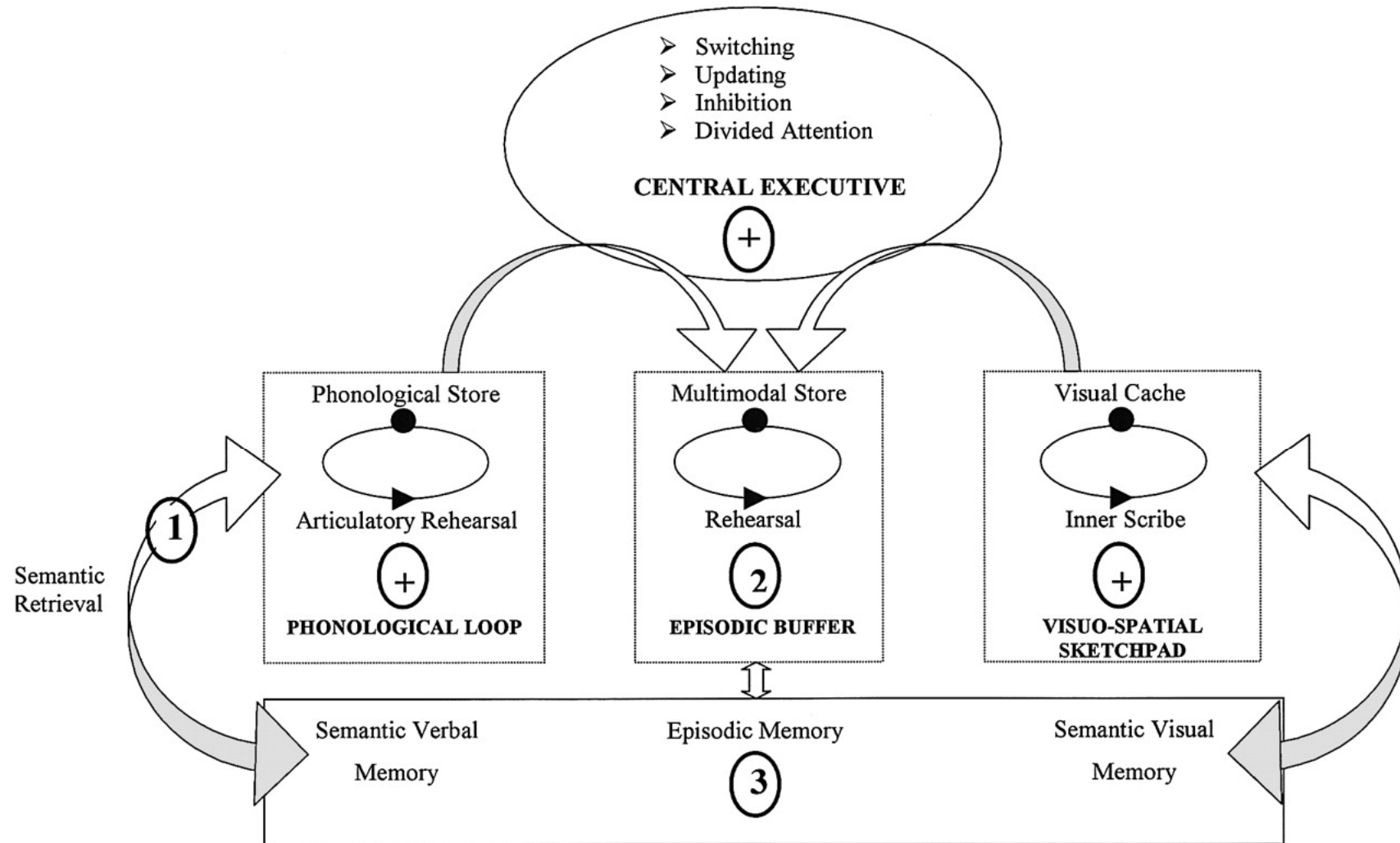
# How does our mind get so much out of so little?



Atkinson, R. C. & Shiffrin, R. M. (1971) *The control processes of short-term memory* (Technical Report 173, April 19, 1971). Stanford, Institute for Mathematical Studies in the Social Sciences, Stanford University.



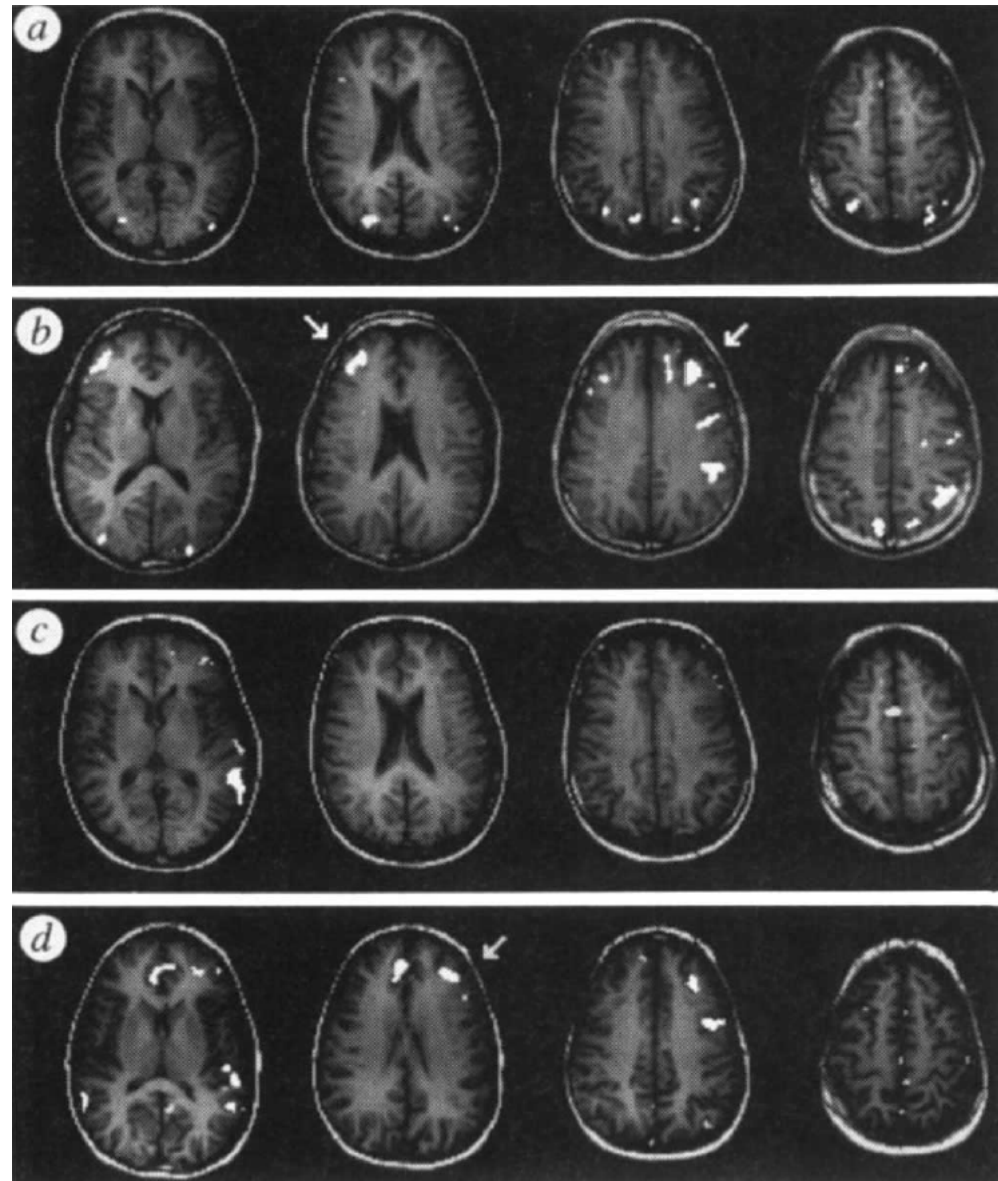
Wickens, C., Lee, J., Liu, Y. & Gordon-Becker, S. (2004) *Introduction to Human Factors Engineering: Second Edition*. Upper Saddle River (NJ), Prentice-Hall.

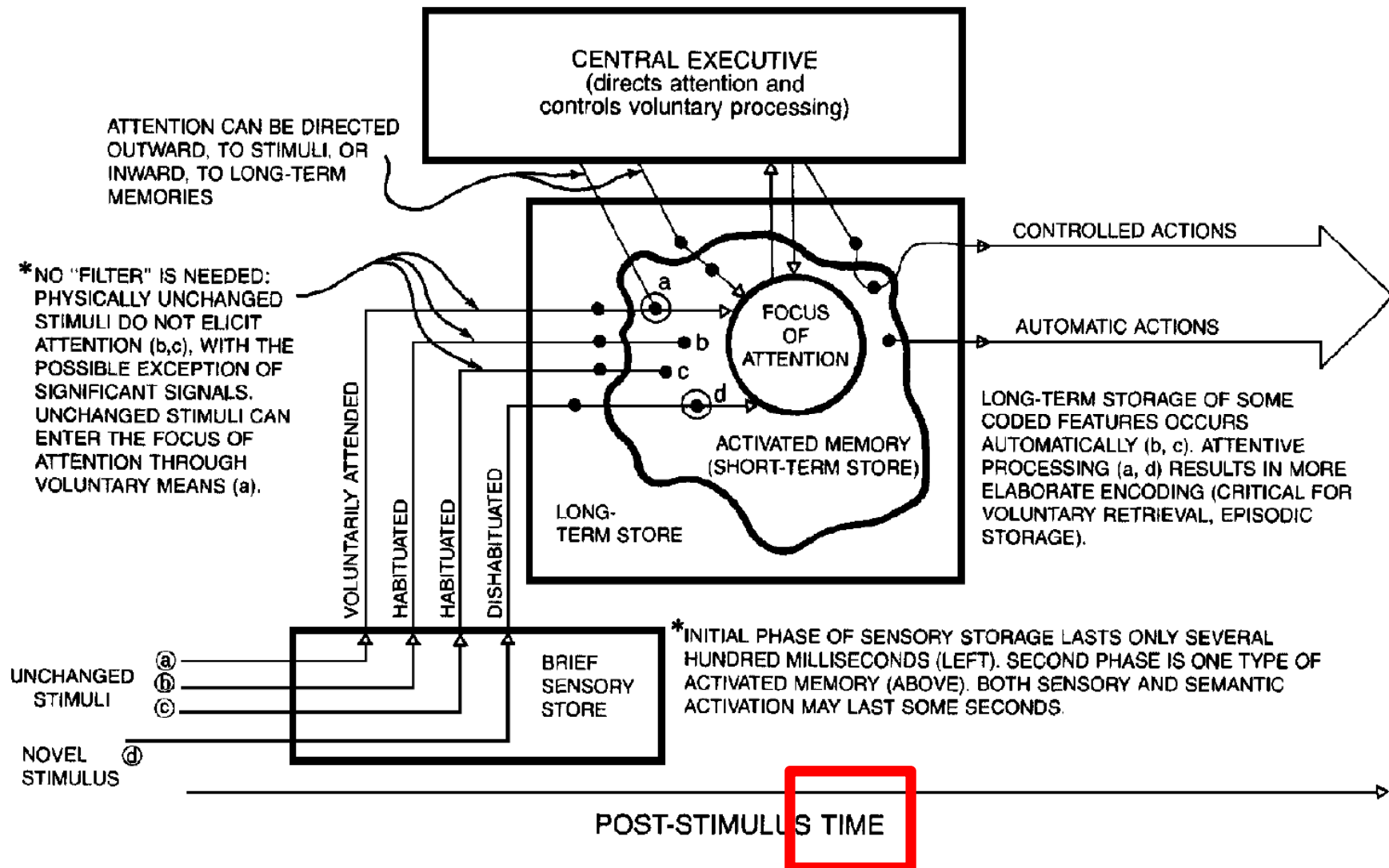


Quinette, P., Guillery, B., Desgranges, B., de la Sayette, V., Viader, F. & Eustache, F. (2003)  
Working memory and executive functions in transient global amnesia. *Brain*, 126, 9, 1917-1934.



D'Esposito, M., Detre, J. A., Alsop, D. C., Shin, R. K., Atlas, S. & Grossman, M. (1995) The neural basis of the central executive system of working memory. *Nature*, 378, 654, 279-281.





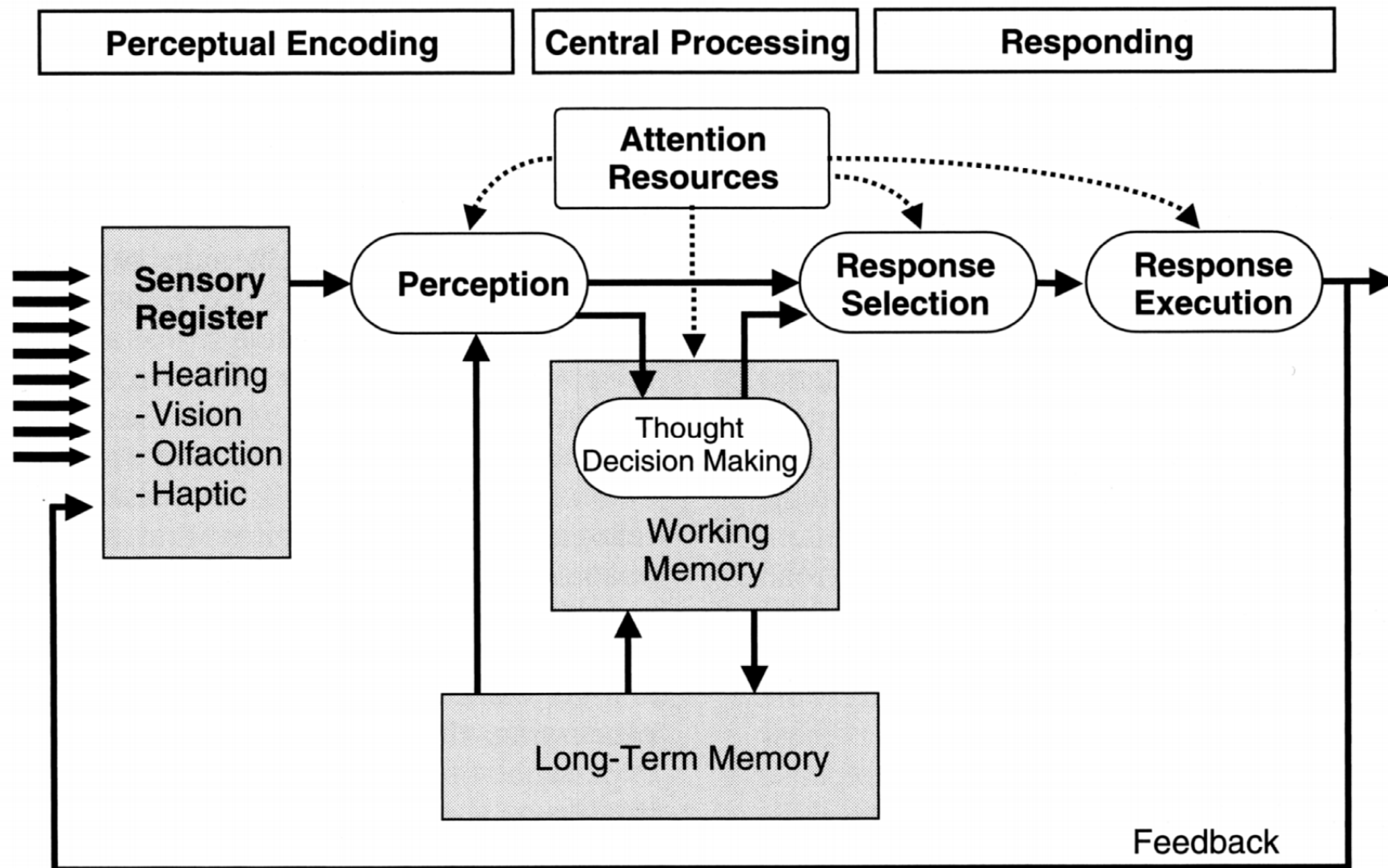
Cowan, N. (1988) Evolving conceptions of memory storage, selective attention, and their mutual constraints within the human information-processing system. *Psychological Bulletin*, 104, 2, 163.



Note: The Test does NOT properly work if you know it in advance or if you do not concentrate on counting

Simons, D. J. & Chabris, C. F. 1999. Gorillas in our midst: sustained inattention blindness for dynamic events. *Perception*, 28, (9), 1059-1074.

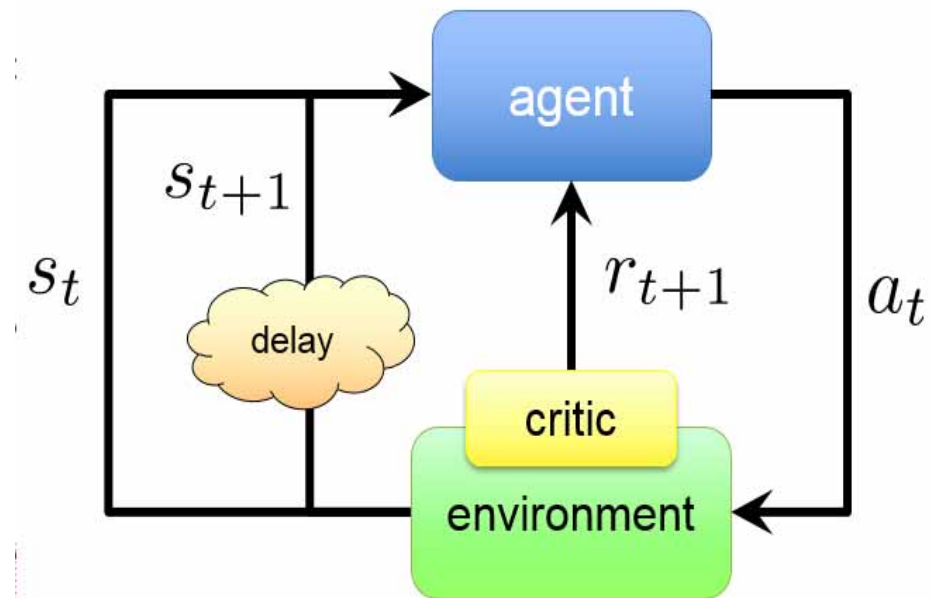




Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus (OH), Charles Merrill.

# 5) The Anatomy of an RL Agent

- Decision-making under uncertainty
- Limited knowledge of the domain environment
- Unknown outcome – unknown reward
- Partial or unreliable access to “databases of interaction”



Russell, S. J. & Norvig, P. 2009. Artificial intelligence: a modern approach (3rd edition), Prentice Hall, Chapter 16, 17: Making Simple Decisions and **Making Complex Decisions**

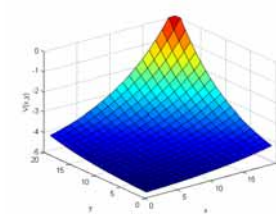
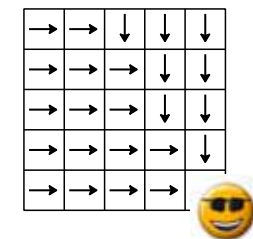
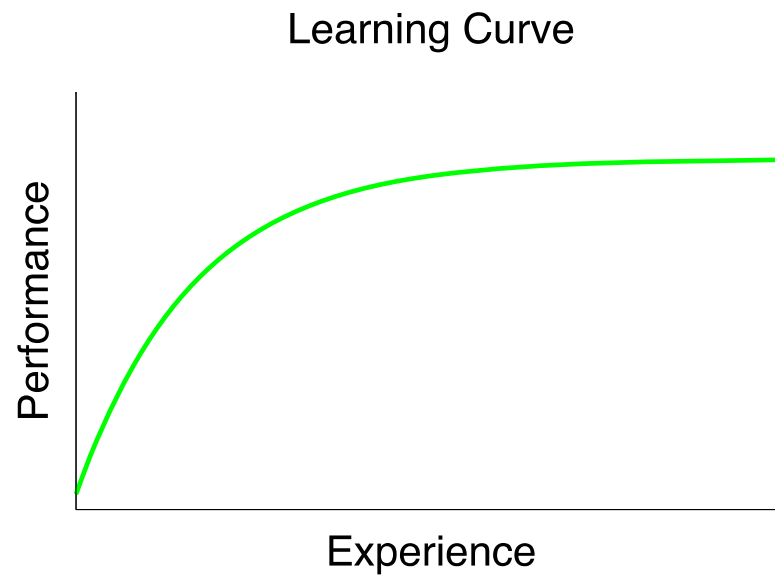
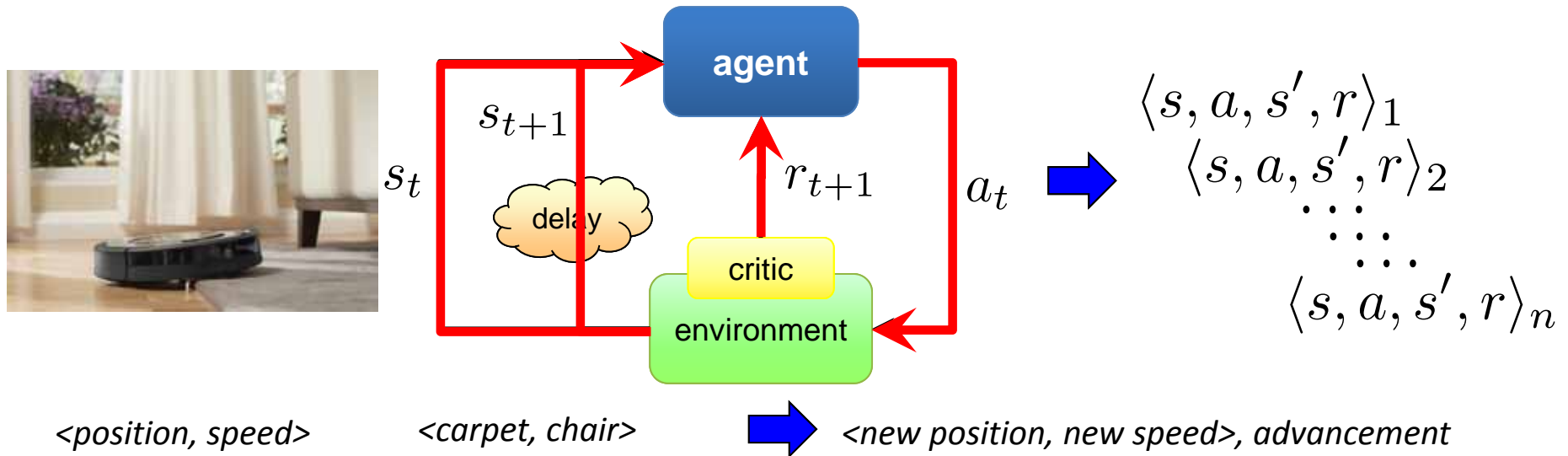


Image credit to Alessandro Lazaric

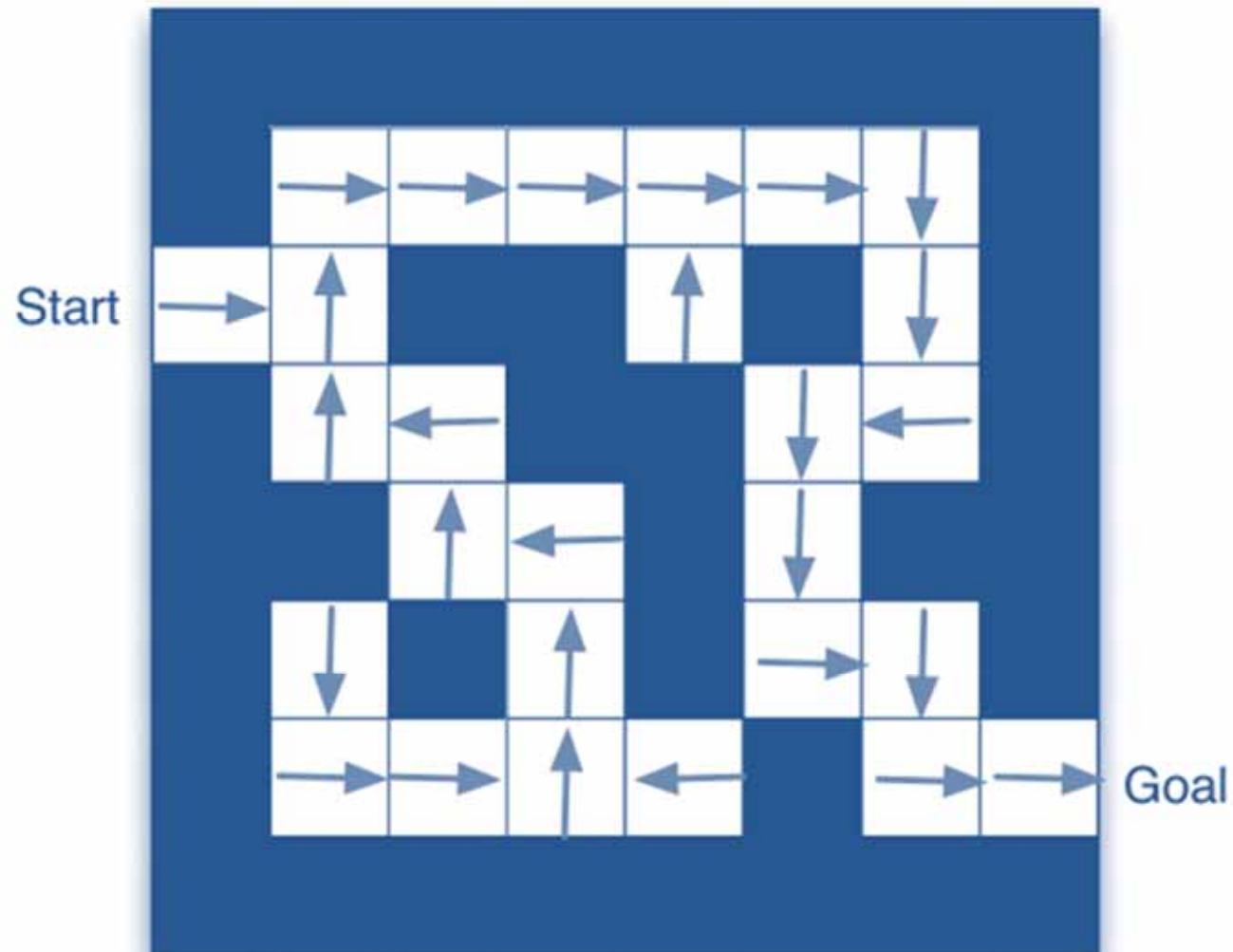
- **Policy:** agent's behaviour function  
e.g. stochastic policy  $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$
- **Value function:** how good is each state and/or action  
e.g.  $v_\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$
- **Model:** agent's representation of the environment  
 $\mathcal{P}$  predicts the next state;  $\mathcal{R}$  the next reward

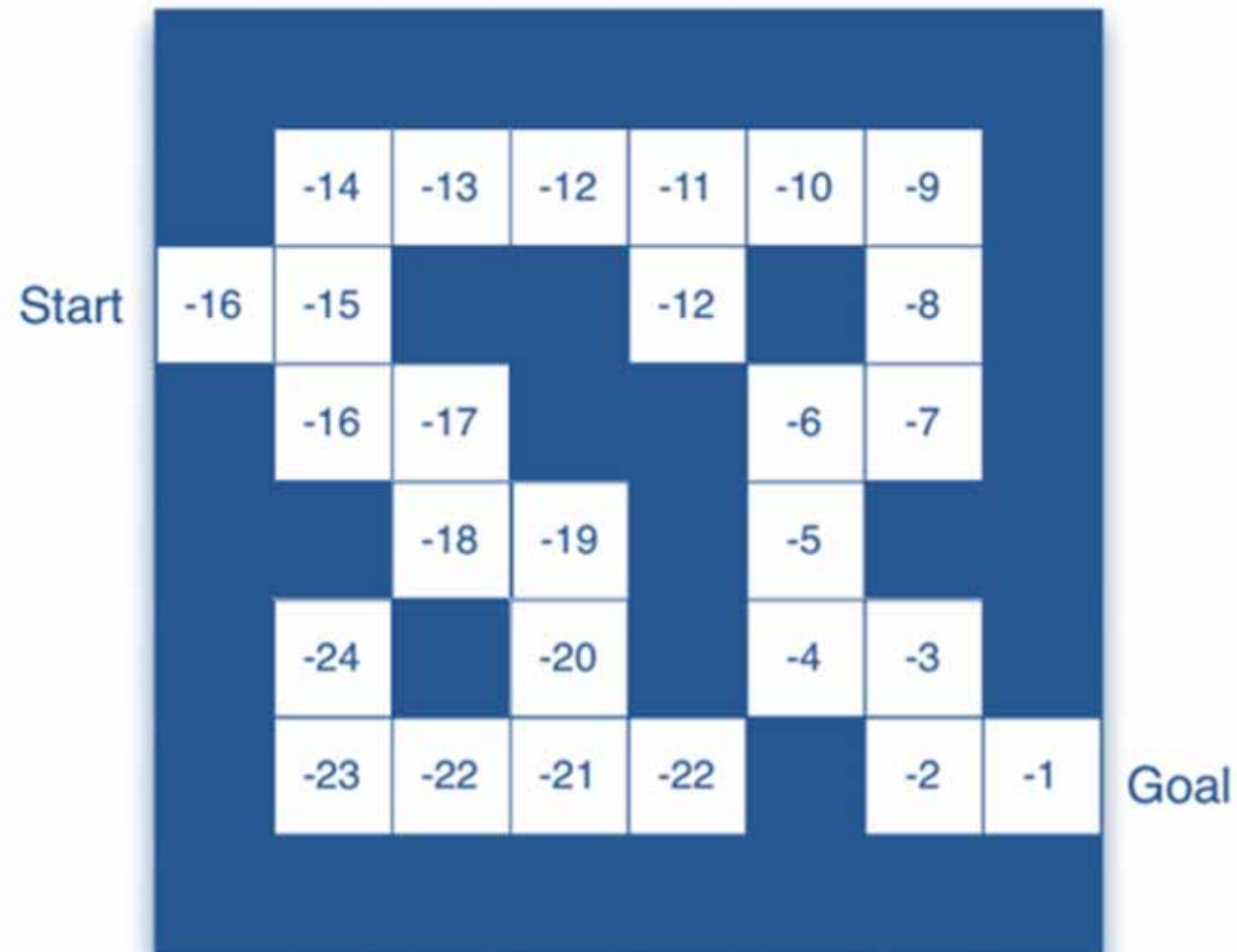
$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

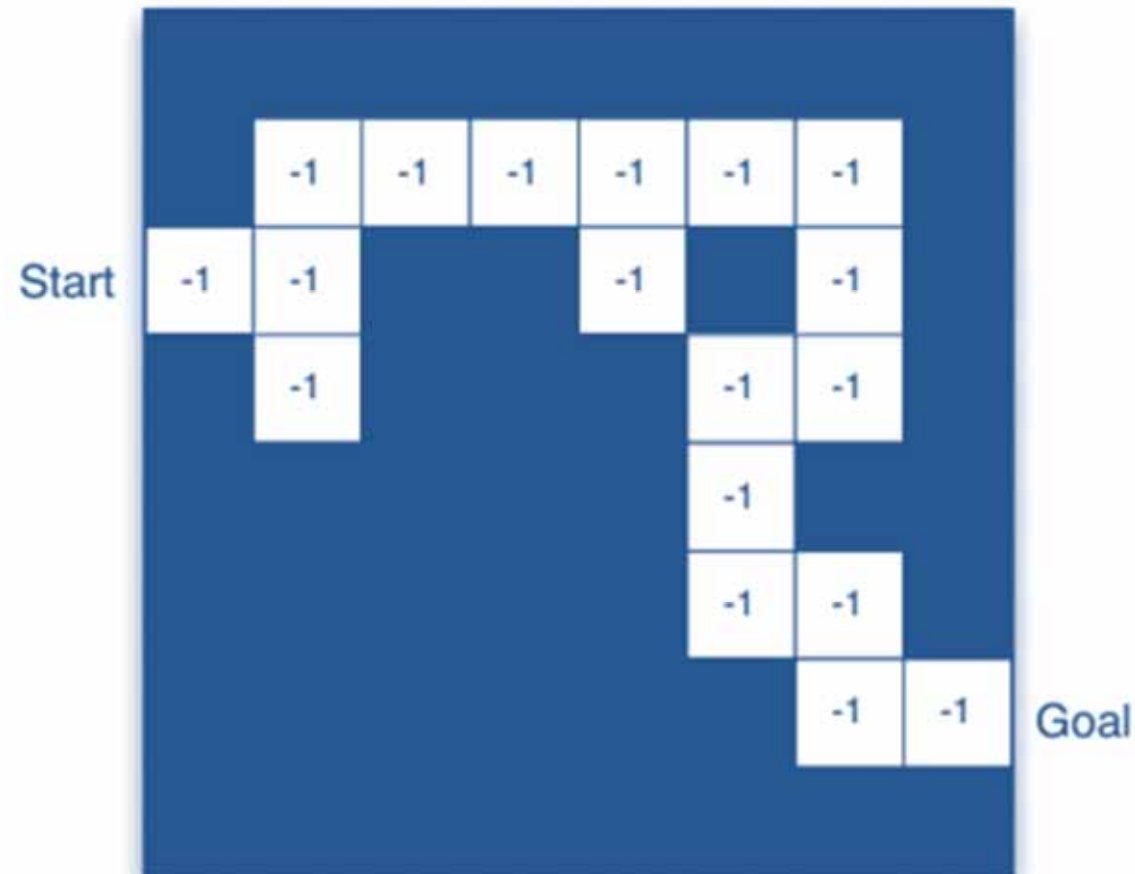
$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

- 1) Value-Based  
(no policy, only value function)
- 2) Policy-Based  
(no value function, only policy)
- 3) Actor-Critic  
(both)
- 4) Model free  
(and/or) – but no model
- 5) Model-based  
(and/or – and model)









- Grid layout represents transition model  $\mathcal{P}_{ss'}^a$
- Numbers represent immediate reward  $\mathcal{R}_s^a$  from each state  $s$  (same for all  $a$ )

Time steps  $t_1, t_2, \dots, t_n$

- Observe the state  $x_t$
- Take an action  $a_t$  (problem of **exploration** and **exploitation**)
- Observe next state and earn reward  $x_{t+1}, r_t$
- Update the policy and the value function  $\pi_t, Q_t$

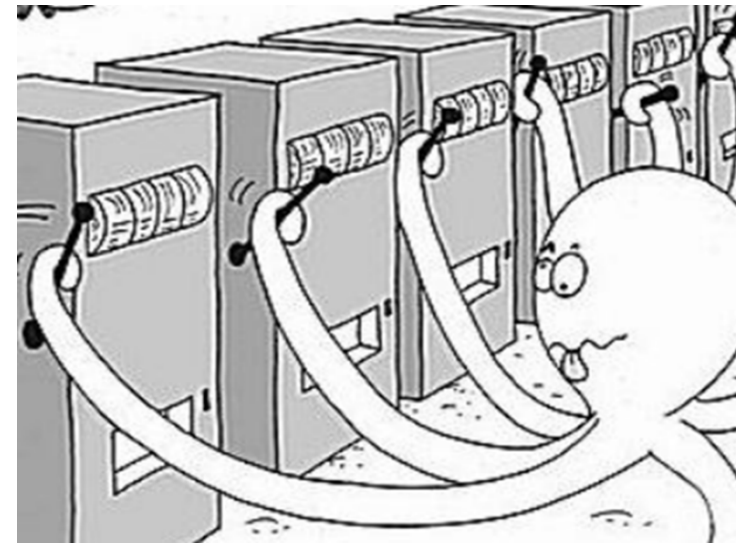
$$Q(x_t, a_t) = Q(x_t, a_t) + \alpha(r_t + \gamma \max_a Q(x_{t+1}, a) - Q(x_t, a_t))$$

$$\pi(x) = \arg \max_a Q(x, a)$$

- Temporal difference learning (1988)
- Q-learning (1998)
- BayesRL (2002)
- RMAX (2002)
- CBPI (2002)
- PEGASUS (2002)
- Least-Squares Policy Iteration (2003)
- Fitted Q-Iteration (2005)
- GTD (2009)
- UCRL (2010)
- REPS (2010)
- DQN (2014)

# 6) Example: Multi-Armed Bandits (MAB)





- There are  $n$  slot-machines (“einarmige Banditen”)
- Each machine  $i$  returns a reward  $y \approx P(y; \Theta_i)$
- Challenge: The machine parameter  $\Theta_i$  is unknown
- Which arm of a -slot machine should a gambler pull to maximize his cumulative reward over a sequence of trials? (stochastic setting or adversarial setting)

Image credit and more information: <http://research.microsoft.com/en-us/projects/bandits>

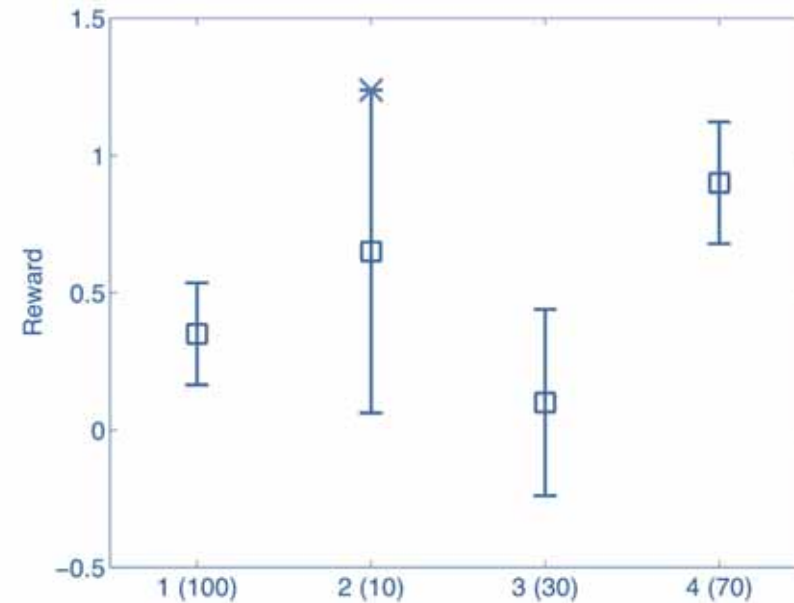
- Let  $a_t \in \{1, \dots, n\}$  be the choice of a machine at time  $t$
- Let  $y_t \in \mathbb{R}$  be the outcome with a mean of  $\langle y_{at} \rangle$
- Now, the given policy maps all history to a new choice:

$$\pi : [(a_1, y_1), (a_2, y_2), \dots, (a_{t-1}, y_{t-1})] \mapsto a_t$$

- The problem: Find a policy  $\pi$  that  $\max \langle y_T \rangle$
- Now, two effects appear when choosing such machine:
  - You collect more data about the machine (=knowledge)
  - You collect reward
- Exploration and Exploitation
  - **Exploration:** Choose the next action  $a_t$  to  $\min \langle H(b_t) \rangle$
  - **Exploitation:** Choose the next action  $a_t$  to  $\max \langle y_t \rangle$
- models an agent that simultaneously attempts to acquire new knowledge (called "exploration") and optimize his or her decisions based on existing knowledge (called "exploitation"). The agent attempts to balance these competing tasks in order to maximize total value over the period of time considered.

More information: <http://research.microsoft.com/en-us/projects/bandits>

$$a_t = \max_{a \in \mathcal{A}} \left( \hat{r}_t(a) + \sqrt{\frac{\log(1/\delta)}{T_t(a)}} \right)$$



$$a_t = \max_{a \in \mathcal{A}} (\text{rew}_t(a) + \text{uncert}_t(a))$$

**Exploitation**

*the higher the (estimated)  
reward the higher the chance  
to select the action*

**Exploration**

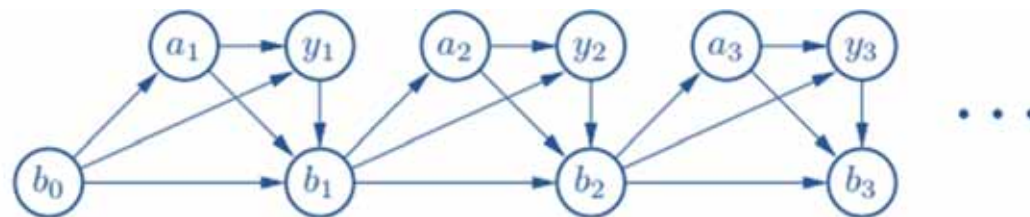
**the higher the (theoretical)  
uncertainty the higher the  
chance to select the action**

Auer, P., Cesa-Bianchi, N. & Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. Machine learning, 47, (2-3), 235-256.

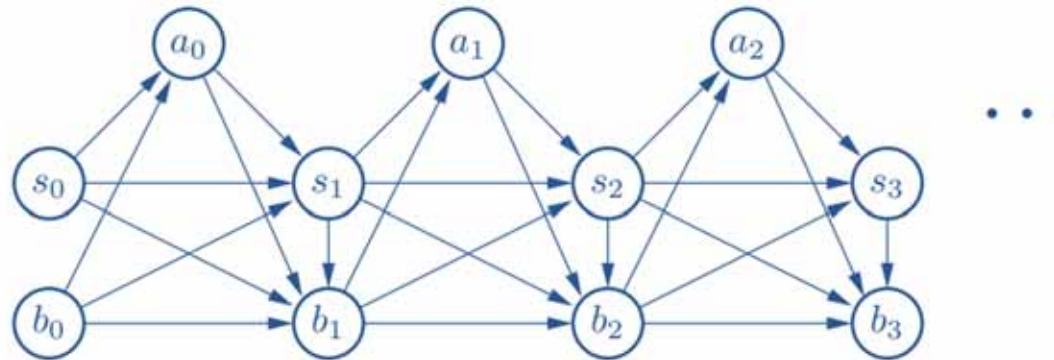
- Knowledge can be represented in two ways:
- 1) as full history  $h_t = [(a_1, y_1), (a_2, y_2), \dots, (a_{t-1}, y_{t-1})]$   
or
- 2) as belief  $b_t(\theta) = P(\theta|h_t)$

where  $\theta$  are the unknown parameters of all machines

The process can be modelled as belief MDP:



$$P(b'|y, a, b) = \begin{cases} 1 & \text{if } b' = b'_{[b,a,y]} \\ 0 & \text{otherwise} \end{cases}, \quad P(y|a, b) = \int_{\theta_a} b(\theta_a) P(y|\theta_a)$$



$$P(b'|s', s, a, b) = \begin{cases} 1 & \text{if } b' = b[s', s, a] \\ 0 & \text{otherwise} \end{cases}, \quad P(s'|s, a, b) = \int_{\theta} b(\theta) P(s'|s, a, \theta)$$

$$V(b, s) = \max_a \left[ \mathbb{E}(r|s, a, b) + \sum_{s'} P(s'|a, s, b) V(s', b') \right]$$

Poupart, P., Vlassis, N., Hoey, J. & Regan, K. An analytic solution to discrete Bayesian reinforcement learning. Proceedings of the 23rd international conference on Machine learning, 2006. ACM, 697-704.

- Clinical trials: potential treatments for a disease to select from new patients or patient category at each round, see:

W. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. Bulletin of the American Mathematics Society, vol. 25, pp. 285–294, 1933.

- Games: Different moves at each round, e.g. GO
- Adaptive routing: finding alternative paths, also finding alternative roads for driving from A to B
- Advertisement placements: selection of an ad to display at the Webpage out of a finite set which can vary over time, for each new Web page visitor



# 7) Applications in Health

YouTube DE health robots



Basically, this robot can move people from a bed to a wheelchair or a wheelchair to a bed, and

Top 9 Medical Robots That Could Change Healthcare

Robots movie

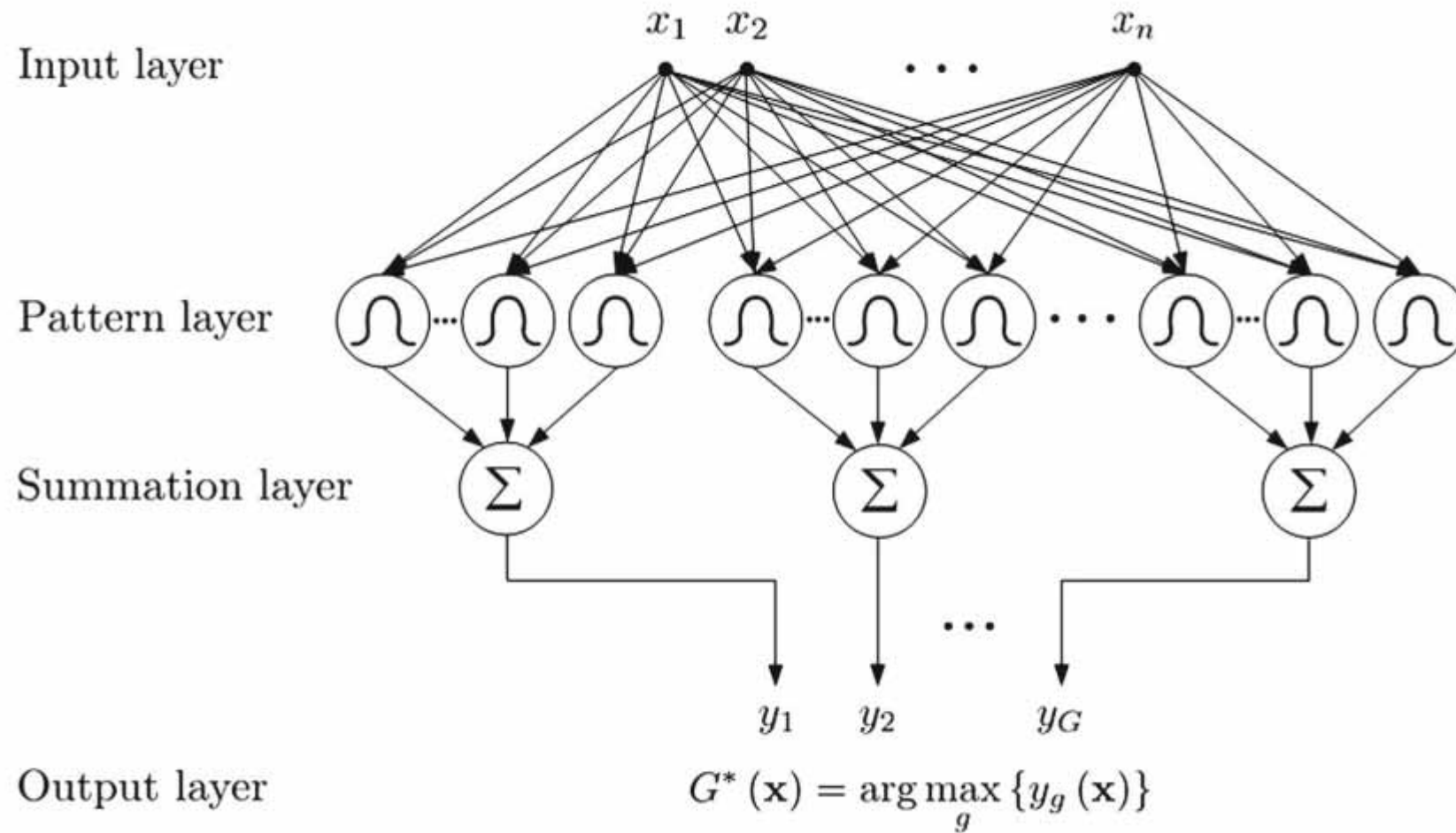
Subscribe

10,653

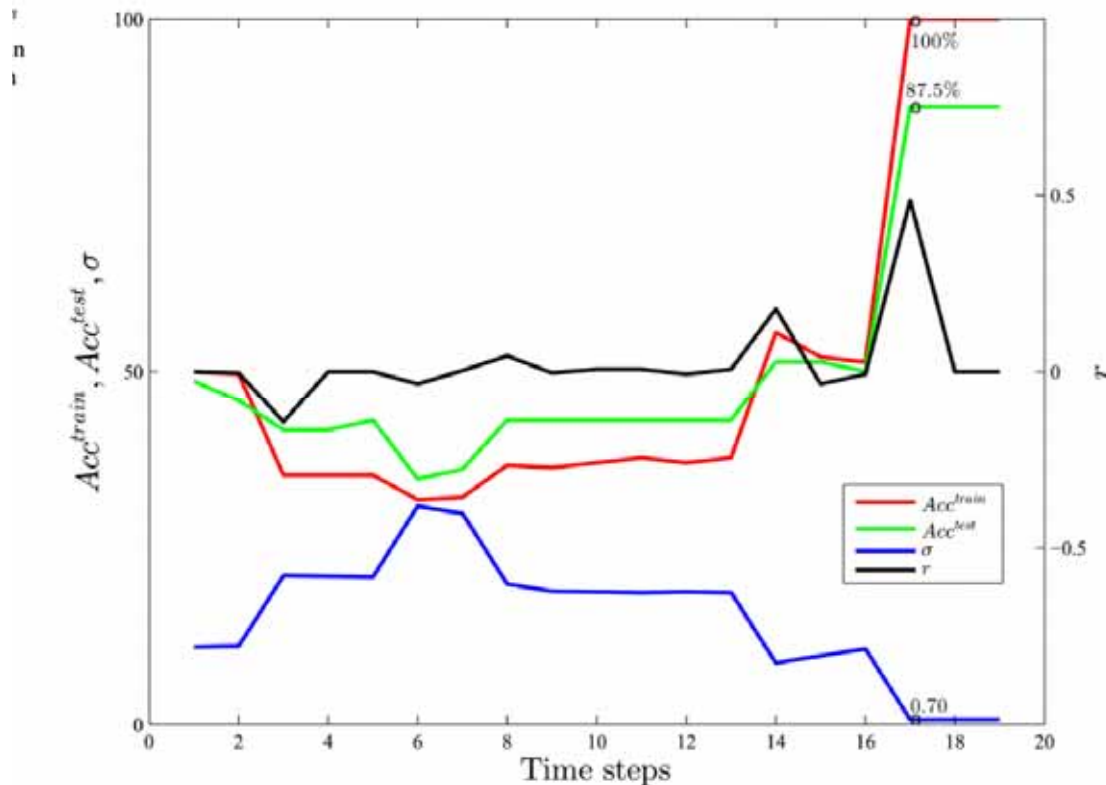
+ Add to Share ... More

Published on Sep 22, 2015  
Top 9 Medical Robots That Could Change Healthcare

<https://www.youtube.com/watch?v=20sj7rRfzm4>



Kusy, M. & Zajdel, R. 2014. Probabilistic neural network training procedure based on Q(0)-learning algorithm in medical data classification. *Applied Intelligence*, 41, (3), 837-854, doi:10.1007/s10489-014-0562-9.

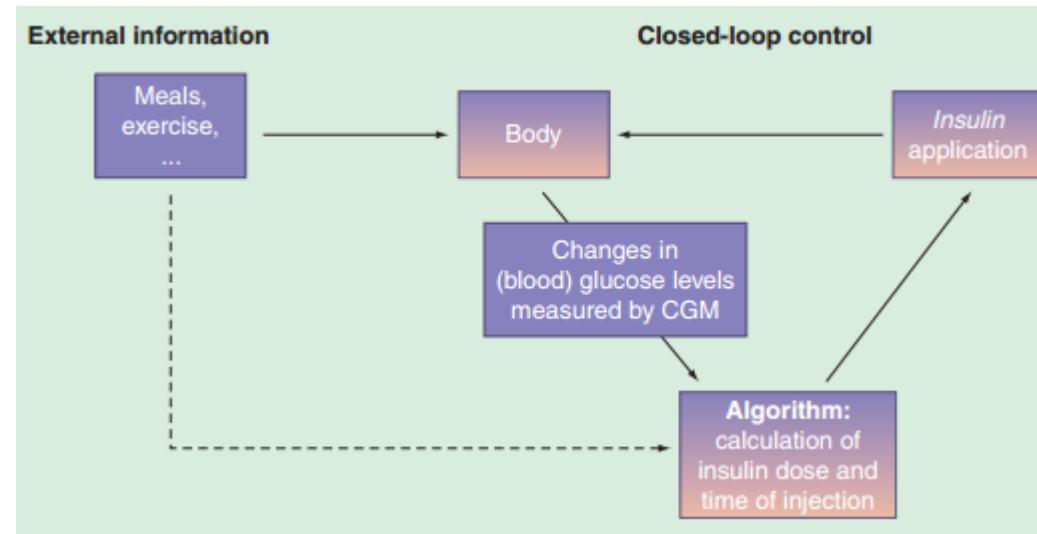
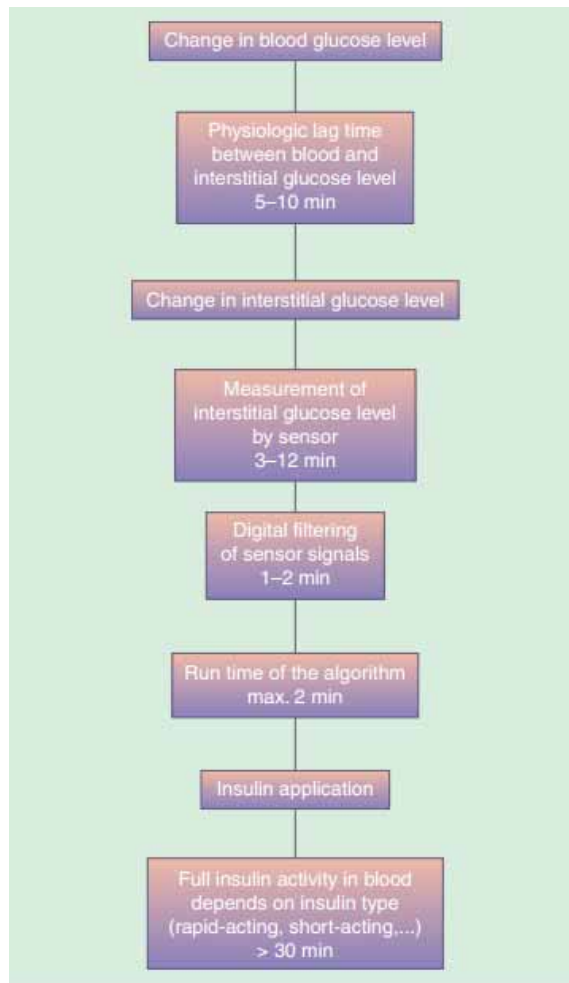


- Wisconsin breast cancer database [24] that consists of 683 instances with 9 attributes. The data is divided into two groups: 444 benign cases and 239 malignant cases. Pima Indians diabetes data set [36] that includes 768 cases having 8 features. Two classes of data are considered: samples tested negative (500 records) and samples tested positive (268 records).

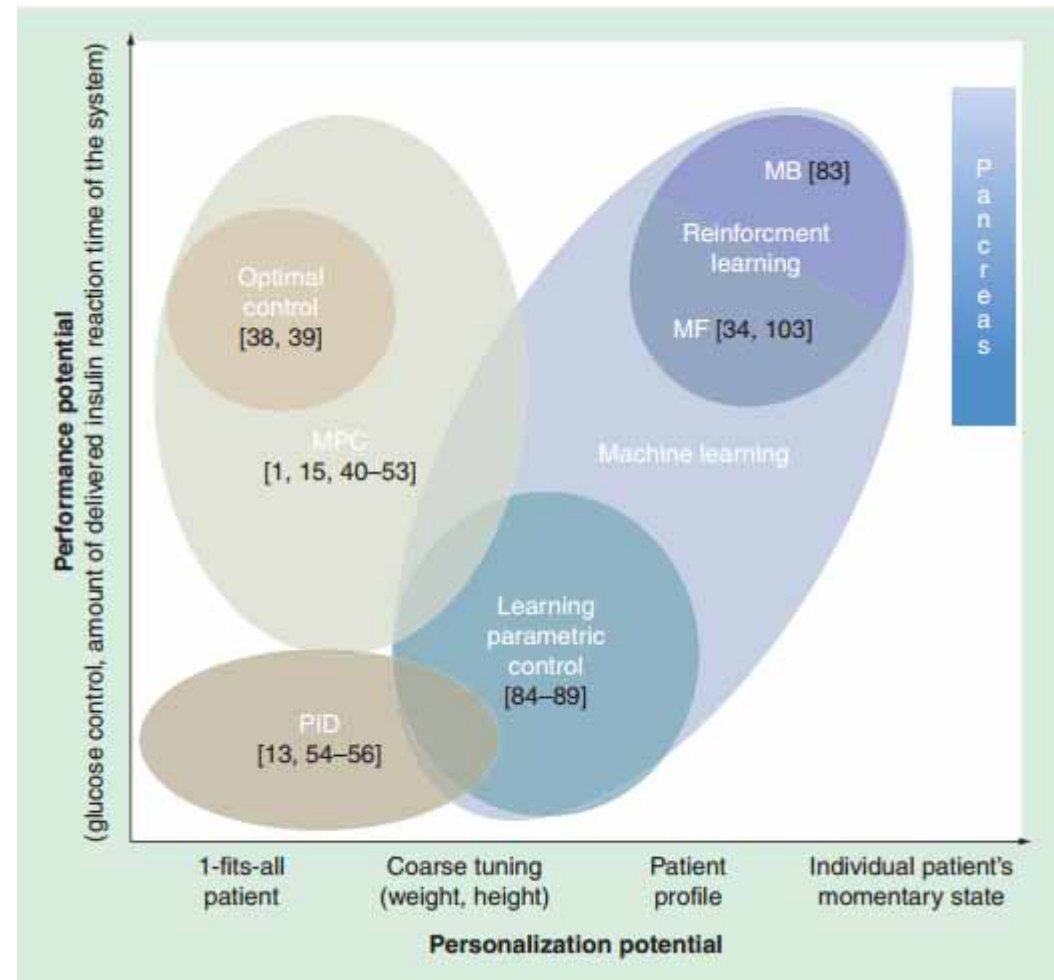
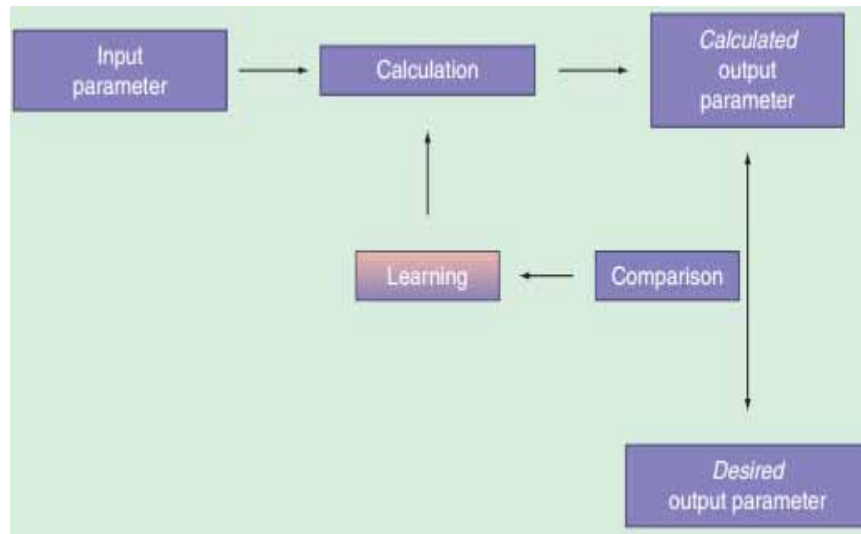
Haberman's survival data [21] that contains 306 patients who underwent surgery for breast cancer. For each instance, 3 variables are measured. The 5-year survival status establishes two input classes: patients who survived 5 years or longer (225 records) and patients who died within 5 years (81 records).

Cardiotocography data set [3] that comprises 2126 measurements of fetal heart rate and uterine contraction features on 22 attribute cardiotocograms classified by expert obstetricians. The classes are coded into three states: normal (1655 cases), suspect (295 cases) and pathological (176 cases).

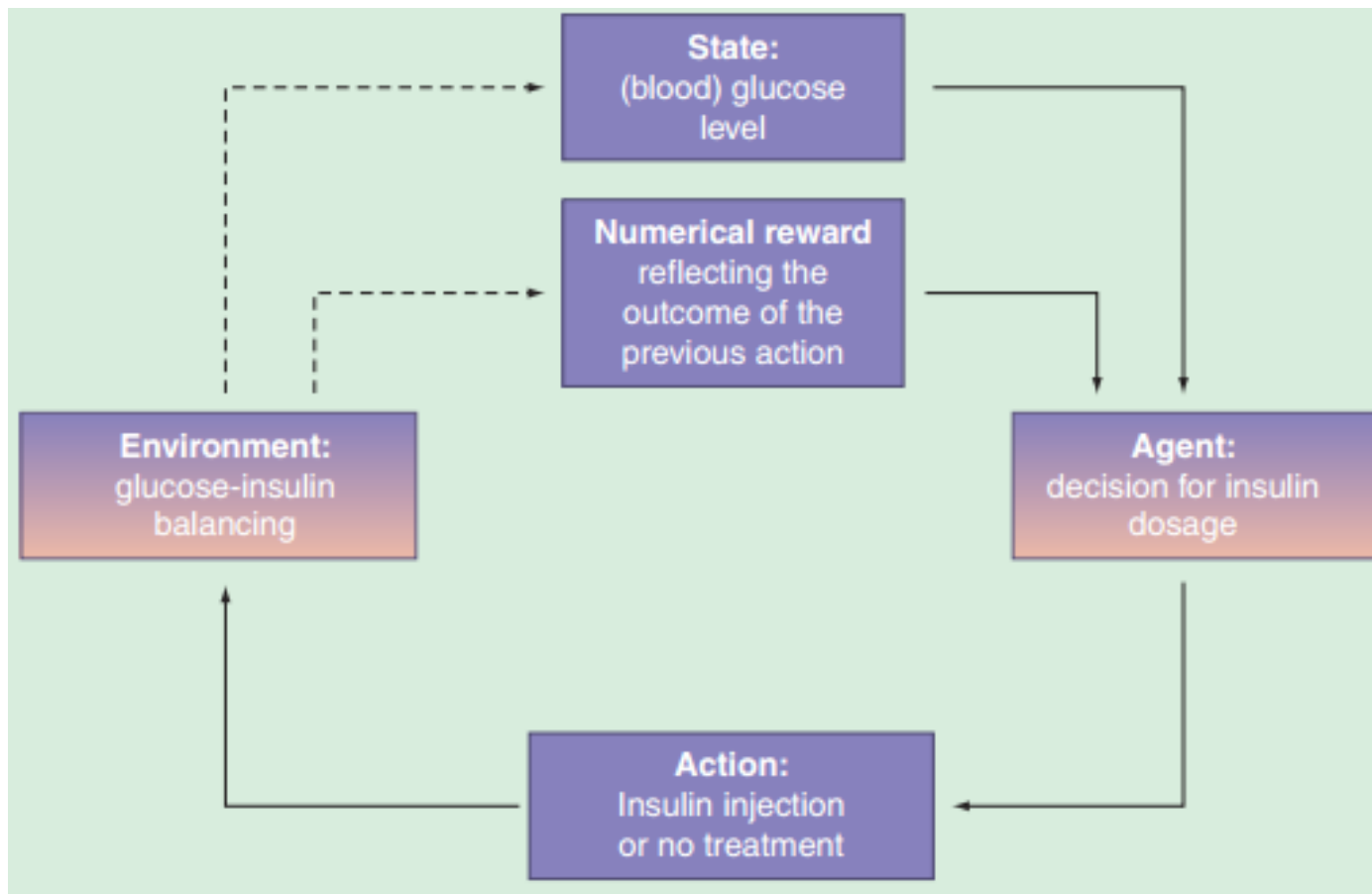
Dermatology data [13] that includes 358 instances each of 34 features. Six data classes are considered: psoriasis (111 cases), lichen planus (71 cases), seborrheic dermatitis (60 cases), chronic dermatitis (48 cases), pityriasis rosea (48 cases) and pityriasis rubra pilaris (20 cases).



Bothe, M. K., Dickens, L., Reichel, K., Tellmann, A., Ellger, B., Westphal, M. & Faisal, A. A. 2013. The use of reinforcement learning algorithms to meet the challenges of an artificial pancreas. Expert Review of Medical Devices, 10, (5), 661-673, doi:10.1586/17434440.2013.827515.

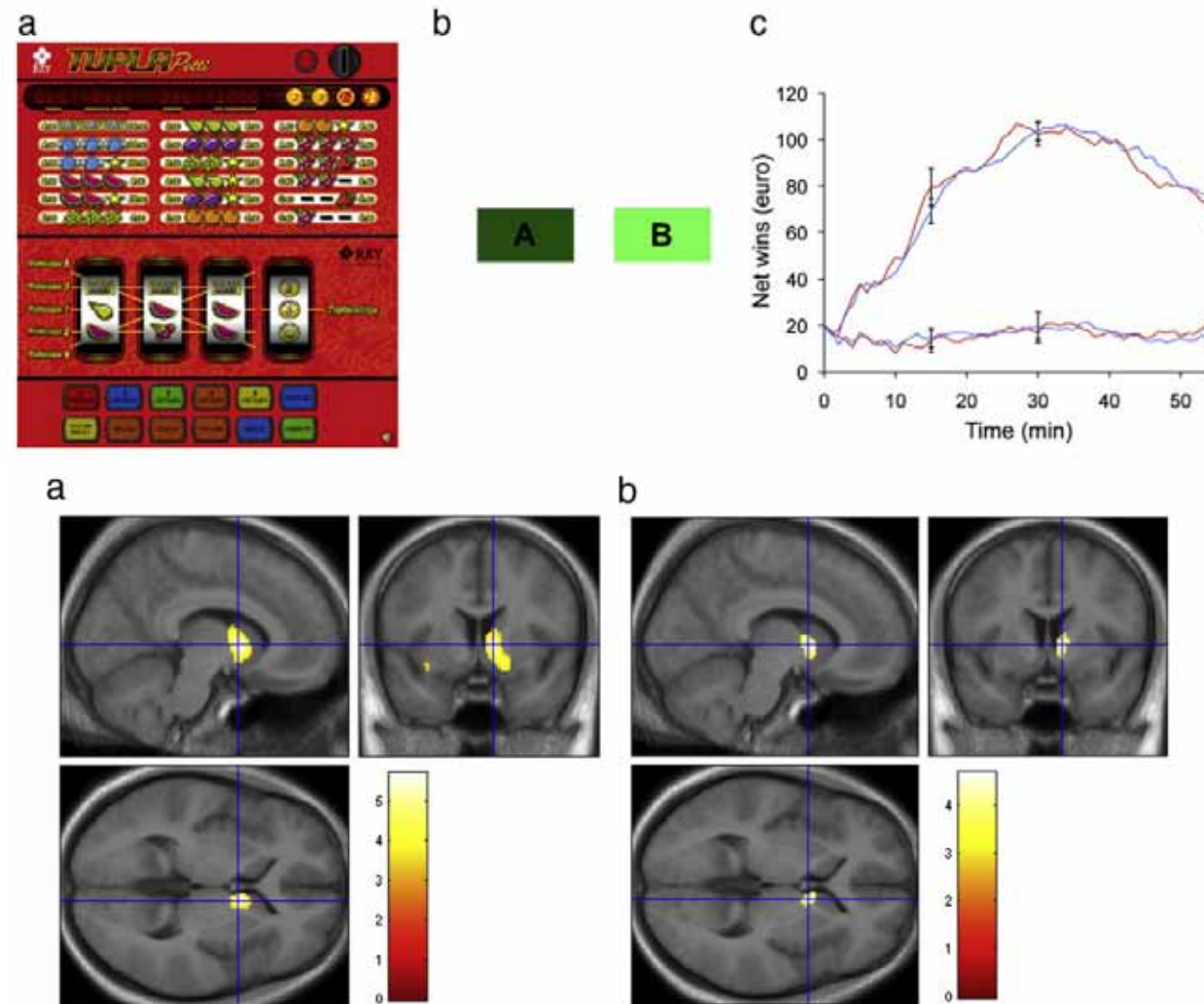


Bothe, M. K., Dickens, L., Reichel, K., Tellmann, A., Ellger, B., Westphal, M. & Faisal, A. A. 2013. The use of reinforcement learning algorithms to meet the challenges of an artificial pancreas. Expert Review of Medical Devices, 10, (5), 661-673, doi:10.1586/17434440.2013.827515.



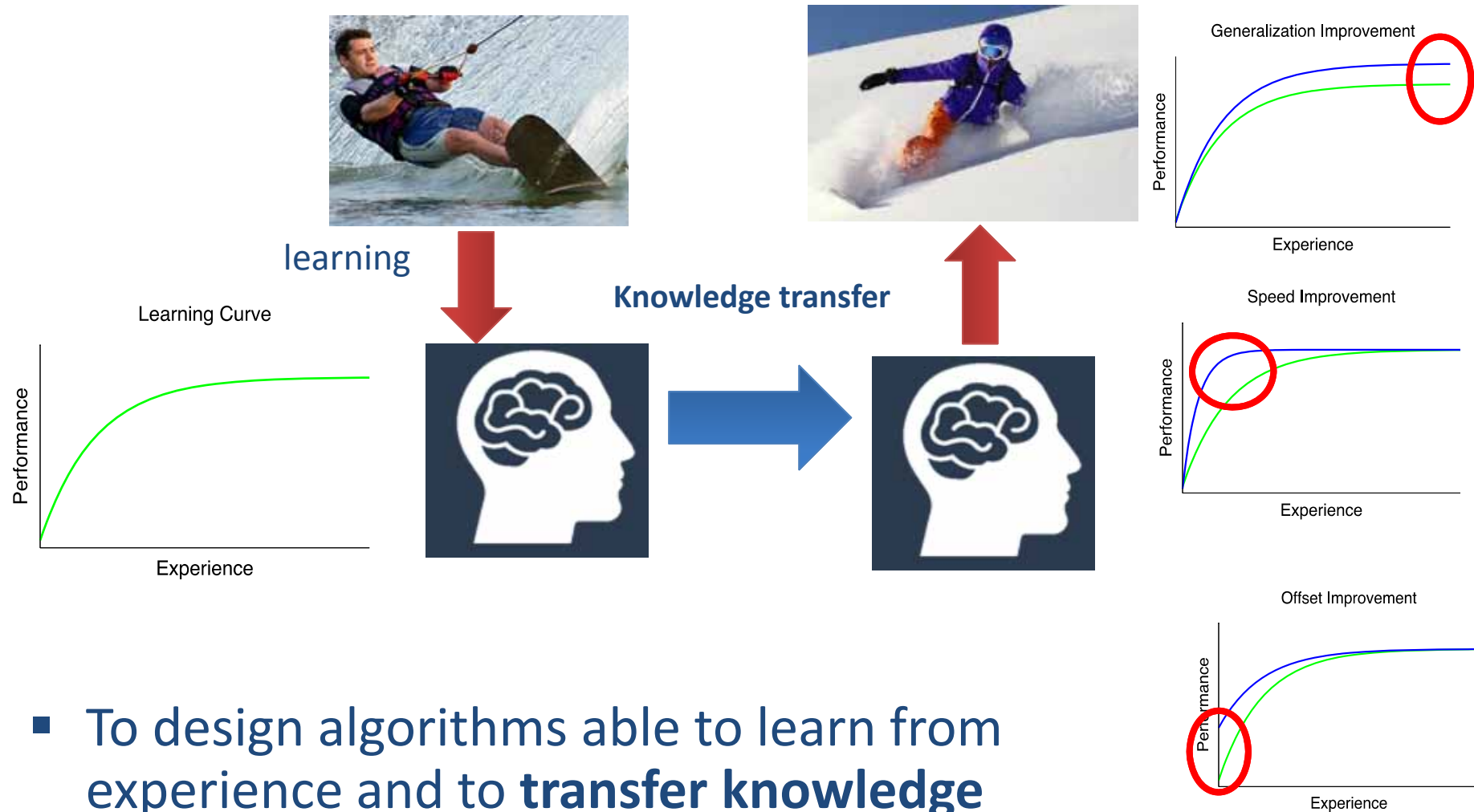
Bothe, M. K., Dickens, L., Reichel, K., Tellmann, A., Ellger, B., Westphal, M. & Faisal, A. A. 2013. The use of reinforcement learning algorithms to meet the challenges of an artificial pancreas. *Expert Review of Medical Devices*, 10, (5), 661-673, doi:10.1586/17434440.2013.827515.



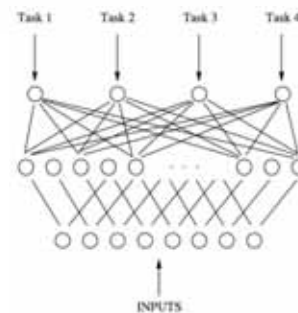
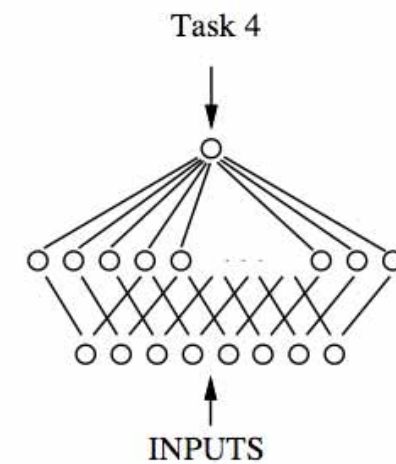
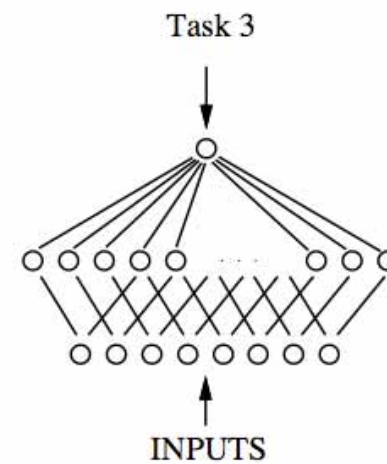
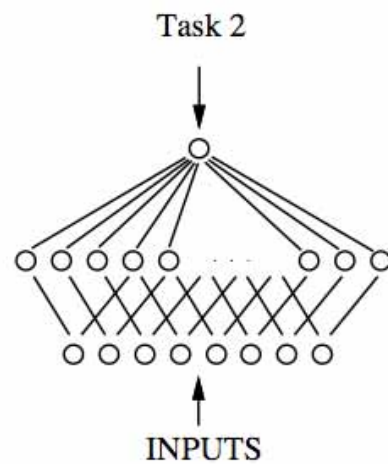
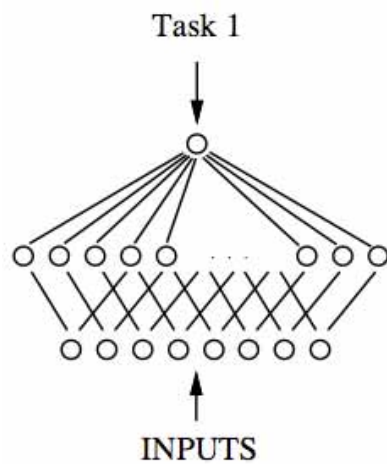
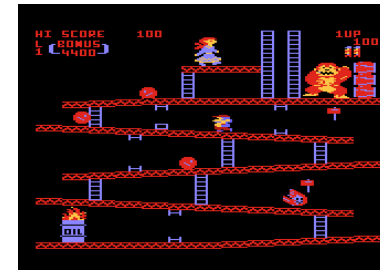


Joutsa et al. Mesolimbic dopamine release is linked to symptom severity in pathological gambling. *NeuroImage*, 60, (4), 1992-1999, doi.org/10.1016/j.neuroimage.2012.02.006.

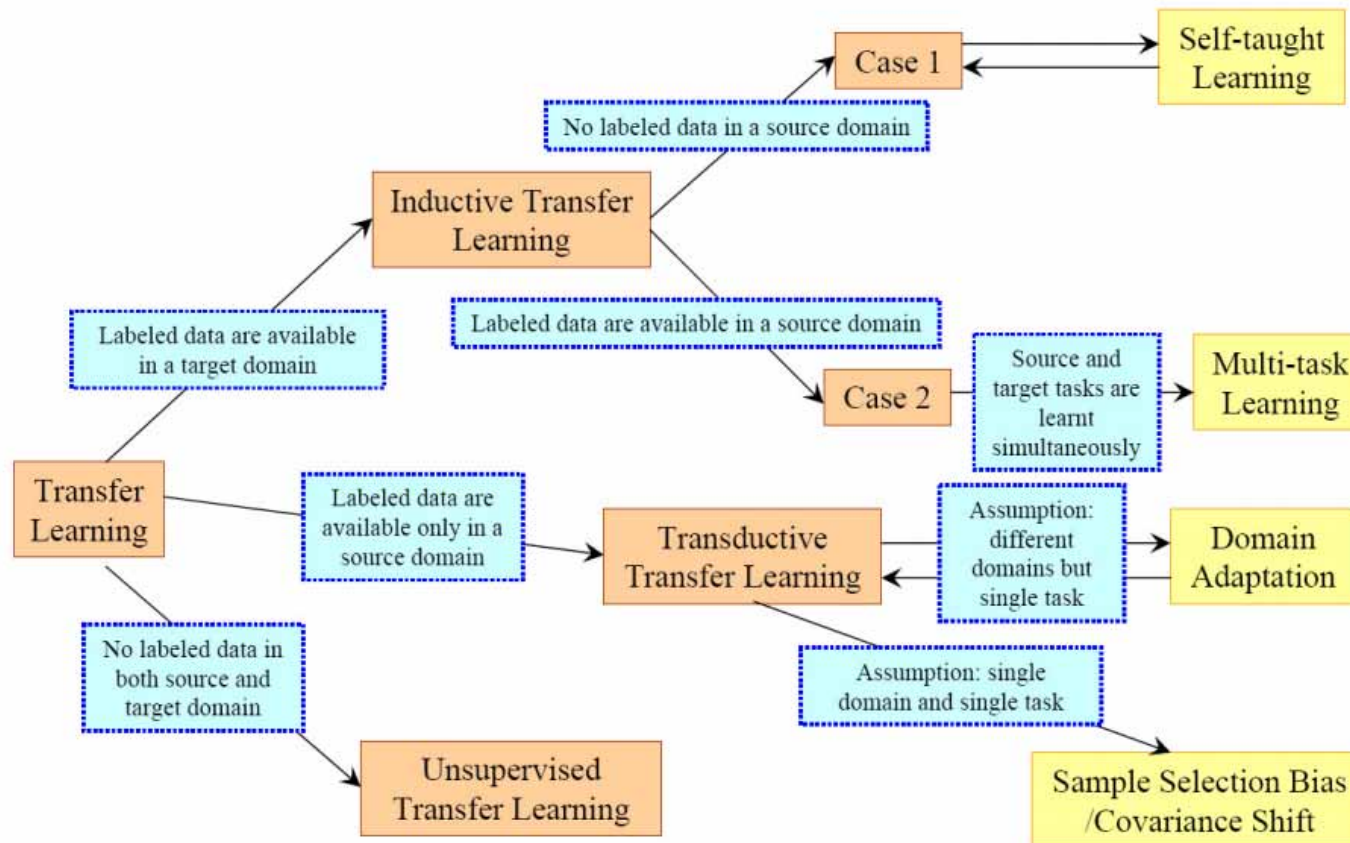
# 8) Future Outlook



- To design algorithms able to learn from experience and to **transfer knowledge across different tasks and domains** to improve their learning performance



V. Mnih *et al.*, "Playing Atari with Deep Reinforcement Learning", Nature (2015)  
 Rich Caruana, "Multi-task Learning", MLJ (1998)



Pan, S. J. & Yang, Q. A. 2010. A Survey on Transfer Learning. IEEE Transactions on Knowledge and Data Engineering, 22, (10), 1345-1359, doi:10.1109/tkde.2009.191.

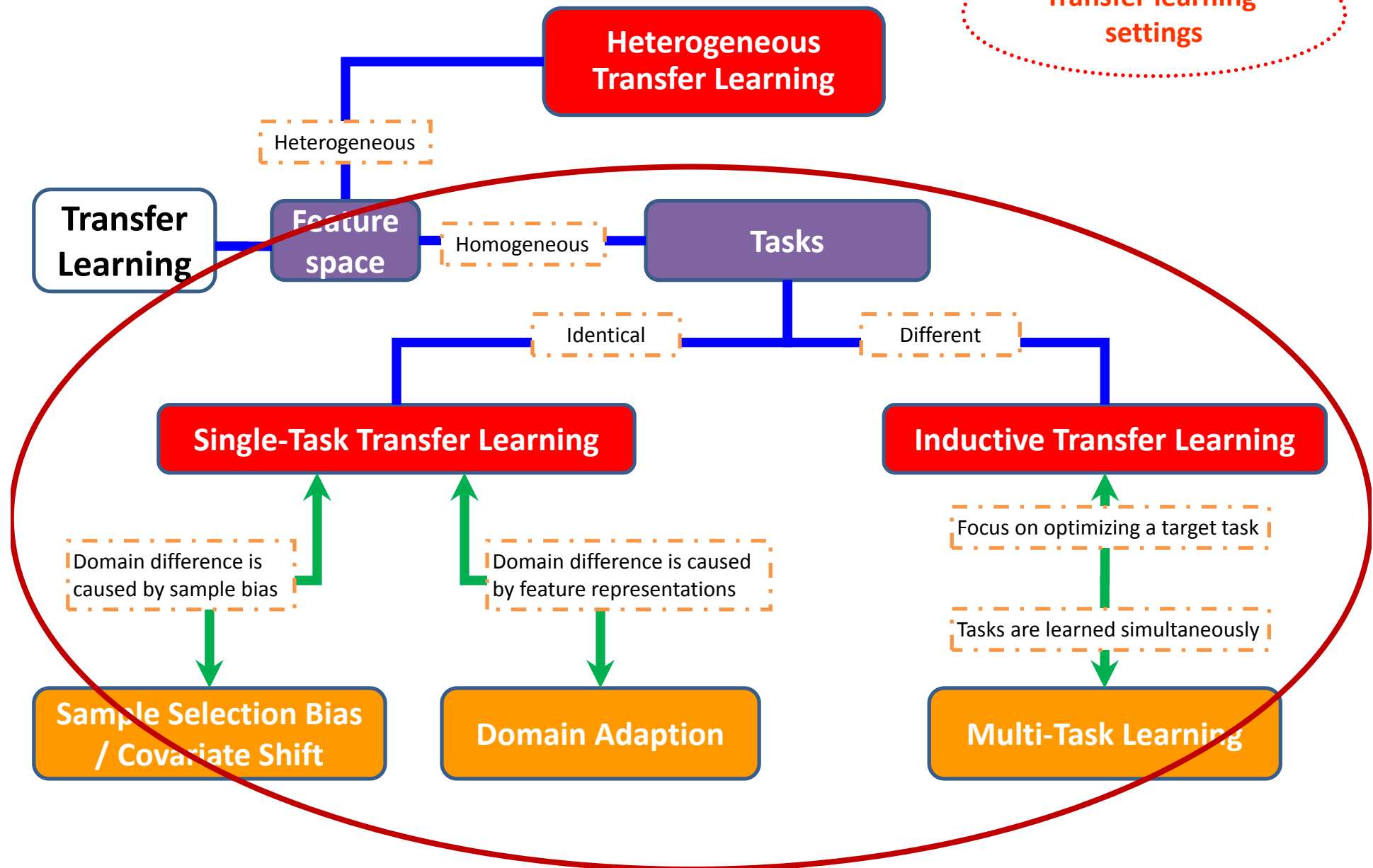
- Thorndike & Woodworth (1901) explored how individuals would transfer in one context to another context that share similar characteristics:
- They explored how individuals would transfer learning in one context to another, similar context
- or how "improvement in one mental function" could influence a related one.
- Their theory implied that transfer of learning depends on how similar the learning task and transfer tasks are,
- or where "identical elements are concerned in the influencing and influenced function", now known as the identical element theory.
- Today example: C++ -> Java; Python -> Julia
- Mathematics -> Computer Science
- Physics -> Economics

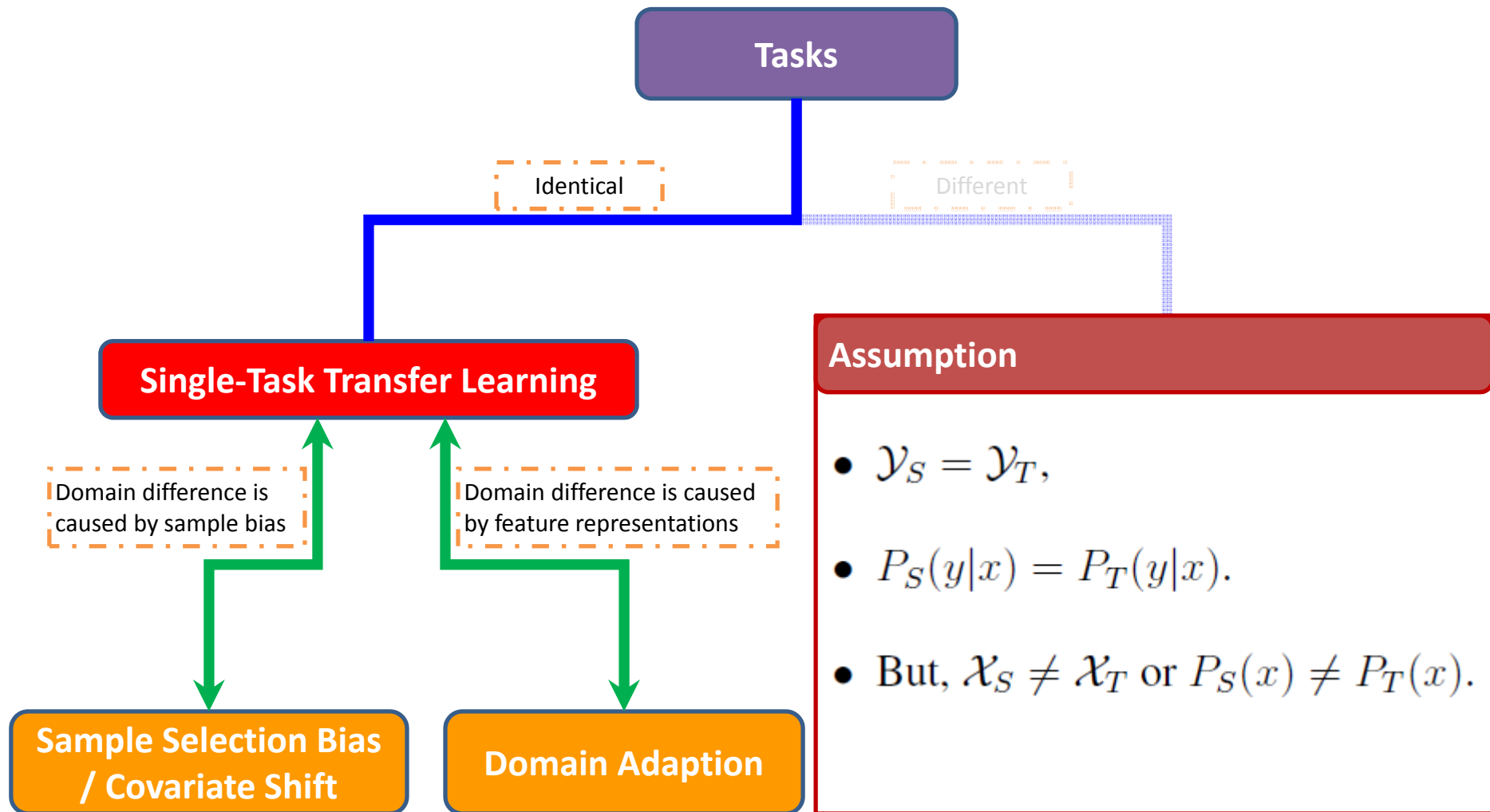
- Feature space  $\mathcal{X}$ ;
- $P(x)$ , where  $x \in \mathcal{X}$ .
- Given  $\mathcal{X}$  and label space  $\mathcal{Y}$ ;
- To learn  $f : x \rightarrow y$ , or estimate  $P(y|x)$ , where  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ .

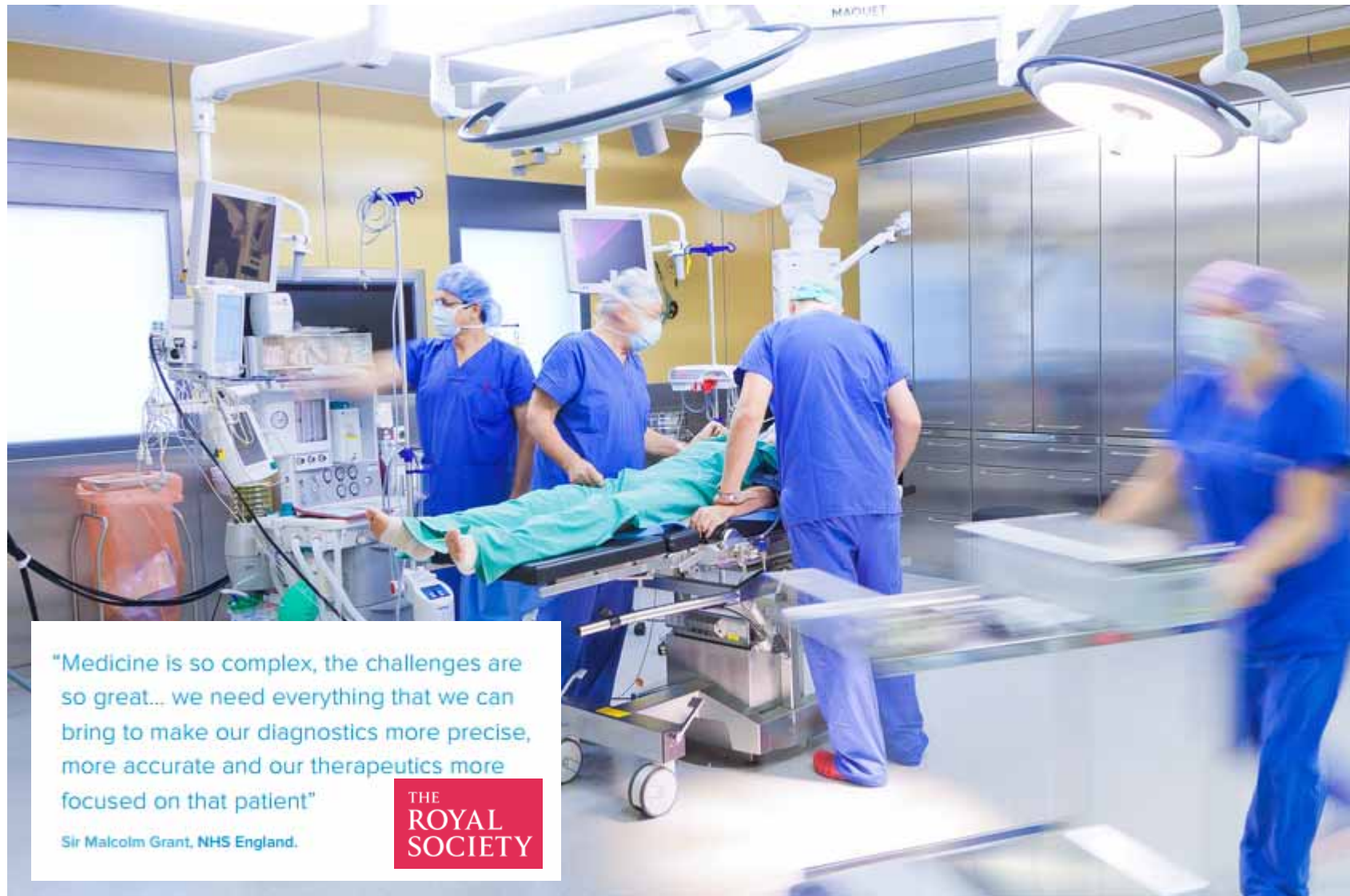
Two domains are different  $\Rightarrow \mathcal{X}_S \neq \mathcal{X}_T$ , or  $P_S(x) \neq P_T(x)$ .  
 Two tasks are different  $\Rightarrow \mathcal{Y}_S \neq \mathcal{Y}_T$ , or  $f_S \neq f_T$  ( $P_S(y|x) \neq P_T(y|x)$ ).

Pan, S. J. & Yang, Q. A. 2010. A Survey on Transfer Learning. IEEE Transactions on Knowledge and Data Engineering, 22, (10), 1345-1359, doi:10.1109/tkde.2009.191.









<https://royalsociety.org/events/2015/05/breakthrough-science-technologies-machine-learning>



# Thank you!

- Why is RL for us in health informatics interesting?
- What is a medical doctor in daily clinical routine doing most of the time?
- Please explain the human decision making process on the basis of the model by Wickens (1984) !
- What is the underlying principle of DQN?
- What is probabilistic inference? Give an example!
- Why is selective attention so important?
- Please describe the “anatomy” of a RL-agent!
- What does policy-based RL-agent mean? Give an example!
- What is the underlying principle of a MAB? Why is it interesting for health informatics?

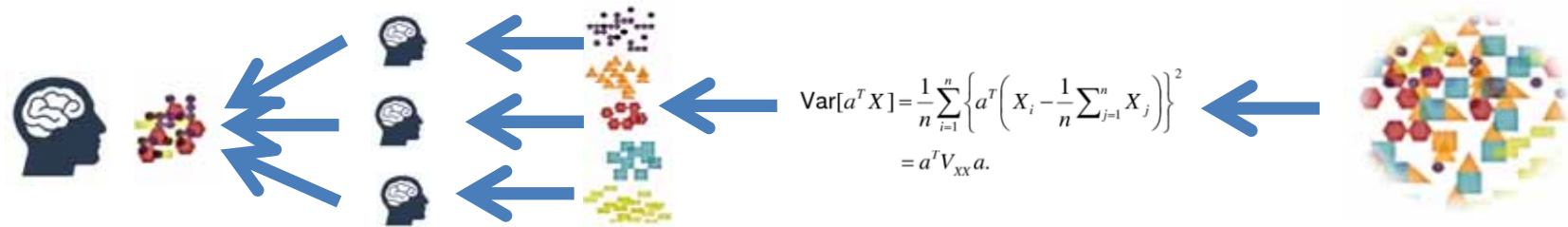
- Reinforcement Learning
- Trial-and-Error Learning
- Markov-Decision-Process
- Utility-based agent
- Q-Learning
- Passive reinforcement learning
- Adaptive dynamic programming
- Temporal-difference learning
- Active reinforcement learning
- Bandit problems



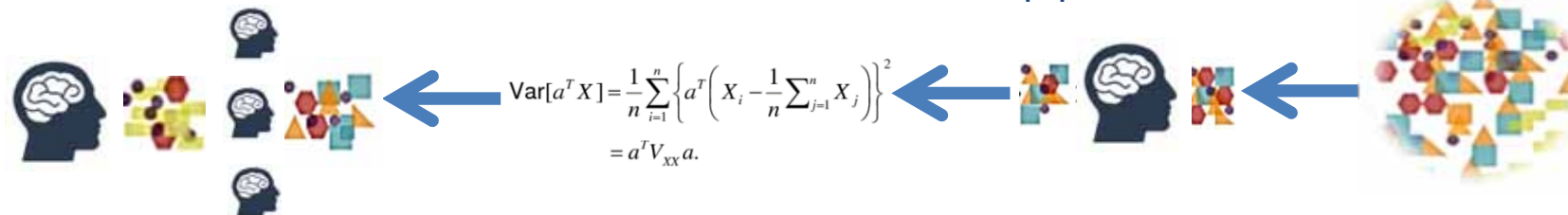
- RL:= general problem, inspired by behaviorist psychology; how software agents learn to make decisions from success and failure, from reward and punishment in an environment – aiming to maximize cumulative reward.
- RL is studied in game theory, control theory, operations research, information theory, simulation-based optimization, multi-agent systems, swarm intelligence, genetic algorithms.
- Aka: approximate dynamic programming.
- The problem has been studied in the theory of optimal control, though most studies are concerned with the existence of optimal solutions and their characterization, and not with the learning or approximation aspects. In economics and game theory, reinforcement learning may be used to explain how equilibrium may arise under bounded rationality.



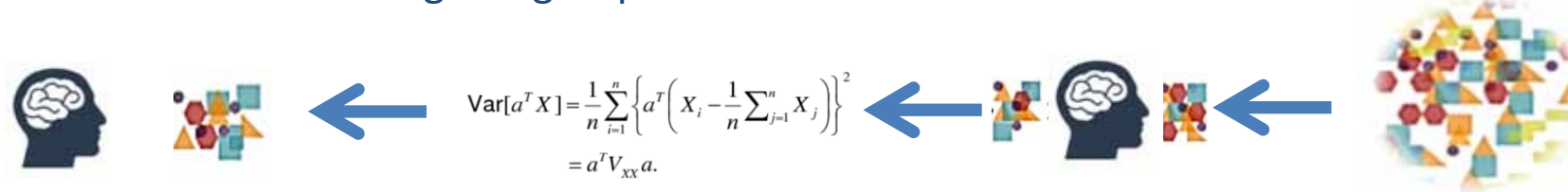
A) Unsupervised ML: Algorithm is applied on the raw data and learns fully automatic – Human can check results at the end of the ML-pipeline



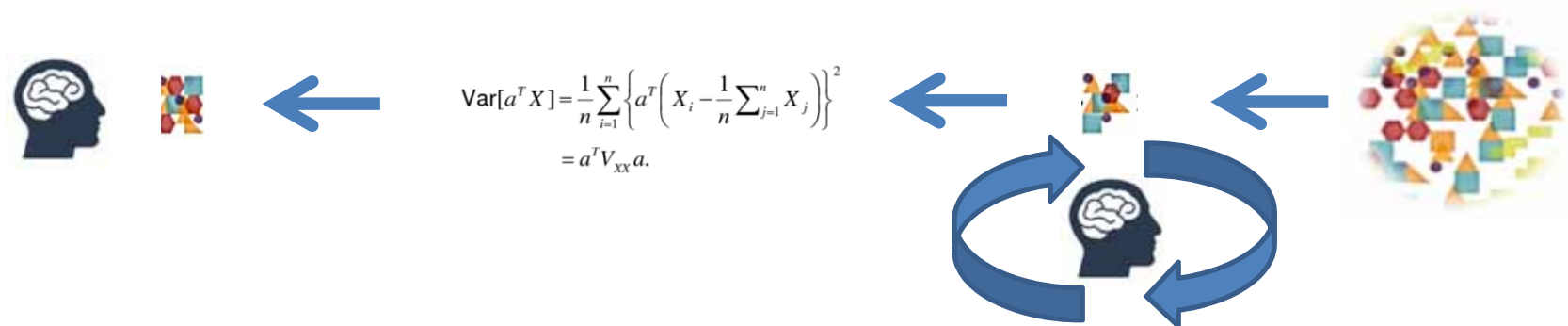
B) Supervised ML: Humans are providing the labels for the training data and/or select features to feed the algorithm to learn – the more samples the better – Human can check results at the end of the ML-pipeline



C) Semi-Supervised Machine Learning: A mixture of A and B – mixing labeled and unlabeled data so that the algorithm can find labels according to a similarity measure to one of the given groups

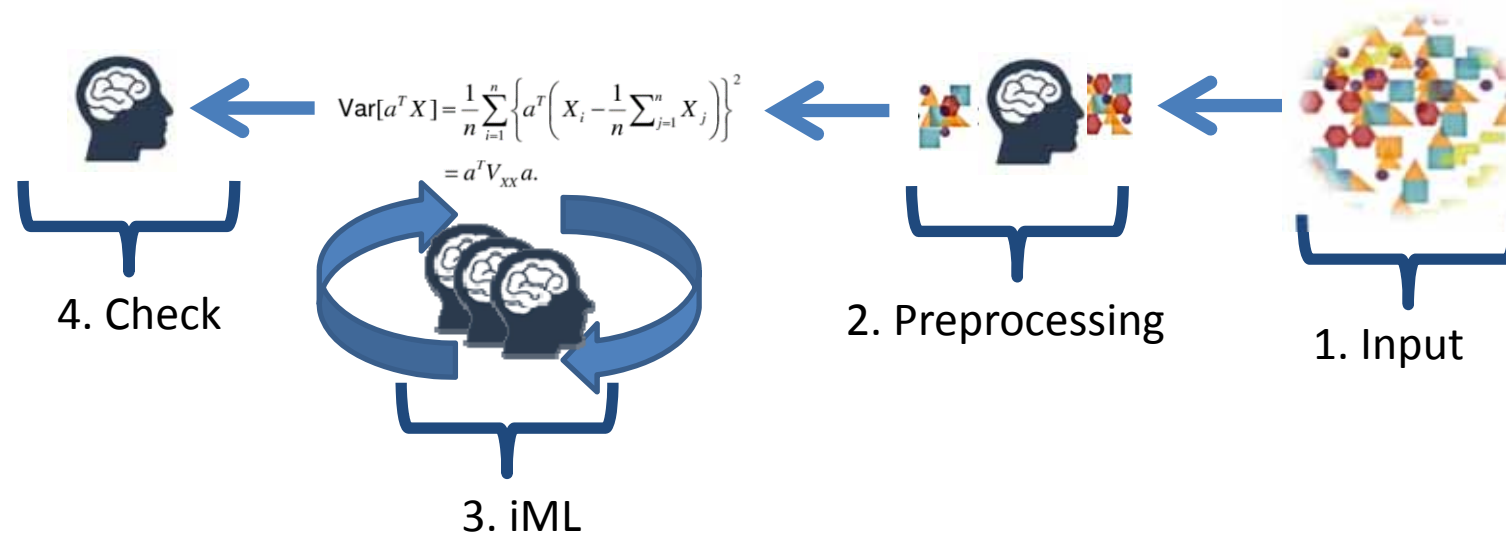


D) Reinforcement Learning: Algorithm is continually trained by human input, and can be automated once maximally accurate



- Advantage: non-greedy nature
- Disadvantage: must learn model of environment

**E) Interactive Machine Learning:** Human is seen as an agent involved in the actual learning phase, step-by-step influencing measures such as distance, cost functions ...



**Constraints of humans:** Robustness, subjectivity, transfer?

**Open Questions:** Evaluation, replicability, ...

Holzinger, A., Plass, M., Holzinger, K., Crisan, G., Pintea, C. & Palade, V. 2016. Towards interactive Machine Learning (iML): Applying Ant Colony Algorithms to solve the Traveling Salesman Problem with the Human-in-the-Loop approach. Springer Lecture Notes in Computer Science LNCS 9817. Heidelberg, Berlin, New York: Springer, pp. in print.