

Presentation at NWO Cognition Programme

<http://www.nwo.nl/cognitie/symposium>

Utrecht 24 Jun 2005

---

# DO MACHINES, NATURAL OR ARTIFICIAL, REALLY NEED EMOTIONS?

**Aaron Sloman**

**School of Computer Science,**

**The University of Birmingham, UK**

<http://www.cs.bham.ac.uk/~axs/>

Many people believe that emotions are required for intelligence. I argue that this is mostly based on (a) wishful thinking and (b) a failure adequately to analyse the variety of types of affective states and processes that can arise in different sorts of architectures produced by biological evolution or required for artificial systems. This work is a development of ideas presented by Herbert Simon in the 1960s in his 'Motivational and emotional controls of cognition'.

---

Updated version of talk at the Birmingham Cafe Scientifique, May 2004

<http://www.cs.bham.ac.uk/research/cogaff/talks/#cafe04>

Last updated (December 5, 2008)

# Key ideas

---

- Our ordinary concepts such as ‘emotion’, ‘consciousness’, ‘feeling’ are really too full of muddle and confusion to be useful in posing scientific questions or formulating explanatory theories
- These ideas evolved for purposes of ordinary communication among lay people and are fine for their original purposes not suitable for use in unambiguous scientific communication.
- We can refine, extend, and subdivide them to produce new more precise, more theoretically-based concepts if we can specify **explanatory architectures** and see what kinds of states and processes they can generate (as happened when physics and chemistry revised and extended our concepts of kinds of matter and kinds of physical and chemical states and processes).
- We need to understand a wide range of architectures – not just humans:  
E.g. [can an insect have emotions?](#) (whatever we think emotions are).  
Or an octopus? Read what the keepers say about Otto the octopus here  
<http://www.telegraph.co.uk/news/newstoppers/howaboutthat/3328480/Otto-the-octopus-wrecks-havoc.html>
- We should resist ‘wishful thinking’ when we try to do science.

For a start, we need to be able to think about architectures of various kinds.

# What is an architecture?

---

A house, a ship a symphony, a computer operating system, a company, a novel, a poem, a mathematical proof, an organisation ... can have an architecture: What is it they have?

- Each of those entities is something complex made of parts, which can also be made of parts made of parts ....
- The parts can have various sorts of relationships to other parts, including being close to or remote from them, having influences in either or both directions, sharing resources, cooperating to perform some function, interacting with external objects, ....
- Some of the parts are physical, like the parts of a house or a ship, while others are more abstract such as the overture of a symphony or the mission of the organisation.
- Talking about the architecture of X is talking about what the components of X are and how they are related to X and to one another, including what they do and what they do it for.

In that sense a mind can have an architecture too: it has components that perform various functions, including perceiving, making inferences, learning, storing information, generating motives, forming plans, forming theories, selecting motives to act on, controlling execution of plans, detecting inconsistencies, resolving conflicts, detecting signs of danger or signs of useful opportunities, and many more.

# The Design-based Approach to study of mind

---

When scientists discuss experimental observations, and evaluate theories, they often formulate questions using language that evolved for informal discourse among people engaged in every day social interaction, like this:

*What does the infant/child/adult/chimp/crow (etc) perceive/understand/learn/intend (etc)?*

*What is he/she/it conscious of?*

*What does he/she/it experience/enjoy/desire?*

*What is he/she/it attending to?*

*What is he/she/it trying to do?*

I suggest that although those questions are very useful in everyday life and some clinical settings, if we wish to gain increased scientific insight we should ask questions like these:

*Which parts of the architecture are involved?*

*What are their functions?*

*What kinds of information do they acquire and use?*

*How do they do this?*

*What is the total architecture in which they function?*

*How is the information represented? (It could be represented differently in different subsystems).*

*What kinds of manipulations and uses of the information occur?*

*What mechanisms make those processes possible?*

*How are the internal and external behaviours selected/controlled/modulated/coordinated?*

*How many different virtual machine levels are involved and how are they related (e.g. physical, chemical, neural, subsymbolic, symbolic, cognitive,...)?*

In other words: The scientific study of mind needs the “design-based” approach.

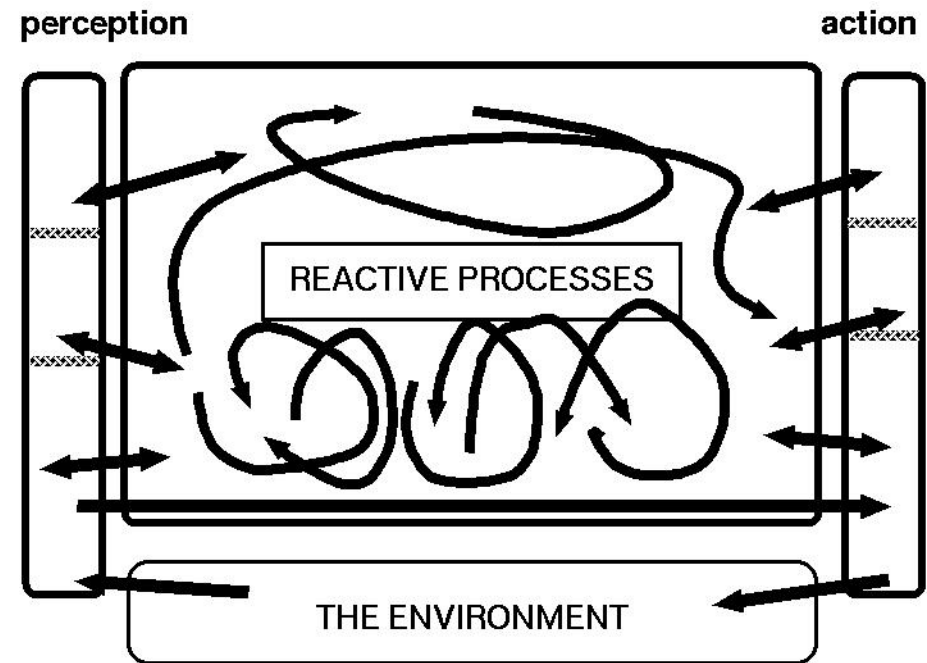
# Example: A simple (insect-like) architecture

A reactive system does not construct complex descriptions of possible futures, evaluate them and then choose one.

**It simply reacts: internally or externally.**

Several reactive sub-mechanisms may operate in parallel. The arrows indicate information flow (including control).

Processing may use a mixture of analog (continuously changing) and discrete mechanisms (changing step-wise).



An adaptive system with reactive mechanisms can be a very successful biological machine.

Some purely reactive species also have a social architecture, e.g. ants, termites, and other insects.

# Most animals are purely reactive

---

The vast majority of animal species are entirely reactive, lacking deliberative capabilities (explained later).

If an ant could consider several possible actions, and for each one imagine several steps ahead, evaluating the options, and then choose one of the multi-step options and carry it out, it would not be purely reactive.

It would include **deliberative** capabilities.

Most organisms (e.g. microbes, invertebrates, including insects) have more or less complex architectures, including perceptual and motor processes functioning at different levels of abstraction, e.g. detecting moisture, detecting an opportunity to mate.

But they don't appear to explore future possibilities in their minds before deciding what to do.

If they did they would have **deliberative** competences: they would not be purely **reactive**.

Purely reactive biological species are **precocial**: they have large amounts of genetically determined capabilities, though minor environmentally driven adaptations are possible.

They lack mechanisms required for rapid creative learning.

But they can be biologically very successful.

# What else is needed for emotions?

---

## CLAIM:

The most general concept that fits our fuzzy and indeterminate mish-mash of uses of the words 'emotion' and 'emotional' is:

**A state in which a monitoring mechanism acquires a tendency (i.e. a disposition, possibly suppressed) to abort, redirect, or modulate some other process or collection of processes.**

*Example: a house-fly consuming food detects something rapidly descending towards it: the 'alarm' mechanism aborts eating and triggers escape behaviour.*

States and processes that people label as 'emotions' vary enormously, involving both evolutionarily old and new mechanisms, producing both short term and long term states (e.g. grief, jealousy, infatuation, ambition) – some of which may occur in future robots.

There need not be any specific **bodily** changes involved: the important things are **control** changes (some in virtual machines)

*E.g. a mathematician working on an important new proof notices the possibility of a fallacy caused by implicit division by zero. This may trigger a disposition to switch to investigating the offending step in the proof.*

Some of these disruptions can be unconscious – like the people who are jealous or infatuated and don't realise it, though it is evident to their friends (a phenomenon exploited by novelists and playwrights).

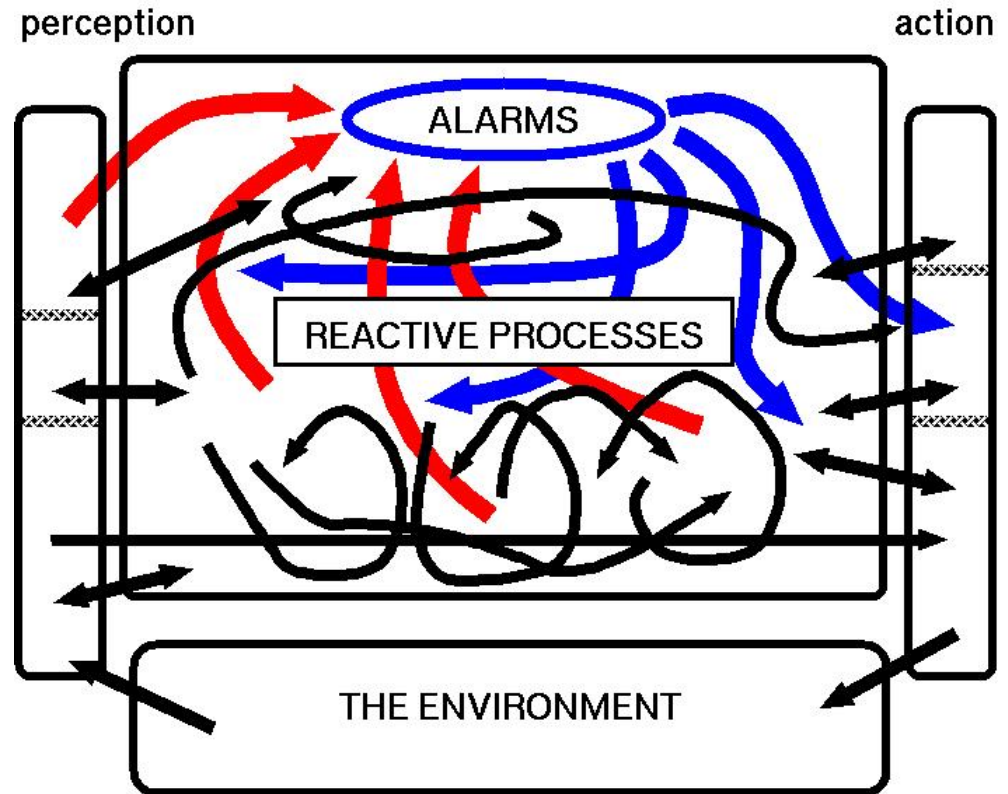
Often the tendency is not resisted and has **immediate** effects. More subtly, the disruptive tendency may be suppressed or overridden, but it is still there, competing for control, and sometimes takes over (e.g. a grieving parent reminded of a dead child months or years after the death).

# Primary emotions in insects?

Even insects may need simplified 'emotions' – e.g. detecting and reacting to unexpected dangers or opportunities, using fast pattern recognition mechanisms.

'Alarm' mechanisms running in parallel with other things, using fast pattern recognition, can detect events that trigger some interference with 'normal' processes, such as

- aborting
- accelerating
- redirecting
- freezing,.....



For humans, and many other animals, far more complex architectures are needed.

NOTE: Most empirical research reveals only shallow behavioural consequences of deep mechanisms and architectures, and studies shallow verbal classifications, not deep cognitive reactions to processes in other people – such as good novelists, playwrights, and poets describe.

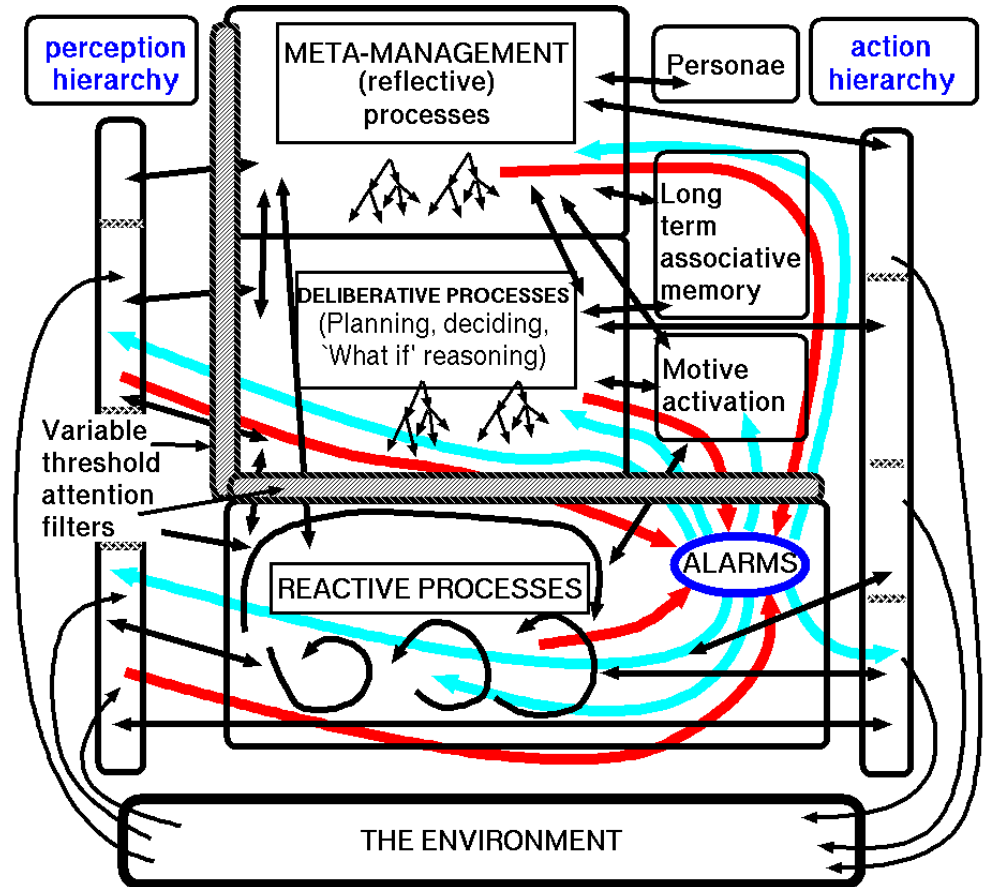
# Something more human: H-Cogaff

In humans, different subsystems operate concurrently, performing different sorts of tasks.

Some are evolutionarily old 'purely reactive' subsystems, and similar to mechanisms in many other kinds of animals.

Other subsystems are newer and do tasks that far fewer animals can perform. (E.g. thinking about past events, remote events, future events, and what another person is thinking.)

Perception and action operate concurrently at different levels of abstraction, in relation to different central sub-systems.



By considering processes in the different layers we can distinguish more kinds of emotion-like states, e.g. primary, secondary, tertiary, and far more, depending on how they differ in detail.

The diagram is explained in more detail in

<http://www.cs.bham.ac.uk/research/projects/cogaff/03.html#200307>

# Damasio's Error

---

- In 1994 Antonio Damasio, a well known neuroscientist, published his book *Descartes' Error*.

He argued that emotions are needed for intelligence, and accused Descartes and many others of not grasping that.

- In 1996 Daniel Goleman published *Emotional Intelligence: Why It Can Matter More than IQ*, quoting Damasio with approval.
- Likewise Rosalind Picard a year later in her book *Affective Computing*.
- Since then there has been a flood of publications and projects echoing Damasio's claim, and many researchers in Artificial Intelligence have become convinced that emotions are essential for intelligence, so they are now producing many computer models containing a module called 'Emotion'.
- Before that, serious researchers had begun to argue that the study of emotions and affect had not had its rightful place in psychology, and cognitive science, but the claims were more moderate.

E.g. a journal called *Cognition and Emotion* was started in 1987.

Even I had a paper in it in the first year. But H.A. Simon's work was mostly ignored.

Alas, all that theorising was not based on a deep understanding of varieties of information-processing architectures required to explain human phenomena.

It did not embrace 'the design-based approach' (defined above).

# Damasio's examples

---

Damasio's argument was partly based on two examples:

- **Phineas Gage**: In 1848, an accidental explosion of a charge he had set blew his tamping iron through his head – destroying the left frontal part of his brain.

“He lived, but having previously been a capable and efficient foreman, one with a well-balanced mind, and who was looked on as a shrewd smart business man, he was now fitful, irreverent, and grossly profane, showing little deference for his fellows. He was also impatient and obstinate, yet capricious and vacillating, unable to settle on any of the plans he devised for future action. His friends said he was No longer Gage.”

<http://www.deakin.edu.au/hbs/GAGEPAGE/Pgstory.htm>

Christopher Green, however, informs me that most popular reports on Gage exaggerate the effects of his injury. See <http://www.nthposition.com/anoddkindoffame.php>

- **Elliot, Damasio's patient** ('Elliot' was not his real name.)  
Following a brain tumor and subsequent operation, Elliot suffered damage in the same general brain area as Gage (left frontal lobe).

Like Gage, he experienced a great change in personality. Elliot had been a successful family man, and successful in business. After his operation he became impulsive and lacking in self-discipline. He could not decide between options where making the decision was important but both options were equally good. He perseverated on unimportant tasks while failing to recognize priorities. He had lost all his business acumen and ended up impoverished, even losing his wife and family. He could no longer hold a steady job. Yet he did well on standard IQ tests.

<http://serendip.brynmawr.edu/bb/damasio/>

# WHAT FOLLOWS FROM THIS?

---

Both patients appeared to retain high intelligence as measured by standard tests, but not as measured by their ability to behave sensibly.

Both had also lost certain kinds of emotional reactions.

WHAT FOLLOWS?

# Damasio's argument

---

Here is the essence of the argument Damasio produced, which many people in many academic disciplines enthusiastically accepted as valid:

There are two factual premises from which a conclusion is drawn.

P1 Damage to frontal lobes impairs certain emotional capabilities

P2 Damage to frontal lobes impairs intelligence

C Those emotional capabilities are required for intelligence

## IS THIS A VALID ARGUMENT?

The conclusion does not follow from the premises.

Whether the conclusion is true is a separate matter, discussed later.

In fairness, Damasio's book did not present the argument in quite such a bare fashion.

However, something like this argument is often approvingly attributed to him.

(An example is mentioned later.)

# Compare this argument

---

We 'prove' that cars need functioning horns in order to start, using two premises on which to base the conclusion:

P1 Damage to a car battery stops the horn working

P2 Damage to a car battery prevents the car starting

C A functioning horn is required for the car to start

DOES C FOLLOW FROM P1 AND P2?

# Why did readers not spot the fallacy?

---

A moment's thought should have reminded Damasio's readers that

- two capabilities **A** and **B** could presuppose some common mechanism **M**, so that
- damaging **M** would damage both **A** and **B**,
- without either of **A** or **B** being **required** for the other.

For instance, even if P1 and P2 are both true, you can damage the starter motor and leave the horn working, or damage the horn and leave the starter motor working!

NOTE:

I am ignoring two points, for the sake of illustration.

- Without a battery some cars can be started by rolling them downhill.
- There may be some cars which have separate batteries for starter motor and horn, in which case P1 and P2 would both be false generalisations.  
For such cars the premisses would be inappropriate because the phrase 'the battery' presupposes that there is only one.

# Why were so many people convinced?

---

Why are so many intelligent people convinced by Damasio's argument?  
I first criticised Damasio's argument in two papers in 1998 and 1999:

A. Sloman, (1998) Damasio, Descartes, Alarms and Meta-management, in *Proceedings International Conference on Systems, Man, and Cybernetics (SMC98)*, San Diego, IEEE, pp. 2652–7,

Available online: <http://www.cs.bham.ac.uk/research/cogaff/0-INDEX96-99.html#36>

A. Sloman, (1999) Review of Affective Computing by R.W. Picard, 1997, in *The AI Magazine*, 20, 1, pp. 127–133

Available online: <http://www.cs.bham.ac.uk/research/cogaff/0-INDEX96-99.html#40>

I have never seen these criticisms of Damasio's arguments made by other authors.

My criticisms were repeated in several subsequent publications.

Nobody paid any attention to the criticism and even people who had read those papers continued to refer approvingly to Damasio's argument in their papers.

Some even quote me as agreeing with Damasio!

Very intelligent people keep falling for the argument.

## WHY? What's wrong with emotion researchers?

E.g. Susan Blackmore did not notice the fallacy when she approvingly summarised Damasio's theories. See page 285 of her excellent recent book *Consciousness: An Introduction (2003)*. She has now informed me that she agrees that the argument used is fallacious.

# A sociological conjecture

---

The best explanation I can offer for the surprising fact that so many intelligent people are fooled by an obviously invalid argument is sociological:

they are part of a culture in which people **want** the conclusion to be true.

There seems to be a wide-spread (though not universal) feeling, even among many scientists and philosophers, that intelligence, rationality, critical analysis, problem-solving powers, are over-valued, and that they have defects that can be overcome by emotional mechanisms.

This leads people to **like** Damasio's conclusion. They **want** it to be true.

And this somehow causes them to accept as valid an argument for that conclusion, even though they would notice the flaw in a structurally similar argument for a different conclusion (e.g. the car horn example).[\*]

**A research community with too much wishful thinking does not advance science.**

Instead of being wishful thinkers, scientists trying to understand the most complex information-processing system on the planet should learn how to think (some of the time) as designers of information-processing systems do.

---

[\*] *This is a general phenomenon: consider how many people on both sides of the evolution/creation debate or both sides of the debate for and against computational theories of mind tend to accept bad arguments for their side.*

# A personal note

---

I find it curious that highly intelligent AI researchers who read my papers concerned with emotions apparently interpret my words as saying  
**what they wish I had said.**

For example:

- My mention of Damasio is taken as an endorsement of his arguments: people either don't read or don't understand what I write about his arguments (e.g. that such theories cannot account for enduring emotions such as long term grief).
- A paper I wrote with Monica Croucher in 1981, entitled 'Why robots will have emotions' is reported as if it had been 'Why robots **should** have emotions'.

See <http://www.cs.bham.ac.uk/research/cogaff/81-95.html#36>

- An online presentation attributes to me the **completely daft** theory that

**In each reasoning cycle:**

**A goal is selected as the focus of attention**

**The one that generates the strongest emotions (Sloman)**

I found that presentation when browsing the section on emotions in the euCognition web site.

[http://www.eucognition.org/affect\\_emotion\\_articles.htm](http://www.eucognition.org/affect_emotion_articles.htm)

We need to counter this wishful thinking, sloppy science, and sloppy reporting?

# To be fair ....

---

In fact Damasio produced additional theoretical explanations of what is going on, so, in principle, even though the quoted argument is invalid, the conclusion might turn out to be true and explained by his theories.

However:

- His theory of emotions as based on 'somatic markers' is very closely related to the theory of William James, which regards emotions as a form of awareness of bodily changes. This sort of theory is incapable of accounting for the huge subset of socially important emotions in humans which involve rich **semantic content** which would not be expressible within somatic markers (e.g. admiring someone's courage while being jealous of his wealth) and emotions that endure over a long period of time while bodily states come and go (such as obsessive ambition, infatuation, or long term grief at the death of a loved one).
- The key assumption, shared by both Damasio and many others whose theories are different, is that all choices depend on emotions, and especially choices where there are conflicting motives. If that were true it would support a conclusion that emotions are needed for at least intelligent conflict resolution.
- Although I will not argue the point here, I think it is very obvious from the experience of many people (certainly my experience) that one can learn how to make decisions between conflicting motives in a totally calm, unemotional, even cold way simply on the basis of having preferences or having learnt principles that one assents to. Many practical skills require learning which option is likely to be better. A lot of social learning provides conflict resolution strategies for more subtle decisions: again without emotions having to be involved.
- **A terminological decision to label all preferences, policies, and principles 'emotions' would trivialise Damasio's conclusion.**

**So, let's start again: what are emotions, and how do they work?**

# Does a Crow need emotions in order to be intelligent?

---

## SHOW BETTY MAKING A HOOK

<http://users.ox.ac.uk/~kgroup/tools/photos.shtml>

See the videos

<http://users.ox.ac.uk/~kgroup/tools/movies.shtml>

There are many more reports of this research in many web sites:

<http://www.google.co.uk/search?num=30&hl=en&q=betty+crow+hook+&btnG=Search&meta=>

### WARNING:

The BBC reporter in one of the reports misquotes the researchers as saying that 9 out of 10 female crows can solve the problem, when what the researchers actually said was that Betty solved the problem 9 out of 10 times.

Beware of reporters, even from the BBC.

# Does being emotional help a child solve his problem?

## SHOW A CHILD FAILING TO UNDERSTAND HOOKS AND GETTING EMOTIONAL

See the video (two versions: small and big)

[http://www.cs.bham.ac.uk/~axs/fig/josh34\\_0096.mpg](http://www.cs.bham.ac.uk/~axs/fig/josh34_0096.mpg) [4.2MB]

[http://www.cs.bham.ac.uk/~axs/fig/josh34\\_0096\\_big.mpg](http://www.cs.bham.ac.uk/~axs/fig/josh34_0096_big.mpg) [11 MB]

What changes when a child who cannot solve a problem later becomes able to solve it?

One possibility: being trained with rewards and punishment, like a circus animal.

Another possibility – learning (with or without emotions):

- to perceive new affordances
- to acquiring a richer ontology
- to use new forms of representation
- to use new procedures or algorithms for making use of the above
- to recognise their relevance to particular problems
- to think instead of getting emotional!

Is having emotions **necessary** for such learning?

# In order to explain all this

---

## We have to think about

- information
- representations
- mechanisms
- architectures

## And a host of related questions

- How many different kinds are there?
- What are the implications of the differences?
- How and why did the different sorts evolve?

In comparison, talking about emotions explains very little: for learning prompted or aided by emotions requires a vast amount of machinery that can also work without emotions, and often does work without emotions when highly intelligent people are coming to understand something new and complex.

Weak students may need many emotional props and motivators, however.

# Some old ways to study emotions

---

There are many ways to study emotions and other aspects of human minds:

- **Reading plays, novels, poems** will teach much about how people who have emotions, moods, attitudes, desires, etc. think and behave, and how others react to them — because many writers are very shrewd observers!
- **Studying ethology** will teach you something about how emotions and other mental phenomena vary among different animals.
- **Studying psychology** will add extra detail concerning what can be triggered or measured in laboratories, and what correlates with what.
- **Studying developmental psychology** can teach you how the states and processes in infants differ from those in older children and adults.
- **Studying neuroscience** will teach you about the physiological brain mechanisms that help to produce and modulate mental states and processes.
- **Studying therapy and counselling** can teach you about ways in which things can go wrong and do harm, and some ways of helping people.
- **Studying philosophy** with a good teacher may help you discern muddle and confusion in attempts to say what emotions are and how they differ from other mental states and processes.

**There's another way that complements those ways.**

## Another way to learn: do some engineering design

Suppose you had to design animals (including humans) or robots capable of living in various kinds of environments, including environments containing other intelligent systems.

What sorts of information-processing mechanisms, including control mechanisms, would you need to include in the design, and how could you fit all the various mechanisms together to produce all the required functionality, including:

- perceiving,
- learning,
- acquiring new motives,
- enjoying some activities and states and disliking others,
- selecting between conflicting motives,
- planning,
- reacting to dangers and opportunities,
- communicating in various ways
- reproducing, **and so on...**

If we combine this “design standpoint” with the previously listed ways to study mental phenomena, we can learn much about all sorts of mental processes: what they are, how they can vary, what they do, what produces them, whether they are essential or merely by-products of other things, how they can go wrong, etc.

The result could be both deep new insights about what we are, and important practical applications.

# The design-based approach – too fragmented now

The design-based approach is not new: over the last half century, researchers in Computational Cognitive Science, and in Artificial Intelligence have been pursuing it.

- Because the work was so difficult and because of the pressures of competition for funding and other aspects of academic life (e.g. lack of time for study), the field fragmented, and as more people became involved the research community became more fragmented, with each group investigating only a small subset of the larger whole, and talking only to members of that group.
- Deep, narrowly focused, research on very specific problems is a requirement for progress, but if **everybody** does only that, the results will be bad.
  - People working on natural language without relating it to studies of perception, thinking, reasoning, and acting may miss out on important aspects of how natural languages work.
  - Likewise those who study only a small sub-problem in perception may miss out ways in which the mechanisms they study need to be modified to fit into a larger system.
  - The study of emotions also needs to be related to the total system.

The European Community's recent initiative in 'Cognitive Systems' is an attempt to remedy this by requiring researchers to think about integrated multi-component systems.

One of the projects funded (including Birmingham) under that initiative is described here:

<http://www.cs.bham.ac.uk/research/projects/cosy/>

A UK grand challenge proposal to put all the pieces together again in a long term research programme is described here <http://www.cs.bham.ac.uk/research/cogaff/gc/>

See also the tutorial on [Representation and learning in robots and animals](#) at IJCAI in Edinburgh, 30-31st July 1005

<http://www.cs.bham.ac.uk/research/projects/cosy/conferences/edinburgh-05.html>

# Example demos

---

Some 'toy' examples of this design-based approach were shown during the talk.

They included

- The simulated 'harassed nursemaid' having to look after too many 'babies' in an environment presenting various opportunities and dangers
- Two simulated 'emotional' individuals trying to get to their 'targets' and becoming glum, surprised, neutral, or happy depending on what happened in their toy world: these have knowledge of their own states (unlike the nursemaid) and express the state both in a change of facial expression and a verbal report.
- A simulated sheepdog which fetches sheep and herds them into a pen (one at a time) in a world in which its plans can be blocked (e.g. because a tree is moved to block its path, or it or one of the sheep can be forcibly moved to a new location, requiring it to abandon its current plan and form a new one), and in which new opportunities can turn up unexpectedly (e.g. because a barrier that required a long detour suddenly acquires a gap, allowing the dog to use a short-cut). This dog has no anger or frustration when things go wrong, or joy when new opportunities suddenly appear: but it is able to detect new developments and react to them appropriately.

There are movies showing these programs online here

<http://www.cs.bham.ac.uk/research/poplog/figs/simagent/>

However, these are all toy systems: all they do is illustrate the design-based approach (for nursery school learners??).

Anyone who wishes to acquire and play with the software tools can fetch them from here

<http://www.cs.bham.ac.uk/research/poplog/freepoplog.html>

The tools require a linux or solaris system.

Linux running under vmware on a windows or mac pc may work.

# Conceptual confusions: how to make progress

---

- We may be able to come up with clear, useful **design-based concepts** for describing what is happening in a certain class of complex information- processing systems, if we study the architecture, mechanisms and forms of representations used in that type of system, and work out the states and processes that can be generated when the components interact with each other and the environment.
- If the system is one that we had previously encountered and for which we already have a rich and useful pre-scientific vocabulary, then the new design-based concepts will not necessarily **replace** the old ones but may instead **refine** and **extend** them, e.g. making new sub-divisions and bringing out deep similarities between previously apparently different cases.
- This happened to our concepts of physical stuff (air, water, iron, copper, salt, carbon, etc.) as we learnt more about the underlying architecture of matter and the various ways in which the atoms and sub-atomic particles could combine and interact. So we now define water as  $\text{H}_2\text{O}$  and salt as  $\text{NaCl}$ , rather than in terms of how they look, taste, feel, etc., and we know that there are different isotopes of carbon with different numbers of neutrons, e.g.  $\text{C}_{12}$ ,  $\text{C}_{13}$  and  $\text{C}_{14}$ .
- As we increase our understanding of the architecture of mind (what the mechanisms are, how they are combined, how they interact) our concepts of mind (e.g. 'emotion', 'consciousness', 'learning', 'seeing', etc.) will also be refined and extended.

In the meantime, muddle and confusion reign.

# Varieties of definitions of emotion

---

Part of the problem is that many of the words we use for describing human mental states and processes (including 'emotion', 'learning', 'intelligence', 'consciousness') are far too ill-defined to be useful in scientific theories.

Not even professional scientists are close to using agreed definitions of 'emotion'.

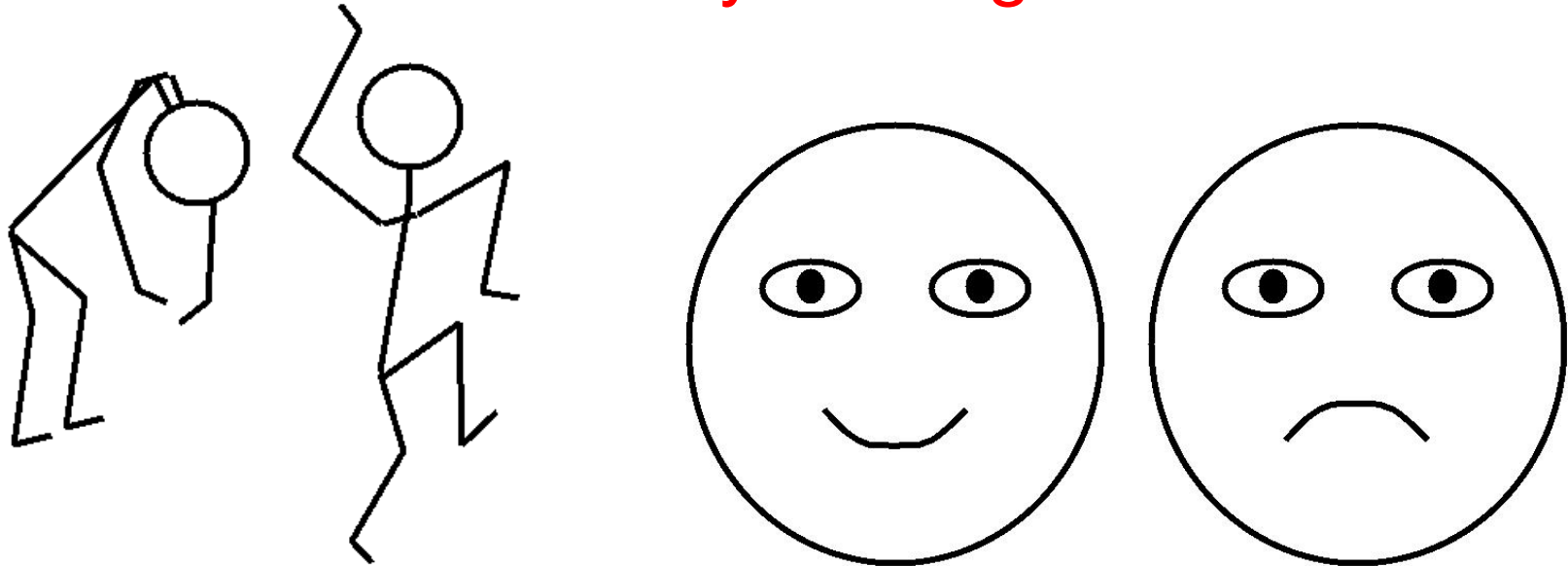
In the psychological literature, for instance, there are attempts to define emotions in terms of

- social relations and interactions between people
  - the kinds of states in the environment that produce them
  - the kinds of behaviours they produce
  - kinds of input-output relations (combining both the above)
  - 'skin-level' and 'sub-skin-level' states and processes, e.g. whether hairs stand on end, galvanic skin responses, blood pressure, muscular tension, etc.
  - the experience of the above bodily changes due to proprioceptive feedback mechanisms (the James/Lange definition, revived in a slightly new form by Damasio's theory of 'somatic markers')
  - which bits of the brain produce them (e.g. amygdala, ...)
  - 'how it feels' to have them
  - how they affect other aspects of mental life
- .....etc.....

All this conceptual confusion and definitional disagreement makes it unclear what question we are asking when we ask whether emotions are needed for intelligence.

# What do we mean by “having an emotion”?

---



- Is it **enough** to produce certain behaviours that people interpret as emotional?
- Do actors actually **have** the states they **portray** so effectively — e.g. despondency, joy, jealousy, hatred, grief...? Not when such states include beliefs and intentions, as despondency, joy, jealousy, hatred, grief etc., often do.
- Behaviour is not enough to define any **mental** state, since
- In principle any behaviour, observed over any time period, can be produced by indefinitely many different mechanisms, using very different internal states and processes. Hence the Turing test is of no use here.
- We need to understand the variety of types of mental states better.  
Then we can define scientific concepts for classifying such states.

# METHODOLOGICAL POINT

---

The concept of **emotion** is but one of a large family of intricately related, but somewhat confused, everyday concepts, including many affective concepts.

E.g. moods, attitudes, desires, dislikes, preferences, values, standards, ideals, intentions, etc., the more enduring of which (along with various skills and knowledge) can be thought of as making up the notion of a “personality”.

Models that purport to account for ‘emotion’ **without accounting for others in the family** are bound to be shallow **though they may have practical applications**.

(See <http://www.cs.bham.ac.uk/research/cogaff/talks/#talk3>)

**A “periodic table” for affective concepts can be based on an architecture, in something like the way the periodic table of elements was based on an architecture for physical matter.**

The analogy is not exact: there are many architectures for minds, each providing its own family of concepts.

**So we need many periodic tables  
generating different sets of concepts.**

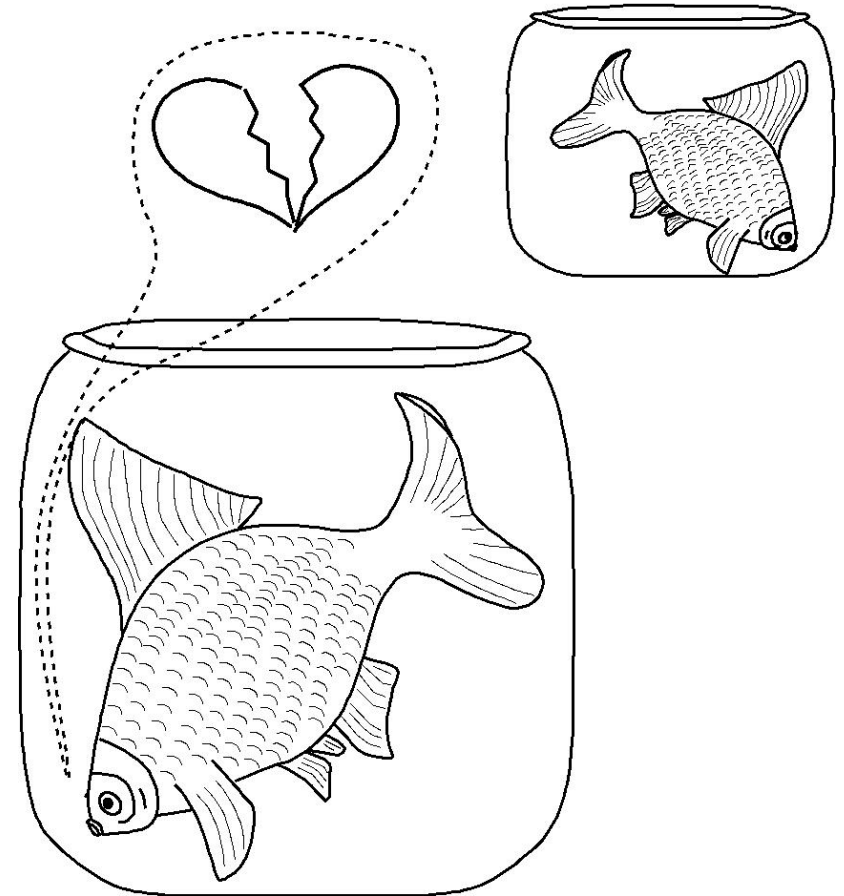
**There may be some concepts applicable across architectures**

# What's wrong with the concepts?

- Everyday concept of 'emotion' mixes up motivations, attitudes, preferences, evaluations, moods, and other affective states and processes.
- There's not even agreement on what sorts of things can have emotions
  - A fly?
  - A woodlouse?
  - A fish?
  - An unborn human foetus?
  - An operating system?
  - A nuclear power plant warning system?

- E.g. some people who argue that emotions are needed for intelligence are merely defending David Hume's truism that **motivation** is needed for action (though not in the case of tornadoes), and **preferences** are needed for selecting between options. Does a tornado select a direction to move in? Does a paramoecium?

WHY CAN'T A GOLDFISH  
LONG FOR ITS MOTHER?



# Towards a general framework

---

We need to talk about “information-using systems” — where “information” has the everyday sense, **not the Shannon technical sense**. This notion is being used increasingly in biology.

## What are information-using systems?

- They acquire, store, manipulate, transform, derive, apply information.
- The information must be expressed or encoded somehow, e.g. in simple or complex structures – possibly in virtual machines.  
(The use of *physical* symbol systems is often too restrictive.)
- These structures may be within the system or in the environment.
- The information may be more or less explicit, or implicit.

A theory of meaning as we normally understand “meaning” in human communication and thinking should be seen as a special case within a general theory of information-using animals and machines.

These ideas are explained in more detail, including the notion of information processing in [virtual machines](#) here:

<http://www.cs.bham.ac.uk/research/projects/cogaff/talks/#models>

<http://www.cs.bham.ac.uk/research/projects/cogaff/talks/#inf>

# Examples of types of processes involving information

- Acquisition
- Filtering/selecting
- Transforming/interpreting/disambiguating
- Compressing/generalising/abstracting
- Deriving (making inferences, but not only using propositions)
- Storage/Retrieval (many forms: exact, pattern-based, fuzzy)
- Training, adaptation (e.g. modifying weights, inducing rules)
- Constructing (e.g. descriptions of new situations or actions)
- Comparing and describing information (meta-information)
- Reorganising (e.g. formation of new ontologies)
- Testing/interrogating (is X in Y, is A above B, what is the P of Q?)
- Copying/replicating
- Syntactic manipulation of information-bearing structures
- Translating between forms, e.g. propositions, diagrams, weights
- Controlling/triggering/modulating behaviour (internal, external)
- Propagating (e.g. in a semantic net, or neural net)
- Transmitting/communicating
- .... (many more)

The differences involve types of content, types of medium used, and the causal and functional relations between the processes and their precursors and successors.

# Control information vs factual information

---

A feature of ordinary language that can confuse discussions of information-processing is that we normally think of information as something that is true or false: e.g. information about when the train will arrive, whereas much information is **control** information which instead of being a potential answer to a question about what is the case is a potential answer to a question about what to do (or not do).

Gilbert Ryle (*The Concept of Mind* 1949) distinguished **knowing that** and **knowing how**, and we could add **knowing what to do, or avoid, or refrain from, or...**

Examples include:

- recipes and instruction manuals
- the ten commandments
- books on etiquette
- commands given by superiors to subordinates
- advice given by parents to children or helpers to friends
- learnt skills that enable us to do things, ....

Control information is more fundamental to intelligent action, or any kind of action, than factual information, since control information can generate action without factual information, whereas the converse is not true (as David Hume and others noted).

Having motives, having preferences, having values, having attitudes, having ideals, having dislikes, all involve control information – in the virtual machines constituting minds – but there's no reason to regard them all as 'emotions'.

# The importance of virtual machines

During the 20th century computer scientists and software engineers came to realise this important truth:

In addition to physical machines, whose components and behaviours can be described using the language of the physical sciences, e.g. physics and chemistry, there are also **virtual** machines whose components and behaviour require a quite different vocabulary for their description, and whose laws of behaviour are not like physical laws.

For more on this see

<http://www.cs.bham.ac.uk/research/cogaff/talks/#inf>

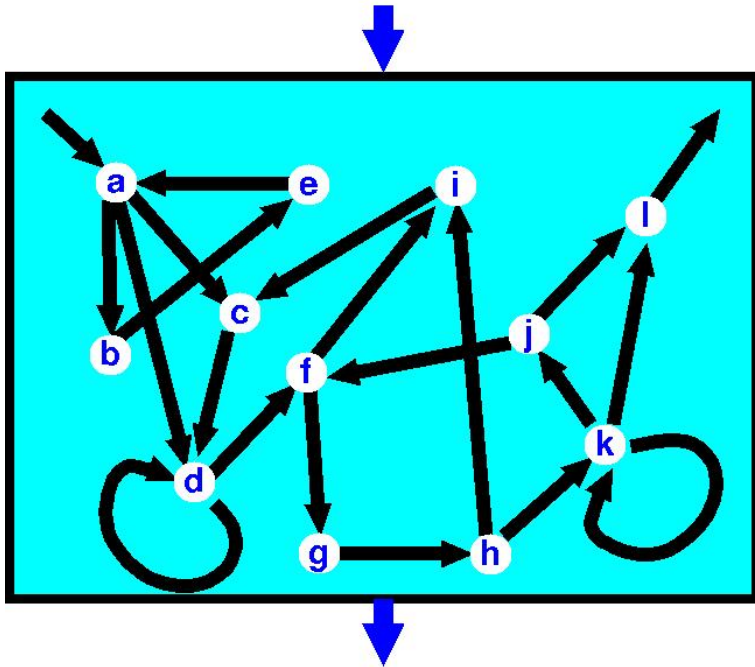
Virtual machines have many advantages over physical machines for various kinds of tasks, e.g.

- They can change their structures without having to rebuild the underlying physical machinery
- They can switch between states containing different structures very quickly, far more quickly than physical structures can be reorganised
  - This is needed for instance when what you see changes rapidly, or while you are having a rapid succession of complex thoughts, e.g. while reading a story or this text.
- Conflicts between inconsistent control processes can be resolved by deliberation and reasoning, instead of being restricted to averaging or vector addition, as is the case with most physical forces pulling in different directions.

It is clear that evolution 'discovered' the benefits of virtual machines long before human scientists and engineers did!

# Functionalism ?

Functionalism is one kind of attempt to understand the notion of virtual machine, in terms of states defined by a state-transition table.



This is how many people think of functionalism: there's a total state which affects input/output contingencies, and each possible state can be defined by how inputs determine next state and outputs.

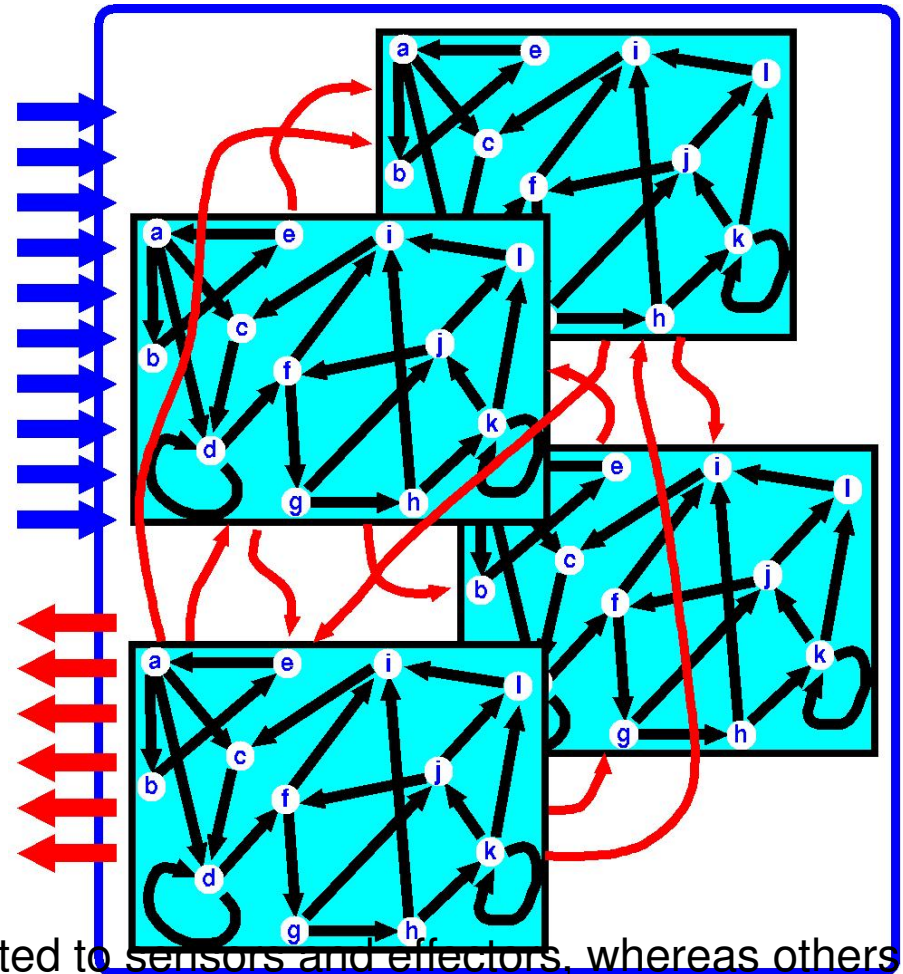
(E.g. see Ned Block's accounts of functionalism.)

HOWEVER THERE'S A RICHER, DEEPER NOTION OF FUNCTIONALISM

# Another kind of Functionalism ?

Instead of a **single** (atomic) state which switches when some input is received, a virtual machine can include **many** sub-systems with their own states and state transitions going on concurrently, some of them providing inputs to others.

- The different states may **change on different time scales**: some change very rapidly others very slowly, if at all.
- They can vary in their **granularity**: some sub-systems may be able to be only in one of a few states, whereas others can switch between vast numbers of possible states (like a computer's virtual memory).
- Some may change **continuously**, others only in **discrete** steps.



Some sub-processes may be **directly** connected to **sensors and effectors**, whereas others have no direct connections to inputs and outputs and may only be affected very **indirectly** by sensors or affect motors only very **indirectly** (if at all!).

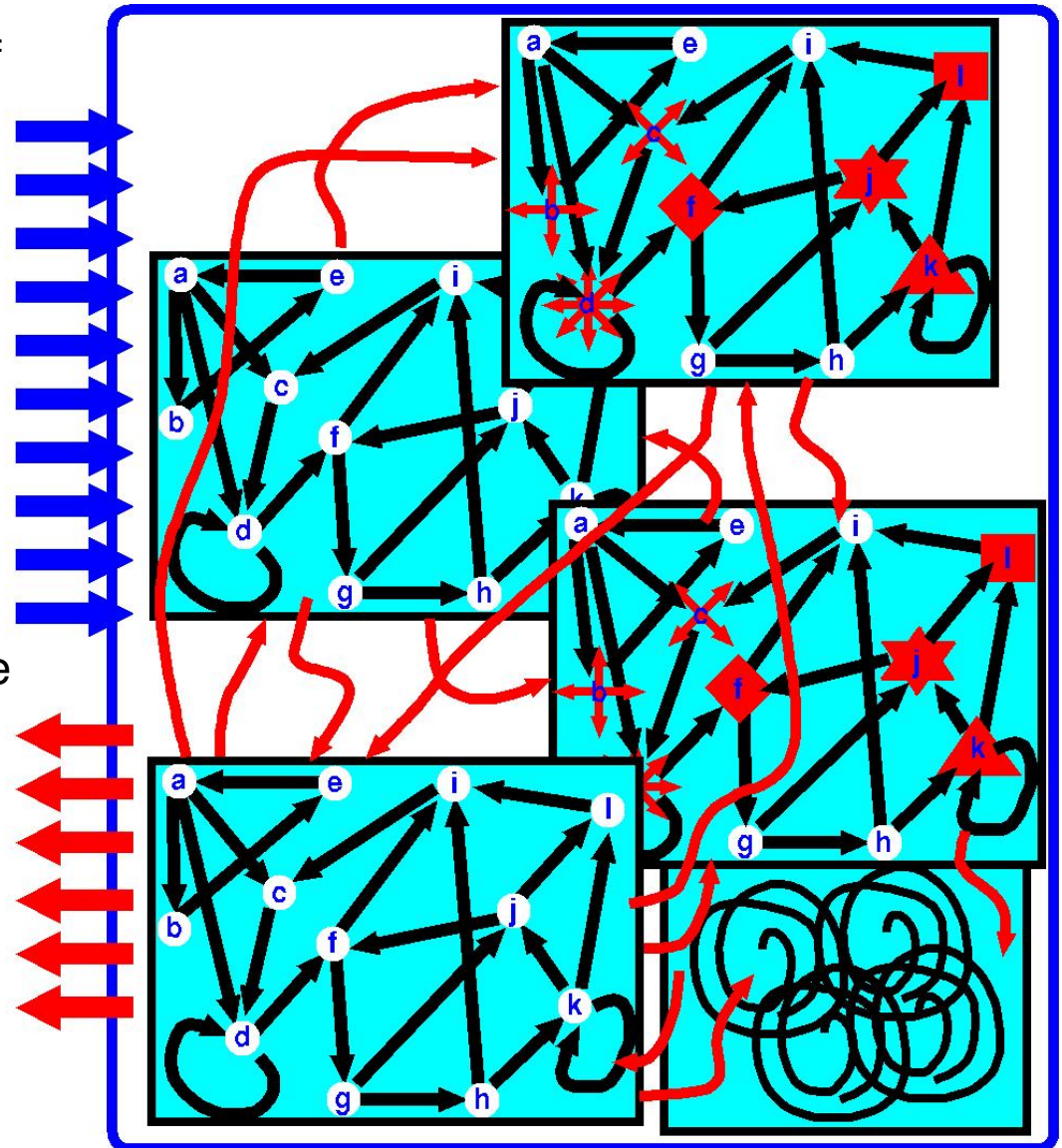
# The previous picture is misleading

Because it suggests that the total state is made up of a **fixed** number of **discretely varying** sub-states:

We also need to allow systems that can grow structures whose complexity varies over time, as crudely indicated on the right, e.g. trees, networks, algorithms, plans, thoughts, etc.

And systems that can change **continuously**, such as many physicists and control engineers have studied for many years, as crudely indicated bottom right e.g. for controlling movements.

The label '**dynamical system**' is trivially applicable to all these types of sub-system and to complex systems composed of them: but it explains nothing.



# VMF: Virtual Machine Functionalism

---

We use “Virtual Machine Functionalism” (VMF) to refer to the more general notion of functionalism, in contrast with “Atomic State Functionalism” (ASF) which is generally concerned with finite state machines that have only **one** state at a time.

- VMF allows multiple concurrently active, interactive, sub-states changing on different time scales (some continuously) with varying complexity.
- VMF also allows that the Input/Output bandwidth of the system with multiple interacting internal states may be too low to reveal everything going on internally.
- There may still be real, causally efficacious, internal virtual machine events and processes that cannot be directly observed and whose effects may not even be **indirectly** manifested externally.

Even opening up the system may not make it easy to observe the VM events and processes (decompiling can be too hard). See

<http://www.cs.bham.ac.uk/research/cogaff/talks/#inf>

<http://www.cs.bham.ac.uk/research/projects/cogaff/talks/#wpe08>

- VMF allows some processes to have the effect of providing control information for others, and for different processes to compete for control.
- If all control is dubbed ‘emotional’ the label becomes vacuous: but it may be useful to recognize some special cases as emotional, namely some of the cases **where where one process disrupts, aborts, suspends, or otherwise interferes with another — e.g. when we are ‘moved’ by something.**

# Our own work in Birmingham

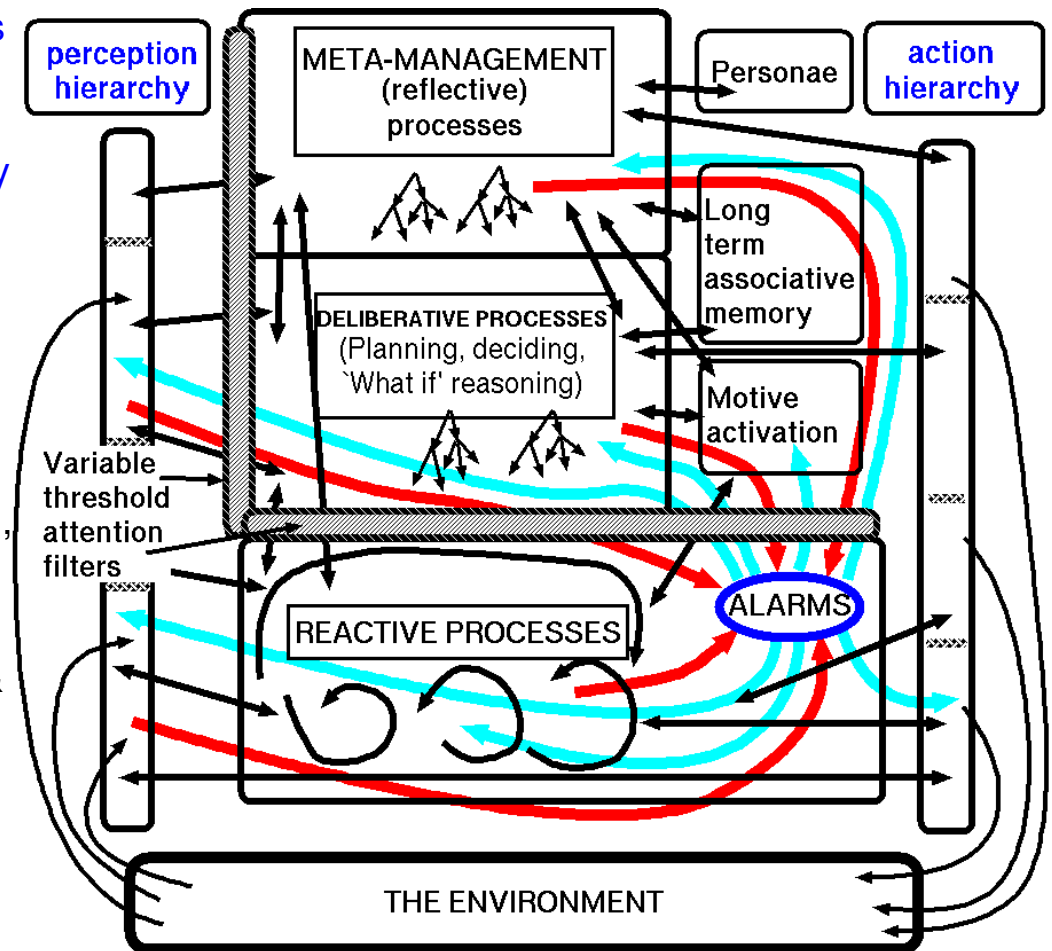
## The architecture of a human mind

(very sketchy draft – see <http://www.cs.bham.ac.uk/research/cogaff/>)

The H-Cogaff (Human Cogaff) architecture is a (conjectured) special case of the CogAff architecture schema, containing many different sorts of concurrently active, mutually interacting components.

It includes 'old' reactive components shared with many other animals (most species are purely reactive) 'newer' deliberative mechanisms (for considering non-existent possibilities) and relatively rare meta-management capabilities for inspecting, evaluating, and influencing internal information-processing.

Papers and presentations on the Cognition & Affect web site give more information about the functional subdivisions in the (still very sketchy) H-Cogaff architecture, and suggest that many familiar kinds states (e.g. several varieties of emotions) could arise in such an architecture, in animals or robots.



See other Cognition and Affect papers and talks for details

# CogAff: A schema for a variety of architectures.

'CogAff' is our label, not for an architecture (like 'H-Cogaff'), but for a way of specifying architectures in terms of which sorts of components they include and how they are connected: H-Cogaff is a special case of the schema.

Think of a grid of **co-evolved** types of **sub-organisms**, each contributing to the niches of the others, each performing different functions, using different mechanisms, etc.

We could add lots of arrows between boxes indicating possible routes for flow of information (including control signals) – in principle, mechanisms in any two boxes can be connected in either direction.

However, not all organisms will have all the kinds of components, or all possible connections.

E.g. insects are purely reactive, and perhaps also all reptiles and fish. A few species have deliberative capabilities in a simple form and perhaps even fewer have meta-management. **Many kinds need "alarm" mechanisms.**

For a survey of varieties of deliberative systems from 'proto-deliberative' to 'fully deliberative' see <http://www.cs.bham.ac.uk/research/projects/cosy/papers/#dp0604>

Perception	Central Processing	Action
	<b>Meta-management (reflective processes) (newest)</b>	
	<b>Deliberative reasoning ("what if" mechanisms) (older)</b>	
	<b>Reactive mechanisms (oldest)</b>	

# As processing grows more sophisticated, so it can become slower, to the point of danger

REMEDY: FAST, POWERFUL, “GLOBAL ALARM SYSTEMS”

Resource-limited alarm mechanisms must use fast pattern-recognition and will therefore inevitably be stupid, and capable of error!

Many variants are possible. E.g. purely innate, or trainable.

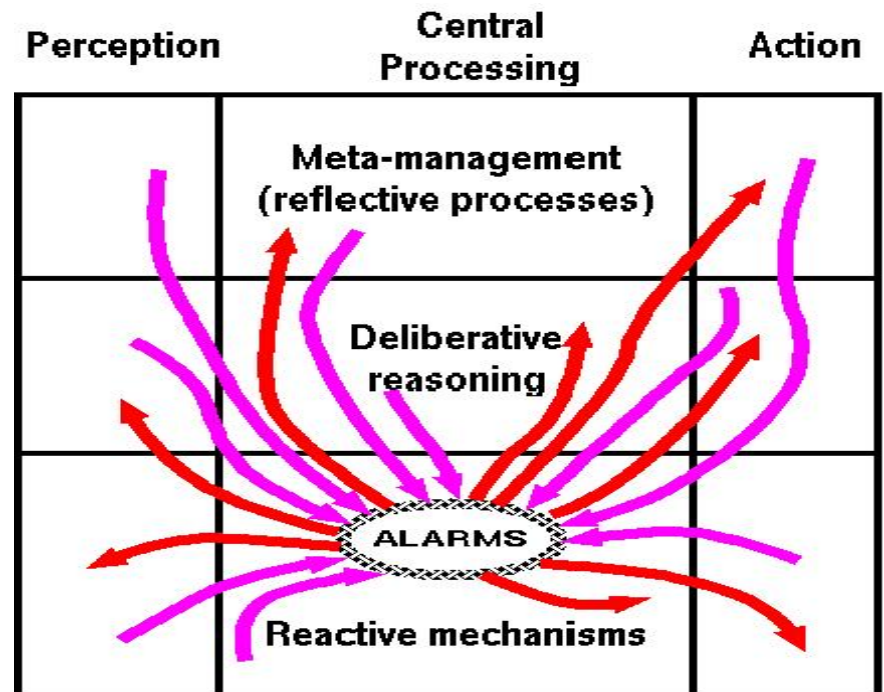
E.g. one alarm system or several?  
(Brain stem, limbic system, ...???)

See Cogaff papers and talks

<http://www.cs.bham.ac.uk/research/cogaff/>

[http:](http://)

[//www.cs.bham.ac.uk/research/cogaff/talks/](http://www.cs.bham.ac.uk/research/cogaff/talks/)



Many different kinds of emotional state can be based on such an alarm system, depending on what else is in the architecture.

Don't confuse the alarms (and emotions they produce) with the evaluations that trigger them, or the motives, preferences, policies, values, attitudes that have different sorts of functional roles – different sorts of control functions (including conditional control in many cases).

# Emotions and control mechanisms

---

What is there in common between

- a crawling woodlouse that rapidly curls up if suddenly tapped with a pencil,
- a fly on the table that rapidly flies off when a swatter approaches,
- a fox squealing and struggling to escape from the trap that has clamped its leg,
- a child suddenly terrified by a large object rushing towards it,
- a person who is startled by a moving shadow when walking in a dark passageway,
- a rejected lover unable to put the humiliation out of mind
- a mathematician upset on realising that a proof of a hard theorem is fallacious,
- a grieving parent, suddenly remembering the lost child while in the middle of some important task?

Proposed Answer:

in all cases there are at least two sub-systems at work in the organism, and one or more specialised sub-systems, somehow interrupt or suppress or change the behaviour of others, producing some alteration in (relatively) global (internal or external) behaviour of the system — which could be in a virtual machine.

Some people would wish to emphasise a role for *evaluation*: the interruption is based at least in part on an assessment of the situation as good or bad.

Is a fly capable of evaluation? Can it have emotions? [Evaluations are another bag of worms.](#)

Some such 'emotional' states are useful, others not: they are not required for all kinds of intelligence — only in a **subset** of cases where the system is too slow or too uninformed to decide intelligently what to do — they can often be disastrous!

# Emotions are a subclass of “affective” states

---

Affective states are of many kinds. They include not only what we ordinarily call emotions but also states involving desires, pleasures, pains, goals, values, ideals, attitudes, preferences, and moods.

The general notion of “affective state” is very hard to define but very roughly it involves using some kind of information that is compared (explicitly or implicitly) against what is happening, sensed either internally or externally.

- When there’s a discrepancy some action is taken, or tends to be taken to remove the discrepancy by acting on the sensed thing: affective states involve a *disposition* to change reality in some way to reduce a mismatch, or preserve a match.
- In contrast, if the information is part of a percept or a belief, then detecting a discrepancy tends to produce a change in the stored “reference” information.

**The two cases differ in what has been labelled “direction of fit”.**

Without **affect** there is no reason to do anything.

Affect: whatever initiates, preserves, prevents, selects between, modulates, actions.

So I am NOT arguing that knowledge and powers of reasoning suffice for behaving intelligently: in particular, **without motivation** nothing will be done (except in purely reactive systems).

Hume: **Reason is and ought to be the slave of the passions.**

NOTE: Some affective states are derivative on others

(e.g. wanting X because it is conducive to, or prevents or preserves Y, etc.)

# This is just the beginning

---

- I have tried to give some of the flavour of the kind of thinking involved in the design-based approach to thinking about minds of humans, other animals or machines.
- When we start investigating what could happen in an architecture as rich as H-Cogaff (which is still much simpler than a normal adult human architecture) we see that many more kinds of states and processes are possible than we have convenient labels for.
- So we can start classifying in a more precise way than ever before various classes of states and processes.
- We'll see that a subset of the things we call being in an emotional state (e.g. being startled, frightened of a cliff-edge, joyful at recognition of a loved one) may involve operations of something like the 'alarm' mechanism, though not all cases will be alike.
- Some of the long-term cognitively rich emotions including grief or jealousy may not depend on alarm mechanisms, likewise many attitudes often confused with emotions, e.g. dedication to one's job, love of one's family or country.

## The periodic table of human mental states still has far to grow.

The ideas sketched here are a development of ideas that can be found in

H. A. Simon, (1967) *Motivational and emotional controls of cognition*, Reprinted in *Models of Thought*, Yale University Press, 29–38, 1979

# For more on this approach SEE THE COGNITION AND AFFECT PROJECT AND THE COSY ROBOT PROJECT

---

## OVERVIEW, INCLUDING PAPERS & DISCUSSION NOTES:

[http:](http://www.cs.bham.ac.uk/research/projects/cogaff/misc/emotions-questions.html)

[//www.cs.bham.ac.uk/research/projects/cogaff/misc/emotions-questions.html](http://www.cs.bham.ac.uk/research/projects/cogaff/misc/emotions-questions.html)

<http://www.cs.bham.ac.uk/research/projects/cogaff/>

<http://www.cs.bham.ac.uk/research/projects/cosy/papers/>

(References to other work can be found in papers in both directories)

## TOOLS:

<http://www.cs.bham.ac.uk/research/poplog/freepoplog.html>

<http://www.cs.bham.ac.uk/research/projects/poplog/packages/simagent.html>

(the SIM\_AGENT toolkit)

## DEMO-MOVIES:

<http://www.cs.bham.ac.uk/research/poplog/figs/simagent/>

## SLIDES FOR TALKS:

<http://www.cs.bham.ac.uk/research/projects/cogaff/talks/>

<http://www.cs.bham.ac.uk/research/projects/cosy/papers/#presentations>

(Including several on emotions)

## ALSO STRONGLY RECOMMENDED:

Barrett, L. F. (2006). Emotions as natural kinds? *Perspectives on Psychological Science*, 1, 28-58.

<http://www2.bc.edu/~barretli/pubs/2006/Barrett2006kinds.pdf>

# THANKS

---

I am very grateful to  
the developers of Linux  
and other free, open-source,  
platform-independent, software systems.

LaTeX was used to produce these slides. (I should switch to LaTeX Beamer).

Diagrams are created using `tgif`, freely available from

`http://bourbon.cs.umd.edu:8001/tgif/`

Demos built using `Poplog`

`http://www.cs.bham.ac.uk/research/poplog/freepoplog.html`