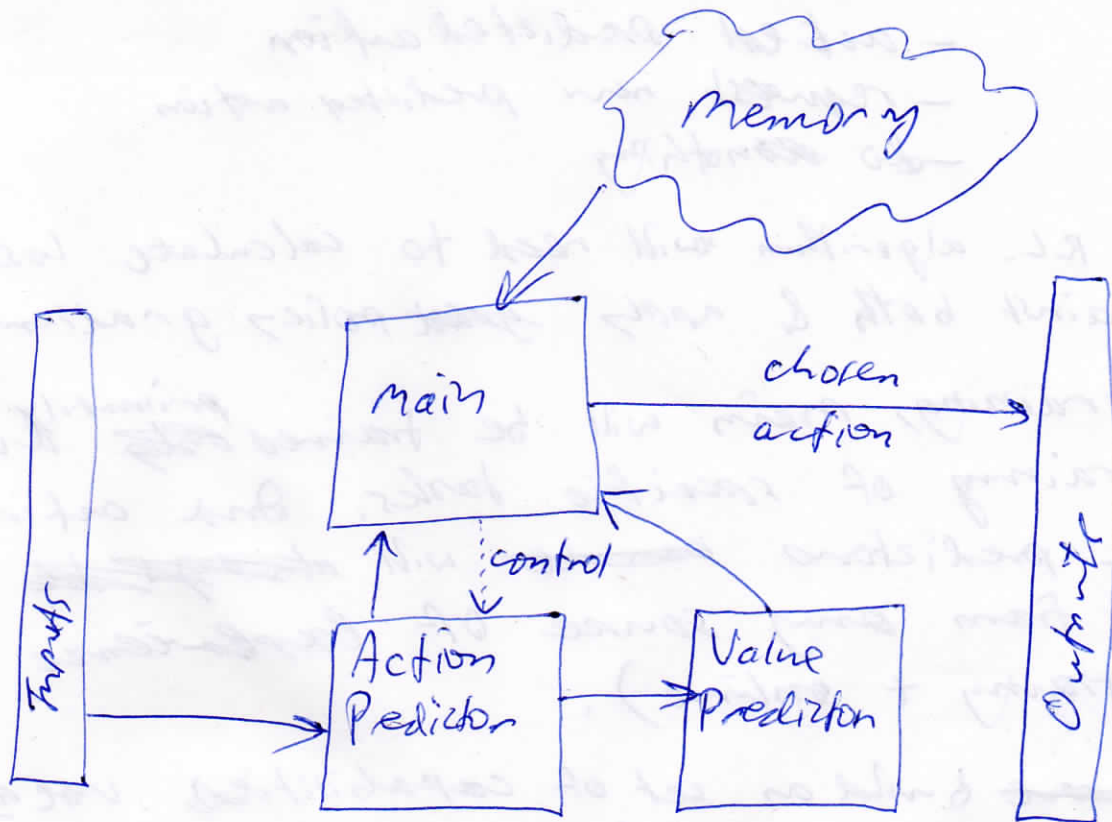


# Deconstructed R.L. Agent as basic for AGI



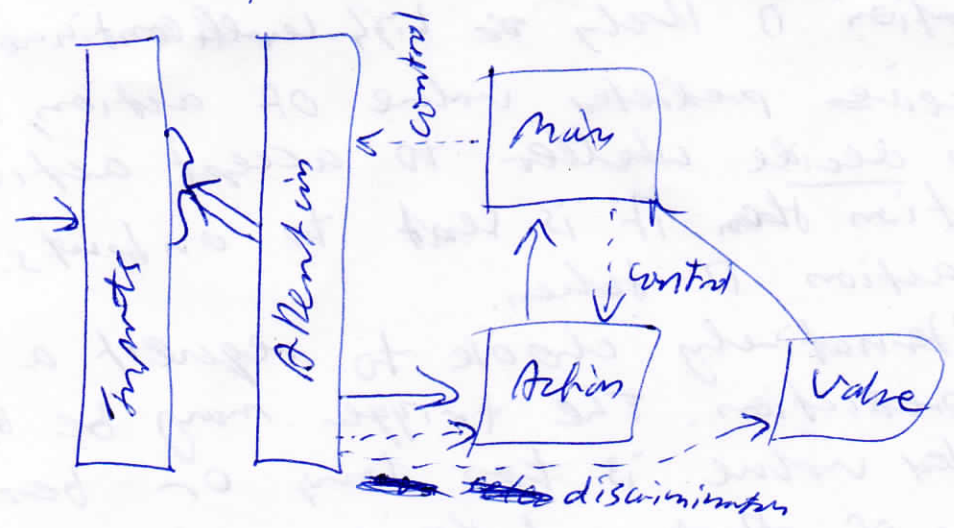
- main receives a suggested action from Action predictor. Action representation is likely a high-level continuous.
- main also receives predicted value of action, and uses that to decide whether to accept action. If it accepts action then it is sent to outputs. Otherwise no action is taken.
- main can alternatively choose to request a new action prediction. The trigger may be because the ~~exp~~ predicted value is too low, or because of ~~too~~ memory of ~~that~~ ~~and~~ the current predicted action already being attempted & not working. The action predictor could be used in a continuous stochastic way, and thus merely requesting a new prediction could produce a new action.
- In a more advanced form, some sort of inhibition of the previous prediction(s) would be applied.

2/ main now has actions itself. These are simple categorized actions:

- accept predicted action
- request new predicted action
- do something.

So RL algorithm will need to calculate loss against both & apply ~~just~~ policy gradients.

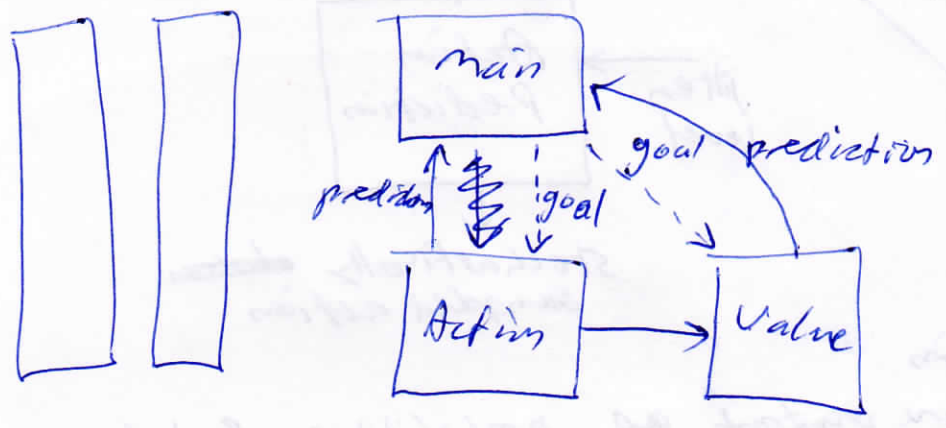
- For training, main will be trained <sup>primarily</sup> ~~only~~ during pre-training of specific tasks. And actions & value predictions ~~can be~~ will ~~always be~~ be trained from any source of experience (pre-training + online).
- To ~~more~~ build on set of capabilities, use attention to act as a ~~selector~~ discriminator or variable that selects particular ~~traits~~ sub-networks.



- Action & Value networks have 'input + attention' as inputs so that they can ~~produce~~ operate against different scenarios. Ideally needs catastrophic forgetting protection.
- main can control attention.

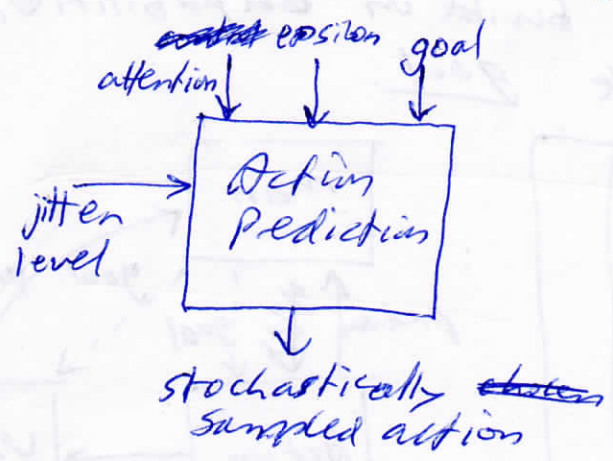
• ~~And~~

- To further build on capabilities, main needs to control the goal.



- Main can now learn its own higher-level strategies for driving this system. As new capabilities are added it will have more options for the sources of predicted actions (eg: multiple action prediction blocks, or choose different action groups through goal ~~set~~ selection). And it will develop different strategies for explaining possibilities.
- Ideally also need a modelling capability to build up novel strategies, to build theories of external systems to aid in action/value prediction, and to discover & control its own learning agenda.

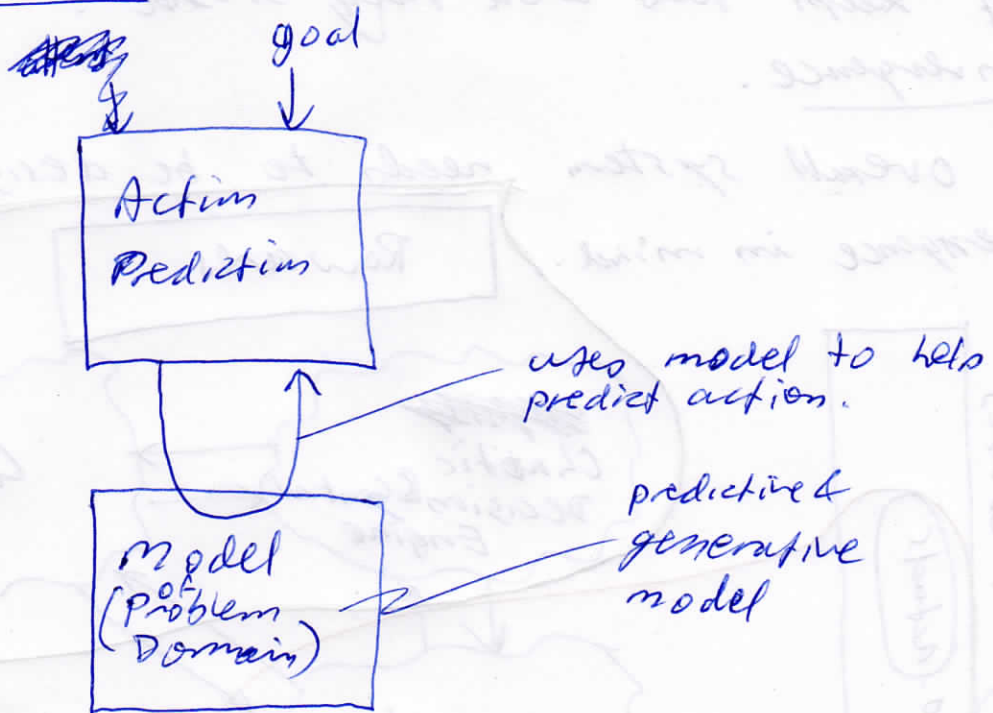
4/ • Emulation baby "jitter learning" & stochastic experimentation:



~~Let main~~

- Treat min. output as probabilities & take a stochastic sample.
- Let main ~~not~~ directly control experimentation see pg 7. no. exploit through control of epsilon.  
(choose completely random action out out with probability  $\epsilon$ )  
or better yet: through the level of noise (std. dev. increased according to ~~the~~ level of requested noise).  
→ this makes the agent occasionally make small mistakes when it knows what it's doing, but can also experiment widely when it doesn't know.
- Lastly, human babies learn <sup>initial</sup> models of their limbs through enforced jitter, that acts as a source of initial ~~experimentation~~ learning.

- Combine with model.

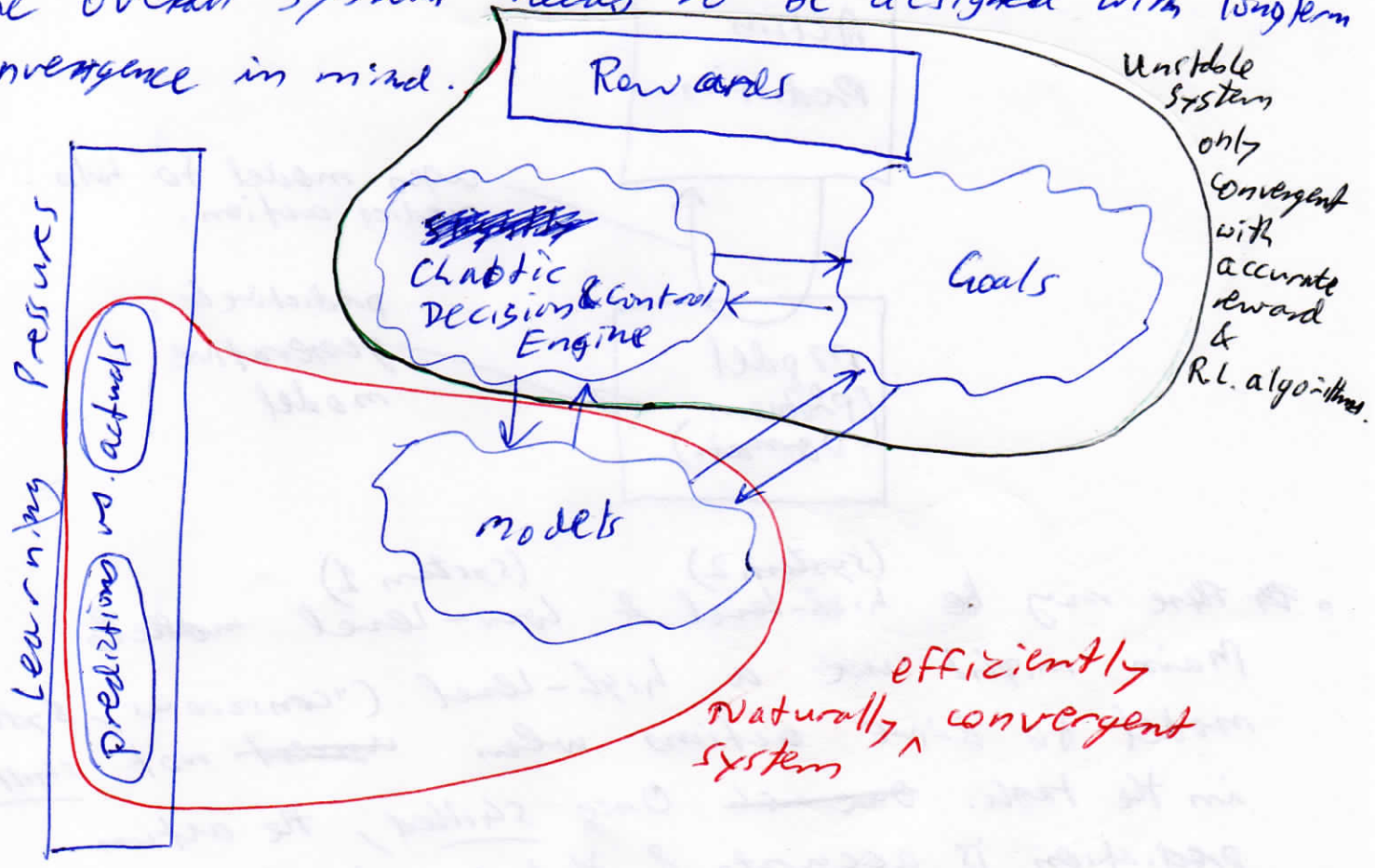


- There may be <sup>(system 2)</sup> high-level & <sup>(system 1)</sup> low-level models. Brain might use a high-level ("conscious", system 2) model to drive actions when ~~unskilled~~ not skilled in the task. ~~Once sk~~ Once skilled, the action prediction is accurate & that can be blindly trusted ~~via~~ as system 1 thought, without using ~~the~~ a system 2 model.

6/ What keeps this whole thing stable?

A: Convergence.

The overall system needs to be designed with long term convergence in mind.

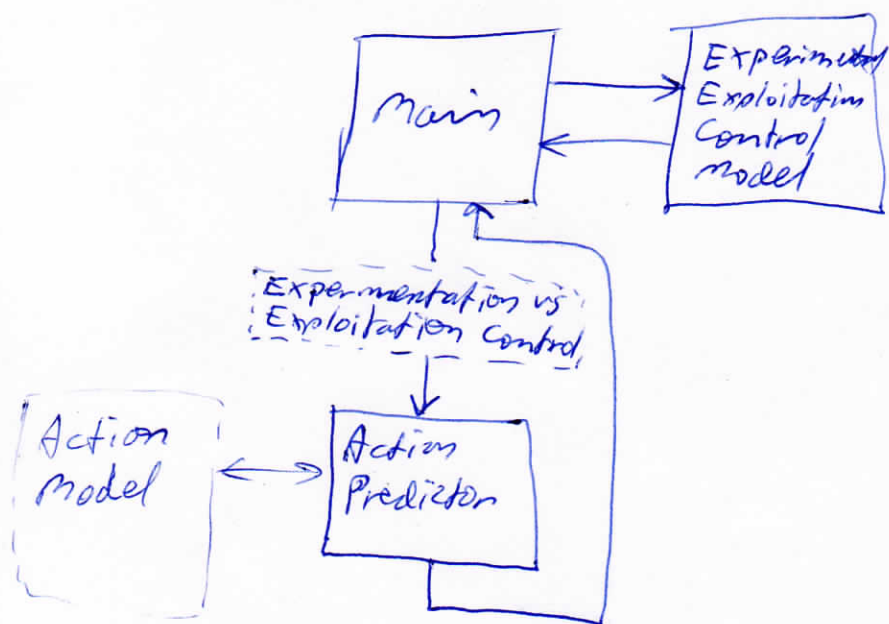


The naturally convergent model + prediction + actuals system is really important to keep the main decision engine in check

~~Making Self-controlled Learning~~

Self-controlled Epsilon/Noise Level (see p 4).

To make this work, can't just rely on implicit pre-trained handling. Need explicit online handling of experimentation vs. exploitation control. It needs to learn a model that represents its benefits & behaviours.



Learning to use this model becomes one of many "strategies" that the agent can employ.

The E.E. control model will itself need to be learned over time. It is probably a bayesian inference model. And some sort of bootstrapping system will be needed to effectively bypass executive control during initial agent training when the E.E. control model is too inaccurate to use.