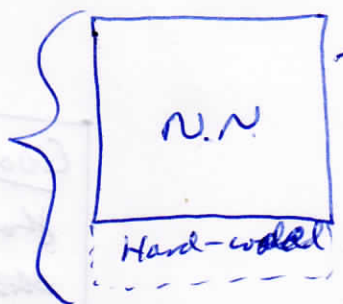


Internalised Reward Function

Value Functions:

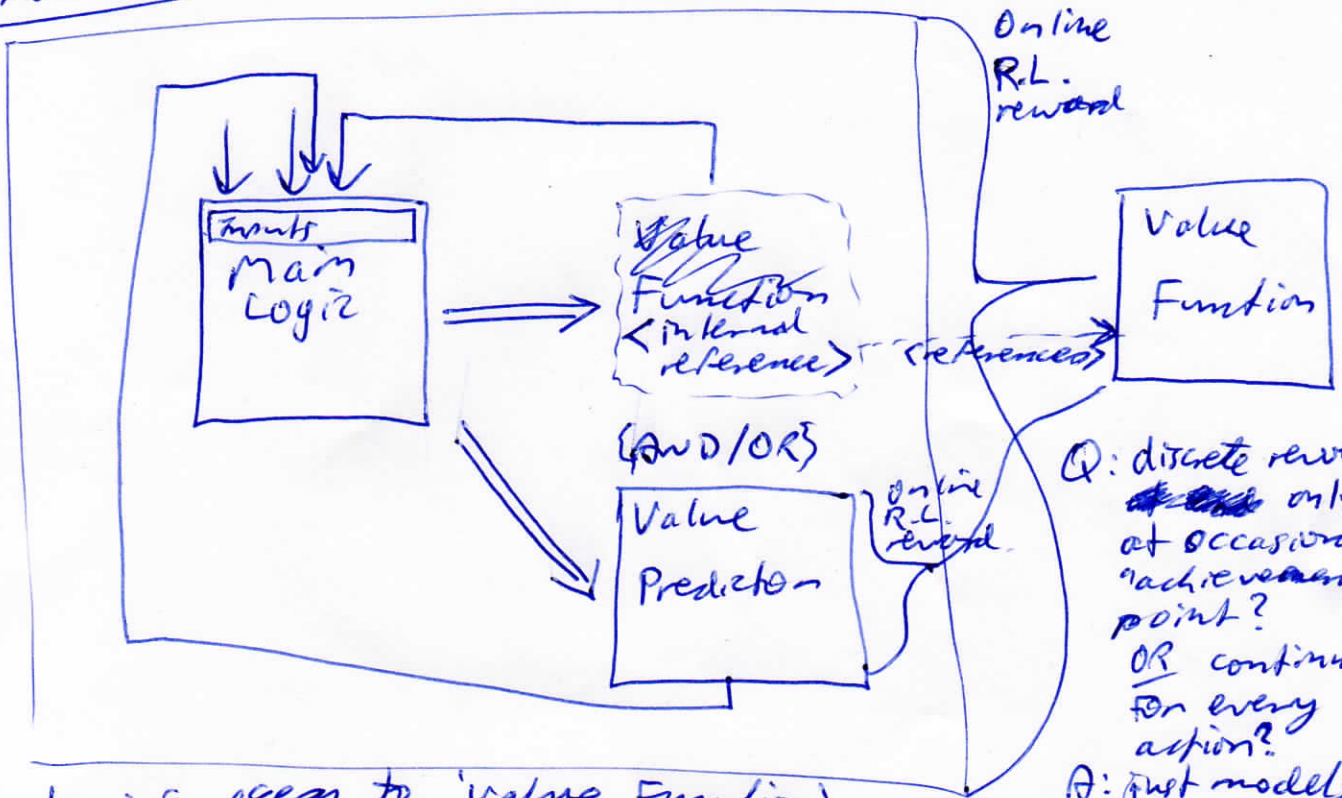


initially trained simply to predict some value as reward from pre-training.

optional "evolved" hardwired value.
eg: that ~~discourages~~ discourages too much ~~eff~~ physical or mental effort → act to improve efficiency.

based on simple hard-wired measures such as number of ~~main~~ main loop iterations.

Usage architecture



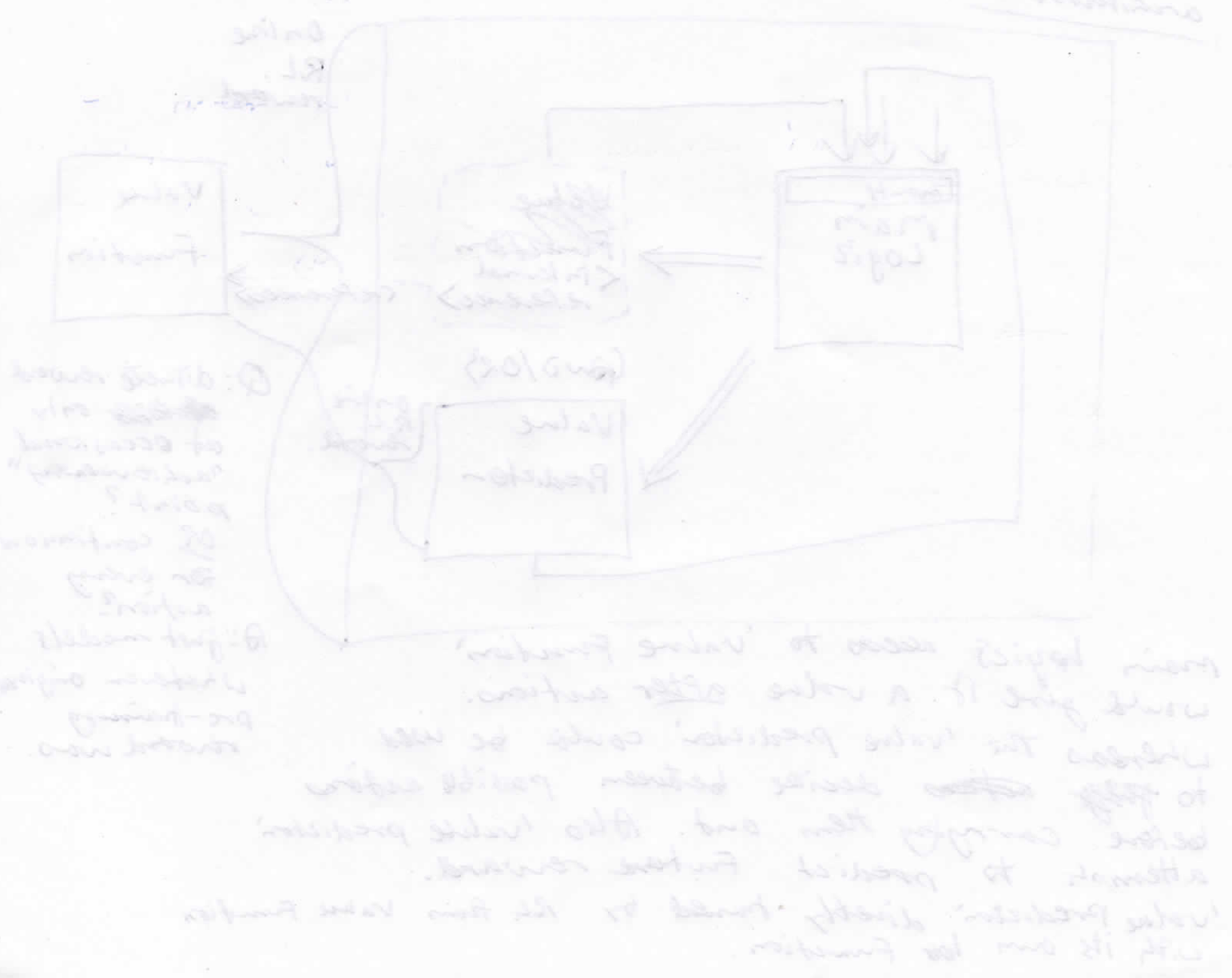
Q: discrete reward ~~only~~ only at occasional "achievement" point?
OR continuous for every action?

A: just models whatever original pre-training reward was.

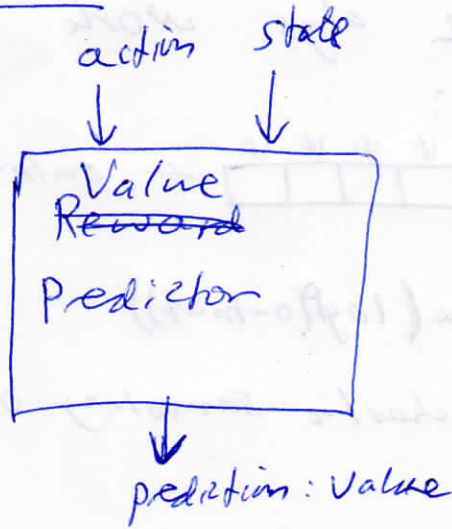
main logic's access to 'value Function' would give it a value after actions.
Whereas the 'value predictor' could be used to ~~judge~~ ~~actions~~ decide between possible actions before carrying them out. Also 'value predictor' attempts to predict future reward.
'value predictor' directly trained by RL from value function with its own loss function.

2/ Q: what can provide continual training to Value Function? without that, there is no hope to learn new skills?

Goals: an agent that can intentionally ~~choose~~ take a course of actions that leads to it learning something from these actions.



Learnable Predictors



NN. representation



Predict likelihood of achieving goal ~~given a particular action~~ if taking a particular action from state.

~~It~~ works off an AI's prior chosen goal, as represented in whatever way it uses.

4) Continuous ~~Act~~ vs. Categorized Actions

PPO and other RL algs work best of categorized actions:



action: $\text{argmax}(\log P(\text{output}))$

And RL uses stochastic sampling of that probability distribution.

My AI needs to output a "model" rather than an action selection.



action: output.

In order to apply stochastic sampling, treat each output node as an independent probability distribution (node value \Rightarrow mean?) and sample from each.

Why ~~from~~ distribution sampling so important?
Because it provides the source of "novel" attempts that may be better than the policy's current preference.