Executive
control
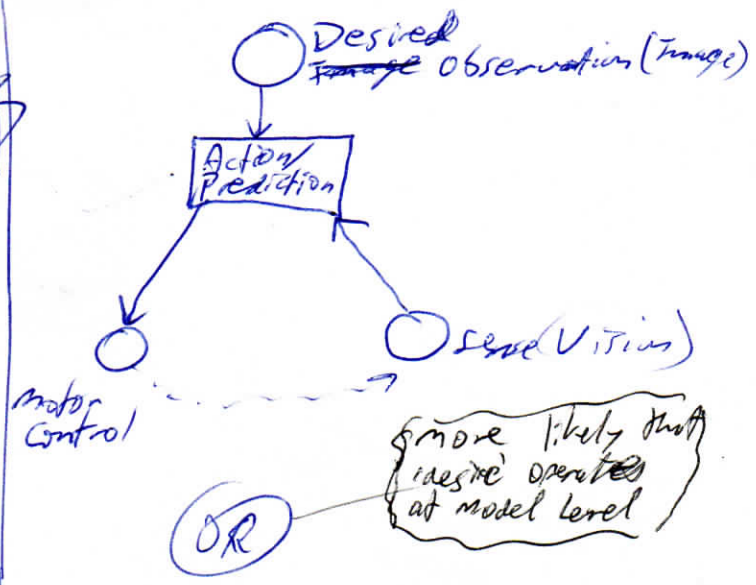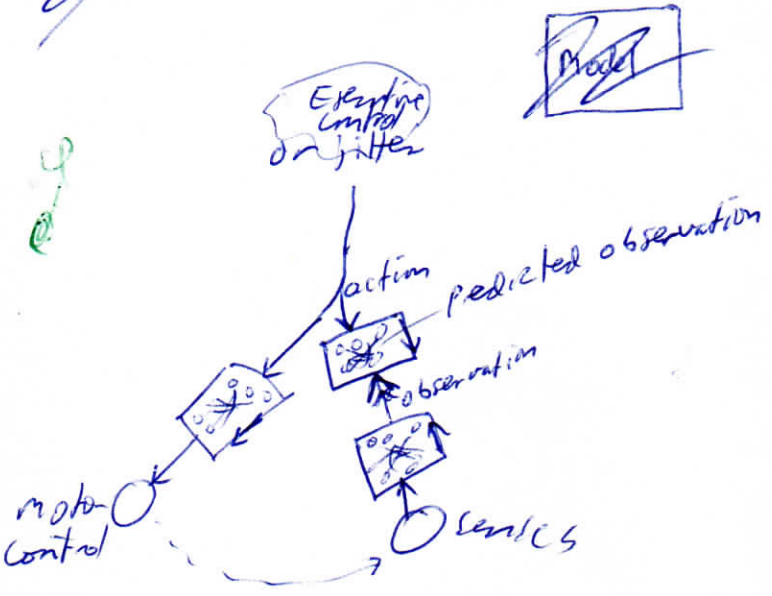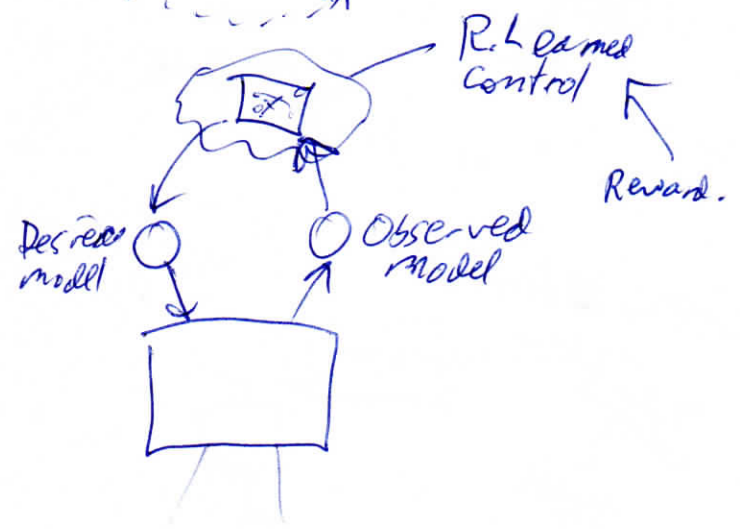or jitter

Model

action    predicted observation

observation

motor
control

Senses

---

See arm and desire to see
it move to particular location
(in order to pick something up).
(ignore executive control side of it
for now).

Desired
~~Image~~ Observation (Image)

Action/
Prediction

motor
Control

Issue (Vision)

OR

more likely that
desire operates
at model level

Desired
Observation (High-level
model)

Action/
modelling/
prediction

Observation
(model)

R.Learner
Control

Reward.

Desired
model

Observed
Model

---

Needs jitter to bootstrap
& lead to convergence
(stability).

Jitter 2 ~~Mots:~~ Random signals
to Desired Model.

Jitter 1: Random motor control
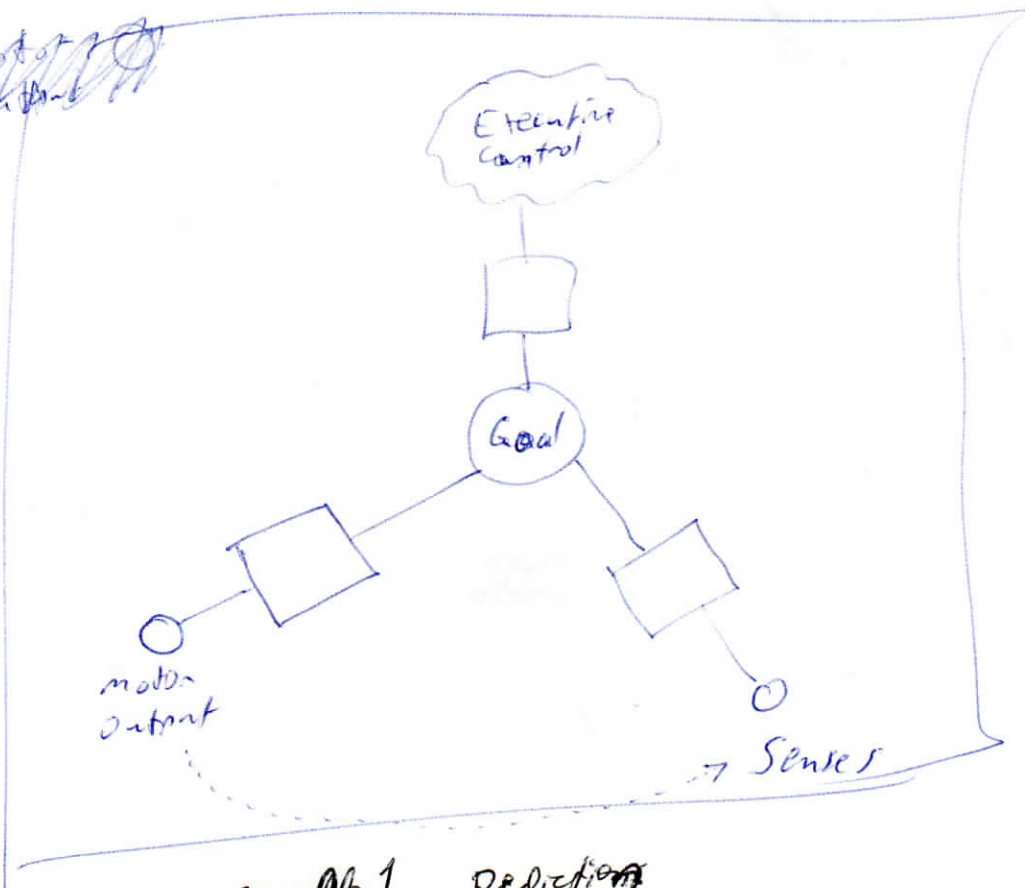signals. Use touch as
shared sense that builds up
model.

(a la Daniel Kahneman)

Executive Control

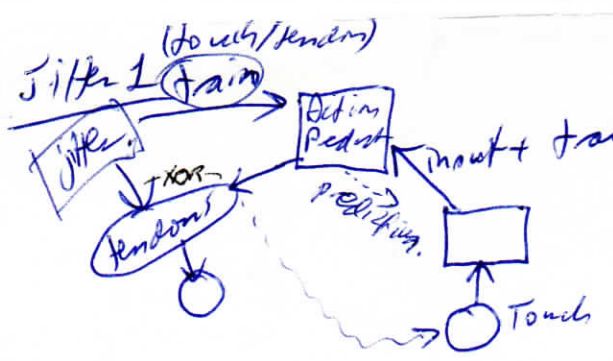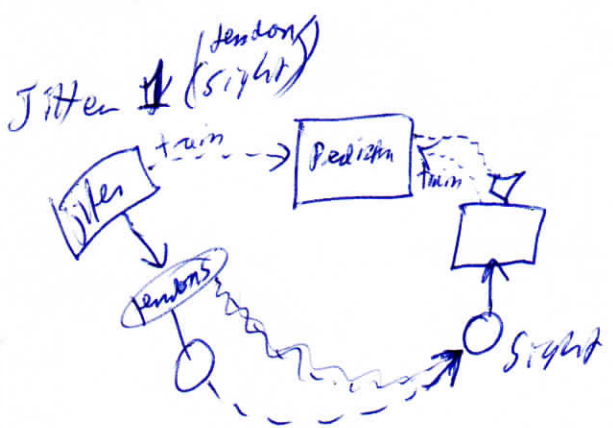Goal

motor output

→ Senses

Unconscious **Model 1** Prediction

Jitter 1 (touch/tendon) train
Jitter
XOR
tendons
Action Pedict
input + train (expected)
predicting.
Touch

source: jitter on control
tendon
action
(motor control)
Prediction
→ touch result
(predicted v.s. actual)

Jitter 1 (tendon)(sight)
Jitter → train → Prediction train
tendons
Sight

tendon action → Prediction → sight result.

Now somehow layer this up. This prediction becomes input to jitter 2 circuit.

Ctrl

Tend/Obj

Motor

dff.
(@Raw)

Touch/
surf

(@model)
diff.

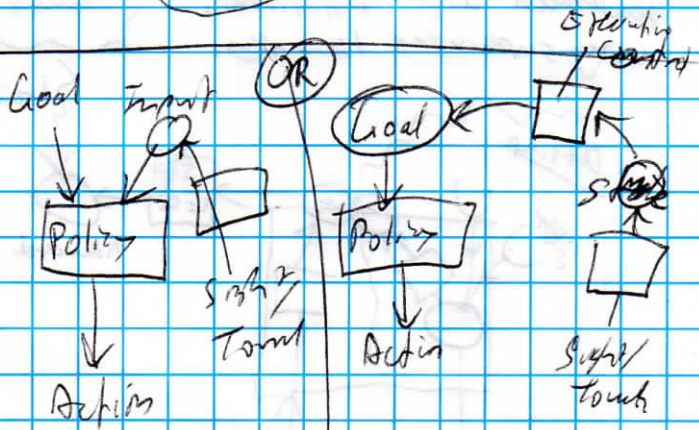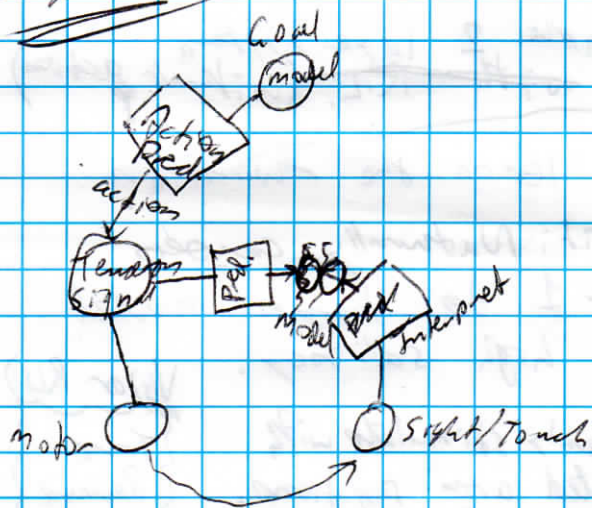Basically an adversarial model

auto-tran
loss Fn

Risk of learning evolved soln: always blank. Evolutionary pressure will maximise contrast/ salience/utility: w.N and other pressure/mechanics.
eg: some sort of white-balance or normalisation.

4) Benefit of this layering approach is that it enforces convergence. The low-level system is inherently convergent, so it counteracts any chaos from higher layers.

# Layer 2

Goal
Model

Action
Pred

action

Tenags
sitn

Pedi
Model
BK
interpret

motor

Sight/Touch

---

Goal      Input      OR      Goal      Executin
Control

Policy

Policy

Sight/
Touch

Sight/
Touch

Action

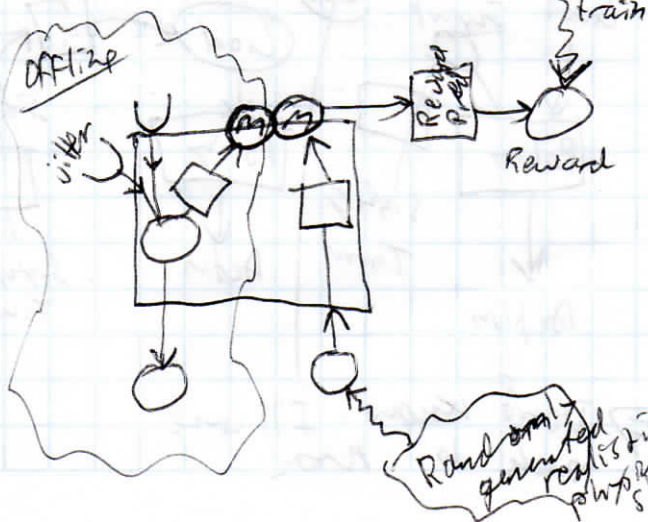Action

Sight/
Touch

⇒ Don't know. I'll ane
back to this.

6/5 Simple 2 layer system with~~ R.L. (without action)

Trick: learn the reward fn

Benefit: Naturally causes layer 1 to learn model with high saliency.

Randomly stimulate with generated arm positions. using supervised learning.



$V(s)$ or $R(s)$
True Reward

↓train

Reward

OFFLINE

jitter

Randomly generated realistic photoreal states

Result:

Learns a model representation with some utility, but biased towards just predicting return.

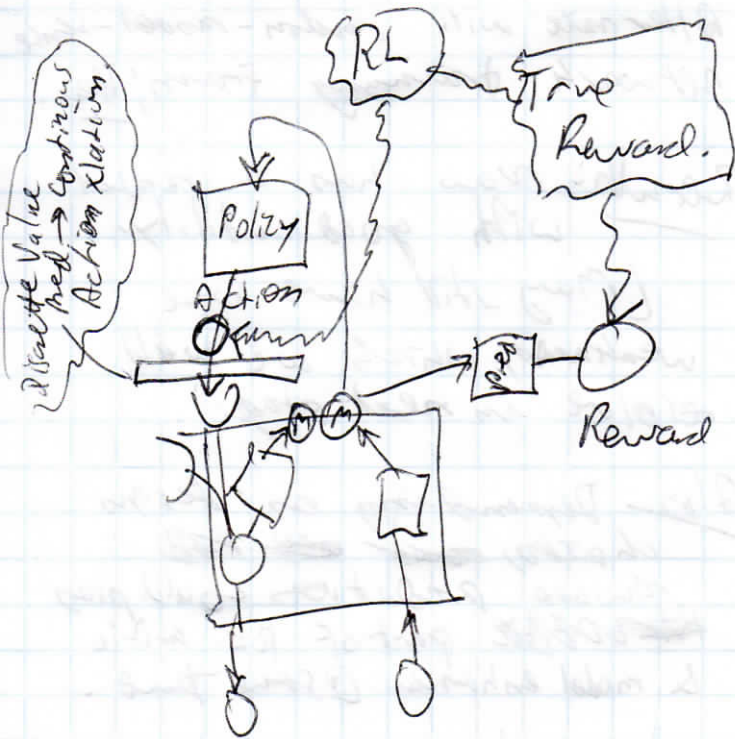Alternate with tendon-model-since network training from jitter.

Result: Now has a model with good utility.

(May still have some weakness, which we will resolve in next stage).

Also: Depending on design choices, ~~could~~ ~~add~~ ~~of~~ reward prediction could play R.L. ~~robbie~~ part of R.L. critic & model enhancer @ same time.

Alternatively, go straight to full RL (as per next page). But I the idea of internalised reward.

Simple 2 layer with action \8
Policy & RL.



RL

True Reward.

Policy

Action

Discrete Value → Continuous
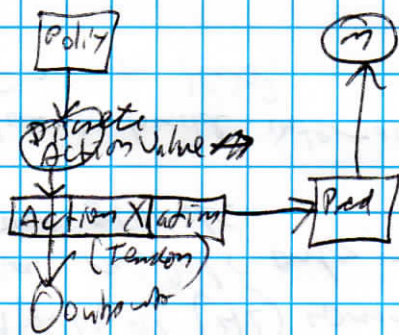Pred → Continuous
Action xlation

Reward

Result:

model with full utility.

Q: how does continuous action
policy learning work.

However, maybe we can work
with any solution, in order to
train tendon → model prediction.



Next: layer 3 produces the
goal ← and this is what
we are conscious of.
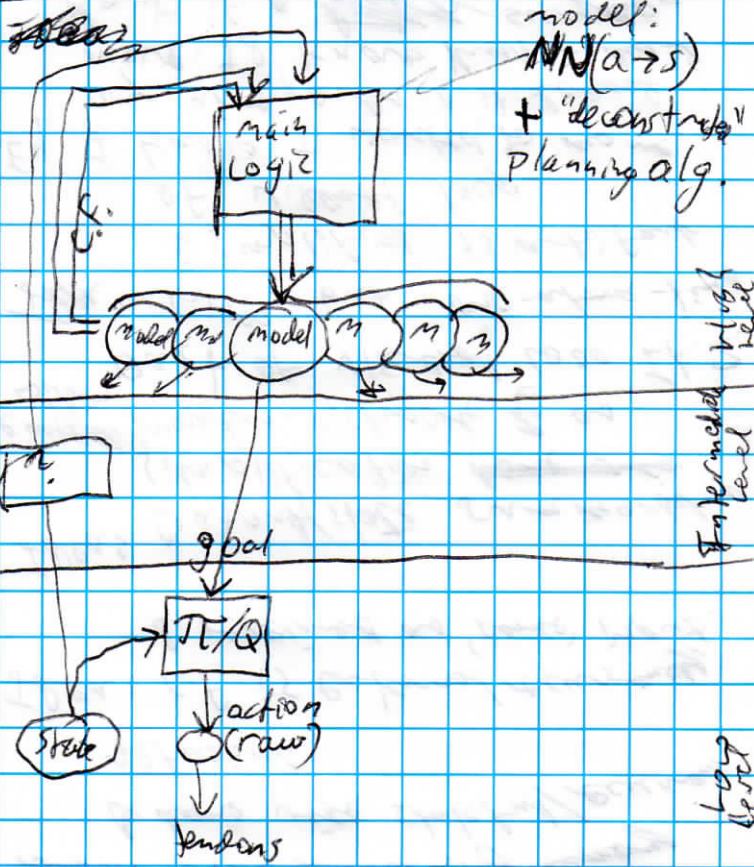
# For Experiment below

Idea:

High Level (HL) communicates
via models only — no low level
raw signals.

Q: How to build internalised reward
fn into this?

Future addition:
Add "suggestion" prediction block
NN(s, g → a) as input that learns
best approach from past experience
& optimises action search space.

Experiment — from High level
ideas

Main
Logic

model: model model M M M 3

model: NN(a → s)
+ "deconstruct"
Planning alg.

"High level"

"Intermediate level"

goal

π/Q

State

action
(raw)

tendons

"low level"

# For below

Idea: state ≡ working memory
& done via stateful/recurrent
network.
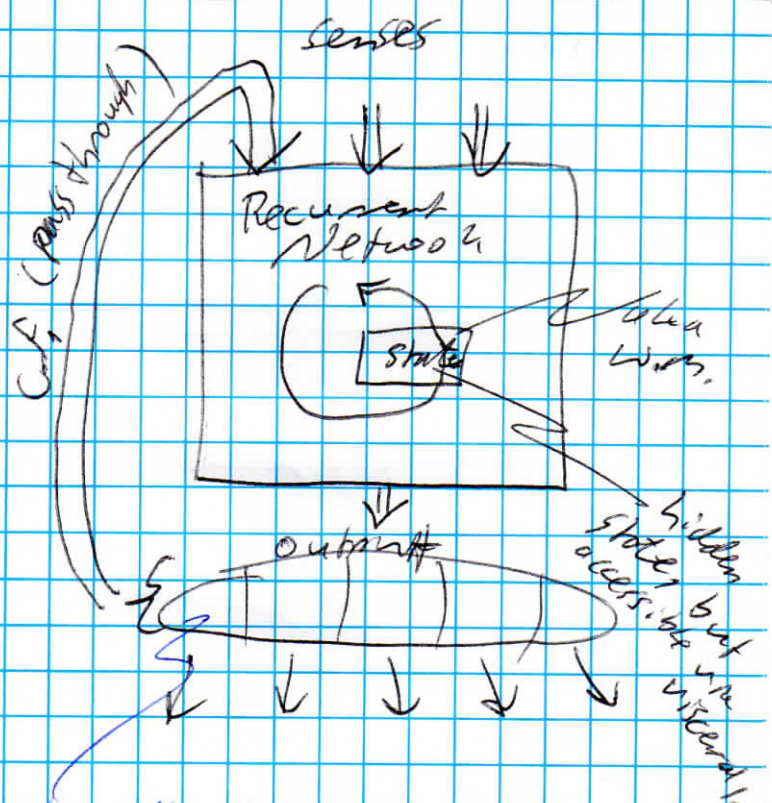
Idea: c.f. is external recurrency
& received as 'sense' input.

Idea: Output state summarising
simplification ~~done via~~
embedded/within meta network & as
part of visceral loop I.t.2.

Idea: W.M. but one-at-a-time
& simplified is artifact
of visceral loop.

Eg: I know I wanted to move
my arm, even if it doesn't
move. I'd know that I didn't
do it if a ~~tested~~ surgeon
moved my tendon.

# Review High Level State + C.F.

Senses

C.F. (pass through)

Recurrent Network

state

idea worm.

output

hidden states but accessible via internal or.

Physically structured output array sections allocated & diverted to specific targets.

~~Important~~
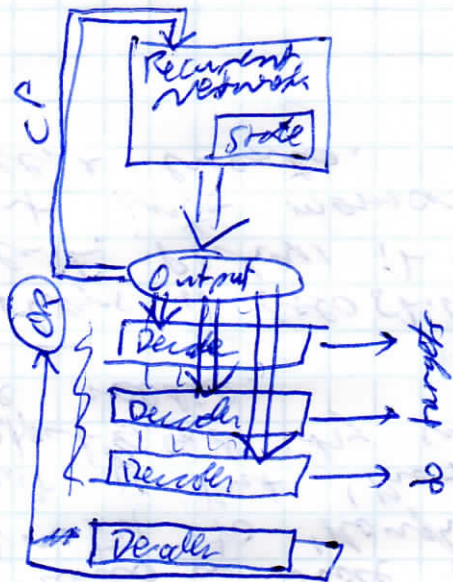
- Above doesn't quite fit
with receiving own thoughts
as cf input ·· it's
redundant — internal & diss
state fully accounts for it
& main output ~~as needed~~
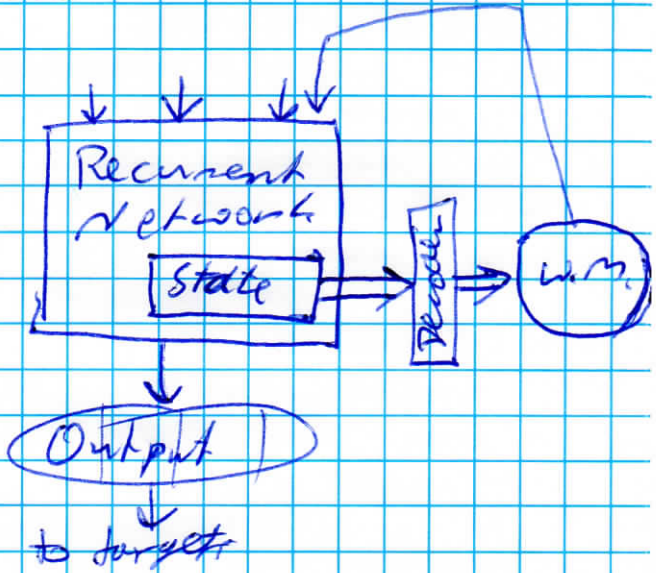~~#~~ needs a whole section
devoted to 'thought'.
what is that 'thought'
output & and why would it
be needed.

~~#~~

Perhaps it does still have
merit. Perhaps it ~~becomes~~
input into memory of
recent events.

# Variant

# Variant



Recurrent
Network

State → Decoder → $w.m.$
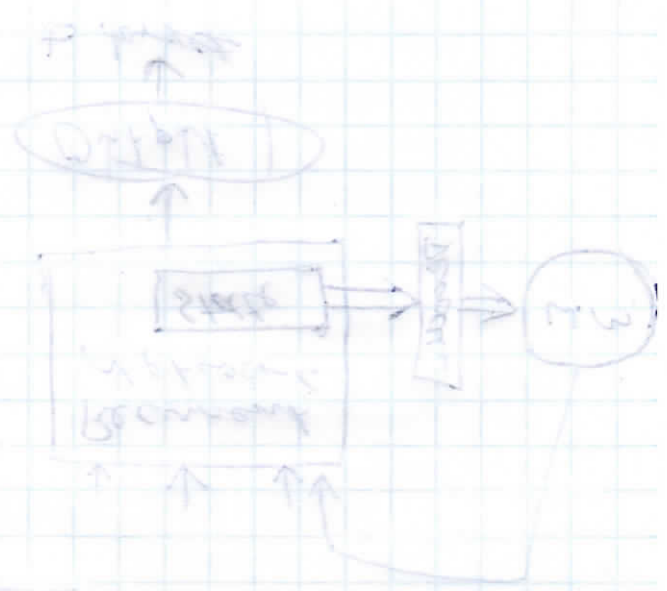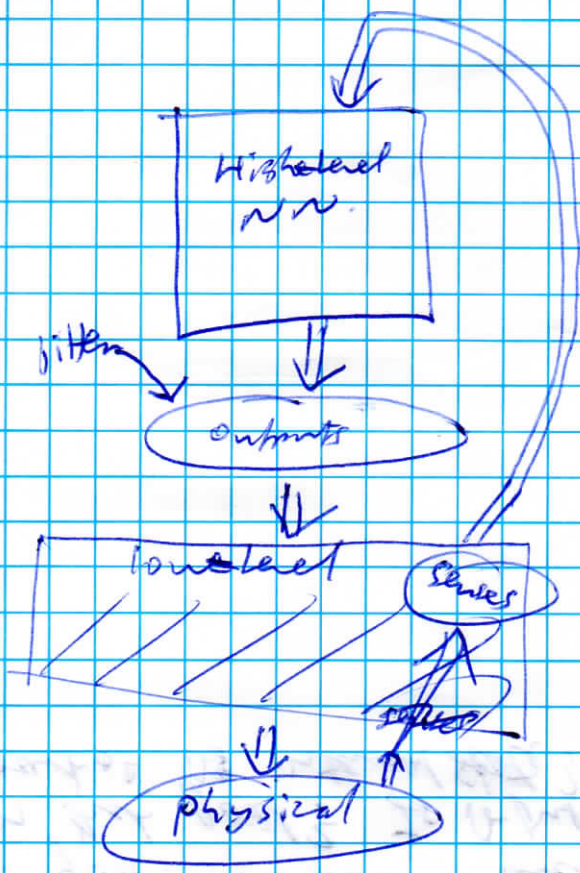
Output

to targets

# Variant

- BY Original external
  recurrency ~~~~~ ~~~~~ for
  both state & cs

# Executive control

Needs to "discover" that it can
control low-level. For now,
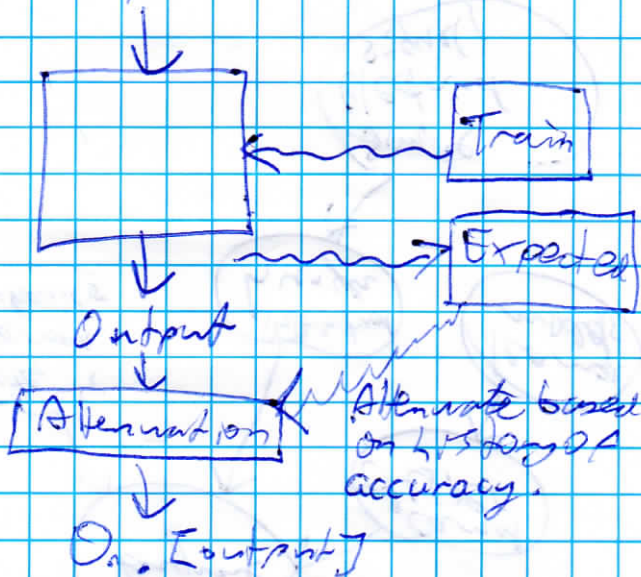re-use 'filler' approach.

High-level
NN

bitten

Outputs

low-level

Senses

Physical

## Goal

Now Executive Control needs
a goal. How does it decide
on that goal? Is it just
part of its internal state?

# Aspec : Attenuated networks

Input
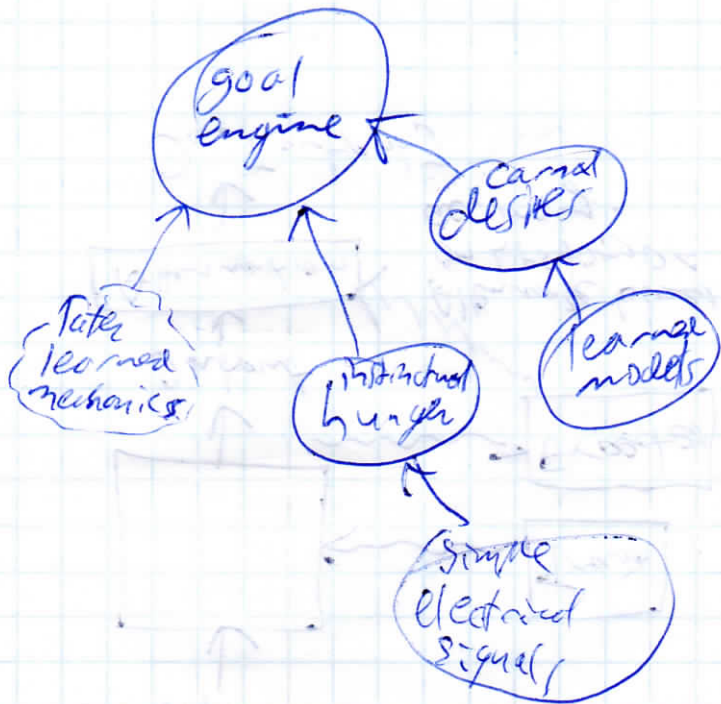
↓



Train

Expected

Output

↓

Attenuation ← Attenuate based on % Accuracy.

↓

$O_n$. [output]

# Self Driven Goal — Basics.

goal engine

carnal desires

later learned mechanics

instinctual hunger

learned models

simple electrical signals

- Rewards in humans.
- external reward is vague & sparse. This allows for a great variability in internal goal/reward systems.

- The internalized reward systems give much finer grained rewards and thus can drive our style of behaviour. As long as they are consistent with the long term reward.

- much of our learning is governed by hard-coded internal rewards (eg: hunger, ~~prediction~~ surprise), so there is plenty of room to ignore external rewards. This is quite different to current RL techniques, which are 100% external.

This makes complete sense, because:

|| There is no external
|| reward, until the agent
|| can grok the existence
|| of an externality.

So, all rewards start as
internal rewards. This explains
why we have so much leeway
to choose to ignore external
rewards — we are hard-wired
for internal rewards, but using
external rewards is ~~an~~ a
~~optional add-on~~. learned
and optional add-on.