-> Single class manages environment

-> task class sets up scene, process actions,
compute rewards, observations,


Prims: basic building blocks of a scene
- light       -> transform
-> mesh
-> can have other prims/ objects
under it.
Attributes:  key - value pair
e.g   color: red


relationships: pointers
-> relationship
-> mesh can be related to
light for shading.

Sim. Spawners  is  a wrapper for USD
API

When simulation starts only after properties
of prim.


The first path in prims is the
objects heirchy in USD.


world                    -> this is
                            how prims
  ├— Cone I                are organize
  └——— Xform

— gibson    in the
   down      world.

```python
# spawn a green cone with colliders and rigid body
cfg_cone_rigid = sim_utils.ConeCfg(
    radius=0.15,
    height=0.5,
    rigid_props=sim_utils.RigidBodyPropertiesCfg(),
    mass_props=sim_utils.MassPropertiesCfg(mass=1.0),
    collision_props=sim_utils.CollisionPropertiesCfg(),
    visual_material=sim_utils.PreviewSurfaceCfg(diffuse_color=(0.0, 1.0, 0.0)),
)
cfg_cone_rigid.func(
    "/World/ConeRigid", cfg_cone_rigid, translation=(-0.2, 0.0, 2.0), orientation=(0.5, 0.0, 0.5, 0.0)
)
```

① ②

Path in USD
heirchy

Config

—> then

translations -----

① create config for shape/
   mesh

   reference the API in isaaclab
   in this case it was
   the ConeCfg method

②

POLICY : takes State of the world and produces an action.

-> its a function
-> because its a function it can be a neural network

-> The policy can take in images instead of velocities, position etc . . . .
-> each neuron will adapt to one pattern of the same (ball / paddle position)
-> This is called policy gradient.

Bellman equation: relationships between The value of a state and the value of the next state.

state: gas runs out, tire wears out

gas is proportional to speed ①
② tire is determined by type of tire
    -> different tires allow for different top speed

① options: make a pit stop
        -> 10 sec penalty
    run out of gas:
        -> DNF

② options: choose between type 1, 2, 3

    1) higher speed, lower life

    2) med speed, med life

    3) lower speed, highest life

    -> pit stop
        -> 10 sec penalty
    don't stop
    -> DNF

high level model:                          -> linked to reward

    at each time step: input $\begin{bmatrix} \text{gas life} \\ \text{Tire life} \\ \text{state} \end{bmatrix}$

output at each time
        step = $\begin{bmatrix} \text{Action} \\ \text{space} \end{bmatrix}$

Action space = [ steering, gas, break, pit stop ]

it will be penalized if it DNFs, or gets
a slower lap time against adverary

reward = 100 ~~=~~ $\frac{10000}{100}$ ft if it take
$\overline{100}$ (0.000 steps
its 0

driving score

reward function = 100 $- \frac{t}{100} + \frac{\theta}{100} + \Omega$

Total          reward          for every
score          lost            second
for            for how         its lap time
               long it         is faster than
being          falls to        opponent.
a good
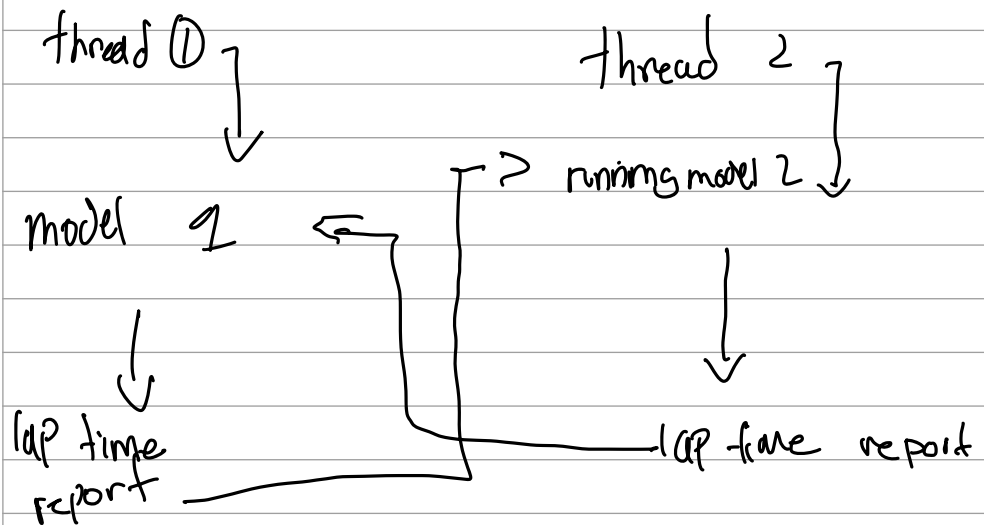driver         drive           $\Omega = -10000000$
                               if DNF

for calculating lap time:

(current time in lap)

30sec    to complete

80%   &

300sec
(100%)

thread ①                         thread 2

model   2  ⇐              ⇒ running model 2

↓

lap time                        lap time report
report

$$\text{Reward function} = \text{Total reward} - \frac{t}{100} + \frac{\theta}{100} - \frac{P}{100} + W$$

total reward = how perfect can you drive

$t =$ penalizing for how perfect you can drive (time
in lap?)

$\theta =$ reward for every step which your time is
less than your competitors

$P =$ penalty for fuel/tire being less than 40%
and increased gradtly when bellow 0.

$w =$ reward/penalty for winning or losing.

1. trained the high level ✓
   -> generate graphs

2. get single agent racing ✓

3.

End goals: assuming there is a pit stop method, agent to be competetive against human.

⤷ We need a robust pit stop method.
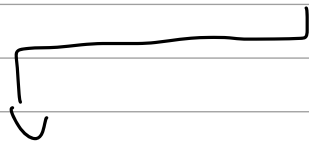
pit stop method;

currently : Total fuel,

   Total tire level,

   fixed rate of consumption

   —> you can race

   a whole lap without

   pit stopping

increase rate

⤷ doesn't matter because
both agents will just have to
pit stop at least once.

   ⤷ randomize the rate of
      consumption [clip it]

↳ Just adding a randomized rate of consumption will not add any value or complexity to the system because the model is learning to pit stop based on fuel left, other in words rate of consumption is related to total and we are only observing the related variable.

Essentially the problem we are trying to solve is least number of pit stops in order to complete a lap.

in real life or in sim, multi agent or not the objective is to make it to a pit stop with least amount

of fuel possible.

and the skill comes in

Where how well can you
drive to make your car
last as long as possible
and minimize pit stops while
minimizing time taken

(pit stop)    minimizing use of supplies

which is related to

howell you drive / racing line, speeds/
breaking,
acceleration

while you minimize time around

minimize  a track which means you
breathing ⟵ want to maximize speed

and minimize racing line

with in rules,
very recursive,

reward function of higher level model should minimize pit stops.

– The reason why we removed the randomization in consuption rate, is because it will add variability to the threshold at which the model makes a decision to take a pitstop.

↳ Eg: Let's say it's at 0.3 rate and we want the model to take a pitstop at fuel level 0.2, in time step 1 it's at 0.7, then next it's at 0.4 then 0.1, → This forces the model to take a pit stop at 0.4 or 0.1, both at which are not a ideal condition. Because the way the reward function is setup it will get punished either option. We saw af

We saw after training, that the reward wasn't converging and was continuisly going into the negatives, and the model was continously getting punished and wasn't learning. Regardless of randomizingf the consumption rate, the variable that is being observed by the model is the total fuel and tire trend level left, therefore changing the consumption rate, only changes the number of timestamps required for the fuel level and tire trend level to reach the threshold for the model to make a decision, therefore it added only unessessary complexity and variance on how close the level of the fuel and tire trend will reach the threshold. Realistically it only makes sense to have the consumption rate to be related to the driving

behaviour to the lower models, but it was not possible to do this because we did not have the ability to train the lower model because of the computational limitation, thus we used a pretrteianed lower level model.`