Casting Sim2Real as
Meta-Reinforcement Learning

# Agenda
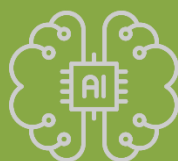
Motivation

What's new

PEARL

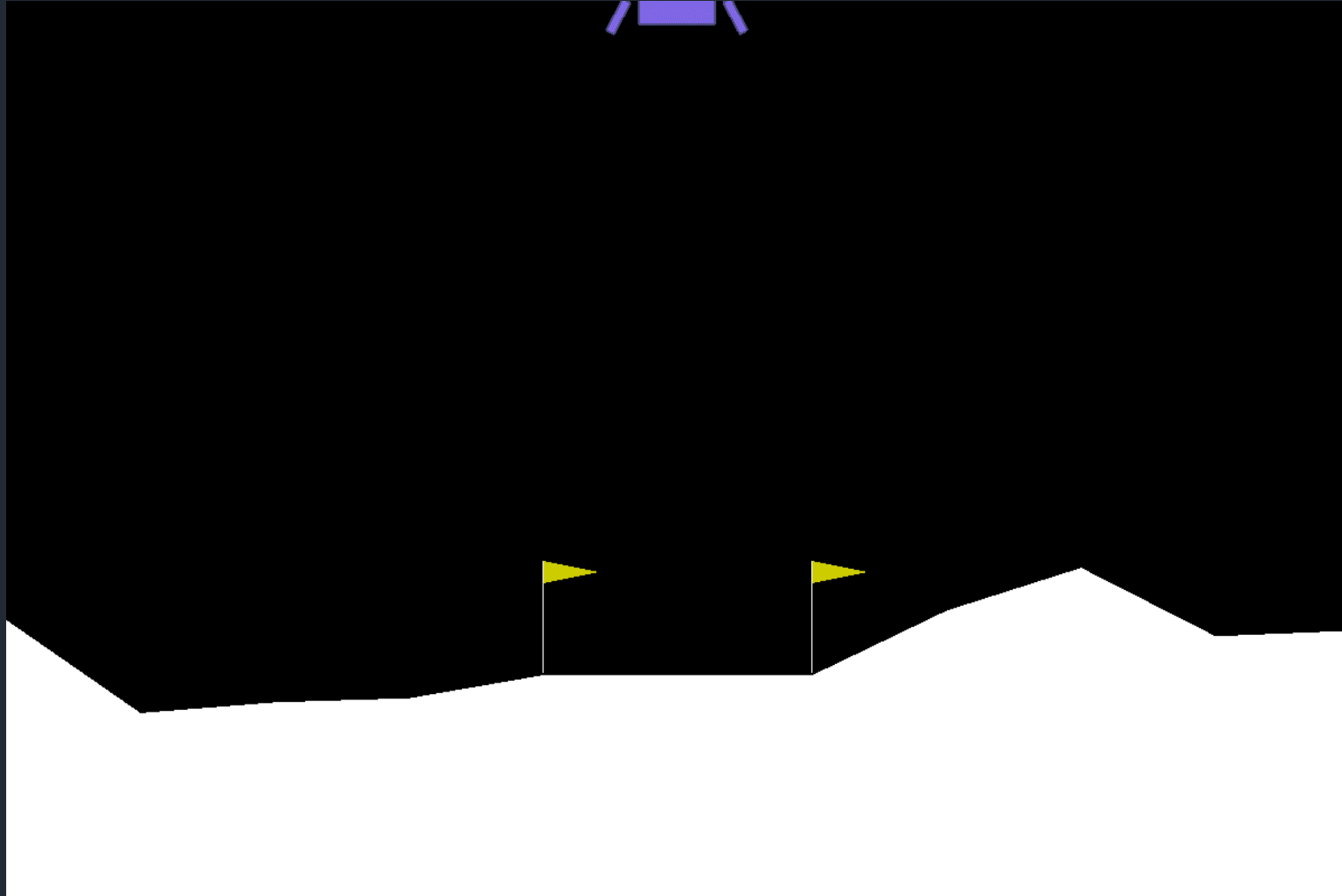PEARL2

Experiments

# Motivation

1. Sim2Real describes the problem of transferring a policy learned in simulation to the real world

2. Simulations are great because they are low risk and low cost

3. Ultimate goal is to train agent in simulation and adapt the policy in real world in a sample efficient way

4. Investigate meta-reinforcement learning on distribution of simulated tasks for sample efficient adoption in new tasks
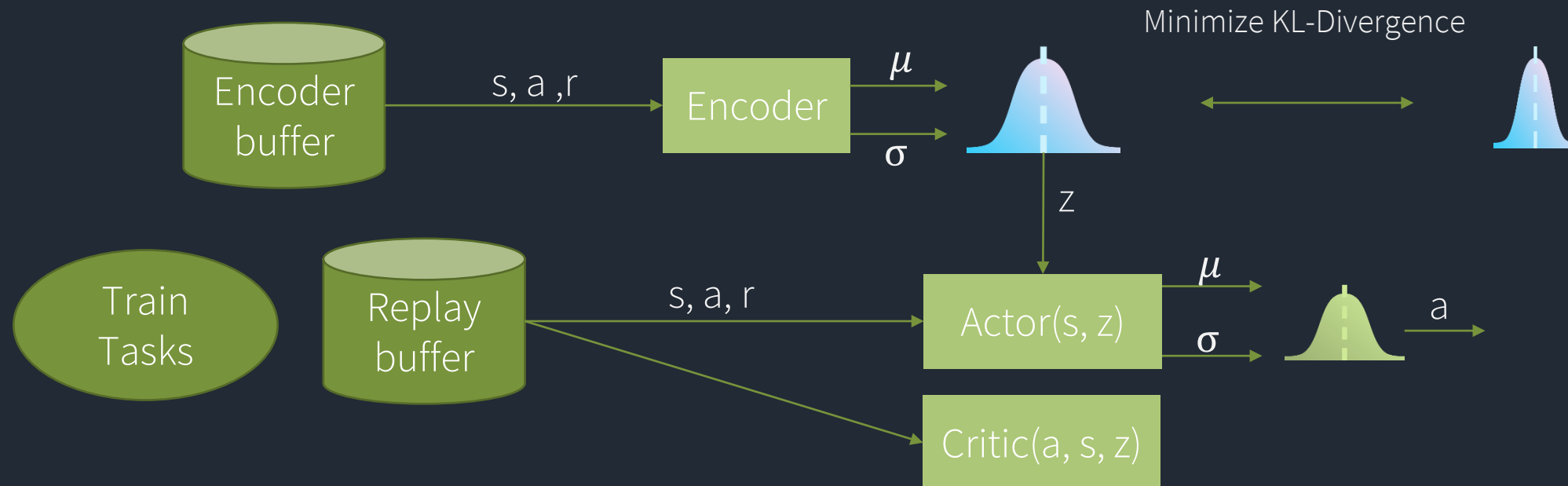
# Our Environement

# What's new

- Implemented PEARL

- Implemented variation of PEARL

- Out of distribution experiments

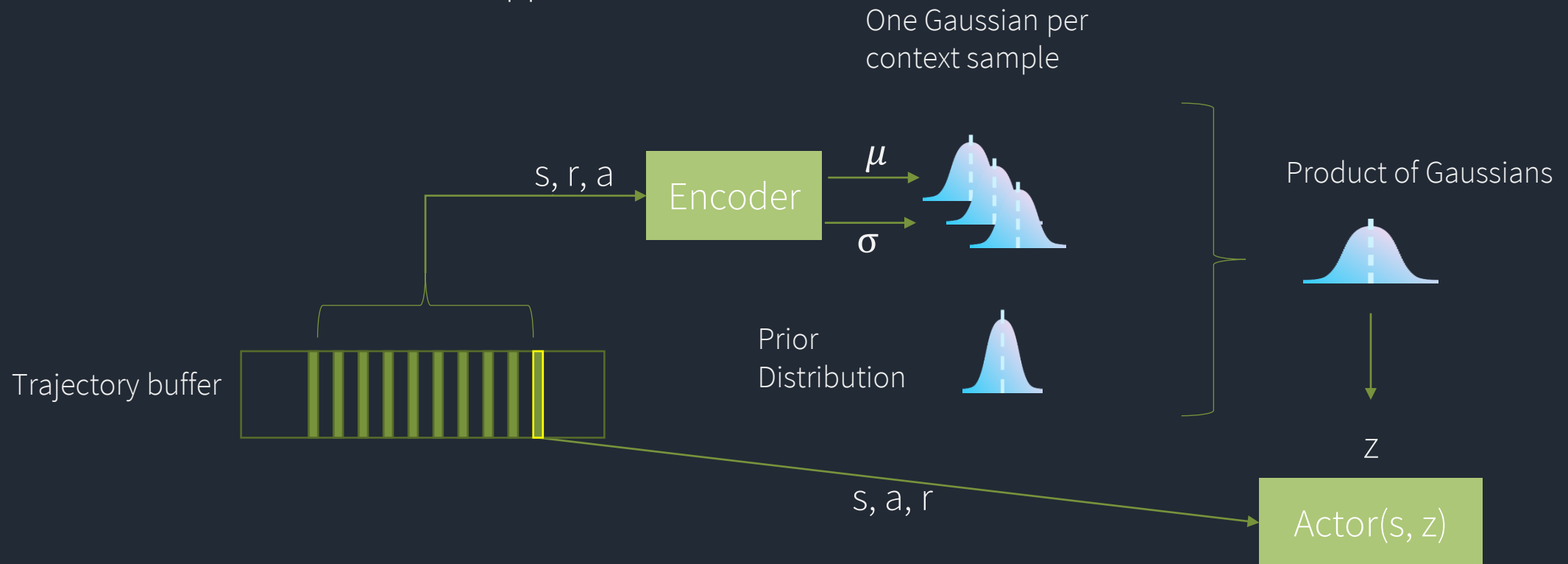- Random wind dynamics

- Many experiments

# PEARL

- Off policy meta reinforcement learning

- Based on SAC

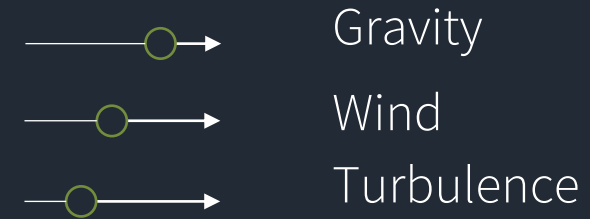- Probabalistic context model to condition actor and critic

Minimize KL-Divergence

Encoder buffer → s, a ,r → Encoder → $\mu$ / $\sigma$

z

Train Tasks

Replay buffer → s, a, r → Actor(s, z) → $\mu$ / $\sigma$ → a

Critic(a, s, z)

# PEARL2

- Only use context from current trajectory

- Infer posterior distribution during trajectory execution

- More realistic approach

One Gaussian per context sample

s, r, a → Encoder

$\mu$

$\sigma$

Product of Gaussians

Prior Distribution

Trajectory buffer

z

s, a, r

Actor(s, z)
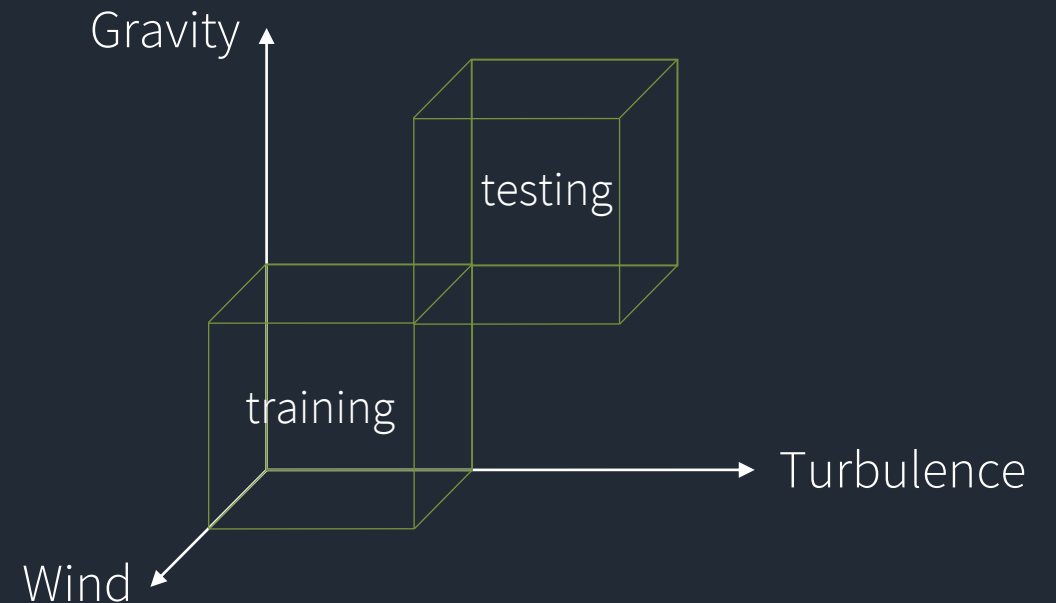
# Experiments

- Domain randomization
  - Random training parameters
  - Validation on grid

- Out of distribution
  - Random training parameters
  - Validation on grid

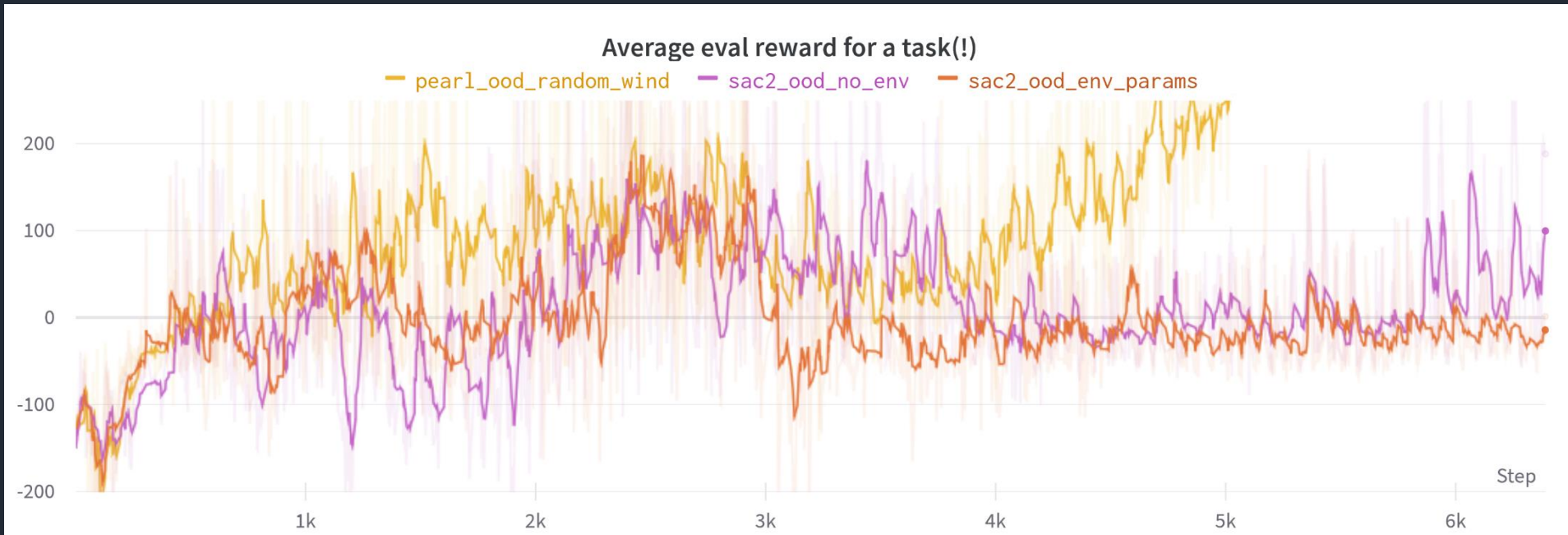# Experiments

- Comparing models
  - SAC
  - SAC2
  - PEARL
  - PEARL2
- Different modes
  - Inside distribution
  - Out of distribution
  - Passing parameters
  - Not passing parameters
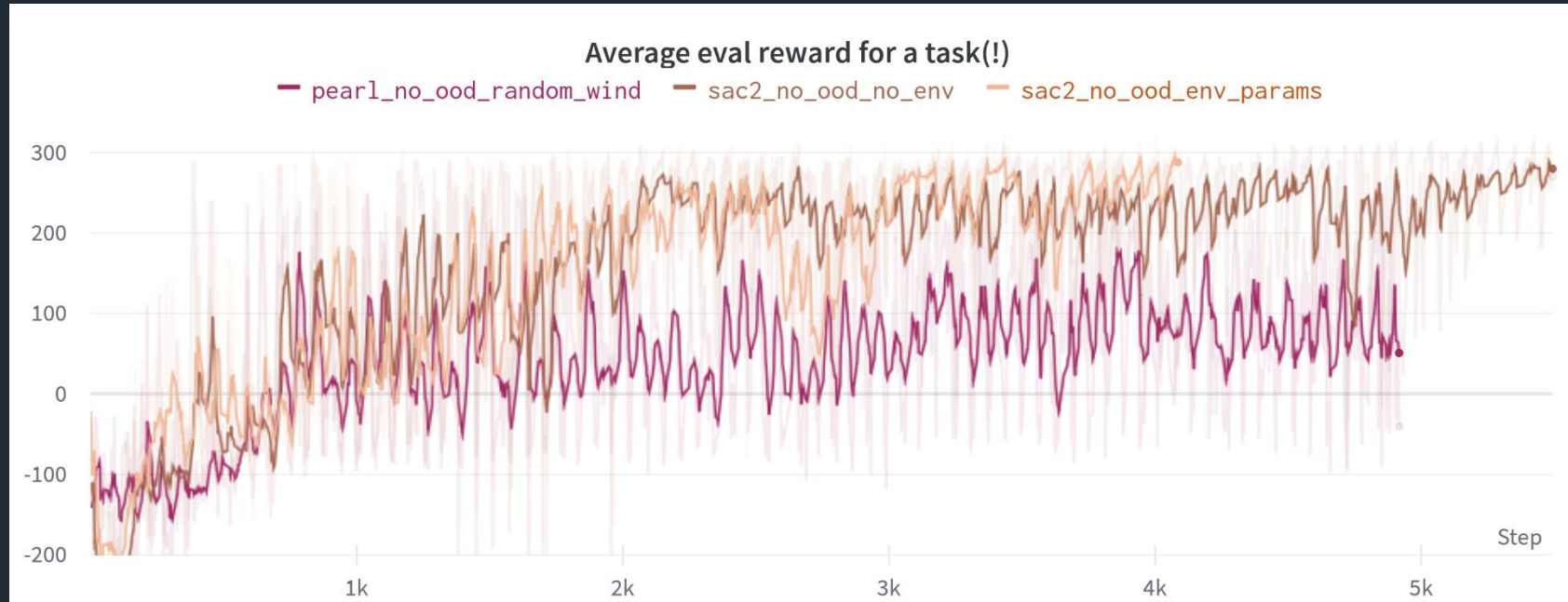  - Fixed wind
  - Random wind

uninformed

| SAC | PEARL | SAC | PEARL |
| --- | --- | --- | --- |
| SAC2 | PEARL2 | SAC2 | PEARL2 |

Inside distribution                                      Out of distribution

| SAC | | SAC | |
| --- | --- | --- | --- |
| SAC2 | | SAC2 | |

informed

# Experiments - OOD



**Average eval reward for a task(!)**
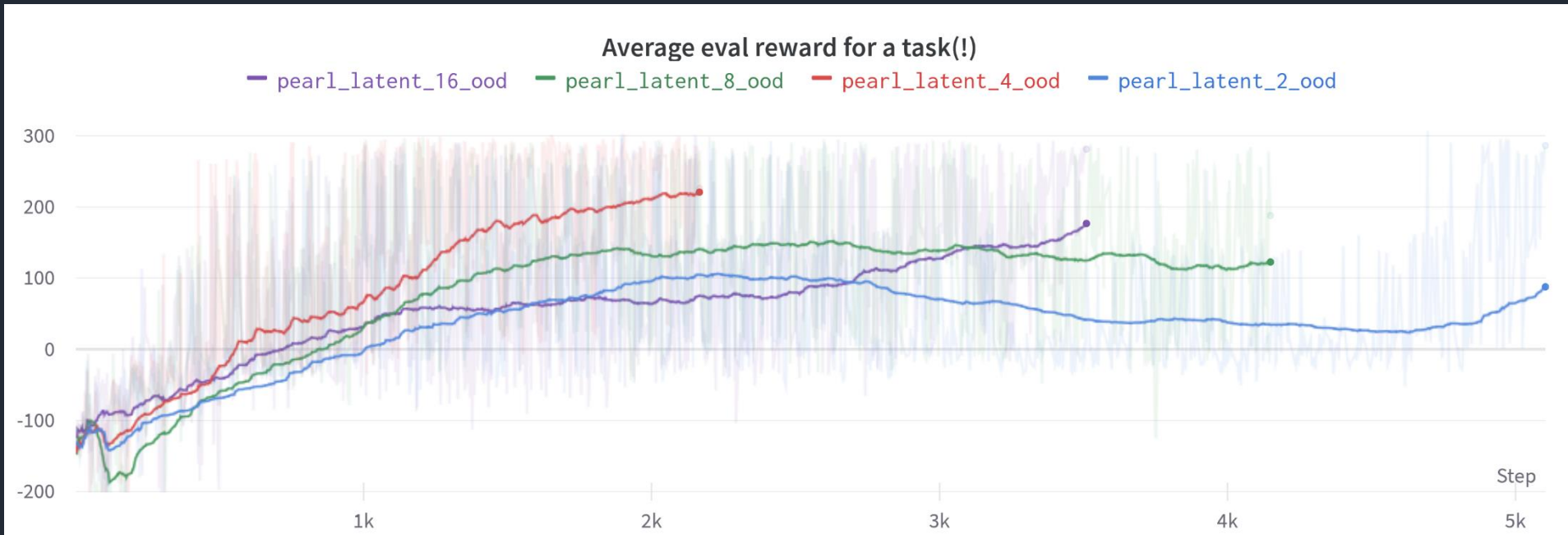— pearl_ood_random_wind  — sac2_ood_no_env  — sac2_ood_env_params

- PEARL outperforms SAC & SAC2 in OOD case

- PEARL solved all 27 evaluation tasks after 5k steps

- (in that set of experiments PEARL has 5 latent dimensions)
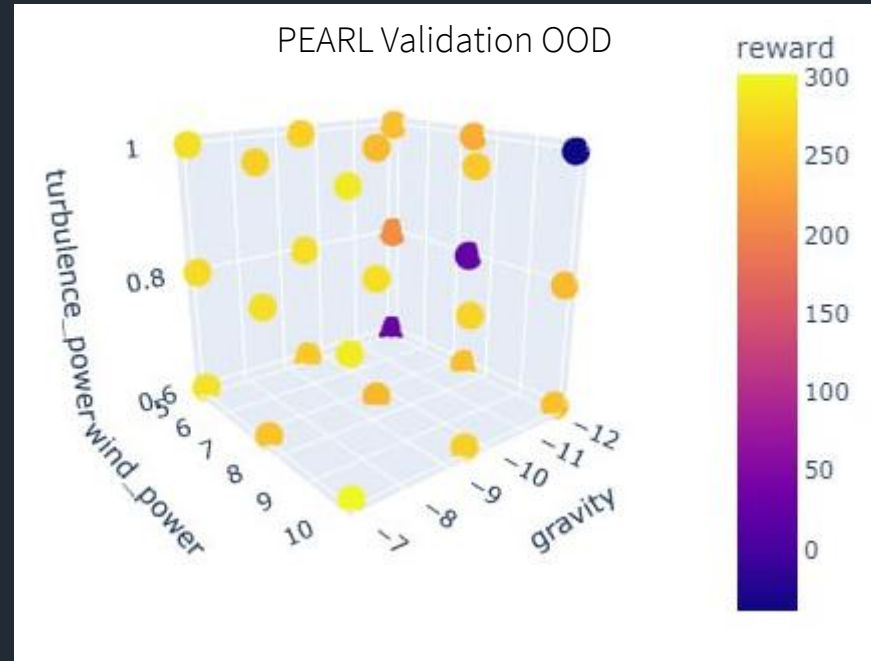
# Experiments - ID



Average eval reward for a task(!)

— pearl_no_ood_random_wind  — sac2_no_ood_no_env  — sac2_no_ood_env_params

- SAC & SAC2 outperform for PEARL inside distribution tests

- Latent size of 5

# Experiments – PEARL latent size



**Average eval reward for a task(!)**
— pearl_latent_16_ood    — pearl_latent_8_ood    — pearl_latent_4_ood    — pearl_latent_2_ood

- PEARL performs best for latent size of 4
- More than 4 overfits, less than 4 underfits

# Experiments – Validation Hypercube



PEARL Validation OOD

- In OOD validation PEARL performs worst for high gravity setting

- Gravity has strongest impact on reward from all three parameters

- PEARL was trained on low gravity

# Experiments – Latent Correlation Map



PEARL, Correlation map between environment parameters and latent space means in OOD setting

- How does PEARL encode environment parameters?
- Explored for the best model with 4 latens variables
- Correlation map shows:
  - Latent variable 1 & 3 have same correlations
  - Latent variable 2 is orthogonal
  - variable 4 slightly different to 1 & 2

# Perfomance of PEARL2



Average eval reward for a task(!)

— pearl2_context_16   — pearl_latent_4_ood   — sac2_ood_no_env   — sac2_ood_env_params

- PEARL2 was trained in hurry with small batches, still outperforms SAC

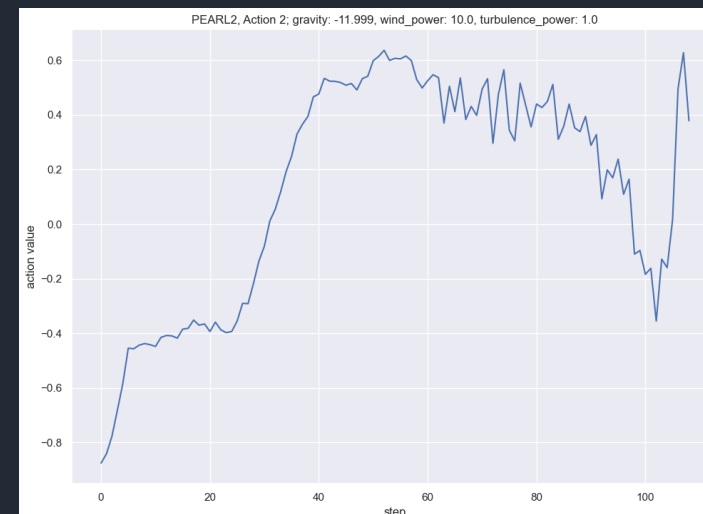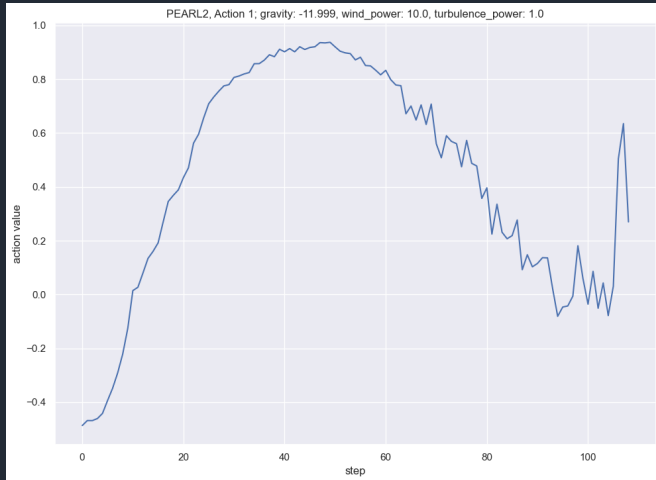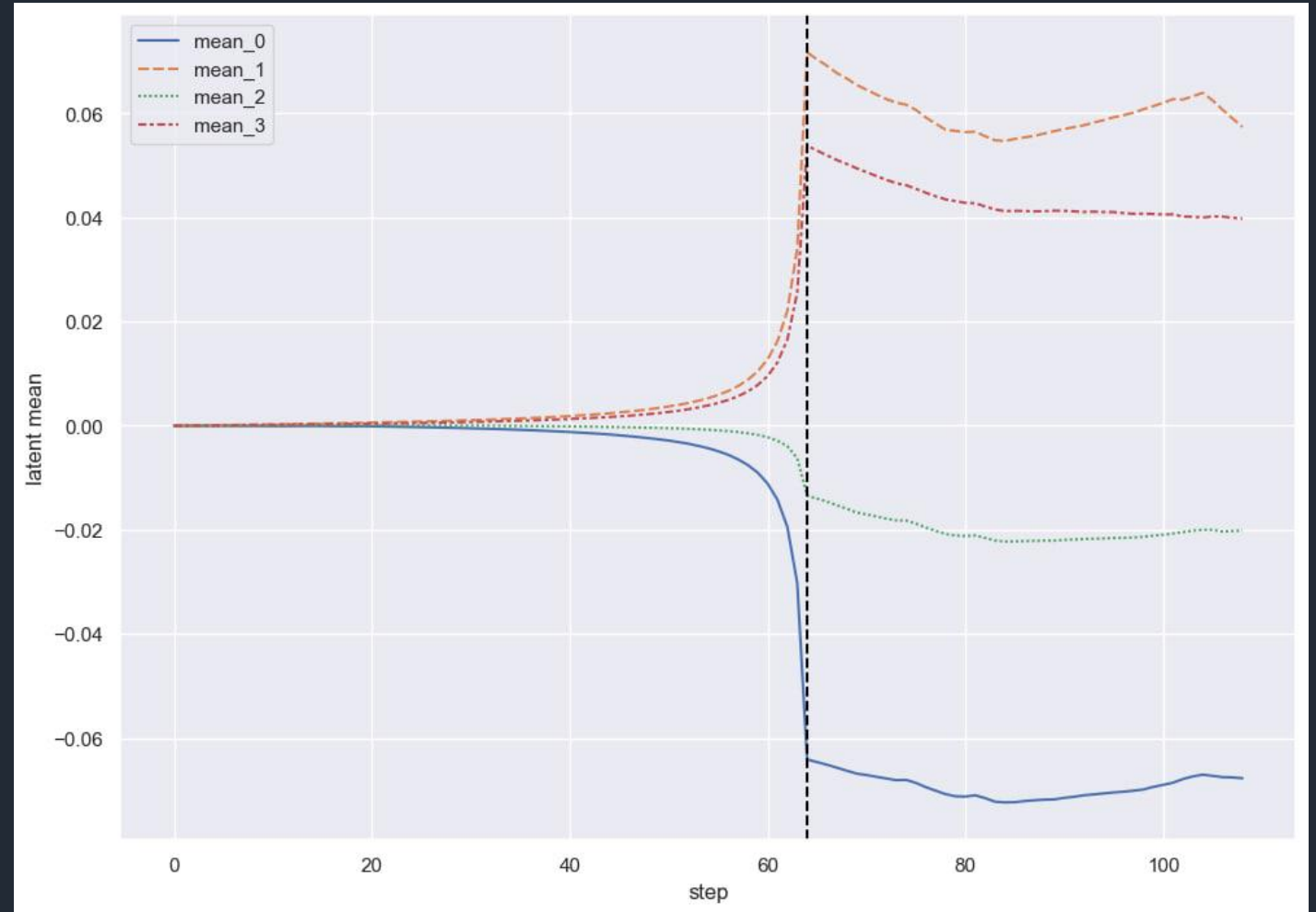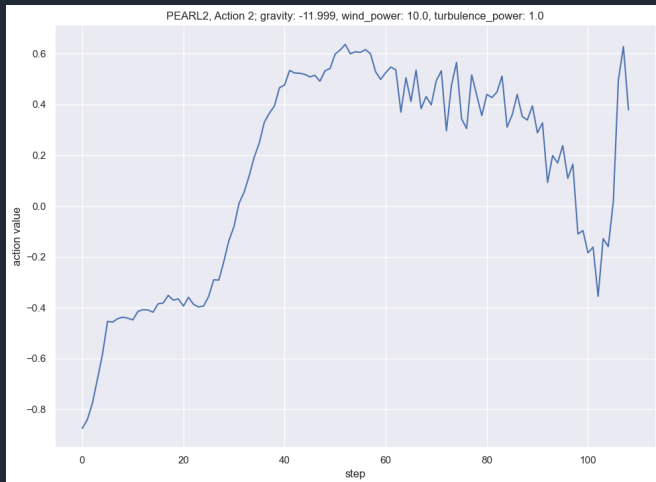# Demonstration – the hardest case

PEARL

SAC2

PEARL2

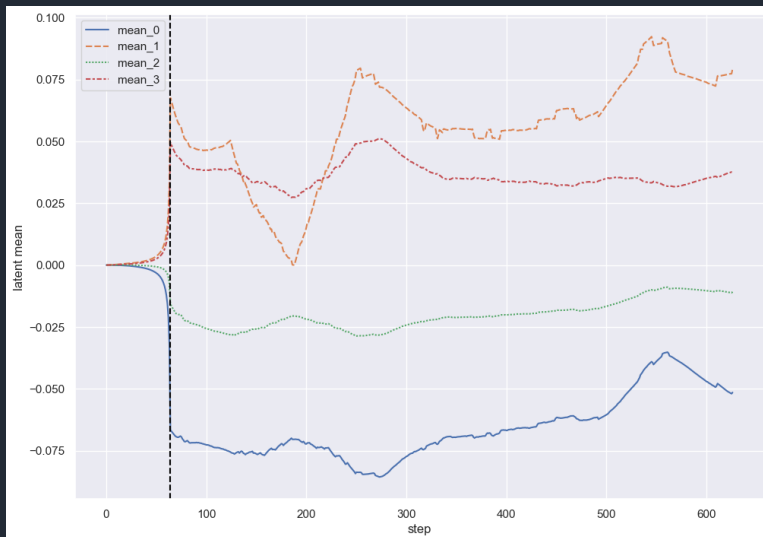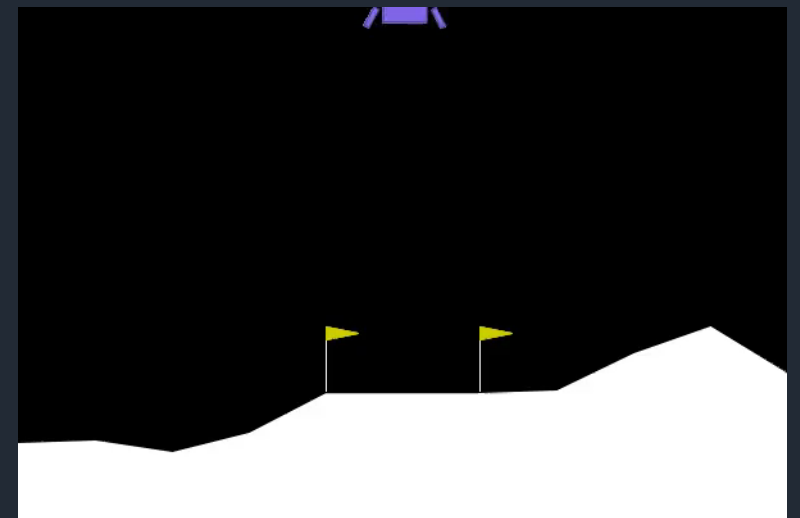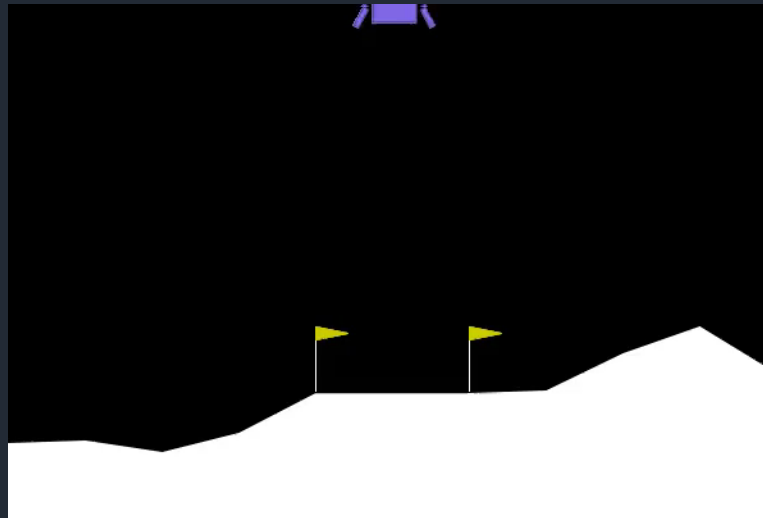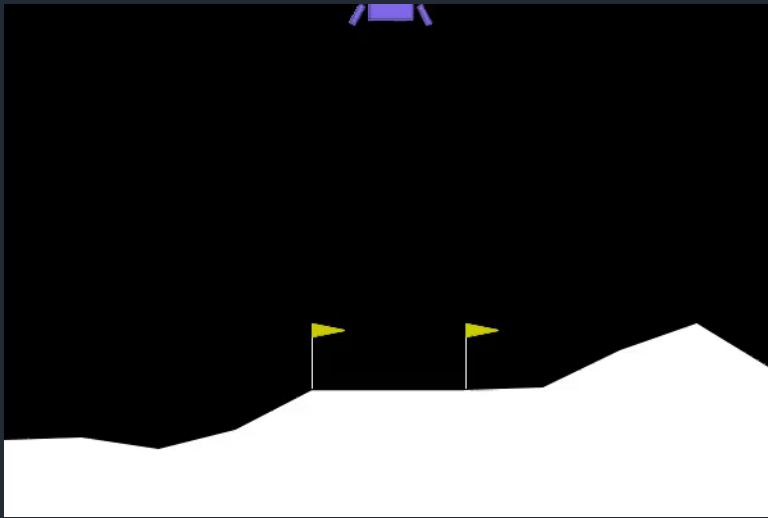# Action comparisons

# Dynamics of latent variables in PEARL2



After latent variables become "saturated", actions begin to become more "flicking", implicitly indicating uncertainty about the environment.

# PEARL 2 with different wind trajectories

(but same environmental parameters)

# Engineering results to show off:

- The training process (>= 500 epochs) was completed at least 551 times

-  =~ 300 hours of compute time (4 cores, 16 gb RAM)

- ~3000 lines of code, debugged with pain and tears, including config files with total 200 options

# Possible continuations

- Smarter way to encode latent dimensions, possibly enforcing orthogonality between hidden features (as in modern GANs)

- Incorporate uncertainty into latent variables in more formal way

- Combined with our online latent variable generation approach, can be used to make models more safe

  - For example, high wind -> uncertainty about dynamics -> switch from learned approach to backoff classical control