

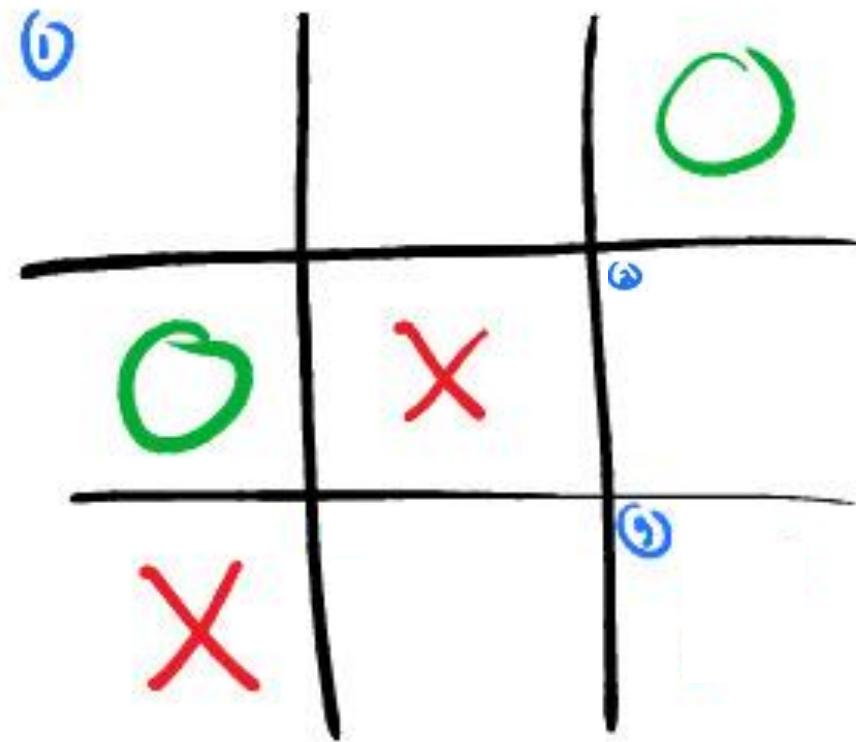
Decisions from data Controlling complex systems with reinforcement learning

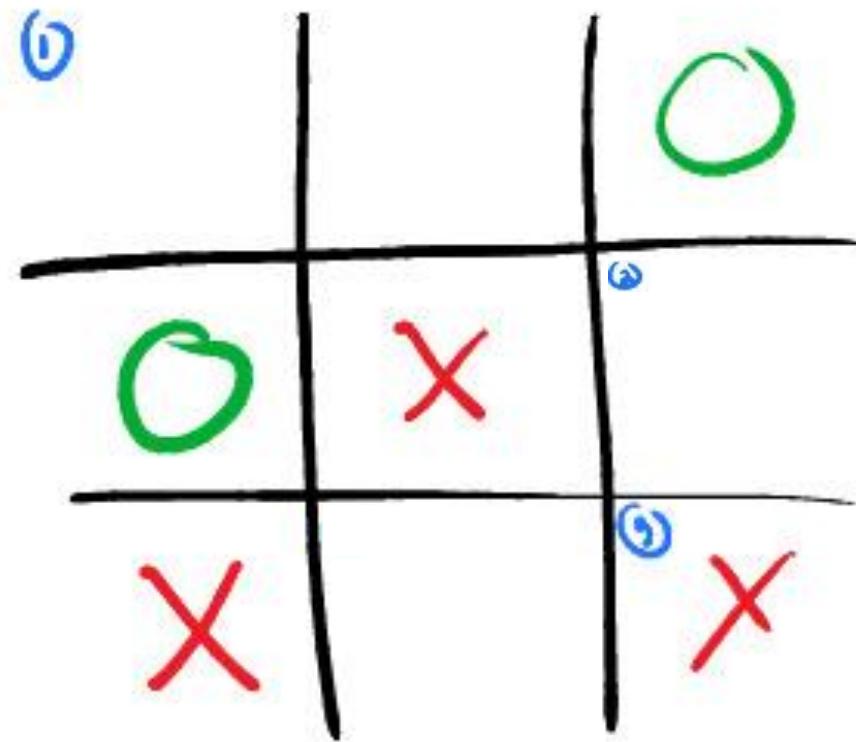
Marc G. Bellemare

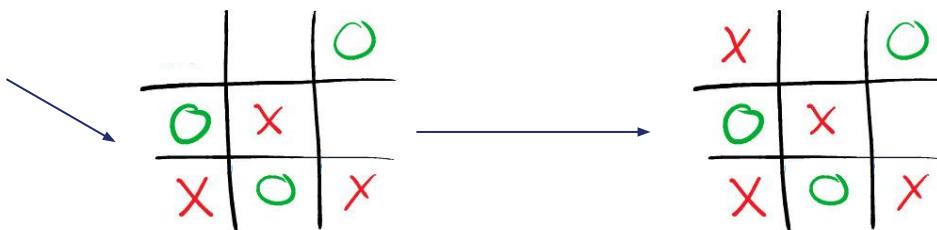
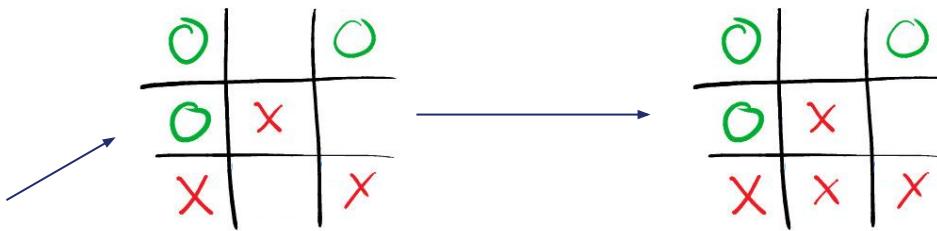
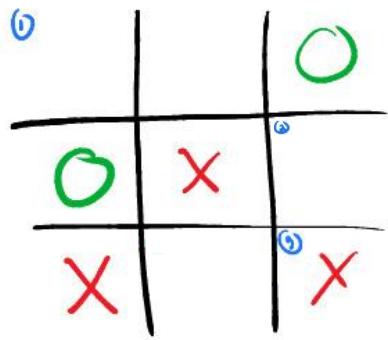
Google Brain, Montreal; Mila

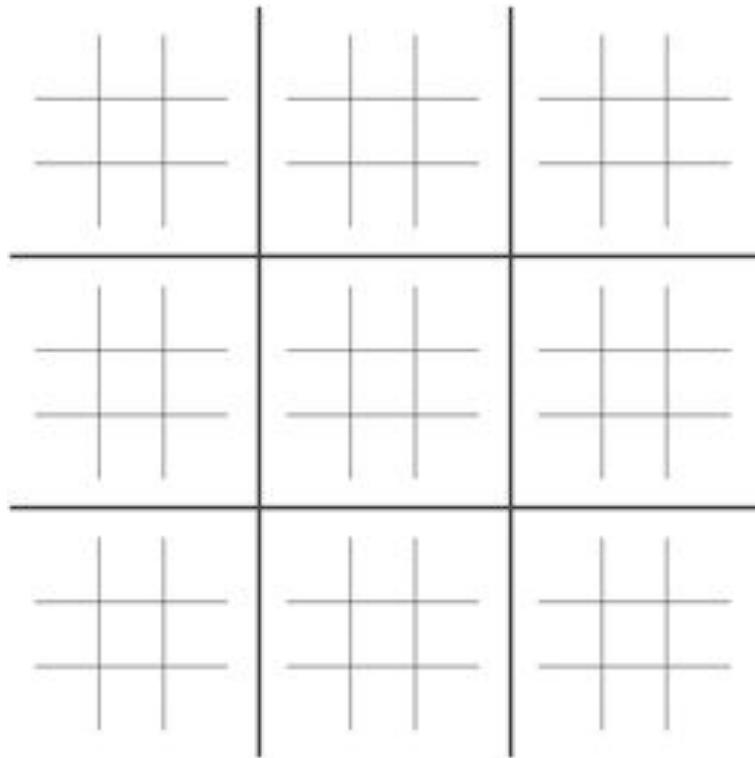
With: Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C. Machado, Subhodeep Moitra, Sameera S. Ponda, Ziyu Wang













Reinforcement learning = trial and error

data → decisions

...

Reinforcement learning = trial and error
data → decisions

Trial and error and cats

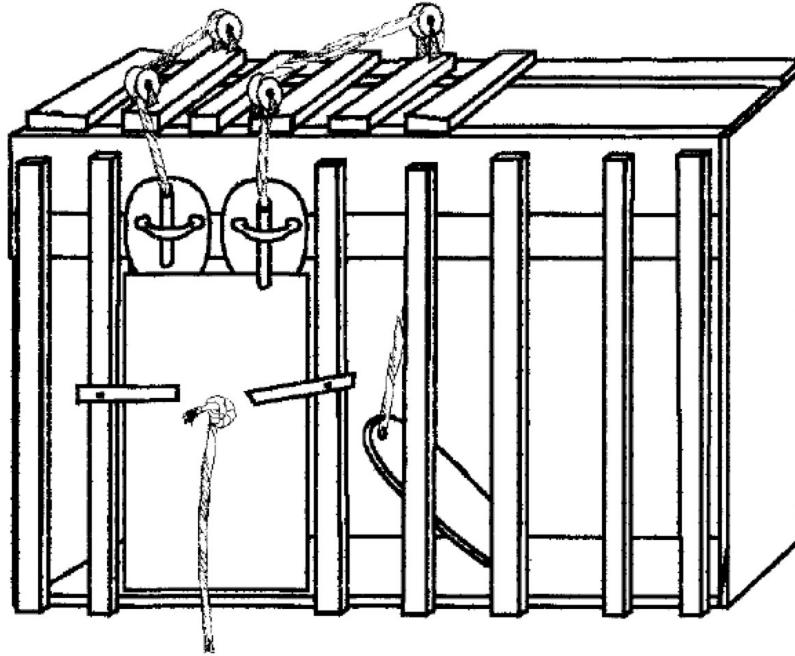
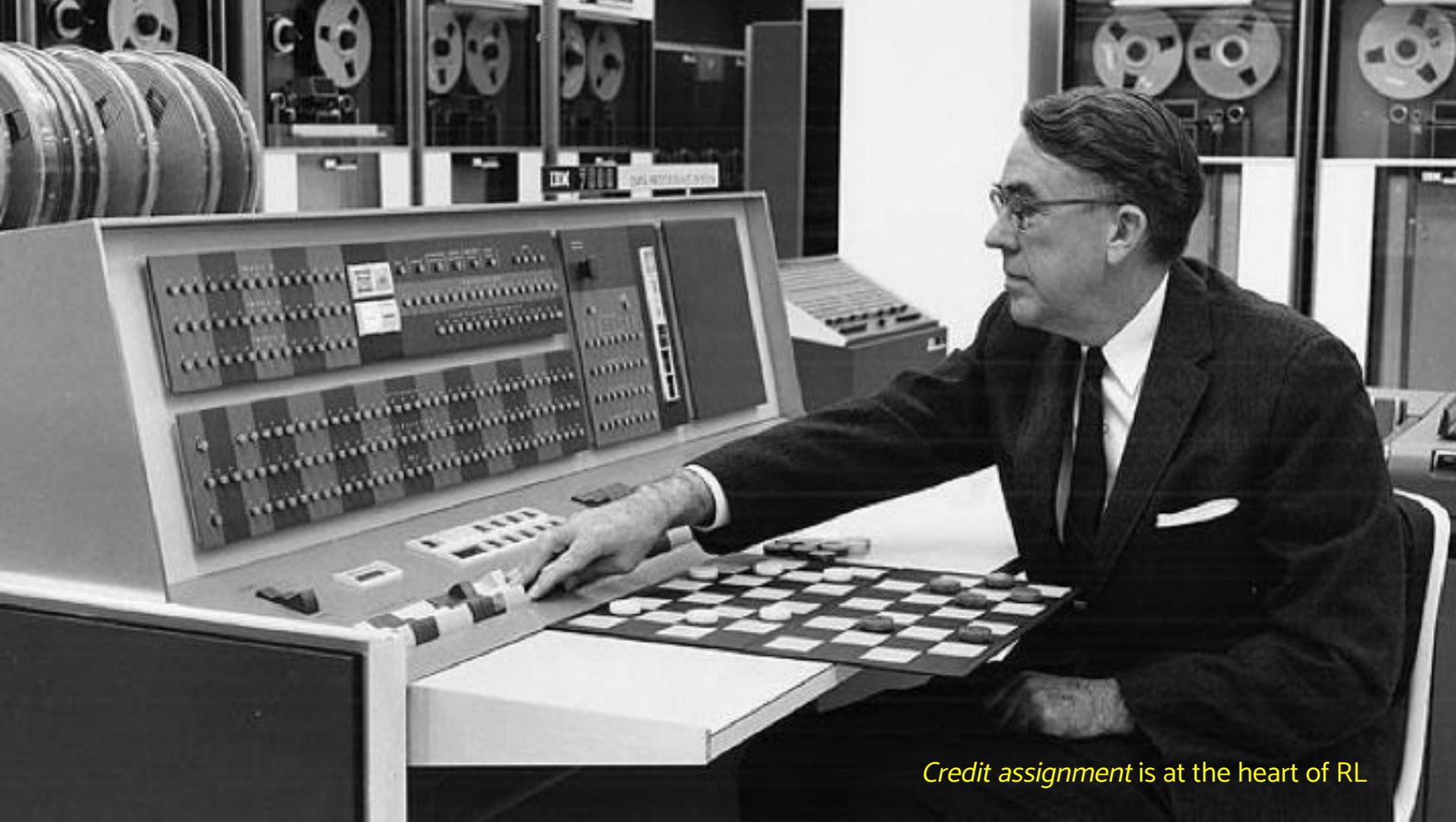


Fig. 4. Box K. The door is held in place by a weight suspended by a string. To open the door, a cat had to depress a treadle, pull on a string, and push a bar up or down. (After Thorndike, 1898, Figure 1, p. 8.)



Credit assignment is at the heart of RL

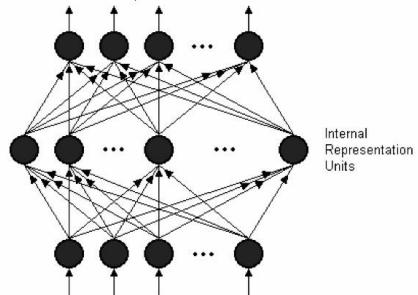
Credit assignment via the Bellman equation

$$Q^*(x, a) = R(x, a) + \gamma \mathbf{E}_{x' \sim P} \left[\max_{a' \in \mathcal{A}} Q^*(x', a') \right]$$

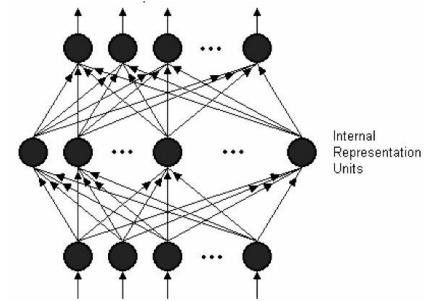
Markov decision process

Implemented as a
Deep neural network

Deep reinforcement learning



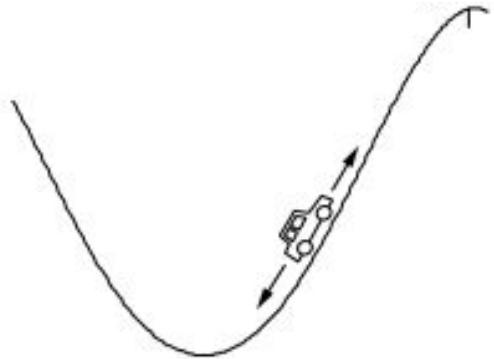
80 units; Tesauro (1995)



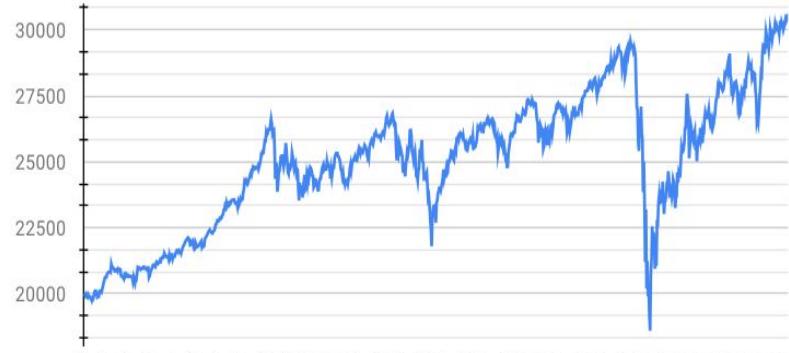
40 x 256 convolutional filters;
Silver et al. (2017)

Many RL problems are...

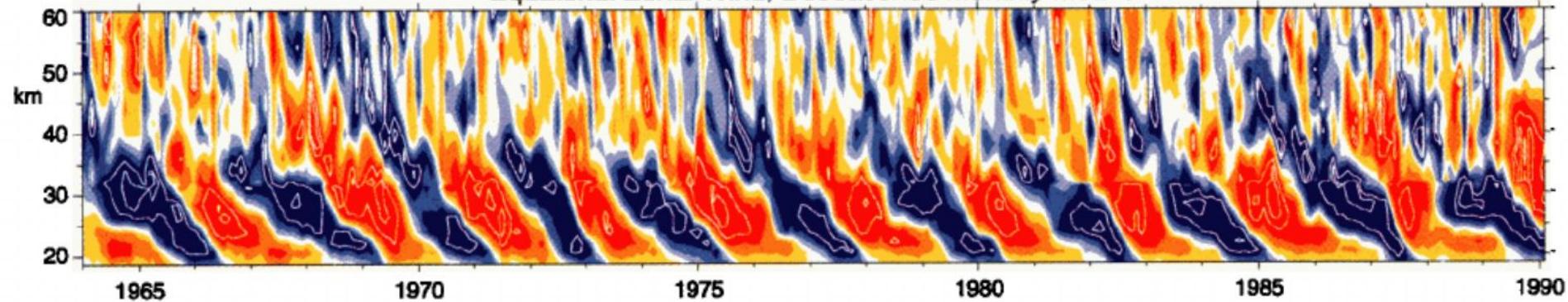
Underactuated



Partially observable



Equatorial Zonal Wind, Deseasoned Monthly Means



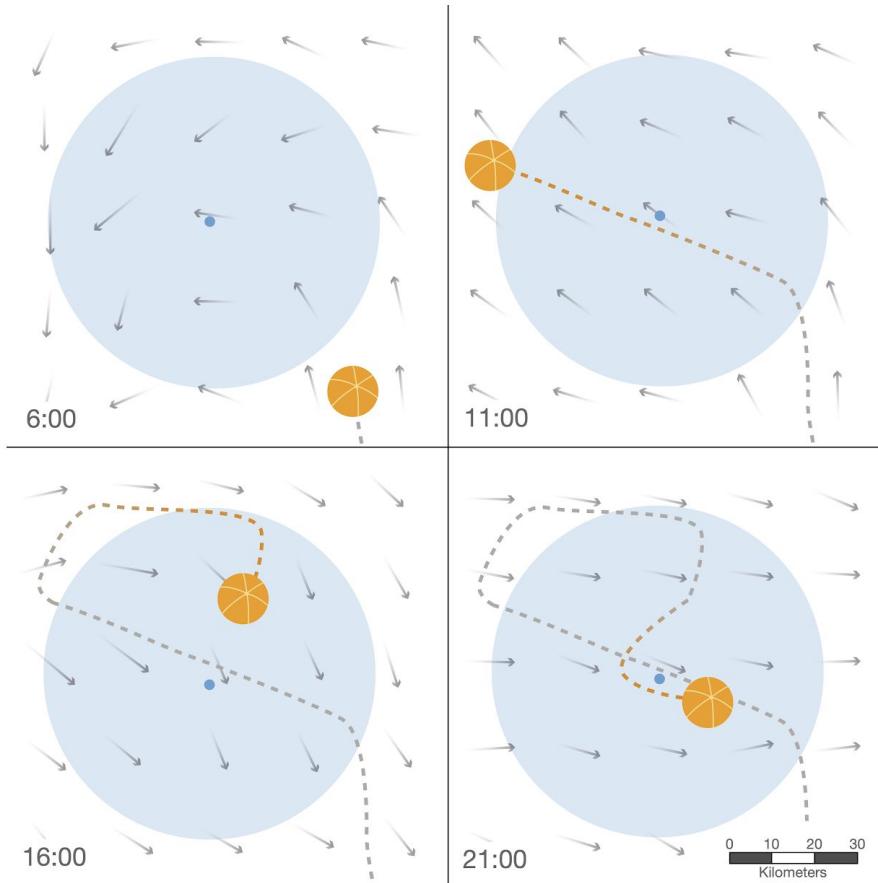
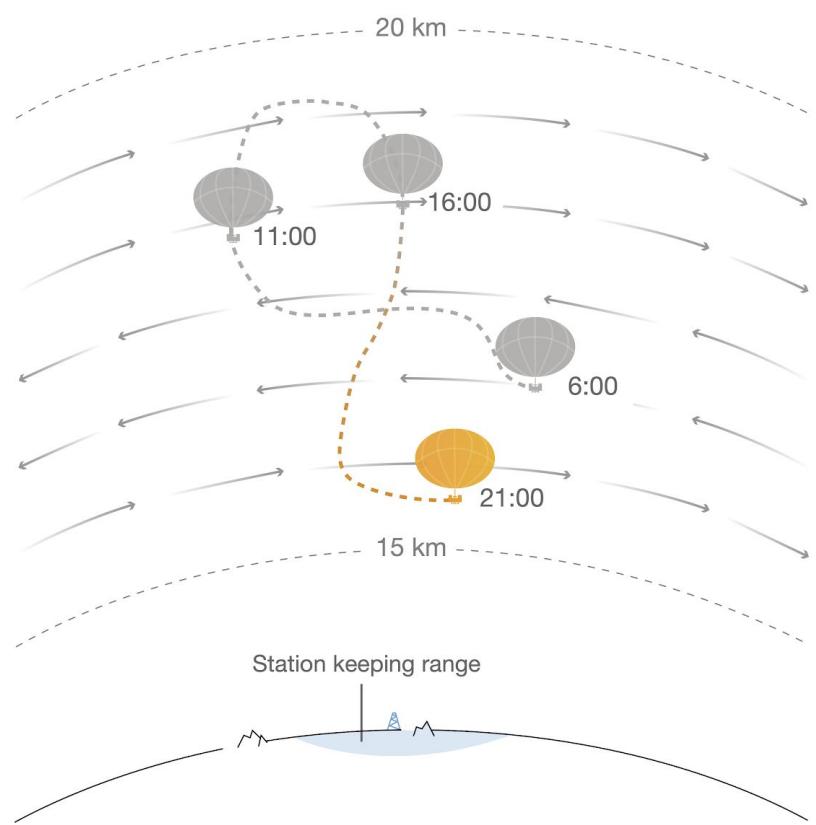
The quasi-biennial oscillation, Baldwin et al. (2001)



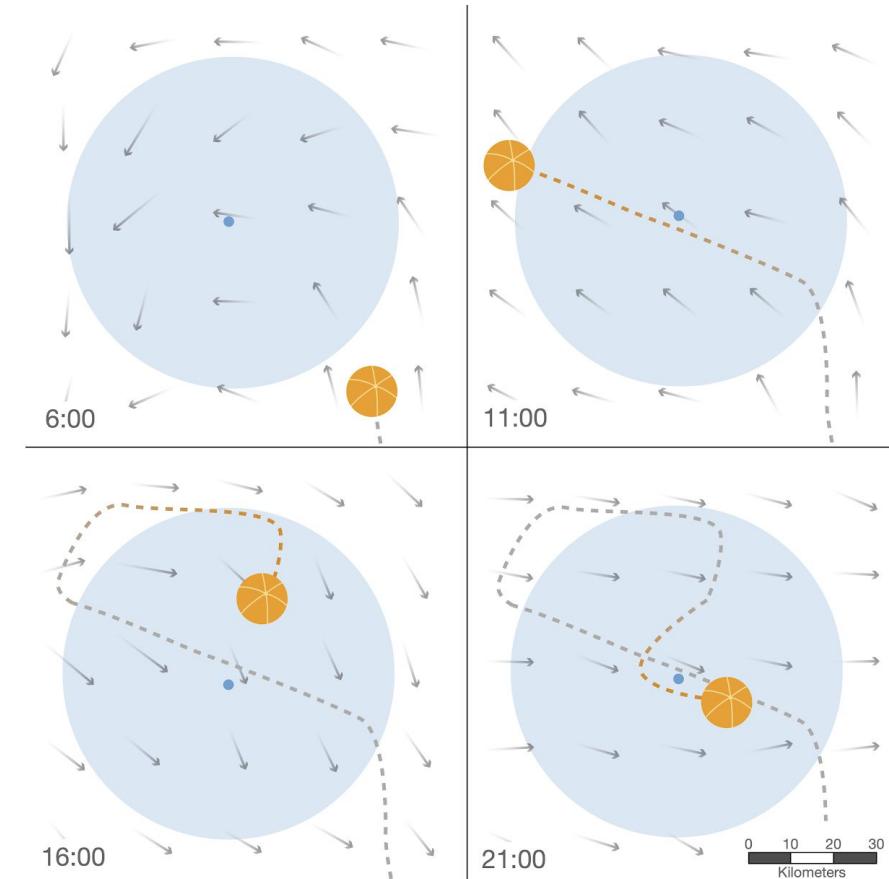
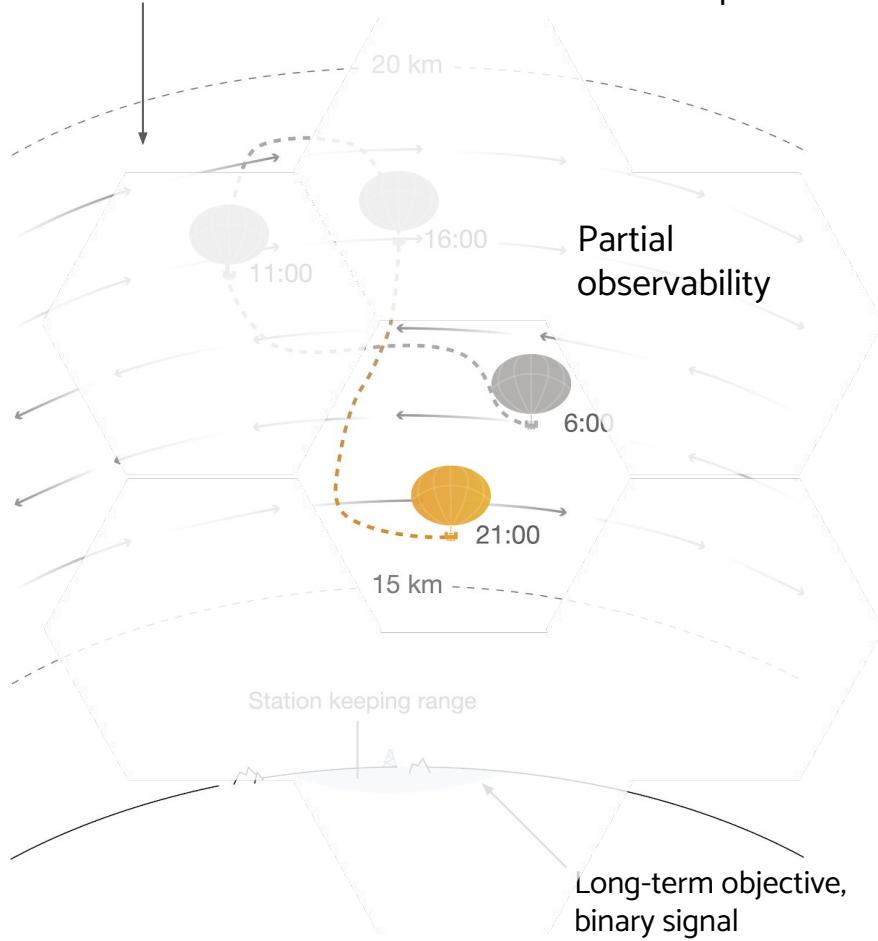


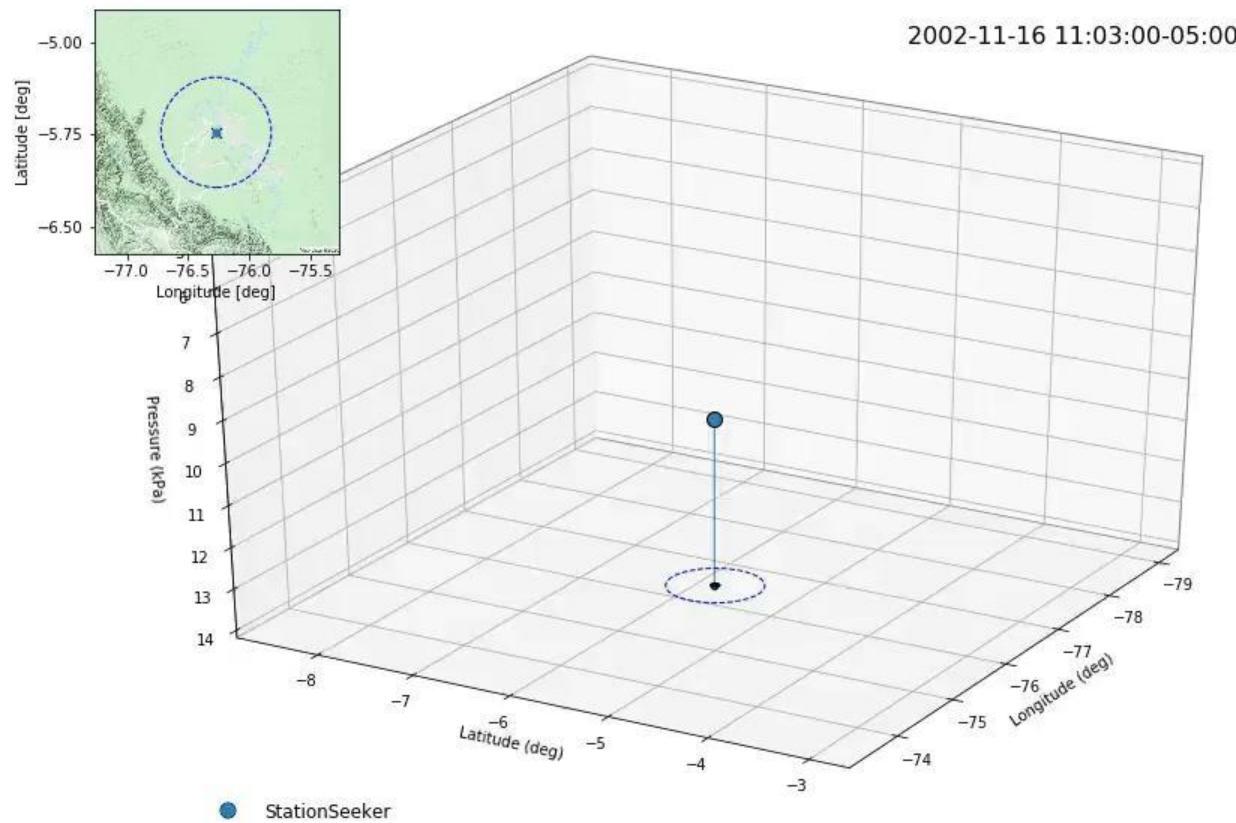


[312 Days in the Stratosphere](#), Loon, Oct 28 2020.



Underactuated system,
stochastic dynamics





StationSeeker in equations

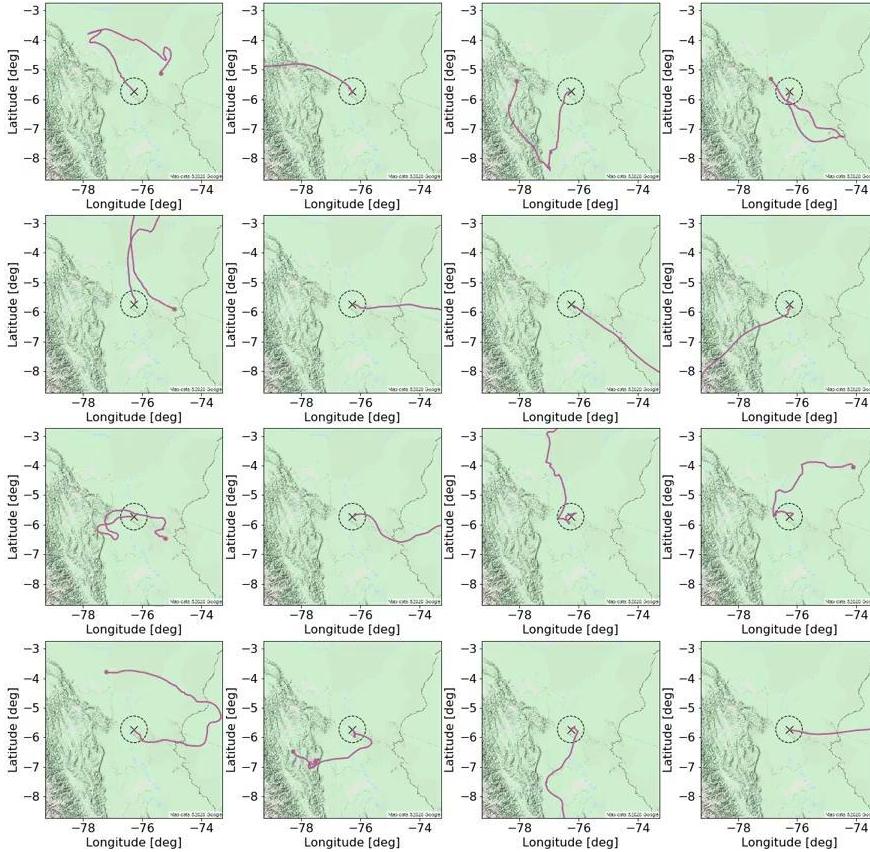
1) Wind score.

$$g_l = (1 - \alpha_\Delta)e^{-w_\Delta \theta_l} + \alpha_\Delta e^{-k_1 \mu_l}$$

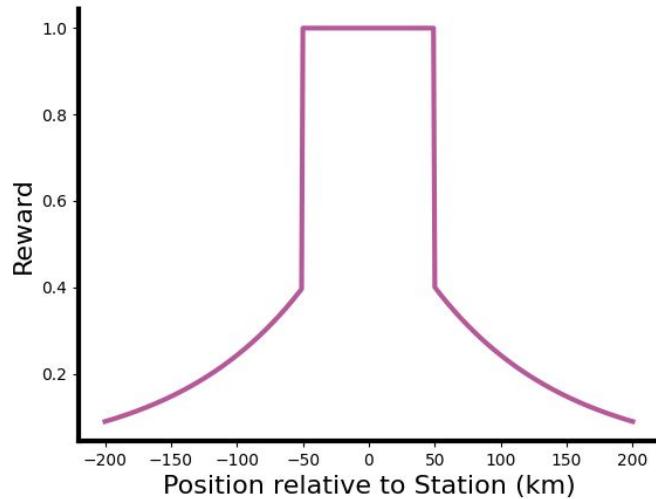
2) Per-altitude score.

$$s_l = (1 - u_l)g_l + u_l g_{\text{unknown}} + k_2 e^{-k_3 |l - l_{\text{current}}|}$$

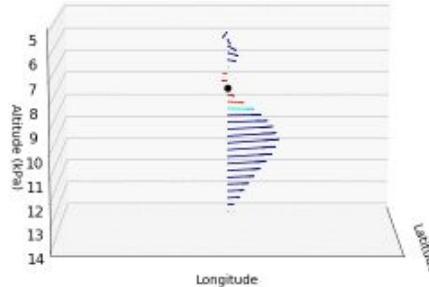
3) Setpoint to max. scoring altitude.



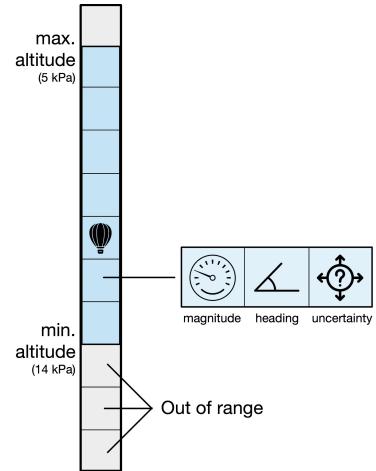
Deep reinforcement learning for balloon navigation



$$r(\Delta) = c_{\text{CLIFF}} 2^{-(\Delta-\rho)/\tau} \times \text{POW}(\omega)$$



Forecast +
measurements +
Gaussian process =
wind column



+16 ambient variables

The simulator

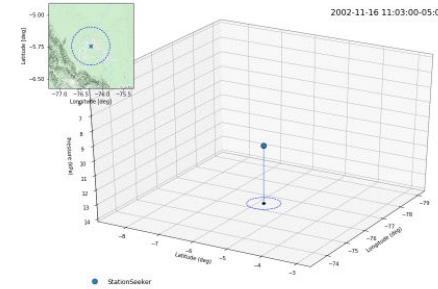
The ERA5 reanalysis (dataset) provides **baseline winds**

- Like real, but
- Low resolution.

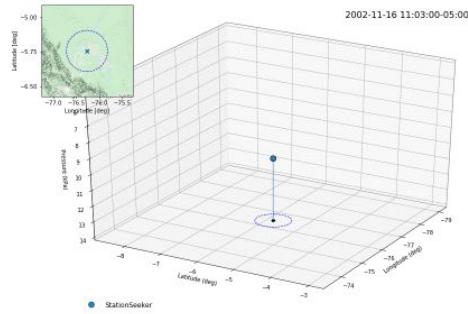
Baseline winds are **upsampled** using procedural noise:

- Statistically plausible
- High resolution
- Effectively infinite supply

An **episode** is defined by **initial conditions** (*latitude, longitude, altitude, time*), station location, and random seed



Design and training



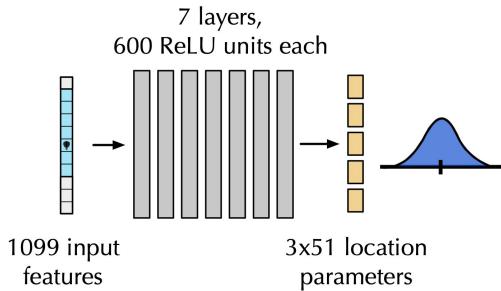
2-day training simulations

In the tropics (+/- 25 lat.)

Starting up to 200km away

Light filtering of “impossible” conditions

Combined exploration + greedy behaviour

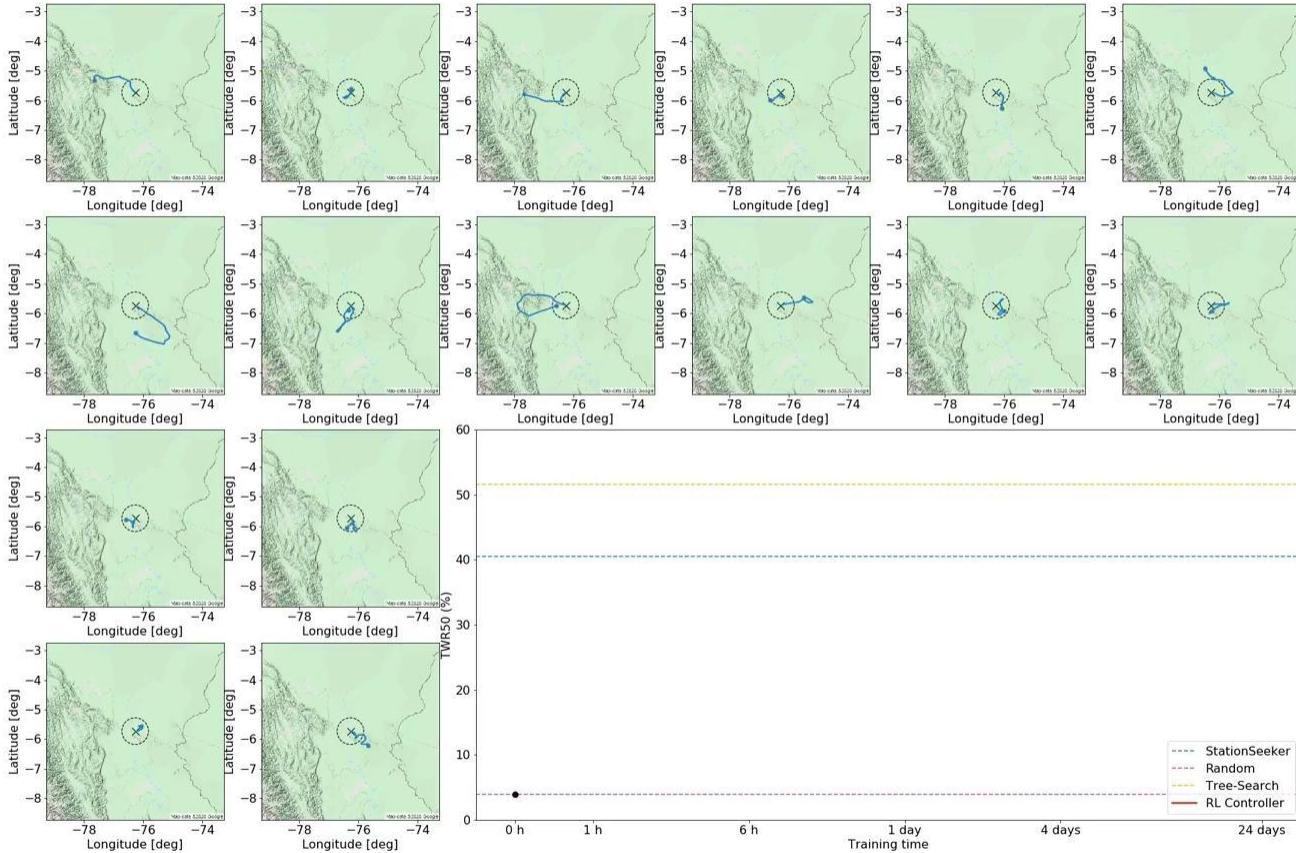


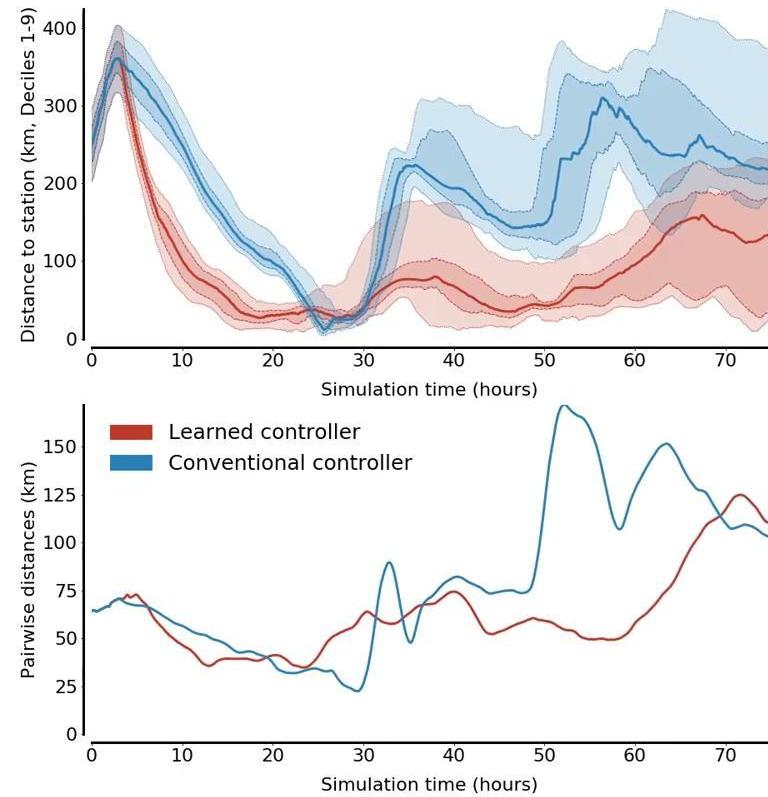
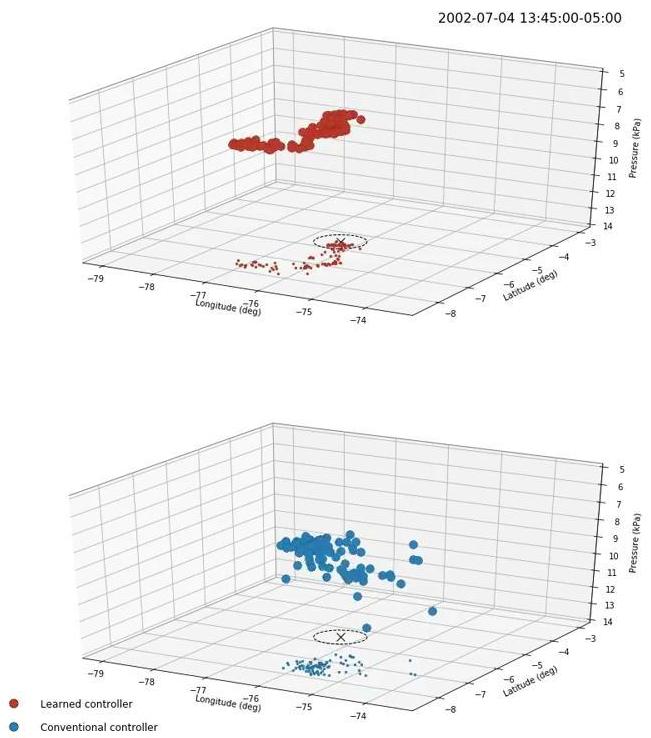
Distributional predictions (QR-DQN)

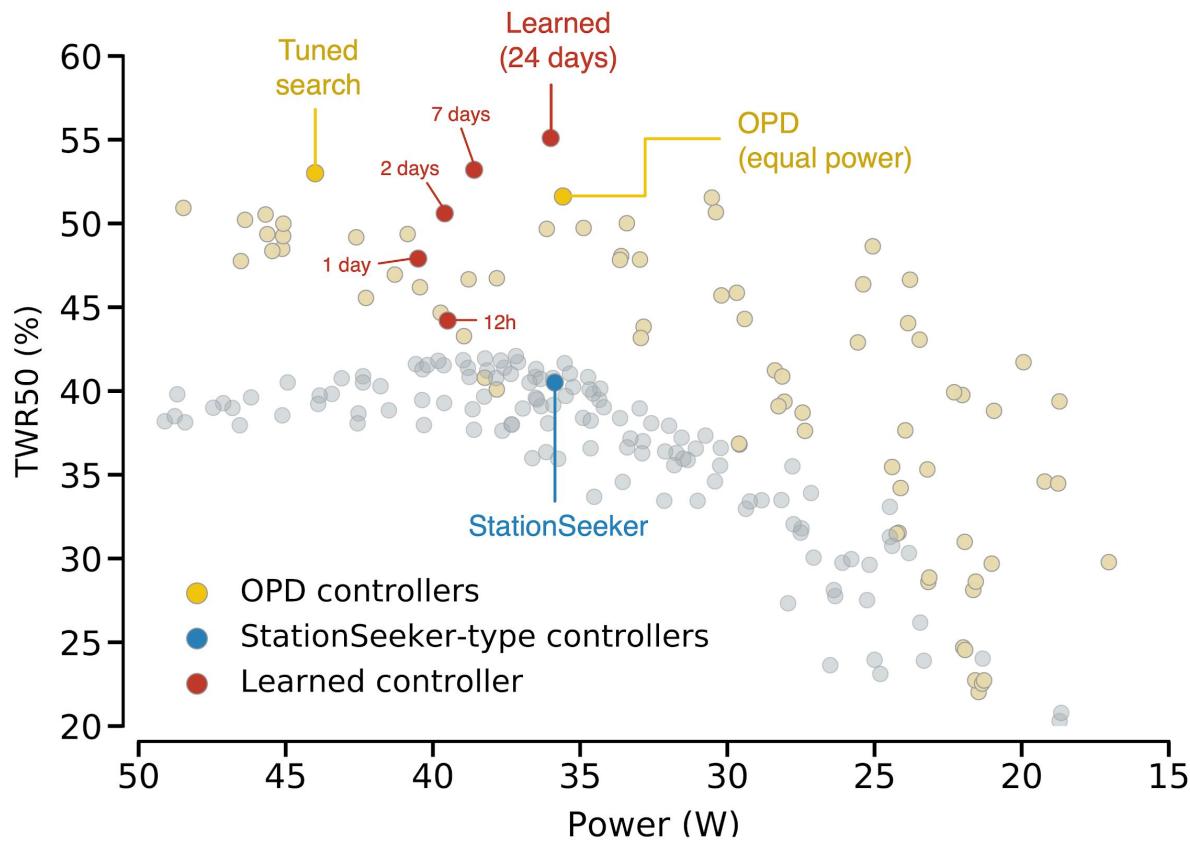
Distributed training:

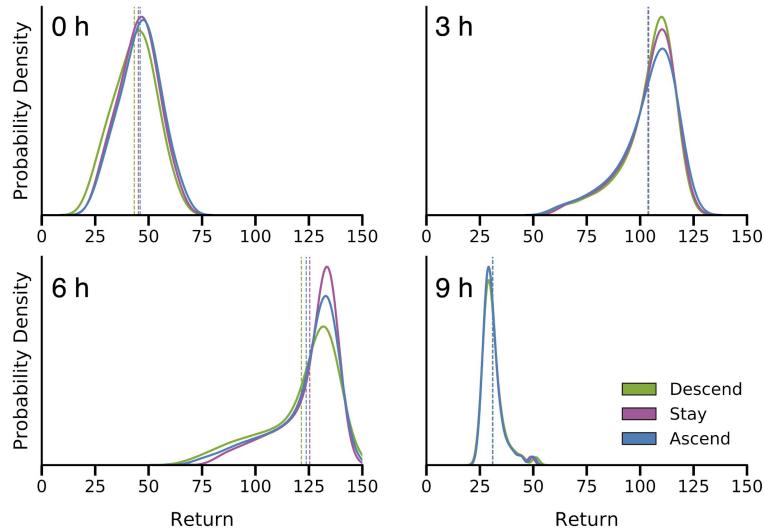
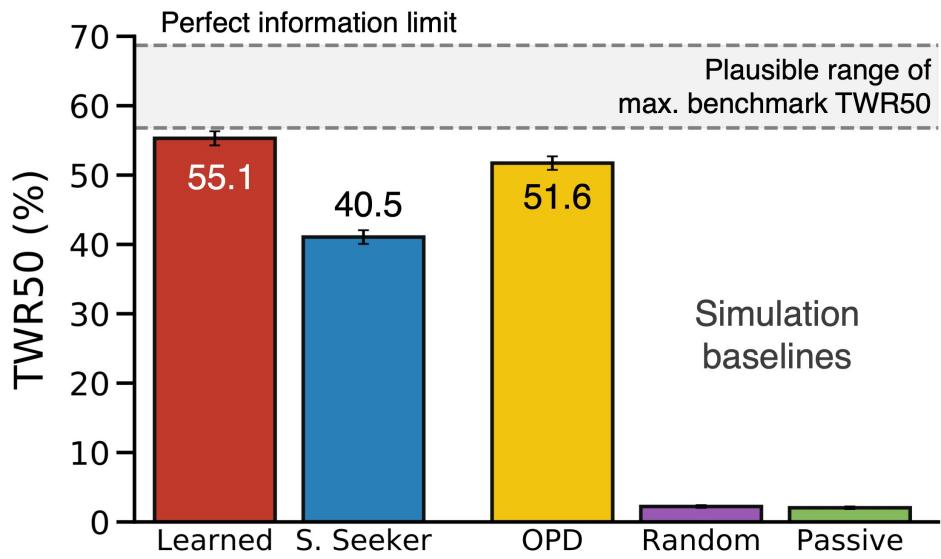
- 100 actors
- 4 replay buffers
- 1 GPU

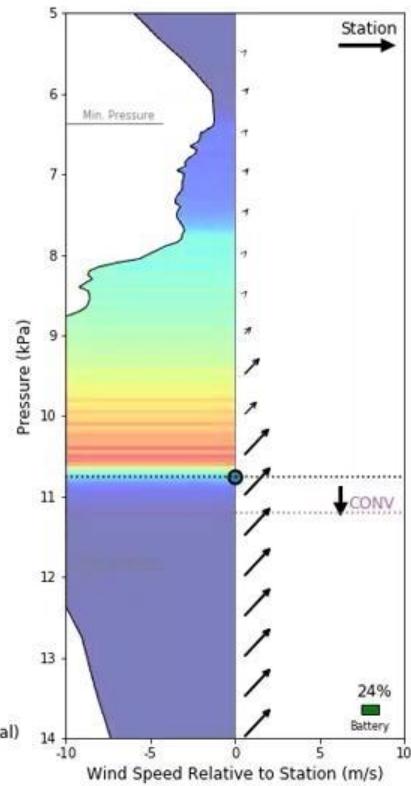
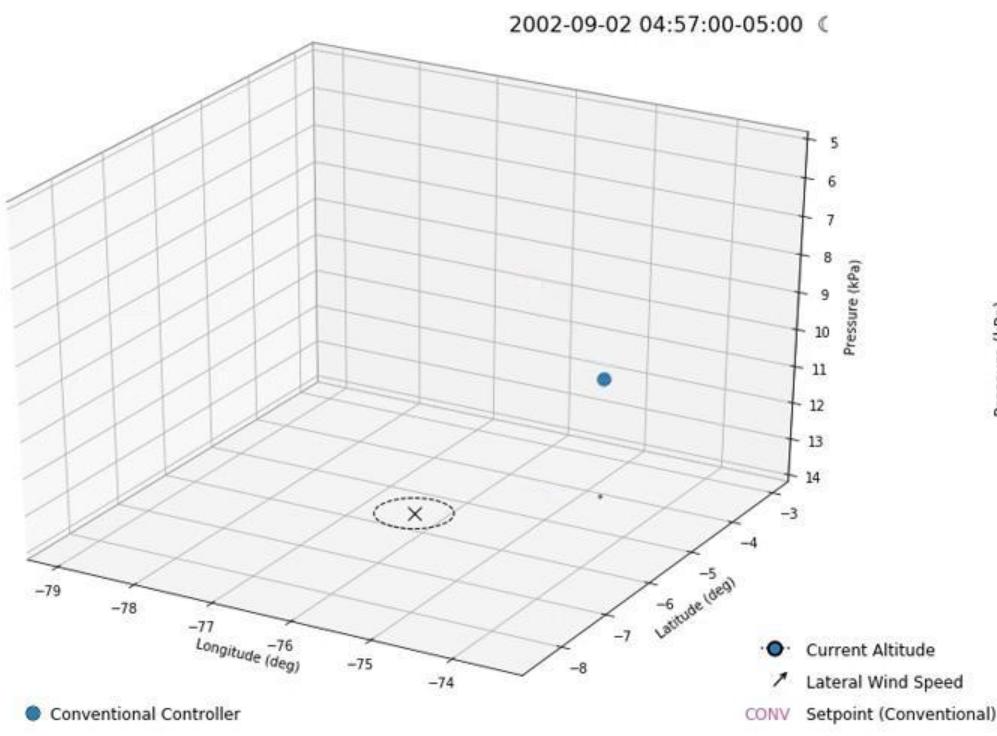
1.1B training steps (~30 days wall time)













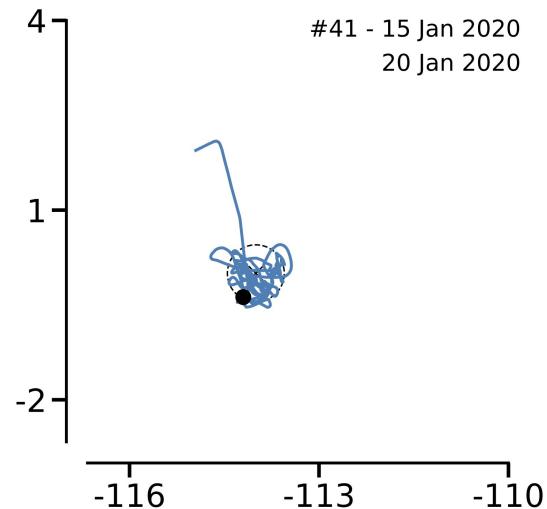
Pacific Ocean Experiment

26 Oct 2019 – 25 Jan 2020

13 balloons

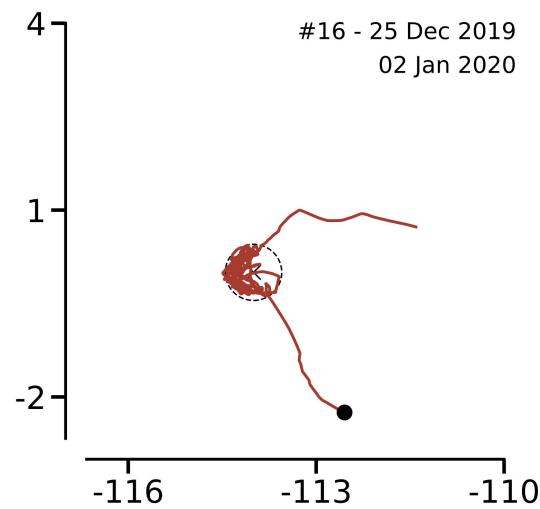
Total **2884** RL flight hours
Longest RL flight ~**16** days

StationSeeker



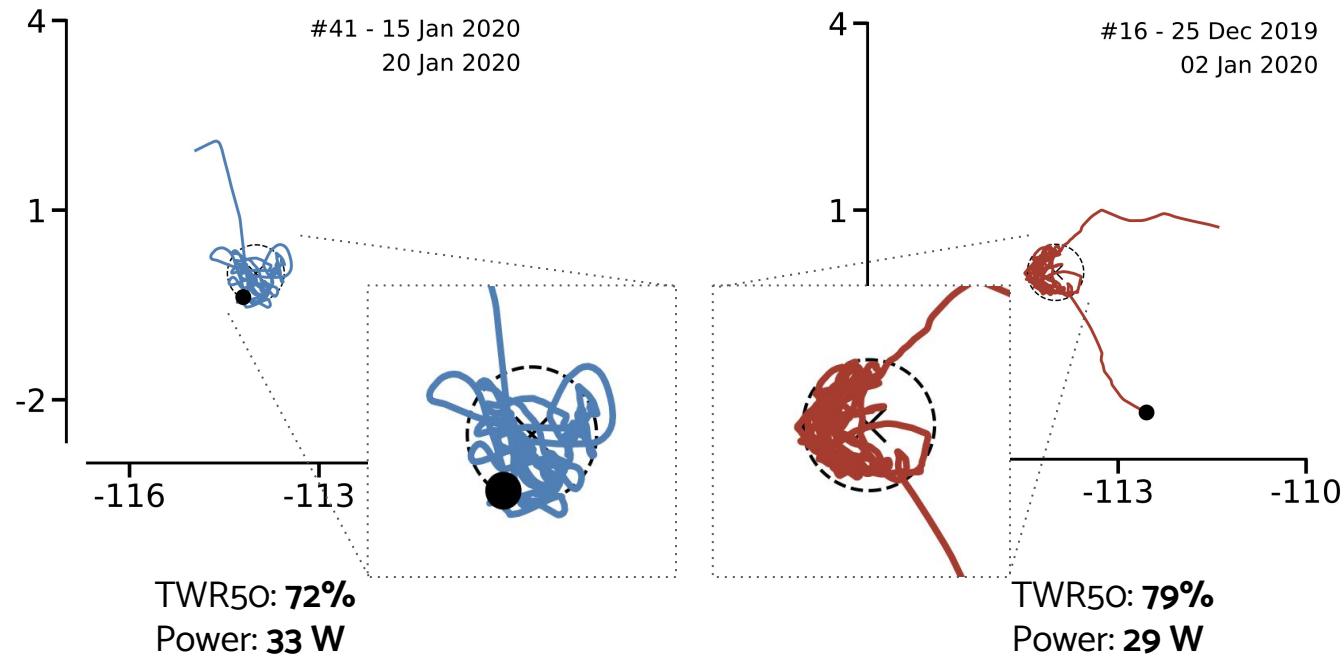
TWR50: **72%**
Power: **33 W**

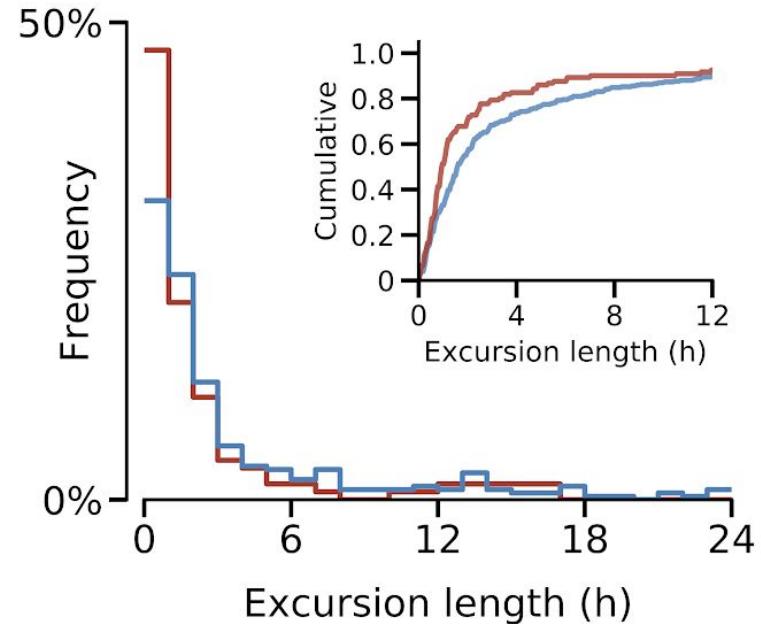
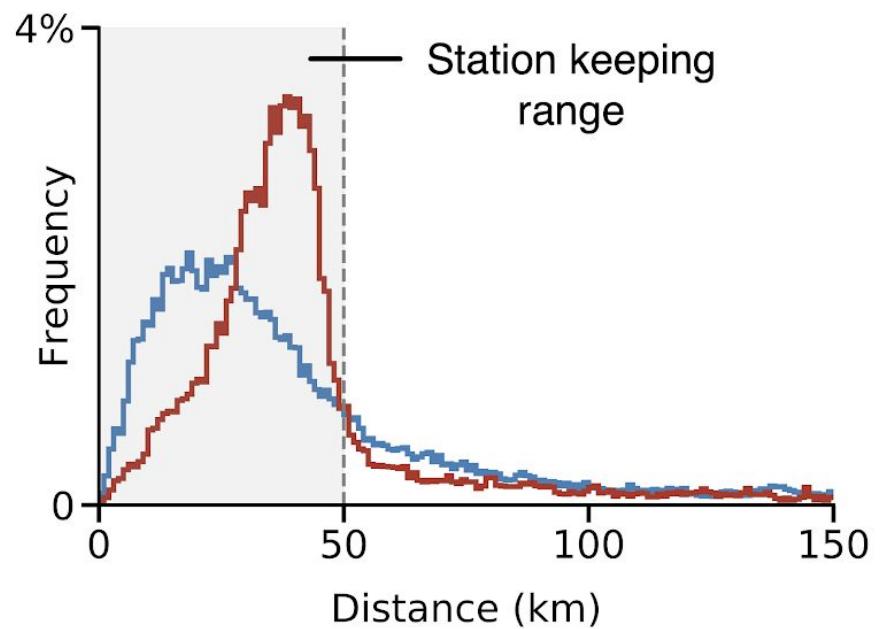
Perciatelli

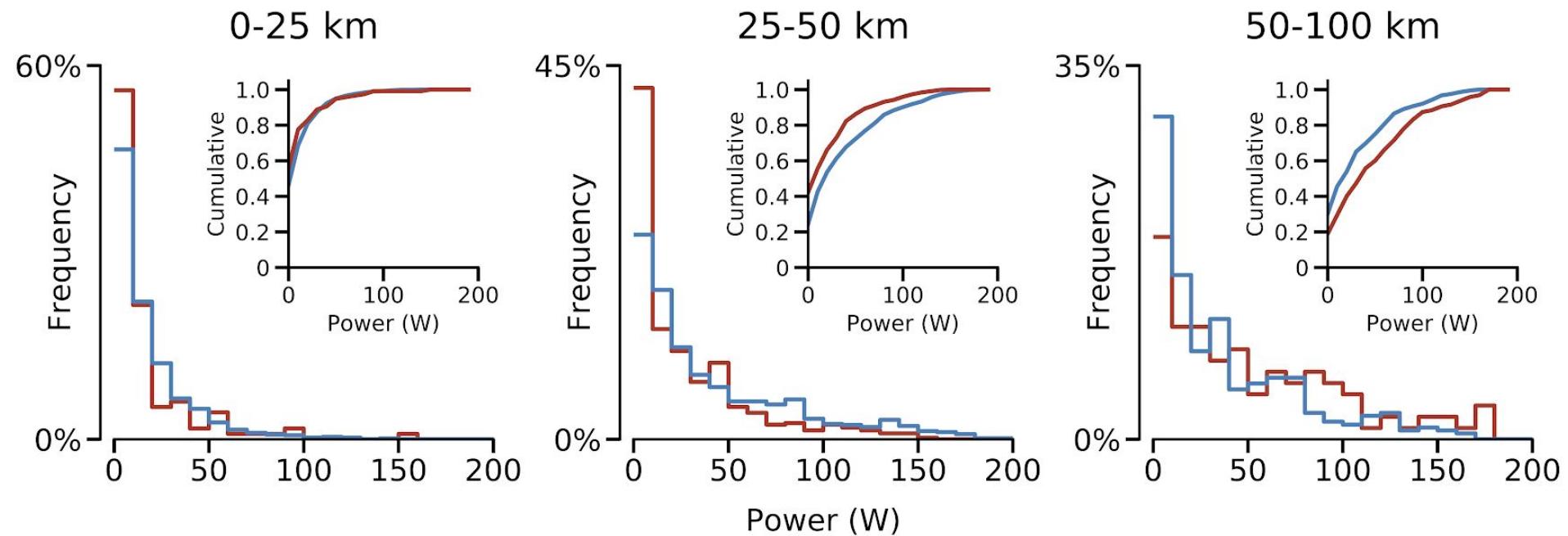


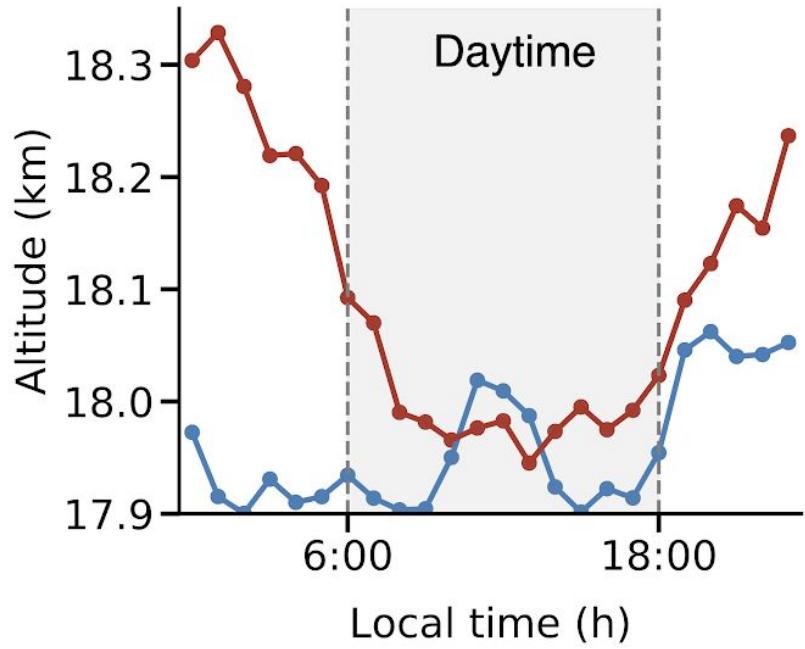
TWR50: **79%**
Power: **29 W**

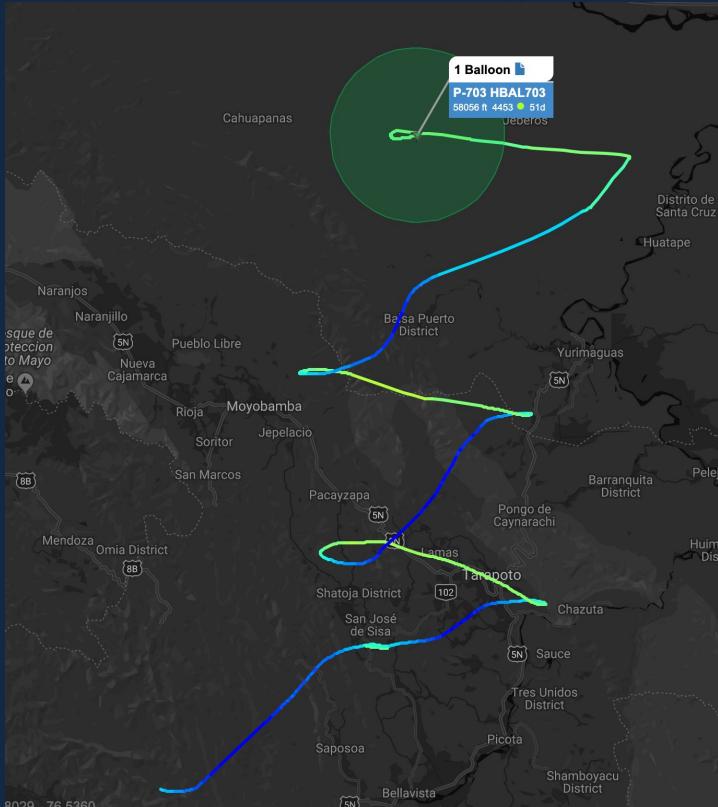
StationSeeker





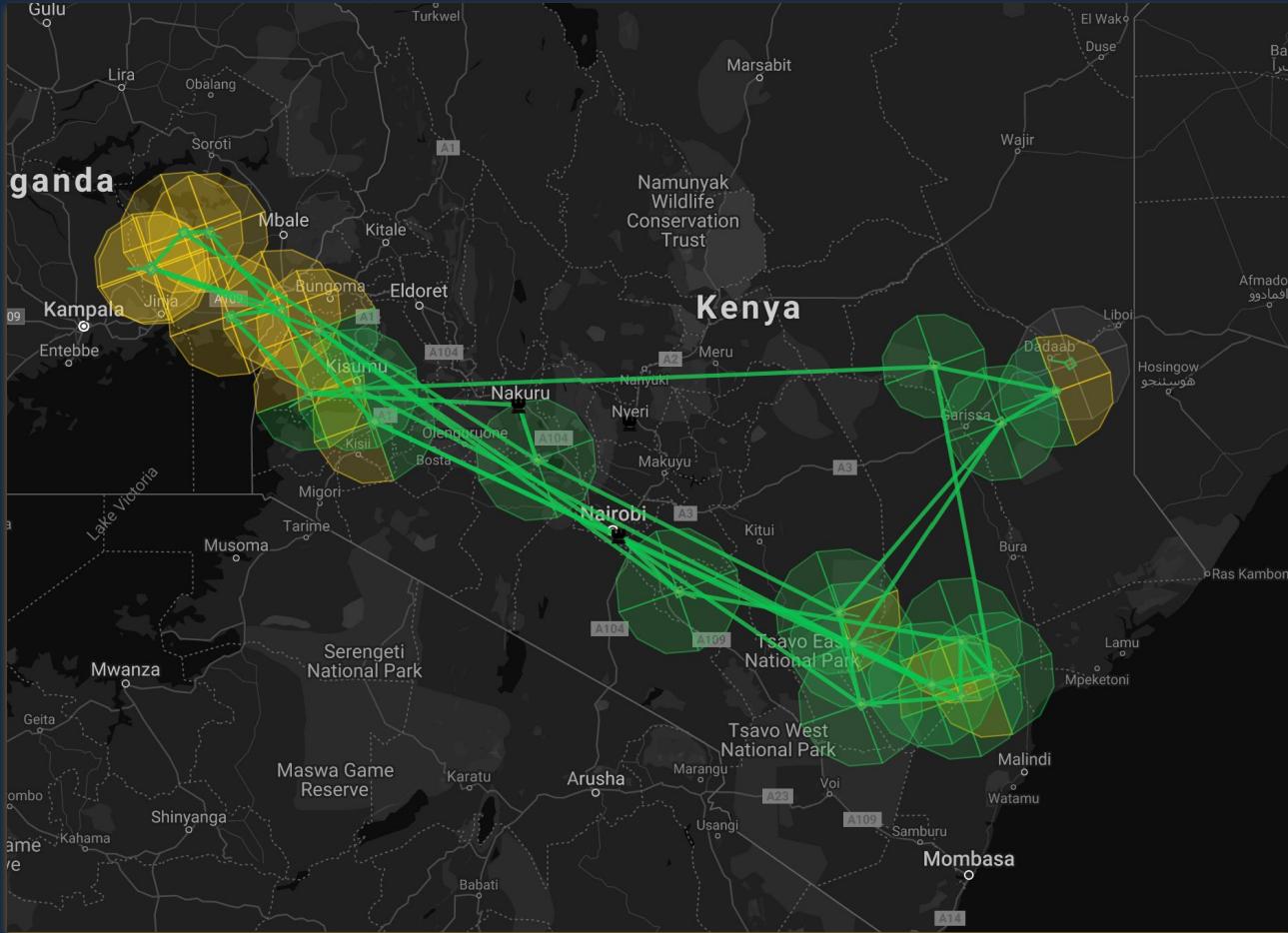


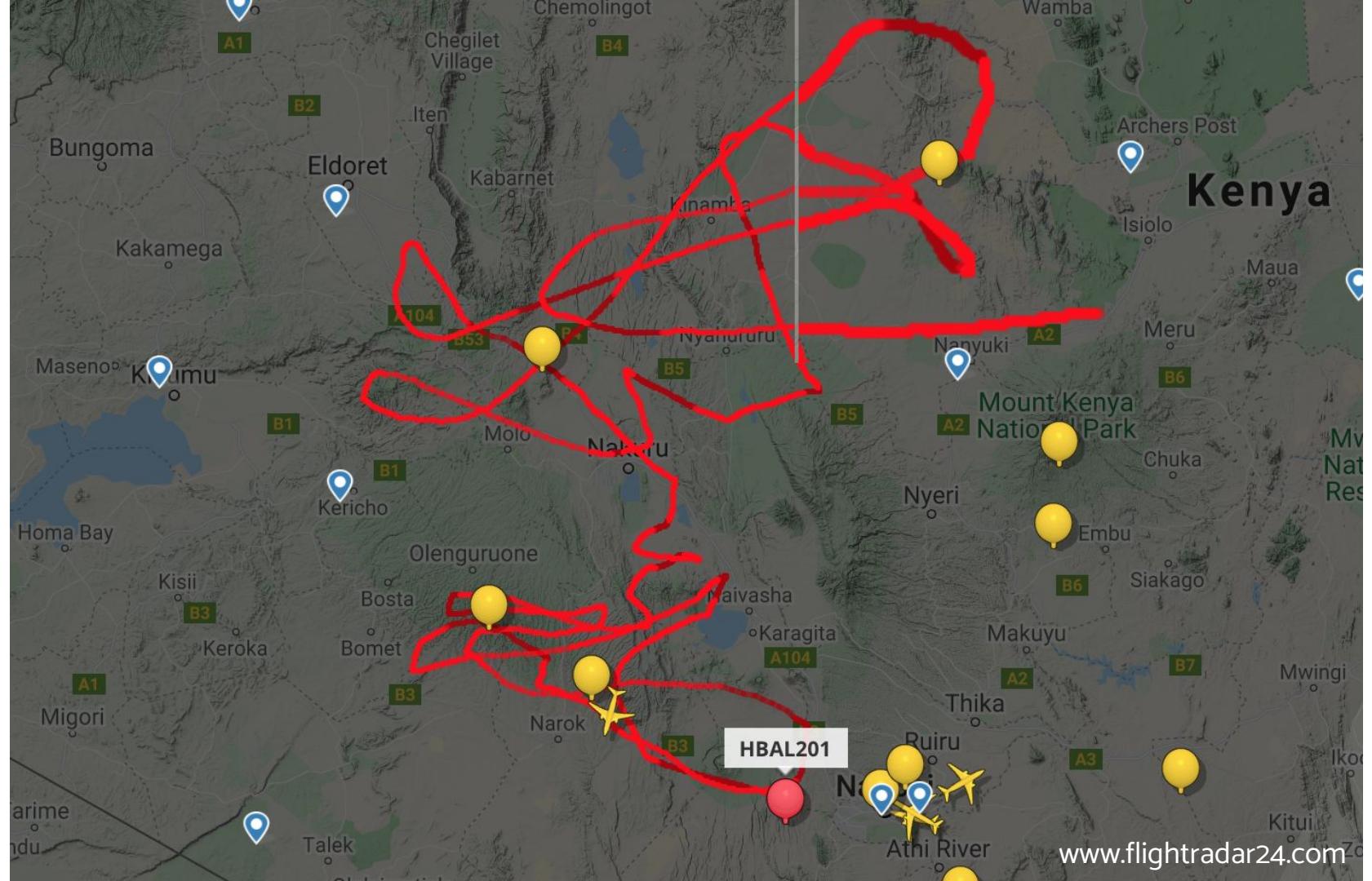






[312 Days in the Stratosphere](#), Loon, Oct 28 2020.





Why does this matter?

Deep reinforcement learning



Silver, Huang, Maddison, et al. (2016, 2017)



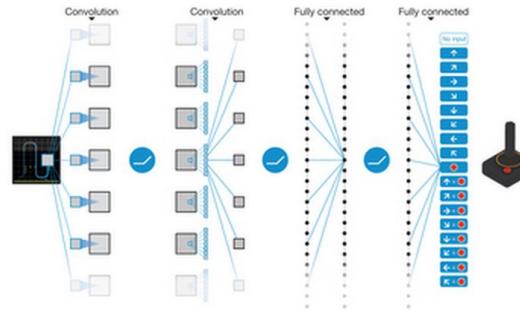
Bellemare, Naddaf, Veness, Bowling (2013)



Bard, Foerster, Chandar, et al. (2020)



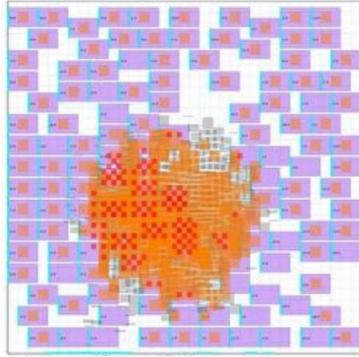
Levine et al., (2016)
Kalashnikov, Irpan, et al. (2018)



Mnih, Silver, Kavukcuoglu, Rusu, Veness, Bellemare, et al. (*Nature*, 2015)



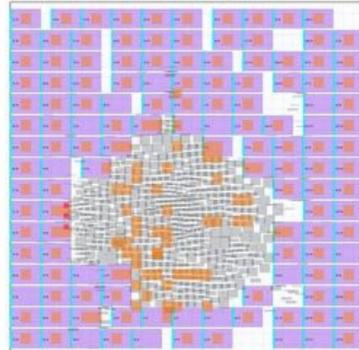
OpenAI et al. (2019)



Mirhoseini, Goldie, et al. (arXiv, 2020)



Won et al., 2020

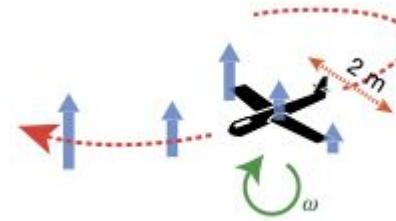


Managing Fragmented Fire-Threatened Landscapes with Spatial Externalities

Christopher J. Lauer, Claire A. Montgomery, and Thomas G. Dietterich



Glavic et al., 2017



Reddy et al. (2018)

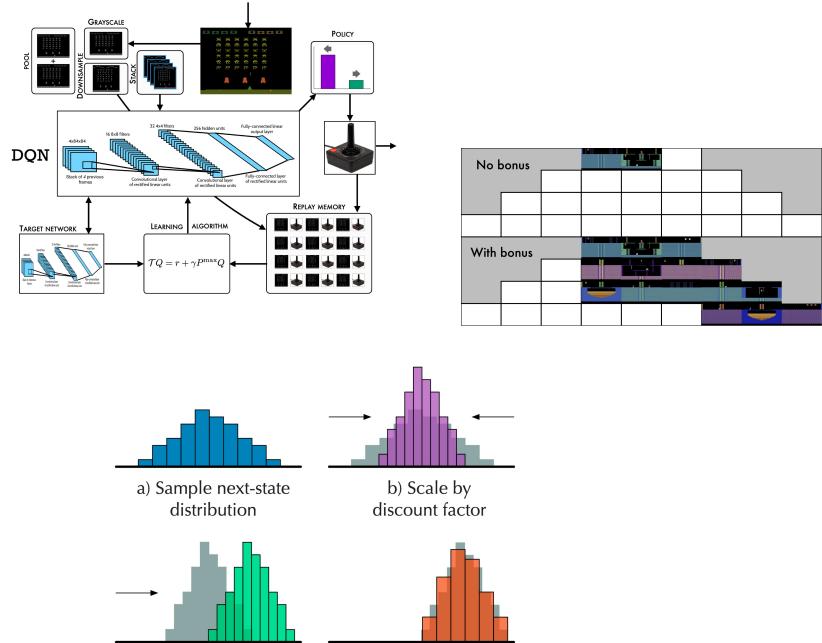


Ie et al., 2019

Challenging our paradigms



The Arcade Learning Environment, Bellemare, Naddaf, Veness, Bowling (2013).



“Real-world applications drive fundamental research questions”

- Chris Bishop

Decisions from data Controlling complex systems with reinforcement learning

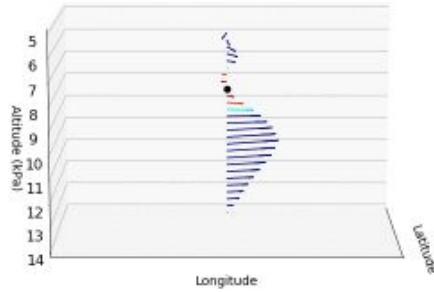
Marc G. Bellemare

Google Brain, Montreal; Mila

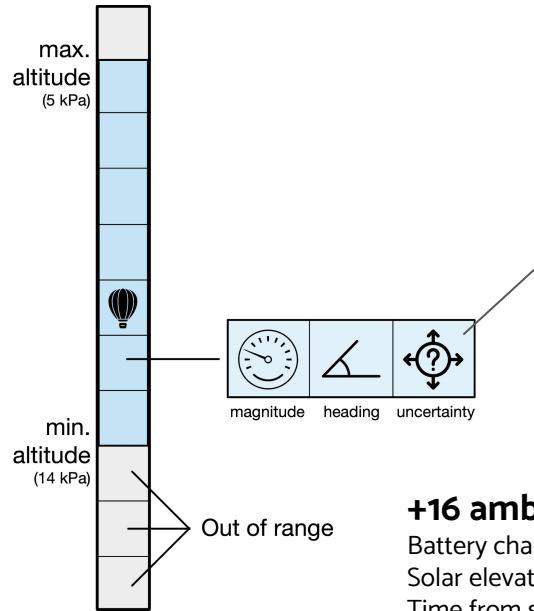
With: Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C. Machado, Subhodeep Moitra, Sameera S. Ponda, Ziyu Wang



RL for balloon navigation: inputs

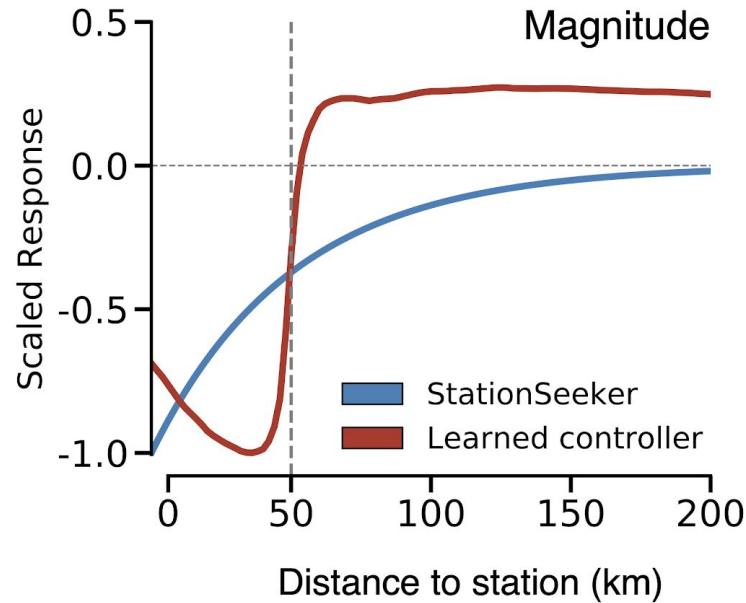
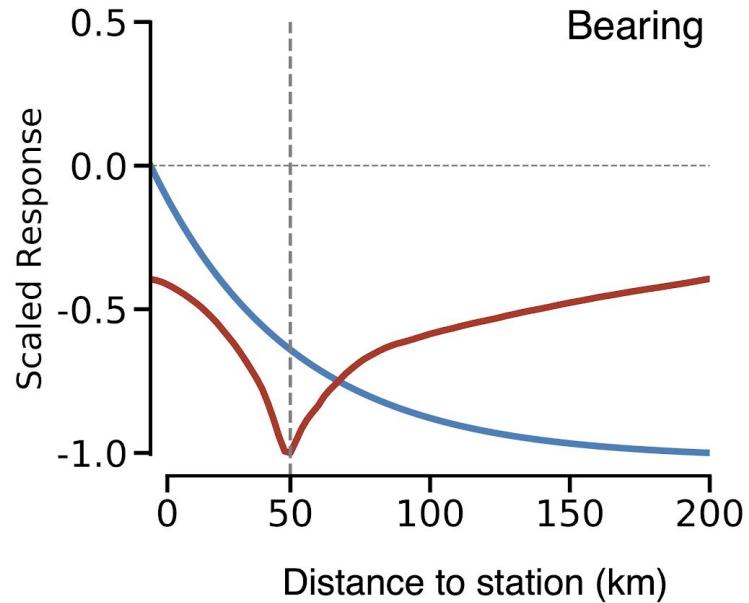


Forecast +
measurements +
Gaussian process =
wind column



The uncertainty term at
each altitude level acts as a
belief feature

+16 ambient variables:
Battery charge
Solar elevation
Time from sunrise
Distance to station
Internal pressure ratio
etc.



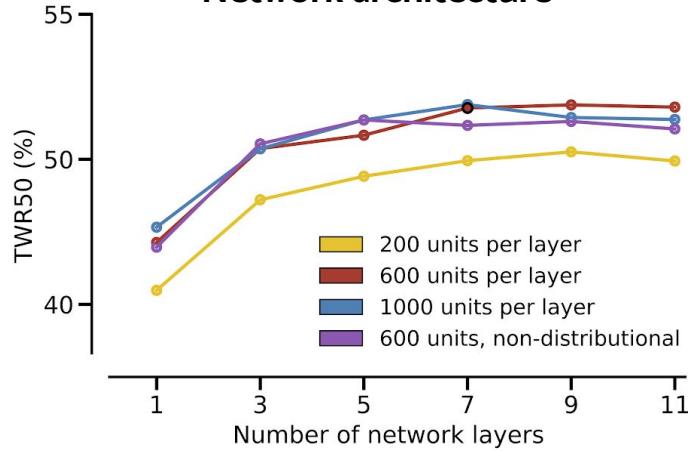
Sensitivity analysis

Twelve 2-day flights

Estimate **gradient of value function** w.r.t. wind characteristics (finite differences)

Plot average gradient **as a function of distance**

Network architecture



Reward parameters

