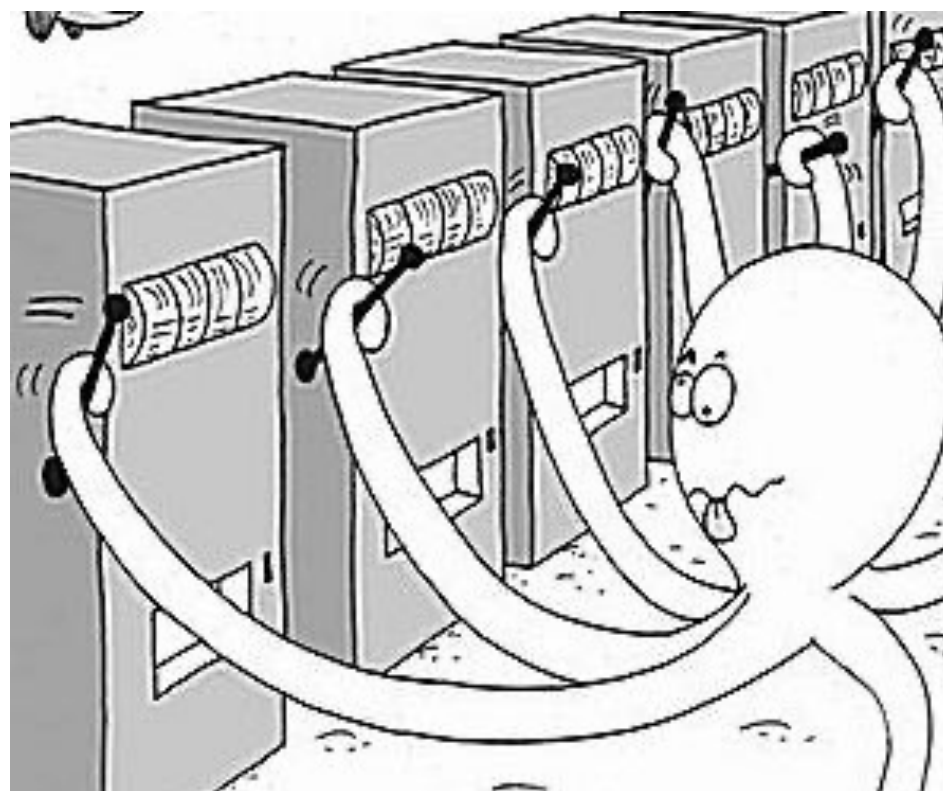# CSE 574 Planning and Learning Methods in AI
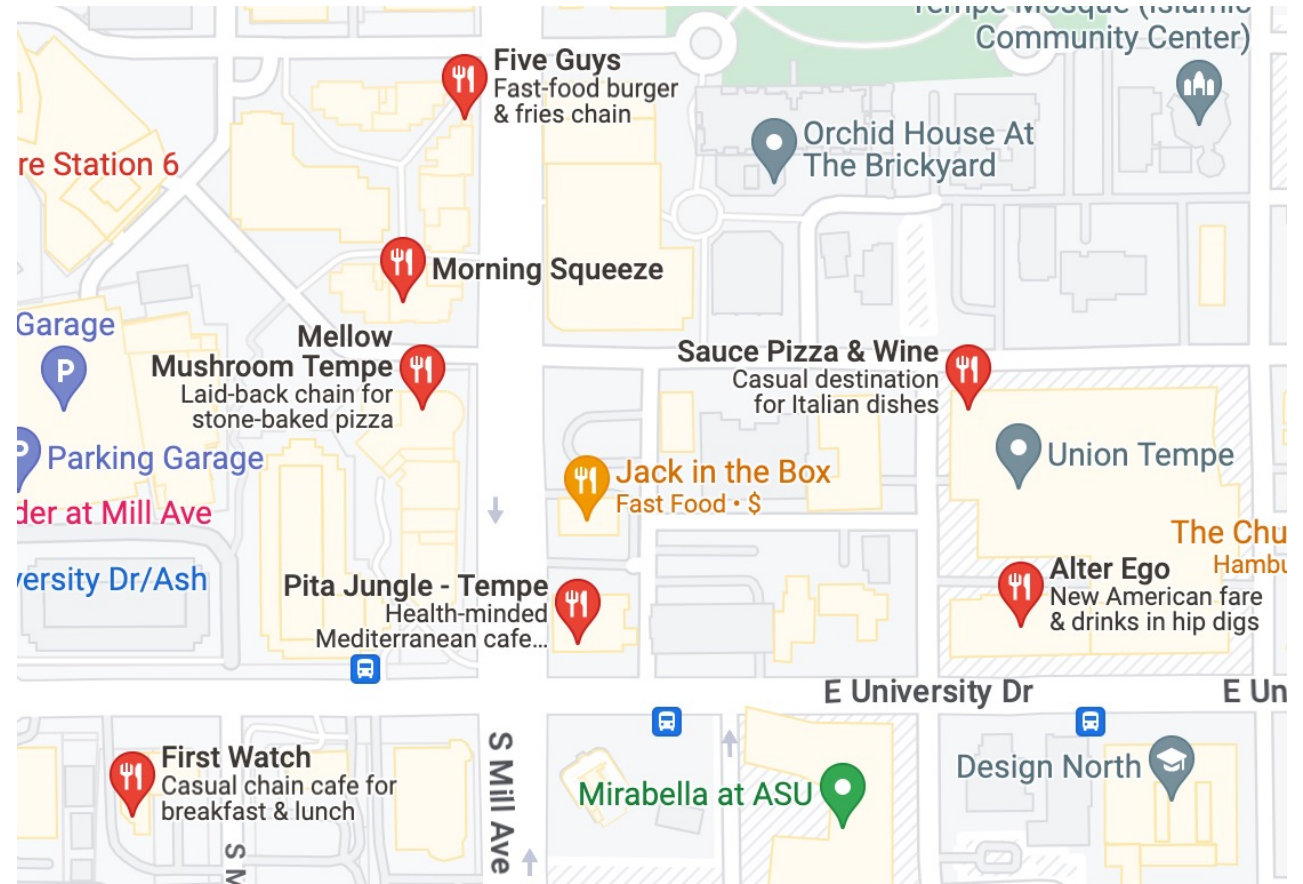
Ransalu Senanayake

Week 3

# Multi-Arm Bandits

# Multi-Arm Bandits Solutions

- Exploitation only
- Exploration only (greedy)
- $\epsilon$-first (exploration-first)
- $\epsilon$-greedy
- UCB

# $\epsilon$-greedy

$$\text{arm}_t = \begin{cases} \text{arm that maximizes reward, with probability } 1 - \epsilon \\ \text{random arm, with probability } \epsilon \end{cases}$$

- What is the effect of $\epsilon$?
- Fixed $\epsilon$ (e.g., $\epsilon$=10%), decreasing $\epsilon$, adaptive $\epsilon$, etc.
- Can we utilize more information than the average?

# UCB1 Algorithm

Randomly pull arms **k={1,…,K}** several times $(n)$ to get an initial estimate of <u>expected</u> rewards $\bar{r}_k$

For iteration **t=1,…,T**

    Play machine $k_{t+1} = argmax\left(\bar{r}_k + \alpha\sqrt{\frac{2\,logN_t}{n_k}}\right)$

end

**<u>Expected regret</u>**

Initial phase (figuring out the reward from each arm):   $\mathcal{O}(\sqrt{KTlogT})$

Later phase (when we get to know about arms/$\delta r_k$):   $\mathcal{O}\left(\sum_k \frac{1}{\delta r_k} logT\right)$

$\delta r_k$ is the reward gap of the kth arm compared to the arm with the best reward

# Bayesian Bandits and Thompson Sampling

```
Assume parameterized distributions for the prior and likelihood

For T iterations
   Compute the posterior p(θ|D) ∝ p(D|θ)p(θ)
   Sample parameters from each arm
   Compute the reward for each sample
   Pick the arm that maximizes the reward
   Append the dataset with D={(arm,reward)}
end
```
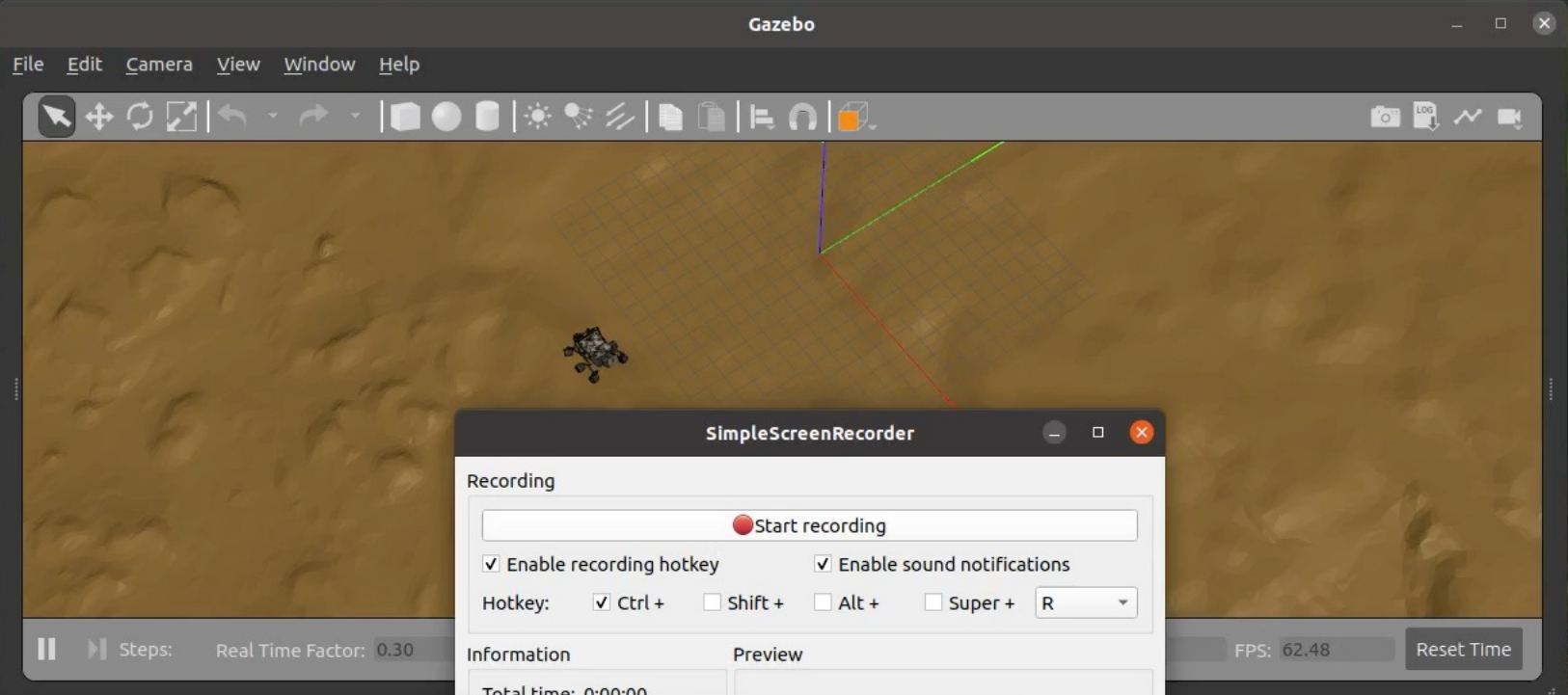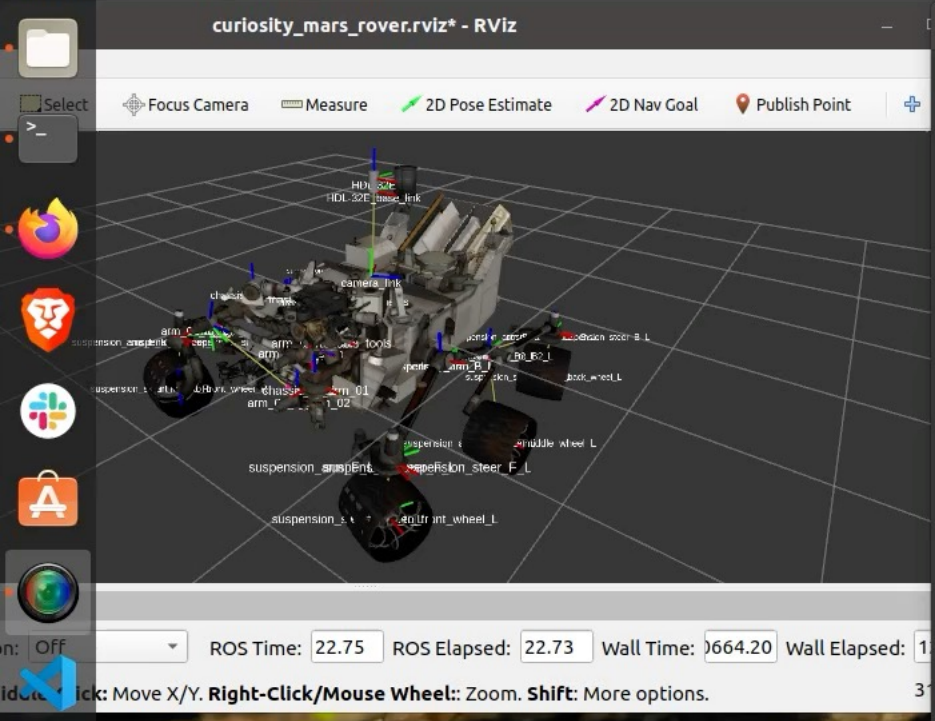
- Non-informative/uniform/flat/broad prior. Conjugate prior.

# Multi-Objective Optimization
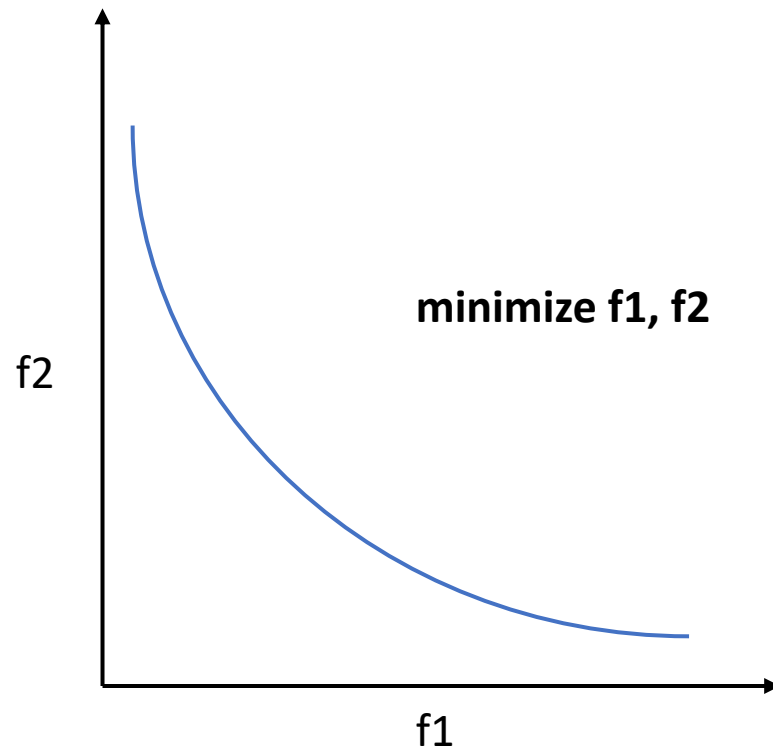
# Multi-Objective Optimization

- A choice is *Pareto optimal* if it is impossible to improve in one objective without worsening at least one other objective

**minimize f1, f2**

f2

f1

- $wf_1(x) + (1 - w)f_2(x)$
- Possible solutions
- What if we maximize
- Hypervolume

# Multi-Objective Bayesian Optimization (MOBO)

## Differentiable Expected Hypervolume Improvement for Parallel Multi-Objective Bayesian Optimization

**Samuel Daulton**
Facebook
sdaulton@fb.com

**Maximilian Balandat**
Facebook
balandat@fb.com

**Eytan Bakshy**
Facebook
ebakshy@fb.com