

ECE 558

Lecture 1 August 25th

Focus:

(1) Centralised Stochastic Control

- Stochastic dynamical system (!)
- One Controller (!)
- Controller has perfect recall (!)

(2) Structural and Foundational aspects

Why?

Algorithmic ideas will be developed but explored only to some extent.

Algorithms Focus of Reinforcement Learning (RL)

CSE 598 Satinder Banja FA2025

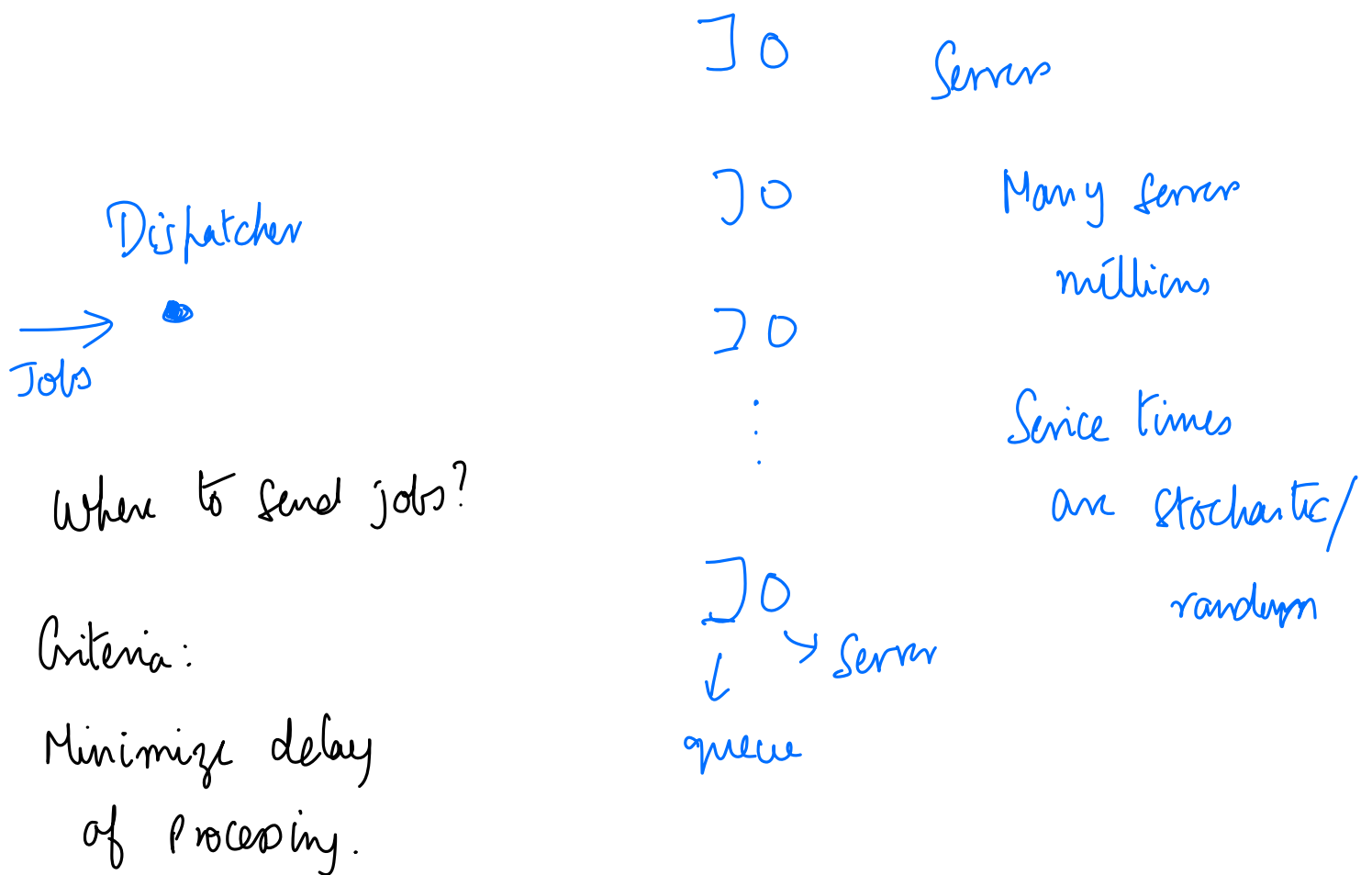
ECE 602 RL LEI YING WN 2026

Motivation: lots of applications Controls, Communications, Networks, Signal Processing, OR, Math, Finance, Economics,

Statistics, Manufacturing, Robotics, AI, ..., Quantum

Examples:

1. Scheduling in Cloud Systems/ Data Centers



Perspective A: Dispatcher sees full state, know
ability of servers

Based on job type, send the job to a specific server!

MARKOV DECISION PROCESS (MDP)

Perspective B: SENDING INFORMATION IS COSTLY

Dispatcher polls for information

\Rightarrow Poll a fixed number of servers

Partial observability - imperfect information

PARTIALLY OBSERVED MDPs (POMDPs)

Poll a fixed number of servers

(i) Soft constraint - On the average meet this

(ii) Hard constraint - Never allowed to exceed this.

CONSTRAINED MDPs + POMDPs (More advanced, could be a project topic)

2. Controlling the trajectory of a spacecraft,
autonomous car, robot

System State known - MDP (State "feedback")

If sensors used - POMDP (Output "feedback")

3. Portfolio management

Have some money on day t - x_t

Have K instruments to choose from.

$K+1$ options - Choosing distribution π_t on $K+1$ options

Option k gets $\pi_t^k x_t$ amount of money

y_{t+1}^k - a stochastic outcome

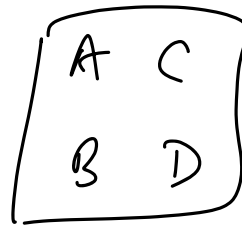
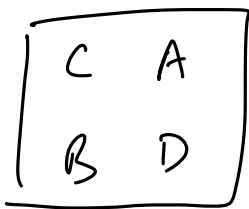
Function of $\pi_t^k x_t$ + Market behaviour

$$x_{t+1} = \sum_{k=0}^K y_{t+1}^k \quad (\text{Stochastic})$$

Goal: T periods, π_t 's each day to maximize wealth at the end of the period, but minimizing the chance of going bankrupt.

4. Active Hypothesis Testing

Compare 2 images to find differences (small differences)



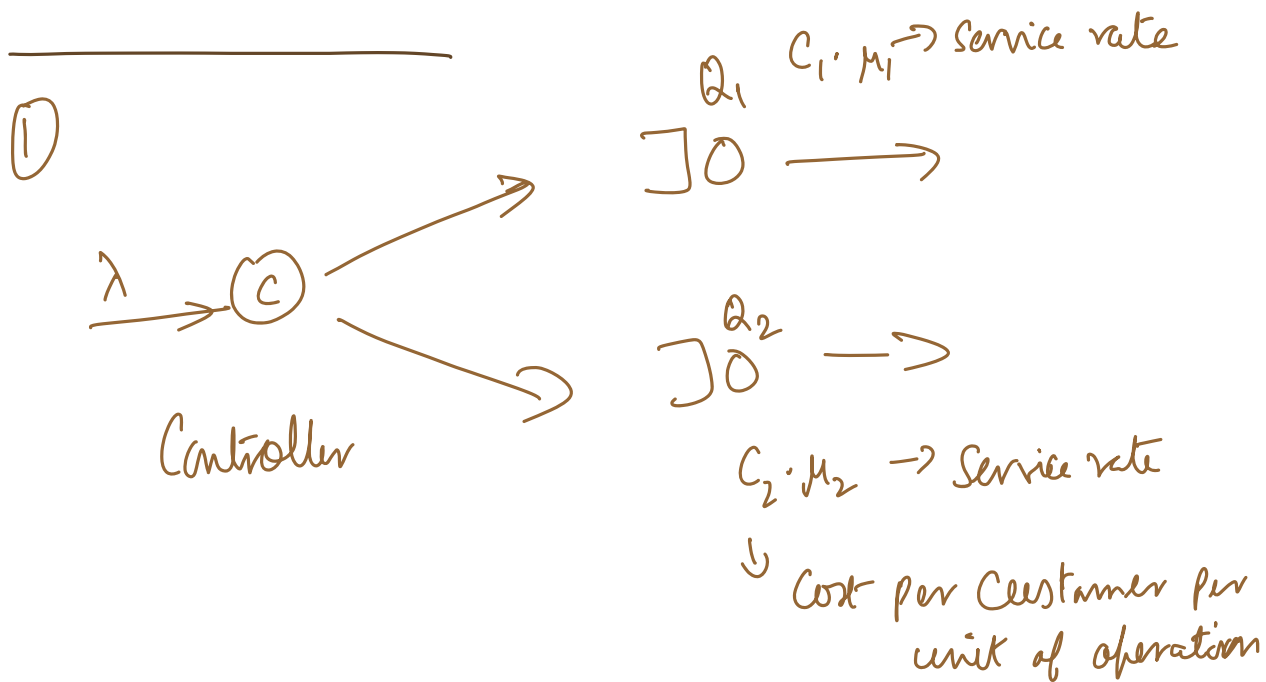
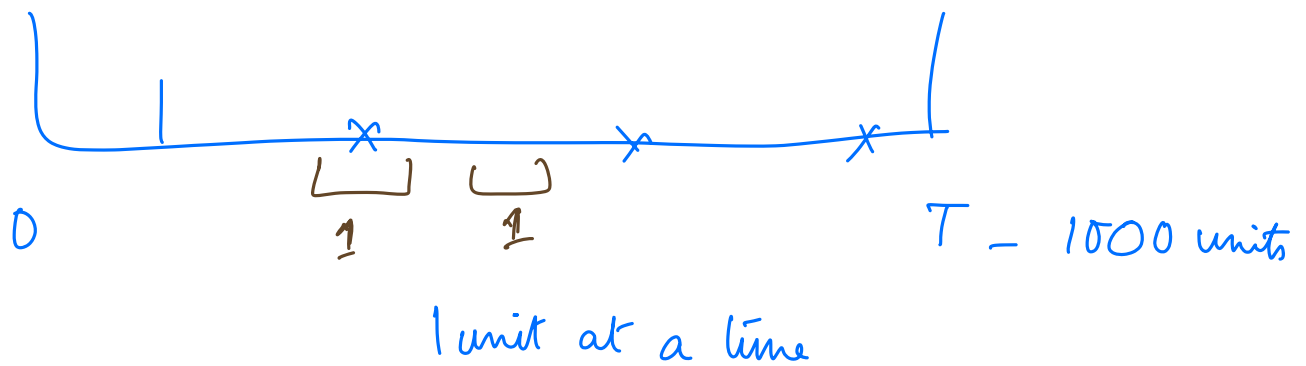
(KNOWN MODEL)

Goal: Tell apart the images (if different) + identify differences. FAST

How? Take different views at different granularities

Choice { High resolution over a small region
Low resolution over a bigger region

Choosing rent resolution & region based on past choices & results.



Network operates until time T

Occupancy of each queue at time $t \in \{1, 2, \dots, T\}$
 $\triangleq [T]$

x_t^1 in Q_1

x_t^2 in Q_2

$A_t = \begin{cases} \text{Route to } Q_1 & Q_1 \\ \text{Route to } Q_2 & Q_2 \end{cases}$

$A_t \in \{Q_1, Q_2\}$

Goal: Find a routing policy

$$g = (g_1, g_2, g_3, \dots, g_T)$$

that minimizes

$$E^g \left[\sum_{t=1}^T (C_1 x_t^1 + C_2 x_t^2) \right]$$

$$g_t = (x_1^1, x_1^2), (x_2^1, x_2^2), \dots, (x_t^1, x_t^2)$$

(Past states values including the present)

$$A_1, A_2, \dots, A_{t-1}$$

(Past action values)

Controller using all the past information

PERFECT RECALL

Assmt. Queues are finite — max B

$$x_t^1 \in \{0, 1, 2, \dots, B\} \quad (BH)$$

$$g_t: I_t \rightarrow A_t \quad A_t = \{a_1, a_2\}$$

I_t takes values in a space that has size

$$((BH)^2)^t \times 2^{t-1}$$

$$\text{Number of functions} = 2^{((B+1)^2)^T \times 2^{T-1}}$$

- Paring down to simpler functions to choose from while ensuring good performance.
- DP - breaks the time dependence

② Spacecraft - Vector state

$$t \in \{0, 1, \dots\} \quad x_{t+1} = x_t + \underset{\substack{\uparrow \text{State} \\ \downarrow \text{Control}}}{u_t} + \underset{\uparrow \text{noise}}{w_t}$$

Observation $y_t = C x_t + v_t$ \rightarrow noise in sensors.

Define $y_{0:t} = (y_0, y_1, \dots, y_t)$

$$u_{0:t} = (u_0, u_1, \dots, u_t)$$

$$u_t = g_t(y_{0:t}, u_{0:t-1}) \quad \text{and at time } t$$

"Cost" is $x_t^* Q x_t + u_t^* R u_t \quad t < T$

$$x_T^* Q x_T \quad \text{if } t = T$$

X^* - Transpose of X (Hermitian transpose)

$$X = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}, \text{ then } X^* = \begin{bmatrix} 0 & 1 & 2 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}, \quad A^* = \begin{bmatrix} 0 & 2 \\ 1 & 3 \end{bmatrix}$$

Goal: Choose g to minimize

$$\mathbb{E}^g \left[\sum_{t=0}^{T-1} (X_t^* Q X_t + U_t^* R U_t) + X_T^* Q X_T \right]$$

$$Q = I$$

$$R = I$$

$$X_t^* Q X_t = \sum_{k=1}^K (X_t^k)^2$$

Covered - Single controller with perfect recall.

Discrete-time, but state discrete or continuous.

Not covered

- (1) Continuous time problems (Technically more demanding, conceptually same ideas apply)
- (2) Multi-controller problems (Conceptual leap)
 - (a) Team
 - (b) Games
 - (c) Games of teams.