

Topics in Reinforcement Learning:  
AlphaZero, ChatGPT, Neuro-Dynamic Programming,  
Model Predictive Control, Discrete Optimization, Applications  
Arizona State University  
Course CSE 691, Spring 2025

Links to Class Notes, Videolectures, and Slides at  
<http://web.mit.edu/dimitrib/www/RLbook.html>

Prof. Dimitri P. Bertsekas (dimitrib@mit.edu)  
and

Dr Yuchao Li (yuchaoli@asu.edu)

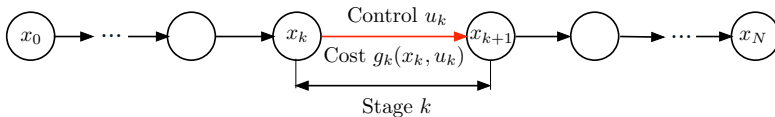
Video Credit: Alejandro P. Riveiros   William Emanuelsson   Pratyusha Musunuru

Lecture 7  
Applications of Approximation in Value Space to  
Multiagent and Multiple Objects Tracking Problems

- 1 Review: Approximation in Value Space - Truncated Rollout
- 2 Multiagent Rollout for Warehouse Robots Path Planning
- 3 Approximation in Value Space for Multiple Object Tracking/Data Association Problem

- 1 Review: Approximation in Value Space - Truncated Rollout
- 2 Multiagent Rollout for Warehouse Robots Path Planning
- 3 Approximation in Value Space for Multiple Object Tracking/Data Association Problem

## Review: DP Algorithm for Deterministic Problems



Go backward to compute the optimal costs  $J_k^*(x_k)$  of the  $x_k$ -tail subproblems (**off-line training** - involves lots of computation)

Go forward to construct optimal control sequence  $\{u_0^*, \dots, u_{N-1}^*\}$  (**on-line play**)

- Start with

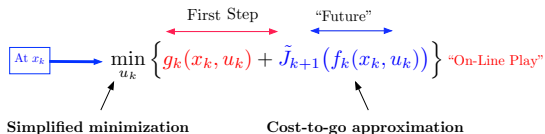
$$u_0^* \in \arg \min_{u_0 \in U_0(x_0)} \left[ g_0(x_0, u_0) + J_1^*(f_0(x_0, u_0)) \right], \quad x_1^* = f_0(x_0, u_0^*).$$

- Sequentially, going forward, for  $k = 1, 2, \dots, N-1$ , set

$$u_k^* \in \arg \min_{u_k \in U_k(x_k^*)} \left[ g_k(x_k^*, u_k) + J_{k+1}^*(f_k(x_k^*, u_k)) \right], \quad x_{k+1}^* = f_k(x_k^*, u_k^*).$$



# Review: Approximation in Value Space with One-Step Lookahead



We replace  $J_k^*(x_k)$  with an approximation  $\tilde{J}_k$  during on-line play

- Start with

$$\tilde{u}_0 \in \arg \min_{u_0 \in U_0(x_0)} \left[ g_0(x_0, u_0) + \tilde{J}_1(f_0(x_0, u_0)) \right].$$

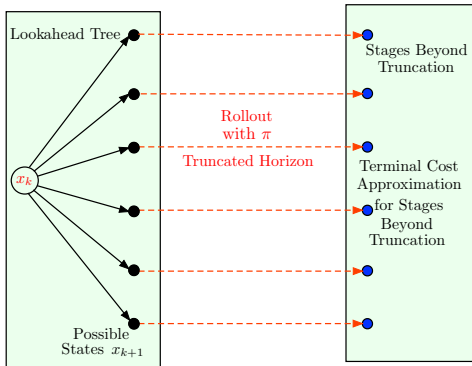
- We set  $\tilde{x}_1 = f_0(x_0, \tilde{u}_0)$
- Sequentially, going forward, for  $k = 1, 2, \dots, N - 1$ , set

$$\tilde{u}_k \in \arg \min_{u_k \in U_k(\tilde{x}_k)} \left[ g_k(\tilde{x}_k, u_k) + \tilde{J}_{k+1}(f_k(\tilde{x}_k, u_k)) \right], \quad \tilde{x}_{k+1} = f_k(\tilde{x}_k, \tilde{u}_k).$$

## Two Challenges

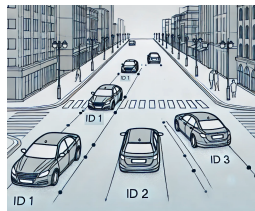
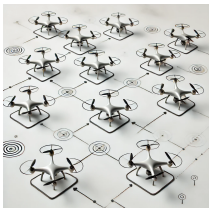
- How to construct the approximation  $\tilde{J}_{k+1}$  **offline**?
- How to solve the minimization problem **online**?

## Review: Truncated Rollout with One-Step Lookahead



- Truncated rollout encompasses all key ingredients: a **base policy  $\pi$**  and a **terminal cost approximation  $\hat{J}$** :  $\tilde{J}_{k+1}(x_{k+1}) = \hat{J}(x_{k+m+1}) + \sum_{\ell=k+1}^{k+m} g_{\ell}(x_{\ell}, \mu_{\ell}(x_{\ell}))$ .
- When no rollout is involved, the **terminal cost approximation  $\hat{J}$**  serves as  $\tilde{J}$
- When there is no truncation, we use  $J_{\pi}$  as  $\tilde{J}$
- However, the challenge of **computing the minimization remains!**

# Focus of This Lecture: Solving the Minimization in Approximation in Value Space Efficiently



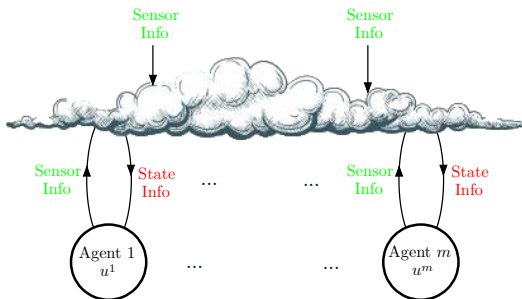
$$\boxed{\text{At } x_k} \longrightarrow \min_{u_k} \left\{ \overset{\text{First Step}}{g_k(x_k, u_k)} + \overset{\text{"Future"}}{\tilde{J}_{k+1}(f_k(x_k, u_k))} \right\} \text{ "On-Line Play"}$$

Simplified minimization                      Cost-to-go approximation

- We will consider **multiagent coordination** and **multi-object tracking** problems.
- The computational demand can be reduced by **leveraging the problem structure**, as in **multiagent rollout**, still yielding **a legitimate Newton's step**.
- The terminal cost approximation  $\tilde{J}$  can be suitably designed to **enhance the minimization calculation**.

- 1 Review: Approximation in Value Space - Truncated Rollout
- 2 Multiagent Rollout for Warehouse Robots Path Planning**
- 3 Approximation in Value Space for Multiple Object Tracking/Data Association Problem

# Review: A Classical Information Pattern



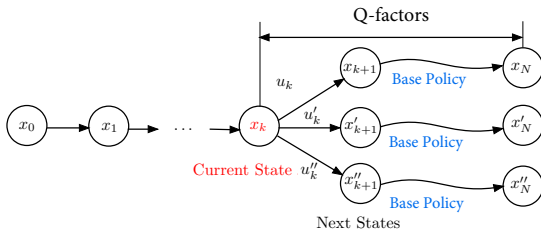
At each time: Agents have exact state info; choose their controls as a function of state

Model: A discrete-time (possibly stochastic) system with state  $x$  and control  $u$

- Decision/control has  $m$  components  $u = (u^1, \dots, u^m)$  corresponding to  $m$  "agents"
- "Agents" is just a metaphor - the important math structure is  $u = (u^1, \dots, u^m)$
- For every policy  $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$ , the functions  $\mu_k$ ,  $k = 0, \dots, N-1$  take the form

$$\mu_k(x) = (\mu_k^1(x), \mu_k^2(x), \dots, \mu_k^m(x))$$

# Standard (Truncated) Rollout for Multiagent Problem



## Standard (truncated) rollout for multiagent problem

- Suppose for simplicity that control constraint set  $U_k(x_k)$  has a Cartesian product structure:

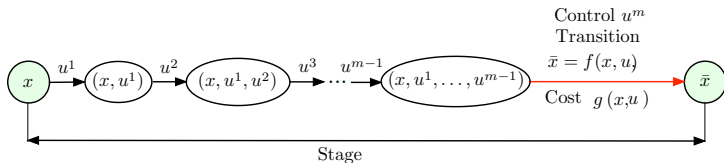
$$U_k(x_k) = U_k^1 \times U_k^2 \times \cdots \times U_k^m.$$

- Standard rollout, which **considers all agents at once**, involves the following minimization:

$$(\tilde{u}_k^1, \dots, \tilde{u}_k^m) \in \arg \min_{(u_k^1, \dots, u_k^m) \in U_k^1 \times U_k^2 \times \cdots \times U_k^m} Q_k(x_k, u_k^1, \dots, u_k^m).$$

- Computing each Q-factor requires simulation: **Infeasible even for modest  $m$ !**

# Multiagent Rollout - One Agent at a Time



Multiagent rollout requires **much less computation**

- At  $x_k$ , multiagent rollout **solves sequentially  $m$  minimization problems**:

$$\tilde{u}_k^1 \in \arg \min_{u_k^1 \in U_k^1} Q_k(x_k, u_k^1, \mu_k^2(x_k), \dots, \mu_k^m(x_k))$$

$$\tilde{u}_k^2 \in \arg \min_{u_k^2 \in U_k^2} Q_k(x_k, \tilde{u}_k^1, u_k^2, \mu_k^3(x_k), \dots, \mu_k^m(x_k))$$

$\vdots$

$$\tilde{u}_k^m \in \arg \min_{u_k^m \in U_k^m} Q_k(\tilde{u}_k^1, \dots, \tilde{u}_k^{m-1}, u_k^m)$$

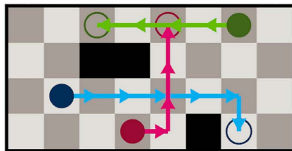
- The computational demand is reduced from  $n^m$  to  $n \cdot m$ !
- Multiagent rollout is standard rollout applied to the reformulated problem - **A legitimate Newton step!**

## A Fifteen-Minute Break

Catch our breath and think about issues relating to the first half of the lecture. Ask questions when you return.



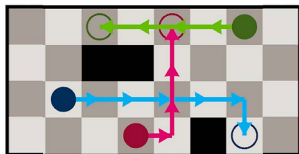
# Multiagent Path Finding Example: Modeling



warehouse robots path planning  $\implies$  grid world representation

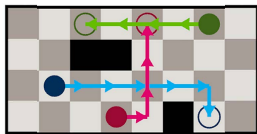
- There are  $m = 3$  agents (solid circles) moving in 4 directions or standing still **with perfect vision**
- The agents have been assigned to some targets (open circles with the same color).
- The objective: **reaching their respective targets in minimum time while avoiding collision with each other**
- Simple heuristic: each agent follows the shortest path to the respective target, assuming other agents are not present (arrows in the figure)

# Multiagent Path Finding Example: DP Formulation and Standard Rollout

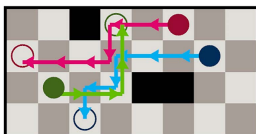


- States: current positions of all agents and their respective targets
- Control: each agent has at most 5 choices, their combination **grows exponentially with  $m$**
- Stage cost: related to the number of collisions and the number of reached targets
- When applying standard rollout, we must evaluate  $\approx 5^m$  approximate Q-factors

# Multiagent Rollout for Multiagent Path Finding



case 1



case 2

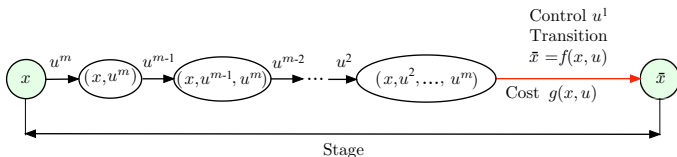


case 3

- Multiagent rollout reduces it to  $5 \cdot m$  (while maintaining good properties)
- Key idea: Break down the control into a sequence of one-agent-at-a-time moves
- Each stage involves the following sequence of operations:
  - ▶ Minimizing Q-factors associated with the first agent, while the remaining two agents follow base heuristics
  - ▶ Minimizing Q-factors associated with the second agent, while the last agent follows base heuristics
  - ▶ Minimizing Q-factors associated with the last agent
- We allow a change of the order in which the agents are selecting their controls, at every stage

Will be presented in class

# Implementation Variants of Multiagent Rollout



- Reshuffling the order of agents results in **a different, yet still equivalent** problem
- Multiagent rollout allows parallel computation of Q-factors
- Multiple base heuristics can be applied to enhance the performance further
- **All those ideas are independent of each other and can be combined**
- See textbook for additional material for **order optimization** and other variants

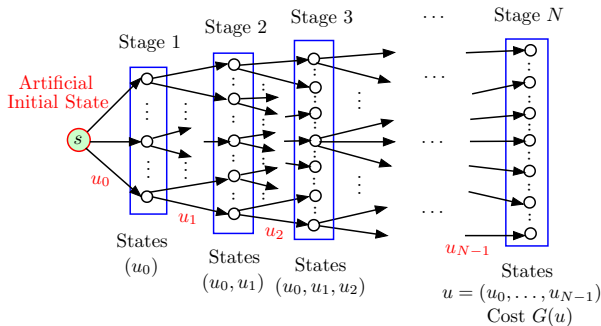
# Multiagent Rollout for Multiagent Path Finding: Animation For A Large Scale Problem

Will be presented in class

- Base policy is computed offline and **stored in a lookup table**
- Can adapt to changing environment: **some robots may breakdown halfway**
- Tested up to 200 robots in simulation, required computational time **less than 1 sec**
- Code can be found at <https://github.com/will-em/multi-agent-rollout>
- Paper: “Multiagent Rollout with Reshuffling for Warehouse Robots Path Planning”, by W. Emanuelsson et al., IFAC World Congress. Also see arXiv:2211.08201

- 1 Review: Approximation in Value Space - Truncated Rollout
- 2 Multiagent Rollout for Warehouse Robots Path Planning
- 3 Approximation in Value Space for Multiple Object Tracking/Data Association Problem

# General Discrete Optimization



Minimize  $G(u)$  subject to  $u \in U$

- Assume that **each solution  $u$  has  $N$  components:  $u_0, \dots, u_{N-1}$**
- View the components as the controls of  $N$  stages
- Define  $x_k = (u_0, \dots, u_{k-1})$ ,  $k = 1, \dots, N$ , and introduce artificial start state  $x_0 = s$
- Define just terminal cost as  $G(u)$ ; all other costs are 0

**This formulation typically makes little sense for exact DP, but often makes a lot of sense for approximate DP/approximation in value space**



## DP solution to the discrete optimization problem

- Start with

$$J_N^*(x_N) = G(x_N) = G(u_0, \dots, u_{N-1}) \quad \text{for all } x_N \in U$$

- For  $k = 0, \dots, N - 1$ , let

$$J_k^*(x_k) = \min_{u_k \in U_k(x_k)} J_{k+1}^*(x_k, u_k) \quad \text{for all } x_k$$

where  $U_k(x_k)$  need to be suitably defined.

- Construct the optimal solution  $(u_0^*, \dots, u_{N-1}^*)$  by forward calculation

$$u_k^* \in \arg \min_{u_k \in U_k(x_k)} J_{k+1}^*(x_k, u_k) \quad \text{for all } x_k$$

## Approximation in value space

- Use some  $\tilde{J}_{k+1}$  in place of  $J_{k+1}^*$
- Starting from the artificial initial state, for  $k = 0, \dots, N - 1$ , set

$$\tilde{u}_k \in \arg \min_{u_k \in U_k(x_k)} \tilde{J}_{k+1}(x_k, u_k) \quad \text{for all } x_k$$

# Multiple Object Tracking

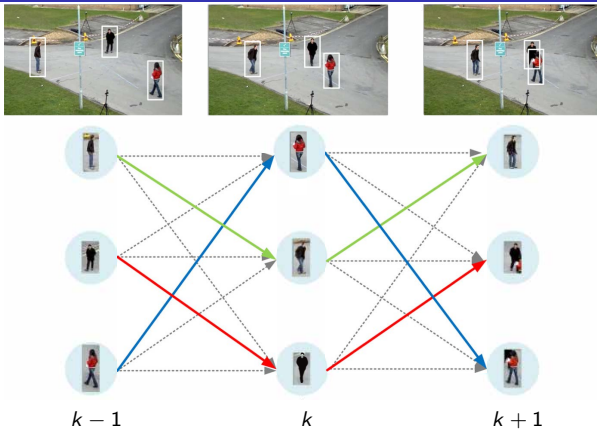
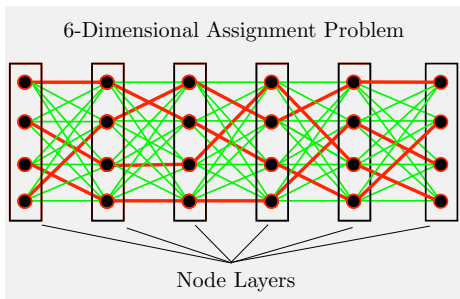


Figure source: [CGI22]<sup>1</sup>

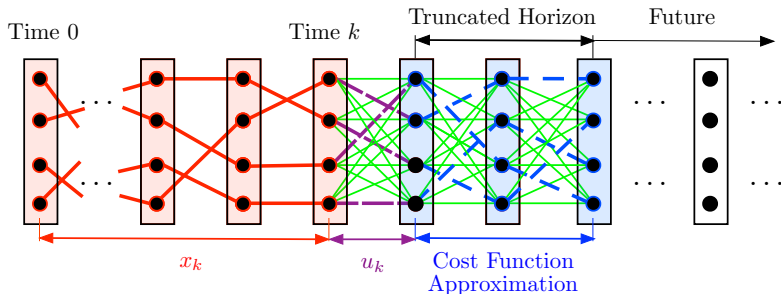
- Multiple object tracking (MOT) aims to match the same objects over various frames
- Nontrivial: **occlusion, changes in object appearance, and real-time computation constraint**
- Important problem with many applications: **traffic monitoring, robotics, consumer analytics, augmented and virtual realities ...**

# Multidimensional Assignment Problem



- MOT can be modeled as a multidimensional assignment problem
- There are  $(N + 1)$  layers (frames) of nodes
- A **grouping** consists of  $N + 1$  nodes  $(i_0, \dots, i_N)$  where  $i_k$  belongs to  $k$ th layer, and  $N$  corresponding arcs
- For each grouping, there is an associated cost **depending on the entire grouping**
- Our goal: find  $m$  groupings so that each node belongs to **one and only one** grouping and **the sum of the costs of the groupings is minimized**

# Approximation in Value Space



- Approximation in value space involves the following key ideas:
  - ▶ One-step lookahead minimization
  - ▶ Truncated rollout
  - ▶ Cost approximation  $\tilde{J}$  with structure that matches the assignment problem
- Q-factor minimization reduces to **solving a 2-dimensional assignment problem**
- Results: **robust and consistent matching against occlusion**
- Paper: “An Approximate Dynamic Programming Framework for Occlusion-Robust Multi-Object Tracking”, by P. Musunuru et al., arXiv:2405.15137

Will be presented in class

## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



## MOT Example: Base Heuristic





## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



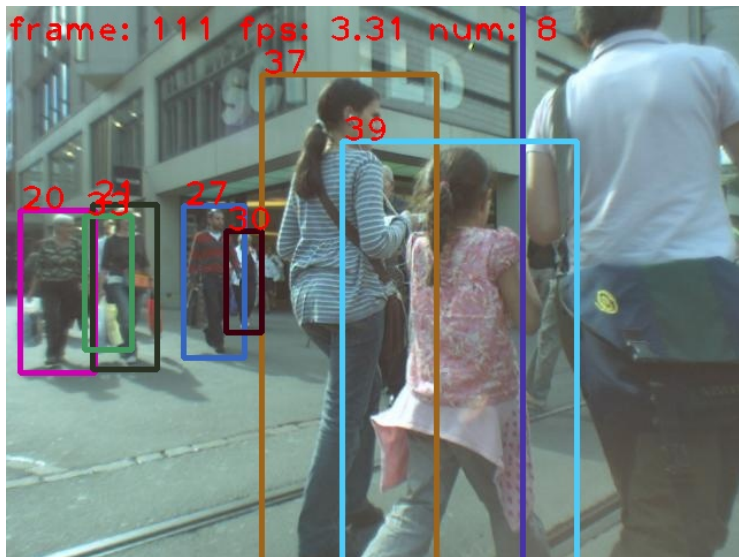
## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



## MOT Example: Base Heuristic

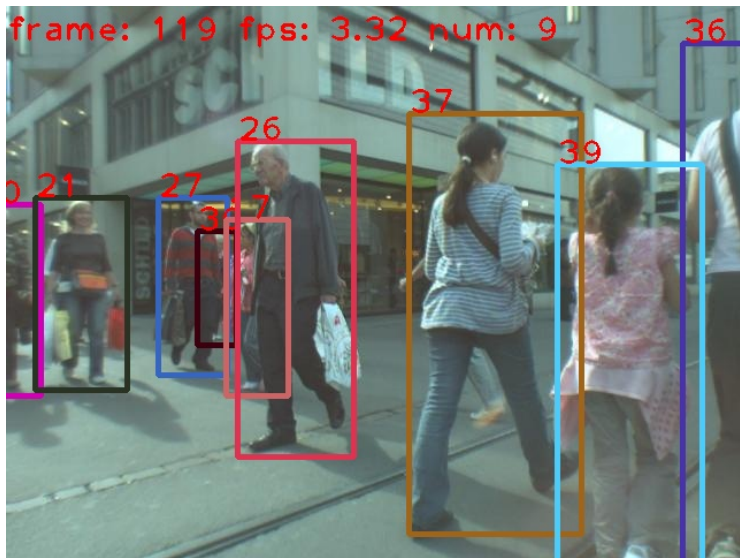


## MOT Example: Base Heuristic

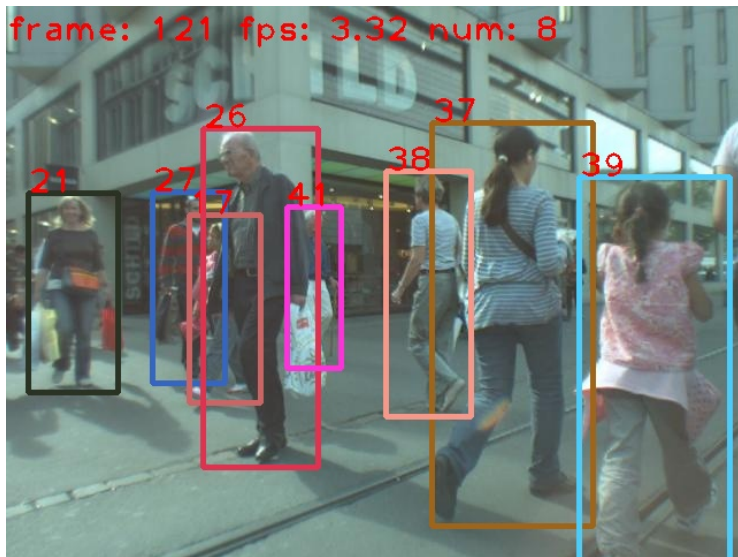




## MOT Example: Base Heuristic



## MOT Example: Base Heuristic



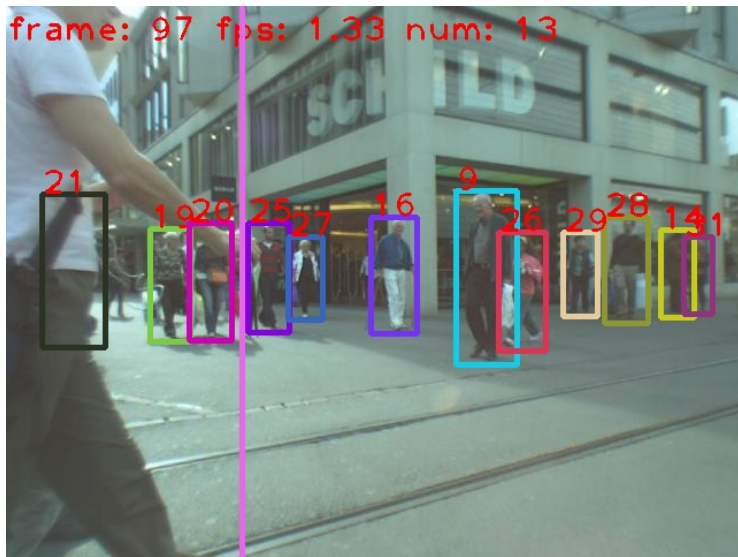
## MOT Example: Base Heuristic



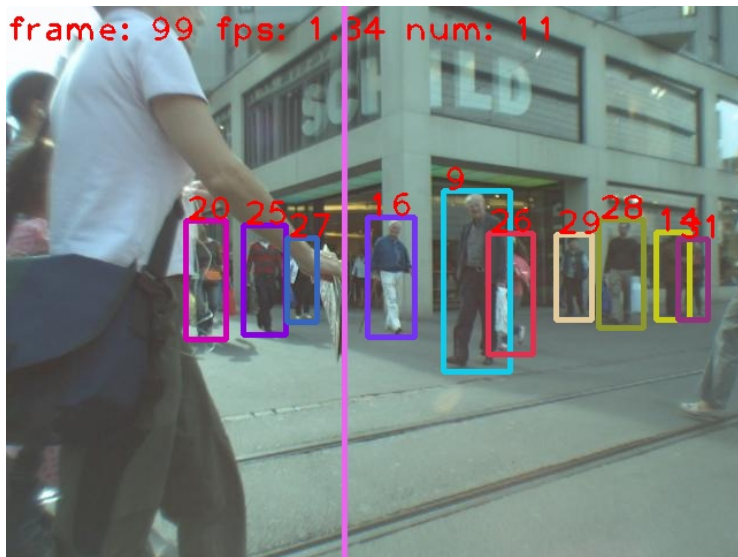
## MOT Example: Base Heuristic



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space





## MOT Example: Approximation in Value Space



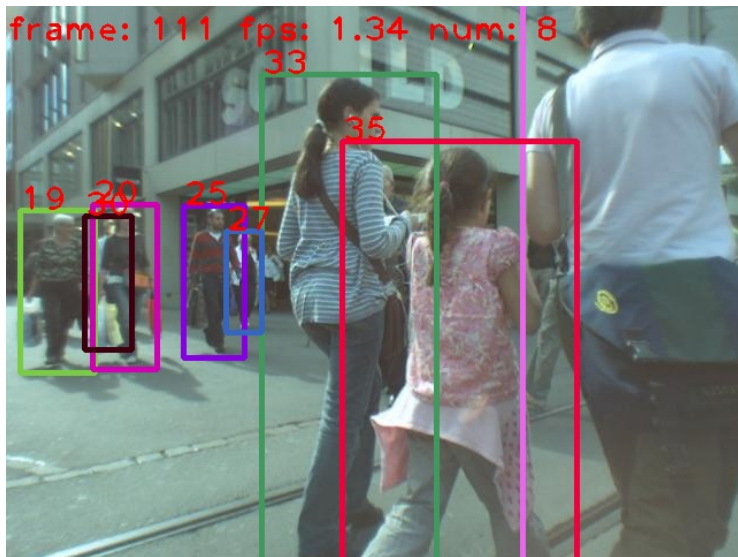
## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space

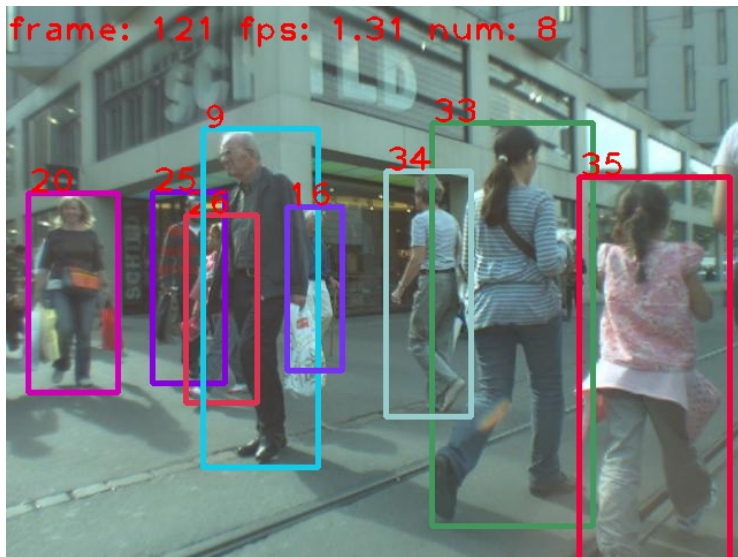


## MOT Example: Approximation in Value Space





## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



## MOT Example: Approximation in Value Space



In the next lecture we will cover:

- Most likely sequence generated by  $n$ -grams and HMM
- Rollout algorithm for approximate solution
- Application to large language models