



Deep Reinforcement Learning for Double Auction Processes

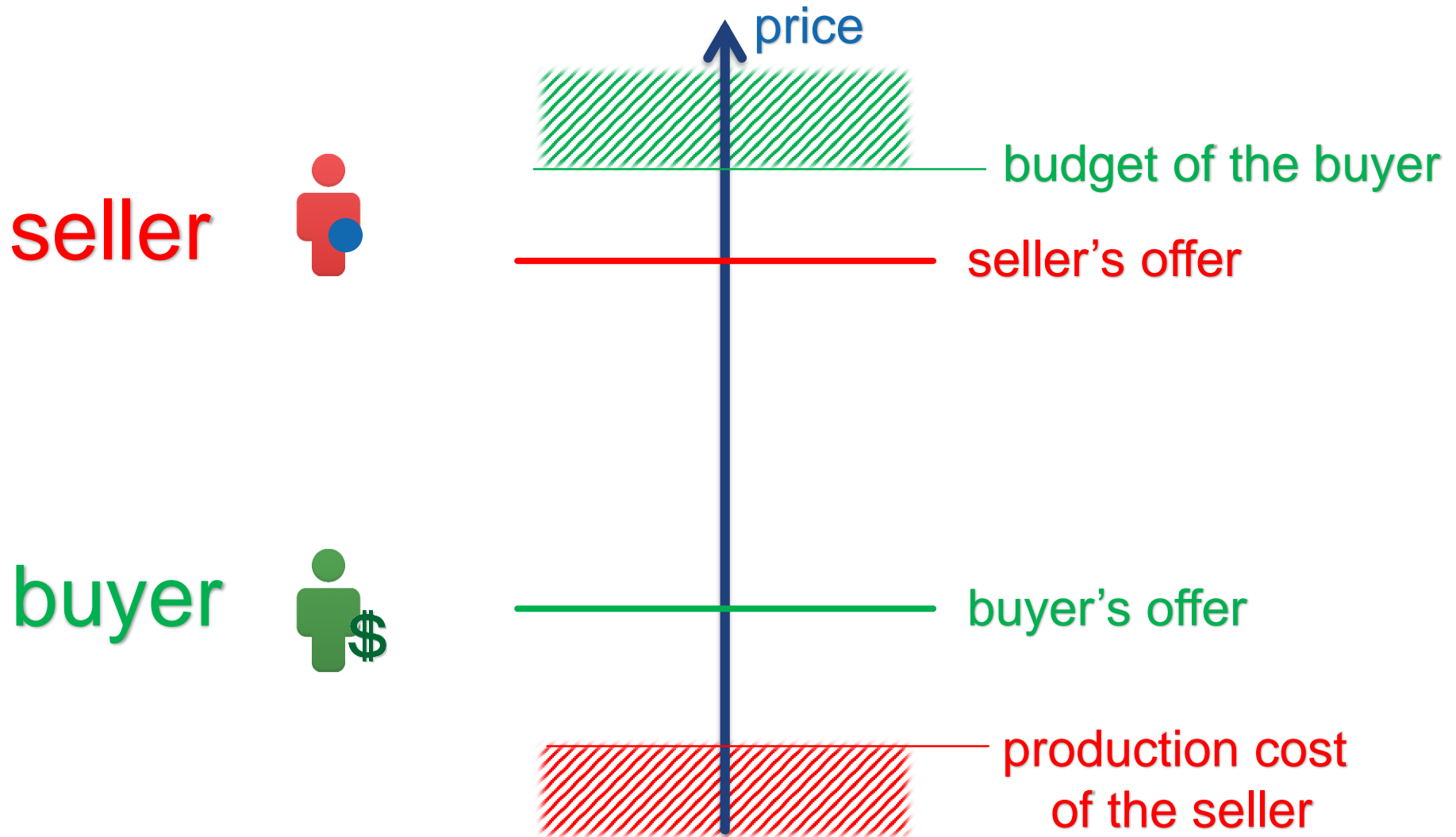
Batuhan Yardim

Aleksei Khudorozhkov

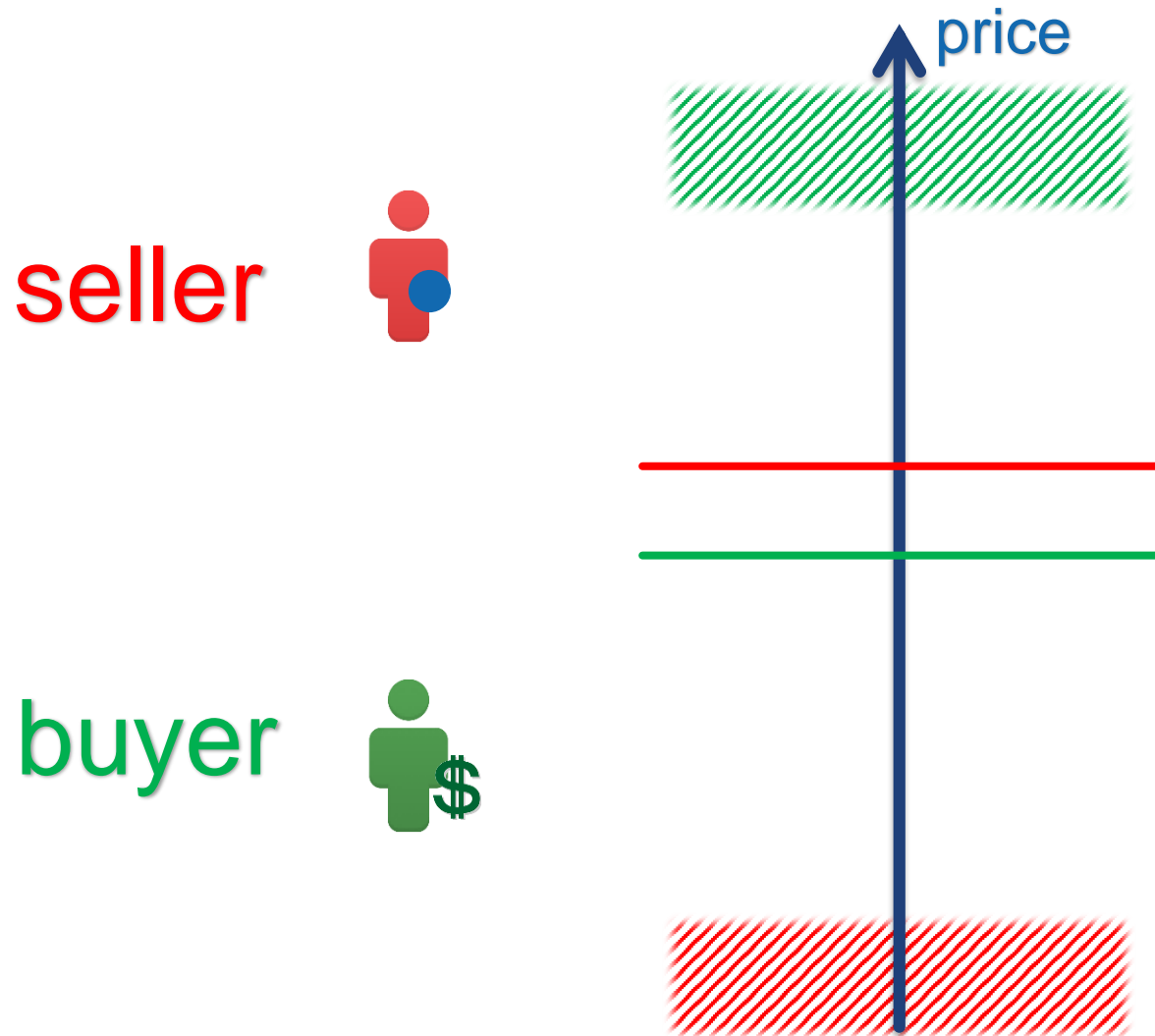
Ka Rin Sim

Neri Passaleva

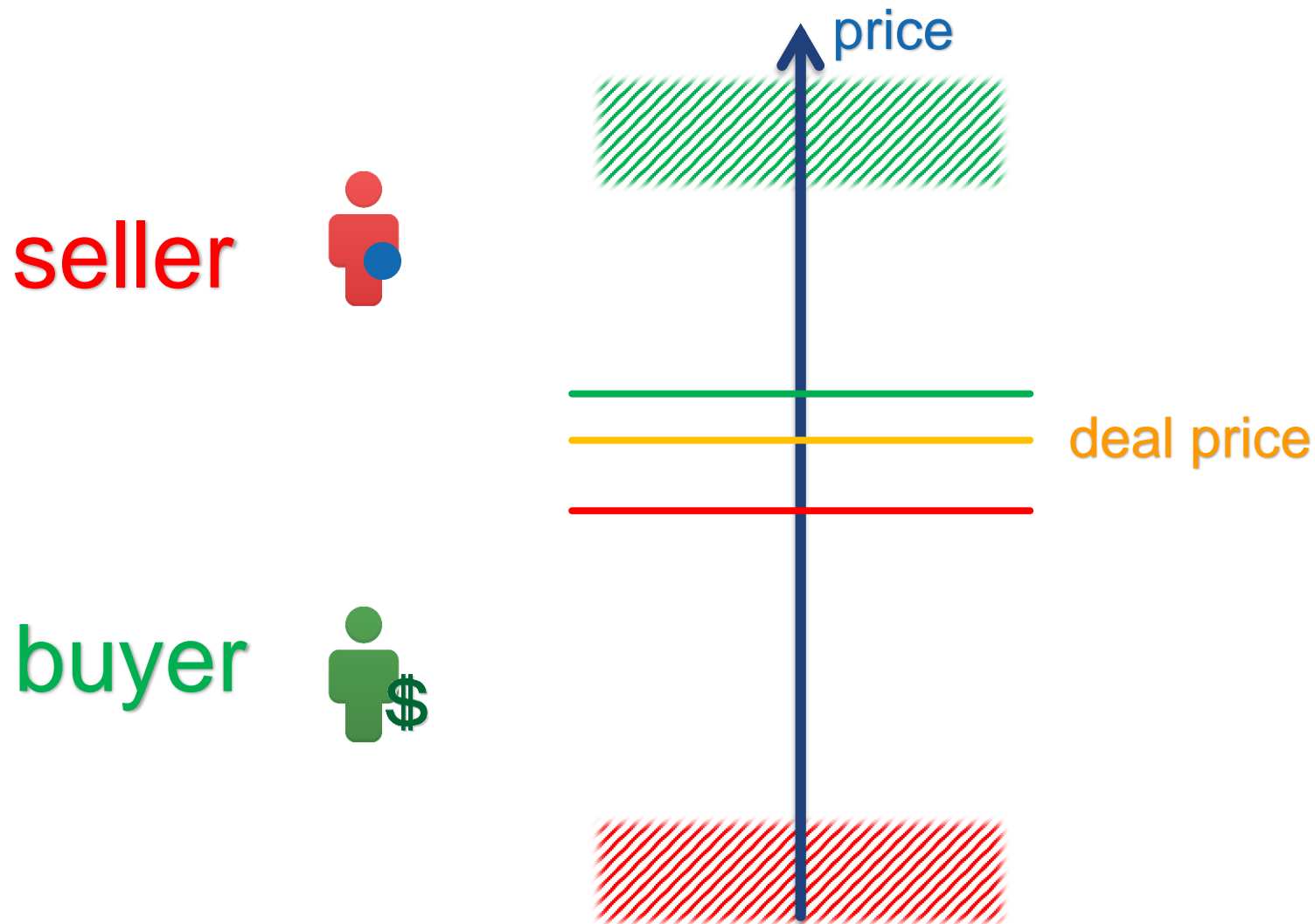
Double auction



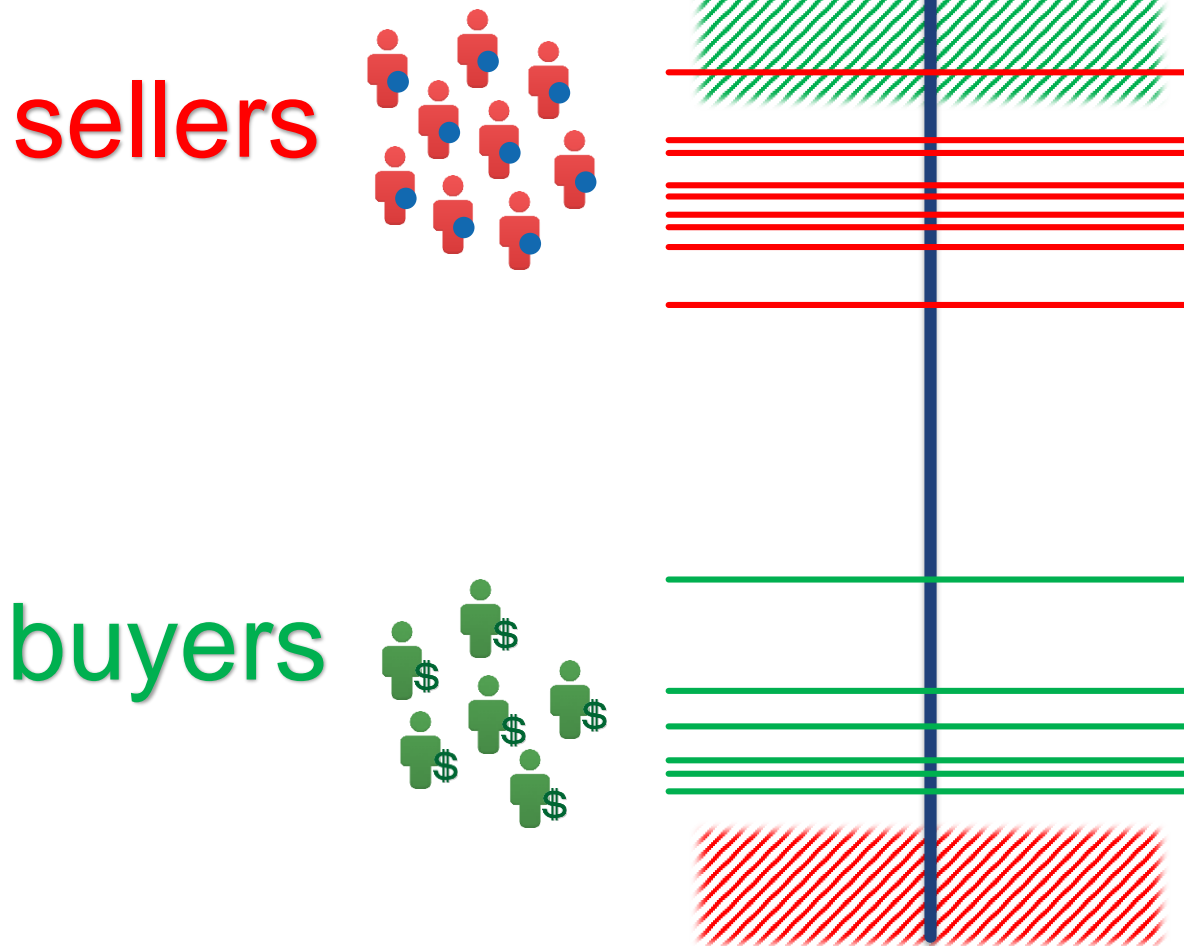
Double auction



Double auction



Double auction



Each agent wants to maximize the reward

For this they can choose different strategies

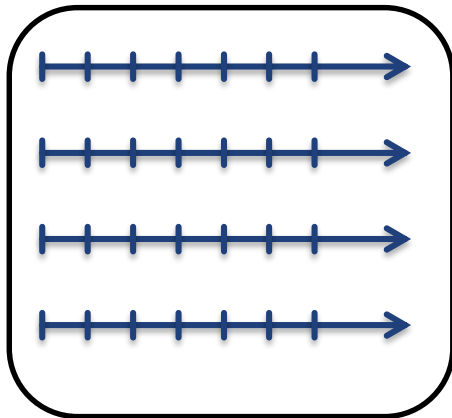
Market environment

- Each round consists of time steps



- Round terminates when T_{max} is reached or no more deals can be made

- Each game consists of rounds



- Agents can have memory about the previous rounds
- Between the games agents can learn and adjust their strategy

Observations of agents

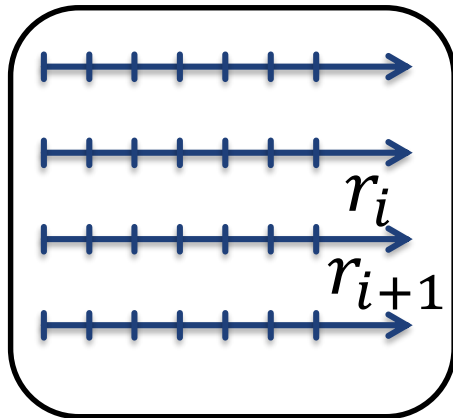
- After each time step an agent receives observations from the market environment.
- Core observations are:
 - The last offer of the agent
 - Current time step
 - bool: if the agent managed to make a deal in the previous round
- Other observations might be included:
 - Last offers of agents of the same/opposite side
 - Reservation prices of agents of the same/opposite side
 - Information about completed deals in the current round
 - The maximum time steps in a round
 - The total number of buyers/sellers
 - The number of buyers/sellers who hasn't made a deal yet

Reward Mechanism

- The agent's **reward** is the absolute difference between the reservation price and the agent's deal offer.

$$r_i = |p_i - a_i| \quad \text{reward for round } i$$

- Reward is cumulative throughout the rounds



$$r = \sum_i r_i \quad \text{total reward}$$

Zero-Intelligence Agent (ZIA)

The agent randomly chooses the next offer according to the exponential distribution around the reservation price



No observations are needed in order to decide on a new offer

Linear Markov Decision Agent (LMDA)

- The new demand is a linear combination of the agent's observations

$$d_i = \alpha d_{i-1} + \beta s + n_i$$

Demand at current time step

Demand at previous time step

Boolean indicator of previous round outcome

Noise

$$s = \begin{cases} 0 & \text{if unsuccessful} \\ 1 & \text{if successful} \end{cases}$$

Price Aggressive Agents (PAA)

s is used as an indicator for whether the agent should be aggressive or not

If **s** is **TRUE**:

$$d_i = \alpha d_{i-1} + n_i$$

If **s** is **FALSE**

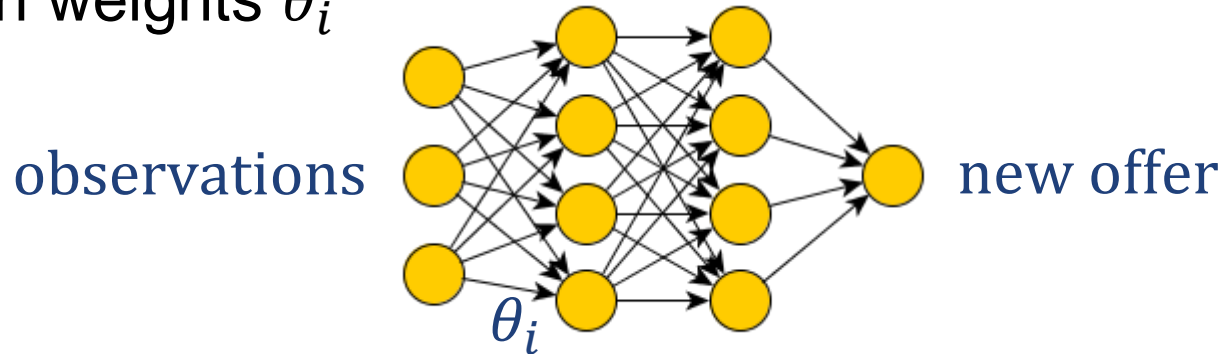
$$d_i = (\alpha + \varepsilon) d_{i-1} + n_i$$

 ε is the agent's level of aggressiveness

The agent becomes aggressive after an unsuccessful round and tries to make a deal even with low reward

Deep RL Agents

- The new offer decision mechanism is a neural network with weights θ_i



$$d_i = \pi_{\theta_i}(o_i) + \mathcal{N}(0, \sigma_i)$$

$\pi_{\theta_i}(o_i)$ → Parametrized Agent's Policy

$\mathcal{N}(0, \sigma_i)$ → Gaussian noise

- Reinforcement learning through Deep Deterministic Policy Gradient (DDPG) framework

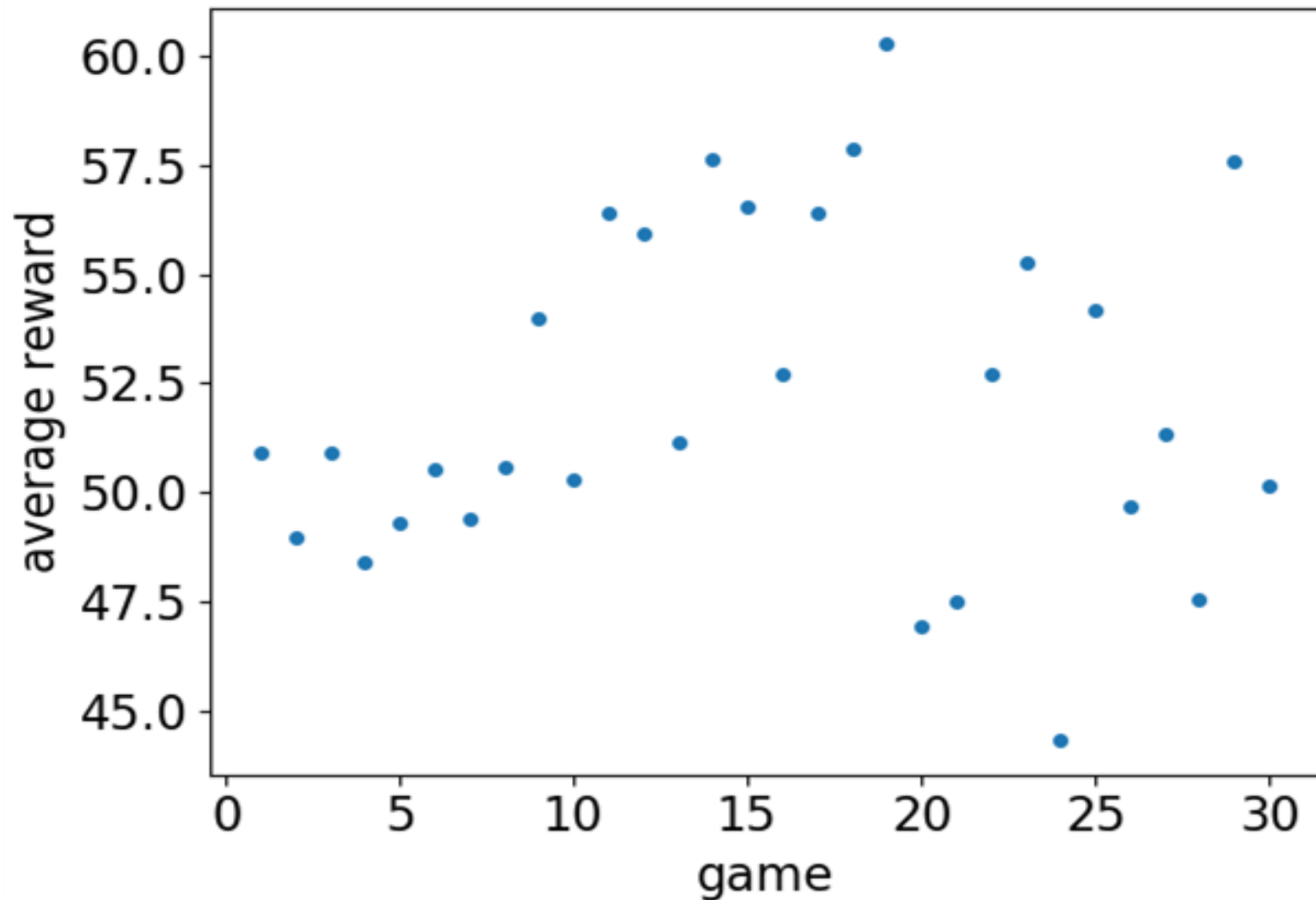
Results:

Non-Intelligent Agents

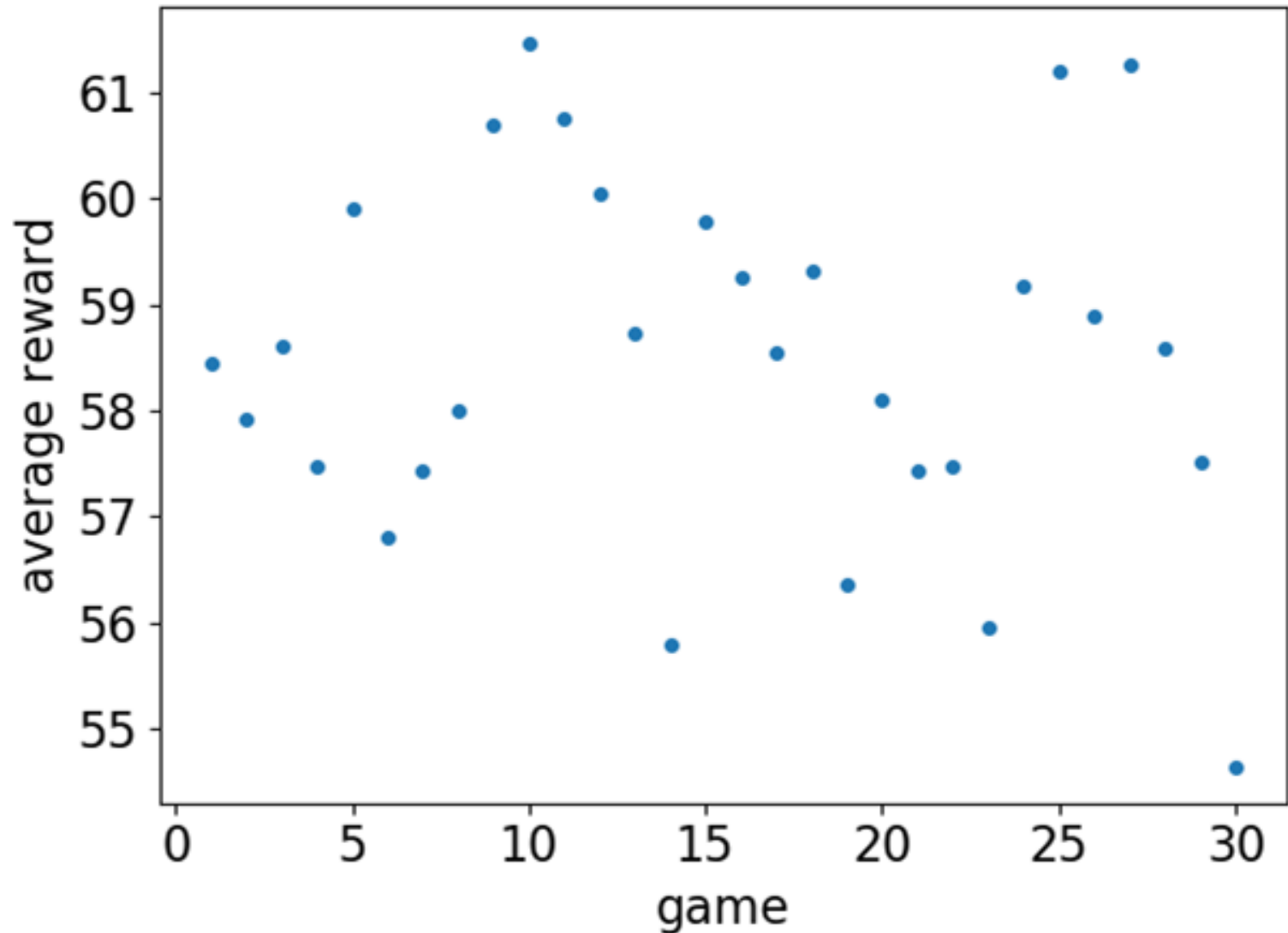
(ZIA, LMDA, PAA)

- ★ No learning
- ★ No correlation

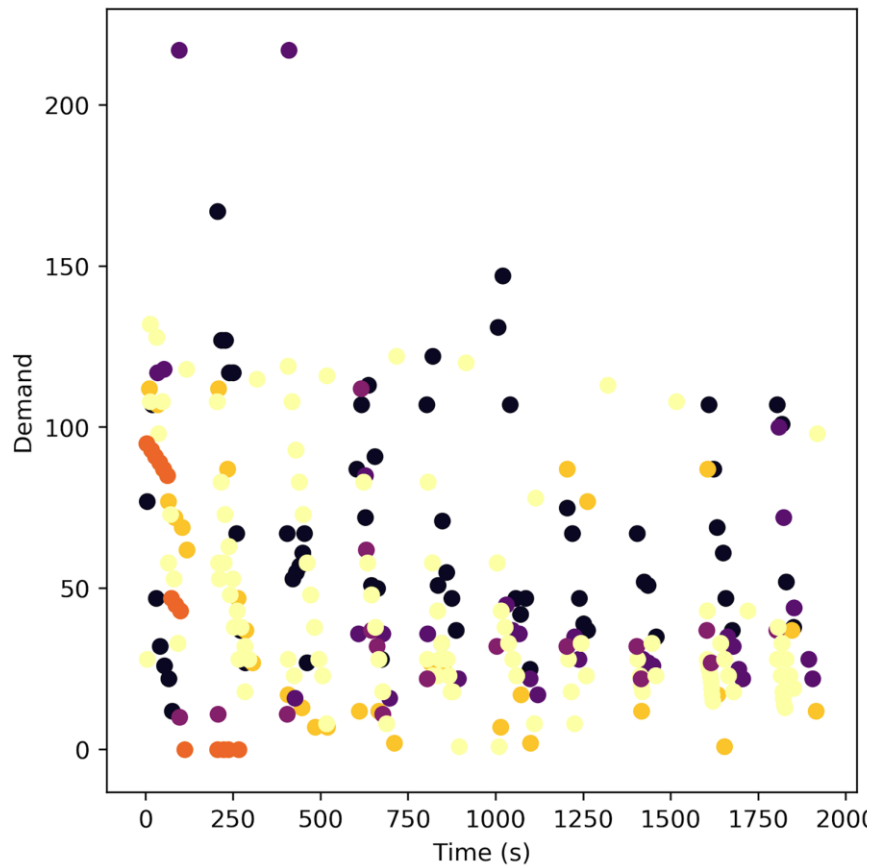
Zero Intelligence Agents (ZIA)



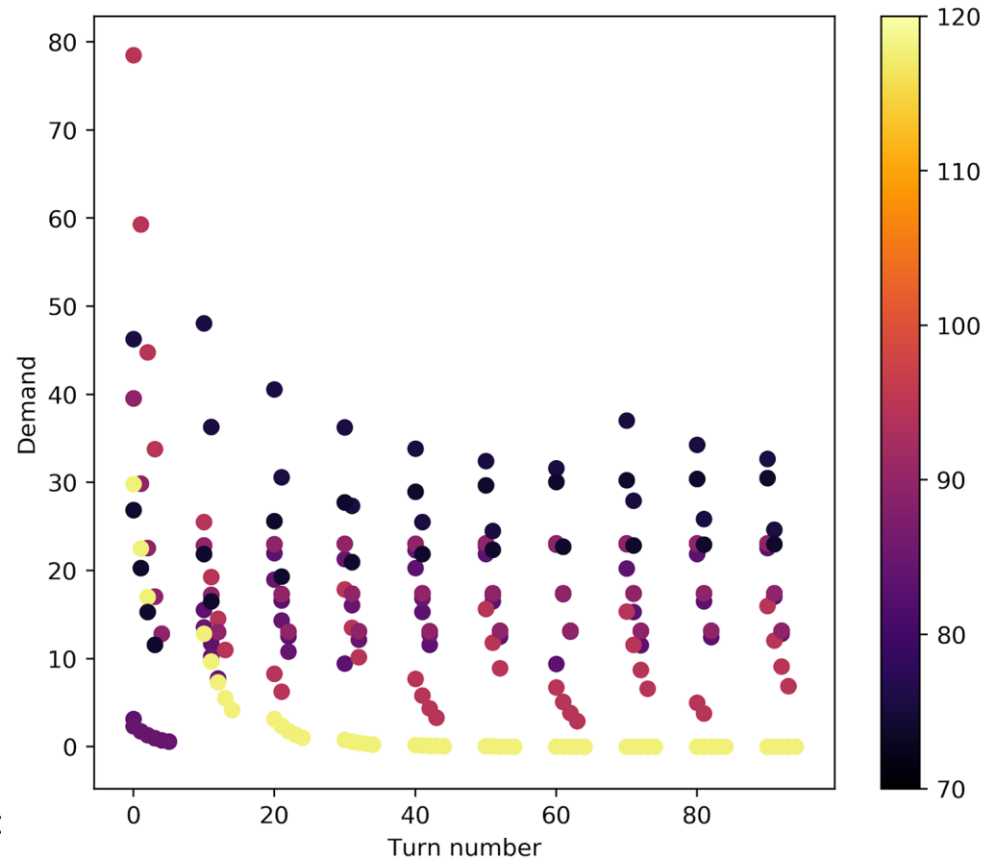
Linear Markov Decision Agents (LMDA)



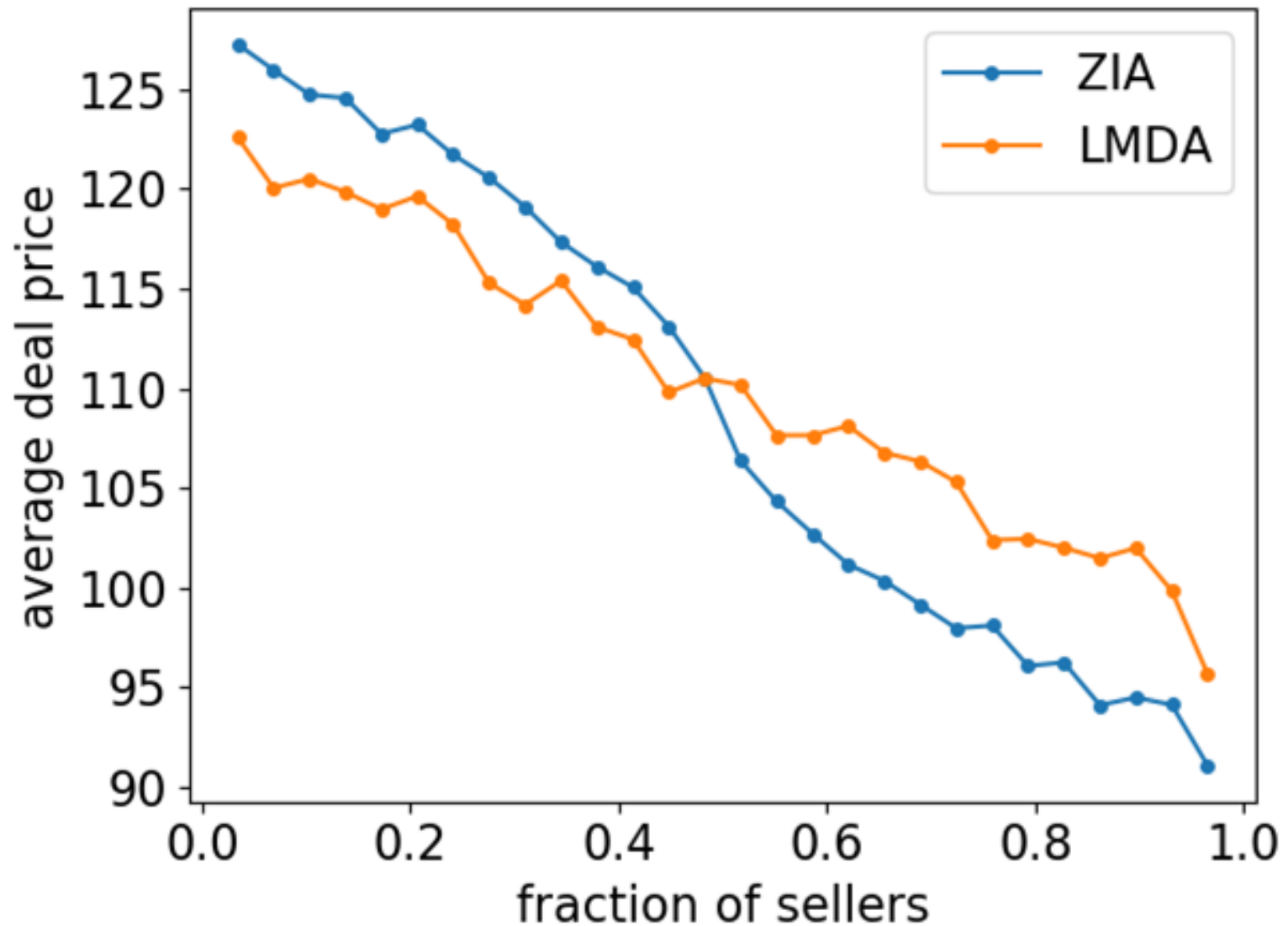
Demands -- Experiments vs LMDA



Experiments



LMDA



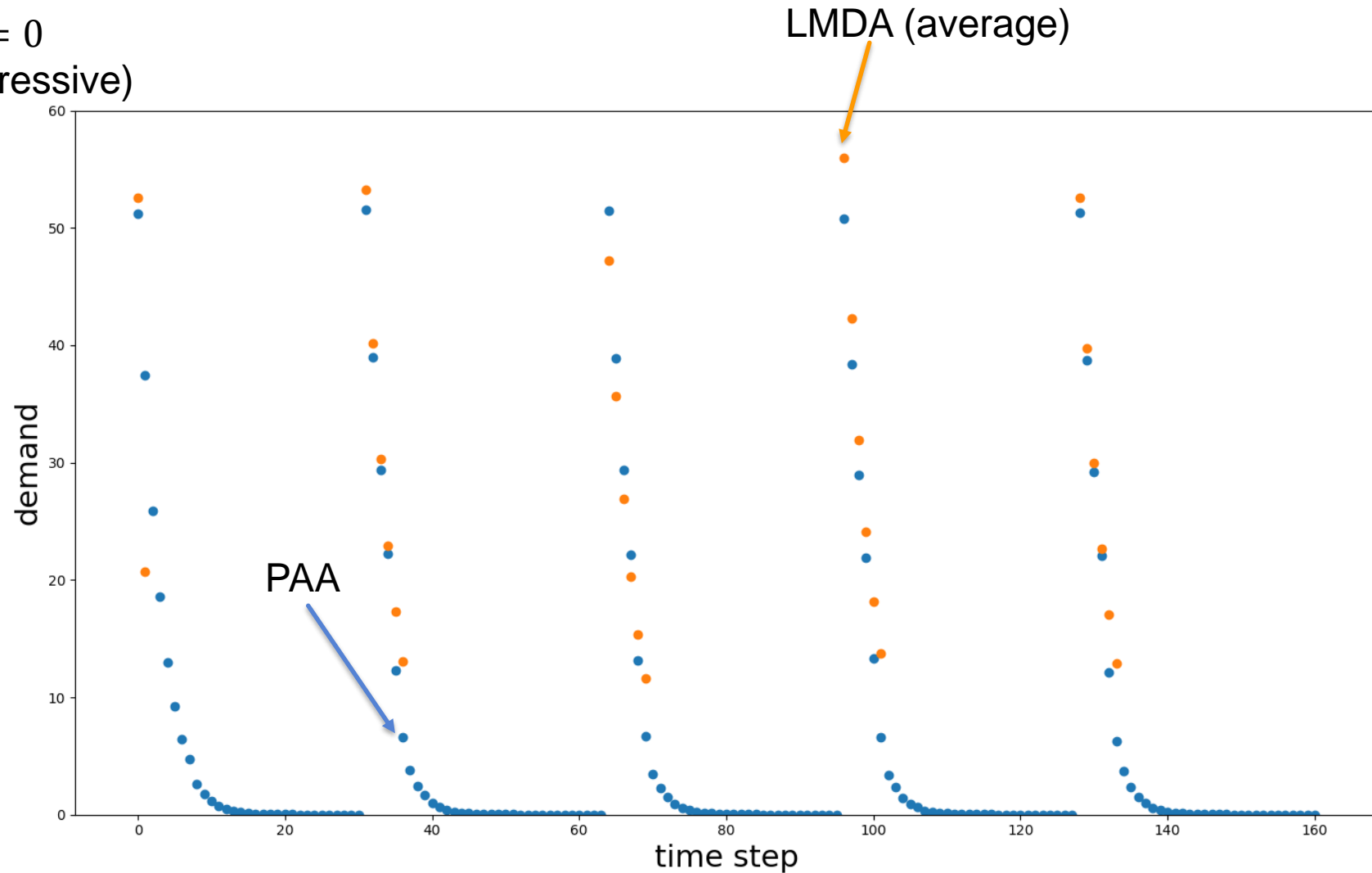
Deal Price vs Fraction of Sellers

- ★ Seller fraction increases → deal price decreases
- ★ Competition
- ★ Slope → competition
 - ★ ZIA faces tougher competition than LMDA

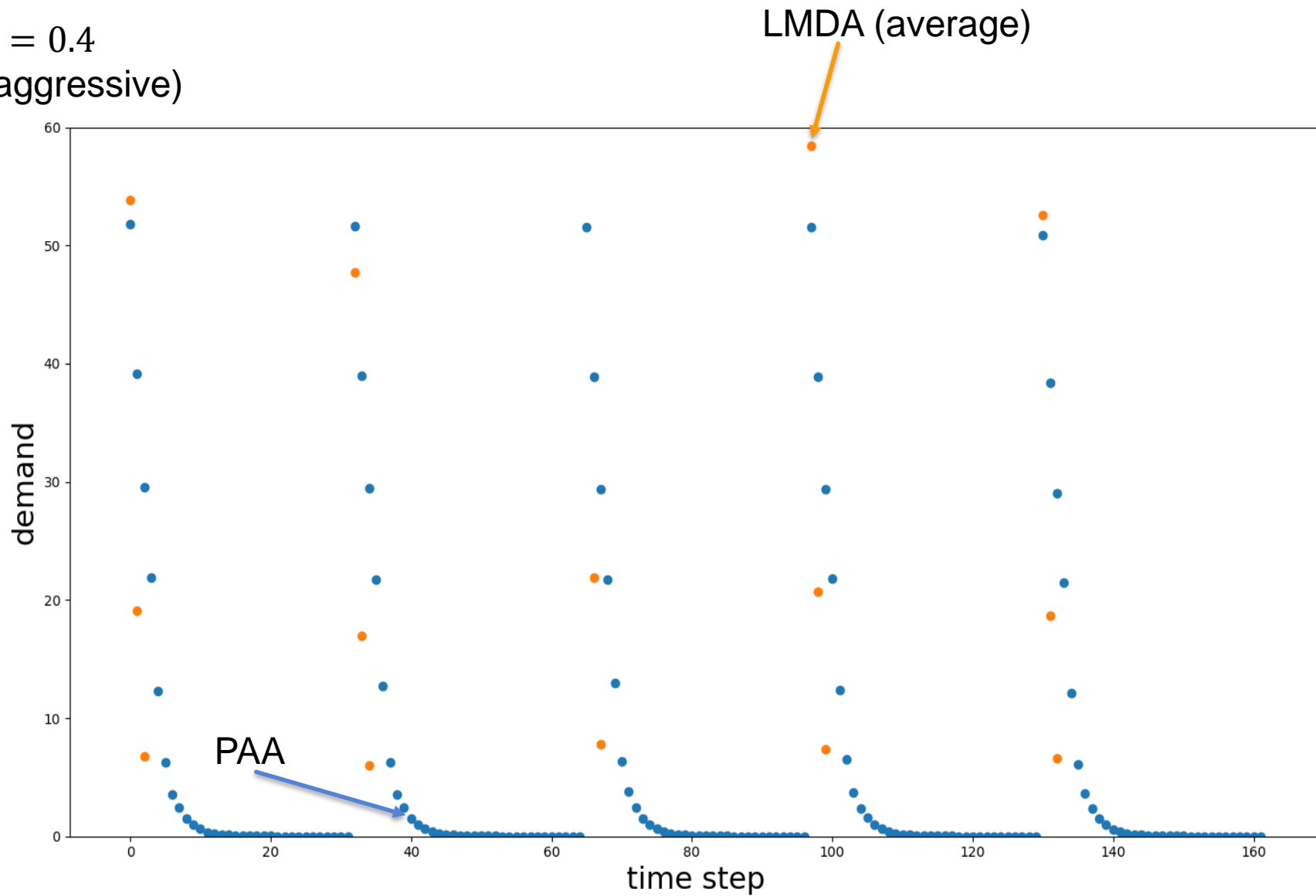
Price Aggressive Agents (PAA)

 $\varepsilon = 0$

(non-aggressive)



$\varepsilon = 0.4$
(quite aggressive)

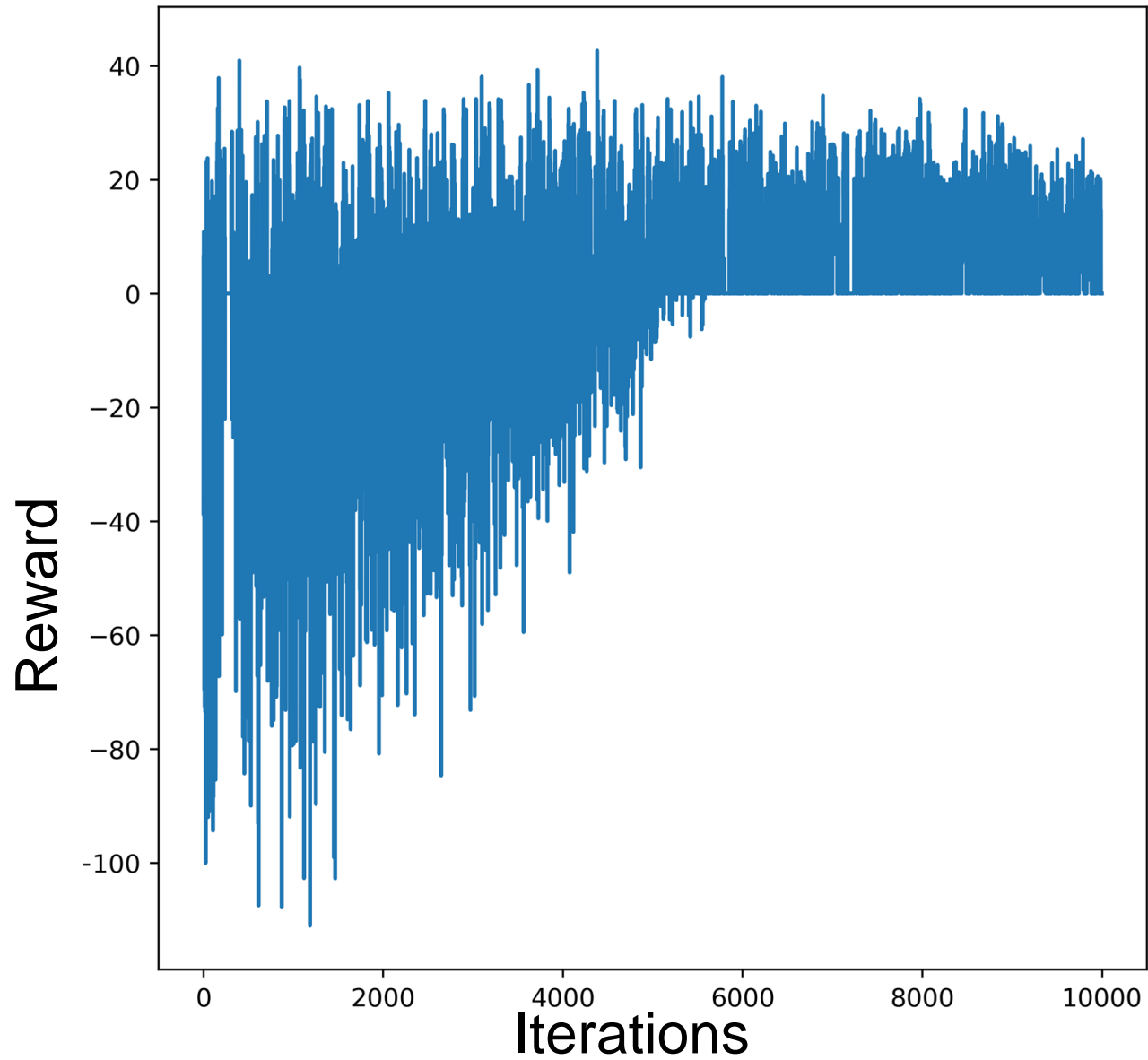


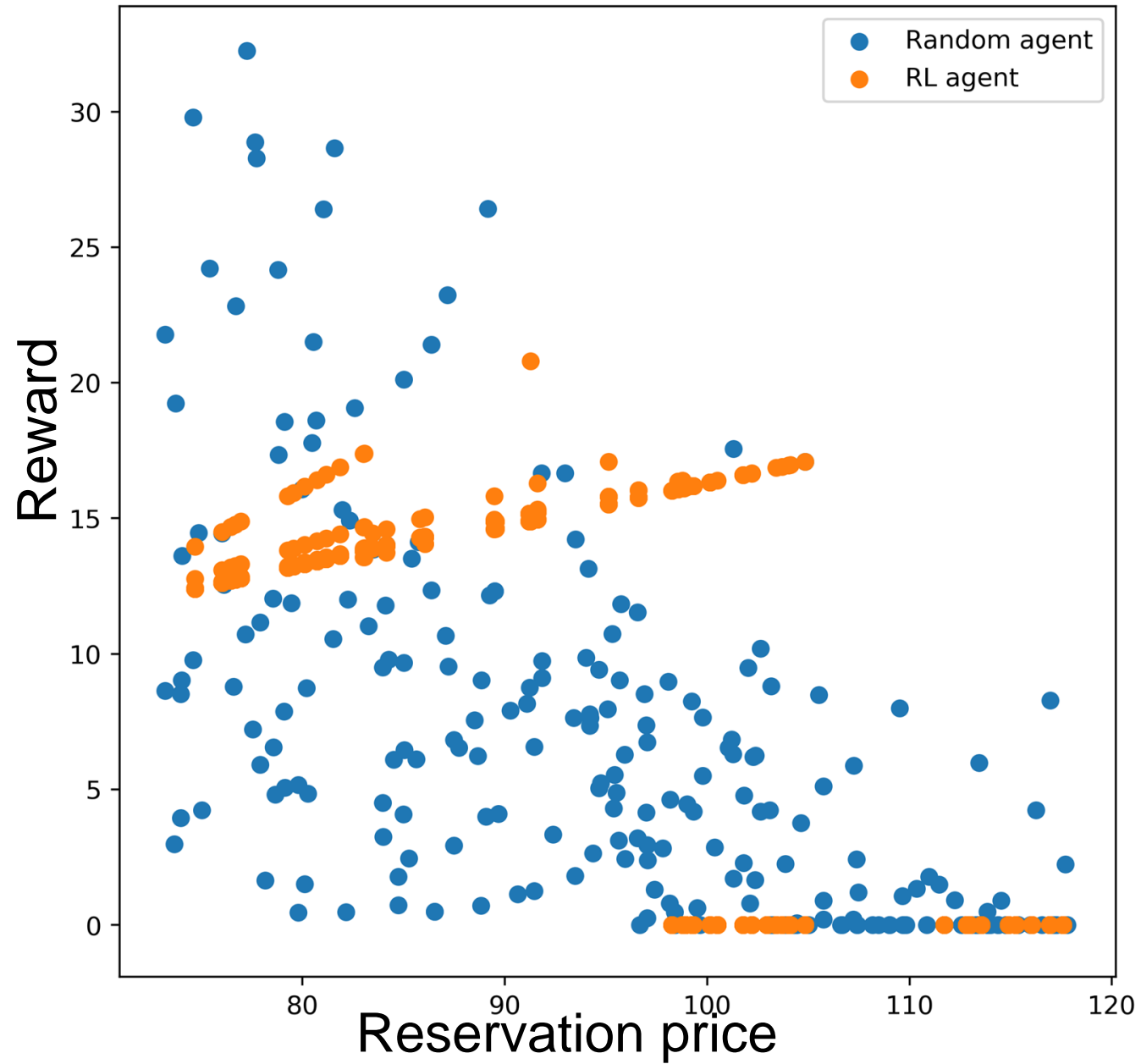
- PAA always has a lower offer price
- Deal is achieved faster by PAA

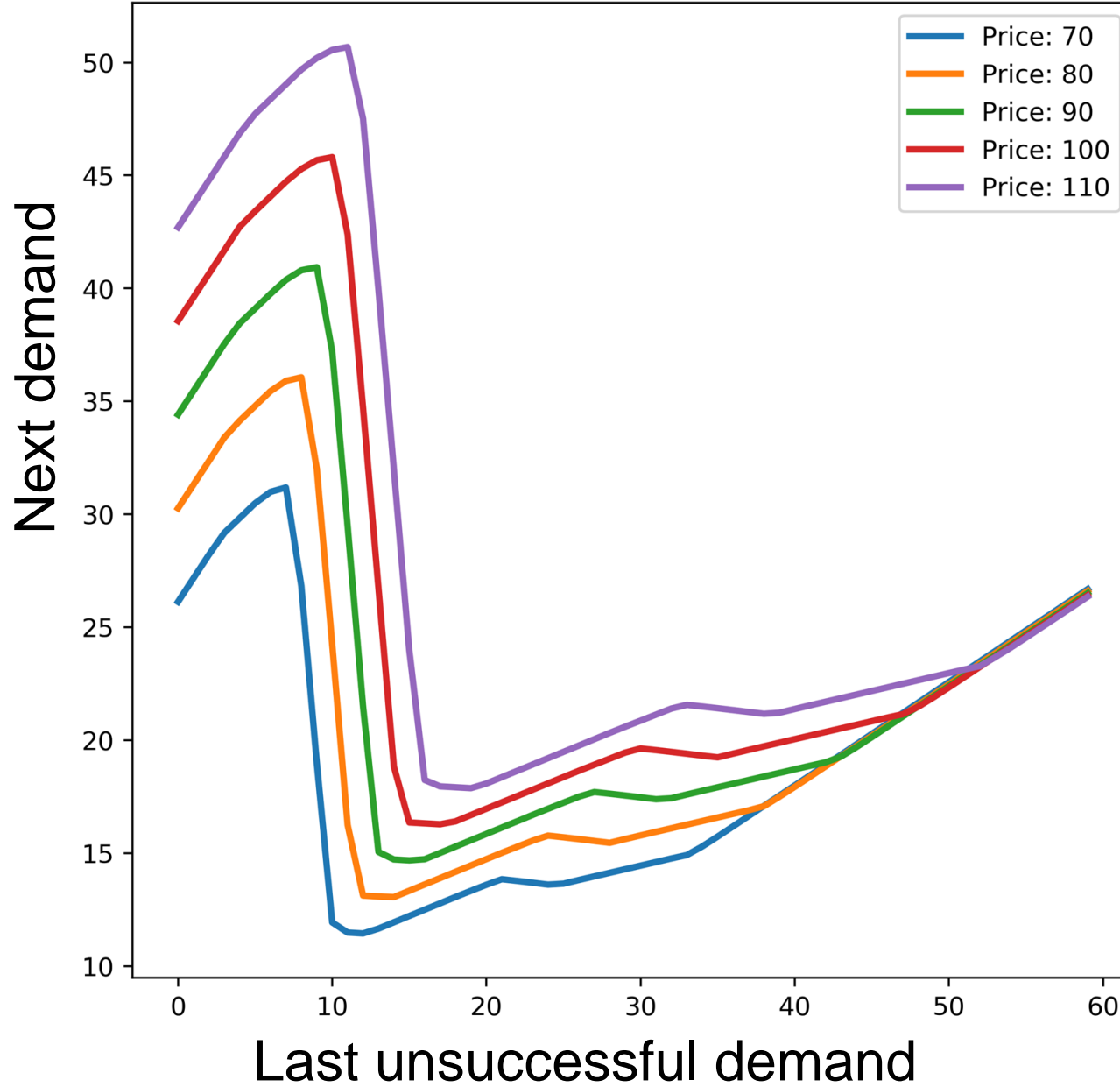
Deep Reinforcement Learning Model

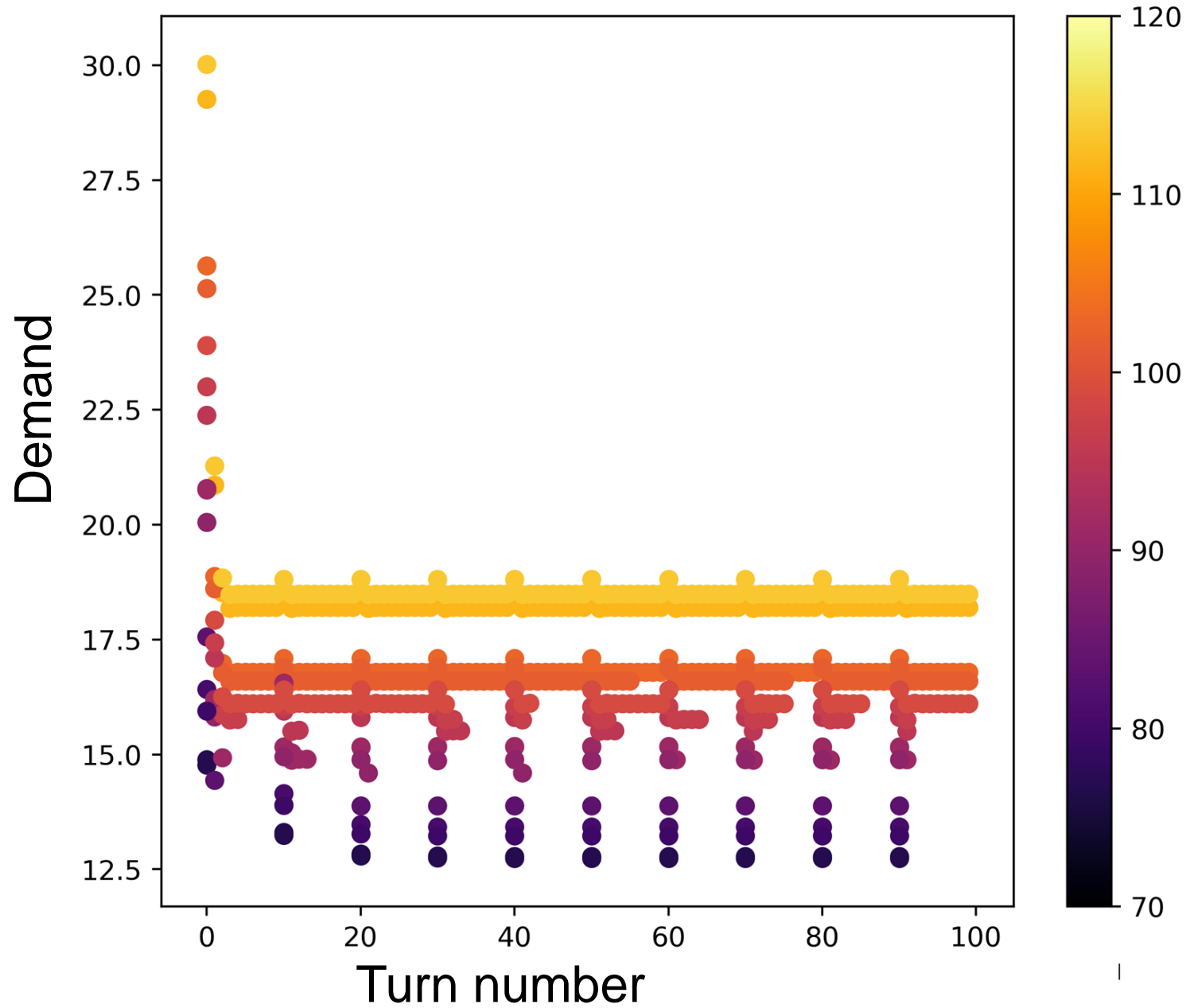
- ★ 2 different exploration policies:
 - ★ Gaussian
 - ★ Ornstein-Uhlenbeck (OU)
- ★ 1 intelligent agent + a pool of non-intelligent agents (e.g. ZIA, LMDA)

Gaussian Exploration Policy

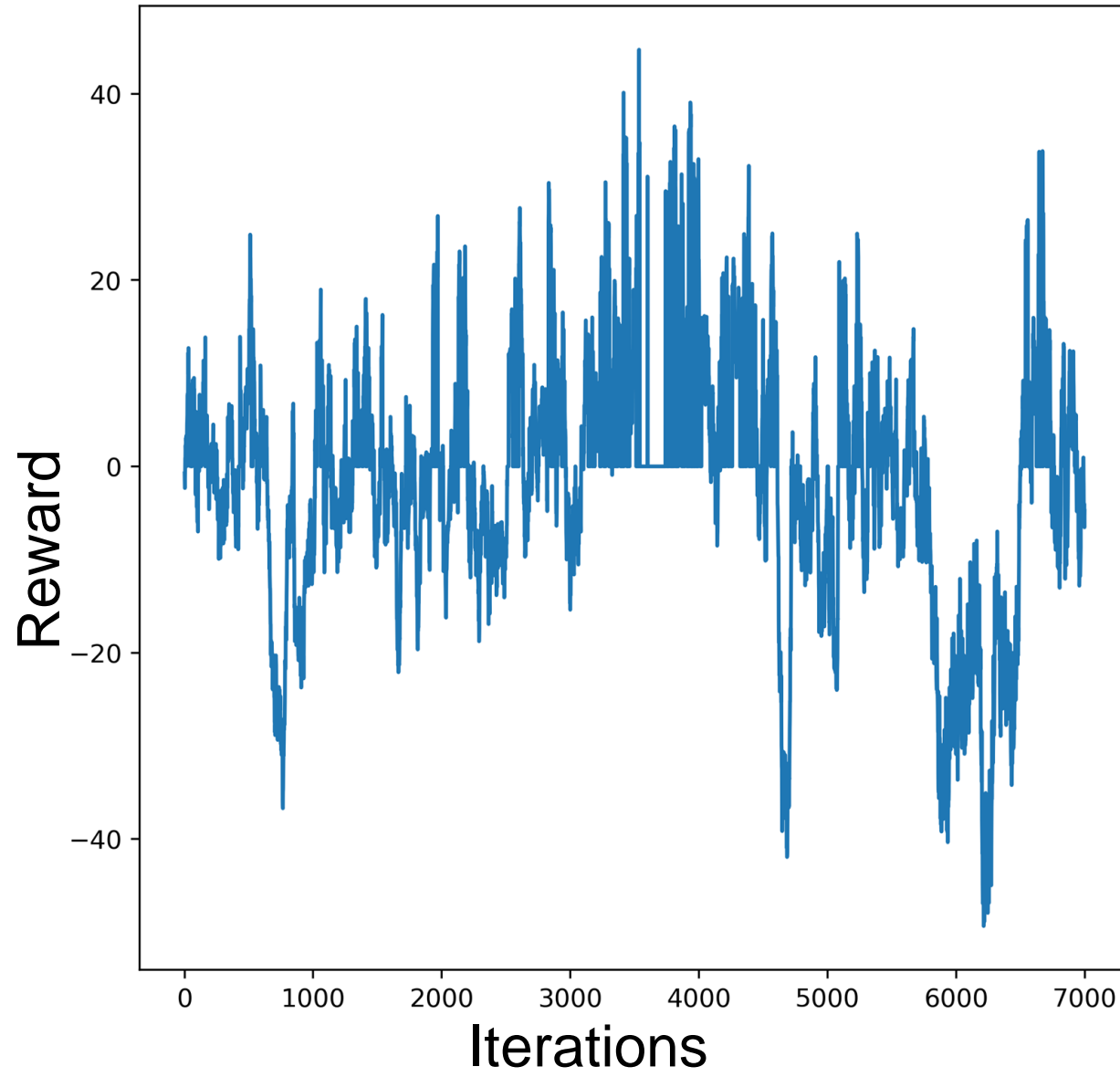


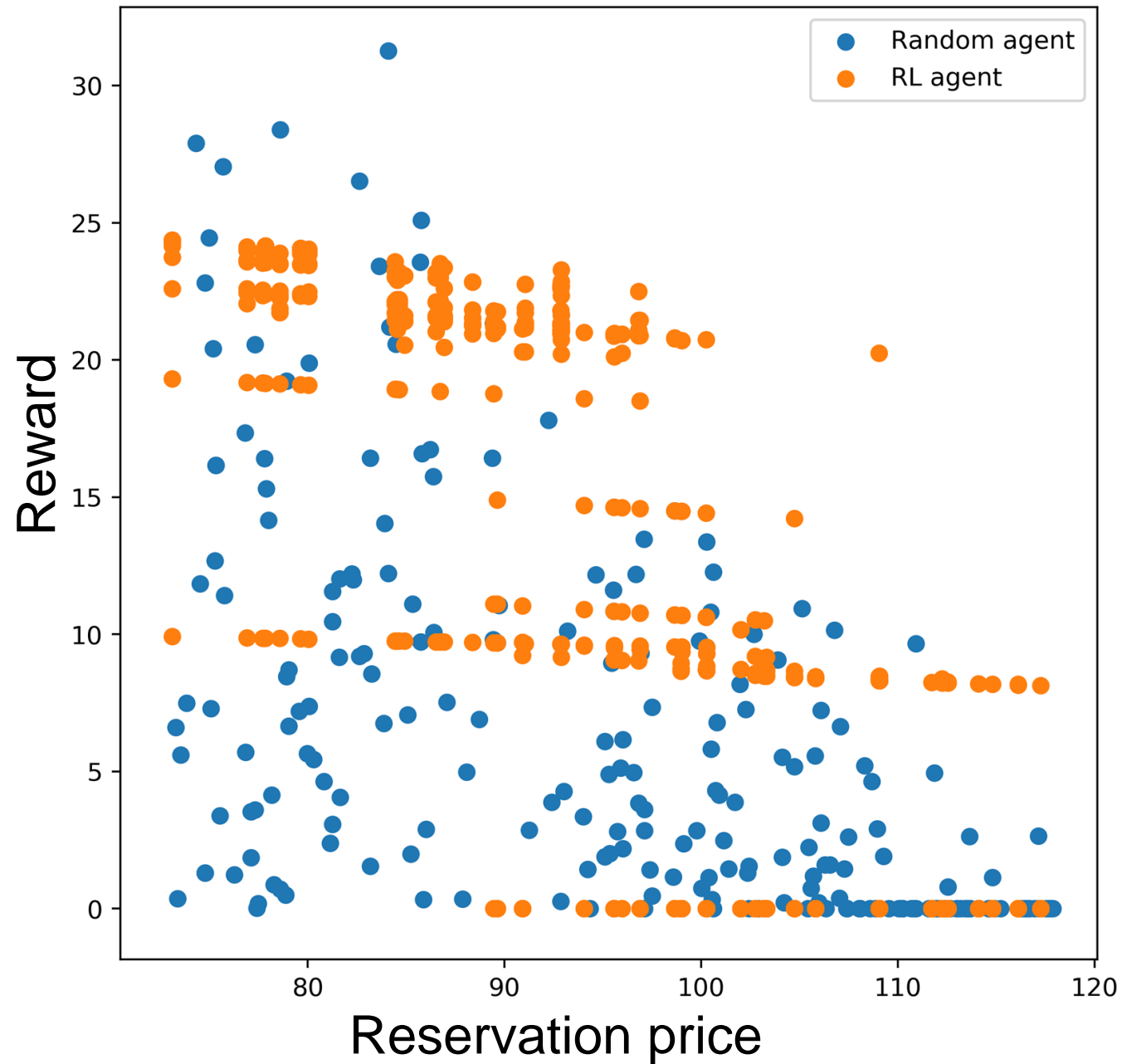


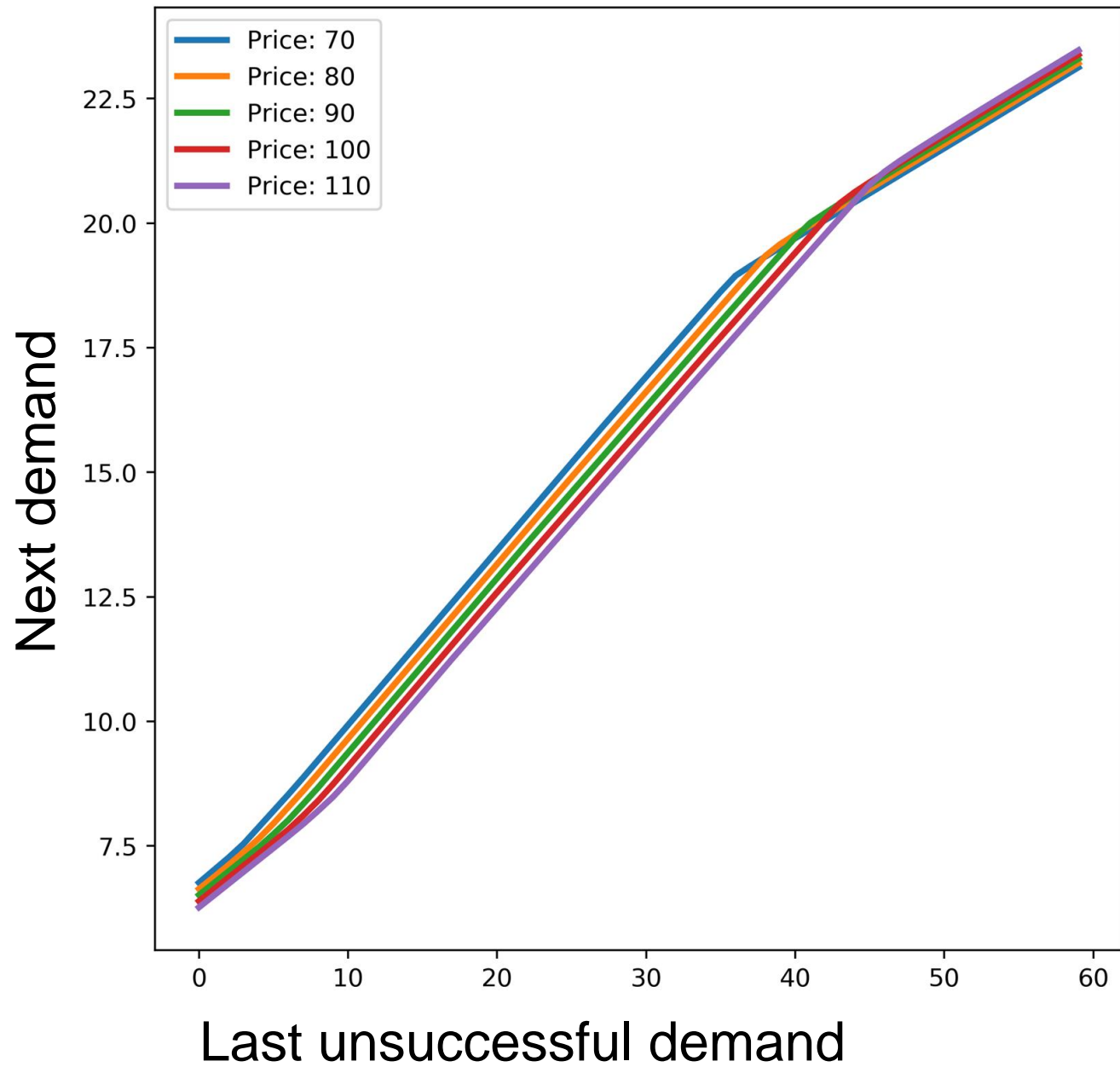


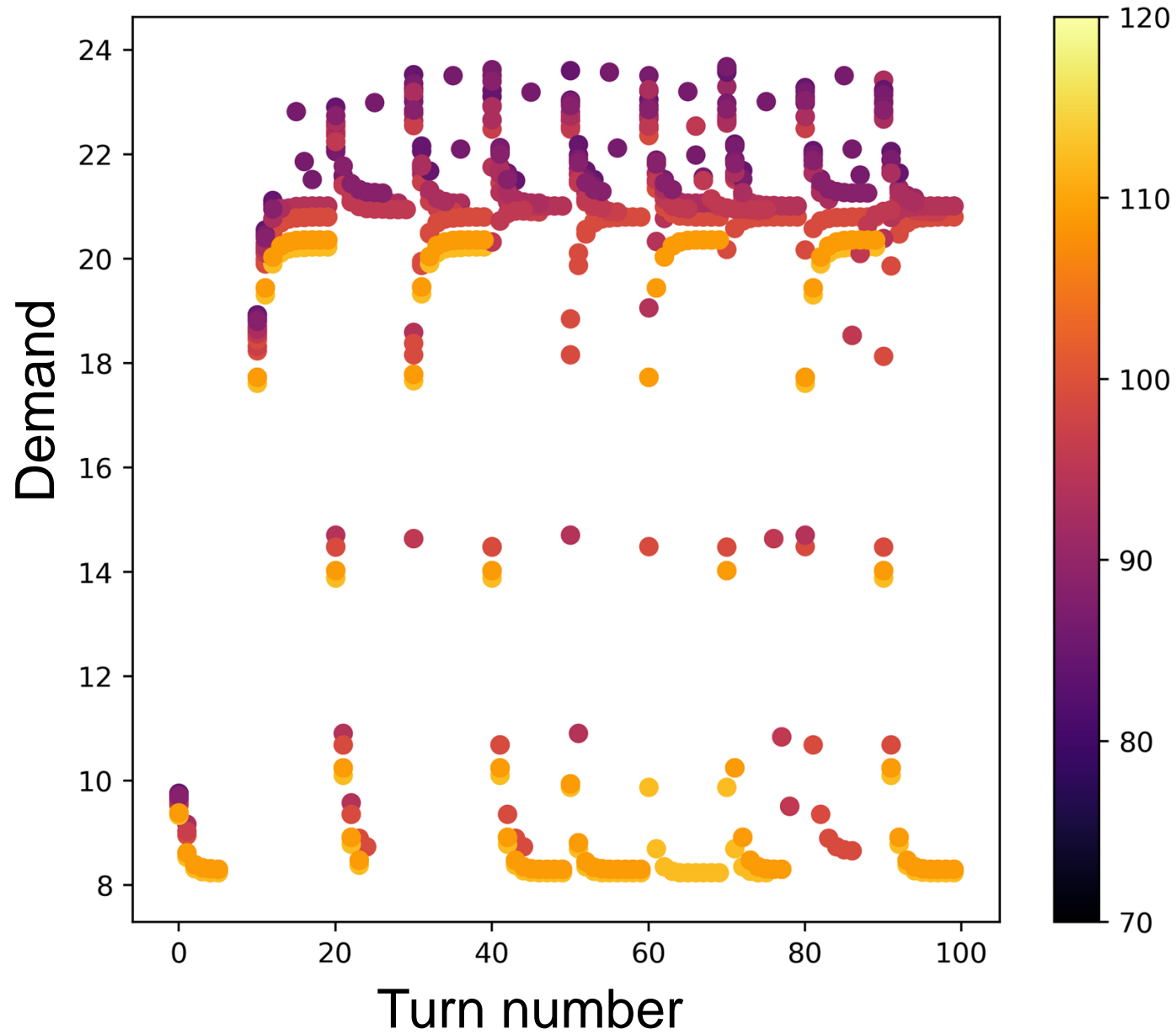


Ornstein-Uhlenbeck (OU) Exploration Policy









- ★ reward is negative at early stages
- ★ Gaussian agent learns quickly to avoid negative reward
- ★ higher reservation price
 - ★ more difficult to sell
 - ★ market performs worse
- ★ Gaussian agent performs better than OU agent

Agent Pool	Pool Earnings	RLA Earnings	Learning Agent
ZIA	5.23	9.27	RLA+Gaussian
ZIA	7.01	8.39	RLA+OU
ZIA	5.27	12.32	RLA+OU+anneal
LMDA	7.19	-	RLA+Gaussian
LMDA	-	-	RLA+OU
PAA	-	-	RLA+OU

Conclusion & Outlook

- RL agent outperformed the pool (ZIA, LMDA, PAA)
- RL agent (Gaussian) learned to avoid negative rewards; reward increased with iterations
- Gaussian agents have a larger profit margin than OU agents for higher res. price
- Humans are much more conservative than RL
- Extension: include more information into observation space
 - Currently using black box setting