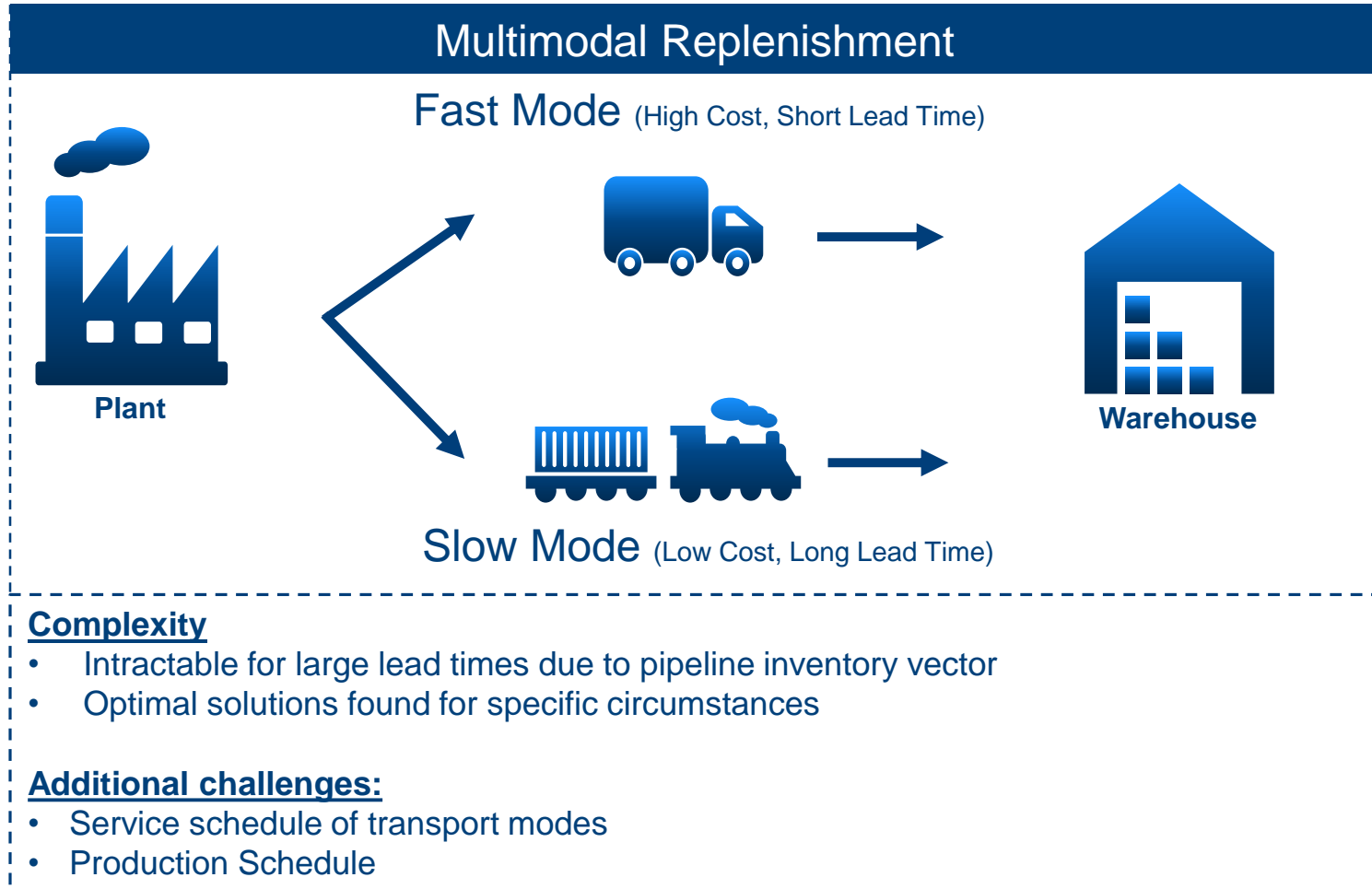**KU LEUVEN**

# A deep reinforcement learning approach for synchronized multi-modal replenishment

joren.gijsbrechts@kuleuven.be
robert.boute@kuleuven.be
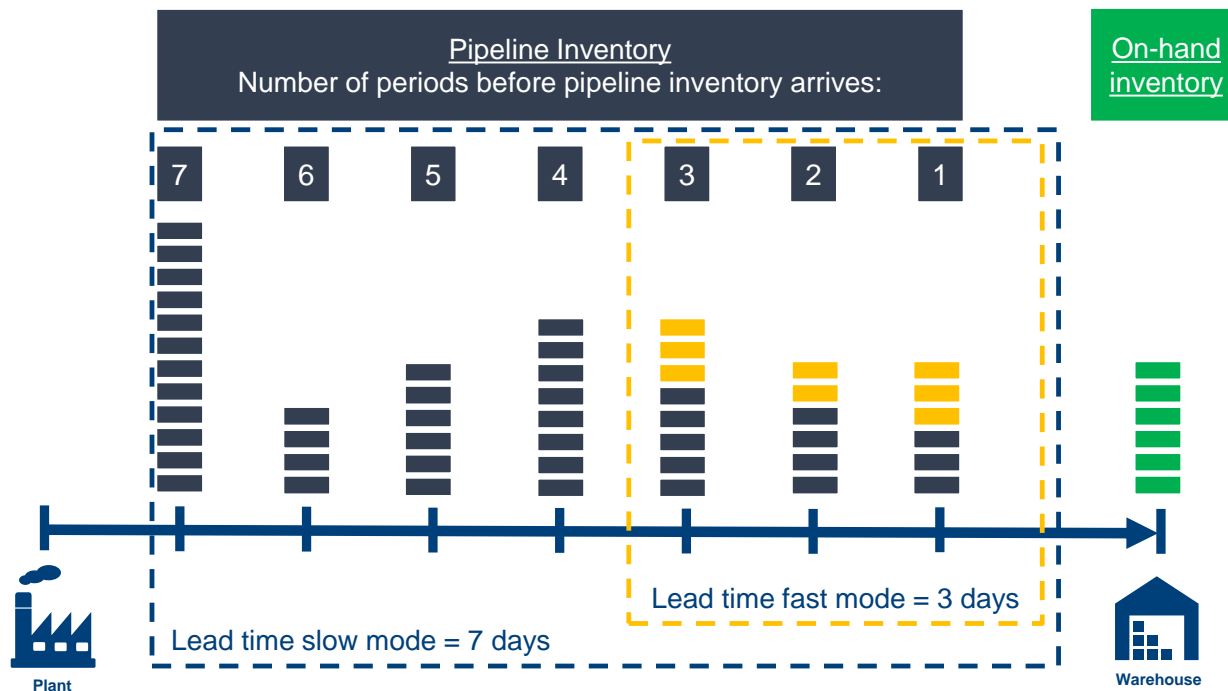
# Problem Statement
## Multimodal replenishment as a dual sourcing problem



**Multimodal Replenishment**

Fast Mode (High Cost, Short Lead Time)

Plant

Warehouse

Slow Mode (Low Cost, Long Lead Time)

**Complexity**
- Intractable for large lead times due to pipeline inventory vector
- Optimal solutions found for specific circumstances

**Additional challenges:**
- Service schedule of transport modes
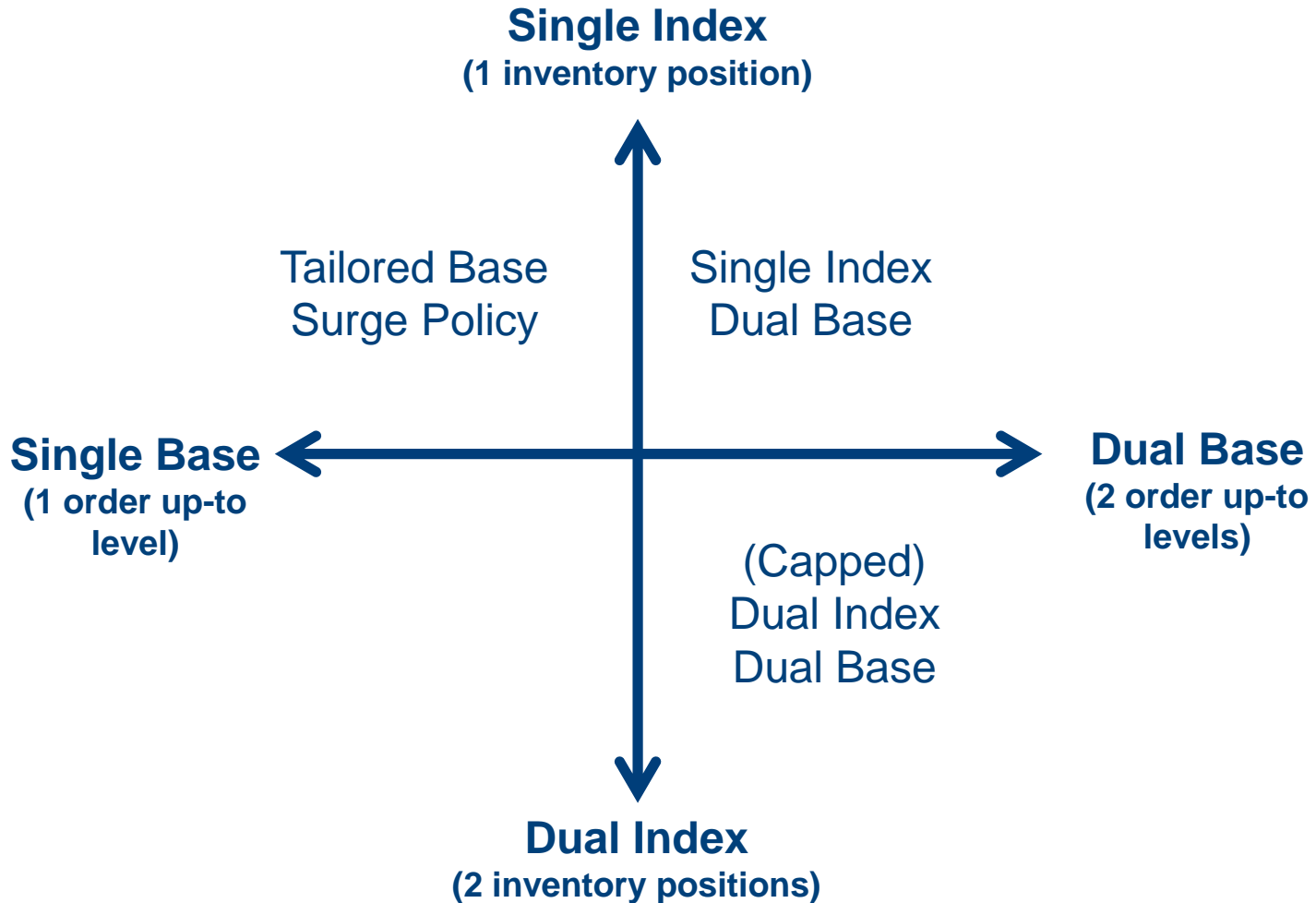- Production Schedule

KU LEUVEN

# State of the art
## Heuristic policies group parts of pipeline inventory vector

- Complexity arises from pipeline inventory vector
- State-of-the-art heuristic policies hence use:
  - 1 or 2 inventory positions
  - 1 or 2 order up-to levels



Pipeline Inventory
Number of periods before pipeline inventory arrives:

On-hand inventory

Lead time fast mode = 3 days

Lead time slow mode = 7 days

Plant

Warehouse

KU LEUVEN

# State of the art
## Heuristic policies work around complexity

**Single Index**
**(1 inventory position)**

Tailored Base
Surge Policy

Single Index
Dual Base

**Single Base**
**(1 order up-to level)**

**Dual Base**
**(2 order up-to levels)**

(Capped)
Dual Index
Dual Base

**Dual Index**
**(2 inventory positions)**

KU LEUVEN

# Motivation
## Can Artificial Intelligence be used to solve the dual sourcing problem?

"I would say, a lot of the value that we're getting from **machine learning** is actually happening kind of **beneath the surface**. It is things like improved search results, improved product recommendations for customers, improved forecasting for inventory management, and literally hundreds of other things beneath the surface," Bezos said.
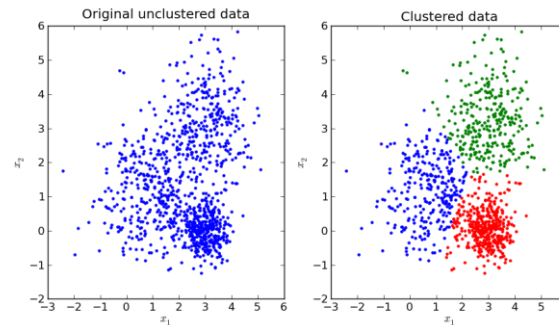


KU LEUVEN

# Machine Learning Overview

# Reinforcement Learning, no new field!

# But … major recent breakthroughs!

# Contribution

- Smart algorithm learns itself a replenishment policy based on full pipeline inventory vector

- Suitable for complex settings
  - Non-linear ordering cost
  - Include ordering/delivery/production schedules
    - E.g. non-daily train/boat schedule

- First application of deep reinforcement learning in dual sourcing

KU LEUVEN

# Methodology

- Problem modeled as a **Markov Decision Process** *(S,A,R(s,a), γ)*
  - *State space (S) = Inventory Vector + Day of week*

$$S = \begin{bmatrix} I^{(0)} & I^{(1)} & \dots & I^{(L_s)} & D \end{bmatrix},$$

  - *Action space (A) = Ordering Vector (Fast + Slow)*

$$A = \begin{bmatrix} a^{(f)} & a^{(s)} \end{bmatrix},$$

  - *Rewards (R(s,a)) = Reward realization (ordering + inventory cost)*

$$r(s_t, a_i) = c^f a_i^{(f)} + c^s a_i^{(s)} + h[I_{t+1}^{(0)}]^+ + b[I_{t+1}^{(0)}]^- + \sum_{i=1}^{L_s} p I_{t+1}^{(i)}.$$

  - *γ = discount factor*
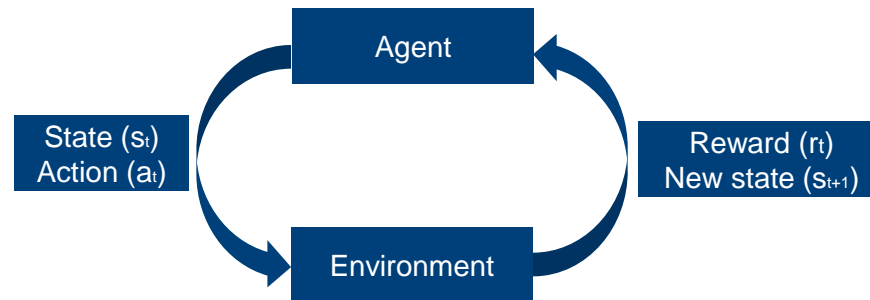
- *Objective: minimize future discounted costs*

$$\min r_1 + \gamma r_2 + \gamma^2 r_3 + \cdots + \gamma^{t-1} r_t + \dots,$$

# Methodology

- Dynamic Programming intractable ➜ Approximate Dynamic Programming
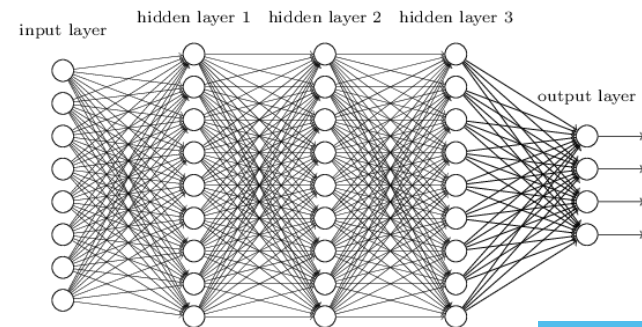  - Reinforcement Learning – *Q*-learning

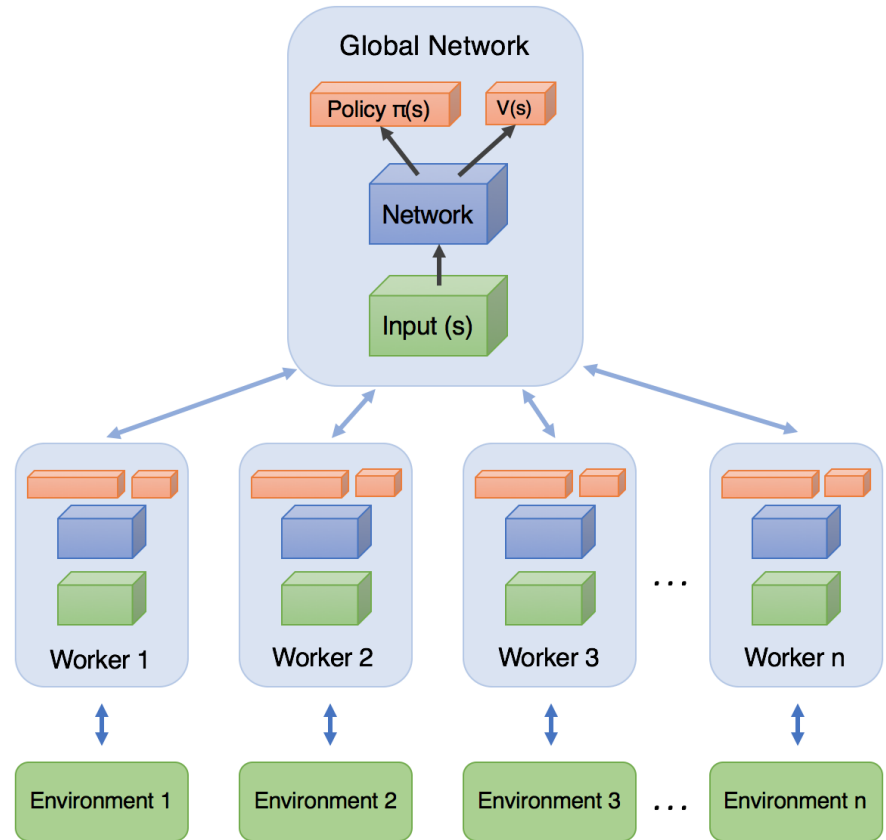$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t))$$



- Q-learning slow for large state space ➜ Deep *Q*-learning
  - o Input: states
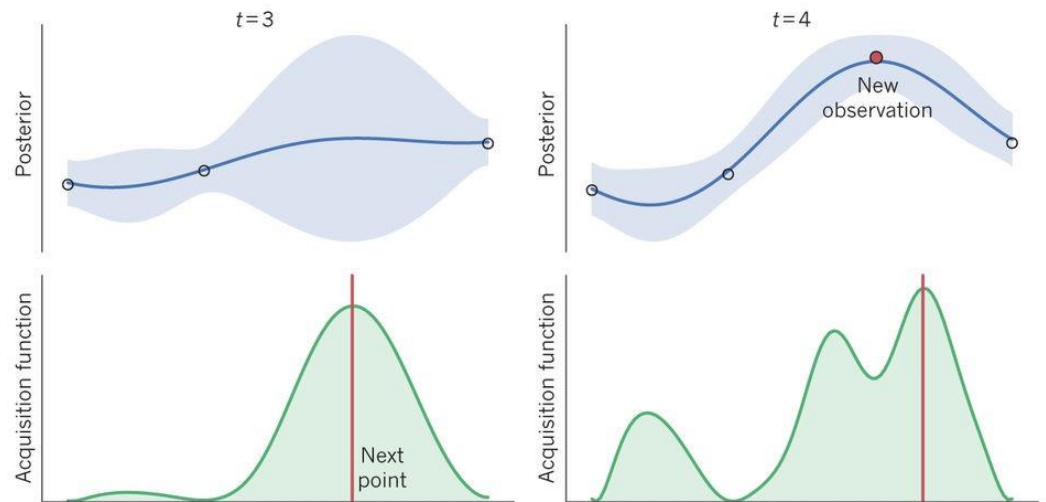  - o Output: *Q*-value for each action

# Methodology

- Asynchronous Advantage Actor-Critic (A3C)
  - Actor develops policy
  - Critic evaluates policy

# Hyperparameter Tuning

- Grid Search
- Random Search
- Bayesian Optimization

# Results
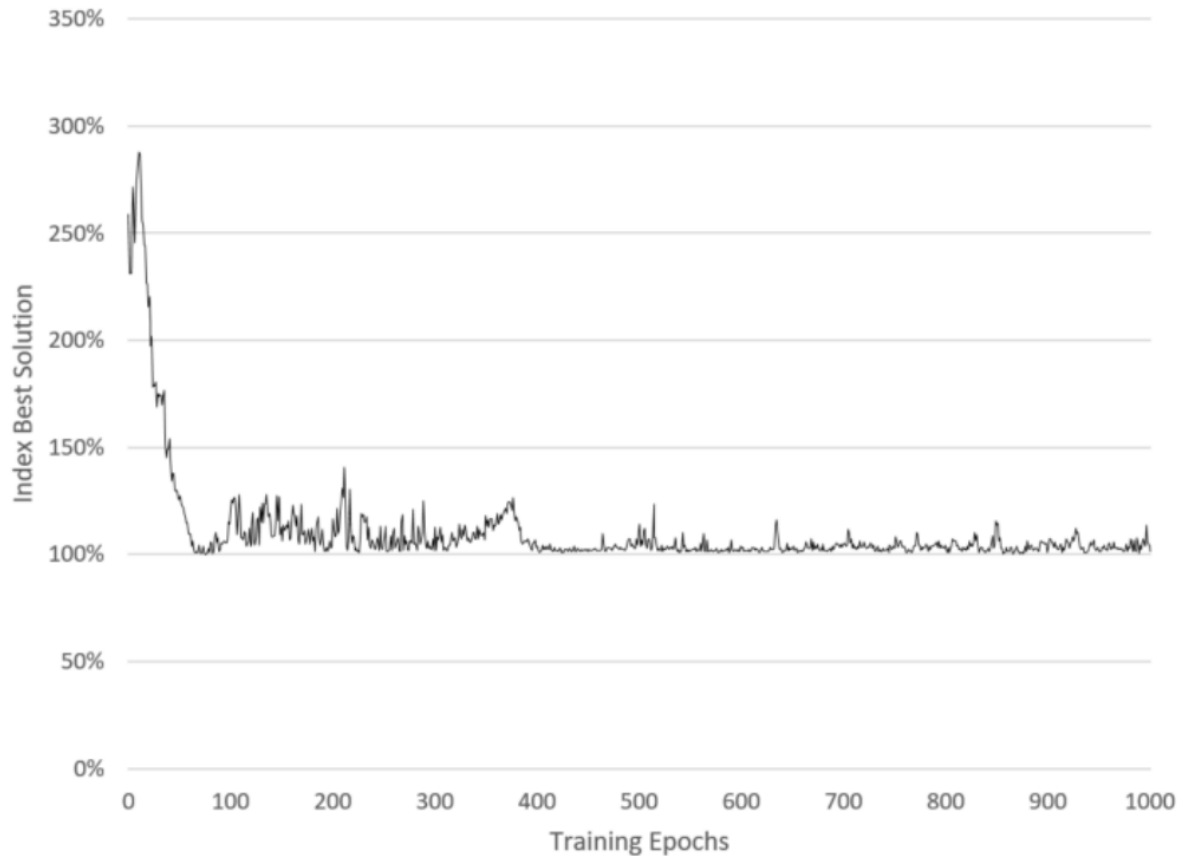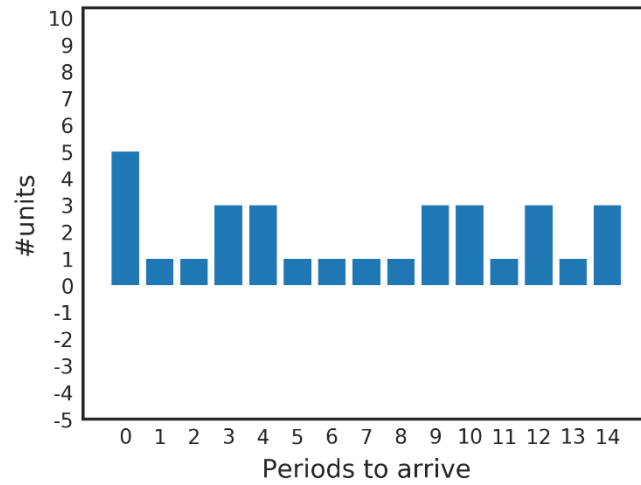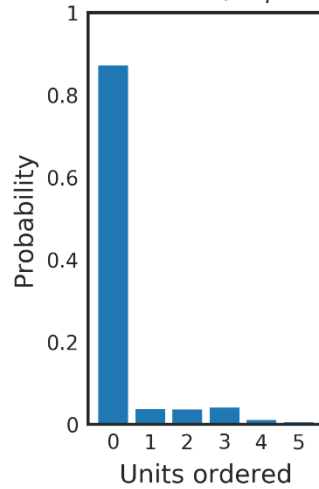## Deep Q-Learning algorithm learns itself a 'smart' replenishment policy
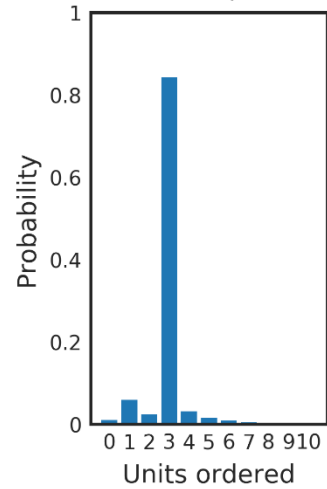


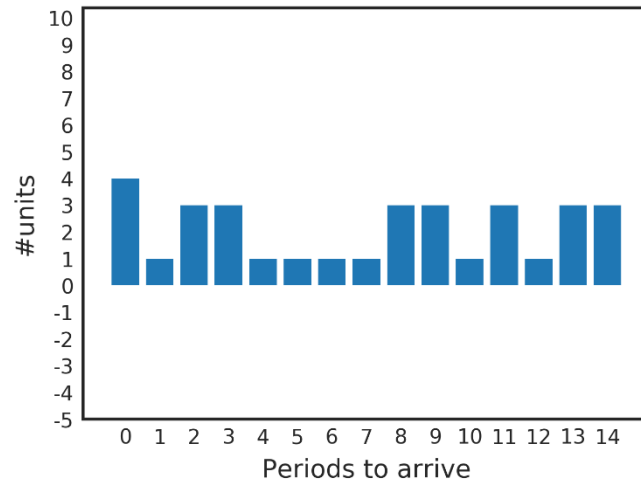Figure 1: Cost performance during training

KU LEUVEN

State: Inventory Pipeline

Action: Fast/Express

Action: Slow/Remote

State: Inventory Pipeline

Action: Fast/Express

Action: Slow/Remote

# Results
## Matching performance state-of-the-art dual sourcing policies with daily frequencies



DRL = Deep Reinforcement Learning

TBS = Tailored Base Surge

SIDB = Single Index Dual Base

DIDB = Dual Index Dual Base

CDIDB = Capped Dual Index Dual Base

Figure 2: Benchmark with state-of-the-art dual sourcing policies

- Equal performance typical dual sourcing setting (daily frequencies)

- We do not lose performance when extending problem to:
  - Non-daily ordering frequencies
    - E.g. including rail schedule or production schedule
  - Non-linear ordering cost (e.g. per container instead of per unit)

KU LEUVEN

# Backslides

# Methodology
## Smart replenishment algorithm – Deep Q-learning

---

**Algorithm 1** Deep Reinforcement Learning (DRL) algorithm

---

1: Initialize replay memory D to capacity $\phi$
2: Initialize Q-network with random weights $\theta$
3: Initialize Target Network $\hat{Q}$ with weights $\theta^- = \theta$
4: Choose Initial State $s_1$
5: **for** t $= 1, T$ **do**
6:     **for** n$=1, N$ **do**
7:         $a_t = \begin{cases} \text{random action,} & \text{with probability } (1 - \epsilon) \\ \text{argmin}_{a \in A}\{Q(s_t, a; \theta)\}, & \text{else} \end{cases}$
8:         Simulate using $a_t$ and observe reward $r_t$ and next state $s_{t+1}$
9:         Store $(s_t, a_t, r_t, s_{t+1})$ in replay memory D
10:     **end for**
11:     Sample random minibatches of size K from D
12:     **for** j$= 1, K$ **do**
13:         $y_j = \begin{cases} r_j + \gamma Q(s, \text{argmin}_{a' \in A} Q(s_{j+1}, a'; \theta^-), \theta), & \text{if } (s, a) \text{ in sample } j \\ Q(s_j, a; \theta), & \text{else} \end{cases}$
14:     **end for**
15:     Minimize loss $= \sum_{j \in K}(y_j - Q(s_j, a_j; \theta))^2$ using Adam optimizer
16:     Every $z$ steps, update Target network $\hat{Q} = Q$
17: **end for**

---

# Methodology
## Smart replenishment algorithm – A3C algorithm

**Algorithm S3** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

*// Assume global shared parameter vectors $\theta$ and $\theta_v$ and global shared counter $T = 0$*
*// Assume thread-specific parameter vectors $\theta'$ and $\theta'_v$*
Initialize thread step counter $t \leftarrow 1$
**repeat**
    Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.
    Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$
    $t_{start} = t$
    Get state $s_t$
    **repeat**
        Perform $a_t$ according to policy $\pi(a_t|s_t; \theta')$
        Receive reward $r_t$ and new state $s_{t+1}$
        $t \leftarrow t + 1$
        $T \leftarrow T + 1$
    **until** terminal $s_t$ **or** $t - t_{start} == t_{max}$
$$R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t \text{// Bootstrap from last state} \end{cases}$$
    **for** $i \in \{t-1, \ldots, t_{start}\}$ **do**
        $R \leftarrow r_i + \gamma R$
        Accumulate gradients wrt $\theta'$: $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$
        Accumulate gradients wrt $\theta'_v$: $d\theta_v \leftarrow d\theta_v + \partial(R - V(s_i; \theta'_v))^2 / \partial\theta'_v$
    **end for**
    Perform asynchronous update of $\theta$ using $d\theta$ and of $\theta_v$ using $d\theta_v$.
**until** $T > T_{max}$

KU LEUVEN

# Results

## Within 2% of optimal solution in a simple setting



Performance versus optimal for different lead times regular mode assuming Linear Ordering Cost

(LT_e = 0 // c_r = 100 //c_e = 150 // h=5 //b=495)

|  | LT_r = 1 | LT_r = 2 | LT_r = 3 | LT_r = 4 |
|---|---|---|---|---|
| ■ A3C | 102.34% | 101.74% | 100.06% | 100.01% |
| ■ SIDB | 100.00% | 118.71% | 131.60% | 128.35% |
| ■ DIDB | 100.00% | 102.80% | 105.06% | 108.18% |
| ■ CDIDB | 100.00% | 101.84% | 104.61% | 107.35% |
| ■ TBS | 119.06% | 136.07% | 131.86% | 128.94% |

■ A3C  ■ SIDB  ■ DIDB  ■ CDIDB  ■ TBS

**KU LEUVEN**

# Results
## Maintaining performance in a more complex setting



Performance versus optimal for different lead times regular mode assuming Stepwise Ordering Cost
(LT_e = 0 // c_r = 100 //c_e = 150 // h=5 //b=495//cap_r = 2 //cap_e =2)

|  | LT_r = 1 | LT_r = 2 | LT_r = 3 | LT_r = 4 |
|---|---|---|---|---|
| A3C | 102.49% | 101.76% | 101.03% | 108.16% |
| SIDB | 113.64% | 147.14% | 145.03% | 140.52% |
| DIDB | 112.73% | 122.49% | 129.38% | 132.49% |
| CDIDB | 111.91% | 114.88% | 126.23% | 122.54% |
| TBS | 123.41% | 148.17% | 146.13% | 142.28% |