

DeepBike: A Deep Reinforcement Learning Based Model for Large-scale Online Bike Share Rebalancing

Zhuoli Yin

zhuoliyin@purdue.edu

Purdue University

Zhaoyu Kou

Purdue University

Hua Cai

Purdue University

Research Article

Keywords: Bike share system, large-scale, dynamic rebalancing, deep reinforcement learning

Posted Date: March 4th, 2024

DOI: <https://doi.org/10.21203/rs.3.rs-3998473/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: No competing interests reported.

32 1. Introduction

33 Rebalancing bikes to accommodate real-time customer demands is becoming critical for
34 bike share system (BSS) operation as BSS is scaling up in major cities worldwide (Yin,
35 Hardaway, et al., 2023). The accessibility, convenience, and low cost of BSS have attracted lots
36 of customers. Additionally, BSS is regarded as a promising mobility solution for cities to help
37 mitigate congestion and reduce greenhouse gas (GHG) emissions (Fishman, 2016; Shaheen et
38 al., 2010). Due to the popularity and benefits of BSSs, BSS companies have continued to expand
39 their existing systems or launch new ones (Fishman, 2016; Shaheen et al., 2010). Nevertheless,
40 spatiotemporal demand heterogeneity across a city often leads to the imbalanced distribution of
41 bike share ridership. This imbalance will then cause bikes to be overstocked or depleted in
42 stations and impede the returning or renting of bikes, which thereby results in customer loss, and
43 revenue decrease, as well as a GHG emission increase (Ghosh et al., 2017). Consequently, to
44 avoid potential customer loss, the system operator needs to regularly rebalance bikes among
45 stations to maintain a reasonable distribution across the service region (Fishman, 2014). Among
46 various rebalancing strategies, dynamic bike share rebalancing becomes a major solution for
47 BSS operation. This dynamic strategy is crucial because it monitors the system in real-time to
48 maintain the stations to be well-functioning throughout the day, in contrast to static rebalancing
49 which merely optimizes the initial distribution of bikes in a system (Gleditsch et al., 2022; Shui
50 & Szeto, 2020; Vallez et al., 2021). Such a strategy is typically modeled as a Dynamic Bike
51 Share Rebalancing Problem (*DBSRP*). In specific, the rebalancing is usually conducted by a fleet
52 of automobiles, such as trucks or vans (referred to as *rebalancing vehicles* hereafter) (Citi Bike,
53 2023). The rebalancing process is sequential and involves two entangled decisions: (1) *routing*
54 decisions, which choose a series of bike stations that a rebalancing vehicle is to visit, and (2)
55 *repositioning* decisions, which determine the number of bikes to be collected from or distributed
56 into stations decided by the *routing* decisions.

57 *DBSRP* has been widely studied in the literature, but the prior works still have major
58 limitations in addressing the online¹ and large-scale needs of *DBSRP* (detailed discussions of the
59 existing literature are provided in Section 2). Traditional approaches primarily solved *DBSRP*
60 through Mixed Integer Programming (MIP) (Contardo et al., 2012; Ghosh et al., 2015, 2017; T.
61 Wang, 2014; Zheng et al., 2021). However, they are only applicable to offline scenarios. Their

¹ We will refer to online and real-time interchangeably.

62 solutions rely on complete data of the entire planning horizon, but in dynamic rebalancing,
63 customer demands for the upcoming planning horizon are unknown to the decision-makers. In
64 recent years, to generate online rebalancing solutions, heuristic algorithms and artificial
65 intelligence-based approaches (e.g., deep learning and reinforcement learning) have been
66 employed (Hu et al., 2021; Li et al., 2018; Shui & Szeto, 2018). However, these algorithms were
67 only designed for and evaluated on small-scale BSSs (typically less than 300 bike stations) or
68 they divided the study region into small clusters and dispatched only one rebalancing vehicle in
69 each cluster to simplify the problem (Vallez et al., 2021). Additionally, solutions from heuristic
70 algorithms are myopic by looking forward only a few time steps (Lowalekar et al., 2017). In
71 contrast, the size of BSS in major cities, such as Chicago and New York City in the U.S., has
72 increased to be large-scale to cover wider urban regions, including more than 500 stations and 10
73 rebalancing vehicles in service (General Bikeshare Feed Specification, 2023). Thus, the scale of
74 current online solutions no longer suffices the needs of these expanded systems.

75 Models that are capable of generating online and far-sighted solutions for large-scale
76 BSSs are essential to support effective rebalancing operations. The advancement of Deep
77 Reinforcement Learning (DRL) has demonstrated the potential opportunity to reach this goal.
78 Integrating Q-learning with a deep neural network, DRL presents a strong representation
79 learning capability for high-dimensional and large-scale system observations to estimate the
80 long-term (i.e., the period covering the current time to the end of the planning horizon) benefits
81 of decisions (de Bruin et al., 2018; Mnih et al., 2013, 2015). Once well trained, it has promising
82 effectiveness in online sequential decision-making based on real-time observations (J. Wang &
83 Sun, 2020). Leveraging these benefits, DRL has been widely adopted in urban mobility
84 management, such as on-demand ride sharing dispatching (Qin et al., 2019). Thus, given the
85 sequential decision-making nature of *DBSRP* (the formulation is detailed in Section 3), DRL
86 could also be employed to generate online and far-sighted rebalancing solutions for this problem.

87 In this study, we proposed a DRL-based framework—*DeepBike*— to generate the optimal
88 online solutions for large-scale *DBSRP*. *DeepBike* enables real-time dispatching for multiple
89 rebalancing vehicles to rebalance bikes. In the decision-making module within the *DeepBike*
90 framework (referred to as *agent* hereafter), we designed a Deep Q-Network (DQN), whose
91 inputs are real-time large-scale BSS state observations, and whose outputs are the long-term
92 quality values (Q-value) of rebalancing actions of each rebalancing vehicle. Through the

93 interaction with the BSS simulator, the agent iteratively learns to approximate the Q-values of
94 rebalancing actions and selects the highest-valued action for individual rebalancing vehicles. The
95 comparison with an MIP-based algorithm and a heuristic-based algorithm demonstrates that our
96 algorithm outperforms these benchmarks in reducing customer loss and improving the overall net
97 profits.

98 The remainder of this paper is organized as follows. Section 2 discusses the related
99 literature, research gaps, and the main contributions of this study. Section 3 provides a full
100 description and mathematical formulation of *DBSRP*. Section 4 discusses input data, and the
101 *DeepBike* framework including the simulator, demand prediction, and DRL model. Additionally,
102 it introduces the baseline algorithms and metrics that measure the performance of different
103 algorithms. In Section 5, the model training results, and model performance comparisons are
104 described. Last, Section 6 summarizes the findings and discusses the limitations of this study and
105 future research directions.

106 **2. Literature review**

107 Existing studies have proposed various algorithms to efficiently solve *DBSRP* since the
108 emergence of BSSs. The algorithms could mainly be classified into three groups: Mixed Integer
109 Programming (MIP), heuristic, and artificial intelligence. Table 1 summarizes the models in
110 existing literature and also compares them with our proposed model. As an NP-hard problem, the
111 complexity of *DBSRP* is governed by factors such as the scale of the bike stations (i.e. the
112 number of stations and capacity of each station), the size of the rebalancing fleet (i.e. the number
113 of rebalancing vehicles and capacity of each vehicle), and the length of planning horizon (i.e. the
114 number of decisions to be made over the entire planning horizon) (Lowalekar et al., 2017;
115 Schuijbroek et al., 2017; Shui & Szeto, 2020).

116 Traditionally, researchers modeled *DBSRP* through MIP or linear programming (LP),
117 treating it as a variant of the vehicle routing problem (VRP). However, these models only offer
118 offline solutions. For example, Contardo et al. (2012) proposed a mathematical formulation of
119 *DBSRP* for peak hours and derived the lower and upper bounds via Dantzig-Wolfe
120 decomposition and Benders decomposition. Shu et al. (2013) addressed the *DBSRP* based on the
121 stochastic network flow model and solved it via LP. (Ghosh et al., 2015, 2017; Zheng et al.,
122 2021) handled *DBSRP* by an MIP model with the objective of maximizing the overall profit. To

123 accelerate the computation due to the NP-hard nature of the problem, they first aggregated all
124 individual stations into different clusters and viewed each cluster as an abstract station. Then
125 they solved the *DBSRP* at the abstract station level by Lagrangian dual decomposition. Similarly,
126 Wang (2014) proposed an MIP formulation for *DBSRP* and solved the model through two
127 heuristic methods and one exact approach. These traditional mathematical algorithms yield stable
128 and exact solutions but require complete information on the considered planning horizon,
129 generating only offline solutions. They are unable to adapt to upcoming real-time planning
130 horizons with unknown information. Consequently, there exists a research gap in making online
131 rebalancing decisions during the daytime to accommodate dynamically changing bike usage
132 demands.

133 To facilitate online rebalancing decisions, many studies developed heuristic algorithms
134 but their solutions are often myopic. For example, Lowalekar et al. (2017) proposed a Greedy
135 Online Anticipatory Heuristic approach capable of offering online solutions within a short
136 computation time (e.g., one minute). This approach picks the best rebalancing action for each
137 rebalancing vehicle by anticipating several timeslots ahead and approximating the value of each
138 valid action. Shui and Szeto (2018) employed a rolling horizon method to break down the
139 planning period, solving a static bike repositioning sub-problem within each segment.
140 Additionally, within their model, an enhanced artificial bee colony (EABC) algorithm and a
141 route truncation heuristic were integrated to optimize the route design in each stage. Chiariotti et
142 al. (2018) exploited historical data to predict future bike station conditions to decide the desired
143 inventory level of each bike station. Then they selected the rebalancing actions that maximize the
144 “survival time” of stations. Even though the aforementioned approaches could generate online
145 rebalancing solutions, the selection of actions is myopic because they heuristically approximated
146 the benefits of rebalancing actions by looking ahead only a few steps within the entire planning
147 horizon, which is likely to reach the local optimum. However, *DBSRP* typically operates on an
148 hourly basis within a day so each rebalancing action not only has an immediate effect on the BSS
149 but also impacts subsequent system dynamics (O’Mahony & Shmoys, 2015; Li et al., 2018). For
150 example, stations in the downtown area have high morning check-in demands for commute
151 purposes and have symmetric high afternoon check-out demands (Zhou, 2015). Such stations
152 could potentially achieve self-rebalancing status without the need for extra rebalancing effort. An
153 online rebalancing decision that greedily moving bikes out of the self-rebalancing stations in the

154 morning might cause greater customer loss and rebalancing cost because it overlooks the path
155 dependence in the full planning horizon and further results in supply scarcity in those stations in
156 the afternoon (D. Chen & Sakai, 2022). Therefore, developing online as well as farsighted
157 solutions that evaluate the payoff of rebalancing decisions over the entire planning horizon is still
158 a gap for *DBSRP*.

159 Prior works also leverage the capability of artificial intelligence in sequential decision-
160 making to solve *DBSRP* and provide online solutions (Yin, Kou, et al., 2023). However, their
161 models together with previous works with MIP and heuristic algorithms were only designed on
162 small-scale BSSs or a small subset of BSSs with the deployment of a small size of rebalancing
163 fleet in BSSs (less than 300 bike stations and fewer than 5 rebalancing vehicles in the fleet)
164 (Vallez et al., 2021). For example, Li et al. (2018) dealt with *DBSRP* through a spatial-temporal
165 reinforcement learning framework aiming to minimize long-term customer loss. They first
166 proposed an interdependent inner-balance clustering algorithm and then designated only one
167 vehicle to each cluster. Nevertheless, their approach was only applicable to a subset of the Citi
168 Bike system and overlooks spatiotemporal imbalance of customer demand among those clusters.
169 Assigning one vehicle to each cluster is inadequate to handle this disparity. Given the
170 heterogeneity of customer demands across different regions, such as downtown versus suburban
171 areas, the allocation of rebalancing vehicles should be smartly and efficiently adjusted for global
172 demands and overall system dynamics in the entire planning horizon. As such, their model is
173 overly simplified for the operation of multi-rebalancing vehicles and large-scale stations in
174 current BSSs. Notably, (Zhu et al., 2022) also employed reinforcement learning to tackle large-
175 scale *DBSRP*. However, they employed a user incentivized crowd-sourcing rebalancing strategy
176 which is different from the *DBSRP* with centralized rebalancing managed by the operators. Chen
177 et al. (2021) applied deep learning to learn from historical solution instances and utilized the
178 trained model for future cases. The solutions heavily rely on the quality of solution instances
179 from which their model learned. Additionally, the model can only generate rebalancing decisions
180 for one rebalancing vehicle, regardless of the coordination between multiple rebalancing vehicles
181 in service. Additionally, in terms of the scale of solutions in previous studies, such as (Contardo
182 et al., 2012; Shu et al., 2013; Kloimüller et al., 2014; T. Wang, 2014; Regue & Recker, 2014),
183 they proposed the *DBSRP* models for BSSs of fewer than 100 stations. Other works such as
184 (Ghosh et al., 2015; Lowalekar et al., 2017; D. Zhang et al., 2017; Ghosh et al., 2017; Brinkmann

185 et al., 2020) dealt with *DBSRP* with fewer than 400 stations in the BSSs. Furthermore, in
186 (Chiariotti et al., 2018; Ghosh et al., 2016; T. Wang, 2014), the problem was simplified by
187 allocating a single rebalancing vehicle to the service region. However, current station-based
188 BSSs in major cities have grown to possess more than 500 stations in their service regions. For
189 instance, as of 2023, Citi Bike in New York City and Velib’ Metropole in Paris maintained
190 stations of 2,000 and 1,472, respectively (*General Bikeshare Feed Specification, 2023*).
191 Moreover, current operators typically deploy multiple rebalancing vehicles in lieu of a single one
192 to meet rebalancing needs. For example, (Sun, 2021) reported that BSSs in Chicago, Boston and
193 Los Angeles need 16, 10 and 5 rebalancing vehicles, respectively. Therefore, existing algorithms
194 fall short of accommodating the expanding scale of BSSs and there is a gap in developing
195 models adequate for generating high-quality online rebalancing solutions for large-scale BSSs
196 with multiple rebalancing vehicles.

197 To address the above-mentioned research gaps, we proposed a DRL-based model—
198 *DeepBike*— to generate the optimal strategy for the dynamic rebalancing of large-scale BSS in
199 real-time. Specifically, in the proposed framework, an agent interacts with the BSS simulator
200 (which simulates the dynamics of the BSS) iteratively to learn the optimal policy, whose
201 objective is to maximize the system’s overall profits. We designed a Deep Q-Network (DQN) for
202 the agent to estimate the long-term (from the current timeslot to the end of the planning horizon)
203 benefits of the rebalancing actions, given the global system observations as network inputs. The
204 agent selects the rebalancing action with the highest benefit for each rebalancing vehicle. Our
205 model differs from the existing work in an important aspect that we incorporated large-scale,
206 real-time, and farsighted rebalancing decision-making simultaneously, catering to BSS
207 operations in large-scale cities nowadays. In the context of the existing literature, we summarize
208 the main contributions of this study as follows:

- 209 (1) We proposed a DRL-based framework (*DeepBike*) tailored for *DBSRP*. Different from
210 existing literature, our model focuses on optimizing large-scale *DBSRP* in real-time,
211 without partitioning the service region into smaller clusters.
- 212 (2) We constructed a convolutional neural network (i.e., DQN) which is capable of
213 processing large-scale system observations covering the entire service region. Once the
214 network converges, the *DeepBike* model generates instantaneous rebalancing decisions
215 for each rebalancing vehicle every hour in a day based on real-time observation.

216 (3) We trained the DQN through the Double Deep Q-Network (DDQN) algorithm, enforcing
217 the network to approximate the long-term benefits of rebalancing decisions from the
218 current step to the end of the planning horizon to reach the far-sighted optimum.

219 (4) We evaluated the *DeepBike* model on the real-world dataset of Divvy Bike Share in
220 Chicago which has more than 500 bike stations and 16 rebalancing vehicles. The
221 empirical results demonstrate the effectiveness of our model in improving the overall
222 system's profits.

223 Table 1. Summary of the related literature on dynamic bike share rebalancing problem and comparison with this study.

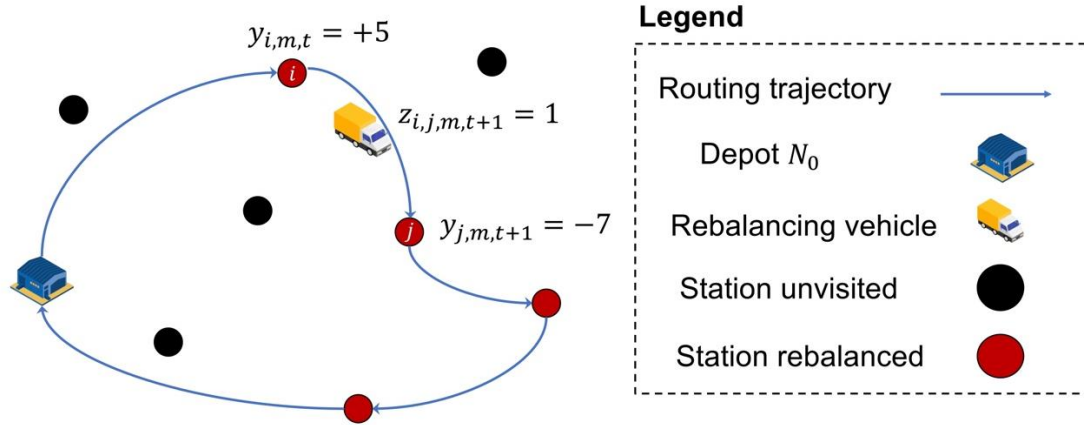
Paper	Offline/ Online	# of stations	Clustering/ # of clusters	# of rebalancing vehicles	Planning horizon	Algorithm
Contardo et al. (2012)	Offline	25/50/100	Yes/9	5	2 hours	MIP
Shu et al. (2013)	Online	100	No	/	24 hours	LP
<u>Ghosh et al. (2015)</u>	Offline	95	Yes/25	4	7 hours	MIP
<u>Ghosh et al. (2017)</u>	Offline	95/305	Yes/25/30	3/5	24 hours	MIP
Zheng et al. (2021)	Offline	60/95/305	Yes/12	9/13/45	19 hours	MIP
Wang (2014)	Offline	12/24/36/48	No	1	1 hour/2 hours	MIP/Heuristic
Kloimüllner et al. (2014)	Offline	30/60/90	No	1/2/3/4/5	4 hours/8 hours	Heuristic
Regue and Recker (2014)	Offline	61	No	2	24 hours	Heuristic
<u>Lowalekar et al. (2017)</u>	Online	95/305	No	3/6	6 hours	MIP/Heuristic
Zhang et al. (2017)	Online	100/104/200	No	2	9 hours	Heuristic
Shui and Szeto (2018)	Online	30/60/90/120/180	No	1	6/12/18 time steps	Heuristic
Chiariotti et al. (2018)	Online	280	No	1	24 hours	Heuristic
<u>Li et al. (2018)</u>	Online	389	Yes/4	4	24 hours	DRL
Brinkmann et al. (2020)	Online	35/169	No	1/2/3/4	24 hours	MIP+heuristic
<u>Chen et al. (2021)</u>	Online	149	No	1	8 hours	DL
This study (DeepBike)	Online	578	No	16	24 hours	DRL

224 Note: MIP: Mixed Integer Programming; LP: Linear Programming; DL: Deep Learning; DRL: Deep Reinforcement Learning

225 **3. Problem formulation**

226 In this section, we begin by illustrating the *DBSRP* with an example. Then we describe
 227 the generic MIP formulation of *DBSRP*. Notations employed throughout this paper are detailed
 228 in Appendix Table A-1.

229 As demonstrated in Fig. 1, the rebalancing process involves a rebalancing vehicle,
 230 typically a truck or van, initiating its routes from a depot at the start of the planning horizon and
 231 ending its route at the same depot. The primary objective of this rebalancing vehicle is to
 232 reposition shared bikes across various stations, either by picking up or dropping off bikes. This
 233 rebalancing process ensures a reasonable inventory level of bikes at each station, thereby
 234 facilitating uninterrupted bike usage for customers. Meanwhile, customers rent and return bikes
 235 at these stations throughout the day. In the event that a station either runs out of bikes or docking
 236 space, customers fail to rent or park bikes at their visited stations. These service failures are
 237 considered as customer loss.



238
 239 Fig. 1. An illustration of *DBSRP* studied in this work. In this example, the rebalancing vehicle m
 240 departs from the depot N_0 . It routes to bike station i and moves 5 bikes into the station in
 241 timeslot t and then travels to bike station j and moves 7 bikes out from the station in timeslot $t +$
 242 1. The routing decision that vehicle m travels from station i to station j is denoted as $z_{i,j,m,t+1} =$
 243 1, while the repositioning decisions at station i in timeslot t and at station j in timeslot $t + 1$ are
 244 denoted as $y_{i,m,t}$ and $y_{j,m,t+1}$, respectively. After visiting the subsequent stations, the rebalancing
 245 vehicle returns to the depot at the end of the planning horizon.

246 We depict the comprehensive dynamics of *DBSRP* in Eqs. (1) – (12) based on MIP
 247 formulations proposed by (Ghosh et al., 2017; Luo et al., 2020). These formulations describe
 248 customer interaction with BSS and operational constraints involved in the rebalancing process,
 249 where decision variables involve repositioning variables $y_{i,m,t}$ and routing variables $z_{i,j,m,t}$:

$$250 \quad \min_{z,y} (\alpha \cdot \sum_{i \in N, t \in T} CL_{i,t} + \beta \cdot \sum_{i,j \in N, m \in M, t \in T} P_{i,j} \cdot z_{i,j,m,t}) \quad (1)$$

251 s.t.

$$252 \quad S_{i,t+1}^{Bike} = \min \left\{ \max \left\{ 0, S_{i,t}^{Bike} + \sum_{m \in M} y_{i,m,t} + D_{i,t} \right\}, C_i^{station} \right\} \quad \forall i \in N, t \in T \quad (2)$$

$$253 \quad 0 \leq S_{i,t}^{Bike} + \sum_{m \in M} y_{i,m,t} \leq C_i^{station} \quad \forall i \in N, t \in T \quad (3)$$

$$254 \quad CL_{i,t} = \left| S_{i,t+1}^{Bike} - (S_{i,t}^{Bike} + \sum_{m \in M} y_{i,m,t} + D_{i,t}) \right| \quad \forall i \in N, t \in T \quad (4)$$

$$255 \quad V_{m,t+1}^{Bike} = V_{m,t}^{Bike} + \sum_{m \in M} -y_{j,m,t} \quad \forall m \in M, t \in T \quad (5)$$

$$256 \quad \sum_{j \in N} z_{i,j,m,t} - \sum_{j \in N} z_{j,i,m,t-1} = 0 \quad \forall i \in N, m \in M, t \in T \quad (6)$$

$$257 \quad \sum_{i \in N, m \in M} z_{i,j,m,t} \leq 1 \quad \forall j \in N, t \in T \quad (7)$$

$$258 \quad \sum_{i,j \in N} z_{i,j,m,t} \leq 1 \quad \forall m \in M, t \in T \quad (8)$$

$$259 \quad -C_m^{vehicle} \cdot \sum_{i \in N} z_{i,j,m,t} \leq y_{j,m,t} \leq C_m^{vehicle} \cdot \sum_{i \in N} z_{i,j,m,t} \quad \forall j \in N, m \in M, t \in T \quad (9)$$

$$260 \quad S_{i,0}^{Bike} = C_0^{station} \quad \forall i \in N, \quad V_{m,0}^{Bike} = C_0^{vehicle} \quad \forall m \in M \quad (10)$$

$$261 \quad 0 \leq S_{i,t}^{Bike} \leq C_i^{station}, 0 \leq V_{m,t}^{Bike} \leq C_m^{vehicle} \quad \forall i \in N, m \in M, t \in T \quad (11)$$

$$262 \quad z_{i,j,m,t} \in \{0,1\}, y_{i,m,t} \text{ is integer} \quad \forall i, j \in N, m \in M, t \in T \quad (12)$$

263 The objective function that is described in Eq. (1) minimizes the total costs of customer
 264 loss (the sum of the customer loss from each station i in each timeslot t , $CL_{i,t}$, the definition of
 265 which is detailed in Section 4.4) and fuel cost of rebalancing vehicles that route to target stations
 266 (the sum of the routing distance between stations i and j , $P_{i,j} \cdot z_{i,j,m,t}$). Additionally, in Eq. (1), α
 267 represents the average price paid by customers for a single bike share trip and β represents the
 268 average fuel cost per mile of vehicle use. The objective function also considers the tradeoff
 269 between the reduced customer loss (to improve the revenue from trips) and the cost of routing
 270 rebalancing vehicles that is required.

271 For constraints, Eqs. (2)-(3) conserve the in- and out-flows of bikes in the station. The
 272 number of bikes at the station either after being rebalanced or after the rent/return of customers
 273 should always stay within the station's capacity, $C_i^{station}$. Eq. (4) calculates the customer loss in

274 each station, which is the difference between the expected inventory level of the station if all
 275 customer demands could be fulfilled after being rebalanced (that is, $S_{i,t}^{Bike} + \sum_{m \in M} y_{i,m,t} + D_{i,t}$) and
 276 the actual inventory level constrained by Eq. (2) (that is, $S_{i,t+1}^{Bike}$). Eq. (5) defines the transition of
 277 bikes stocked in the rebalancing vehicles over time. The inflow and outflow of the bikes on the
 278 vehicles are the result of the repositioning decision $\sum_{m \in M} -y_{j,m,t}$. Eq. (6) enforces that the
 279 rebalancing vehicles departing from a station $\sum_{j \in N} z_{i,j,m}^t$ could only be the ones that arrive at this
 280 station in the previous timeslot $\sum_{j \in N} z_{j,i,m}^{t-1}$. Note that staying in a station could be viewed as
 281 moving out of the station and moving into it again. Intuitively, $\sum_{j \in N} z_{N_0,j,m,0}$ and $\sum_{i \in N} z_{i,N_0,m,|T|}$
 282 are set to be 1 for every rebalancing vehicle m to ensure it to depart from the depot initially and
 283 return to it at the end. Eqs. (7) and (8) guarantee that, in any timeslot, no more than one vehicle is
 284 visiting the same station, and each vehicle can only visit at the most one station, respectively.
 285 These constraints avoid the risk of multiple vehicles rebalancing the same station, thereby
 286 preventing offsetting effects. Eq. (9) couples the repositioning decision and routing decision. It
 287 enforces that rebalancing vehicles can only pick up or drop off bikes in physically visited
 288 stations. That is, $y_{j,m,t}$ will be forced to be zero if no rebalancing vehicles $\sum_{i \in N} z_{i,j,m,t}$ are
 289 assigned to visit station j departing from each station i . Eq. (10) sets the initial inventory level of
 290 bikes among stations (constant $C_0^{station}$) and the initial capacity of the rebalancing vehicles
 291 (constant $C_0^{vehicle}$). Eq. (11) confines the maximum capacities of bike stations and rebalancing
 292 vehicles, denoted by the constants $C_i^{station}$ and $C_m^{vehicle}$, respectively. Eq. (12) defines that the
 293 routing decision variable $z_{i,j,m,t}$ is a binary and repositioning decision variable $y_{i,m,t}$ is an
 294 integer. The variable $z_{i,j,m,t}$, when setting to 1, represents vehicle m in the timeslot t will travel
 295 from station i to station j . When setting $y_{i,m,t}$ to an integer $K/-K$, it indicates rebalancing vehicle
 296 m moves K bikes into/out of the station i in the timeslot t .

297 4. Data and methods

298 This section introduces our proposed model—*DeepBike*—designed for the *DBSRP*, as
 299 well as the associated data used to train and evaluate the model. Specifically, Section 4.1
 300 describes the Divvy BSS data from Chicago and the preprocessing to align the data with our
 301 framework. In Section 4.2, we detail the *DeepBike* framework based on the RL paradigm,
 302 involving the environment, agent, state, action, and reward. More specifically, Section 4.2.1

303 establishes a BSS simulator (i.e., the environment)², which replicates BSS dynamics and
304 provides information on BSS status for the agent to make rebalancing decisions. Following this,
305 we develop a prediction model in Section 4.2.2, which forecasts next-hour customer demands
306 across the stations based on the BSS status information. This model offers predicted demands for
307 the agent to assist in generating rebalancing actions. In Section 4.2.3, we devise the state space,
308 action space, and reward function within the framework of DRL. Then, the Double-DQN
309 algorithm and DQN architecture are introduced in Sections 4.2.4 and 4.2.5, respectively. Section
310 4.3 introduces baseline models for comparative analysis. Last, Section 4.4 outlines evaluation
311 metrics for measuring *DBSRP* algorithm performance.

312 4.1. Case study data and preprocessing

313 We chose Divvy BSS in Chicago as our case study system because Divvy is a well-
314 established station-based BSS and has been actively used since 2013. Divvy’s trip records and
315 station information were employed to build the simulator to train the *DeepBike* model. The
316 historical trip records are publicly available, spanning from the year 2013 to the year 2023
317 (*Divvy System Data*, 2023). For our case study, we selected four consecutive weeks of historical
318 trip data from Sept. 2nd, 2019 (Monday) to Sept. 29th, 2019 (Sunday). This period includes over
319 449,000 trip records. Specifically, we used the start and end time of a bike trip, and the trip’s
320 origin and destination station IDs from the trip records. Following this, the data from the first
321 three weeks was employed for model training and the data of the remaining one week was
322 reserved for evaluation purposes. The station information of Divvy in the year 2019 was
323 collected through GBFS (General Bikeshare Feed Specification) API. The station information
324 includes the station ID, the corresponding location (latitude & longitude), and capacity. Those
325 station IDs are consistent with the trips’ origin and destination station IDs. In addition to the
326 above data, we also utilized Divvy’s third quarter data in 2018 to train a separate demand
327 prediction model, prior to *DeepBike* model training (detailed in SI-S2).

328 In the preprocessing phase, we trimmed 40 inactive stations that displayed no trip
329 activities in September 2019. This results in 578 active physical stations across Chicago. We
330 considered a rectangle region bounded by the locations of active stations as the service region.
331 Additionally, following the gridding approach used by Al-Abbasi et al. (2019) and Oda and Joe-

² In the following content, the terms “simulator” and “environment” are used interchangeably.

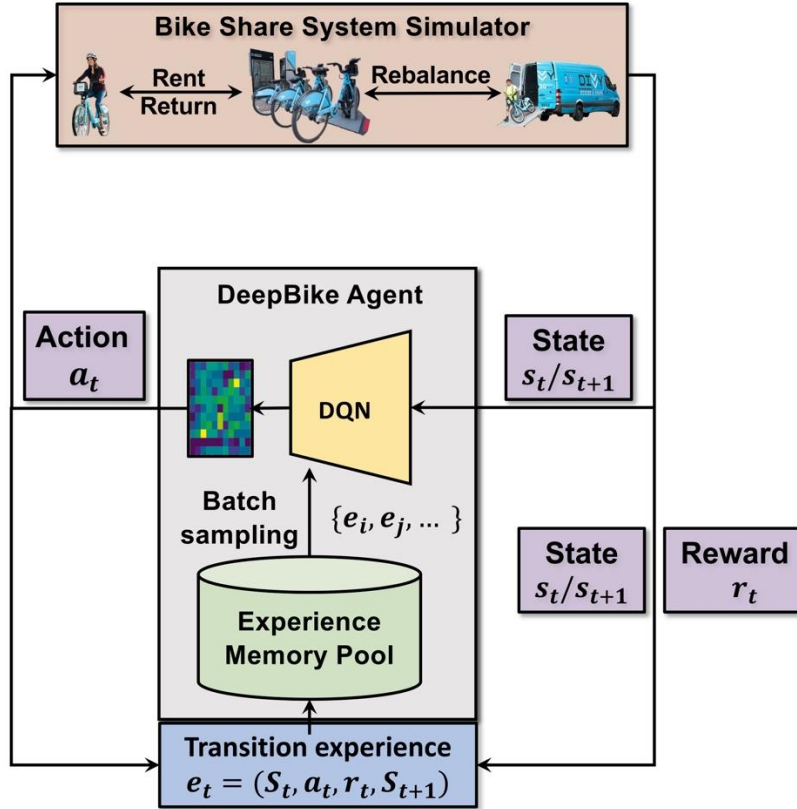
332 Wong (2018), we divided the service region into 51×51 grids with each grid measuring
333 $400\text{ m} \times 700\text{ m}$. This gridding approach aligns the service region with the input requirements of
334 the DQN (detailed further in Section 4.2.2). Stations residing within the same grid are
335 consolidated as one abstract station³, leading to 451 abstract active stations in total. This results
336 in the number of physical stations within a single abstract station ranging from zero to seven.
337 Additionally, considering the grid size, physical stations within one abstract station are close
338 enough to each other such that the subsequent allocation of rebalancing within an abstract station
339 can be intuitively determined by the operators (Ghosh et al., 2017; Osorio et al., 2021). The
340 capacity of the abstract station is therefore the cumulative capacities of all physical stations
341 within that grid. The estimated route between two abstract stations is calculated based on the
342 average distance between all physical station pairs within the respective grids. For each station
343 pair, the route distance is approximated using road network distance from the Open Street Map
344 dataset for cars in Chicago based on 578 stations’ physical locations.

345 The physical address of the bike service warehouse of Divvy (i.e., the depot) is located at
346 (41.89002, -87.658054) on Google Maps, which corresponds to the grid (27, 26) in this study.
347 The rebalancing fleet of Divvy is estimated to include at least 16 automobiles in service, as
348 reported by Sun (2021). Considering the work shift, we assumed that these vehicles operate daily
349 from 6:00 am to 11:00 pm.

350 4.2. *DeepBike* framework

351 The overall framework of *DeepBike* is illustrated in Fig. 2 and each component within the
352 framework is detailed in the subsections. A *DeepBike* agent interacts with the simulator (Section
353 4.2.4) through three components—states, actions and rewards (Section 4.2.1). The components
354 are continuously stored as experiences. Based on the observed states, the *DeepBike* estimates the
355 Q-values (long-term benefits) of actions via a Deep Q-Network (Section 4.2.2). The *DeepBike*
356 iteratively improves its Q-value estimation by learning from the collected experiences based on
357 the Double Deep Q-Network algorithm (Section 4.2.3).

³ In the following context, the abstract station will be referred to as the station or grid.



358
 359 Fig. 2. Illustration of the *DeepBike* learning framework following the RL paradigm. The agent
 360 iteratively interacts with the simulator to collect experiences to learn the relationship between the
 361 input states and the long-term benefits of actions. The agent generates an action by DQN based
 362 on the observed state. The simulator executes the action and simulates the bike usage. It then
 363 sends the updated state and reward to the agent as feedback. States, actions, and rewards are
 364 continuously stored in the memory pool, from which the agent updates the parameters of its
 365 DQN based on the experiences sampled.

366 4.2.1. State, action, and reward

367 In the interactions between the *DeepBike* agent and the BSS simulator, the transition from
 368 timeslot t to $t + 1$ is characterized by three fundamental components: state, action, and reward.
 369 Given a specific rebalancing policy π (i.e., the decision-making rule), each transition is
 370 formulated by the tuple $(s_t, a_t, r_t, s_{t+1}, \pi)$ (Sutton & Barto, 2018).

371 “State” depicts the observed status of the environment. In this study, the state is populated
 372 from the real-time information of the BSS environment on bike stations, rebalancing vehicles,
 373 and predicted future customer demands. Mathematically, we denoted the state at the beginning of
 374 timeslot t as $s_t = (S_{i,t}^{Bike}, S_{i,t}^{Dock}, V_{m,t}^{Bike}, V_{m,t}^{Slot}, V_{m,t}^{Loc}, D_{i,t:t+1}^{pred}), i \in \{1, 2, \dots, N\}, m \in \{1, 2, \dots, M\}$,
 375 where $S_{i,t}^{Bike}$ and $S_{i,t}^{Dock}$ are the number of parked bikes and available docks in station i of

376 timeslot t , respectively; $V_{m,t}^{Bike}$, $V_{m,t}^{Slot}$ and $V_{m,t}^{Loc}$ represent the number of bikes onboard, available
 377 slots, and the location of the rebalancing vehicle m in timeslot t ; $D_{i,t:t+1}^{pred}$ is the forecasted
 378 customer demands across the grids in the next timeslot, which is detailed in S3 in SI).

379 “Action” is the decision generated based on the current state according to a specific
 380 policy, leading to the transition of the environment to a subsequent state. Here, the rebalancing
 381 actions align with the decision variables in the MIP model Section 3. Thus, the rebalancing
 382 actions generated by *DeepBike* are denoted as $a_{m,t} = (z_{i,j,m,t}, y_{j,m,t})$ for rebalancing vehicle m
 383 in timeslot t , where $z_{i,j,m,t}$ is the routing from station i to j and $y_{j,m,t}$ is the number of bikes to
 384 be collected/deposited at the target station j . To reduce the computation time, the action space
 385 for routing action $z_{i,j,m,t}$ is limited to stations reachable within ten horizontal or vertical grids
 386 from the current station, which totals 441 discrete actions. Such a reachability constraint also
 387 ensures that a rebalancing vehicle can arrive at a target station within one timeslot. Meanwhile,
 388 we aggregated the repositioning actions into seven discrete actions to reduce the computational
 389 complexity and avoid inefficient rebalancing: (1) do not move any bike in the target bike station,
 390 (2,3) move five in/out of the target station, (4,5) move ten bikes into/out of the target station, (6)
 391 move maximum feasible bikes out of the station (“greedy all out”), and (7) move maximum
 392 feasible bikes from the rebalancing vehicle into the station (“greedy all in”).

393 “Reward” is the signal received after transitioning from the current state to the next one
 394 by taking an action. The overarching aim of an agent in DRL is to select actions that maximize
 395 the expected total future rewards until the episode (i.e. the entire planning horizon) terminates
 396 (Mnih et al., 2015; Oda & Joe-wong, 2018) as described in Eq. (13):

$$397 \quad \max \sum_{t'=t}^{\infty} \gamma^{t'-t} r(a_{t'}, s_{t'}) \quad (13)$$

398 where $r(a_{t'}, s_{t'})$ is the reward function calculating a scalar reward value given the state $s_{t'}$ and
 399 action $a_{t'}$; $\gamma \in [0, 1]$ is the time discount factor which discounts future rewards; t is the current
 400 timeslot; and t' is the future timeslot. In this study, we designed the reward function for a
 401 rebalancing vehicle m that executes rebalancing action $a_{m,t}$ in timeslot t given the state s_t as:

$$402 \quad r(a_{m,t}, s_t) = \alpha \cdot (CL_{j,t} - \widehat{CL}_{j,t}) - \beta \cdot P_{i,j} \cdot z_{i,j,m,t} \quad (14)$$

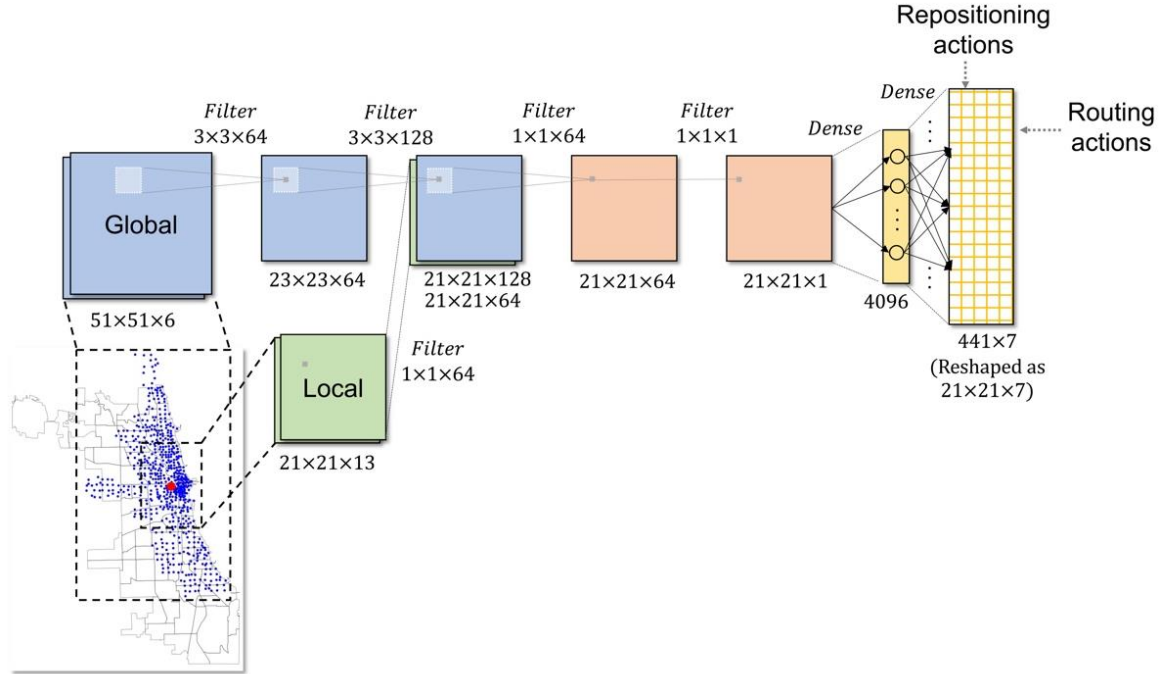
403 Eq. (14) is a modification to the objective function Eq. (1) to reflect an immediate reward
 404 in current timeslot, where $(CL_{j,t} - \widehat{CL}_{j,t})$ is the delta customer loss in the target station j in

405 timeslot t and the $P_{i,j} \cdot z_{i,j,m,t}$ is the routing distance from station i to j . Specifically, the delta
406 customer loss is the difference between the customer loss $CL_{j,t}$ occurring with executed
407 rebalancing under a given policy and the baseline customer loss $\widehat{CL}_{j,t}$ (i.e., the loss if no
408 rebalancing was performed). This modification results in a scalar reward that can range between
409 positive and negative values, thereby better indicate what *DeepBike* is supposed to accomplish
410 and speed up reaching optimality (Sutton & Barto, 2018). Considering that the *DeepBike* is
411 trained offline before deploying online, baseline customer loss is always attainable following the
412 simulation setup in this work.

413 4.2.2. DQN architecture

414 The design of a neural network is crucial in effectively estimating the long-term benefits
415 of actions in a given state (Silver et al., 2016). Such a neural network is known as a Deep Q-
416 Network (DQN) while the long-term benefits of actions corresponding to a state are denoted as
417 the Q-value $Q_{\pi}(s_t, a_t)$. In this study, we designed a convolutional neural network architecture as
418 the DQN, motivated by (Oda & Joe-wong, 2018; Silver et al., 2016; Zhao et al., 2021).

419 As presented in Fig. 3, we designed two input channels for the DQN to gather both global
420 and local observations to capture the spatial-temporal information of the entire BSS service
421 region. Specifically, the global observation is viewed from the entire BSS system’s aspect, while
422 the local observation is viewed from each rebalancing vehicle’s aspect centering the rebalancing
423 vehicle in the plane. The DQN output is structured as a two-dimensional 441×7 table to
424 approximate the Q-values for different rebalancing actions. This tabular approach is designed
425 due to the intertwined nature of routing and repositioning. Each table dimension aligns with the
426 action spaces of repositioning and routing as described in Section 4.2.1, respectively. Thus,
427 selecting an action from this table effectively maps to corresponding rebalancing actions in both
428 dimensions simultaneously. Detailed descriptions of the input features and the process by which
429 the convolutional neural network estimates the outputs are provided in the Supplementary
430 Information (SI) Sections S-1 and S-2, respectively.



431
432 Fig. 3. Deep Q-Network architecture of *DeepBike*.

433 We further applied an *action mask* on the output Q-value table to eliminate infeasible
434 actions. This *a priori* constraint excludes specific actions from selection, preventing infeasible
435 actions from entering the experience pool and improving the exploration efficiency. Specifically,
436 pixels in the Q-value table (i.e., individual rebalancing actions) will be masked as zero under three
437 conditions: (1) if a grid lacks an installed station; (2) the station has already been assigned to be
438 visited by another rebalancing vehicle (as per Eq. (7).); and (3) if available bikes parked at the
439 intended stations or bike inventory in a rebalancing vehicle are insufficient for fulfilling the
440 rebalancing action associated with this pixel.

441 4.2.3. Action selection and DQN learning

442 We first introduce how an action is chosen from the action space and then describe how a
443 DQN is trained to become a better estimator. Essentially, to avoid the overestimation of Q-
444 values, two DQNs with identical architecture (that is, a primary network with weights θ and a
445 target network with weights θ^-) were employed when selecting actions and improving DQN.

446

447

448 During the interaction with the environment, the rebalancing action to be executed is
 449 determined by the primary network. For a given state s_t , this network estimates the Q-values of
 450 the rebalancing actions for individual rebalancing vehicles (Mnih et al., 2013; Silver et al.,
 451 2016). Formally, the Q-value of an action a_t in a state s_t is defined in Eq. (15) as the expected
 452 cumulative future rewards yielded from selecting the immediate action and future actions
 453 following a policy thereafter:

$$454 \quad Q(s_t, a_t; \theta) \equiv \mathbb{E} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} r_k \mid S = s_t, A = a_t, \pi \right] \quad (15)$$

455 where $\gamma \in [0,1]$ is a discount factor that compensates for the importance of immediate and future
 456 rewards, which was set to be 0.9. Based on Eq. (13), we always choose the highest-valued action
 457 $a_t^* = \underset{a_t}{\operatorname{argmax}} Q^*(s_t, a_t; \theta)$ to maximize the cumulative rewards to be received, where
 458 $Q^*(s_t, a_t; \theta)$ denotes the optimal Q-value (Van Hasselt et al., 2016).

459 Prior to the DQN being fully trained and reaching convergence, the selected action may
 460 only be sub-optimal because the Q-values can be poorly estimated. Therefore, the *exploration-*
 461 *exploitation* scheme is applied to trade off gathering a new observation in an unknown state-action
 462 space and making the best decision based on the current DQN (Silver, 2014). We implemented the
 463 ϵ -greedy policy to select rebalancing actions during the training stage. Specifically, in each
 464 iteration, with probability $1 - \epsilon$, we select the highest valued action from the DQN as described
 465 above (i.e., exploitation); while with probability ϵ , we choose a random action from the action
 466 space (i.e., exploration) (Mnih et al., 2013). Thus, the exploration will lead to discovering more
 467 rewarding spaces that are unseen from the existing experiences. In contrast, as the DQN learns to
 468 correctly identify the actions yielding higher rewards, exploitation will reinforce the current DQN.
 469 In this study, we set the initial ϵ to be the commonly adopted value of 0.99. It decreases linearly
 470 to 0.05 over 40,000 training timeslots and remains unchanged afterward.

471 As the interaction progresses, the experiences attained over time, represented as $e_{t'} =$
 472 $(s_{t'}, a_{t'}, r_{t'}, s_{t'+1})$, are streamlined as a replay buffer, denoted as $\mathcal{D} = \{e_1, \dots, e_{t'}, \dots, e_R\}$ (Lin,
 473 1992). The buffer \mathcal{D} retains the most recent R experiences, from which DQN will learn. In this
 474 study, we set R to 40,000.

475 DQN progressively improves its Q-value estimation by learning from accumulated
 476 experiences, which is parallel to ongoing interaction with the environment. We adopted the
 477 Double-Deep Q-Network (DDQN) algorithm proposed by Van Hasselt et al. (2016) to train our
 478 designed DQN. The DQN is improved (i.e., the weights are tuned) by minimizing the error
 479 between the predicted Q-value of the selected action and a target Q-value updated based on the
 480 true reward when performing this action. Specifically, in each timeslot t , we uniformly retrieve a
 481 batch of experiences at random according to a pre-determined batch size (set to 128 in this study).
 482 These experiences are used to update the weights θ_t of the primary network by minimizing the
 483 mean-squared error (MSE) and performing backpropagation. This error is expressed as $L_t(\theta_t)$ in
 484 Eq. (16):

$$485 \quad L_t(\theta_t) = \mathbb{E}_{s_{t'}, a_{t'}} \left[\left(y_{t'}^{DDQN} - Q(s_{t'}, a_{t'}; \theta_t) \right)^2 \right] \quad (16)$$

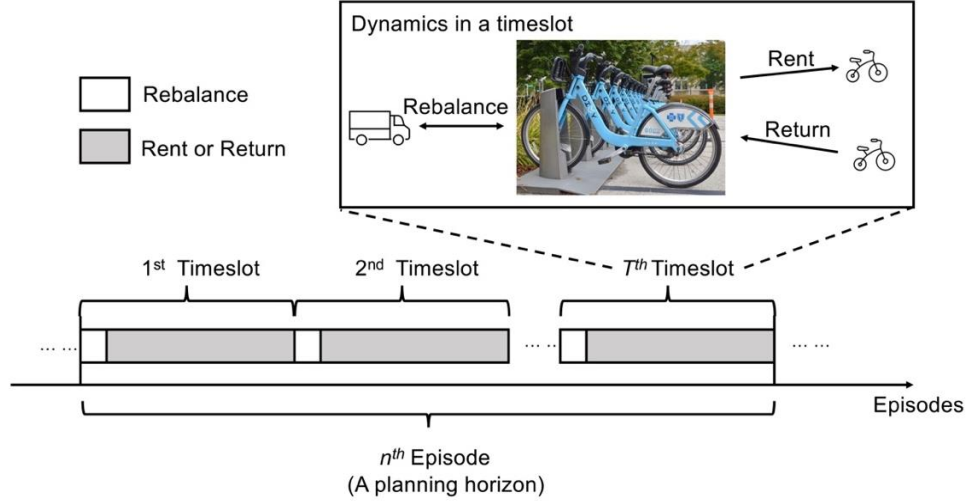
486 In this equation, $y_{t'}^{DDQN} = r_{t'} + \gamma Q(s_{t'+1}, \max_{a_{t'+1}} Q(s_{t'+1}, a_{t'+1}; \theta_t); \theta_t^-)$ represents the
 487 target Q-value for the selected action in the sampled timeslot t' while $Q(s_{t'}, a_{t'}; \theta_t)$ is the
 488 predicted Q-values for the selected action based on the state $s_{t'}$ by using primary network (Mnih
 489 et al., 2015). The target Q-value is calculated as the sum of the true reward ($r_{t'}$) received for the
 490 selected action $a_{t'}$ and the discounted Q-values in the subsequent state $s_{t'+1}$. Decoupling the
 491 action selection and estimation by involving two DQN networks has been proven to mitigate
 492 DQN's overestimation bias in large-scale sequential decision-making scenarios (Mnih et al.,
 493 2015). The target network's parameters θ^- stay fixed and are periodically updated from copying
 494 the primary network's parameters θ .

495 4.2.4. Bike share system simulator

496 We designed a BSS simulator with which the *DeepBike* agent could interact and collect
 497 transition experiences to improve its DQN. The simulator models the dynamics of the bike
 498 renting and return activities and bike rebalancing operation.

499 Fig. 4 illustrates the simulation process, adhering to constraints from Section 3. We set
500 every 24 hours as one episode (i.e., a simulation day/a planning horizon⁴) and split each episode
501 evenly into $|T| = 24$ timeslots with $\Delta t = 1$ hour. The simulator is re-initialized when a new
502 simulation day starts. Adopting the Divvy initialization from Sun (2021), the number of
503 (homogenous) rebalancing vehicles in the fleet was set to be $|M| = 16$ while each rebalancing
504 vehicle has the capacity $C_m^{vehicle} = 20$ and departs from the depot at the beginning of the
505 working shift. Furthermore, the inventory of both stations and rebalancing vehicles are set to be
506 half full to simplify each re-initialization. Following the sequential decision-making scheme, we
507 assumed that the simulator always performs the rebalancing decision first and then executes the
508 bike usage activities within a timeslot. By executing the decision $(z_{i,j,m,t}, y_{j,m,t})$, the rebalancing
509 vehicle m routes to the target station j by the shortest road network path and changes station's
510 inventory level by moving bikes in or out. The rebalancing vehicle remains at the station until
511 receiving the next timeslot's rebalancing decision. The execution time of rebalancing was
512 omitted in the simulation, considering the time interval Δt is generally long enough to perform
513 the rebalancing for workers, as adopted the same approach from Lowalekar et al. (2017) and
514 Schuijbroek et al. (2017). The bike usage activities (renting and return) were simulated based on
515 Divvy's historical data described in Section 4.1. The rent/return demands are regarded as
516 fulfilled only if the originally targeted stations have sufficient available bikes/docks. In the event
517 that customers fail to return bikes to the intended stations, it is assumed that they will find
518 alternate docks from the nearest feasible neighbor stations that are able to accommodate their
519 trips. If they fail to rent bikes at the intended stations, customers will exit the system unserved
520 (i.e. using other alternative modes instead). This ensures that the number of bikes is conserved in
521 the system (Li et al., 2018; Lowalekar et al., 2017).

⁴ Terminologies across domains often overlap in meaning and are interchangeably used in this study. "step" is more commonly used in RL whereas "timeslot" is used in the VRP. Similarly, "planning horizon" is a common term in VRP while "episode" is frequently used in RL.



522

523 Fig. 4. The bike share system simulator. Each episode consists of $|T|$ timeslot. In each timeslot,
 524 the rebalancing actions are executed first and then the customers rent or return the bikes.

525 *4.3. Baseline algorithms*

526 To evaluate the performance of the proposed *DeepBike* model, we compared it against
 527 four *DBSRP* baseline algorithms. The baselines are selected because they are representative
 528 algorithms derived from either MIP or heuristics discussed in Section 2.

- 529 • *No Rebalancing* -- The operator does not execute any rebalancing in the BSS, which
 530 establishes a base scenario of BSS operation without rebalancing the bike fleet.
- 531 • *Greedy Rebalancing* -- As outlined in (Duan et al., 2018; Li et al., 2018; C. Zhang et al.,
 532 2022), in this algorithm, for a more than half-full rebalancing vehicle, it always travels to
 533 the nearest less than half-full stations and offloads the most possible bikes there until full.
 534 For a less than half-full rebalancing vehicle, it always travels to the nearest more than
 535 half-full stations and collects as many bikes as this rebalancing vehicle could.
- 536 • *Lagrangian Dual Decomposition-based MIP (LDD-MIP)* -- (Ghosh et al., 2017)
 537 exploited the Lagrangian Dual Decomposition and abstraction mechanism to solve the
 538 large-scale MIP model formulated in Section 3. To align with the action space of the
 539 *DeepBike* model defined in Section 4.2.3 for a fair comparison, Eq. (5) in Section 3 was
 540 modified as Eq. (17) such that rebalancing vehicles were constrained to route to reachable
 541 stations:

542
$$\sum_{j \in B_i} z_{i,j,m,t} - \sum_{j \in B_i} z_{j,i,m,t-1} = 0 \quad \forall i \in N, m \in M, t \in T \quad (17)$$

543 where B_i is the set of reachable stations via a maximum of 10 horizontal or vertical
544 moves from the current station i by a rebalancing vehicle m in timeslot t .

- 545 • Greedy Online Anticipatory Heuristic (GOAH) -- (Lowalekar et al., 2017) provided a
546 heuristic approach to produce online solutions for *DBSRP* by computing policy for each
547 individual rebalancing vehicle. Based on looking ahead to three timeslots as well as
548 historically observed demands, this approach estimates the extra bike availability in
549 stations in the current timeslot t . Then, it computes the weight of each valid rebalancing
550 action based on the potential reduction in customer loss if executing this action. Finally,
551 each rebalancing vehicle selects and executes the action that has the maximum weight.

552 4.4. Evaluation metrics

553 We used the following three metrics to measure the performance of different rebalancing
554 methods. The metrics were defined in accordance with the terms in the objective function in Eq.
555 (1).

- 556 • Customer loss and customer loss reduction. *Customer loss* is the total number of
557 customers who failed to rent bikes at their initial departure stations or return bikes at the
558 original arrival stations. *Customer loss reduction* is defined as the delta between customer
559 loss with *No Rebalancing* and customer loss when using a rebalancing strategy.
- 560 • Vehicle routing distance (VRD). To account for the cost of conducting rebalancing, we use
561 VRD, which is the sum of the distance that a rebalancing vehicle travels among visited
562 stations in an episode. This is also known as vehicle-miles-traveled (VMT).
- 563 • Overall net profit. The net benefit of rebalancing is then measured by the overall net profit,
564 which is the sum of the increased ride revenue earned from *customer loss reduction* and
565 the cost (negative revenue) of operating the rebalancing vehicles. In this work, the revenue
566 from each ride is equivalent to the average price for a trip (\$3.3/trip) (Divvy, 2021). The
567 operation cost of the rebalancing vehicles is estimated by multiplying the mileage rate (58
568 cents/mile in 2019, (*Home: Internal Revenue Service, 2023*)) and the routing distance of
569 rebalancing vehicles (in miles).

570 5. Experiment results

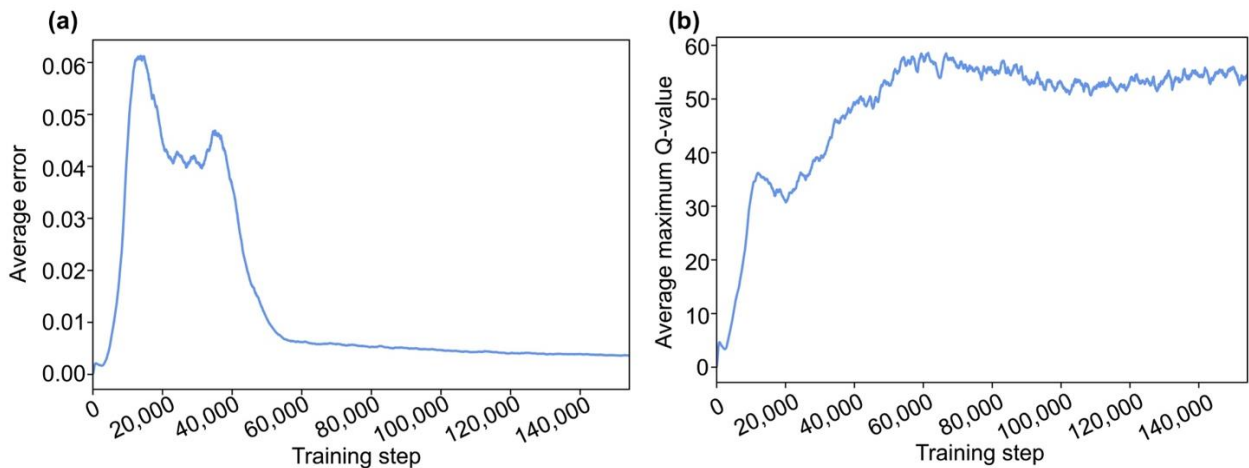
571 In this section, we first present the training results of our proposed *DeepBike* model
572 (Section 5.1). Then we show the empirical evaluation results of comparing the performance of

573 the *DeepBike* model with different rebalancing baselines, measured by the defined metrics
574 (Section 5.2), in order to examine the effectiveness of our proposed model for *DBSRP*.

575 During the training phase, the *DeepBike* model was trained on the Gilbreth community
576 cluster with 256 GB of RAM and two P100 Nvidia Tesla GPUs. The model underwent training
577 across 6,500 episodes, totaling 156,000 training steps. The DQN reached full convergence after
578 168 hours of physical training time. As described in Section 4.1, three consecutive weeks of data
579 served as training data and the subsequent one-week data was utilized to test the performance of
580 the *DeepBike* model. Furthermore, to alleviate the uncertainty of rebalancing models, we ran the
581 experiments ten times and calculated their average performance as well as the standard deviation.

582 5.1. Model training results

583 Fig. 5 presents the training curves of the DQN in the *DeepBike* framework. In Fig. 5(a),
584 the average error of DQN training at each step represents the mean MSE of DQN (as defined in
585 Eq. (16)) over the sampled batch (i.e., 128 experiences). Meanwhile, Fig. 5(b) displays the mean
586 highest Q-value over the batch for the selected action. Both the average loss and the average
587 maximum Q-value converged after around 60,000 steps. This observation implies that the agent
588 has learned a rebalancing policy and is able to well approximate the long-term benefits (i.e., Q-
589 values) of rebalancing actions.



590
591 Fig. 5. Training curves of the *DeepBike* model. (a) Average error over the training steps; (b)
592 Average maximum Q-value over the training steps.

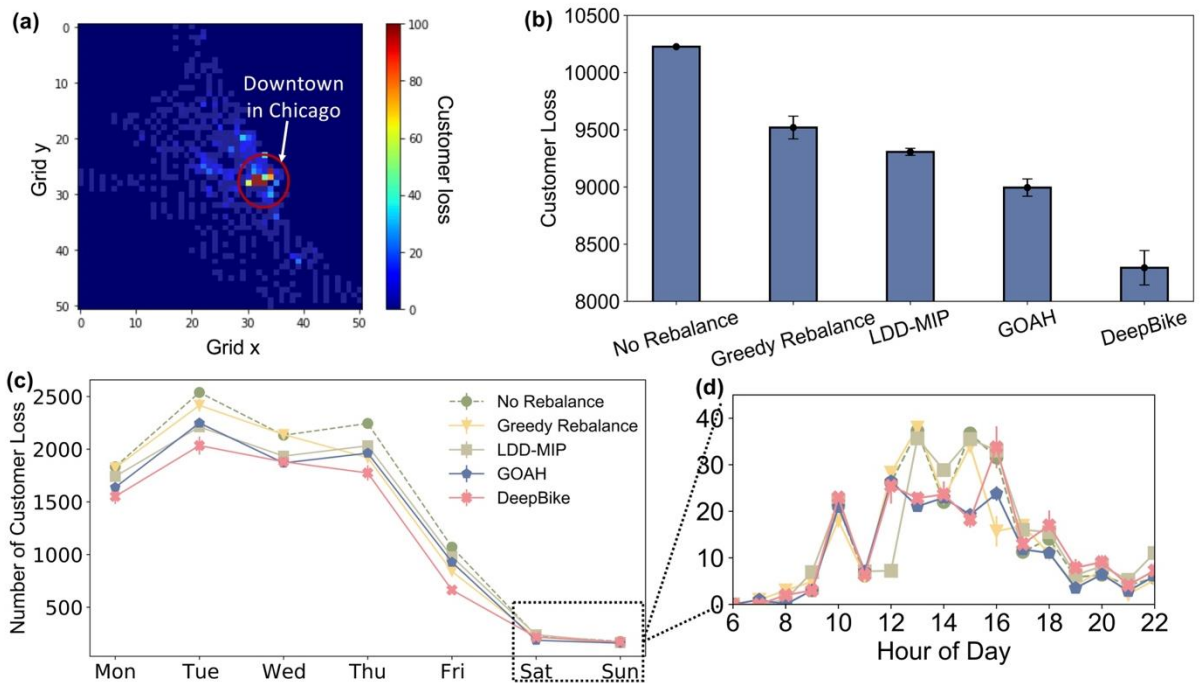
593

594

595 5.2. Performance results and analysis

596 5.2.1. Customer loss reduction

597 We first present the customer loss under different rebalancing strategies, reflecting how
598 the rebalancing strategy promotes the success of bike rentals and returns for potential customers.
599 Fig. 6(a) illustrates that without rebalancing, the majority of daily customer losses happened in
600 downtown Chicago, while notable customer losses can be observed in suburban regions as well.
601 This implies the potential locations for rebalancing vehicles to target and subsequently mitigate
602 customer loss.

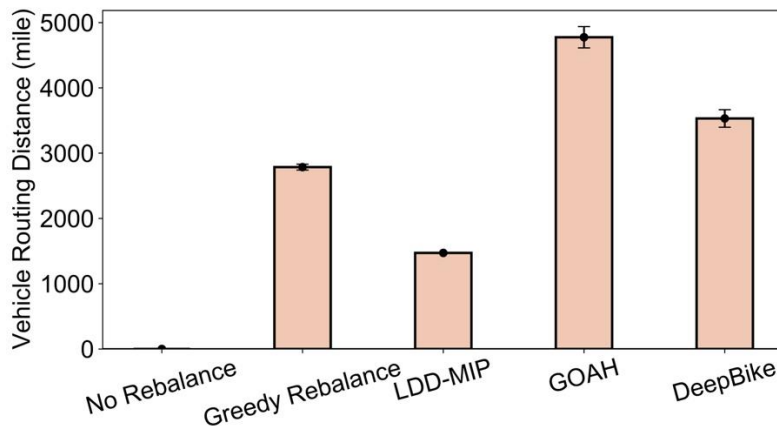


603
604 Fig. 6. Comparison of different rebalancing strategies in customer loss. (a) Average daily
605 customer loss distribution across grids in Chicago without rebalancing. (b) Comparison of
606 average customer loss in a week. (c) Comparison of average daily customer loss. (d) Comparison
607 of average hourly customer loss over Saturday and Sunday.

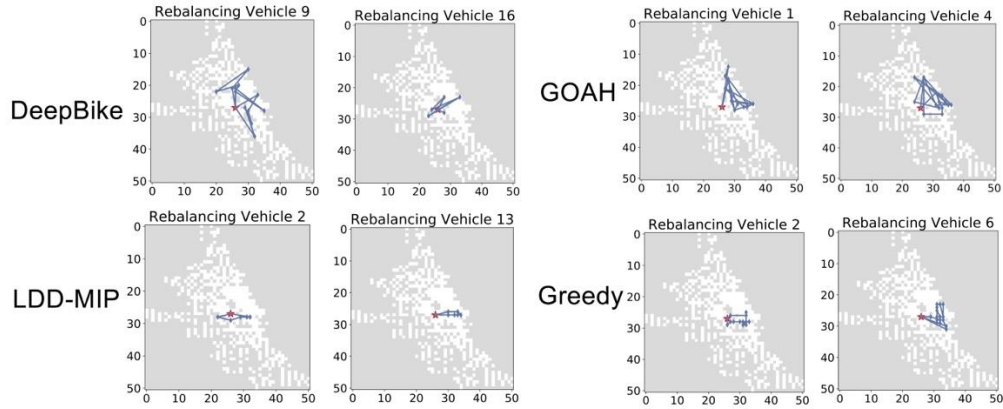
608 As demonstrated in Fig. 6(b), the *DeepBike* model minimizes total customer loss over the
609 one-week test horizon, outperforming all other strategies. Compared to the *No Rebalance*
610 strategy, *DeepBike* achieves an average reduction of 1,942 customer losses in a week. In
611 contrast, *GOAH*, *LDD-MIP*, and *Greedy Rebalance* yield average weekly reductions of 1,232,
612 920, and 707 customer losses, respectively. That is, *DeepBike* surpasses *GOAH*, *LDD-MIP*, and
613 *Greedy Rebalance* by reducing additional customer loss by 57.6%, 111.09%, and 174.7%,
614 respectively.

615 Fig. 6(c) shows that the *DeepBike* model consistently outperforms all other strategies in
 616 daily customer loss reduction throughout the weekdays. Both *DeepBike* and *GOAH*
 617 outperformed *LDD-MIP* in reducing customer loss due to their adaptability to BSS’s real-time
 618 customer demands, whereas *LDD-MIP* only produces offline rebalancing solutions. *Greedy*
 619 *Rebalance* underperforms as it neglects upcoming demands and lacks coordination among
 620 rebalancing vehicles. However, these rebalancing strategies demonstrate no significant
 621 disparities in reducing customer loss during weekends. This can be attributed to the
 622 comparatively lower customer demands on weekends as opposed to weekdays. A closer look at
 623 the average hourly customer loss distribution on the weekend, depicted in Fig. 6(d), reveals that
 624 each rebalancing strategy excels the others only at certain times of the day, exhibiting minor
 625 overall gaps in customer loss. The customer loss without rebalancing reaches a maximum of 40
 626 even during peak hours, implying adequate bike inventory at most stations and minimal customer
 627 loss due to low usage on weekends. As a result, all rebalancing strategies tend to yield similar
 628 rebalancing decisions on weekends, which leads to comparable customer loss reduction. *Vehicle*
 629 *routing distance*

630 Fig. 7 shows that, from the perspective of VRD, the *DeepBike* model’s rebalancing
 631 vehicles rank second in routing activity. Its vehicles travel to bike stations more frequently than
 632 those of the *Greedy Rebalance* and *LDD-MIP* but not as often as the *GOAH*. Additionally, we
 633 show the representative routing trajectories under different strategies in Fig. 8 (Complete
 634 trajectories for all rebalancing vehicles are available in the SI from S-2 to S-5).
 635



636
 637 Fig. 7. Comparison of cumulative routing distance (in miles) in a week of different rebalancing
 638 strategies.
 639



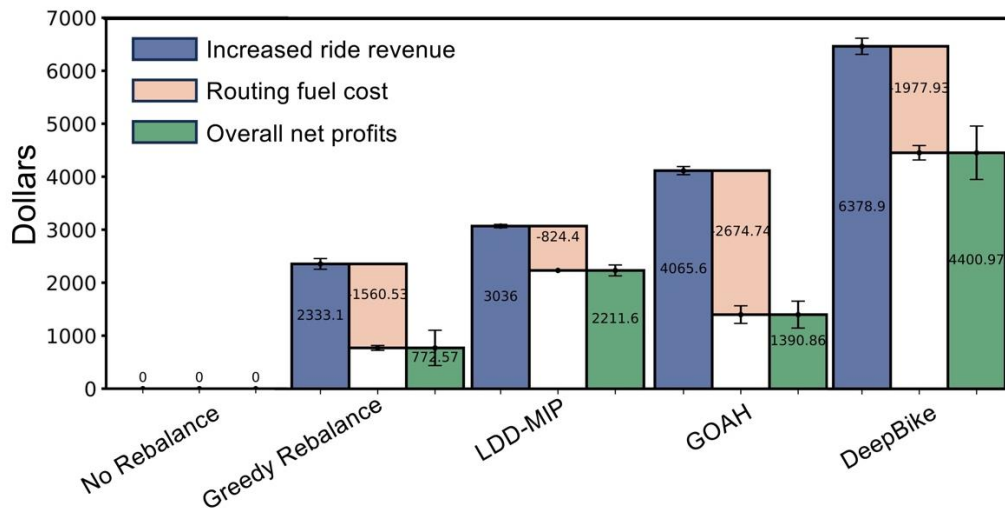
640
641 Fig. 8. Comparison of representative trajectories under different rebalancing strategies

642 The *DeepBike* model proactively dispatches the rebalancing vehicles not only to the
643 downtown area but also to suburban areas to overcome customer loss. Its VRD is 137.57% and
644 25.45% higher than that of *LDD-MIP* and *Greedy Rebalancing*, respectively. Despite that the
645 rebalancing vehicles are constrained to move at most ten grids either vertically or horizontally in
646 a single timeslot in Section 4.2.3.1, the *DeepBike* model has far sight to maximize this limit. For
647 instance, for rebalancing vehicle #15 in Fig. S-2, it opts for the maximum move in 61% (11 out
648 of 18) timeslots. This strategy enables the vehicles to reach suburban areas over multiple
649 timeslots. Additionally, *DeepBike* learns to coordinate the routings among rebalancing vehicles.
650 This is observed as rebalancing vehicle #16 conservatively circulates around the downtown area
651 contrasting with other rebalancing vehicles that are more aggressive in routing. It suggests that
652 the *DeepBike* model has a far-sighted strategy to balance the customer loss reduction and fuel
653 cost efficiency for the rebalancing fleet to reach long-term optimal rebalancing performance as
654 well as cover a wider range of bike stations in the service region. Contrarily, the routing
655 distances under the *LDD-MIP* strategy are the shortest over other strategies and rebalancing
656 mainly concentrates on the downtown area as indicated in Fig. 7 and Fig. 8, respectively. Even
657 though *GOAH* employs a more active rebalancing strategy as depicted in Fig. 8, it falls short of
658 achieving a higher reduction in customer loss, as discussed earlier. Such a result of the *GOAH*
659 echoes the limitation addressed in (Lowalekar et al., 2017) that their model did not specifically
660 consider the long-term fuel cost in evaluating the marginal benefits of the rebalancing actions in
661 the algorithm.

662

663 5.2.2. Overall net profit

664 Last, to comprehensively assess both the benefit and cost of rebalancing, we compare the
665 overall net profits achieved under different rebalancing strategies shown in Fig. 9. The results
666 show that the proposed *DeepBike* model attained the highest overall net profits compared with all
667 other evaluated rebalancing strategies. The overall net profits represent the tradeoff between
668 reducing customer loss and incurring fuel costs. Despite the longer routing distances from
669 aggressive rebalancing to service both downtown and suburban areas, *DeepBike's* strategy
670 markedly diminishes potential customer loss and maximizes profits. This implies that incurring
671 higher routing costs can be justified by the greater profit returns achieved through better
672 customer demand fulfillment.



673
674 Fig. 9. Comparison of overall net profits in a week with the breakdown of increased ride revenue
675 and routing fuel cost of different rebalancing strategies.

676 To sum up, our proposed *DeepBike* rebalancing strategy effectively addresses the large-
677 scale *DBSRP*. It generates specific rebalancing decisions for each individual rebalancing vehicle
678 based on real-time system states. Our experiments show that, by trading off between the routing
679 distance and customer loss reduction, rebalancing decisions generated by the *DeepBike* model
680 are beneficial for customers to improve the success of using bikes and for operators to boost the
681 overall net profits.

682 **6. Discussions and conclusion**

683 In this study, we proposed a deep reinforcement learning based model—*DeepBike*—for
684 large-scale *DBSRP*. In *DeepBike*, we designed a neural network (i.e., the DQN) to process the

685 real-time system states and estimate the long-term benefits (Q-values) of the rebalancing actions
686 for each rebalancing vehicle. *DeepBike* updates its DQN parameters through collected
687 experiences from the interaction with the BSS environment (i.e., learning from trial and error),
688 such that DQN can accurately approximate the Q-values. The numerical results presented in
689 Section 5 demonstrate that *DeepBike* mitigates the most customer loss and maximizes the overall
690 net profit over the one-week test horizon in a large-scale BSS with over 500 bike stations and 16
691 rebalancing vehicles. We also observed that the *DeepBike* model tends to balance routing fuel
692 cost against customer loss to achieve the highest overall net profits. The resulting customer loss
693 reduction and overall net profit improvement are beneficiaries for both customers and BSS
694 operators.

695 Our proposed model, *DeepBike*, provides an advancement over existing solutions to
696 *DBSRP*. The complexity of *DBSRP* escalates with the expansion of BSS, marked by more bike
697 stations and a larger fleet size of rebalancing vehicles. In online decision making, only historically
698 observed data is available. By employing a DRL scheme, *DeepBike* trains a neural network that,
699 once converged, can generate online rebalancing decisions for large-scale BSS, contrasting with
700 the MIP approaches that require the complete data for the planning horizon and solve the model
701 from scratch. Our model does not necessitate partitioning the service region; instead, it is able to
702 take the entire region as input for the DQN. Additionally, the output of the neural network
703 approximates the long-term benefits of the rebalancing actions till the end of the planning horizon,
704 rather than myopically making the decisions. The model creates individualized rebalancing actions
705 for each rebalancing vehicle, therefore it is able to adapt to changes in fleet size. It is also notable
706 that, even though our *DeepBike* model was experimented on a station-based BSS because of data
707 availability, our model is transferable to free-floating BSS because zoning the service region into
708 $n \times n$ grids and aggregating free-floating bikes parked within neighbor grids as abstract stations
709 are still applicable to align with the inputs of DQN.

710 Although the *DeepBike* model contributes to building a more efficient and profitable BSS
711 that benefits both customers and operators, it has the following limitations that future works
712 could address to further improve the model. First, the model relies on publicly available trip data
713 from BSS operators, representing observed demands. Such data underestimates station check-in
714 rates as it overlooks potential customers deterred by the unavailability of bikes at empty stations.
715 A true demand estimation model would better represent the latent customer demands and

716 therefore refine the simulator. Second, all rebalancing vehicles share the same DQN to generate
717 rebalancing decisions, leading to a uniform rebalancing policy across all vehicles. However,
718 adopting a multi-agent decision-making approach, where each rebalancing vehicle operates on
719 distinct policies, could further foster cooperation among vehicles, potentially better maximizing
720 overall revenue. Third, this work neglected the labor cost in the objective, presuming
721 unconditional completion of rebalancing assignments by workers alone. Building a model that
722 accounts for not only the revenues from customers but also the varying cost schemes of
723 conducting rebalancing tasks (such as hiring workers to fulfill all rebalancing needs versus
724 incentivizing customers to ride bikes to a designated location) could lead to a more operationally
725 realistic and economical representation of the BSS.

726 Appendix

727 Table A-1. List of notations.

Decision Variables	
$z_{i,j,m,t}$	Binary, the routing decision variable that is set to 1 for vehicle m in the timeslot t if traveling from station i to station j
$y_{j,m,t}$	Integer, the repositioning decision variable for vehicle m . The positive/negative denotes the number of bikes to drop off/pick up when visiting target station j in timeslot t
Intermediate Variables	
$S_{i,t}^{Bike}$	The number of bikes in the station i at the beginning of the timeslot t
$S_{i,t}^{Dock}$	The number of docks in the station i at the beginning of timeslot t
$V_{m,t}^{Bike}$	The number of bikes in the rebalancing vehicle m at the beginning of the timeslot t
$V_{m,t}^{Slot}$	The number of slots in the rebalancing vehicle m at the beginning of the timeslot t
$V_{m,t}^{Loc}$	The location of vehicle m at the beginning of timeslot t
$D_{i,t}$	The net demands of the station i in the timeslot t . The net demand equals to the value that the return demands minus the rent demands
$D_{i,t:t+n}^{pred}$	The predicted future net demands of the station i in the next n timeslots from the beginning of the timeslot t
$CL_{i,t}$	The customer loss occurred in the station i in timeslot t .
Deep Reinforcement Learning Components	
S_t	BSS environment state in timeslot t
r_t	The reward earned in timeslot t
e_t	A sample consists of S_t, a_t, r_t and S_{t+1} , which records the transition from the timeslot t to the timeslot $t + 1$
a_t	Rebalancing action in timeslot t , a tuple consists of the routing decision and repositioning decision $a_t = (z_{i,j,m,t}, y_{j,m,t})$
Sets and Constants	
T	The set of discrete timeslots, $t \in \{1, 2, \dots, T \}$
M	The set of rebalancing vehicles in the fleet, $m \in \{1, 2, \dots, M \}$
N	The set of aggregated bike stations in BSS, $i, j \in \{1, 2, \dots, N \}$

$C_m^{vehicle}$	The capacity of vehicle m – the maximum number of bikes that can be carried in the vehicle $m, C_m^{vehicle} = v_{m,t}^{Bike} + v_{m,t}^{Slot}$
α	The average price of a single trip riding the shared bike (dollar/trip)
β	The mileage rate of rebalancing vehicles (dollar/mile)
Δt	The duration of each timeslot
ϵ	The probability of selecting a random action during the training process of <i>DeepBike</i>
γ	Time discount factor
ρ	The probability of doing no rebalancing action during the initial training process

728

729

730

References

- 731 Brinkmann, J., Ulmer, M. W., & Mattfeld, D. C. (2020). The multi-vehicle stochastic-dynamic
732 inventory routing problem for bike sharing systems. *Business Research*, 13(1), 69–92.
733 <https://doi.org/10.1007/s40685-019-0100-z>
- 734 Chen, D., & Sakai, K. (2022). A User-Based Bike Return Algorithm for Docked Bike Sharing
735 Systems. *Workshop Proceedings of the 51st International Conference on Parallel*
736 *Processing*, 1–8. <https://doi.org/10.1145/3547276.3548443>
- 737 Chen, J., Yang, Z., Shu, Y., & Cheng, P. (2020). Rebalance Bike-Sharing System With Deep
738 Sequential Learning. *IEEE Intelligent Transportation Systems Magazine*, January 2020,
739 2–8. <https://doi.org/10.1109/MITS.2019.2926252>
- 740 Chiariotti, F., Pielli, C., Zanella, A., & Zorzi, M. (2018). A dynamic approach to rebalancing
741 bike-sharing systems. *Sensors (Switzerland)*, 18(2), 1–22.
742 <https://doi.org/10.3390/s18020512>
- 743 Citi Bike. (2023). *January 2023 Monthly Report*. [https://citibikenyc.com/system-data/operating-](https://citibikenyc.com/system-data/operating-reports)
744 [reports](https://citibikenyc.com/system-data/operating-reports)
- 745 Contardo, C., Morency, C., & Rousseau, L.-M. (2012). Balancing a dynamic public bike-sharing
746 system. *Cirrelet*.
- 747 de Bruin, T., Kober, J., Tuyls, K., & Babuška, R. (2018). Integrating State Representation
748 Learning Into Deep Reinforcement Learning. *IEEE Robotics and Automation Letters*,
749 3(3), 1394–1401. <https://doi.org/10.1109/LRA.2018.2800101>
- 750 Divvy. (2021). *Single Ride*. <https://www.divvybikes.com/pricing/single-ride>
- 751 *Divvy System Data*. (2023, May 20). Divvy Bikes. <https://ride.divvybikes.com/system-data>
- 752 Duan, Y., Wu, J., & Zheng, H. (2018). A Greedy Approach for Vehicle Routing When
753 Rebalancing Bike Sharing Systems. *2018 IEEE Global Communications Conference*
754 *(GLOBECOM)*, 1–7. <https://doi.org/10.1109/GLOCOM.2018.8647755>
- 755 Fishman, E. (2014). *Bikeshare: Barriers, facilitators and impacts on car use*. Queensland
756 University of Technology.
- 757 Fishman, E. (2016). Bikeshare: A review of recent literature. *Transport Reviews*, 36(1), 92–113.
- 758 *General Bikeshare Feed Specification*. (2023, July 22). GitHub.
759 <https://github.com/MobilityData/gbfs>

760 Ghosh, S., Trick, M., & Varakantham, P. (2016). *Robust repositioning to counter unpredictable*
761 *demand in bike sharing systems.*

762 Ghosh, S., Varakantham, P., Adulyasak, Y., & Jaillet, P. (2015). Dynamic redeployment to
763 counter congestion or starvation in vehicle sharing systems. *Proceedings of the 8th*
764 *Annual Symposium on Combinatorial Search, SoCS 2015, 2015-Janua*, 230–231.

765 Ghosh, S., Varakantham, P., Adulyasak, Y., & Jaillet, P. (2017). Dynamic repositioning to
766 reduce lost demand in bike sharing systems. *Journal of Artificial Intelligence Research*,
767 *58*, 387–430. <https://doi.org/10.1613/jair.5308>

768 Gleditsch, M. D., Hagen, K., Andersson, H., Bakker, S. J., & Fagerholt, K. (2022). A column
769 generation heuristic for the dynamic bicycle rebalancing problem. *European Journal of*
770 *Operational Research*. <https://doi.org/10.1016/j.ejor.2022.07.004>

771 *Home: Internal Revenue Service*. (2023, May 20). Internal Revenue Service | An Official
772 Website of the United States Government. <https://www.irs.gov/>

773 Hu, R., Zhang, Z., Ma, X., & Jin, Y. (2021). Dynamic Rebalancing Optimization for Bike-
774 Sharing System Using Priority-Based MOEA/D Algorithm. *IEEE Access*, *9*, 27067–
775 27084. <https://doi.org/10.1109/ACCESS.2021.3058013>

776 Kloimüller, C., Papazek, P., Hu, B., & Raidl, G. R. (2014). Balancing bicycle sharing systems:
777 An approach for the dynamic case. *Lecture Notes in Computer Science (Including*
778 *Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*,
779 *8600*(June), 73–84. https://doi.org/10.1007/978-3-662-44320-0_7

780 Li, Y., Zheng, Y., & Yang, Q. (2018). Dynamic bike reposition: A spatio-temporal reinforcement
781 learning approach. *Proceedings of the ACM SIGKDD International Conference on*
782 *Knowledge Discovery and Data Mining*, 1724–1733.
783 <https://doi.org/10.1145/3219819.3220110>

784 Lin, L.-J. (1992). *Reinforcement learning for robots using neural networks*. Carnegie Mellon
785 University.

786 Lowalekar, M., Varakantham, P., Ghosh, S., Jena, S. D., & Jaillet, P. (2017). Online
787 repositioning in bike sharing systems. *Proceedings International Conference on*
788 *Automated Planning and Scheduling, ICAPS*, 200–208.

789 Luo, H., Zhao, F., Chen, W. Q., & Cai, H. (2020). Optimizing bike sharing systems from the life
790 cycle greenhouse gas emissions perspective. *Transportation Research Part C: Emerging*
791 *Technologies*, *117*(September 2019), 102705. <https://doi.org/10.1016/j.trc.2020.102705>

792 Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller,
793 M. (2013). *Playing Atari with Deep Reinforcement Learning*. 1–9.

794 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A.,
795 Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A.,
796 Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., & Hassabis, D. (2015).
797 Human-level control through deep reinforcement learning. *Nature*, *518*(7540), 529–533.
798 <https://doi.org/10.1038/nature14236>

799 Oda, T., & Joe-wong, C. (2018). MOVI: A Model-Free Approach to Dynamic Fleet
800 Management. *IEEE INFOCOM 2018 - IEEE Conference on Computer Communications*,
801 2708–2716.

802 O’Mahony, E., & Shmoys, D. B. (2015). Data analysis and optimization for (Citi)bike sharing.
803 *Proceedings of the National Conference on Artificial Intelligence*, 1, 687–694.

804 Osorio, J., Lei, C., & Ouyang, Y. (2021). Optimal rebalancing and on-board charging of shared
805 electric scooters. *Transportation Research Part B: Methodological*, 147, 197–219.
806 <https://doi.org/10.1016/j.trb.2021.03.009>

807 Qin, Z., Tang, J., & Ye, J. (2019). Deep reinforcement learning with applications in
808 transportation. *Proceedings of the 25th ACM SIGKDD International Conference on*
809 *Knowledge Discovery & Data Mining*, 3201–3202.

810 Regue, R., & Recker, W. (2014). Proactive vehicle routing with inferred demand to solve the
811 bikesharing rebalancing problem. *Transportation Research Part E: Logistics and*
812 *Transportation Review*, 72, 192–209. <https://doi.org/10.1016/j.tre.2014.10.005>

813 Schuijbroek, J., Hampshire, R. C., & van Hoes, W. J. (2017). Inventory rebalancing and vehicle
814 routing in bike sharing systems. *European Journal of Operational Research*, 257(3),
815 992–1004. <https://doi.org/10.1016/j.ejor.2016.08.029>

816 Shaheen, S., Guzman, S., & Zhang, H. (2010). Bikesharing in Europe, the Americas, and Asia.
817 *Transportation Research Record*, 2143, 159–167. <https://doi.org/10.3141/2143-20>

818 Shu, J., Chou, M. C., Liu, Q., Teo, C. P., & Wang, I. L. (2013). Models for effective deployment
819 and redistribution of bicycles within public bicycle-sharing systems. *Operations*
820 *Research*, 61(6), 1346–1359. <https://doi.org/10.1287/opre.2013.1215>

821 Shui, C. S., & Szeto, W. Y. (2018). Dynamic green bike repositioning problem – A hybrid
822 rolling horizon artificial bee colony algorithm approach. *Transportation Research Part*
823 *D: Transport and Environment*, 60, 119–136. <https://doi.org/10.1016/j.trd.2017.06.023>

824 Shui, C. S., & Szeto, W. Y. (2020). A review of bicycle-sharing service planning problems.
825 *Transportation Research Part C: Emerging Technologies*, 117(April 2019), 102648.
826 <https://doi.org/10.1016/j.trc.2020.102648>

827 Silver, D. (2014). Lecture 9: Exploration and Exploitation. *Computer Science Department*,
828 *University of London*.

829 Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser,
830 J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J.,
831 Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., &
832 Hassabis, D. (2016). Mastering the game of Go with deep neural networks and tree
833 search. *Nature*, 529(7587), 484–489. <https://doi.org/10.1038/nature16961>

834 Sun, R. (2021). *Bike Share System-Rebalancing Estimation and System Optimization* [PhD
835 Thesis]. Purdue University Graduate School.

836 Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.

837 Vallez, C. M., Castro, M., & Contreras, D. (2021). Challenges and opportunities in dock-based
838 bike-sharing rebalancing: A systematic review. *Sustainability*, 13(4), 1829.

839 Van Hasselt, H., Guez, A., & Silver, D. (2016). Deep reinforcement learning with double Q-
840 Learning. *30th AAAI Conference on Artificial Intelligence, AAAI 2016*, 2094–2100.
841 <https://doi.org/10.1609/aaai.v30i1.10295>

842 Wang, J., & Sun, L. (2020). Dynamic holding control to avoid bus bunching: A multi-agent deep
843 reinforcement learning framework. *Transportation Research Part C: Emerging*
844 *Technologies*, 116(April), 102661. <https://doi.org/10.1016/j.trc.2020.102661>

845 Wang, T. (2014). *Solving Dynamic Repositioning Problem for Bicycle Sharing Systems: Model,*
846 *Heuristics, and Decomposition.*

847 Yin, Z., Hardaway, K., Feng, Y., Kou, Z., & Cai, H. (2023). Understanding the demand
848 predictability of bike share systems: A station-level analysis. *Frontiers of Engineering*
849 *Management*, 10(4), 551–565. <https://doi.org/10.1007/s42524-023-0279-8>

850 Yin, Z., Kou, Z., & Cai, H. (2023). A Deep Reinforcement Learning Model for Large-Scale
851 Dynamic Bike Share Rebalancing with Spatial-Temporal Context. *The 12th International*
852 *Workshop on Urban Computing*. [http://urban-](http://urban-computing.com/urbcomp2023/file/UrbComp2023_paper_7.pdf)
853 [computing.com/urbcomp2023/file/UrbComp2023_paper_7.pdf](http://urban-computing.com/urbcomp2023/file/UrbComp2023_paper_7.pdf)

854 Zhang, C., Wu, F., Wang, H., Tang, B., Fan, W., & Liu, Y. (2022). A Meta-Learning Algorithm
855 for Rebalancing the Bike-Sharing System in IoT Smart City. *IEEE Internet of Things*
856 *Journal*, 9(21), 21073–21085. <https://doi.org/10.1109/JIOT.2022.3176145>

857 Zhang, D., Yu, C., Desai, J., Lau, H. Y. K., & Srivathsan, S. (2017). A time-space network flow
858 approach to dynamic repositioning in bicycle sharing systems. *Transportation Research*
859 *Part B: Methodological*, 103, 188–207. <https://doi.org/10.1016/j.trb.2016.12.006>

860 Zhao, J., Mao, M., & Zhao, X. (2021). A hybrid of deep reinforcement learning and local search
861 for the vehicle routing problems. *Ieeexplore.Ieee.Org*.
862 [https://ieeexplore.ieee.org/abstract/document/9141401/?casa_token=joKmbnG-](https://ieeexplore.ieee.org/abstract/document/9141401/?casa_token=joKmbnG-H1kAAAAA:XSSUCwt_NeOCg0twxOJeS3yl3uwiz1o5te-NkOxMisQunkE8_zq3ordTkIzNona4NOzpCrE12kA)
863 [H1kAAAAA:XSSUCwt_NeOCg0twxOJeS3yl3uwiz1o5te-](https://ieeexplore.ieee.org/abstract/document/9141401/?casa_token=joKmbnG-H1kAAAAA:XSSUCwt_NeOCg0twxOJeS3yl3uwiz1o5te-NkOxMisQunkE8_zq3ordTkIzNona4NOzpCrE12kA)
864 [NkOxMisQunkE8_zq3ordTkIzNona4NOzpCrE12kA](https://ieeexplore.ieee.org/abstract/document/9141401/?casa_token=joKmbnG-H1kAAAAA:XSSUCwt_NeOCg0twxOJeS3yl3uwiz1o5te-NkOxMisQunkE8_zq3ordTkIzNona4NOzpCrE12kA)

865 Zheng, X., Tang, M., Liu, Y., Xian, Z., & Zhuo, H. H. (2021). Repositioning bikes with carrier
866 vehicles and bike trailers in bike sharing systems. *Applied Sciences (Switzerland)*, 11(16).
867 <https://doi.org/10.3390/app11167227>

868 Zhou, X. (2015). Understanding spatiotemporal patterns of biking behavior by analyzing massive
869 bike sharing data in Chicago. *PLoS ONE*, 10(10), 1–20.
870 <https://doi.org/10.1371/journal.pone.0137922>

871 Zhu, H., Shou, T., Guo, R., Jiang, Z., Wang, Z., Wang, Z., Yu, Z., Zhang, W., Wang, C., &
872 Chen, L. (2022). Redpacketbike: A graph-based demand modeling and crowd-driven
873 station rebalancing framework for bike sharing systems. *IEEE Transactions on Mobile*
874 *Computing*. <https://ieeexplore.ieee.org/abstract/document/9693278/>

875

876

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [DeepBikeSIfinal.docx](#)