# A semi-Markov decision model for the optimal control of an emergency medical service system

## Giannis A. Kechagias and Alexandros C. Diamantidis\*

School of Economics, Faculty of Economics and Political Sciences, Aristotle University of Thessaloniki, Greece Email: ioannika@econ.auth.gr Email: adiama@econ.auth.gr \*Corresponding author

### Theodosis D. Dimitrakos

School of Sciences, Department of Mathematics, University of the Aegean, Karlovassi, Samos, Greece Email: dimitheo@aegean.gr

Abstract: A mathematical model for the analysis of an emergency medical service (EMS) system with a specific number of advanced life support units (ALS) and a specific number of basic life support (BLS) units is presented in this paper. The system admits incoming emergency calls which are divided into two classes: 1) urgent, high-priority calls for which the patient's life is potentially at risk; 2) less urgent low-priority calls. Under a suitable cost structure, the system is modelled using an appropriate Markov decision process in continuous time for which we seek a stationary policy that minimises a predefined optimality criterion for vehicle mixes over a set of candidate ambulance fleets. Based on this formulation, it is possible to implement standard Markov decision algorithms, such as the standard value-iteration algorithm and the standard policy-iteration algorithm. A sensitivity analysis of some model parameters is provided to examine their effect in the vehicle mix and in the cost of the system. An integer programming formulation is also provided for the corresponding location-allocation problem of the model. Numerical results are also presented for the examined problem.

**Keywords:** emergency medical service system; EMS; vehicle mix; Markov decision process; MDP; integer programming.

**Reference** to this paper should be made as follows: Kechagias, G.A., Diamantidis, A.C. and Dimitrakos, T.D. (2024) 'A semi-Markov decision model for the optimal control of an emergency medical service system', *Int. J. Industrial and Systems Engineering*, Vol. 46, No. 2, pp.169–194.

**Biographical notes:** Giannis A. Kechagias is a PhD candidate at the School of Economics of Aristotle University of Thessaloniki. He is a graduate of the Department of Mathematics of the University of Ioannina and the Department of Economics of the Aristotle University of Thessaloniki. He holds an MSc in

#### 170 G.A. Kechagias et al.

Theoretical Mathematics from the Mathematics Department of the Aristotle University of Thessaloniki. His research activity focuses on the design, analysis and optimisation of industrial systems, and on stochastic models in operational research with an emphasis on Markov decision-making processes and queues.

Alexandros C. Diamantidis is an Assistant Professor in Quantitative Methods in Management at the School of Economics of Aristotle University of Thessaloniki. He holds a PhD in the Analysis of Manufacturing Systems with Linear and Nonlinear Flow of Material from the University of the Aegean, MSc in the analysis of Job Shop Models from Dublin City University and BSc in Mathematics from the University of the Aegean. His research interests include the analysis and optimisation of manufacturing systems, supply chains and service management systems. He is a referee to many leading international scientific journals.

Theodosis D. Dimitrakos is an Assistant Professor in Applied Probabilities and in Stochastic Operations Research at the School of Sciences of University of the Aegean. He holds a PhD in the Optimal Control of Stochastic Processes from the University of the Aegean, MSc in Statistics from Athens University of Economics and Business and BSc in Mathematics from the University of the Aegean. His research interests include stochastic models in operations research, single vehicle routing problems and stochastic dynamic programming. He is a referee to many leading international scientific journals such as *European Journal of Operational Research, Operational Research: An International Journal* and *Communications in Statistics-Theory and Methods*.

This paper is a revised and expanded version of a paper entitled 'A vehicle mix Markov decision model for an emergency medical service system' presented at 6th Stochastic Modeling Techniques and Data Analysis International Conference, Barcelona, 2–5 June 2020.

#### **1** Introduction and literature review

Emergency medical service (EMS) system, most commonly known as EMS, is a system that provides emergency medical pre-hospital care. It is activated by an incident that causes serious illness or injury. The primary goal of EMS is to provide patients with emergency medical care and to transfer them in the hospital. The process of an EMS system is depicted in Figure 1.



Figure 1	The proces	s of an	EMS	system
				- /

There are two general types of EMS models used in many countries; the Anglo-American and the Franco-German system (see Dick, 2003). The Franco-German system is a Physician-EMS-based model that enables a doctor and an EMS to evaluate and treat a patient on the scene of a medical emergency. The patient can be taken to a hospital or clinic if further evaluation is required. It is widely adopted in European countries such as Germany, France, Greece, Malta and Austria (see Al-Shaqsi, 2010). The Anglo-American model, on the other hand, consists of ambulances staffed with emergency medical technicians (EMTs) and paramedics trained in basic, intermediate and advanced life support (ALS). It could be found in the USA, Canada, New Zealand, Oman and Australia (see Al-Shaqsi, 2010).

The basic component in these systems is the ambulances, which are responsible for providing pre-hospital care. There are two basic types of ambulances: the ALS ambulances and the basic life support (BLS) ambulances. BLS is staffed by EMTs and it is associated with the 'load and go' philosophy providing non-invasive basic interventions and rapid transport to definitive health care facility. ALS has a paramedic on board (some countries require a higher level of care and they employ a physician among the staff) along with a EMTs and it fits more with the 'stay and stabilise' approach. It includes all the BLS procedures with the addition of invasive procedures such as intravenous line placement, fluid replacement, needle-chest decompression and others. Several reports have been published by researchers comparing the effectiveness of ALS and BLS in many medical situations. Indicatively, we mention Nguyen-Van-Tam et al. (1997), Rainer et al. (1997), Eisen and Dubinsky (1998), Stiell et al. (2004, 2007) and Di Bartolomeo et al. (2005).

According to the type of fleet and the chosen dispatch, the EMS systems can be grouped into two categories:

- a the All-ALS systems in which an ALS ambulance is always dispatched
- b the tiered or mixed systems.

The tiered system uses a combination of ALS and BLS units and the dispatch depends on the severity of the emergency call. For each system type there are advantages and disadvantages. Thus, in all-ALS system the triage of a call is not complicating and that leads to decreased response times. Additionally, there is always an advanced practitioner on scene (paramedic). A disadvantage is that the paramedic loses their skills because they serve more frequent less urgent calls. As for tiered systems, the majority of emergency calls do not need an ALS ambulance, so the mixed system has the advantage of freeing up ALS ambulances for the acute care of seriously ill patients. Also, ALS ambulances are more expensive to operate, so mixed fleets can be larger and as consequence the system has shorter response time. As a disadvantage we can refer that there is a risk of sending a BLS ambulance to a call requiring paramedic support. Consequently, there is a reasonable dilemma if the EMS fleet should have only ALS ambulances or a combination of ALS and BLS ambulances. The selection of a combination of ALS and BLS ambulances in an EMS system is known as the vehicle mix problem. There are many studied results by advocators of each system type over time. We can refer as advocators of all-ALS systems, (Ornato et al., 1990; Wilson et al., 1992) and as advocators of tiered systems (Braun et al., 1990; Clawson, 1989; Slovis et al., 1985; Stout et al., 2000).

The choice of vehicle mix is crucial for an EMS system due to the effect both in level of cost and in level of service. In this paper, an EMS system is modelled using an appropriate Markov decision process (MDP) in continuous time for which we can decide about the vehicle mix of the system over a set of candidate ambulance fleets. The use of Markov models in the analysis of EMS systems problems and of systems with similarly characteristics (such as fire departments) is a usual practice. Jarvis (1975) introduced a MDP for determining optimal dispatching policies for a single type of server. Berman (1981a, 1981b, 1981c) and Zhang (2012) formulated the ambulance redeployment problem as a MDP and then solved for an optimal policy using exact dynamic programming. Swersey (1982) developed a Markov model for determining how many fire vehicles to dispatch to a call that balances the costs associated with dispatching too few or too many. McLay and Mayorga (2012) presented a MDP to dispatch distinguishable EMS vehicles to prioritised calls and considered the fact that errors in the classification of patient priorities might occur. In this model, the dispatching policy that maximises the expected coverage of true high-risk calls was determined. Alanis et al. (2013) represented the EMS system as a two-dimensional Markov chain model to evaluate the system given a compliance table. Ramirez-Nafarrate et al. (2014) studied optimal ambulance diversion control policies using a MDP formulation that minimises the average time that patients wait beyond their recommended safety time threshold. Chong et al. (2015) constructed Markov decision models that allow for quantitative comparisons between feasible vehicle mixes. They studied the problem of dispatching into a mixed fleet system when there are a number of busy ALS and BLS ambulances. Our approach is closely related to this article. Lee and Lee (2018) formulated a finite horizon MDP model for the study of the admission control problem for patients arriving at an emergency department in the aftermath of a mass casualty incident. They took into consideration a policy restriction that immediate-patients should be admitted as long as there are available beds. Nasrollahzadeh et al. (2018) developed an optimisation framework for real-time ambulance dispatching and relocation. They formulated the problem as an infinite-horizon MDP and implemented that framework on an EMS system in Mecklenburg County, NC. Park and Lee (2019) proposed a two-tiered ambulance system, consisting of advanced and BLS units for emergency and non-emergency patient care, providing a cost-efficient medical service. They formulated their dynamic decision-making problem as a semi-MDP and proposed a mini-batch monotone-approximate dynamic programming (ADP) algorithm to solve the problem within a reasonable computation time. Most recently, Mengyu et al. (2020) formulated a discrete time, infinite time horizon, discounted MDP model to determine when it is advantageous to send appropriate patients to out-of-region Emergency departments, which have longer transport times but shorter offload times. DuBois and Albert (2021) studied how to optimally dispatch ambulances to prioritised patients during mass casualty incidents. They formulated the ambulance dispatching problem as a MDP model with patients prioritised by the benefit, they will receive from ambulance care and with two classes of ambulances.

A rich literature also exists on ambulance location-allocation models for EMS systems. We can mention some of the most recent reports on this issue. Pouraliakbari et al. (2018) formulated and solved a probabilistic maximal covering location model for determining the optimal location of facilities in congested EMS systems with referral hierarchical structure. Their goal was to minimise the total amount of demand that is lost in the system. To solve the model, two meta-heuristic algorithms, including population-based simulated annealing (PBSA) and ant colony optimisation (ACO) have been executed. Benabdouallah and Bojji (2018) presented a review on coverage models

applied to emergency location. They analysed these models classified under three classes. The static class focuses on the earlier models about emergency location coverage, the probabilistic class defines the ambulance unavailability ratio, and the dynamic class describes how to reassign dynamically the ambulances. Van Den Berg and Van Essen (2019) compared several ambulance locations models on coverage and response time criteria. They showed that the maximum expected covering location problem (MEXCLP) and the expected response time model (ERTM) perform the best overall considered criteria. Ji et al. (2020) proposed a data-driven real-time ambulance redeployment approach that redeploys an ambulance to a proper station after it becomes available, to optimise the transporting capability of an EMS system. In this paper, we give the IP formulation for the geographical distribution of the system's fleet at the stations of a service area. We consider heterogeneous fleet (ambulances are of different types) and different types of calls.

The contribution of this paper in the investigation of EMS problems is the following. An EMS system is modelled using an appropriate MDP in continuous time for which we seek a stationary policy that minimises a predefined optimality criterion for vehicle mixes over a set of candidate ambulance fleets. The point of view for our choices is associated with the number of emergency calls being in the system waiting to be served. We can claim that our model is an 'open system' in the sense that the EMS can redirect calls to other EMS systems when the system is in red alert. For our model we make the following assumption: the BLS units cannot serve adequately the high-priority calls and it is undesirable that the ALS units serve low-priority calls due to the highest service cost. We treat two models: a simplified version for the above description and a model which is closer to reality with two classes of emergency calls.

The rest of the paper is organised as follows. In Section 2, the simplified model of an EMS system and its semi-MDP formulation is described. In Section 3, the generalised model is presented with two types of ambulances and two classes of emergency calls. In Section 4, the semi-MDP formulation of the generalised model is presented. In Section 5, we provide numerical results which are based on information of Virginia Beach EMS system (VBEMS). In Section 6, we present a sensitivity analysis of some model parameters to see the effect of those in the vehicle mix and in the cost of the system. Finally, in Section 7, a possible integer programming formulation of the system is given, as a further research direction, for a suitable location-allocation problem of the model. The conclusions of the paper are also provided.

#### 2 A simplified version of the model and its semi-MDP formulation

We consider an EMS system which consists of N life support units. In fact, during each shift only a number of these vehicles are available for service. Suppose that k < N ambulances are available for service, we assume that at any given time some of the k-vehicles either remain at the base for refuelling and maintenance or they are a safety fleet for propulsion in some extreme unforeseen situations. The EMS admits incoming emergency calls. Let i be the number of calls waiting to be served at any given time, and t < k, the vehicles available for these calls (at ambulance stations remain k - t vehicles refuelling or preparedness in extreme cases). We allow the system to choose not to temporarily serve calls, which during the evaluation it seems to be unreliable or which do

not need immediate medical care, with the aim of the subsequent utilisation of resources in more urgent situations. This logic makes sense during rush hour when the system should be able to handle calls that require immediate emergency response.

Consider a Markov process in continuous-time  $\{(X(t), t \ge 0)\}$ , where the random variable X(t) denotes the state of the process at time  $t \ge 0$ . More precisely, the random variable X(t) denotes the number of emergency calls, which are in the system waiting to be served at time  $t \ge 0$ . Thus, if X(t) = i, the system state can be described by the value *i* of emergency calls, being in the system waiting to be served at time  $t \ge 0$ . We assume that the system can admit a maximum number Q of emergency calls. We further assume that the maximum number of emergency calls Q is greater than the maximum number N of the available ambulance units, i.e., Q > N. Thus, the state space of the system is defined by the set:

$$S = \{i \mid 0 \le i \le Q\}.$$

The decision epochs include the epochs at which an emergency call arrives to the system waiting to be served by the system.

For states  $i \in S$ , such that  $0 \le i \le N$ , actions 0 and 1 are possible. It is reasonable to assume that in state 0 the only possible action is action 0. According to action 0, emergency calls arrive to the system according to a Poisson process with rate equal to  $\lambda$ . The system makes the following transition:

Transition	Rate
$i \rightarrow i + 1$	$\lambda, i \ge 0$

According to action 1, emergency calls are served by the available t ambulance units with rate equal to  $\mu$ . The system makes the following transition:

Transition	Rate
$i \rightarrow i - 1$	$t\mu, t \leq i \leq N$
$i \rightarrow i - 1$	$i\mu$ , $1 \le i \le t-1$

That is, if  $i \ge t$ , the system sends t available ambulances when i calls are waiting to be served and the system serves with i ambulances if i < t calls are waiting to be served. In what follows, for reasons of simplicity, in the formulation of the simplified model as a semi-Markov decision model that we present below, we consider that t = 1.

For states  $i \in S$ , such that  $N + 1 \le i \le Q - 1$ , the possible actions are actions 0 and 2. According to action 2, k emergency calls are redirected to available neighbour systems with a rate equal to  $\gamma$ . If a redirection of calls is chosen, when the number of emergency calls is equal to i,  $N + 1 \le i \le Q - 1$ , the number of emergency calls after the calls redirection is reduced to i - k,  $1 \le k \le i - 1$ , with probability  $p(1-p)^{i-k-1}$  and is reduced to zero emergency calls with probability  $(1 - p)^{i-1}$ , where  $p \in (0, 1)$ . In state i = Q, it is reasonable to assume that the only feasible action is action 2. The system makes the following transitions:

Transition	Rate
$i \rightarrow i - k$	$\gamma p(1-p)^{i-k-1}, N+1 \le i \le Q-1, 1 \le k \le i-1$
$i \rightarrow 0$	$\gamma(1-p)^{i-1}, N+1 \le i \le Q-1$
$Q \rightarrow Q - k$	$\gamma p(1-p)^{Q-k-1}, 1 \le k \le Q-1$
$Q \rightarrow 0$	$\gamma(1-p)Q^{-1}$

The rate  $\lambda$  is necessarily equal to zero when i = Q, since Q is the maximum number of emergency calls that the system can admit. The action space and a transition graph of the model for a specific case of some input model parameters values are depicted in Figure 2.





We assume that the system incurs a holding cost at a rate equal to  $h_i > 0$  when there are *i* emergency calls being in the system waiting to be served where  $h_i$  is bounded increasing and non-negative function with respect to *i*. We also assume that there is a service cost at a rate equal to C > 0, whenever an emergency call is served by the system. A redirection cost is also incurred at a rate equal to R > 0, whenever a redirection of calls is chosen. Our goal is to consider an appropriate MDP in continuous time for which we seek a decision rule that minimises a predefined optimality criterion. We shall consider the criterion of minimising the long-run expected average cost per unit time. For this criterion, a semi-Markov decision model is in fact determined by the following three characteristics:

- a the probability  $p_{ij}(a)$  that at the next decision epoch the system will be in state *j* if action *a* is chosen in the present state *i*
- b the expected time  $T_a(i)$  until the next decision epoch if action *a* is chosen in the present state *i*
- c the expected cost  $C_a(i)$  until the next decision epoch if action *a* is chosen in the present state *i*.

Let  $p_{ij}(a)$  be the transition probabilities from state *i* to state *j* if action *a* is selected in state *i* and let  $T_a(i)$  and  $C_a(i)$  be the one-step expected transition time and cost, respectively, when action *a* is chosen in state *i*. These quantities can be computed for each state  $i \in S$  and for each possible action *a*. These quantities are specified below.

• Non-zero one-step transition probabilities

$$\begin{split} p_{01}(0) &= 1, \ p_{i,i+1}(0) = 1, 1 \leq i \leq N, \\ p_{i,i+1}(1) &= \lambda(\lambda + u)^{-1}, \ p_{i,i-1}(1) = \mu(\lambda + u)^{-1}, 1 \leq i \leq N. \\ p_{i,i+1}(0) &= 1, \ p_{i,i-k}(2) = \gamma p(1-p)^{i-k-1}(h+\gamma)^{-1}, N+1 \leq i \leq Q-1, 1 \leq k \leq i-1, p \in (0, 1). \\ p \in (0, 1). \\ p_{i0}(2) &= \gamma(1-p)^{i-1}(\lambda + \gamma)^{-1}, \ p_{i,i+1}(2) = \lambda(\lambda + \gamma)^{-1}, N+1 \leq i \leq Q-1. \\ p_{Q,Q-k}(2) &= p(1-p)^{Q-k-1}, 1 \leq k \leq Q-1, p_{Q0}(2) = (1-p)^{Q-1}, \ p \in (0, 1). \end{split}$$

• One-step expected times

$$T_0(0) = \lambda^{-1}, 0 \le i \le Q - 1, T_1(i) = (\lambda + \mu)^{-1}, 1 \le i \le N,$$
  
$$T_2(i) = (\lambda + \gamma)^{-1}, N + 1 \le i \le Q - 1, T_2(Q) = y^{-1}.$$

• One-step expected costs

$$C_{0}(i) = h_{i}\lambda^{-1}, 0 \le i \le Q - 1^{\circ}, C_{1}(i) = (C + h_{i})(\lambda + \mu)^{-1}, 1 \le i \le N.$$
  
$$C_{2}(i) = (R + h_{i})(\lambda + \gamma)^{-1}, N + 1 \le i \le Q - 1^{\circ}, C_{2}(Q) = (R + h_{Q})\gamma^{-1}.$$

A computational treatment of the problem is possible by applying various standard Markov decision algorithms such as the standard value-iteration algorithm and the standard policy-iteration algorithm.

#### **3** Generalisation of the simplified model

In this section, we generalise the model described in the previous section. We consider an EMS system which consists of  $N_A$  ALS units and  $N_B$  BLS units which are available to serve patients. The EMS admits incoming emergency calls which are divided into two classes:

- a urgent, high-priority calls for which the patient's life is potentially at risk
- b less urgent, low-priority calls.

Consider the two-dimensional Markov process in continuous-time  $\{(X(t), Y(t)), t \ge 0\}$ , where the random variables X(t) and Y(t) denote the state of the process at time  $t \ge 0$ . More precisely, the random variables X(t) and Y(t) denote the number of high-priority calls and the number of low-priority calls, respectively, which are in the system waiting to be served at time  $t \ge 0$ . Thus, if X(t) = i and Y(t) = j, the system state can be described by two values:

- 1 the number *i* of high-priority calls
- 2 the number *j* of low-priority calls, being in the system waiting to be served at time  $t \ge 0$ .

We assume that the system can admit a maximum number  $Q_H$  of high-priority calls and a maximum number  $Q_L$  of low-priority calls. We further assume that the maximum number of high-priority calls and the maximum number of low-priority calls, respectively, are greater than the number of the available ALS units and of the available BLS units. Thus, the state space of the system is defined by the set:

$$S = \{(i, j) \mid 0 \le i \le Q_H, 0 \le j \le Q_L\}.$$

The decision epochs include the epochs at which a call, either a high-priority one or a low-priority one, arrives to the system. We assume that the EMS system has an annual operating budget equal to D. We also assume that the annual operating costs of a single ALS and of a single BLS ambulance are equal to  $C_A$  and to  $C_B$ , respectively. The rates  $\lambda_H$  and  $\lambda_L$  are necessarily equal to zero when  $i = Q_H$  and  $j = Q_L$ , since  $Q_H$  and  $Q_L$  are the

maximum number of high-priority and low-priority calls that the system can admit, respectively. The action space and a transition graph of the model for a specific case of some input model parameters values are depicted in Figure 3.

Figure 3 The action space and the transition graph of the process for  $N_A = 4$ ,  $N_B = 2$ ,  $Q_H = 6$ ,  $Q_L = 4$  (see online version for colours)



We assume that the system incurs a holding cost at a rate equal to  $h_i > 0$  (or to  $\tilde{h}_j > 0$ ) when there are *i* high-priority calls (or *j* low-priority calls) being in the system waiting to be served where  $h_i$  (or  $\tilde{h}_j$ ) is bounded increasing and non-negative functions with respect to *i* (and with respect to *j*), respectively. We also assume that there is a service cost at a rate equal to  $C_H > 0$  (or to  $C_L > 0$ ) whenever a high-priority (or a low-priority call) is served by the system. A redirection cost is also incurred at a rate equal to  $R_H > 0$  (or to  $R_L > 0$ ) whenever a redirection of high-priority calls (or a redirection of low-priority calls) is chosen. It is reasonable to assume that  $C_H > C_L$ .

Let  $f(N_A, N_B)$  denote some measure of system performance associated with a fleet operating  $N_A$  ALS units and  $N_B$  BLS units, where  $N_AC_A + N_BC_B \leq D$ . Our goal is to consider an appropriate MDP in continuous time for which we seek a decision rule that minimises a predefined optimality criterion for vehicle mixes  $(N_A, N_B)$  over a set of candidate fleets under the above cost restriction. We shall consider the criterion of minimising the long-run expected average cost per unit time. In next section, we provide the semi-MDP formulation for the generalised version of the model.

### 4 Formulation as a semi-Markov decision model for the generalised system

Let  $(N_A, N_B)$  be the total number of ambulances of each category that has an EMS system in its fleet. As in the simplified model, during each shift only a number of these vehicles are available for service. Suppose that *k* ALS ambulances and  $\theta$  BLS ambulances are available for service during each shift. We also assume that at any given time some of the vehicles  $(k, \theta)$  either remain at the base for refuelling and maintenance or are a safety fleet for propulsion in some extreme unforeseen situations. Let (i, j) be the number of calls waiting to be served at any given time and t < k,  $\rho < \theta$ , the vehicles available for these calls (at ambulance stations remain k - t ALS and  $\theta - \rho$  BLS vehicles for refuelling or preparedness in extreme cases). The system with the forwarding of *t*-ALS and  $\rho$ -BLS vehicles makes the following transition:

Rate

Transition

$(i,j) \rightarrow (i-1,j)$	$t\mu, t \leq i \leq N_A$
$(i,j) \rightarrow (i-1,j)$	$i\mu$ , $1 \le i \le t - 1$
$(i,j) \rightarrow (i,j-1)$	$ ho\mu,  ho \leq j \leq N_B$
$(i,j) \rightarrow (i,j-1)$	$j\mu$ , $1 \leq j \leq \rho - 1$

The most realistic option for an EMS system is to forward for (i, j) calls waiting to be served, a number of ALS units and a number of BLS units for these calls. We could consider the impact on the system's selection to send  $(t_1, \rho_1)$  ambulances compared to  $(t_2, \rho_2)$  ambulances, but in this paper, we want to compare the dynamics of ALS ambulances in service and at the cost of the system compared to BLS ambulances. So, we send either *k*-ALS ambulances to handle *i*-ALS calls (and  $\theta$ -BLS vehicles stay at the station) or  $\theta$ -BLS vehicles to handle *j*-BLS calls (and *k*-ALS vehicles stay at the station). Consequently, the system must choose between the following transitions:

Transition 1	Rate
$(i,j) \rightarrow (i-1,j)$	$k\mu, k \leq i \leq N_A$
$(i,j) \rightarrow (i-1,j)$	$i\mu$ , $1 \le i \le k-1$
Transition 2	Rate
$(i,j) \rightarrow (i,j-1)$	$\theta\mu, \theta \leq j \leq N_B$
$(i, j) \rightarrow (i, j-1)$	$j\mu$ , $1 \le j \le \theta - 1$

In what follows, for reasons of simplicity, in the model that we describe below, we consider that k = 1 and  $\theta = 1$ .

Let  $p_{(i,j)(i',j')}(a)$  be the transition probabilities from state (i, j) to state (i', j') if action a is selected in state (i, j) and let  $T_a(i)$  and  $C_a(i, j)$  be the one-step expected transition time and cost, respectively, when action a is chosen in state (i, j). These quantities can be computed for each state  $(i, j) \in S$  and for each possible action a. For each state  $(i, j) \in S$ , we distinct four cases.

#### 4.1 First case

We consider the states  $(i, j) \in S$ , for which  $0 \le i \le N_A$  and  $0 \le j \le N_B$ . It is intuitively reasonable to assume that in state (0, 0) the only possible action is action 0. For states (0, j),  $1 \le j \le N_B$ , the possible actions are actions 0 and action 2 (serve a low-priority call). For states (i, 0),  $1 \le i \le N_A$ , the possible actions are actions 0 and 1 (serve a high-priority call). For the states (i, j),  $1 \le i \le N_A$ ,  $1 \le j \le N_B$ , the possible actions are actions 0, 1 (serve a high-priority call) and 2 (serve a low-priority call).

Action 0. High-priority calls and low-priority calls arrive to the system according to independent Poisson processes with rates equal to  $\lambda_H$  and to  $\lambda_L$ , respectively. The system makes the following transitions:

Transition	Rate
$(i,j) \rightarrow (i+1,j)$	$\lambda_H$
$(i,j) \rightarrow (i+1,j)$	$\lambda_L$

Action 1 (serve a high-priority call). A high-priority call is served by one of the available  $N_A$  ALS units with rate equal to  $\mu$ . The system makes the following transition:

TransitionRate
$$(i,j) \rightarrow (i-1,j)$$
 $\mu, 1 \le i \le N_A$ 

Action 2 (serve a low-priority call). A low-priority call is served by one of the available  $N_B$  BLS units with rate equal to  $\mu$ . The system makes the following transition:

TransitionRate
$$(i,j) \rightarrow (i,j-1)$$
 $\mu, 1 \le j \le N_B$ 

In Case 1, the non-zero one-step transition probabilities from each state to another state, if an action is chosen, and the one-step expected transition times and costs, for each possible action, can be specified. For example, for states (i, j),  $1 \le i \le N_A$ ,  $1 \le j \le N_B$ , where action a belongs to the set  $\{0, 1, 2\}$ , we have:

Non-zero one-step transition probabilities

$$\begin{aligned} p_{(i,j)(i+1,j)}(0) &= \lambda_H \left( \lambda_H + \lambda_L \right)^{-1}, \ p_{(i,j)(i+1,j)}(0) = \lambda_L \left( \lambda_H + \lambda_L \right)^{-1}, 1 \le i \le N_A, 1 \le j \le N_B. \\ p_{(i,j)(i+1,j)}(1) &= \lambda_H \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \ p_{(i,j)(i+1,j)}(1) = \lambda_L \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \\ p_{(i,j)(i-1,j)}(1) &= \mu \left( \lambda_H + \lambda_L + \mu \right)^{-1}, 1 \le i \le N_A, 1 \le j \le N_B. \\ p_{(i,j)(i+1,j)}(2) &= \lambda_H \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \ p_{(i,j)(i,j+1)}(2) = \lambda_L \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \\ p_{(i,j)(i,j-1)}(2) &= \mu \left( \lambda_H + \lambda_L + \mu \right)^{-1}, 1 \le i \le N_A, 1 \le j \le N_B. \end{aligned}$$

One-step expected times

$$T_{0}(i, j) = (\lambda_{H} + \lambda_{L})^{-1}, T_{1}(i, j) = (\lambda_{H} + \lambda_{L} + \mu)^{-1},$$
  
$$T_{2}(i, j) = \lambda_{H} + \lambda_{L} + \mu()^{-1}, 1 \le i \le N_{A}, 1 \le j \le N_{B}.$$

One-step expected costs

$$C_{0}(i, j) = (h_{i} + \tilde{\lambda}_{j})(h_{H} + \lambda_{L})^{-1}, C_{1}(i, j) = (h_{i} + \tilde{\lambda}_{j} + C_{H})(\lambda_{H} + \lambda_{L} + \mu)^{-1},$$
  

$$C_{2}(i, j) = (h_{i} + \tilde{\lambda}_{j} + C_{L})(\lambda_{H} + \lambda_{L} + \mu)^{-1}, 1 \le i \le N_{A}, 1 \le j \le N_{B}.$$

#### 4.2 Second case

We consider the states  $(i, j) \in S$ , for which  $0 \le i \le N_A$  and  $N_B + 1 \le j \le Q_L$ . For states (0, j) $N_B + \le j \le Q_L - 1$ , the possible actions are action 0 and action 3 (redirect *k* low-priority calls). For states (i, j),  $1 \le i \le N_A$ ,  $N_B + 1 \le j \le Q_L - 1$ , the possible actions are actions 0, 1 (serve a high-priority call) and 3 (redirect *k* low-priority calls). For the state  $(0, Q_L)$ , the possible actions are Actions 0 (with  $\lambda_L = 0$ ) and 3 (redirect *k* low-priority calls). For the states  $(i, Q_L)$ ,  $1 \le i \le N_A$ , the possible actions are actions 0 (with  $\lambda_L = 0$ ), 1 (serve a high-priority call) and 3 (redirect k low-priority calls). Action 0 and action 1 are the same as in the first case.

Action 3 (redirect k low-priority calls). According to action 3, if the number of low-priority calls is j, then k low-priority calls can be redirected to available neighbour systems with a constant rate equal to  $\gamma$ .

If calls redirection is chosen, when the number of low-priority calls is *j*, the number of low-priority calls is reduced to j - k,  $1 \le k \le j - 1$ , with probability  $p(1 - p)^{j-k-1}$ , and is reduced to zero low-priority calls with probability  $(1 - p)^{j-1}$ , where  $p \in (0, 1)$ . The system makes transitions to the states (i, j - k),  $1 \le k \le j - 1$  and to the states (i, 0).

Transition	Rate
$(i,j) \rightarrow (i,j-k)$	$\gamma(1-p)^{j-k-1}, 1 \le k \le j-1$
$(i,j) \rightarrow (i,0)$	$\gamma(1-p)^{j-1}$

In Case 2, the non-zero one-step transition probabilities from each state to another state, if an action is chosen, and the one-step expected transition times and costs, for each possible action, can be specified. For example, for states (i, j),  $1 \le i \le N_A$ ,  $N_B + 1 \le j \le Q_L - 1$ , where the possible actions are actions 0, 1 and 3, we have:

Non-zero one-step transition probabilities

$$\begin{aligned} p_{(i,j)(i+1,j)}(0) &= \lambda_H \left( \lambda_H + \lambda_L \right)^{-1}, \ p_{(i,j)(i,j+1)}(0) = \lambda_L \left( \lambda_H + \lambda_L \right)^{-1}, 1 \le i \le N_A, N_B + 1 \le j \le Q_L - 1. \\ p_{(i,j)(i+1,j)}(1) &= \lambda_H \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \ p_{(i,j)(i,j+1)}(1) = \lambda_L \left( \lambda_H + \lambda_L + \mu \right)^{-1}, \\ p_{(i,j)(i-1,j)}(1) &= \mu \left( \lambda_H + \lambda_L + \mu \right)^{-1}, 1 \le i \le N_A, N_B + 1 \le j \le Q_L - 1. \\ p_{(i,j)(i+1,j)}(3) &= \lambda_H \left( \lambda_H + \lambda_L + \gamma \right)^{-1}, \ p_{(i,j)(i,j+1)}(3) = \lambda_L \left( \lambda_H + \lambda_L + \gamma \right)^{-1}, \\ 1 \le i \le N_A, N_B + 1 \le j \le Q_L - 1. \\ p_{(i,j)(i,j-k)}(3) &= \gamma p (1-p)^{j-k-1} \left( \lambda_H + \lambda_L + \gamma \right)^{-1}, 1 \le k \le j - 1, \\ p_{(i,j)(i,0)}(3) &= \gamma (1-p)^{j-1} \left( \lambda_H + \lambda_L + \gamma \right)^{-1}, 1 \le i \le N_A, N_B + 1 \le j \le Q_L - 1, p \in (0,1) \end{aligned}$$

One-step expected times

$$T_{0}(i, j) = (\lambda_{H} + \lambda_{L})^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1.$$
  

$$T_{1}(i, j) = (\lambda_{H} + \lambda_{L} + \mu)^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1.$$
  

$$T_{3}(i, j) = (\lambda_{H} + \lambda_{L} + \gamma)^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1.$$

One-step expected costs

$$C_{0}(i, j) = (h_{i} + \tilde{\lambda}_{j})(h_{H} + \lambda_{L})^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1.$$

$$C_{1}(i, j) = (h_{i} + \tilde{\lambda}_{j} + C_{H})(\mu + h_{H} + \lambda_{L})^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1$$

$$C_{3}(i, j) = (R_{L} + \lambda_{i} + \tilde{\lambda}_{j})(h_{H} + \lambda_{L} + \gamma)^{-1}, 1 \le i \le N_{A}, N_{B} + 1 \le j \le Q_{L} - 1.$$

#### 4.3 Third case

We consider the states  $(i, j) \in S$ , for which  $N_A + 1 \le i \le Q_H$  and  $0 \le j \le N_B$ . For states (i, 0) $N_A + 1 \le i \le Q_H - 1$ , the possible actions are action 0 and action 3 (redirect *l* high-priority calls). For states (i, j),  $N_A + 1 \le i \le Q_H - 1$ ,  $1 \le j \le N_B$ , the possible actions are actions 0, 2 (serve a low-priority call) and 3 (redirect *l* high-priority calls). For the state  $(Q_H, 0)$ , the possible actions are actions 0 (with  $\lambda_H = 0$ ) and 3 (redirect *l* high-priority calls). For the states  $(Q_H, j)$ ,  $1 \le j \le N_B$ , the possible actions are actions 0 (with  $\lambda_H = 0$ ), 2 (serve a low-priority call) and 3 (redirect *l* high-priority calls). Action 0 and action 2 are the same as in the first case.

Action 3 (redirect l high-priority calls). According to action 3, if the number of low-priority calls is *i*, then *l* low-priority calls can be redirected to available neighbour systems with a constant rate equal to  $\gamma$ . If calls redirection is chosen, when the number of low-priority calls is *i*, the number of low-priority calls is reduced to i - l,  $1 \le l \le i - 1$ , with probability  $p(1 - p)^{i-l-1}$ , and is reduced to zero low-priority calls with probability  $(1 - p)^{i-l-1}$ , where  $p \in (0, 1)$ . The system makes transitions to the states (l, j),  $1 \le l \le i - 1$  and to the states (0, j).

TransitionRate
$$(i,j) \rightarrow (i-l,j)$$
 $\gamma(1-p)^{i-l-1}, 1 \le l \le i-1$  $(i,j) \rightarrow (0,j)$  $\gamma(1-p)^{i-l}$ 

In Case 3, the non-zero one-step transition probabilities from each state to another state, if an action is chosen, and the one-step expected transition times and costs, for each possible action, can be specified. For example, for states (i, j),  $N_A + 1 \le i \le Q_H - 1$ ,  $1 \le j \le N_B$ , where the possible actions are actions 0, 2 and 3, we have:

Non-zero one-step transition probabilities

$$\begin{split} p_{(i,j)(i+1,j)}(0) &= \lambda_{H} \left( \lambda_{H} + \lambda_{L} \right)^{-1}, \ p_{(i,j)(i,j+1)}(0) = \lambda_{L} \left( \lambda_{H} + \lambda_{L} \right)^{-1}, \\ N_{A} + 1 &\leq i \leq Q_{H} - 1, 1 \leq j \leq N_{B}. \\ p_{(i,j)(i+1,j)}(2) &= \lambda_{H} \left( \lambda_{H} + \lambda_{L} + \mu \right)^{-1}, \ p_{(i,j)(i,j+1)}(2) = \lambda_{L} \left( \lambda_{H} + \lambda_{L} + \mu \right)^{-1}, \\ p_{(i,j)(i+1,j)}(2) &= \mu \left( \lambda_{H} + \lambda_{L} + \mu \right)^{-1}, \ N_{A} + 1 \leq i \leq Q_{H} - 1, 1 \leq j \leq N_{B}. \\ p_{(i,j)(i+1,j)}(3) &= \lambda_{H} \left( \lambda_{H} + \lambda_{L} + \gamma \right)^{-1}, \ p_{(i,j)(i,j+1)}(3) = \lambda_{L} \left( \lambda_{H} + \lambda_{L} + \gamma \right)^{-1}, \\ N_{A} + 1 \leq i \leq Q_{H} - 1, 1 \leq j \leq N_{B}. \\ p_{(i,j)(i-l,j)}(3) &= \gamma p (1-p)^{i-l-1} \left( \lambda_{H} + \lambda_{L} + \gamma \right)^{-1}, \ N_{A} + 1 \leq i \leq Q_{H} - 1, 1 \leq l \leq i-1, \ p \in (0, 1). \\ p_{(i,j)(0,j)}(3) &= \gamma (1-p)^{i-1} \left( \lambda_{H} + \lambda_{L} + \gamma \right)^{-1}, \ N_{A} + 1 \leq i \leq Q_{H} - 1, \ p \in (0, 1). \end{split}$$

One-step expected times

$$T_0(i, j) = (\lambda_H + \lambda_L)^{-1}, N_A + 1 \le i \le Q_H - 1, 1 \le j \le N_B.$$

$$T_{2}(i, j) = (\mu + \lambda_{H} + \lambda_{L})^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, 1 \le j \le N_{B}.$$
  
$$T_{3}(i, j) = (\lambda_{H} + \lambda_{L} + \gamma)^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, 1 \le j \le N_{B}.$$

• One-step expected costs

$$C_{0}(i, j) = (h_{i} + \tilde{\lambda}_{j})(h_{H} + \lambda_{L})^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, 1 \le j \le N_{B}.$$

$$C_{0}(i, j) = (h_{i} + \tilde{\lambda}_{j} + C_{L})(\mu + h_{H} + \lambda_{L})^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, 1 \le j \le N_{B}.$$

$$C_{3}(i, j) = (h_{i} + \tilde{\lambda}_{j} + R_{H})(h_{H} + \lambda_{L} + \gamma)^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, 1 \le j \le N_{B}.$$

#### 4.4 Fourth case

We consider the states  $(i, j) \in S$ , for which  $N_A + 1 \le i \le Q_H$  and  $N_B + 1 \le j \le Q_L$ . For the states (i, j), such that  $N_A + 1 \le i \le Q_H - 1$  and  $N_B + 1 \le j \le Q_L - 1$ , two actions are possible, action 0 and action 3 (redirect *l* high-priority calls or *k* low-priority calls). For the states  $(Q_H, j)$ , such that  $N_B + 1 \le j \le Q_L - 1$ , two actions are possible, action 0 (with  $\lambda_H = 0$ ) and action 3 (redirect *l* high-priority calls or *k* low-priority calls). For the states  $(i, Q_L)$ , such that  $N_A + 1 \le i \le Q_H - 1$ , two actions are possible, Action 0 (with  $\lambda_L = 0$ ) and action 3 (redirect *l* high-priority calls or *k* low-priority calls). For the states  $(i, Q_L)$ , such that  $N_A + 1 \le i \le Q_H - 1$ , two actions are possible, Action 0 (with  $\lambda_L = 0$ ) and action 3 (redirect *l* high-priority calls or *k* low-priority calls). For the states  $(Q_H, Q_L)$ , the only possible action is action 3 (redirect *l* high-priority calls or *k* low-priority calls). action 0 is the same as in the first case.

Action 3 (redirect l high-priority calls or k low-priority calls). According to action 3, if the number of high-priority calls is *i* and the number of low-priority calls is *j*, then *l* high-priority calls or *k* low-priority calls can be redirected to available neighbours systems with a constant rate for each kind of redirection equal to *y*. If calls redirection is chosen, when the number of high-priority calls is *i*, the number of high-priority calls is reduced to i - l,  $1 \le l \le i - 1$ , with probability  $p(1 - p)^{i-l-1}$  and is reduced to zero high-priority calls with probability  $(1 - p)^{i-l}$ , where  $p \in (0, 1)$ . If calls redirection is chosen, when the number of low-priority calls is *j*, the number of low-priority calls is reduced to j - k,  $1 \le k \le j - 1$ , with probability  $p(1 - p)^{i-k-1}$  and is reduced to zero low-priority calls with probability  $(1 - p)^{j-1}$ , where  $p \in (0, 1)$ . The system makes transitions to the states (i - l, j) and to the states (0, j) or to the states (i, j - k) and to the states (i, 0).

Transition	Rate
$(i, j) \rightarrow (i - 1, j)$	$\gamma(1-p)^{i-l-1}, 1 \leq l \leq i-1$
$(i,j) \rightarrow (0,j)$	$\gamma(1-p)^{i-1}$
$(i,j) \rightarrow (i,j-k)$	$\gamma(1-p)^{j-k-1}, 1 \le k \le j-1$
$(i,j) \rightarrow (i,0)$	$\gamma(1-p)^{j-1}$

In Case 4, the non-zero one-step transition probabilities from each state to another state, if an action is chosen, and the one-step expected transition times and costs, for each possible action, can be specified. For example, for the states (i, j), such that  $N_A + 1 \le i \le Q_H - 1$  and  $N_B + 1 \le j \le Q_L - 1$ , where the possible actions are actions 0 and 3, we have:

• Non-zero one-step transition probabilities

$$\begin{split} p_{(i,j)(i+1,j)}(0) &= \lambda_{H} \left( \lambda_{H} + \lambda_{L} \right)^{-1}, \ p_{(i,j)(i,j+1)}(0) = \lambda_{L} \left( \lambda_{H} + \lambda_{L} \right)^{-1}, \\ N_{A} + 1 &\leq i \leq Q_{H} - 1, N_{B} + 1 \leq j \leq Q_{L} - 1. \\ p_{(i,j)(i+1,j)}(3) &= \lambda_{H} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1}, \\ p_{(i,j)(i,j+1)}(3) &= \lambda_{L} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, N_{B} + 1 \leq j \leq Q_{L} - 1. \\ p_{(i,j)(i,j+1)}(3) &= \gamma p(1-p)^{i-l-1} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, N_{B} + 1 \leq j \leq Q_{L} - 1, \\ 1 \leq l \leq i-1, p \in (0, 1). \\ p_{(i,j)(0,j)}(3) &= \gamma (1-p)^{j-l} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, p \in (0, 1). \\ p_{(i,j)(i,j-k)}(3) &= \gamma p(1-p)^{j-k-1} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, 1 \leq k \leq j - 1, p \in (0, 1). \\ p_{(i,j)(i,j-k)}(3) &= \gamma p(1-p)^{j-k-1} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, 1 \leq k \leq j - 1, p \in (0, 1). \\ p_{(i,j)(i,j-k)}(3) &= \gamma p(1-p)^{j-k-1} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, p \in (0, 1). \\ p_{(i,j)(i,j-k)}(3) &= \gamma p(1-p)^{j-k-1} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, p \in (0, 1). \\ p_{(i,j)(i,j-k)}(3) &= \gamma (1-p)^{j-l} \left( \lambda_{H} + \lambda_{L} + 2\gamma \right)^{-1} N_{A} + 1 \leq i \leq Q_{H} - 1, \\ N_{B} + 1 \leq j \leq Q_{L} - 1, p \in (0, 1). \end{split}$$

• One-step expected times

$$T_0(i, j) = (\lambda_H + \lambda_L)^{-1}, N_A + 1 \le i \le Q_H - 1, N_B + 1 \le j \le Q_L - 1.$$
  
$$T_3(i, j) = (\lambda_H + \lambda_L + 2\gamma)^{-1}, N_A + 1 \le i \le Q_H - 1, N_B + 1 \le j \le Q_L - 1.$$

One-step expected costs

$$C_{0}(i, j) = (h_{i} + \tilde{\lambda}_{j})(h_{H} + \lambda_{L})^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, N_{B} + 1 \le j \le Q_{L} - 1.$$
  

$$C_{3}(i, j) = (h_{i} + \tilde{\lambda}_{j} + R_{H} + R_{L})(h_{H} + \lambda_{L} + 2\gamma)^{-1}, N_{A} + 1 \le i \le Q_{H} - 1, N_{B} + 1 \le j \le Q_{L} - 1.$$

Based on the above semi-MDP formulation, a direct implementation of the standard Markov decision algorithms is possible for the numerical computation of an optimal stationary policy. For a detailed description of these algorithms we refer, for example, to the books of Puterman (1994), Heyman and Sobel (2003) and Tijms (2003).

#### 5 Numerical results for the generalised model

In this section, we implement numerically the generalised model using information and making assumptions for Virginia Beach EMS (VBEMS) which is the largest volunteer-based EMS system in the USA. It provides EMS to the residents and visitors of

Virginia Beach, a city in USA with population 450.000 inhabitants and classified as tiered system with 37 ambulances.

The system provides training to volunteers either in EMTs skills either in paramedic skills and from 1940s it operates as all-volunteer system. In 2004, the system hired career staff to satisfy the growing demand and it converted from all-volunteer to volunteer-based system. We suppose that the system hires both career EMTs and career paramedics, so there is a distinction between the members of the staff ambulances which is:  $V_{par}$  (volunteer paramedic),  $V_{EMT}$  (volunteer EMTs),  $C_{par}$  (career paramedic),  $C_{EMT}$  (career EMTs). We suppose additionally that an ALS ambulance staffed by one paramedic and one EMTs and a BLS ambulance staffed by two EMTs. From the annual report of 2018, we drew that 85% of all ambulance crews were staffed by volunteers. So, for an ALS ambulance which requires three crews to operate 24 hours per day, we have the following crews: ( $C_{par}$ ,  $V_{EMT}$ ) for one shift and ( $V_{par}$ ,  $V_{EMT}$ ) for two shifts and for a BLS ambulance the crews are ( $C_{EMT}$ ,  $V_{EMT}$ ) for one shift and ( $V_{EMT}$ ,  $V_{EMT}$ ) for two shifts.

The annual salary for career paramedic and EMTs can vary depending on many factors including education, certification, additional skills and the number of years spent in the profession. We may assume that the average paramedic salary is \$40.000 per year and the average EMTs salary is \$35.000 per year. We also suppose that the cost to equip and operate an ALS ambulance in 24/7 is \$210.000 per year and for BLS ambulance is \$185.000 per year. This cost includes supplies, operating cost, maintenance and other related costs. Therefore, in our approach, we have that  $C_A = $250.000$  and  $C_B = $220.000$ . Also, the annual operating budget for VBEMS, according to the annual report of 2018, is D = \$8.000.000. Then, we have that the set of candidate fleets for VBEMS is:

$$\Gamma = \{ (N_A, N_B) : 25N_A + 22N_B \le 800 \}.$$

In an EMS system, there are two service procedures, the service of the call canter and the ambulance service. As far as the call centre service is concerned, we suppose that a high priority call, which requires an ALS ambulance has duration equivalent to the arrival of ambulance on scene (it is a reasonable assumption because the call taker may provide instruction to the caller, which in many cases is essential to stabilising or saving a life) and a low priority call which requires an BLS ambulance has duration equivalent to the ambulance departure. For VBEMS the time up to the arrival on scene is nine-fold to the time up to the ambulance departure. So, we suppose that  $R_H =$ \$ 45 per hour and  $R_L =$ \$5 per hour as the redirection costs for high and low priority calls, respectively.

There are many research papers about the maximum number of calls allowed to be admitted in a call centre, see for example, Akl et al. (2005). In this paper we do not take that factor into account, but we aim to examine all the feasible vehicle mixes of VBEMS, so we assume that the dispatcher of the system can handle the maximum  $Q_H = 40$  of high priority calls and  $Q_L = 40$  of low priority calls.

Some EMS systems redirect calls to an advice line that can be handled telephonically with a clinician to develop an alternative to dispatching an ambulance. There are studies that determined that 15%–20% of all EMS calls could be diverted to an advice line. Some other systems redirect calls to other EMS systems when they enter in red alert. For VBEMS we assume that it redirects 20% of incoming EMS calls when the number of high priority calls is greater than the number of ALS ambulances. That percentage allocates equally to high and low priority calls, so we have  $\gamma = 0.385$  calls per hour. This

redirection rate remains stable when the number of high priority calls is greater than the number of ALS ambulances and the number of low priority calls is less than the number of BLS ambulances and when the number of high priority calls is less than the number of ALS ambulances and the number of low priority calls is greater than the number of BLS ambulances. If a redirection of calls is chosen, we define the discount factor p = 0.5 for the transition probabilities to fewer calls.

Using information from Centres for Medicare and Medicaid Services (CMS), we assume that the service cost whenever a high-priority call is served is  $C_H = $420$  per hour and the service cost whenever a low-priority call is served is  $C_L = $360$  per hour. Furthermore, following the logic of Chong et al. (2015), we assume that the holding costs of the system are determined by the service costs in the following way: the patient do not feel obliged to pay if he waits more than one hour to being served so, we have  $h_i = $420 \cdot i$  the holding cost of the system when there are *i*-high priority calls waiting to be served per hour and  $h_j = $420 \cdot j$  the holding cost of the system when there are *j*-low priority calls waiting to be served per hour. Finally, from the annual report of 2018, we set the values  $\lambda_H = 1.18$  calls per hour,  $\lambda_L = 2.67$  calls per hour and  $\mu = 1.2$  customers per hour.

We solve the case of VBEMS using the above values with policy iteration algorithm and with value iteration algorithm. In Figure 4, we observe the minimum long run average cost per hour on the attainable combinations of VBEMS system. The best choice for VBEMS is to operate as All-ALS system.





Virginia Beach EMS

In Table 1, the successive policies at each iteration generated by the policy iteration algorithm and their average costs are presented, for the vehicle mix (32, 0). We denote by r(i, j) the action selected in state (i, j). The algorithm is stopped after three iterations with the value of the minimum long-run expected average cost approximately equal to 24 thousand dollars. Note that, as it is empirically verified in policy iteration algorithm, the average costs of the policies generated by the algorithm have shown great improvements in their values in the first few iterations. The Markov decision algorithms are implemented by running the corresponding programs, using MATLAB 2016a software

on a personal computer equipped with an AMD Quad Core A12-9720P, 2.7 GHz processor and 8GB of RAM. The computation time (CPU time) of the policy iteration algorithm for the combination (32, 0) was 9.62 seconds.

In Figure 5, we see the optimal policy for VBEMS for the combination  $N_A = 32$  and  $N_B = 0$ , where we have the minimum cost among all vehicle mixes. We notice that the optimal policy for VBEMS is to service the high priority calls at any state where it is possible and to redirect to the rest of the states.

**Figure 5** Optimal policy for VBEMS in combination  $N_A = 32$  and  $N_B = 0$  (see online version for colours)

| 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1                     | 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1                | 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1 | 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1    | 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1 | 1<br>1<br>1<br>1<br>1<br>1              | 1<br>1<br>1<br>1<br>1<br>1 | 1<br>1<br>1<br>1<br>1<br>1 | 1<br>1<br>1<br>1<br>1   | 1<br>1<br>1<br>1  
   
  | 1<br>1<br>1 | 1<br>1<br>1  
   
   
  | 1<br>1<br>1   | 1 1   | 1<br>1   
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
|--|---|--|---|--------------------------------------|---|----------------------------|----------------------------|---
--
--
--|-------------
--
--
---|---|---
---	---	---	---	---	---	---	---	---	---
---	---	---	---	---	---				
- 1 1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1	1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1	1 1 1 1 1 1	1 1 1 1 1	1 1 1 1 1 1
   
  | 1<br>1<br>1 | 1 1 1 1  
   
   
  | 1   | 1   | 1  
  | ÷.  |   | -   | _   | _   |   |   |   |   | _   | _  
  | _   | _   | _   | _   | -   | -   | С.  | ×   | ×.  | ×.  | ×.                 
  | ×.  | ×.  | . T   |
|  | 1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1<br>1                | 1<br>1<br>1<br>1<br>1<br>1<br>1                | -<br>1<br>1<br>1<br>1<br>1<br>1<br>1    | 1<br>1<br>1<br>1<br>1<br>1<br>1      | 1<br>1<br>1<br>1<br>1                   | 1<br>1<br>1<br>1<br>1      | 1<br>1<br>1<br>1<br>1      | 1<br>1<br>1<br>1  | 1 1 1 1   
   
  | 1           | 1  
   
   
  | 1   | Ξ.  |  
  |   | 1   | 1   | 1   | 1   | 1   | ĩ   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| -<br>1<br>1<br>1<br>1<br>1<br>1<br>1                     | 1<br>1<br>1<br>1<br>1<br>1<br>1                               | 1<br>1<br>1<br>1<br>1<br>1                     | 111111111111111111111111111111111111111 | 1<br>1<br>1<br>1<br>1<br>1           | 1<br>1<br>1<br>1<br>1                   | 1 1 1 1 1 1 1              | 1<br>1<br>1<br>1           | 1<br>1<br>1   | 1   
   
  | 1           | 1  
   
   
  | _   |   | 1  
  | 1   | ĩ   | 1   | 1   | 1   | ī.  | 1   | ĩ   | ĩ   | ĩ   | ĩ  
  | ĩ   | 1   | ĩ   | 1   | ī   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1                  | 1<br>1<br>1<br>1<br>1<br>1                                    | 1<br>1<br>1<br>1<br>1                          | 1<br>1<br>1<br>1<br>1                   | 1<br>1<br>1<br>1<br>1                | 1<br>1<br>1<br>1                        | 1<br>1<br>1                | 1<br>1<br>1                | 1<br>1  | 1   
   
  |             |  
   
   
  | 1   | ĩ   | 1  
  | 1   | ĩ   | 1   | 1   | 1   | ĩ   | ĩ   | ĩ   | ĩ   | ĩ   | ĩ  
  | ĩ   | ĩ   | ĩ   | ĩ   | ĩ   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1                  | 1<br>1<br>1<br>1<br>1   | 1<br>1<br>1<br>1                               | 1 1 1 1 1 1                             | 1<br>1<br>1<br>1                     | 1 1 1                                   | 1 1 1                      | 1                          | 1   |   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | ĩ   | 1   | 1   | 1   | ĩ   | 1   | ĩ   | ĩ   | ĩ   | 1  
  | 1   | ĩ   | ĩ   | 1   | ĩ   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1<br>1<br>1<br>1<br>1                                    | 1<br>1<br>1<br>1  | 1<br>1<br>1<br>1                               | 1<br>1<br>1                             | 1<br>1<br>1                          | 1<br>1                                  | 1                          | 1                          |   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1<br>1<br>1<br>1   | 1<br>1<br>1   | 1<br>1<br>1                                    | 1                                       | 1<br>1                               | 1                                       | 1                          |                            | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1<br>1<br>1<br>1   | 1<br>1<br>1   | 1<br>1   | 1                                       | 1                                    |   |                            | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1<br>1<br>1  | 1<br>1  | 1  | 1                                       |                                      | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | з   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1<br>1   | 1   |  | -                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  |   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | з   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
|  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | з   | 3   | з   | з   | 3                  
  | 3   | з   | з   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | з   | 3   | 3   | з   | 3                  
  | 3   | з   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | з   | 3   | 3   | з   | 3                  
  | 3   | з   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
| 1  | 1   | 1  | 1                                       | 1                                    | 1                                       | 1                          | 1                          | 1   | 1   
   
  | 1           | 1  
   
   
  | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1   | 1  
  | 1   | 1   | 1   | 1   | 1   | 1   | 3   | 3   | 3   | 3   | 3                  
  | 3   | 3   | 3   |
|  |   |  |   |                                      |   |                            |                            | 1       1 | 1     1 <td></td> <td>1     1<td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1
      1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1    
  1       1       1       1       1       1       1       1       1       1       1       1       1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td><td>1       1</td></td> |             | 1     1
    1 <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1 
     1       1       1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1
      1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> <td>1       1</td> | 1       1 | 1       1 | 1       1  
    1       1       1       1       1       1       1       1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1      
1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1 | 1       1       1  
    1       1 | 1       1 | 1       1 | 1       1 |

We observe, in Figure 4, that, as  $N_A$  increases, the cost decreases. For values up to the combination  $N_A = 10$  and  $N_B = 25$  the cost is almost constant and then it decreases. That behaviour is determined by the holding costs.

In Figures 6(a), 6(b) and 6(c), we present the long run proportion of time that the system spends in each state for the vehicle mixes (5, 30), (10, 25), (32, 0), respectively. We notice that as ALS ambulances increases, the system spends less time in states with high holding costs.

Iteration	Successive policies	Average cost
1	$r(i, j) = \begin{cases} 0, & otherwise \\ 3, & i = j = 40 \end{cases}$	\$31,616.37
2	$r(i, j) = \begin{cases} 0, & i = j = 0 \\ 3, & i = 0, 1 \le j \le 40 \\ 1, & 1 \le i \le 29, 0 \le j \le 40 \\ 3, & 33 \le i \le 40, 0 \le j \le 40 \\ 3, & i = 30, j = 4,, 6 \\ 1, & i = 30, j \ne 4,, 6 \\ 3, & i = 31, j = 3,, 10 \\ 1, & i = 31, j \ne 3,, 10 \\ 3, & i = 32, j = 2,, 14 \\ 1, & i = 32, j \ne 2,, 14 \end{cases}$	\$23,844.96
3	$r(i, j) = \begin{cases} 0, & i = j = 0\\ 3, & i = 0, 1 \le j \le 40\\ 1, & 1 \le i \le 32, 0 \le j \le 40\\ 3, & 33 \le i \le 40, 0 \le j \le 40 \end{cases}$	\$23,844.31

Table 1The successive policies generated by the policy iteration algorithm and their average<br/>costs for the combination (32, 0)

In Figure 7, we can see the sum of long run proportion of time, for the Markov chain which induced by the optimal policy, in states  $\{(i, j) | 17 \le i \le 40, 30 \le j \le 40\}$  where we have high holding costs.

We notice that in combinations with a few ALS ambulances, we have high sum of long run proportion of time. That results in high holding costs. In combinations with many ALS ambulances, we have the opposite result.

This particular behaviour of the system is a function of two things: the level  $\gamma$  of the redirection rate and the traffic intensities  $\alpha_A = \frac{\lambda_H}{\mu}$  and  $\alpha_B = \frac{\lambda_L}{\mu}$ . As long as  $N_A$  increases

and  $\alpha_A < 1$ ,  $\alpha_B > 1$ , the system services the ALS calls wherever it is possible. Also, the system needs at least ten ALS units in this level of  $\gamma$  to service effectively the remaining ALS calls. As long as  $N_A$  continues to increase and after the combination (10, 25) the long run average cost decreases because the system serves more ALS calls for which  $\alpha_A < 1$ . This seems to be reasonable since as  $N_A$  increases and the traffic intensity  $\alpha_A$  is lower than 1, the system does not collect high holding costs from the serving of ALS calls, it gathers significant holding costs only from BLS calls. On the other hand, due to  $\alpha_B > 1$ , the system collects pricey holding costs when it serves those calls because it gathers notable holding costs from both ALS and BLS calls.

Figure 6 Long run proportion of time in each state for the combinations (5, 30), (10, 25), (32, 0) (see online version for colours)



Figure 7 Long run proportion of time in states with high holding costs (see online version for colours)



#### 6 Sensitivity analysis

In this section, we present a sensitivity analysis for VBEMS in two model parameters: the redirection rate  $\gamma$  and service rate  $\mu$ . We select these parameters because they are the two main options with which an EMS system manages incoming calls. Our goal is to clarify how they affect vehicle mix and the cost of the system.





Figure 9 The effect of  $\mu$  in vehicle mix and in the cost of the system (see online version for colours)



In Figure 8, five curves are depicted for values of  $\gamma$  ranging from 0.192 to 0.770. We observe that for low values of  $\gamma$  the system needs larger number of ALS ambulances to service the ALS calls and for high values of  $\gamma$  the system needs smaller number of ALS ambulances. Also, we observe that the cost is lower in higher level of redirection rate for each vehicle mix. This could be explained by the fact that the more ALS calls the system

redirects, the less ALS ambulances it needs to handle the remaining calls. The option to redirect calls for high values of  $\gamma$  is more preferable than the option to service these calls. Regarding the reduction of cost that implies the increase of  $\gamma$ , we can note that it seems reasonable due to the reduction of holding costs from the redirection of calls.

In Figure 9, five curves for values of  $\mu$  ranging from 0.4 to 4 are depicted. We observe that for  $\mu = 0.4$  we have  $1 < a_A < a_B$  and the minimum cost is appeared in the combination (0, 36). The low price for  $\mu$  pushes the system to redirect the calls and not to serve them. As  $N_A$  increases after the combination (0, 36) the system serves the high priority calls. This increases the holding costs due to  $a_A < 1$ . For  $\mu = 1.2$  and  $\mu = 2$  we have  $a_A < 1 < a_B$ . The minimum cost for  $\mu = 1.2$  is appeared in the combination (32, 0). In this case the system serves all the ALS calls wherever it is possible and redirects the calls in the rest states. For  $\mu \mu = 2$  the system needs 6 ALS units to serve effectively the ALS calls and for  $N_A > 19$  the cost is constant due to very low value of  $a_A$ . This is intuitively reasonable since the high level of service stabilises the cost after a certain number of ambulances. For  $\mu = 3$  and  $\mu = 4$  we have  $a_A < a_B < 1$ . For  $\mu = 3$  the cost increases up to the combination (6, 29). The exchange of service between BLS calls and ALS calls increases the cost up to this point. For  $N_A > 6$ , the system serves effectively the ALS calls, so the cost decreases and becomes constant for  $N_A > 13$ . We reach similar conclusions for  $\mu = 4$ . We also notice that as  $\mu$  increases the cost decreases and it became constant after a specific number of  $N_A$  for large values of  $\mu$ . These observations are compatible with reality.

#### 7 A further research direction and conclusions

The geographical distribution of ambulances in the service area is another important issue regarding the EMS system effectiveness. We assume that an area is divided into N zones and every zone  $i \in N$  has demand  $\lambda_i^H$  for high priority calls and  $\lambda_i^L$  for low priority calls. Also, there is a station of  $\alpha$  ALS ambulances and  $\beta$  BLS ambulances in zone *i*. The call centre of the system can admit the maximum  $\gamma$  high priority calls and  $\delta$  low priority calls from the zone *i* waiting to be served. We also assume that  $\alpha < \gamma$ ,  $\beta < \delta$  and each call is serviced by only one station. In the following figure, an area which is divided into seven zones is depicted. Each zone has one station which services high priority calls, denoted by red dot and low priority calls, denoted by blue dot.

Figure 10 An area which is divided into seven zones (see online version for colours)



Our future pursuit is to determine the values  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$  in each zone for an EMS system which services an area with N zones. We assume that it consists of  $N_A$  ALS ambulances,

 $N_B$  BLS ambulances and can admit a maximum number  $Q_H$  of high-priority calls and a maximum number  $Q_L$  of low-priority calls. We formulate this problem as an integer program as follows.

Whenever there is an incoming high priority call from some zone in the system, we can consider the service procedure as a  $M/M/\alpha/\alpha + \gamma$  queue. There are two cases for a call. The call is either serviced directly by one of  $\alpha$  ALS ambulances, if the queue system is in state k with  $0 \le k \le \alpha - 1$  or the call remains on hold and after some time it is serviced or redirected, if the queue system is in state k with  $\alpha \le k \le \alpha + \gamma$ . In the second case, when there are k calls in the system, there are  $k - \alpha$  calls waiting to be served. The same applies for low priority calls and we can consider the service procedure as a  $M/M/\beta/\beta + \delta$  queue. Considering this structure, we can explore the approach of the system for the calls waiting to be served with a Markovian model with parameters  $Q_H = \gamma$ ,  $Q_L = \delta$ ,  $N_A = \alpha$ ,  $N_B = \beta$ . In next figure, we see an example of this correspondence, for some values of the parameters.





In the Markovian model we represent the proportion of action induced by the optimal policy where the system services the high priority calls with  $\phi_{ALS}$ . We also represent the proportion of action where the system services the low priority calls with  $\phi_{BLS}$ , the proportion of action where the system redirects the high priority calls with  $\theta_{ALS}$  and the proportion of action where the system redirects the low priority calls with  $\theta_{BLS}$ . Finally, we represent the proportion of time when there are *k* customers in the system with  $P_k$  and the holding cost rates for high and low priority calls with *h* and  $\tilde{h}$ , respectively.

So, the expected cost from a high priority call is:

$$\mathcal{C}_{ALS} = \sum_{k=0}^{\alpha-1} P_k \cdot C_H + \sum_{k=\alpha}^{\alpha+y} P_k \left( C_H \cdot \varphi_{ALS} + R_H \cdot \theta_{ALS} + h \right)$$

And the expected cost from a low priority call is:

$$\mathcal{C}_{\text{BLS}} = \sum_{k=0}^{\beta-1} P_k \cdot C_L + \sum_{k=\beta}^{\beta+\delta} P_k \left( C_L \cdot \varphi_{BLS} + R_L \cdot \theta_{BLS} + \tilde{h} \right)$$

For the ambulance fleet and the maximum admitted calls allocation we want to minimise the function:

$$\sum_{\substack{i \in N \\ \gamma = 1}} \lambda_i^H \cdot \sum_{\alpha = 0}^{N_A} \sum_{\beta = 0}^{N_B} \cdot \sum_{\substack{\gamma = 1 \\ \gamma > \alpha}}^{Q_H} \sum_{\substack{\delta = 1 \\ \delta > \beta}}^{Q_L} \cdot y_{i\alpha B\gamma \delta} \cdot \mathcal{C}_{ALS} + \sum_{i \in N} \lambda_i^L \cdot \sum_{\alpha = 0}^{N_A} \sum_{\beta = 0}^{N_B} \cdot \sum_{\beta = 0}^{Q_H} \sum_{\substack{\gamma = 1 \\ \delta > \beta}}^{Q_H} \sum_{\substack{\delta = 1 \\ \delta > \beta}}^{Q_L} \cdot y_{i\alpha B\gamma \delta} \cdot \mathcal{C}_{BLS}$$

The variable  $y_{i\alpha\beta\gamma\delta}$  takes the value 1 if zone *i* uses  $\alpha$  ALS ambulances,  $\beta$  BLS ambulances, the call centre of EMS system can admit a maximum number  $\gamma$  of high priority calls and a maximum number  $\delta$  of low priority calls and takes the value 0, otherwise. The constraints of the integer programming problem can be defined such that the sum of alphas equals  $N_A$ , the sum of betas equals  $N_B$ , the sum of gammas equals  $Q_H$  and the sum of deltas equals  $Q_L$ .

The inspiration for this paper came from the paper of Chong et al. (2015). We researched the EMS systems from the perspective of the calls waiting to be served at any given time. Using a MDP in continuous time, we found the decisions that should be made by the system and the most economical vehicle mix based on its cost data and its features. The main options of the system are the service of calls and their redirection. We examined how changes in the values of these parameters affect the options and costs for the system. Unfortunately, the complexity of the model makes mathematical evidence difficult to predict and to analyse the characteristics of the system. Nevertheless, the application of the model in VBEMS offered results that are characterised as rational. Finally, we gave the integer programming formulation for the geographical distribution of the system's fleet at the stations of a service area which is another important issue for these systems and has a great impact on their efficiency.

#### References

- Akl, R.G., Hegde, M. and Naraghi-Pour, M. (2005) 'Mobility-based CAC algorithm for arbitrary call-arrival rates in CDMA cellular systems', *IEEE Transactions on Vehicular Technology*, Vol. 54, No. 2, pp.639–651.
- Alanis, R., Ingolfsson, A. and Kolfal, B. (2012) 'A Markov chain model for an EMS system with repositioning', *Production and Operations Management*, Vol. 22, No. 1, pp.216–231.
- Al-Shaqsi, S. (2010) 'Model of international emergency medical services (EMS) systems', Oman Medical Journal, Vol. 25, No. 4, pp.320–323.
- Benabdouallah, M. and Bojji, C. (2018) 'A review on coverage models applied to emergency location', *International Journal of Emergency Management*, Vol. 14, No. 2, pp.180–199.
- Berman, O. (1981a) 'Dynamic repositioning of indistinguishable service units on transportation networks', *Transportation Science*, Vol. 15, No. 2, pp.115–136.
- Berman, O. (1981b) 'Repositioning of distinguishable urban service units on networks', *Computers and Operations Research*, Vol. 8, No. 2, pp.105–118.
- Berman, O. (1981c) 'Repositioning of two distinguishable service vehicles on networks', *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. 11, No. 3, pp.187–193.
- Braun, O., McCallion, R. and Fazackerley, J. (1990) 'Characteristics of midsized urban EMS systems', *Annals of Emergency Medicine*, Vol. 19, No. 5, pp.536–546.
- Chong, K.C., Henderson, S.G. and Lewis, M.E. (2015) 'The vehicle mix decision in emergency medical service systems', *Manufacturing & Service Operations Management*, Vol. 18, No. 3, pp.347–360.
- Clawson, J.J. (1989) 'Emergency medical dispatching', in Roush, W.R. (Ed.): *Principles of EMS Systems: A Comprehensive Text for Physicians*, pp.119–133, American College of Emergency Physicians, Dallas.

- Di Bartolomeo, S., Sanson, G., Nardi, G., Michelutto, V. and Scian, F. (2005) 'HEMS vs. ground-BLS care in traumatic cardiac arrest', *Prehospital Emergency Care*, Vol. 9, No. 1, pp.79–84.
- Dick, W. (2003) 'Anglo-American vs. Franco-German emergency medical services system', *Prehospital and Disaster Medicine*, Vol. 18, No. 1, pp.29–35.
- DuBois, E. and Albert, L.A. (2021) 'Dispatching policies during prolonged mass casualty incidents', Journal of the Operational Research Society, pp.1–15, Taylor & Francis, https://doi.org/10.1080/01605682.2021.1999181.
- Eisen, J.S. and Dubinsky, I. (1998) 'Advanced life support vs. basic life support field care: an outcome study', *Academic Emergency Medicine*, Vol. 5, No. 6, pp.592–598.
- Heyman, D.P. and Sobel, M.J. (2003) Stochastic Models in Operations Research, Stochastic Optimization, Vol. II, Dover, New York
- Jarvis, J.P. (1975) Optimization in Stochastic Service Systems with Distinguishable Servers, PhD thesis, Massachusetts Institute of Technology, Cambridge, MA
- Ji, S., Zheng, Y., Wang, W. and Li, T. (2020) 'Real-time ambulance redeployment: a data-driven approach', *IEEE Transactions on Knowledge and Data Engineering*, Vol. 32, No. 11, pp.2213–2226.
- Lee, H-R. and Lee, T. (2018) 'Markov decision process model for patient admission decision at an emergency department under a surge demand', *Flexible Services and Manufacturing Journal*, Vol. 30, No. 1, pp.98–122, Springer.
- McLay, L.A. and Mayorga, M.E. (2012) 'An optimal dispatching model for server-to-customer systems with classification errors', *IIE Transactions*, Vol. 45, No. 1, pp.1–24.
- Mengyu, L., Carter, A., Goldstein, J., Hawco, T., Jensen, J. and Vanberkel, P. (2020) 'Determining ambulance destinations when facing offload delays using a Markov decision process, *Omega*, Article 102251, Vol. 101, No. 3.
- Nasrollahzadeh, A.A., Khademi, A. and Mayorga, M.E. (2018) 'Real-time ambulance dispatching and relocation', *Manufacturing & Service Operations Management*, Vol. 20, No. 3, pp.467–480.
- Nguyen-Van-Tam, J.S., Dove, A.F., Bradley, M.P., Pearson, J.C., Durston, P. and Madeley, R. (1997) 'Effectiveness of ambulance paramedics versus ambulance technicians in managing out of hospital cardiac arrest', *Journal of Accident and Emergency Medicine*, Vol. 14, No. 3, pp.142–148.
- Ornato, J.P., Racht, E.M., Fitch, J.J. and Berry, J.F. (1990) The need for ALS in urban and suburban EMS systems, *Annals of Emergency Medicine*, Vol. 19, No. 12, pp.1469–1470.
- Park, S.H. and Lee, Y.H. (2019) 'Two-tiered ambulance dispatch and redeployment considering patient severity classification errors', *Journal of Healthcare Engineering*, Vol. 2019, Article ID: 6031789, 14p.
- Pouraliakbari, M., Mohammad, M. and Mirzazadeh, A. (2018) 'Analysis of maximal covering location-allocation model for congested healthcare systems in user choice environment', *International Journal of Industrial and Systems Engineering (IJISE)*, Vol. 28, No. 2, pp.240–274.
- Puterman, M.L. (1994) Markov Decision Processes: Discrete Stochastic Dynamic Programming, Wiley, New York.
- Rainer, T.H., Houlihan, K.P., Robertson, C.E., Beard, D., Henry, J.M. and Gordon, M.W. (1997) 'An evaluation of paramedic activities in prehospital trauma care', *Injury*, Vol. 28, Nos. 9–10, pp.623–627.
- Ramirez-Nafarrate, A., Gel, E.S., Fowler, J.W. and Hafizoglu, A.B. (2014) 'Optimal control policies for ambulance diversion', *European Journal of Operational Research*, Vol. 236, No. 1, pp.298–312.

- Slovis, C.M., Carruth, T.B., Seitz, W.J., Thomas, C.M. and Elsea, W.R. (1985) 'A priority dispatch system for emergency medical services', *Annals of Emergency Systems*, Vol. 14, No. 11, pp.1055–1060.
- Stiell, I.G., Spaite, D.W., Field, B., Nesbitt, L.P., Munkley, D., Maloney, J., Dreyer, J., Toohey, L.L., Campeau, T., Dagnone, E., Lyver, M. and Wells, G.A. (2007) 'OPALS study group: advanced life support for out-of-hospital respiratory distress', *New England Journal of Medicine*, Vol. 356, No. 21, pp.2156–2164.
- Stiell, I.G., Wells, G.A., Field, B., Spaite, D.W., Nesbitt, L.P., De Maio, V.J., Nichol, G., Cousineau, D., Blackburn, J., Munkley, D., Luinstra-Toohey, L., Campeau, T., Dagnone, E. and Lyver, M. (2004) 'Ontario prehospital advanced life support study group: advanced cardiac life support in out-of-hospital cardiac arrest', *New England Journal of Medicine*, Vol. 351, No. 7, pp.647–656.
- Stout, J., Pepe, P.E. and Mosesso, V.N. (2000) 'All-advanced life support vs tiered-response ambulance systems', *Prehospital Emergency Care*, Vol. 4, No. 1, pp.1–6.
- Swersey, A.J. (1982) 'A Markovian decision model for deciding how many fire companies to dispatch', *Management Science*, Vol. 28, No. 4, pp.352–365.
- Tijms, H.C. (2003) A First Course in Stochastic Models, Wiley, Chichester, New York.
- Van Den Berg, P.L. and Van Essen, T. (2019) 'Comparison of static ambulance location models', International Journal of Logistics Systems and Management, Vol. 32, Nos. 3/4, pp.292–321.
- Wilson, B., Gratton, M.C., Overton, J. and Watson, W.A. (1992) 'Unexpected ALS procedures on non-emergency ambulance calls: the value of a single-tier system', *Prehospital and Disaster Medicine*, Vol. 7, No. 4, pp.380–382.
- Zhang, L. (2012) Simulation Optimization and Markov Models for Dynamic Ambulance Redeployment, PhD thesis, The University of Auckland, New Zealand.