

ASYMPTOTIC FLUID OPTIMALITY AND EFFICIENCY OF THE TRACKING POLICY FOR BANDWIDTH-SHARING NETWORKS

KONSTANTIN AVRACHENKOV,* *INRIA*
ALEXEY PIUNOVSKIY ** *** AND
YI ZHANG,** **** *University of Liverpool*

Abstract

Optimal control of stochastic bandwidth-sharing networks is typically difficult. In order to facilitate the analysis, deterministic analogues of stochastic bandwidth-sharing networks, the so-called fluid models, are often taken for analysis, as their optimal control can be found more easily. The tracking policy translates the fluid optimal control policy back to a control policy for the stochastic model, so that the fluid optimality can be achieved asymptotically when the stochastic model is scaled properly. In this work we study the efficiency of the tracking policy, that is, how fast the fluid optimality can be achieved in the stochastic model with respect to the scaling parameter. In particular, our result shows that, under certain conditions, the tracking policy can be as efficient as feedback policies.

Keywords: Bandwidth-sharing network; fluid model; optimal control; tracking policy; rate of convergence

2010 Mathematics Subject Classification: Primary 60K25; 68M20

1. Introduction

Let us start this article with the following adopted conventions.

- By \mathbb{R}^L , \mathbb{R}_+^L , and \mathbb{Z}_+^L we respectively indicate the sets of L -vectors of real, nonnegative real, and nonnegative integer numbers. By $D[0, \infty)$ we denote the space of right-continuous left-limit functions on $[0, \infty)$. By $\{e_l\}$, e , \cdot^T , $\mathbf{1}\{\cdot\}$ we respectively denote the natural basis of \mathbb{R}^L , the element in \mathbb{R}^L with all its components equal to 1, the transposition, and the indicator function. For all $x, y \in \mathbb{R}$, $x \wedge y := \min\{x, y\}$. By $O(1/n)$ and $O(1/\sqrt{n})$ we respectively mean that there exists some constant C such that $O(1/n)/(1/n) \rightarrow C$ and $O(1/\sqrt{n})/(1/\sqrt{n}) \rightarrow C$ as $n \rightarrow \infty$.
- Unless stated otherwise, by a vector we always mean a column vector.
- When we say that $X_1 \leq X_2$, where X_1 and X_2 are vectors, ' \leq ' is understood component-wise, and likewise for other vector inequalities.

A bandwidth-sharing network can be described as follows. Through a network of J resources (links), L flows are routed in a predetermined way. The L flows are assumed to be distinct in

Received 21 December 2009; revision received 27 November 2010.

* Postal address: INRIA, MAESTRO Team, 2004 Route des Lucioles - BP 93 FR-06902 Sophia Antipolis Cedex, France. Email address: k.avrachenkov@sophia.inria.fr

** Postal address: Department of Mathematical Sciences, University of Liverpool, Liverpool L69 7ZL, UK.

*** Email address: piunov@liverpool.ac.uk

**** Email address: zy1985@liverpool.ac.uk

the sense that each one is routed differently. The flows represent aggregate streams of finitely sized files sent according to Poisson processes. Thus, the files are classified according to which flow they belong to. Files belonging to flow l are called files of type l . Each resource has a finite capacity shared among the flows passing through it. Let us introduce some notation. A network is described by a configuration (A, Z) , where A is a $(J \times L)$ -matrix with $A_{jl} = 1$ if resource j participates in serving files of type l and $A_{jl} = 0$ otherwise. We assume that each column of A has at least one nonzero component, meaning that each file (flow) must be served somewhere. Let Z be a J -vector with $Z_j > 0$ indicating the maximal capacity of resource j . If the current number of files of each type is $Y \in \mathbb{Z}_+^L$ and the current time is $t \geq 0$, then let $U(Y, t)$, an L -vector satisfying $0 \leq AU(Y, t) \leq Z$, represent the instantaneous (feasible) allocations of resources for each type of files. Here, we are restricted to the class of deterministic Markov control policies, so that $U(Y, t)$ will be assumed to be a mapping measurable in $t \geq 0$. For brevity, below, when we say a control policy (for the stochastic model), we shall often omit its argument, but simply indicate U .

We shall model the network dynamics as a Markovian system. To the best of the authors' knowledge, this bandwidth-sharing network model originated in [15] and [25]. It has subsequently been extensively studied. An interested reader can find a thorough literature review about the model in [27].

Quite formally, given a fixed U , our stochastic model is a regular Q -process $\{Y_t, t \geq 0\}$ in the state space \mathbb{Z}_+^L with the conservative and stable Q -matrix given by $q_{Y, Y+e_l}(t) = \lambda_l$ and $q_{Y, Y-e_l}(t) = \mu_l U_l(Y, t) \mathbf{1}\{Y^l > 0\}$, where Y^l is the l th component of Y . See Chapter 2 and Appendix B of [11] and [29] for more details. Here, we have decided not to indicate any of the nonpositive 'diagonal' components and the other zero components of the Q -matrix: this does not generate any confusion, with the conservativeness of the Q -matrix in mind. Note also that the dependence of the Q -matrix on U will not be signified for simplicity. Put differently, in our model we assume that the files of flow l arrive according to a Poisson process and have exponentially distributed service times. In greater detail, denoting by $Y_t \in \mathbb{Z}_+^L$ the instantaneous number of files of each type in the network, we suppose that $\{Y_t, t \geq 0\}$ is a continuous-time Markov chain with λ_l the transition rate from Y_t to $Y_t + e_l$ and $\mu_l U_l$ the transition rate from Y_t to $Y_t - e_l$ (if the corresponding component of Y_t is positive). This is equivalent to saying that files of type l arrive in a Poisson process with intensity λ_l , and each file is of an exponentially distributed size with mean $1/\mu_l$; hence, the mean time for its service to be completed will be the ratio of U_l , its allocated capacity, against $1/\mu_l$. Here, we have assumed that all the exogenous arrivals of files are independent. Furthermore, if the fixed U means that $U_l(Y, t)$ can be written as $U_l(Y^l, t)$ for all $l = 1, \dots, L, t \geq 0$, and $Y \in \mathbb{Z}_+^L$, i.e. the allocation of resources for files of type l depends only on the current time as well as the number of files of type l , then it will be our standing assumption that, under this U , $\{Y_t^l, t \geq 0\}$ are independent continuous-time Markov chains with transition rates λ_l and $U_l(Y_t^l, t)\mu_l$, where Y_t^l is the l th component of Y_t , i.e. each file is served independently, and the presences of different types of files at a common resource do not interfere with each other. Below, we shall frequently refer to this standing assumption of independence.

For such a system, the performance measure of our interest is the expected total holding cost, and the optimization problem of our concern is in the form of

$$E_{Y_0} \left[\int_0^T e^{-\tau} Y_\tau \, d\tau \right] \rightarrow \min_U$$

where E_{Y_0} is the expectation with the initial state Y_0 (fixed and deterministic) and $T > 0$ is a given finite horizon. Note that, without leading to confusion, we have decided not to signify the dependence of E_{Y_0} on U . In words, we aim at finding the optimal resource allocation to minimize this performance measure. The optimization with this criterion can be interpreted as the minimization of the total workload. More discussions about this criterion can be found in [27].

On the other hand, we can consider a deterministic analogue of the above stochastic (Markovian) bandwidth-sharing network model, its so-called fluid model, which we now describe. Let $y(t)$ and $u(y(t), t)$ be the analogues of Y_t and $U(Y_t, t)$. In fact, below we shall only write $u(t)$ for the fluid control policy, which, in such a deterministic system, is a plausible thing to do. If we further agree on the notation $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_L)^\top$ and

$$M = \text{diag}(\mu_1, \mu_2, \dots, \mu_L) := \begin{pmatrix} \mu_1 & 0 & \cdots & 0 \\ 0 & \mu_2 & \cdots & 0 \\ 0 & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \mu_L \end{pmatrix},$$

then the fluid model can be written as the following linear program:

$$\int_0^T e^\top y(t) dt \rightarrow \min_{u,y} \tag{1}$$

such that $\frac{dy}{dt} = -Mu(t) + \Lambda, \quad y(0) = Y_0, \quad Au(t) \leq Z, \quad y(t), u(t) \geq 0.$

We emphasize that in the fluid model the state space is \mathbb{R}_+^L instead of \mathbb{Z}_+^L .

In spite of some loss of accuracy compared to its corresponding stochastic counterpart, often a fluid model is more amenable for analysis. For example, Verloop and Núñez-Queija [28] (see also [27, Chapter 5]), for a certain set of parameters, found the optimal solution to the above fluid problem for a bandwidth-sharing network with two resources and three flows. Once an optimal policy is obtained for the fluid model, one may translate it into a policy for the stochastic model. It is of interest to see how well the resulting policy performs. In general, directly applying the fluid optimal policy to the stochastic model may be far from optimal. Here, by directly we mean without scaling the stochastic model. To illustrate this point, let us consider the following example.

Example 1. Suppose that we are concerned with a controlled M/M/1 queueing system, where the initial state is $Y_0 = 1$, $q_{Y_t, Y_{t+1}}(t) = \lambda > 0$ (if $Y_t \in \mathbb{Z}_+$), and $q_{Y_t, Y_{t-1}}(t) = \mu - \pi(Y_t)$ (if $Y_t \in \mathbb{Z}_+ \setminus \{0\}$) with $\mu > 0$ and $\lambda > \mu$ (by a sufficiently small difference), where $\pi(\cdot)$ is a mapping from the state space to $\{0, \mu\}$. In words, this means that we are restricted to the class of deterministic stationary policies to choose actions from $\{0, \mu\}$. We shall fix the absorbing state to be some large enough integer, so that the state space is indeed finite. Suppose that the cost rate is $c(x, a) = C \mathbf{1}\{x \in [0, \frac{1}{5}]\} + a \mathbf{1}\{x \in [\frac{1}{2}, 1]\}$, where $C > 0$ is a sufficiently large (in relation to μ) penalty constant. We are concerned with the problem of

$$E_1^\pi \left[\int_0^\infty c(Y_t, \pi(Y_t)) dt \right] \rightarrow \min_\pi,$$

which is well defined for any $\pi(\cdot)$. Clearly, since C is much larger than μ , which is then rather close to λ , the optimal control at state 1 should be $\pi(1) = \mu$, meaning that the value of its

performance functional will be strictly positive. Indeed, on the one hand, if we set $\pi(1) = 0$ then the process visits state 0 after the first jump at a significant probability, but visiting state 0 (and staying there for a while) is heavily penalized. In this case, the expected total cost is more than $(\mu/(\lambda + \mu))C/\lambda$, with $\mu/(\lambda + \mu)$ representing the probability that the first jump is downward. On the other hand, if we set $\pi(x) = \mu, x \geq 1$ (so, in particular, $\pi(1) = \mu$), then the process definitely jumps upward and never visits state 0. In this case, the total cost will be given by μ/λ (i.e. the expected cost incurred up to the first jump), which is close to 1 since μ is close to λ . This quantity is smaller than $(\mu/(\lambda + \mu))C/\lambda$ for sufficiently large C , which is the case here. Hence, the optimal control at state 1 should be given by $\pi(1) = \mu$.

Now consider its fluid model. Since $y(0) = 1$ and $\lambda - \mu > 0, \pi(y) = 0$ is optimal, resulting in $\int_0^\infty c(y, \pi(y)) dt = 0$, y is read as a function of t , which we decide not to indicate explicitly for simplicity. Therefore, we see that in this example both the optimal policies and the values of performance functionals differ significantly between the stochastic model and its fluid model. Another observation of this type can be found in [19, Section 5, p. 427], where the authors considered an epidemic model proposed in [10].

Intuitively, the reason for the discrepancy in the above example lies in the stochastic and discrete nature of the stochastic model versus the deterministic and continuous nature of the fluid model. However, when the stochastic process is scaled properly, in the limiting case, by referring to the ‘functional strong law of large numbers’, the randomness gets eliminated (or alleviated, at least). See [5], [6], and [14] for more details. Let us call the underlying scaling ‘fluid scaling’. Then it is natural to consider the case when the fluid optimal policy is translated to a scaled stochastic model and to examine its limiting case with respect to the fluid scaling parameter. This gives rise to the concept of asymptotic fluid optimality. See [8] and [9]. Basically, a policy is said to be asymptotically fluid optimal (AFO) if applying it to a sequence of scaled stochastic models results in a sequence of performance functionals converging to that optimal for the fluid model. The question then becomes one of asking whether or not naively translating the fluid optimal policy can provide an AFO policy. In fact, generally speaking, translations that look natural can result in policies not AFO, according to [13, Section 2]. See also [20, Example 1] for another example where the fluid model provides inaccurate approximations even in the sense of fluid scaling.

Clearly, to define the asymptotic fluid optimality more accurately, we firstly need to define the fluid scaling. The general idea is to scale the parameters (rates) and to aggregate the space both linearly, while the time scale is kept unchanged: given the (fluid) scaling parameter $n \in \mathbb{N}$ fixed, we consider $\{^n Y_t, t \geq 0\}$ with $^n Y_0 := nY_0$ a continuous-time Markov chain in the state space \mathbb{Z}_+^L with parameters $n\lambda_l$ and $n\mu_l U_l$, and performance measure of the form

$$^n \hat{W}(^n Y_0) := E_{^n Y_0} \left[\int_0^T e^{-\frac{c}{n} Y_t} dt \right].$$

Here, we recall that Y_0 is a constant that has appeared above, and U stands for the control policy. It should also be emphasized that in the scaled stochastic model the parameter n , as the denominator in the above expression, aggregates the space, and, as the multiplicative factor, scales up the transition rates. Then a policy is AFO if it results in the convergence

$$\lim_{n \rightarrow \infty} |^n \hat{W}(^n Y_0) - \hat{w}(Y_0)| = 0,$$

where

$$\hat{w}(Y_0) := \int_0^T e^{-c y(t)} dt, \quad y(0) = Y_0,$$

is the performance functional for the fluid model under its optimal control policy. As mentioned earlier, we have dropped the indication of the control policy, because below we shall always fix a control policy, and, thus, deal with essentially uncontrolled processes. As the initial state ${}^n Y_0$ is fixed, from now on we shall even only write ${}^n \hat{W}$ and \hat{w} instead of ${}^n \hat{W}({}^n Y_0)$ and $\hat{w}(Y_0)$ for simplicity. Furthermore, we shall call $|{}^n \hat{W} - \hat{w}|$ the actual accuracy of the fluid model.

Now we can describe the translation of the fluid optimal policy. Potentially, there are several options. One possibility is a feedback-type translation, which was studied in [18] and [20] for controlled birth-and-death processes, in [17] for a multiclass single-station system, and in [8] and [9] for tandem queues, to mention some. Another possibility is the tracking policy translation, considered in [2], which is also the object of the current work. The tracking policy translation is a naive translation, because, for the scaled stochastic model, the same action as for the fluid model will be taken on the same time interval, unless the state is null in the stochastic model, on which occasion we do not allocate any capacity. More exactly, given the fluid optimal control u^* , the tracking policy U^* is defined via

$$U_l^*({}^n Y_t, t) = U_l^*({}^n Y_t^l, t) = u_l^*(t) \mathbf{1}\{{}^n Y_t^l > 0\}, \quad l = 1, 2, \dots, L, \quad (2)$$

where we recall that u_l^* , U_l^* , and ${}^n Y_t^l$ are the l th components of u^* , U^* , and ${}^n Y_t$, respectively. In [2], the author considered a scheduling problem for a multiclass queueing network. There, for a discount model, the tracking policy was proved to be AFO, which was done by showing the convergence of performance functionals. However, the author revealed no information about the rate of convergence. Given the preliminary belief that the tracking policy translation also results in an AFO policy in our case, which is indeed the case as shown below, it is natural to ask the following question.

- How fast can we achieve the fluid optimality? Or, in other words, how efficient is the tracking policy?

This question is important, as the efficiency helps compare different AFO translations of fluid optimal policies, and at the same time provides the accuracy of the fluid approximations.

In the present work, for the concerned finite-horizon problem, we provide two answers to this question with upper boundary estimates for the rate of convergence of performance functionals. Specifically, we show that the tracking policy can be as efficient as or less efficient than the feedback-type translation, depending on the parameters of the underlying fluid model. While one example in [9] showed that the tracking policy was less efficient than the feedback-type translation considered there, we shall provide an opposite example to show that the tracking policy can also be efficient, and, hence, favoured due to the less information it requires to be implemented. The result obtained for the efficiency implies that the tracking policy is AFO. In addition, we discuss a modification of the tracking policy towards the improved efficiency. To the best of the authors' knowledge, the question about the efficiency of translations of fluid optimal policy has not yet been studied intensively. Relevant works include [9], [17], [18], [20], and [21], all of which focus on feedback-type translations, with the exception of [9], which mainly focuses on feedback-type translations, but also has a short discussion on the tracking policy.

The rest of the work is organized as follows. In Section 2 we state the main results, which are verified with two examples in Section 3. Finally, we sum up this work with a conclusion. The proofs of the main statements, alongside some auxiliary lemmas, are given in Appendix A.

2. Main statements

In this section we shall provide two answers to the question about the efficiency of the tracking policy raised at the end of Section 1. The first answer given in Theorem 1 below corresponds to the general case, where we do not impose extra assumptions on the parameters of the network; the second answer given in Theorem 2 below looks more interesting, but is based on some extra assumptions.

Let us start with a lemma, giving the form of the optimal control policy for the fluid model, following which the main notation used in this section, in addition to those given in Section 1, can be introduced.

Lemma 1. *There is an optimal policy $u^*(t)$ for fluid model (1) in the form of a piecewise-constant function in time $t \geq 0$ with N subintervals $[T_{i-1}, T_i]$, where $i = 1, \dots, N < \infty$, with $T_0 := 0$ and $T_N := T$.*

Owing to the form of the tracking policy U^* (see (2)) and our standing assumption about independence mentioned in Section 1, for all $l = 1, 2, \dots, L$, we refer to files of type l as the l th (Markovian) ‘queue’. Let ${}^i\mu_l := \mu_l u_l^*(t)$ on $[T_{i-1}, T_i]$ for all $l = 1, 2, \dots, L$ and $i = 1, 2, \dots, N$ be the potential service rate of the l th queue on the i th subinterval.

2.1. Efficiency of the tracking policy for the general case

Now we are in the position to state the following theorem.

Theorem 1. *Let us set $\gamma := \max_{i=1, \dots, N, l=1, \dots, L} \{\lambda_l, {}^i\mu_l\}$ and $\bar{T} := \max_{i=1, \dots, N} \{T_i - T_{i-1}\}$. For any fixed initial state and network configuration, the above quantities are fixed. Then*

$$|{}^n\hat{W} - \hat{w}| \leq \frac{TL}{\sqrt{n}} \{12\sqrt{N(N+1)\gamma\bar{T}}\}.$$

In particular, $\lim_{n \rightarrow \infty} |{}^n\hat{W} - \hat{w}| = 0$.

It will be confirmed by an example in Section 3 that, as far as the order is concerned, the obtained $O(1/\sqrt{n})$ is correct and cannot be improved. Given the uniformity in the values of parameters, Theorem 1 does not require extra assumptions on the primary data. Then, with Theorem 1, we say that the efficiency of the tracking policy is $O(1/\sqrt{n})$. It is then interesting to compare the efficiency of the tracking policy with that of the feedback-type translation, as considered in, for example, [17], [18], and [20]. There, for an (absorbing) birth-and-death process (or, say, M/M/1 queue), the author(s) considered a feedback-type translation, which was shown to be of efficiency $O(1/n)$. If we consider a corresponding bandwidth-sharing network with only one resource and flow, then Theorem 1 somehow suggests that the tracking policy is in general less efficient than feedback-type translations. This is in line with the result in [9], where, for a different model, the authors also observed that the tracking policy could be less efficient than feedback-type translations.

On the other hand, suppose that we still consider the simple case of the M/M/1 queue, i.e. the bandwidth-sharing network is of one resource and one flow. In addition, we impose some assumptions on the parameters so that, for the given $Y_0 = y(0) > 0$, the time horizon T is sufficiently small. Then, with the absorbing case considered in [20] in mind, instead of $O(1/\sqrt{n})$, we might expect $|{}^n\hat{W} - \hat{w}|$ to converge as fast as $O(1/n)$, because we now count the deviation of the stochastic model from the fluid model for less time. This motivates our study of a preabsorbing case in the next subsection, where the time horizon is small enough so that over it, for each type of file, the fluid model does not reach 0.

2.2. Efficiency of the tracking policy for a preabsorbing case

Assumption 1. Define $\underline{\lambda} := \min_{l=1, \dots, L} \{\lambda_l\}$, $\bar{\lambda} := \max_{l=1, \dots, L} \{\lambda_l\}$, and $\bar{\mu} := (\max_{l=1, \dots, L} \{\mu_l\})$ ($\max_{l=1, \dots, L} \{Z_l\}$), where Z_l is the l th component of Z . Then $T > 0$, the time horizon, satisfies the condition that there exists a $\underline{y} > 0$ such that, for all $l = 1, \dots, L$, $Y_0^l - |\underline{\lambda} - \bar{\mu}|T \geq \underline{y} > 0$. Here we emphasize that the above inequality holds uniformly in $l = 1, \dots, L$, \underline{y} is a (deterministic) constant, and we recall that Y_0^l is the l th component of Y_0 , which is then a fixed (deterministic) constant (see Section 1).

With the fixed $T > 0$, satisfying Assumption 1 or not, we can always fix two (deterministic) constants $\bar{y} > 0$ and $K > 0$ such that $\max_{l=1, \dots, L} \{Y_0^l + \bar{\lambda}T\} \leq \bar{y} < K$. In fact, below K , as can always be fixed, will be regarded as a sufficiently large constant. In words, Assumption 1 means that $T > 0$ is so small that over the time horizon, the number of files of any type never reaches either 0 or K in the fluid model. If we view state 0 as ‘absorbing’ then, naturally, we call a bandwidth-sharing network model a ‘preabsorbing’ case, if Assumption 1 is satisfied.

Theorem 2. Under Assumption 1,

$$|^n \hat{W} - \hat{w}| \leq O\left(\frac{1}{n}\right).$$

Here, the exact expression of the upper boundary estimate ($O(1/n)$) can be obtained from the primary data by scanning the proof of this theorem.

Theorem 2 implies that the tracking policy can be as efficient as some feedback-type translations, at least for a certain set of parameters. While one way of understanding the better efficiency of the tracking policy in the preabsorbing case has been explained at the end of the last subsection (after Theorem 1), we shall give another intuitive explanation in Section 3.

3. Examples

Let us illustrate the obtained results via two examples in this section.

Example 2. (The simplest example.) Let us set $L = J = 1$, $A = [1]$, $Z_1 = 1$, and ${}^n Y_0^1 = n y_1(0) = 0$. Fixing a constant $M > 0$, we suppose that $\mu_1 = 2M$ and $\lambda_1 = M$. Suppose also that some $T > 0$ is fixed. Since $L = 1$, throughout this example, we shall omit the index standing for the type of file under consideration.

For the fluid model (1), clearly, $u^*(t) = \frac{1}{2}$ (on $[0, T)$) is optimal, because, under it,

$$\int_0^T y(t) dt = 0.$$

For the (scaled) stochastic model, under the tracking policy $U^*({}^n Y(t), t) = \frac{1}{2} \mathbf{1}\{{}^n Y_t > 0\}$ (see (2)), we effectively deal with an M/M/1 queue with the arrival and (potential) service rates both equal to nM . Then

$$|^n \hat{W} - \hat{w}| = E_0 \left[\int_0^T \frac{{}^n Y_t}{n} dt \right] = \frac{1}{n} \int_0^T E_0[{}^n Y_t] dt,$$

where $E_0[{}^n Y_t]$ can be computed as follows. For a small enough time increment h , we have, omitting the subscript for the initial position,

$$E[{}^n Y_{t+h} \mid {}^n Y_t > 0] = ({}^n Y_t + 1)(nhM + o(h)) + ({}^n Y_t - 1)(nhM + o(h)) + {}^n Y_t(1 - nhM - nhM + o(h))$$

and

$$E[{}^n Y_{t+h} \mid {}^n Y_t = 0] = nhM + o(h),$$

where $o(h)$ stands for a term of order lower than h in the sense that $\lim_{h \rightarrow 0} o(h)/h = 0$. This leads to

$$\begin{aligned} & E[{}^n Y_{t+h} \mid {}^n Y_t] = {}^n Y_t + nhM \mathbf{1}\{{}^n Y_t = 0\} + o(h) \\ \Rightarrow & E[{}^n Y_{t+h}] = E[{}^n Y_t] + nhM P\{{}^n Y_t = 0\} + o(h) \\ \Rightarrow & \frac{E[{}^n Y_{t+h}] - E[{}^n Y_t]}{h} = nM P\{{}^n Y_t = 0\} + \frac{o(h)}{h} \\ \Rightarrow & E[{}^n Y_t] = \int_0^t nM P\{{}^n Y_s = 0\} ds, \end{aligned}$$

where the last step is a result of first taking $h \rightarrow 0$ and then integrating. Therefore, we have

$$\begin{aligned} |{}^n \hat{W} - \hat{w}| &= E \left[\int_0^T \frac{{}^n Y_t}{n} dt \right] \\ &= \frac{1}{n} \int_0^T \int_0^t nM P\{{}^n Y_s = 0\} ds dt \\ &= M \int_0^T \int_0^t e^{-2nMs} \{I_0(2nMs) + I_1(2nMs)\} ds dt, \end{aligned} \tag{3}$$

where, in accordance with [3, Equation (1)] (see also [1, pp. 110–112]), $P\{{}^n Y_s = 0\} = e^{-2nMs} \{I_0(2nMs) + I_1(2nMs)\}$ with $I_0(2nMs)$ and $I_1(2nMs)$ standing for the modified Bessel functions of the first kind of orders 0 and 1 at $2nMs$, respectively.

For the convenience of numerical evaluations, let us set $M = T = 1$ (for simplicity), and recall that $\bar{T} = \gamma = N = 1$. Then the difference between our upper boundary estimate and the actual accuracy given by (3) is plotted in the Figure 1, and the ratio of our estimate against the actual accuracy is given in Figure 2, both against the scaling parameter. From the figures

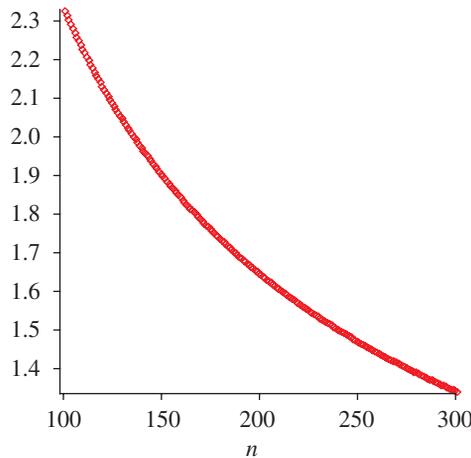


FIGURE 1: The difference between the actual accuracy and the estimated accuracy for the example of the M/M/1 queue. The vertical axis gives the difference, and the horizontal axis gives the scaling parameter, n , ranging from 100 to 300.

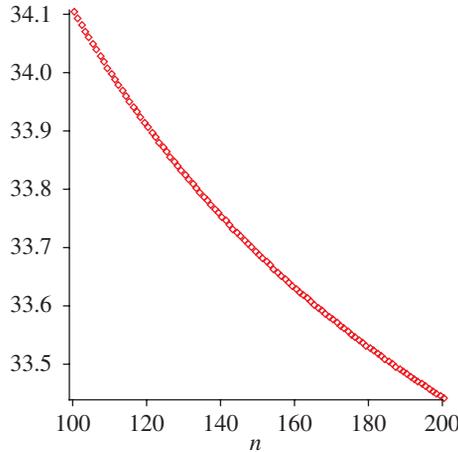


FIGURE 2: The ratio of the estimated accuracy against the actual accuracy for the example of the M/M/1 queue. The vertical axis gives the ratio, and the horizontal axis gives the scaling parameter, n , ranging from 100 to 200.

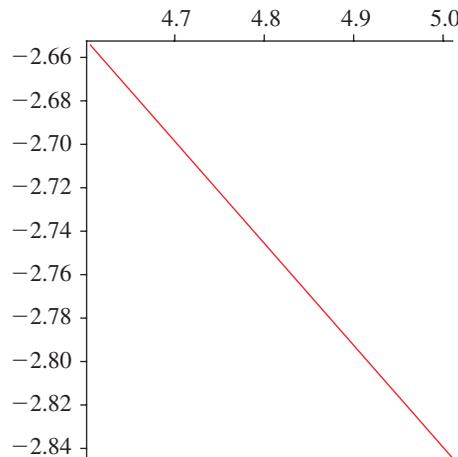


FIGURE 3: A log-scale plot of the actual accuracy for the example of the M/M/1 queue. The vertical axis gives the log of the actual accuracy, and the horizontal axis gives $\log(n)$, with n ranging from 100 to 150.

Note that the curve is very close to a straight line with slope of about $-\frac{1}{2}$.

we see that our estimate is a rather rough estimate. In particular, Figure 2 suggests that the convergence happens around thirty times faster than estimated. In fact, when $n = 10\,000$, the ratio of our estimate against the actual accuracy is about 32, and, when $n = 100\,000\,000$, the ratio is also around 32. On the other hand, Figure 3 shows that the actual rate of convergence is of order $O(1/\sqrt{n})$, which is exactly our estimate. This says, given the uniformity in parameters as in Theorem 1, that we do not have convergence faster than $O(1/\sqrt{n})$.

Let us mention three reasons for the simple settings of Example 2.

The first reason is to have a simple expression for $P\{^n Y_s = 0\}$ and, thus, for the actual rate of convergence (accuracy). In fact, according to [3], for an arbitrary initial state $n_0 \in \mathbb{Z}_+$,

denoting by $\rho := \lambda/\mu$ the traffic intensity for the standard M/M/1 queue, we have

$$P\{^n Y_s = 0\} = e^{-(1+\rho)n\mu s} \left\{ \rho^{-n_0/2} I_{-n_0}(2\rho^{1/2}n\mu t) + \rho^{(-n_0-1)/2} I_{n_0+1}(2\rho^{1/2}n\mu s) + (1 - \rho)\rho^n \sum_{l=n_0+2}^{\infty} \rho^{-l/2} I_l(2\rho^{1/2}n\mu t) \right\},$$

which is difficult to evaluate. Other alternative formulae, though available (see, e.g. [26, Section 1]), are also rather complicated. In addition, when there are at least two subintervals, we deal with a time-dependent M/M/1 queue, for which the transient probability is provided in [30], as a solution to some Volterra-type integral equation, making it very difficult to evaluate also. After all, even in our simplified setting, the obtained expression for $E[\int_0^T ({}^n Y_t/n) dt]$ is still not of a very simple form. However, at least it is easy to compute its numerical values.

The second reason is that this simple setup itself is interesting and typical. Suppose now that we modify the primitives of the above example by letting ${}^n Y_0 = nY_0 > 0$ and $T = Y_0/M$, while we keep all the other primitives unchanged. Then, for the fluid model, obviously, u^* given by $u^*(t) = {}^1\mu = 2M$ on $[0, T/2)$ and $u^*(t) = {}^1\mu = M$ on $[T/2, T)$ is optimal, because, under it, the fluid model decreases at the fastest rate to 0 at $T/2$ and stays at 0 from $T/2$ to the end of the horizon. Therefore, on $[T/2, T)$, we have exactly the situation as in the above example. On the other hand, under the tracking policy U^* , the (scaled) stochastic model can be viewed as a time-dependent M/M/1 queue. As we increase the scaling parameter n , because the trajectory converges (see [4]–[7] and [14]) at $T/2$, we are likely to end up with the stochastic model starting with initial state close to 0 as well as a unit traffic intensity, which, as we have seen in the above example, will result in the rate of convergence $O(1/\sqrt{n})$. This also explains another piece of intuition for avoiding the fluid model reaching 0 in the preabsorbing case (see Subsection 2.2).

Thirdly, in Example 2 if we modify the tracking policy via $U_{\text{modified}}^*({}^n Y_t, t) = \mathbf{1}\{{}^n Y_t > 0\}$ (compared to $U^*({}^n Y_t, t) = \frac{1}{2} \mathbf{1}\{{}^n Y_t > 0\}$ in Example 2), then simulations suggest that U_{modified}^* leads to the better efficiency (compared to the efficiency of $O(1/\sqrt{n})$ obtained in Theorem 1), in the sense of $|{}^n \hat{W}_{\text{modified}} - \hat{w}| \leq O(1/n)$, where ${}^n \hat{W}_{\text{modified}}$ stands for the value of the performance measure for the stochastic model given the fixed initial state and $\hat{w} = 0$. Indeed, Figure 4 is produced by simulations using MATLAB[®], with the diamond symbols standing for the actual accuracy (on a log scale). The solid line (against $\log n$) is the best fitted line coming from the least squares method. In the simulations, for each fixed scaling parameter n , we simulate one hundred trajectories, and, thus, obtain one hundred cost values, whose sample mean is then taken as the estimate for ${}^n \hat{W}_{\text{modified}}$. Here we set $T = M = 1$ for numerical evaluations as above. The slope of the best fitted line is approximately -1 , suggesting the improved efficiency (compared to the slope of $-\frac{1}{2}$ in the case of the original tracking policy). This situation suggests that without influencing the service of other types of file, allocating extra resources than needed (in the fluid model) could potentially lead to the better efficiency.

Example 3. (A linear bandwidth-sharing network.) Let us fix the following primitives: $L = 3$, $J = 2$, $\mu_1 = 4$, $\mu_2 = \mu_3 = 2$, $\lambda_1 = \lambda_2 = \lambda_3 = 1$, $Y_0^1 = 3$, $Y_0^2 = Y_0^3 = 10$, $T = 2$, $Z_1 = Z_2 = 1$, and

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}.$$

This so-called linear bandwidth-sharing network has been studied intensively in [27, Chapter 5] (see also [28]). In particular, by [27, Proposition 5.2.3, Chapter 5], u^* is given by $u_1^*(t) = 1$,

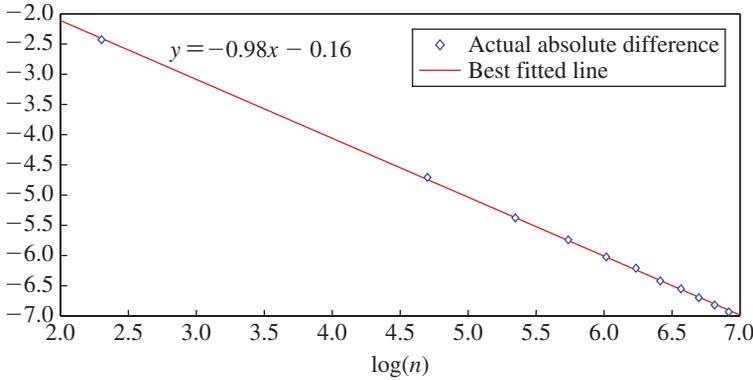


FIGURE 4: The actual accuracy on a log scale for the example of the M/M/1 queue. The vertical axis gives the accuracy (on a log scale), and the horizontal axis gives the scaling parameter (on a log scale), n , ranging from 10 to 1010 in increments of 100.

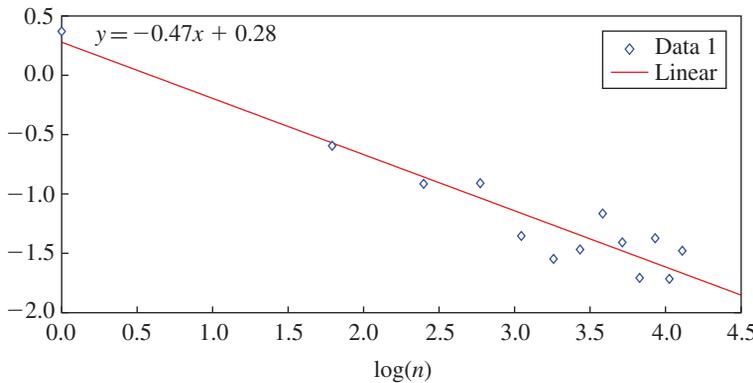


FIGURE 5: The actual accuracy on a log scale for the example of the linear bandwidth-sharing network. The vertical axis gives the accuracy (on a log scale), and the horizontal axis gives the scaling parameter (on a log scale), n , ranging from 1 to 61 in increments of 5. The sample size is 100.

$u_2^*(t) = u_3^*(t) = 0$ on $[0, 1)$ and $u_1^*(t) = \frac{1}{4}$, $u_2^*(t) = u_3^*(t) = \frac{3}{4}$ on $[1, 2)$. Furthermore, the fact that u^* is piecewise constant confirms Lemma 1. Trivial calculations result in $\hat{w} = 44$. Unfortunately, the calculation for ${}^n\hat{W}$ tends to be overwhelmingly difficult. However, Figure 5, with the slope of the line in it being approximately -0.5 , obtained from simulations, confirms the efficiency of $O(1/\sqrt{n})$. The simulations are carried out using MATLAB in the same way as explained at the end of Example 2, and, thus, we shall omit the explanations to avoid repetition. Here we would like to mention that, under u^* , files of type 1 in the fluid model will reach 0 at $t = 1$ and stay there for the rest of the horizon with the unit traffic intensity, while files of the other two types never reach 0 on the horizon. This accounts for the order of $O(1/\sqrt{n})$ (compared with Example 2). In fact, with the slope of the line in it being approximately -1 , Figure 6 suggests that $|{}^n\hat{W}^2 - \hat{w}^2|$ is of order $O(1/n)$, confirming Theorem 2.

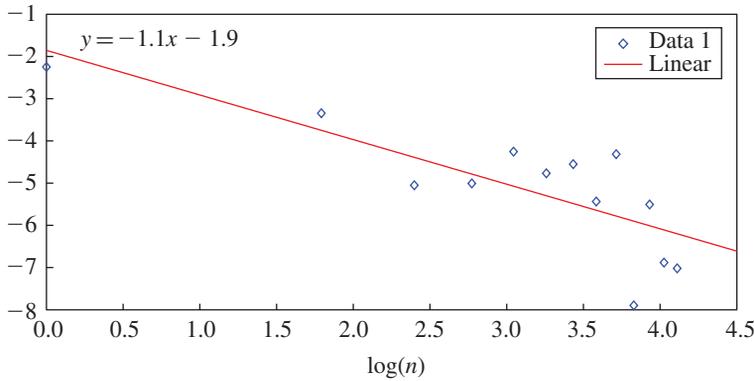


FIGURE 6: The actual accuracy on a log scale for the files of type 2 for the example of the linear bandwidth-sharing network. The vertical axis gives the accuracy (on a log scale), and the horizontal axis gives the scaling parameter (on a log scale), n , ranging from 1 to 61 in increments of 5. The sample size is 1500.

4. Conclusion

To sum up, in this work we have studied the efficiency of the tracking policy for the bandwidth-sharing network model. While it was known for some other networks in the literature that the tracking policy is AFO, its efficiency, to the best of the authors’ knowledge has not been studied intensively. Indeed, [9] is the only work we know that contains a short discussion on that topic. We have shown, in terms of explicit upper boundary estimates for the rate of convergence for performance functionals, that the tracking policy could be as efficient as or less efficient than the feedback-type translation, depending on the parameters of the underlying fluid model. In particular, our result is in favour of the tracking policy over the feedback-type translations, at least for short enough time horizons, owing to its good efficiency as well as to the small amount of required information. The current work contributes new insights about the accuracy of fluid approximations. It appears that the existing knowledge on the accuracy of fluid approximations is very scarce and needs significant development. We hope that the present work becomes an important step in this development.

Appendix A. Proofs of the main statements

A.1. Proof of Lemma 1

For the linear program (1), integrating by parts results in an objective function of the form

$$\begin{aligned} \int_0^T e^\top y(t) dt &= \int_0^T e^\top dt y(T) - \int_0^T \left(\int_0^t e^\top ds \right) (-Mu(t) + \Lambda) dt \\ &= \int_0^T e^\top dt \left(y(0) + \int_0^T (-Mu(t) + \Lambda) dt \right) \\ &\quad - \int_0^T \left(\int_0^t e^\top ds \right) (-Mu(t) + \Lambda) dt \end{aligned}$$

(write $\vec{t} := \int_0^t e^\top ds = (t, \dots, t)$)

$$= \vec{T} y(0) - \int_0^T \vec{T} Mu(s) ds + \vec{T} \Lambda T + \int_0^T \vec{t} Mu(t) dt - \int_0^T \vec{t} \Lambda dt.$$

Here \vec{T} denotes the vector (T, \dots, T) , and we similarly understand \vec{s} below. Therefore, by ignoring the uncontrolled terms in the above expression, fluid model (1) is equivalent to considering

$$-\int_0^T \vec{T} M u(s) ds + \int_0^T \vec{s} M u(s) ds = \int_0^T (\vec{s} - \vec{T}) M u(s) ds$$

as the objective function. Now we note that $(\vec{s} - \vec{T})M$, Λ , and Z , being respectively linear, constant, and constant, are respectively piecewise analytic, piecewise linear, and piecewise constant on $[0, T]$, and the problem has a feasible region, nonempty and bounded, which altogether validate [22, Theorem 3.3], implying that the optimal control for the fluid model is a function $u(t)^*$, constant on intervals $[T_{i-1}, T_i)$, where $i = 1, \dots, N < \infty$ with $T_0 := 0$ and $T_N := T$.

A.2. Proof of Theorem 1

Let us recall that in this proof, though we do not indicate them explicitly, u^* and U^* are the fixed policies under consideration. Then, in addition to our standing assumption of independence made in Section 1, we also observe that

$$\int_0^T e^{\top} \frac{{}^n Y_t}{n} dt = \sum_{l=1}^L \int_0^T \frac{{}^n Y_t^l}{n} dt,$$

meaning that, under the fixed U^* and u^* , we could effectively deal with only one ‘queue’, and apply the same approach to the others. Therefore, from now on we shall focus on files of type l , and, consequently, we have two one-dimensional processes, $\{y_l(t), 0 \leq t \leq T\}$ and $\{{}^n Y_t^l, 0 \leq t \leq T\}$, where we recall that ${}^n Y_t^l$ and $y_l(t)$ are the l th components of ${}^n Y_t$ and $y(t)$. For a real function $x(t)$, we define

$$\|x(t)\|_{[t_1, t_2]} := \sup_{t \in [t_1, t_2]} \{|x(t)|\} \quad \text{and} \quad \|x(t)\|_T := \sup_{t \in [0, T]} \{|x(t)|\}.$$

Lemma 2. *Consider a Poisson process with the counter ${}^n A_t$ and intensity $n\lambda$. Then $\mathbb{P}\{\|{}^n A_t/n - \lambda t\|_T \geq \varepsilon\} \leq \lambda T/n\varepsilon^2$.*

Proof. Clearly, for all $T > 0$, ${}^n M(t) := {}^n A_t - n\lambda t$ gives a martingale with right-continuous trajectories (with the natural filtration) and the index interval $[0, T]$. By the well-known Doob’s L^p -inequality [23, Chapter 2, Theorem 1.7],

$$\mathbb{P}\left\{\left\|\frac{{}^n M(t)}{n}\right\|_T \geq \varepsilon\right\} = \mathbb{P}\left\{\sup_{0 \leq t \leq T} |{}^n M(t)| \geq n\varepsilon\right\} \leq \frac{\mathbb{E}[{}^n M(T)^2]}{(n\varepsilon)^2} = \frac{\lambda T}{n\varepsilon^2},$$

as required.

Lemma 3. *Consider files of type l , that is, the l th ‘queue’ on the first subinterval. Then $\mathbb{P}\{\|{}^n Y_t^l/n - y_l(t)\|_{T_1} \geq \varepsilon\} \leq D_1 := 72\vec{T}\gamma/n\varepsilon^2$. Here we recall that the notation ${}^l \mu_1$, γ , and \vec{T} have been introduced in Section 2.*

Proof. Without generating confusion, throughout this proof, we omit the index for the underlying ‘queue’. Also, without loss of generality (see also Remark 1 below), we assume that ${}^l \mu > 0$. Define ${}^n \tilde{X}_t = {}^n Y_0 + {}^n A_t - {}^n S_t$, where ${}^n S_t$ is a Poisson process with intensity

$n^{-1}\mu$, and nA_t is as in Lemma 2. Of course, $n\tilde{X}_t$, nS_t , and nA_t are defined on the same probability space and independent by assumption. Define also $n\hat{X}_t = \int_0^t \mathbf{1}\{nY_s = 0\} d nS_s$. Then it follows that $n\hat{X}_t$ is nondecreasing and such that $\int_0^\infty \mathbf{1}\{nY_t > 0\} d n\hat{X}_t = 0$. In this way we can write $nY_t = n\tilde{X}_t + n\hat{X}_t$ in the form of the one-dimensional Skorokhod problem, where $n\tilde{X}_t$ is the free process and $n\hat{X}_t$ is the unused capacity process. This representation was also adopted in [14]. See also [16] and [24, Chapter 9]. The solution to this one-dimensional Skorokhod problem is well known (see [14, Appendix A]), so that we have $nY_t = \varphi(n\tilde{X}_t) = n\tilde{X}_t + \sup_{0 \leq s \leq t} [-(n\tilde{X}_s \wedge 0)]$. Here, φ is a Lipschitz mapping with Lipschitz constant 2 on $D[0, \infty)$ in the following sense: for all $x_s, y_s \in D[0, \infty)$,

$$\begin{aligned} \|\varphi(x_t) - \varphi(y_t)\|_{T_1} &= \left\| x_t + \sup_{0 \leq s \leq t} [-(x_s \wedge 0)] - y_t - \sup_{0 \leq s \leq t} [-(y_s \wedge 0)] \right\|_{T_1} \\ &\leq \|x_t - y_t\|_{T_1} + \left\| \sup_{0 \leq s \leq t} [-(x_s \wedge 0)] - \sup_{0 \leq s \leq t} [-(y_s \wedge 0)] \right\|_{T_1}, \end{aligned}$$

where the second term in the last expression is clearly not larger than $\|x_t - y_t\|_{T_1}$. Furthermore, the mapping φ is homogeneous in the sense that, for any $x_s \in D[0, \infty)$, $\varphi(nx_s) = n\varphi(x_s)$ (see [14, Appendix A]). Define $\tilde{x}(t) = y(0) + (\lambda - \mu)t$, the deterministic analogue of $n\tilde{X}_t$. We can write down the Skorokhod problem for the fluid model as well. Then we have

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{nY_t}{n} - y(t) \right\|_{T_1} \geq \varepsilon \right\} &= \mathbb{P} \left\{ \left\| \varphi \left(\frac{n\tilde{X}_t}{n} \right) - \varphi(\tilde{x}(t)) \right\|_{T_1} \geq \varepsilon \right\} \\ &\leq \mathbb{P} \left\{ 2 \left\| \frac{n\tilde{X}_t}{n} - \tilde{x}(t) \right\|_{T_1} \geq \varepsilon \right\} \\ &= \mathbb{P} \left\{ \left\| \frac{n\tilde{X}_t}{n} - \tilde{x}(t) \right\|_{T_1} \geq \frac{\varepsilon}{2} \right\} \\ &\leq \mathbb{P} \left\{ \left\| \frac{nY_0}{n} - y(0) \right\|_{T_1} \geq \frac{\varepsilon}{6} \right\} + \mathbb{P} \left\{ \left\| \frac{nA_t}{n} - \lambda t \right\|_{T_1} \geq \frac{\varepsilon}{6} \right\} \\ &\quad + \mathbb{P} \left\{ \left\| \frac{nS_t}{n} - \mu t \right\|_{T_1} \geq \frac{\varepsilon}{6} \right\} \tag{4} \\ &\leq \frac{36T\lambda}{n\varepsilon^2} + \frac{36T}{n\varepsilon^2} \mu \quad (\text{see Lemma 2}) \tag{5} \\ &\leq \frac{72\bar{T}\gamma}{n\varepsilon^2}, \end{aligned}$$

where the first equality and inequality are due to the homogeneity and Lipschitz continuity of φ , and the last inequality follows from the definitions of γ and \bar{T} . Note also that the first term in (4) vanishes, because $nY_0 = nY_0 = ny(0)$ by definition (see Section 1).

Remark 1. In the proof of Lemma 3, it was assumed that $\mu_l > 0$. Even if this fails to hold, the same result still holds, as then the term in (5) will be absent.

Corollary 1. For all $i = 1, \dots, N$ and $l = 1, \dots, L$,

$$\mathbb{P} \left\{ \left\| \frac{nY_t^l}{n} - y_l(t) \right\|_{[T_{i-1}, T_i]} \geq \varepsilon \right\} \leq iD_1.$$

Proof. Clearly, the statement holds for $i = 1$. Now arguing similarly as in the proof of Lemma 3, we have, for files of type l on the second subinterval, $P\{\|{}^n Y_t^l/n - y_l(t)\|_{[T_1, T_2]} \geq \varepsilon\} \leq D_1 + D_1 = 2D_1$, where the extra D_1 comes from the first term in (4). Here D_1 is accumulated as we consider future intervals, and the statement thus follows. The uniformity in l follows from the universal maximality of γ and \bar{T} , defined in the statement of Lemma 3.

Proof of Theorem 1. According to Corollary 1, we have

$$P\left\{\left\|\frac{{}^n Y_t^l}{n} - y_l(t)\right\|_T \geq \varepsilon\right\} \leq \sum_{i=1}^N P\left\{\left\|\frac{{}^n Y_t^l}{n} - y_l(t)\right\|_{[T_{i-1}, T_i]} \geq \varepsilon\right\} \leq \frac{D_1 N(1 + N)}{2}. \tag{6}$$

Note that the above estimate is uniform in l . Now

$$\begin{aligned} |{}^n \hat{W} - \hat{w}| &= \left| E_{n, Y_0} \left[\int_0^T \sum_{l=1}^L \frac{{}^n Y_t^l}{n} \right] dt - \int_0^T \sum_{l=1}^L y_l(t) dt \right| \\ &\leq \sum_{l=1}^L E_{n, Y_0} \left[\int_0^T \left| \frac{{}^n Y_t^l}{n} - y_l(t) \right| dt \right] \\ &= \sum_{l=1}^L \int_0^T E_{n, Y_0} \left[\left| \frac{{}^n Y_t^l}{n} - y_l(t) \right| \right] dt \\ &= \sum_{l=1}^L \int_0^T \int_0^\infty P_{n, Y_0} \left\{ \left| \frac{{}^n Y_t^l}{n} - y_l(t) \right| \geq \varepsilon \right\} d\varepsilon dt \\ &\leq \sum_{l=1}^L \int_0^T \left\{ \int_0^{6\sqrt{N(N+1)\gamma\bar{T}/\sqrt{n}}} d\varepsilon \right. \\ &\quad \left. + \int_{6\sqrt{N(N+1)\gamma\bar{T}/\sqrt{n}}}^\infty \frac{36N(N+1)\gamma\bar{T}}{n\varepsilon^2} d\varepsilon \right\} dt \quad (\text{see expression (6)}) \\ &= \frac{TL}{\sqrt{n}} \{12\sqrt{N(N+1)\gamma\bar{T}}\}, \end{aligned}$$

where the interchange of integrals is by Fubini’s theorem. Note that in the inequality lines, we set the lower limit and upper limit of the two integrals to be C/\sqrt{n} in order to get the fastest possible convergence of order $O(1/\sqrt{n})$, and set $C = 6\sqrt{N(N+1)\gamma\bar{T}}$ for the same reason.

A.3. Proof of Theorem 2

Let us introduce some additional notation to be used in this subsection. We shall define \hat{w}^l and ${}^n \hat{W}^l$ to be the l th summands of \hat{w} and ${}^n \hat{W}$, i.e. $\hat{w}^l = \int_0^T y_l(t) dt$ and ${}^n \hat{W}^l = E_{Y_0^l}[\int_0^T ({}^n Y_t^l/n) dt]$. Again, the indications of control policies are omitted, as they are fixed. Define, for $t < T$, $w^l(x, t) = \int_t^T \tilde{y}_l(s) ds$, where $\tilde{y}_l(t) = x$, $d\tilde{y}_l/ds = \lambda_l - u_l^*(s)\mu_l$, while $w^l(x, t) \equiv 0$ for $t \geq T$. We shall modify w^l such that it becomes null for all its state arguments larger than K , i.e. we define $v^l(x, t) := w^l(x, t) \mathbf{1}\{x \leq K\}$. In this subsection we are in the framework of Theorem 2.

Lemma 4. *There exist an n -independent positive constant C_l and positive function Φ_l such that, for any $\varepsilon > 0$,*

$$P_{n, Y_0^l} \left\{ \left\| \frac{{}^n Y_t^l}{n} - y_l(t) \right\|_T \geq \varepsilon \right\} \leq C_l e^{-\Phi_l(\varepsilon)n}.$$

Proof. In this proof, it is convenient to omit the subscript indicating the initial position, so that we simply write P instead of $P_{nY_0^l}$. Under the tracking policy (see also Lemma 1), the process ${}^n Y_t^l$ is stationary on subintervals $[0, T_1), \dots, [T_{N-1}, T]$. Moreover, on each subinterval, the process has arrival and potential services which are both renewal processes. Consider now the first subinterval. By Theorem 3.1 of [4] we have an n -independent positive constant ${}^1 C_l$ and positive function ${}^1 \Phi_l$ such that

$$P \left\{ \left\| \frac{{}^n Y_t^l}{n} - y_l(t) \right\|_{T_1} \geq \varepsilon \right\} \leq {}^1 C_l e^{-{}^1 \Phi_l(\varepsilon)n}. \tag{7}$$

In particular, we have

$$P \left\{ \left| \frac{{}^n Y_{T_1}^l}{n} - y_l(T_1) \right| \geq \varepsilon \right\} \leq {}^1 C_l e^{-{}^1 \Phi_l(\varepsilon)n}.$$

Furthermore, Chen [4] showed that, for ${}^n \tilde{Y}_t^l := Y_{nt}^l/n$ and ${}^n \tilde{Y}_0 := nY_0$,

$$P\{\|{}^n \tilde{Y}_t^l - y_l(t)\|_{T_1} \geq \varepsilon\} \leq {}^1 C_l e^{-{}^1 \Phi_l(\varepsilon)n}.$$

However, Mandelbaum and Pats [14, Section 4.4] showed that ${}^n \tilde{Y}_t^l = {}^n Y_t^l/n$ almost surely (see also [6, Section 2.5.3]), and this thus justifies and validates the application of [4, Theorem 3.1] here.

Keeping in mind inequality (7), if we take ${}^n Y_{T_1}$ as the initial state of the process on the next subinterval, then the above reasoning (repeatedly based on [4, Theorem 3.1]) can again be applied, and so on. Eventually, there will be N positive ${}^1 C_l$ -like constants and ${}^1 \Phi_l(\varepsilon)$ -like functions. Finally, we have

$$\begin{aligned} P \left\{ \left\| \frac{{}^n Y_t}{n} - y(t) \right\|_T \geq \varepsilon \right\} &\leq \sum_{i=1}^N P \left\{ \left\| \frac{{}^n Y_t}{n} - y(t) \right\|_{[T_{i-1}, T_i]} \geq \varepsilon \right\} \\ &\leq \sum_{i=1}^N i C_l e^{-i \Phi_l(\varepsilon)n} \\ &\leq C_l e^{-\Phi_l(\varepsilon)n}, \end{aligned}$$

where $\Phi_l(\varepsilon) := \min_{i=1, \dots, N} \{i \Phi_l(\varepsilon)\}$ and $C_l := \sum_{i=1}^N i C_l$.

Remark 2. In Lemma 4, the exact expressions for the constants C_l and the functions Φ_l for all $l = 1, \dots, L$ can be obtained by examining the proofs of [4, Theorems 2.2, 2.3, 3.1]. In fact, we can check from [4, p. 813] that, for all $l = 1, \dots, L$, $\int_0^\infty C_l e^{-\Phi_l(K+x/n-\bar{y})n} x \, dx \leq O(1)$. However, the obtained expressions are much more complicated than those obtained in (6), and are difficult to use in the proof of Theorem 1. On the other hand, the exponentially fast decrease is needed in the proof of Theorem 2 below.

Lemma 5. For all $l = 1, \dots, L$, the following assertions hold on $[0, T)$.

- (a) w^l solves the following equation of dynamic programming type, or the so-called Poisson equation:

$$-\frac{\partial w^l(\tilde{y}_l, t)}{\partial t} = \tilde{y}_l + \frac{\partial w^l(\tilde{y}_l, t)}{\partial \tilde{y}_l} \frac{d\tilde{y}_l(t)}{dt} \tag{8}$$

with boundary condition $w^l(\tilde{y}_l, T) = 0$. Here we take $\partial w^l(\tilde{y}_l, t)/\partial t$ at T_i , $i = 0, \dots, N - 1$, as the right derivatives. In particular, $\hat{w}^l = w^l(Y_0^l, 0)$.

(b) *There exists an $R > 0$ so that*

$$|w^l(\tilde{y}_l, t)|, \quad \left| \frac{\partial w^l(\tilde{y}_l, t)}{\partial \tilde{y}_l} \right|, \quad \left| \frac{\partial w^l(\tilde{y}_l, t)}{\partial t} \right|, \quad \text{and} \quad \left| \frac{\partial^2 w^l(\tilde{y}_l, t)}{\partial \tilde{y}_l^2} \right|$$

are all bounded by R on $[0, K] \times [0, T)$, uniformly in $l = 1, \dots, L$.

(c) *$|v^l(\tilde{y}_l, t)|$ is bounded and has a bounded partial derivative $|\partial v^l(\tilde{y}_l, t)/\partial t|$. Here we take $\partial v^l(\tilde{y}_l, t)/\partial t$ at T_i , $i = 0, \dots, N - 1$, as the right derivatives.*

Remark 3. Here we recall that, for all $l = 1, \dots, L$, $\tilde{y}_l(t)$ defined at the beginning of this subsection differs from $y_l(t)$ in that $\tilde{y}_i(t)$ can be negative while $y_l(t)$ cannot. In other words, on $[0, T)$, $y_l(t) = \tilde{y}_l(t)$ under Assumption 1. This explains the equality of $\hat{w}^l = w^l(Y_0^l, 0)$ in Lemma 5(a).

Proof of Lemma 5. We shall omit the index $l = 1, \dots, L$ for brevity, and, thus, simply write w instead of w^l , and so on throughout this proof. Hopefully, this abuse of notation will not lead to confusion. Clearly, all the assertions hold trivially for the case of only one subinterval, i.e. $N = 1$. Moreover, it is not hard to see that parts (b) and (c) follow from part (a) and the definition of v^l . Therefore, let us verify part (a) only for the case in which $N \geq 2$, by induction.

Consider the case in which $N = 2$, i.e. the case of only two subintervals and, thus, $T_2 = T$. Let us calculate $w(y, t)$ and its partial derivatives. Suppose that $t \in [T_1, T)$. Then we have $\tilde{y}_t = y$ and $\tilde{y}_s = y + (s - t)(\lambda - 2\mu)$ if $s \in [t, T)$, leading to

$$\begin{aligned} w(y, t) &= \int_t^T (y + (s - t)(\lambda - 2\mu)) \, ds \\ &= y(T - t) + (\lambda - 2\mu) \frac{T^2 - t^2}{2} - (\lambda - 2\mu)t(T - t) \end{aligned}$$

and

$$\frac{\partial w}{\partial y} = T - t, \quad \frac{\partial^2 w}{\partial y^2} = 0, \quad \frac{\partial w}{\partial t} = -y - (\lambda - 2\mu)(T - t), \quad \frac{\partial^2 w}{\partial t^2} = \lambda - 2\mu.$$

Suppose that $t \in [0, T_1)$. Then we have $\tilde{y}_t = y$ and $\tilde{y}_s = y + (s - t)(\lambda - 1\mu)$ if $s \in [t, T_1)$, and $\tilde{y}_s = y + (T_1 - t)(\lambda - 1\mu) + (s - T_1)(\lambda - 2\mu)$ if $s \in [T_1, T)$, leading to

$$\begin{aligned} w(y, t) &= \int_t^{T_1} (y + (s - t)(\lambda - 1\mu)) \, ds \\ &\quad + \int_{T_1}^T (y + (T_1 - t)(\lambda - 1\mu) + (s - T_1)(\lambda - 2\mu)) \, ds \\ &= y(T - t) + (\lambda - 1\mu) \frac{T_1^2 - t^2}{2} - (\lambda - 1\mu)(tT_1 - t^2) \\ &\quad + (T_1 - t)(T - T_1)(\lambda - 1\mu) + (\lambda - 2\mu) \frac{T^2 - T_1^2}{2} - (\lambda - 2\mu)T_1(T - T_1) \end{aligned}$$

and

$$\frac{\partial w}{\partial y} = T - t, \quad \frac{\partial^2 w}{\partial y^2} = 0, \quad \frac{\partial w}{\partial t} = -y - (\lambda - 1\mu)(T - t), \quad \frac{\partial^2 w}{\partial t^2} = \lambda - 1\mu.$$

Now, for the case in which $N = 2$, part (a) follows immediately from the direct substitution of the corresponding expressions obtained above into both sides of (8).

Suppose that part (a) holds for $N = k$. Now let us verify that part (a) also holds for the case in which $N = k + 1$. Suppose that $t \in [T_1, T)$. Then, as there are only k subintervals on $[T_1, T)$, part (a) follows from the inductive supposition. Suppose that $t \in [0, T_1)$. Then

$$w(y, t) = \int_t^{T_1} y + (s - t)(\lambda - \mu) \, ds + \int_{T_1}^T (y + (T_1 - t)(\lambda - \mu) + \sum_{i=2}^N (\lambda - \mu)(s - T_1) \mathbf{1}\{s \in [T_{i-1}, T_i)\}) \, ds,$$

based on which, direct calculations result in

$$\frac{\partial w}{\partial y} = T_1 - t + (T - T_1) = T - t$$

and

$$\begin{aligned} \frac{\partial w}{\partial t} &= -y + (-t - (T_1 - t - t))(\lambda - \mu) - (T - T_1)(\lambda - \mu) \\ &= -y - (\lambda - \mu)(T - t) \\ &= -\left(y + \frac{\partial w}{\partial y} \frac{dy}{dt}\right). \end{aligned}$$

Thus, part (a) has been proved by induction.

Since U^* is nonstationary, the resulting process is nonstationary, and we shall apply the extended generator \tilde{L} , whose definition can be found in, for example, [12, Chapter 1]. Here, we refer the reader to the references for this definition to avoid extra definitions and notation.

Lemma 6. For all $l = 1, \dots, L$, v^l is in the domain of \tilde{L} , with the form

$$\tilde{L}v^l(j, t) = \frac{\partial v^l(j, t)}{\partial t} + \sum_{k \in \mathbb{Z}_+} {}^l q_{j,k}(t)v^l(k, t),$$

where $[{}^l q_{jk}(t)]_{j,k \in \mathbb{Z}_+}$ gives the Q -matrix of the continuous-time Markov chain corresponding to the l th queue (recall the standing assumption of independence made in Section 1, the definition of the tracking policy U^* , and Lemma 1).

Proof. Given Lemma 5, this statement directly follows from the proof of [12, Proposition 14.4].

Proof of Theorem 2. Define $\lambda_l(t) := \lambda_l = {}^l q_{j,j+1}(t) \leq \bar{\lambda}$ and $\mu_l(t) := U_l^*({}^n Y^l, t)\mu_l = {}^l q_{n Y^l, n Y^l - 1}(t) \leq \bar{\mu}$, where ${}^l q_{0,-1}(t) = 0$. For simplicity, in what follows, unless necessary, we shall not indicate the control policy, and omit the index $l = 1, \dots, L$ corresponding to the type of file. For example, for an arbitrarily fixed $l = 1, \dots, L$, we shall write ${}^n Y_t$ for ${}^n Y_t^l$, and so on.

Owing to Lemma 6 and [12, Lemma 2.1], the following Dynkin's formula is valid:

$$0 = E_{n Y_0} \left[v \left(\frac{{}^n Y_T}{n}, T \right) \right] = v \left(\frac{{}^n Y_0}{n}, 0 \right) + E_{n Y_0} \left[\int_0^T \tilde{L}v \left(\frac{{}^n Y_t}{n}, t \right) dt \right]. \tag{9}$$

The first equality above follows from the definition of v given at the beginning of this subsection. Adding the finite expectation of $E^{nY_0}[\int_0^T ({}^nY_t/n) dt]$ to the both sides of (9), some arrangements then lead to

$$\begin{aligned} \left| {}^n\hat{W} - v\left(\frac{{}^nY_0}{n}, 0\right) \right| &= \left| E^{nY_0} \left[\int_0^T \left(\frac{{}^nY_t}{n} + \tilde{L}v\left(\frac{{}^nY_t}{n}, t\right) \right) dt \right] \right| \\ &\leq E^{nY_0} \left[\int_0^T \left| \left(\frac{{}^nY_t}{n} + \tilde{L}v\left(\frac{{}^nY_t}{n}, t\right) \right) \right| \mathbf{1}\{0 < {}^nY_t < nK\} dt \right] \\ &\quad + E^{nY_0} \left[\int_0^T \left| \left(\frac{{}^nY_t}{n} + \tilde{L}v\left(\frac{{}^nY_t}{n}, t\right) \right) \right| \mathbf{1}\{{}^nY_t = 0\} dt \right] \\ &\quad + E^{nY_0} \left[\int_0^T \left| \left(\frac{{}^nY_t}{n} + \tilde{L}v\left(\frac{{}^nY_t}{n}, t\right) \right) \right| \mathbf{1}\{{}^nY_t \geq nK\} dt \right] \\ &=: S_1 + S_2 + S_3. \end{aligned} \tag{10}$$

This is justified by the fact that $E^{nY_0}[\int_0^T ({}^nY_t/n) dt]$ is bounded by a function of ${}^nY_0 < \infty$. Recall that the scaled queue is stochastically dominated by a Poisson process with intensity $n\lambda$. In addition, recall that $\hat{w} = w(Y_0, 0) = v(Y_0, 0)$ (see Lemma 5, Remark 3, and Assumption 1). Note also that the fact that ${}^nY_t \in \mathbb{Z}_+$ will often not be indicated explicitly.

Next let us analyze the three terms on the right-hand side (RHS) of (10) case by case.

Case (i): estimating S_1 . Clearly, on the set $\{0 < {}^nY_t < nK\}$, we have the integrand in the above expectation to be of the form

$$\begin{aligned} &\left| \frac{{}^nY_t}{n} + \frac{\partial v({}^nY_t/n, t)}{\partial t} + n\lambda v\left(\frac{{}^nY_t + 1}{n}, t\right) + n\mu(t)v\left(\frac{{}^nY_t - 1}{n}, t\right) - n(\lambda + \mu(t))v\left(\frac{{}^nY_t}{n}, t\right) \right| \\ &= \left| \frac{{}^nY_t}{n} + \frac{\partial v({}^nY_t/n, t)}{\partial t} + n\lambda \left(v\left(\frac{{}^nY_t + 1}{n}, t\right) - v\left(\frac{{}^nY_t}{n}, t\right) \right) \right. \\ &\quad \left. + n\mu(t) \left(v\left(\frac{{}^nY_t - 1}{n}, t\right) - v\left(\frac{{}^nY_t}{n}, t\right) \right) \right|. \end{aligned} \tag{11}$$

Rewriting ${}^nY_t/n$ according to (8), and using the fact that $w = v$ for $0 < {}^nY_t < nY$, which is due to the definition of v given at the beginning of this subsection, we find that

$$\begin{aligned} \text{RHS of (11)} &= \left| n\lambda \left(v\left(\frac{{}^nY_t + 1}{n}, t\right) - v\left(\frac{{}^nY_t}{n}, t\right) \right) + n\mu(t) \left(v\left(\frac{{}^nY_t - 1}{n}, t\right) - v\left(\frac{{}^nY_t}{n}, t\right) \right) \right. \\ &\quad \left. - (\lambda - \mu(t)) \frac{\partial v({}^nY_t/n, t)}{\partial y} \right|. \end{aligned}$$

As a result, we have

$$\begin{aligned} &E^{nY_0} \left[\int_0^T \left| \left(\frac{{}^nY_t}{n} + \tilde{L}v\left(\frac{{}^nY_t}{n}, t\right) \right) \right| \mathbf{1}\{0 < {}^nY_t < nK\} dt \right] \\ &\leq E^{nY_0} \left[\int_0^T \left\{ \lambda \left| n \left(v\left(\frac{{}^nY_t + 1}{n}, t\right) - v\left(\frac{{}^nY_t}{n}, t\right) \right) - \frac{\partial v({}^nY_t/n, t)}{\partial y} \right| \right. \right. \\ &\quad \left. \left. + \bar{\mu} \left| \frac{\partial v({}^nY_t/n, t)}{\partial y} - n \left(v\left(\frac{{}^nY_t}{n}, t\right) - v\left(\frac{{}^nY_t - 1}{n}, t\right) \right) \right| \mathbf{1}\{0 < {}^nY_t < nK\} \right\} dt \right] \\ &\leq \frac{(\bar{\lambda} + \bar{\mu})RT}{2n}, \end{aligned}$$

where the last inequality follows by applying Taylor’s theorem to $v({}^n Y_t/n, t) - v(({}^n Y_t - 1)/n, t)$ and $v(({}^n Y_t + 1)/n, t) - v({}^n Y_t/n, t)$, bounding the indicator by 1, and using Lemma 5.

Case (ii): estimating S_2 . By Lemma 6 and the substitution of ${}^n Y_t/n = 0$, we have

$$\begin{aligned}
 & \mathbb{E}^{n Y_0} \left[\int_0^T \left| \left(\frac{{}^n Y_t}{n} + \tilde{L}v\left(\frac{{}^n Y_t}{n}, t\right) \right) \right| \mathbf{1}\{{}^n Y_t = 0\} dt \right] \\
 &= \mathbb{E}^{n Y_0} \left[\int_0^T \left| \left(\frac{\partial v(0, t)}{\partial t} + n\lambda v\left(\frac{1}{n}, t\right) - n\lambda v(0, t) \right) \right| \mathbf{1}\{{}^n Y_t = 0\} dt \right] \\
 &\leq R(1 + 2n\bar{\lambda}) \mathbb{E}^{n Y_0} \left[\int_0^T \mathbf{1}\{{}^n Y_t = 0\} dt \right] \\
 &= R(1 + 2n\bar{\lambda}) \mathbb{E}^{n Y_0} [\text{time spent by } {}^n Y_t \text{ at } 0 \text{ up to time } T] \\
 &= R(1 + 2n\bar{\lambda}) \mathbb{E}^{n Y_0} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } 0 \text{ up to time } T \right] \\
 &= R(1 + 2n\bar{\lambda}) \mathbb{E}^{n Y_0} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } 0 \text{ up to time } T \mid \frac{{}^n Y_t}{n} \text{ visits } 0 \text{ up to time } T \right] \\
 &\quad \times \mathbb{P}^{n Y_0} \left\{ \frac{{}^n Y_t}{n} \text{ visits } 0 \text{ up to time } T \right\} \\
 &\leq R(1 + 2n\bar{\lambda}) T \mathbb{P}^{n Y_0} \left\{ \left\| \frac{{}^n Y_t}{n} - y(t) \right\|_T \geq \underline{y} \right\} \\
 &\leq R(1 + 2n\bar{\lambda}) T C_I e^{-\Phi_I(\underline{y})n}, \tag{12}
 \end{aligned}$$

where the last inequality follows from Lemma 4. Furthermore, we clearly see that the upper bound for $\mathbb{P}^{n Y_0} \{ \|{}^n Y_t/n - y(t)\|_T \geq \underline{y} \}$ based on (6), which converges to 0 as fast as $O(1/n)$, is not enough for our purpose here.

Case (iii): estimating S_3 . By the definition of v (given at the beginning of this subsection) and Lemma 6, we have

$$\begin{aligned}
 & \mathbb{E}^{n Y_0} \left[\int_0^T \left| \left(\frac{{}^n Y_t}{n} + \tilde{L}v\left(\frac{{}^n Y_t}{n}, t\right) \right) \right| \mathbf{1}\{{}^n Y_t \geq nK\} dt \right] \\
 &= \mathbb{E}^{n Y_0} \left[\int_0^T \left| \left(K + \frac{\partial v(K, t)}{\partial t} + n\mu(t)v\left(\frac{nK-1}{n}, t\right) - n(\lambda + \mu(t))v(K, t) \right) \right| \right. \\
 &\quad \left. \times \mathbf{1}\{{}^n Y_t = nK\} dt \right] \\
 &\quad + \mathbb{E}^{n Y_0} \left[\int_0^T \left| K + \frac{1}{n} + n\mu(t)v(K, t) \right| \mathbf{1}\{{}^n Y_t = nK + 1\} dt \right] \\
 &\quad + \mathbb{E}^{n Y_0} \left[\int_0^T \frac{{}^n Y_t}{n} \mathbf{1}\{{}^n Y_t \geq nK + 2\} dt \right] \\
 &=: S_{31} + S_{32} + S_{33}.
 \end{aligned}$$

Subcase (iii.1): estimating S_{31} . We have

$$\begin{aligned}
 & \mathbb{E}^{n Y_0} \left[\int_0^T \left| \left(K + \frac{\partial v(K, t)}{\partial t} + n\mu(t)v\left(\frac{nK-1}{n}, t\right) - n(\lambda + \mu(t))v(K, t) \right) \right| \mathbf{1}\{{}^n Y_t = nK\} dt \right] \\
 &\leq (K + R + n\bar{\lambda}R + 2n\bar{\mu}R) \mathbb{E}^{n Y_0} \left[\int_0^T \mathbf{1}\left\{ \frac{{}^n Y_t}{n} = K \right\} ds \right]
 \end{aligned}$$

$$\begin{aligned}
 &= (K + R + n\bar{\lambda}R + 2n\bar{\mu}R) E^{n_{Y_0}} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } K \text{ up to time } T \right] \\
 &= (K + R + n\bar{\lambda}R + 2n\bar{\mu}R) \\
 &\quad \times E^{n_{Y_0}} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } K \text{ up to time } T \mid \frac{{}^n Y_t}{n} \text{ visits } K \text{ up to time } T \right] \\
 &\quad \times P^{n_{Y_0}} \left\{ \frac{{}^n Y_t}{n} \text{ visits } K \text{ up to time } T \right\} \\
 &\leq (Y + R + n\bar{\lambda}R + 2n\bar{\mu}R) T P^{n_{Y_0}} \left\{ \left\| \frac{{}^n Y_t}{n} - y(t) \right\|_T \geq K - \bar{y} \right\} \\
 &\leq (K + R + n\bar{\lambda}R + 2n\bar{\mu}R) T C_l e^{-\Phi_l(K-\bar{y})n},
 \end{aligned}$$

where the last inequality follows from Lemma 4.

Subcase (iii.2): estimating S₃₂. Similarly to how we treated subcase (iii.1), we have

$$\begin{aligned}
 &E^{n_{Y_0}} \left[\int_0^T \left| K + \frac{1}{n} + n\mu(t)v(K, t) \right| \mathbf{1}\{{}^n Y_t = nK + 1\} dt \right] \\
 &\leq \left(K + \frac{1}{n} + n\bar{\mu}R \right) E^{n_{Y_0}} \left[\int_0^T \mathbf{1}\{{}^n Y_t = nK + 1\} dt \right] \\
 &= \left(K + \frac{1}{n} + n\bar{\mu}R \right) E^{n_{Y_0}} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } K + \frac{1}{n} \text{ up to time } T \right] \\
 &= \left(K + \frac{1}{n} + n\bar{\mu}R \right) \\
 &\quad \times E^{n_{Y_0}} \left[\text{time spent by } \frac{{}^n Y_t}{n} \text{ at } K + \frac{1}{n} \text{ up to time } T \mid \frac{{}^n Y_t}{n} \text{ visits } K \text{ up to time } T \right] \\
 &\quad \times P^{n_{Y_0}} \left\{ \frac{{}^n Y_t}{n} \text{ visits } K \text{ up to time } T \right\} \\
 &\leq \left(K + \frac{1}{n} + n\bar{\mu}R \right) T P^{n_{Y_0}} \left\{ \left\| \frac{{}^n Y_t}{n} - y(t) \right\|_T \geq K - \bar{y} \right\} \\
 &\leq \left(K + \frac{1}{n} + n\bar{\mu}R \right) T C_l e^{-\Phi_l(K-\bar{y})n},
 \end{aligned}$$

where the last equality follows from that fact that ${}^n Y_t/n$ visits $K + 1/n$ up to time T only if ${}^n Y_t/n$ visits K up to time T , and the last inequality follows from Lemma 4.

Subcase (iii.3): estimating S₃₃. We have

$$\begin{aligned}
 &E^{n_{Y_0}} \left[\int_0^T \frac{{}^n Y_t}{n} \mathbf{1}\{{}^n Y_t \geq nK + 2\} dt \right] \\
 &= E^{n_{Y_0}} \left[\int_0^T \sum_{j=2}^{\infty} \left(K + \frac{j}{n} \right) \mathbf{1}\{{}^n Y_t = nK + j\} dt \right] \\
 &= \sum_{j=2}^{\infty} E^{n_{Y_0}} \left[\int_0^T K \mathbf{1}\{{}^n Y_t = nK + j\} dt \right] + \sum_{j=2}^{\infty} E^{n_{Y_0}} \left[\int_0^T \frac{j}{n} \mathbf{1}\{{}^n Y_t = nK + j\} dt \right],
 \end{aligned}$$

where, similarly to subcases (iii.1) and (iii.2) we have

$$\begin{aligned}
 & \sum_{j=2}^{\infty} E^{nY_0} \left[\int_0^T K \mathbf{1}\{^n Y_t = nK + j\} dt \right] \\
 &= K E^{nY_0} \left[\int_0^T \mathbf{1}\{^n Y_t \geq nK + 2\} dt \right] \\
 &= K E^{nY_0} \left[\text{time spend by } \frac{^n Y_t}{n} \text{ on } \left\{ K + \frac{2}{n}, K + \frac{3}{n}, \dots \right\} \text{ up to time } T \right] \\
 &= K E^{nY_0} \left[\text{time spend by } \frac{^n Y_t}{n} \text{ on } \left\{ K + \frac{2}{n}, K + \frac{3}{n}, \dots \right\} \right. \\
 &\quad \left. \text{up to time } T \mid \frac{^n Y_t}{n} \text{ visits } K \text{ up to time } T \right] \\
 &\quad \times P^{nY_0} \left\{ \frac{^n Y_t}{n} \text{ visits } K \text{ up to time } T \right\} \\
 &\leq K T C_l e^{-\Phi_l(K-\bar{y})n}
 \end{aligned}$$

and

$$\begin{aligned}
 & \sum_{j=2}^{\infty} E^{nY_0} \left[\int_0^T \frac{j}{n} \mathbf{1}\{^n Y_t = K + j\} dt \right] \\
 &= \frac{1}{n} \sum_{j=2}^{\infty} j E^{nY_0} \left[\text{time spent by } \frac{^n Y_t}{n} \text{ at } K + j \text{ up to time } T \mid \right. \\
 &\quad \left. \frac{^n Y_t}{n} \text{ visits } K + \frac{j}{n} \text{ up to time } T \right] \\
 &\quad \times P^{nY_0} \left\{ \frac{^n Y_t}{n} \text{ visits } K + \frac{j}{n} \text{ up to time } T \right\} \\
 &\leq \frac{T}{n} \sum_{j=2}^{\infty} j P^{nY_0} \left\{ \left\| \frac{^n Y_t}{n} - y(t) \right\|_T \geq K + \frac{j}{n} - \bar{y} \right\} \\
 &\leq \frac{T}{n} \sum_{j=2}^{\infty} j C_l e^{-\Phi(K+j/n+\bar{y})n} \\
 &\leq O\left(\frac{1}{n}\right),
 \end{aligned}$$

where the last inequality follows from Remark 2. Hence,

$$E^{nY_0} \left[\int_0^T \frac{^n Y_t}{n} \mathbf{1}\{^n Y_t \geq nK + 2\} dt \right] \leq O\left(\frac{1}{n}\right).$$

Now the above estimates together with (10), (11), and (12) result in $|^n \hat{W}^l - \hat{w}^l| \leq O(1/n)$. Repeating the above reasoning $L - 1$ times and adding the resulting upper boundary estimates together will lead to $|^n \hat{W} - \hat{w}| \leq O(1/n)$, as required.

Acknowledgements

We are grateful to the reviewer for his/her valuable comments, which significantly improved the work. This research was partially supported by the Alliance: Franco-British Research Partnership Programme, project ‘Impulsive Control with Delays and Application to Traffic Control in the Internet’ (PN08.021). During the time of this research, Y. Zhang was supported by the ORSAS award, the University of Liverpool Graduate Association postgraduate scholarship (Hong Kong), and the standard research studentship from the University of Liverpool.

References

- [1] ANDERSON, W. J. (1991). *Continuous-Time Markov Chains*. Springer, New York.
- [2] BÄUERLE, N. (2000). Asymptotic optimality of tracking policies in stochastic networks. *Ann. Appl. Prob.* **10**, 1065–1083.
- [3] CANTRELL, P. (1986). Computation of the transient M/M/1 queue CDF, PDF, and mean with generalized Q-functions. *IEEE Trans. Commun.* **34**, 814–817.
- [4] CHEN, H. (1996). Rate of convergence of the fluid approximation for generalized Jackson networks. *J. Appl. Prob.* **33**, 804–814.
- [5] CHEN, H. AND MANDELBAUM, A. (1991). Discrete flow networks: bottleneck analysis and fluid approximations. *Math. Operat. Res.* **16**, 408–446.
- [6] CHEN, H. AND MANDELBAUM, A. (1994). Hierarchical modeling of stochastic networks. Part I. Fluid models. In *Stochastic Modeling and Analysis of Manufacturing Systems*, ed. D. Yao, Springer, New York, pp. 47–105.
- [7] DAI, J. G. (1995). On positive Harris recurrence of multiclass queueing networks: a unified approach via fluid limit models. *Ann. Appl. Prob.* **5**, 49–77.
- [8] GAJRAT, A. AND HORDIJK, A. (2000). Fluid approximation of a controlled multiclass tandem network. *Queueing Systems* **35**, 349–380.
- [9] GAJRAT, A., HORDIJK, A. AND RIDDER, A. (2003). Large-deviations analysis of the fluid approximation for a controllable tandem queue. *Ann. Appl. Prob.* **13**, 1423–1448.
- [10] GLEISSNER, W. (1988). The spread of epidemics. *Appl. Math. Comput.* **27**, 167–171.
- [11] GUO, X. AND HERNÁNDEZ-LERMA, O. (2009). *Continuous-Time Markov Decision Processes*. Springer, Berlin.
- [12] HERNÁNDEZ-LERMA, O. (1994). *Lectures on Continuous-Time Markov Control Processes*. Sociedad Matemática Mexicana, Mexico City.
- [13] MAGLARAS, C. (2000). Discrete-review policies for scheduling stochastic networks: trajectory tracking and fluid-scale asymptotic optimality. *Ann. Appl. Prob.* **10**, 897–929.
- [14] MANDELBAUM, A. AND PATS, G. (1995). State-dependent queues: approximations and applications. In *Stochastic Networks* (IMA Vol. Math. Appl. **71**), eds F. Kelly and R. Williams, Springer, New York, pp. 239–282.
- [15] MASSOULIÉ, L. AND ROBERTS, J. W. (2000). Bandwidth sharing and admission control for elastic traffic. *Telecommun. Systems* **15**, 185–201.
- [16] PANG, G. AND DAY, M. V. (2007). Fluid limits of optimally controlled queueing networks. *J. Appl. Math. Stoch. Anal.* **2007**, 19 pp.
- [17] PIUNOVSKIY, A. (2009). Controlled jump Markov processes with local transitions and their fluid approximation. *WSEAS Trans. Systems Control* **4**, 399–412.
- [18] PIUNOVSKIY, A. B. (2009). Random walk, birth-and-death process and their fluid approximations: absorbing case. *Math. Meth. Operat. Res.* **70**, 285–312.
- [19] PIUNOVSKIY, A. B. AND CLANCY, D. (2008). An explicit optimal intervention policy for a deterministic epidemic model. *Optimal Control Appl. Meth.* **29**, 413–428.
- [20] PIUNOVSKIY, A. AND ZHANG, Y. (2011). Accuracy of fluid approximations to controlled birth-and-death processes: absorbing case. *Math. Meth. Operat. Res.*, 29 pp.
- [21] PIUNOVSKIY, A. AND ZHANG, Y. (2011). On the fluid approximations of a class of general inventory level-dependent EOQ and EPQ models. To appear in *Adv. Operat. Res.*
- [22] PULLAN, M. C. (1995). Forms of optimal solutions for separated continuous linear programs. *SIAM J. Control. Optimization* **33**, 1952–1977.
- [23] REVUZ, D. AND YOR, M. (1999). *Continuous Martingales and Brownian Motion*, 3rd edn. Springer, Berlin.
- [24] ROBERT, P. (2003). *Stochastic Networks and Queues* (Appl. Math. **52**). Springer, Berlin.
- [25] ROBERTS, J. AND MASSOULIÉ, L. (1998). Bandwidth sharing and admission control for elastic traffic. In *Proc. of ITC Specialist Seminar*, Yokohama, Japan, pp. 185–201.
- [26] SHARMA, O. P. AND TARABIA, A. M. K. (2000). A simple transient analysis of an M/M/1/N queue. *Sankhyā A* **62**, 273–281.

- [27] VERLOOP, I. M. (2009). *Scheduling in Stochastic Resource-Sharing Systems*. Doctoral Thesis, Eindhoven University of Technology.
- [28] VERLOOP, I. M. AND NÚÑEZ-QUEJIA, R. (2009). Assessing the efficiency of resource allocations in bandwidth-sharing networks. *Performance Evaluation* **66**, 59–77.
- [29] YE, L., GUO, X. AND HERNÁNDEZ-LERMA, O. (2008). Existence and regularity of a nonhomogeneous transition matrix under measurability conditions. *J. Theoret. Prob.* **21**, 604–627.
- [30] ZHANG, J. AND COYLE, E. J., JR. (1991). The transient solution of time-dependent M/M/1 queues. *IEEE Trans. Inf. Theory* **37**, 1690–1696.