

A linear programming approach to stability, optimisation and performance analysis for Markovian multiclass queueing networks[★]

Kevin D. Glazebrook^a and José Niño-Mora^b

*^aSchool of Mathematics and Statistics, Newcastle University,
Newcastle-upon-Tyne NE1 7RU, UK*

E-mail: kevin.glazebrook@newcastle.ac.uk

*^bCORE, Université Catholique de Louvain,
B-1348 Louvain-la-Neuve, Belgium*

E-mail: jnimora@alum.mit.edu

Our object of study is a multiclass queueing network (MQNET) which consists of a collection of (connected) single-server stations. Exogenous arrivals into the system form independent Poisson streams, service times are exponential and we have Markovian routing of customers between stations. Recent results concerning linear programming (LP) based approaches enable us to establish a simple and intuitive stability condition. This is of interest in its own right, but also enables us to progress with a study of optimal scheduling and performance analysis. Our methodology here is also based on LP. A primal–dual approach exploits the fact that the system satisfies (approximate) conservation laws to yield performance guarantees for a natural index-based scheduling heuristic. We are also able to analyse the performance of an arbitrary priority policy.

Keywords: achievable region, multiclass queueing network, optimal scheduling, performance guarantee, priority index, stability

1. Introduction

Multiclass queueing networks (MQNETs) provide a rich range of models for complex service systems in application areas that include manufacturing (see Buzacott and Shanthikumar [6]) and computer-communication systems (see Gelenbe and Mitrani [9]). The practical needs to evaluate and improve the performance of such systems

[★] The first author was supported by the Engineering and Physical Sciences Research Council by means of Grant No. GR/K03043. The second author was partially supported by a CORE Research Fellowship and by EC Marie Curie Research Fellowship No. ERBFMBICT961480.

have motivated extensive research efforts on the analysis, optimisation and stability of MQNETs. The model of interest here is a multi-station MQNET with single-server stations and is described in detail in section 2. We use LP-based methodologies to investigate stability (section 3), the approximate optimality of index-based scheduling policies (section 4) and performance analysis (section 5).

Investigation of *stability* is a necessary precursor to consideration of optimisation and performance issues. It is a complex area and has attracted a great deal of recent research effort, much of which has focussed on the development of computational tests for stability either via a fluid model approximation of the MQNET (see, e.g., Dai and Weiss [7]) or via a non-idling performance LP (Kumar and Meyn [15]). We make use of the latter approach to establish in section 3 a simple and intuitive stability condition for our MQNETs. The condition may be stated as follows: if each station has enough service capacity to handle its peak traffic intensity, then the network is stable under the class of stationary nonidling service disciplines. No claim is made that this result is best possible, but it seems to be precisely what is required for addressing subsequent performance evaluation and optimisation issues.

Recent work by Papadimitriou and Tsitsiklis [17] suggests that the *performance optimisation* problem (designing a scheduling policy that optimises a given objective and computing the corresponding optimal system performance) is computationally intractable in most MQNET models. Consequently, researchers have focussed their efforts on designing heuristic scheduling policies with good empirical or asymptotic performance (see, e.g., Harrison and Wein [13]). However, performance guarantees for such policies are rarely available. This is in sharp contrast to the situation in *deterministic scheduling* in which exploitation of the so-called *primal–dual* method has yielded approaches to the design of heuristic algorithms with given performance guarantees. In this method, one seeks to construct simultaneously a heuristic solution to the problem and a feasible solution to the dual of an LP relaxation of it. See the survey paper of Goemans and Williamson [12]. In sections 4 and 5, we use an analogous primal–dual approach to develop a class of priority-index scheduling policies for our MQNETs with associated performance guarantees. These results use the work of Bertsimas and Niño-Mora [3] and Glazebrook and Garbe [11], in which the notion of (approximate) conservation laws plays a central role. Finally, in section 6, this framework yields approaches to the analysis of performance of any given priority policy.

2. The MQNET model

We consider a Markovian open MQNET with N customer classes and M single-server stations. Station $m \in \mathcal{M} = \{1, \dots, M\}$ services a *constituency* of customer classes $C_m \subseteq \mathcal{N} = \{1, \dots, N\}$, where C_1, \dots, C_M is a partition of \mathcal{N} . Class i customers (also called *i-customers*) arrive exogenously at the network as a Poisson process with rate λ_i , and are serviced at station $s(i)$, during an exponential service time with rate μ_i . Upon completion of its service, an i -customer may be *routed* for further service

as a j -job, with probability p_{ij} , or it may *leave* the network, with probability $p_{i0} = 1 - \sum_{j \in \mathcal{N}} p_{ij}$. To ensure that a customer entering the network leaves it with probability one, we require that matrix $\mathbf{I} - \mathbf{P}$ is invertible. We further assume that all arrival, service and routing processes are mutually independent.

We introduce next other parameters of interest. The *total arrival rate of j -customers*, denoted λ_j , is given by the solution of

$$\lambda_j = \lambda_j + \sum_{i \in \mathcal{N}} \lambda_i p_{ij}, \text{ for } j \in \mathcal{N}.$$

The *nominal traffic intensity of j -customers*, denoted ρ_j , is given by

$$\rho_j = \frac{\lambda_j}{\mu_j}, \text{ for } j \in \mathcal{N},$$

and represents the rate at which work brought in by j -customers (external and internal) enters the network. Given a subset of job classes $S \subseteq \mathcal{N}$, we define analogously the *traffic intensity for S -customers*, denoted $\rho(S)$, by

$$\rho(S) = \sum_{j \in S} \rho_j.$$

We further define the *mean S -workload of a j -customer*, denoted V_j^S , as the mean residual service time a current j -customer requires until he first leaves S , where $j \in S$ and $S \subseteq \mathcal{N}$. The V_j^S 's may be computed by solving the linear system

$$V_i^S = \frac{1}{\mu_i} + \sum_{j \in S} p_{ij} V_j^S, \text{ for } i \in S, S \subseteq \mathcal{N}. \quad (0)$$

We also define by extension parameters V_i^S , for $i \in S^c = \mathcal{N} \setminus S$, from the solution of (1). We can now proceed to define the *external traffic intensity for S -customers*, denoted $\rho^0(S)$, by

$$\rho^0(S) = \sum_{j \in S} \lambda_j V_j^S. \quad (2)$$

To explore stability issues for our MQNET model, the above conventional notions of traffic intensity must be supplemented by ideas of peak traffic rates. The *peak traffic intensity from class i into station m* , denoted $R(i, m)$, is the maximum rate at which work brought in by current i -customers can enter that station, i.e.,

$$R(i, m) = \mu_i \sum_{j \in C_m} p_{ij} V_j^{C_m}, \text{ for } i \in \mathcal{N}, m \in \mathcal{M}. \quad (3)$$

The *peak traffic intensity from station m into station m* , denoted $\bar{R}(m, m)$, is the maximum rate at which work can be transferred from station m into m , i.e.,

$$\bar{R}(m, m) = \max_{i \in C_m} R(i, m), \text{ for } m, m \in \mathcal{M}, \text{ with } m = m. \quad (4)$$

Finally, we define the *peak traffic intensity for station m* , denoted $\bar{\rho}^-(m)$, as the maximum rate at which work can be transferred into that station, i.e.,

$$\bar{\rho}^-(m) = \sum_{j \in C_m} V_j^{C_m} + \sum_{m \in \mathcal{M} \setminus \{m\}} \bar{R}(m, m), \text{ for } m \in \mathcal{M}. \quad (5)$$

The network evolution is governed by a *scheduling policy*, which specifies dynamically how servers are allocated to customers. Note that we require of policies that at each time t , each server should be fully allocated to a single customer, should any be available for service. We consider policies that are *nonanticipative* (i.e., decisions may not be based on future information), *stationary* (i.e., decisions depend on the current system state), *nonidling* (i.e., servers cannot idle when they have work to do) and *preemptive* (i.e., the service of a customer may be interrupted at any time and resumed later). Such policies are *admissible*.

The system *state* at time t is given by the following random variables:

- $L_j(t)$: number of j -customers in system at time t .
- $B_j(t)$: 1 if a j -customer is in service at time t ; 0 otherwise.
- $B^m(t)$: 1 if server m is busy at time t ; 0 otherwise.

Our analysis makes extensive use of the following performance measures, defined for each admissible policy:

$$x_j = E[L_j], \quad \text{for } j \in \mathcal{N}, \quad (6)$$

and

$$x_{ij} = E[L_j | B_i = 1], \quad \text{for } i, j \in \mathcal{N} \quad (7)$$

$$x_j^{0m} = E[L_j | B^m = 0], \quad \text{for } j \in \mathcal{N}, m \in \mathcal{M}, \quad (8)$$

where the above expectations are taken under the steady-state distribution of the corresponding stochastic processes. A policy is said to be *stable* if the time-average number of customers in the network is finite. The following key result is due to Bertsimas and Niño-Mora [2]. In its statement, we use the notation $\mathbf{x} = (x_j)_{j \in \mathcal{N}}$, $\mathbf{X} = (x_{ij})_{i,j \in \mathcal{N}}$, $\mathbf{X}^0 = (x_j^{0m})_{j \in \mathcal{N}, m \in \mathcal{M}}$, $\mathbf{1} = (\mathbf{1}_j)_{j \in \mathcal{N}}$ and $\mathbf{D} = \text{Diag}(\mathbf{1})$, where $\mathbf{1} = (\mathbf{1}_j)_{j \in \mathcal{N}}$. We also need $\mathcal{M}(S) = \{m; S \cap C_m \neq \emptyset\}$ and $M(S) = |\mathcal{M}(S)|$.

Theorem 1. Under any admissible stable policy, the performance measures \mathbf{x} , \mathbf{X} and \mathbf{X}^0 satisfy the following:

(a) Conditioning constraints

$$x_j = \sum_{i \in C_m} i x_{ij} + \{1 - M(S)\} x_j^{0m}, \quad \text{for } j \in \mathcal{N}, m \in \mathcal{M}. \quad (9)$$

(b) Flow conservation constraints

$$- \mathbf{x} - \mathbf{x} + (\mathbf{I} - \mathbf{P}) \mathbf{X} + \mathbf{X} (\mathbf{I} - \mathbf{P}) = (\mathbf{I} - \mathbf{P}) \mathbf{1} + (\mathbf{I} - \mathbf{P}) \mathbf{1}. \quad (10)$$

(c) Work decomposition constraints

$$\begin{aligned} \{M(S) - 0(S)\} V_j^S x_j &= \sum_{j \in S} V_j^S + \sum_{i \in S^c \cap (\cup_{m \in \mathcal{M}(S)} C_m)} V_j^S x_{ij} \\ &+ \sum_{i \in S^c} (V_i^S - V_j^S) x_{ij} \\ &+ \sum_{m \in \mathcal{M}(S)} \{1 - (C_m)\} V_j^S x_j^{0m}, \text{ for } S \subseteq \mathcal{N}. \end{aligned} \quad (11)$$

See Bertsimas and Niño-Mora [2] for a full account of theorem 1. Equation (9) is elementary, and arises from consideration of the possible states of station m in the steady state. Constraint (10) was obtained by Bertsimas et al. [4] and by Kumar and Kumar [16]. It may be deduced from application of the flow conservation law $L^- = L^+$ of queueing theory (discovered by Burke [5] and Finch [8]) to all subsystems of the MQNET consisting of single customer classes and also to pairs of classes. The work decomposition result (11) is due to Bertsimas and Niño-Mora [2]. It is of considerable interest in its own right and is deduced from (9) and (10) via algebraic manipulation.

3. Stability of the MQNET model

Before proceeding to issues of optimisation and performance analysis of the system, a fundamental issue concerns whether the MQNETs of interest to us are such that all policies in a given class are stable. This issue has been the subject of a huge research effort. Two kinds of results have emerged: (1) computational tests which seek to construct a Lyapunov function that implies stability, for specific model parameters, from the numerical solution of a linear programming (LP) problem (see, e.g., Kumar and Meyn [15]), and (2) qualitative results, which establish the stability of a family of policies for a restricted network topology under the traditional stability condition (see, e.g., Dai and Weiss [7]).

In this section, we shall present a simple intuitive sufficient stability condition for general Markovian MQNETs: if each station has enough service capacity to cope with its peak traffic intensity, then the network is stable under any admissible policy. No claim is made that this result is best possible, but (i) it is insightful, (ii) it is close to sharp under given conditions and, most importantly, (iii) it seems to be precisely what is required for addressing analytically subsequent performance evaluation and optimisation issues (see sections 4 and 5). In particular, the condition implies a closed-form upper bound on the mean number of customers in the network (and hence on mean customer delay via Little's theorem), uniformly valid under all admissible policies. Our proof draws on theorem 1 and recent results concerning the relation between stability and performance via LP (see Kumar and Meyn [15]).

We state our main result.

Theorem 2 (Global stability condition and performance bound). If the condition

$$(C_m) < 1 \text{ and } \bar{\rho}(m) < 1, \text{ for } m \in \mathcal{M}, \quad (12)$$

holds, then:

- (i) The network is stable under any admissible policy, and the Markov chain that represents its evolution has a geometrically converging exponential moment.
- (ii) Under any admissible policy, the mean number in the network is bounded as follows:

$$E[L_j] \leq \frac{V_j^{C_m}}{(1 - \bar{\rho}(m))V^{C_m}}, \quad (13)$$

where

$$V^{C_m} = \min_{j \in \mathcal{C}_m} V_j^{C_m}.$$

Remarks.

- (1) It is easily shown that if $(C_m) < 1$, then $(C_m) \geq \bar{\rho}(m)$, and so the condition $\bar{\rho}(m) < 1$ is then at least as strong as $(C_m) < 1$. Since also $\bar{\rho}(m) - (C_m) \geq 0$ as $\max_{i \in \mathcal{N} \setminus \mathcal{C}_m, j \in \mathcal{C}_m} p_{ij} \geq 0$ and, moreover, $(C_m) < 1$ is necessary for stability, it follows that condition (12) is close to sharp when the flow between stations is light.
- (2) Notice from theorem 2(i) that condition (12) implies a very strong form of stability.

As indicated above, the LP stability condition of Kumar and Meyn [15] plays a central role in the proof of theorem 2 (given below).

The LP stability condition. The stability condition of Kumar and Meyn [15] refers to a *nonidling performance LP* (shown to be feasible if $(C_m) < 1$ for $m \in \mathcal{M}$), which was introduced independently by Bertsimas et al. [4] and by Kumar and Kumar [16]:

$$Z = \max x_1 + \cdots + x_N$$

subject to $-\mathbf{x} - \mathbf{X} + (\mathbf{I} - \mathbf{P})\mathbf{X} + \mathbf{X}(\mathbf{I} - \mathbf{P}) = (\mathbf{I} - \mathbf{P})\mathbf{1} + (\mathbf{I} - \mathbf{P})\mathbf{1}$ (14)

$$x_j - \sum_{i \in \mathcal{C}_m} i x_{ij} = 0, \quad \text{for } j \in \mathcal{N} \text{ and } m = s(j) \quad (15)$$

$$x_j - \sum_{i \in \mathcal{C}_m} i x_{ij} \leq 0, \quad \text{for } j \in \mathcal{N} \text{ and } m \in \mathcal{M} \setminus \{s(j)\} \quad (16)$$

$$\mathbf{x}, \mathbf{X} \geq 0. \quad (17)$$

The variables $\mathbf{x} = (x_j)_{j \in \mathcal{N}}$ and $\mathbf{X} = (x_{ij})_{i,j \in \mathcal{N}}$ in the performance LP correspond to the performance measures in (6) and (7) above.

Theorem 3 (Kumar and Meyn [15]). If the nonidling performance LP is bounded, then the network is stable under any stationary nonidling policy. Furthermore, the Markov chain that represents its evolution has a geometrically converging exponential moment.

We can now proceed to the proof of our main result.

Proof of theorem 2. We will show that (12) implies that the nonidling performance LP is bounded. Let \mathbf{x} , \mathbf{X} be a feasible solution to this LP. Let matrix $\mathbf{X}^0 = (x_j^{0m})_{j \in \mathcal{N}, m \in \mathcal{M}}$ be given by

$$x_j^{0m} = \frac{1}{1 - (C_m)} x_j - \sum_{i \in C_m} x_{ij}. \quad (18)$$

Notice that from (15) and (16),

$$\begin{aligned} & x_j^{0m} = 0, \quad \text{for } j \in \mathcal{N}, m \in \mathcal{M} \\ \text{and} & \\ & x_j^{0m} = 0, \quad \text{for } s(j) = m. \end{aligned} \quad (19)$$

By the comments following theorem 1, \mathbf{x} , \mathbf{X} and \mathbf{X}^0 satisfy (11). In the case $S = C_m$, equation (11) together with (19) yields

$$\begin{aligned} & 1 - \sum_{j \in C_m} V_j^{C_m} = \sum_{j \in C_m} V_j^{C_m} x_j \\ & = \sum_{j \in C_m} V_j^{C_m} + \sum_{m \in \mathcal{M} \setminus \{m\}} \sum_{i \in C_m} (V_i^{C_m} - V_j^{C_m}) x_{ij}. \end{aligned} \quad (20)$$

Notice further that

$$\begin{aligned} & \sum_{m \in \mathcal{M} \setminus \{m\}} \sum_{i \in C_m} (V_i^{C_m} - V_j^{C_m}) x_{ij} \\ & = \sum_{m \in \mathcal{M} \setminus \{m\}} \sum_{j \in C_m} V_j^{C_m} \sum_{i \in C_m} R(i, m) x_{ij} \end{aligned} \quad (21)$$

$$\begin{aligned} & \sum_{m \in \mathcal{M} \setminus \{m\}} \bar{R}(m, m) \sum_{j \in C_m} V_j^{C_m} \sum_{i \in C_m} x_{ij} \\ & \leq \sum_{m \in \mathcal{M} \setminus \{m\}} \bar{R}(m, m) \sum_{j \in C_m} V_j^{C_m} x_j, \end{aligned} \quad (22)$$

where identity (21) follows from (1) and inequality (22) follows from (16).

Now, combining identity (20) with inequality (22) and condition (12), we obtain

$$\sum_{j \in C_m} V_j^{C_m} x_j \leq \frac{1}{1 - \bar{r}(m)} \sum_{j \in C_m} V_j^{C_m}, \quad \text{for } m \in \mathcal{M}. \quad (23)$$

From (23), it follows that

$$x_j = \frac{c_j V_j^{C_m}}{(1 - \rho_j^{(m)}) V^{C_m}},$$

where

$$V^{C_m} = \min_{j \in \mathcal{C}_m} V_j^{C_m} \min_{j \in \mathcal{C}_m} \frac{1}{\mu_j} > 0, \text{ for } m \in \mathcal{M}.$$

Hence, the nonidling performance LP is bounded, and theorem 2 now follows from theorem 3. \square

Remark. An alternative route to showing stability under condition (12) makes use of the properties of the fluid model approximation (see Dai and Weiss [7]). Under (12), the fluid model can be shown to have a strictly decreasing amount of fluid at each machine. This is enough to guarantee stability.

In sections 4 and 5, we shall assume that the MQNETs which are the subject of our analysis satisfy (12).

4. Approximate optimality of index policies via approximate conservation laws

We suppose that the global stability condition (12) is satisfied and proceed to consider how to schedule in an optimal fashion. Consider a cost structure in which j -customers incur linear holding costs at a rate c_j per unit time in the system. The *optimal scheduling problem* is concerned with finding an admissible scheduling policy that minimises the time-average holding cost rate,

$$c_1 E[L_1] + \cdots + c_N E[L_N], \quad (24)$$

and with evaluating Z^{OPT} , the corresponding minimum cost.

This problem has been solved exactly only in the special case that the network contains a single server, i.e., $M = 1$. For this case, Klimov [14] first showed that the optimal policy is given by a *priority-index* rule: for each class i , one can compute a corresponding *priority index* π_i , such that it is optimal to service at each time a customer with the largest available index. The vector of optimal priority indices $\boldsymbol{\pi} = (\pi_j)_{j \in \mathcal{N}}$ is computed by running Klimov's *adaptive greedy* algorithm – shown below – on input (\mathbf{c}, \mathbf{V}) , where $\mathbf{c} = (c_j)_{j \in \mathcal{N}}$ and $\mathbf{V} = (V_j^S)_{j \in \mathcal{N}, S \subseteq \mathcal{N}}$. The index π_i thus computed represents the maximum rate of decrease in holding cost rate per unit of network processing time for a current i -customer.

Klimov's adaptive greedy algorithm

Input: (\mathbf{c}, \mathbf{A}) , where $\mathbf{c} = (c_j)_{j \in \mathcal{N}}$ and $\mathbf{A} = (A_j^S)_{j \in \mathcal{N}, S \subseteq \mathcal{N}}$.

Output: $(\pi, \bar{\mathbf{y}})$, where $\pi = (\pi_1, \dots, \pi_N)$ is a permutation of \mathcal{N} , $\bar{\mathbf{y}} = (\bar{y}(S))_{S \subseteq \mathcal{N}}$ and $\pi_1 = \pi_1, \dots, \pi_N$.

Step 0. Set $S_1 = \mathcal{N}$; set $\bar{y}(S_1) = \min\{c_i/A_i^{S_1} : i \in S_1\}$;
pick $\pi_1 = \operatorname{argmin}\{c_i/A_i^{S_1} : i \in S_1\}$;
set $\pi_1 = \bar{y}(S_1)$.

Step k. For $k = 2, \dots, N$:

set $S_k = S_{k-1} \setminus \{\pi_{k-1}\}$; set $\bar{y}(S_k) = \min \frac{c_i - \sum_{j=1}^{k-1} A_i^{S_j} \bar{y}(S_j)}{A_i^{S_k}} : i \in S_k$;

pick $\pi_k = \operatorname{argmin} \frac{c_i - \sum_{j=1}^{k-1} A_i^{S_j} \bar{y}(S_j)}{A_i^{S_k}} : i \in S_k$;

set $\pi_k = \pi_{k-1} + \bar{y}(S_k)$.

Step N + 1. For $S \subseteq \mathcal{N}$: set $\bar{y}(S) = 0$, if $S \neq \{S_1, \dots, S_N\}$.

The natural modification of the index policy for the $M = 1$ case which is appropriate to our more complex MQNET model is one which produces indices by running the adaptive greedy algorithm on input (\mathbf{c}, \mathbf{A}) , where $\mathbf{A} = (A_i^S)_{i \in \mathcal{N}, S \subseteq \mathcal{N}}$ is given by

$$A_i^S = \begin{cases} V_i^{S \cap C_m} & \text{if } i \in S \cap C_m; \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

It follows easily from Klimov's [14] result that such an index policy is optimal when the stations in the MQNET are autonomous – i.e., when there is no flow between them. Note also that the choice of matrix \mathbf{A} in [25] arises naturally from the formulation of the optimal scheduling problem in terms of approximate conservation laws, now to be described.

The index result for the $M = 1$ case referred to above was recently recovered by Bertsimas and Niño-Mora [3]. Their approach was (i) to demonstrate that the system satisfies (so-called) *generalised conservation laws* (GCL), (ii) to use this fact to obtain the *performance space* of the system, i.e. the set $\mathbf{x} = (x_j)_{j \in \mathcal{N}}$ (see (7)) which are achievable under admissible policies, and (iii) to recast the optimal scheduling problem as an LP whose feasible space is the above performance space. We cannot follow this program through in the case of our more complex models since they do not satisfy GCL. However, they come close to doing so in a sense described by the following notion.

Definition 1 (Approximate conservation laws). Let $\mathbf{A} = (A_i^S)_{i \in S, S \subseteq \mathcal{N}}$ be a matrix with all $A_i^S > 0$ and $b, \cdot : 2^{\mathcal{N}} \rightarrow \mathcal{R}^+$ be nonnegative set functions, with b strictly positive always. We say that performance vector \mathbf{x} satisfies approximate conservation laws with parameters \mathbf{A} , b and \cdot if the following conditions hold:

(i) Under any admissible scheduling policy,

$$\sum_{j \in S} A_j^S x_j \leq b(S), \quad \text{for } S \subseteq \mathcal{N}. \quad (26)$$

(ii) Under any admissible policy that gives priority to S -customers over S^c -customers,

$$\sum_{j \in S} A_j^S x_j \leq b(S) + \cdot(S), \quad \text{for } S \subseteq \mathcal{N}. \quad (27)$$

Remarks.

- (1) In the special case that $\cdot = 0$, definition 1 reduces to the *generalised conservation laws* introduced by Bertsimas and Niño-Mora [2].
- (2) Notice that definition 1 differs from the original definition of Glazebrook and Garbe [11] in that a less restrictive assumption is required to hold in part (ii). This modification is essential for the application of their approximate conservation laws framework to the models studied in this paper.

When a performance vector satisfies approximate conservation laws, we can solve approximately optimal scheduling problems on linear performance objectives, as in (24). Re-express (24) as

$$c_1 x_1 + \dots + c_N x_N \quad (28)$$

using (6). We shall address such problems via the *achievable region method* (see e.g., the survey paper by Bertsimas [1]). We thus consider the *performance region* \mathcal{X} achievable by performance vector \mathbf{x} under *all* admissible scheduling policies. We can now formulate the optimal scheduling problem as the mathematical program

$$Z^{OPT} = \min\{c_1 x_1 + \dots + c_N x_N : \mathbf{x} \in \mathcal{X}\}. \quad (29)$$

Suppose that the performance vector \mathbf{x} satisfies approximate conservation laws (26) and (27). Then performance region \mathcal{X} is contained in the polyhedron

$$\mathcal{P} = \{\mathbf{x} \in \mathcal{R}_N^+ : \sum_{j \in S} A_j^S x_j \leq b(S), S \subseteq \mathcal{N}\},$$

which yields the following *linear programming (LP) relaxation* of problem (29):

$$Z^{LP} = \min\{c_1 x_1 + \dots + c_N x_N : \mathbf{x} \in \mathcal{P}\}. \quad (30)$$

We obtain performance guarantees for appropriate index policies in theorem 4 by constructing feasible solutions to (29) and the dual of (30) as follows:

1. Run Klimov's adaptive greedy algorithm on input (\mathbf{c}, \mathbf{A}) to obtain an output $(\pi, \bar{\mathbf{y}})$.
2. The proposed heuristic solution for problem (29) is the performance vector of the priority-index scheduling policy that gives higher priority to classes with higher index π_i . We denote by Z^{INDEX} the value achieved by this policy.
3. The proposed feasible solution for the dual of LP relaxation (30) is given by $\bar{\mathbf{y}}$ (see Bertsimas and Niño-Mora [3] for a proof that $\bar{\mathbf{y}}$ is indeed dual feasible).

In what follows, we shall assume, for ease of exposition, that the permutation returned by Klimov's algorithm is $\pi = (1, \dots, N)$, so that

$$1 \quad \dots \quad N.$$

Bertsimas and Niño-Mora [3] further showed that the value of the dual feasible solution $\bar{\mathbf{y}}$ is given by

$$Z^D = \pi_1 b(\{1, \dots, N\}) + (\pi_2 - \pi_1) b(\{2, \dots, N\}) + \dots + (\pi_N - \pi_{N-1}) b(\{N\}). \quad (31)$$

Notice that

$$Z^D \geq Z^{LP} \geq Z^{OPT} \geq Z^{INDEX}. \quad (32)$$

The next result, which reformulates the original approximate optimality result in Glazebrook and Garbe [11], establishes additive and multiplicative performance guarantees for the value Z^{INDEX} of the priority-index heuristic and the value Z^D of the lower bound.

Theorem 4 (Performance guarantees). Suppose the performance vector \mathbf{x} satisfies approximate conservation laws (26) and (27) with parameters \mathbf{A} , b and r . Then

$$(a) \quad Z^{INDEX} \geq Z^{OPT} + \epsilon, \quad (33)$$

and

$$Z^D \geq Z^{OPT} - \epsilon, \quad (34)$$

where

$$\epsilon = \pi_1 (b(\{1, \dots, N\}) - b(\{2, \dots, N\})) + (\pi_2 - \pi_1) (b(\{2, \dots, N\}) - b(\{3, \dots, N\})) + \dots + (\pi_N - \pi_{N-1}) (b(\{N\}) - b(\{N\})). \quad (35)$$

(b) Furthermore,

$$Z^{INDEX} \geq r Z^{OPT} \quad (36)$$

and

$$Z^D \geq \frac{1}{r} Z^{OPT}, \quad (37)$$

where

$$r = 1 + \max_{S \subseteq \mathcal{N}} \frac{b(S)}{b(S)}. \quad (38)$$

Proof. (a) It is shown in Bertsimas and Niño-Mora [3] that \bar{y} is a feasible solution of the dual LP of (30), given by

$$\begin{aligned} (LP_D) \quad Z^{LP} = \text{maximise} \quad & b(S)y(S) \\ & S \subseteq \mathcal{N} \\ \text{subject to} \quad & A_i^S y(S) \leq c_i, \quad \text{for } i \in \mathcal{N}, \\ & S \ni i \\ & y(S) \geq 0, \quad \text{for } S \subseteq \mathcal{N}, \end{aligned}$$

that satisfies the constraints with equality, i.e.,

$$A_i^S \bar{y}(S) = c_i, \quad \text{for } i \in \mathcal{N}. \quad (39)$$

They also show the relation

$$c_i - c_{i-1} = \bar{y}(S_i), \quad \text{for } i \in \mathcal{N},$$

holds, where $S_i = \{i, i+1, \dots, N\}$.

Now let x^{INDEX} be the performance vector corresponding to the proposed index policy. We have

$$\begin{aligned} Z^{INDEX} &= \sum_{i \in \mathcal{N}} c_i x_i^{INDEX} \\ &= \sum_{i \in \mathcal{N}} \sum_{S \ni i} A_i^S \bar{y}(S) x_i^{INDEX} \\ &= \sum_{S \subseteq \mathcal{N}} \bar{y}(S) \sum_{i \in S} A_i^S x_i^{INDEX} \\ &= \sum_{i=1}^N \bar{y}(S_i) \{b(S_i) + c_i\} \\ &= Z^D + \sum_{i=1}^N c_i, \end{aligned} \quad (40)$$

and, by (32), relations (33) and (34) follow. Part (b) follows simply from (31), (38) and (40). \square

Remark. In the special case that $c_i = 0$, it follows from theorem 4 that the proposed index policy is optimal. See Bertsimas and Niño-Mora [3].

5. Approximate conservation laws for MQNETs

Our goal now is to establish approximate conservation laws for our MQNET model. Application of theorem 4 will then yield performance guarantees for a suitably defined index policy.

To pursue this agenda, we require developments of the notions of peak traffic intensity in (3)–(5) above. Let $S \subseteq \mathcal{N}$, $m \in \mathcal{M}(S)$ and $i \in C_m$, with $m \neq m$. We define the (m, m) -peak traffic intensity from class $i \in C_m$ into $S \cap C_m$, denoted $R_{m, m}(i, S \cap C_m)$, as the maximum rate at which work can be transferred from class i into classes in $S \cap C_m$, i.e.,

$$R_{m, m}(i, S \cap C_m) = \mu_i \max_{j \in S \cap C_m} p_{ij} V_j^{S \cap C_m}. \quad (41)$$

Now let $S \subseteq \mathcal{N} \setminus C_m$. We can now define the m -peak traffic intensity from S into $S \cap C_m$, denoted $\bar{R}_m(S, S \cap C_m)$, as the maximum rate at which work can be transferred from classes in S into classes in $S \cap C_m$, i.e.,

$$\bar{R}_m(S, S \cap C_m) = \max_{m \in \mathcal{M}(S)} \max_{i \in S \cap C_m} R_{m, m}(i, S \cap C_m). \quad (42)$$

Theorem 5 (Approximate conservation laws for the MQNET model). Under the global stability condition (12), the MQNET model satisfies approximate conservation laws with parameters \mathbf{A} , b and \hat{b} given by

$$A_i^S = \begin{cases} V_i^{S \cap C_m}, & \text{if } i \in S \cap C_m, \\ 0, & \text{otherwise;} \end{cases}$$

$$b(S) = \max_{m \in \mathcal{M}(S)} \hat{b}(S \cap C_m);$$

and

$$\rho(S) = \max_{m \in \mathcal{M}(S)} \frac{\bar{R}_m(\mathcal{N} \setminus C_m, S \cap C_m)}{1 - \rho_0(S \cap C_m) - \bar{R}_m(\mathcal{N} \setminus C_m, S \cap C_m)} \hat{b}(S \cap C_m),$$

where

$$\hat{b}(S) = \frac{1}{1 - \rho_0(S)} \max_{j \in S} V_j^S.$$

Proof. Let $S \subseteq \mathcal{N}$ and $m \in \mathcal{M}(S)$. From the work decomposition constraints (11) and the non-negativity of the performance measures, we conclude that under any admissible policy (guaranteed stable by theorem 2),

$$\sum_{j \in S \cap C_m} V_j^{S \cap C_m} x_j \leq \frac{\sum_{j \in S \cap C_m} V_j^{S \cap C_m}}{1 - \rho_0(S \cap C_m)} = \hat{b}(S \cap C_m). \quad (43)$$

By now summing both sides of inequality (43) over $m \in \mathcal{M}(S)$, we infer (26) for the above choices of \mathbf{A} , b and \hat{b} .

To obtain (27), consider some policy π that gives priority to S -customers over S^c -customers. Under π , it follows easily from the definition of the performance variables in (7) and (8) that

$$x_{ij} = 0, \text{ if } i \in S^c \cap C_m, j \in S \cap C_m, m \in \mathcal{M}(S) \quad (44)$$

and

$$x_j^{0m} = 0, \text{ if } j \in S \cap C_m. \quad (45)$$

Using (43) and (44) in (11) we obtain

$$\begin{aligned} \{1 - \rho_0(S \cap C_m)\} \sum_{j \in S \cap C_m} V_j^{S \cap C_m} x_j &= \sum_{j \in S \cap C_m} V_j^{S \cap C_m} \\ &+ \sum_{i \in C_m, j \in S \cap C_m} (\mu_i V_i^{S \cap C_m} - \mu_j) V_j^{S \cap C_m} x_{ij}. \end{aligned} \quad (46)$$

Now, using (1), the second term on the right-hand side of (46) can be re-expressed as

$$\begin{aligned} &(\mu_i V_k^{S \cap C_m} p_{ik}) V_j^{S \cap C_m} x_{ij} \\ &\sum_{m \in \mathcal{M}} \sum_{i \in C_m} \sum_{j, k \in S \cap C_m} \max_{i \in C_m} R_{m, m}(i, S \cap C_m) \sum_{j \in S \cap C_m} V_j^{S \cap C_m} \sum_{k \in C_m} x_{kj} \end{aligned} \quad (47)$$

$$\bar{R}_m(\mathcal{N} \setminus C_m, S \cap C_m) \sum_{j \in S \cap C_m} V_j^{S \cap C_m} x_j. \quad (48)$$

We use (41) to infer (47), and (9) and (42) to infer (48). From (46), (48) and the stability condition (12), we deduce that, under π ,

$$\sum_{j \in S \cap C_m} V_j^{S \cap C_m} x_j \leq \frac{\sum_{j \in S \cap C_m} V_j^{S \cap C_m}}{1 - \rho_0(S \cap C_m) - \bar{R}_m(\mathcal{N} \setminus C_m, S \cap C_m)}. \quad (49)$$

Note that (12) guarantees that the denominator in (49) is positive. By now summing both sides of inequality (49) over $m \in \mathcal{M}(S)$ and using (43), we infer (27) for the above choices of \mathbf{A} , b and π . \square

There are a variety of ways in which theorems 4 and 5 can be deployed to provide performance guarantees for the index policy developed above. We shall content ourselves here with two simple results which are immediate.

Corollary 6 (Performance guarantee for index policy).

$$Z^{INDEX} \leq Z^{OPT} \left(1 + \max_{m \in \mathcal{M}} \frac{\bar{R}_m(\mathcal{N} \setminus C_m, C_m)}{1 - \bar{\rho}(m)} \right).$$

Proof. The proof utilises (36) and (38) of theorem 4, together with the set functions b and γ defined in the statement of theorem 5. \square

We now exploit the performance guarantee in corollary 6 to establish asymptotic optimality for the index policy in a limit as the flow between stations becomes lighter. To this end, consider a sequence of networks indexed by n , all of which satisfy (12) and for which

$$\limsup_n \rho^{(n)}(m) < 1 \quad \text{for } m \in \mathcal{M},$$

and

$$\lim_n \max_{s(i), s(j)} p_{ij}^{(n)} = 0.$$

If we use Z_n^{INDEX} and Z_n^{OPT} to denote the expected costs of interest for the network indexed by n , then the following result is immediate from corollary 6.

Corollary 7 (Asymptotic optimality).

$$\min_n \frac{Z_n^{INDEX}}{Z_n^{OPT}} = 1.$$

6. Performance analysis of priority policies via approximate conservation laws

Our goal here is to analyse the performance of some given priority policy π of interest for the MQNET. Quantities related to π will be identified via a superscript. Examples include performance x^π and the objective Z^π . Reasonable goals of performance analysis are (i) the evaluation of Z^π , and (ii) the development of performance guarantees for π in the form of bounds for $Z^\pi - Z^{OPT}$. Both will be discussed here.

Let $S_k = \{k, \dots, N\}$ be the $(N - k + 1)$ classes of highest priority under π . \mathbf{A} , b and γ are as in theorem 5. We introduce \mathbf{y} as the solution to the system of linear equations

$$\sum_{i=1}^k A_k^{S_i} y_i = c_k, \quad \text{if } 1 \leq k \leq N. \quad (50)$$

In theorem 8, we write $u^+ = \max(u, 0)$ and $u^- = \max(-u, 0)$.

Theorem 8 (Evaluation of Z^π). Under the global stability condition (12),

$$Z^\pi = \sum_{k=1}^N (y_k)^\pi(S_k) = Z - \sum_{k=1}^N y_k b(S_k) + \sum_{k=1}^N (y_k)^+(S_k). \quad (51)$$

Proof. Policy π gives priority to S_k -customers over others for all k , $1 \leq k \leq N$. It follows from theorem 5 that

$$b(S_k) = \sum_{i=k}^N A_i^{S_k} x_i = b(S_k) + \sum_{i=k}^N (S_k)_i, \text{ if } 1 \leq k \leq N. \quad (52)$$

But it follows from (50) that

$$Z = \sum_{k=1}^N c_k x_k = \sum_{k=1}^N y_k = \sum_{i=1}^N A_i^{S_k} x_i. \quad (53)$$

The result follows simply from (52) and (53). \square

Before proceeding to the development of interpretable bounds for $Z - Z^{OPT}$, we consider a modification of the feedback matrix $\mathbf{P} = (p_{ij})$. Matrix $\tilde{\mathbf{P}}$ modifies \mathbf{P} by deeming flows between distinct stations to be departures from the system, i.e.,

$$\tilde{p}_{ij} = \begin{cases} p_{ij}, & \text{if } s(i) = s(j), \\ 0, & \text{otherwise,} \end{cases} \quad (54)$$

with $\tilde{p}_{i0} = 1 - \sum_{j \in \mathcal{N}} \tilde{p}_{ij}$. Note that the feedback mechanism internal to each station is unaffected. We shall use $\tilde{\cdot}$ to denote quantities related to the system modified by replacement of \mathbf{P} by $\tilde{\mathbf{P}}$, but with λ and μ unchanged.

It is easy to check that the A -matrix defined in theorem 5 is unaffected by the modification. It then follows that the Gittins indices obtained from the adaptive greedy algorithm are also unchanged, as is \mathbf{y} in (50) above. From theorem 5, we have

$$\tilde{b}(S) = \sum_{m \in \mathcal{M}(S)} \frac{\sum_{j \in S \cap C_m} \tilde{V}_j^{S \cap C_m}}{\sum_{j \in S \cap C_m} V_j^{S \cap C_m}}, \quad (55)$$

where $\tilde{V}_j = \tilde{\lambda}_j / \mu_j$, $j \in \mathcal{N}$, with $\tilde{\lambda}$ the vector of total arrival rates for the modified system. Note that, by construction, $\tilde{\lambda}_j = \lambda_j$, $j \in \mathcal{N}$. We write

$$(S) := b(S) - \tilde{b}(S) = \sum_{m \in \mathcal{M}(S)} \frac{\sum_{j \in S \cap C_m} (\lambda_j - \tilde{\lambda}_j) V_j^{S \cap C_m}}{1 - \lambda_0(S \cap C_m)}. \quad (56)$$

The quantities (S) are natural measures of the connectedness of the original network. Since the modified network has no flow between stations, it follows from theorem 5 that $(S) \geq 0$. It is then a simple consequence of theorem 8 that

$$\sum_{k=1}^N y_k \tilde{b}(S_k) = \tilde{Z},$$

the cost incurred when controlling the modified network by $\tilde{\lambda}$.

The following is an immediate consequence of theorem 8 and the above observations.

Corollary 9.

$$Z - Z^{OPT} = \alpha_1(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}) - \alpha_2(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}),$$

where

$$\alpha_1(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}) = - \sum_{k=1}^N (y_k)^- (S_k) + \sum_{k=1}^N y_k (S_k)$$

and

$$\alpha_2(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}) = \sum_{k=1}^N (y_k)^+ (S_k) + \sum_{k=1}^N y_k (S_k).$$

Corollary 10 (Suboptimality bound for $\tilde{\mathbf{P}}$).

$$Z - Z^{OPT} = (Z^{INDEX} - Z^{OPT}) + (\tilde{Z} - \tilde{Z}^{INDEX}) + \alpha_2(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}) - \alpha_1(INDEX; \mathbf{P}, \tilde{\mathbf{P}}).$$

Proof.

$$Z - Z^{OPT} = (Z - \tilde{Z}) + (\tilde{Z} - \tilde{Z}^{INDEX}) + (\tilde{Z}^{INDEX} - Z^{INDEX}) + (Z^{INDEX} - Z^{OPT}).$$

Corollary 9 is now used to bound the first and third terms on the right-hand side of (57). The result follows. \square

Remark. Corollary 10 essentially solves the problem of producing interpretable performance guarantees for priority policy $\tilde{\mathbf{P}}$, since

- (i) a bound for $(Z^{INDEX} - Z^{OPT})$ is immediately available from theorems 4 and 5. This bound, aggregated with $\alpha_2(\tilde{\mathbf{P}}; \mathbf{P}, \tilde{\mathbf{P}}) - \alpha_1(INDEX; \mathbf{P}, \tilde{\mathbf{P}})$ in corollary 10, gives a natural measure of the degree of connectedness of the network;
- (ii) by the above construction, $\tilde{Z} - \tilde{Z}^{INDEX}$ has an interpretation as the suboptimality of policy $\tilde{\mathbf{P}}$ when applied to the modified system (consisting of autonomous, unconnected stations) for which our index policy is known to be optimal, by Klimov [14]. An index-based bound for $\tilde{Z} - \tilde{Z}^{INDEX}$ is thus available from the work of Glazebrook and Garbe [10] in the form of a natural measure of the extent to which $\tilde{\mathbf{P}}$ departs from the index policy.

Applying (i) and (ii) to corollary 10 yields a performance guarantee for $\tilde{\mathbf{P}}$ which combines a measure of the degree of connectedness of the network with a measure of the extent to which $\tilde{\mathbf{P}}$ departs from the natural index policy.

References

- [1] D. Bertsimas, The achievable region in the optimal control of queueing systems; formulations, bounds and policies, Queueing Systems 21(1995)337–389.
- [2] D. Bertsimas and J. Niño-Mora, Optimization of multiclass queueing networks with changeover times via the achievable region approach: Part I, the multi-station case. Working Paper, Operations Research Center, M.I.T., 1996.

- [3] D. Bertsimas and J. Niño-Mora, Conservation laws, extended polymatroids and multi-armed bandit problems; a polyhedral approach to indexable systems, *Mathematics of Operations Research* 21 (1996)257–306.
- [4] D. Bertsimas, I. Paschalidis and J. Tsitsiklis, Optimisation of multiclass queueing networks: Polyhedral and nonlinear characterizations of achievable performance, *Annals of Applied Probability* 4(1994) 43–75.
- [5] P.J. Burke, The output of a queueing system, *Operations Research* 4(1956)699–704.
- [6] J.A. Buzacott and J.G. Shanthikumar, *Stochastic Models of Manufacturing Systems*, Prentice-Hall, Englewood Cliffs, NJ, 1993.
- [7] J.G. Dai and G. Weiss, Stability and instability of fluid models for reentrant lines, *Mathematics of Operations Research* 21(1996)115–134.
- [8] P.D. Finch, On the distribution of queue size in queueing problems, *Acta Mathematica Hungarica* 10(1959)327–336.
- [9] E. Gelenbe and I. Mitrani, *Analysis and Synthesis of Computer Systems*, Academic Press, London, 1980.
- [10] K.D. Glazebrook and R. Garbe, Reflections on a new approach to Gittins indexation, *Journal of the Operational Research Society* 47(1996)1301–1309.
- [11] K.D. Glazebrook and R. Garbe, Almost optimal policies for systems which almost satisfy conservation laws, *Annals of Operations Research* 92(1999), this volume.
- [12] M.X. Goemans and D.P. Williamson, The primal–dual method for approximation algorithms and its application to network design problems, *Approximation Algorithms* (1996)144–191.
- [13] J.M. Harrison and L.M. Wein, Scheduling networks of queues: Heavy traffic analysis of a simple open network, *Queueing Systems* 5(1989)265–280.
- [14] G.P. Klimov, Time sharing service systems I, *Theory of Probability and its Applications* 19(1974) 532–551.
- [15] P.R. Kumar and S.P. Meyn, Duality and linear programs for stability and performance analysis of queueing networks and scheduling policies, *IEEE Transactions on Automatic Control* 41(1996)4–16.
- [16] S. Kumar and P.R. Kumar, Performance bounds for queueing networks and scheduling policies, *IEEE Transactions on Automatic Control* 39(1994)1600–1611.
- [17] C.H. Papadimitriou and J.N. Tsitsiklis, The complexity of optimal queueing network control, Working Paper, Laboratory for Intelligent Decision Systems, 1993.