# Dynamic Distributed Scheduling
# in Random Access Networks

Alexander L. Stolyar

Bell Labs, Lucent Technologies

Murray Hill, NJ 07974

stolyar@research.bell-labs.com

## Abstract

We consider a model of random access (Slotted-Aloha-type) communication networks of general topology. Assuming network links receive exogenous arrivals of packets for transmission, we seek dynamic distributed random access strategies whose goal is to keep all network queues stable. We prove that two dynamic strategies, which we collectively call Queue length based Random Access (QRA), ensure stability as long as the rates of exogenous arrival flows are within the network *saturation rate region*.

The first strategy, QRA-I, can be viewed as a random-access-model counterpart of MaxWeight scheduling rule, while the second one, QRA-II, is a counterpart of the Exponential (EXP) rule. The two strategies induce different dynamics of the queues in the fluid scaling limit, which can be exploited for the Quality-of-Service control in applications.

# 1  Introduction

In this paper we consider a model of random access (Slotted-Aloha-type) communication networks of general topology. Assuming network links receive exogenous arrivals of packets for transmission, we seek dynamic distributed random access strategies whose goal is to keep the network queues stable. (Queues are formed by packets waiting for transmission at the links.) We prove that two dynamic strategies, which we collectively call Queue length based Random Access (QRA), ensure stability as long as the rates of exogenous arrival flows are within the network *saturation rate region*.

Our basic network model (formally defined in Section 3) is that of [14], which in turn is a slight generalization of those in [15, 24]. (A closely related - but different - model was considered in [13]. All the models mentioned above generalize the classical Slotted Aloha system [17].) Very informally, there is a finite set $\mathcal{N}$ of nodes $n$, and a finite set of directed communication links $i = (n, m)$ connecting some of the node pairs. Time is slotted. In each slot, node $n$ transmits a "packet" (or, *accesses channel*) with probability $p_n$, and it chooses one of its outgoing links $i = (n, m)$ to transmit on with probabilities $p_{nm}/p_n$ summing up to 1. Link access probabilities $p_i$ are chosen generally speaking dynamically, but once they are chosen for a slot, each node transmits independently of the others. A general interference structure between the nodes is defined, so that some of the simultaneous transmissions on the network links can fail due to the interference from other transmissions. Thus a set of (constant) link access probabilities $p = \{p_i\}$ determines the corresponding set of link throughputs $\mu = \{\mu_i\}$ - this dependence is denoted by $\mu(p)$.

The network *saturation throughput region* $M$ for the above model is the set of all vectors $\mu(p)$ for all possible access probability sets $p$, along with all non-negative vectors $x$ which can be dominated by such $\mu(p)$. In other words, region $M$ is the set of all link throughput vectors which can be achieved by static random access strategies (i.e. those using constant access probabilities), under the assumption that all link queues are "saturated" (i.e. have unlimited number of packets available for transmission at any link at any time). Paper [14] shows that the outer (Pareto) boundary of the saturation throughput region $M$ is character-ized essentially as the set of throughput vectors maximizing (over $M$) the system objective function $\sum_i w_i \log \mu_i$, for all possible sets $w = \{w_i\}$ of non-negative weights assigned to the links. (See [13] for an analogous result for a related model.) An important property of the optimal solution $\mu^*$ for a given set of link weights $w$ is that the corresponding set of access probabilities $p = p(w)$ (such that $\mu^* = \mu(p)$) can be determined in a distributed fashion; namely, to determine "its own" access probabilities each node will only need to "know" the weights of the nodes in its appropriately defined "local neighborhood." (This property was established earlier in [15] for the case when all $w_i = 1$.) It was also suggested in [14] (and [13]) that a network can be controlled by changing link weights *dynamically* to satisfy some desired constraints of the link throughput allocation while keeping such allocation efficient (i.e., Pareto optimal). In fact, a specific algorithm was proposed in [14] which seeks to max-imize system utility $\sum_i \alpha_i \log \mu_i$ subject to the desired lower bounds on the link rates, and shown to have good performance.

In this paper we study the network model described above, but consider the case when each link $i$ receives a flow of exogenous packet arrivals at the average rate $\lambda_i$. The question is how to set access probabilities dynamically so that all link queues are stochastically stable. Let $Q_i(t)$ be the queue length at link $i$ in time slot $t$, and dynamic link weight $w_i(t)$, at time $t$, is a function of $Q_i(t)$. A *QRA algorithm* uses (dynamic) link weights $w_i(t)$, and sets access probabilities in each slot $t$ to $p(w(t))$. We consider two different variants of QRA algorithm. The first, QRA-I, uses weights $w_i(t) = \alpha_i + \gamma_i Q_i(t)^\beta$, with parameters $\alpha_i \geq 0$, $\gamma_i > 0$, $\beta > 0$. The second algorithm, QRA-II, uses weights $w_i(t) = \alpha_i \exp[\gamma_i Q_i(t)]^\kappa$, $\alpha_i > 0$, $\gamma_i > 0$, and $\kappa \in (0, 1)$. Our main result is that both QRA-I and QRA-II algorithms ensures stability of the queues as long as the input rates vector $\lambda = \{\lambda_i\}$ of input rates lies (strictly) within the saturation throughput region $M$. We note that QRA-I algorithm was essentially introduced in [14], but the algorithm stability issue was *not* addressed there; QRA-II algorithm is new.

The reason why we consider two different algorithm is because, although they both ensure stability (as long as $\lambda$ is within $M$), they induce different behavior of the queues. This is demonstrated by our analysis of their fluid sample paths (FSP), which are roughly the limits $q(t)$, $t \geq 0$, of "fluid-scaled" processes $(1/r)Q(rt)$, $t \geq 0$, as $r \to \infty$. One can say that QRA-I is the random-access-model counterpart of much studied MaxWeight-type scheduling algorithms. (MaxWeight algorithms were originally introduced in [23], and then extended and generalized to accommodate a large variety of models and scenarios. Cf. [22, 7] for recent reviews.) The dynamics of an FSP under QRA-I is such that a Lyapunov function of the form $\sum_i c_i \gamma_i q_i(t)^{\beta+1}$ is non-increasing. In contrast, FSPs under QRA-II algorithm are such that $\max_i \gamma_i q_i(t)$ is non-increasing. In this sense, QRA-II is the counterpart of the "exponential" (EXP) scheduling rule [19], which has the same property and is known to keep the values of $\gamma_i Q_i(t)$ roughly equal in the *heavy traffic* limit [20]. Thus, QRA-II algorithm (unlike QRA-I) allows one, to a certain extent, "directly control" the ratios of different queues - a useful feature for *Quality-of-Service* control in applications.

In terms of technique, the stability analysis od QRA-II is more involved than that of QRA-I. Similarly to the situation with EXP algorithm, the conventional fluid scaling, leading to FSPs in the limit, is insufficient, and we additionally need to consider a "local fluid scaling," leading to *local fluid sample paths* (LFSP). Analysis of LFSP dynamics for QRA-II (in the proof of Theorem 4) is substantially different from previous analyses of LFSPs under EXP algorithm in that it requires a completely different Lyapunov function. We believe that this part of our work is of independent interest.

We have to clarify why one would need QRA algorithms at all, and not simply use previously known MaxWeight and EXP algorithms for our model. Both MaxWeight and EXP algorithms would have the following form: in each time slot choose access probabilities so as to

$$\text{maximize} \sum_i w_i(t)\mu_i, \tag{1}$$

with weights $w_i(t)$ defined as for QRA-I and QRA-II, respectively. Such an algorithm will ensure stability of the queues as long as input rate vector $\lambda$ lies within system *maximum stability region $V$*. Region $V$ is typically strictly larger than our saturation rate region $M$,

because it is defined as the largest region within which stability is feasible at all under any strategy, including strategies utilizing global and instantaneous state sharing and coordination between network nodes. However, due to the fact that "instantaneous service rates" $\mu_i$ enter the sum in (1) *linearly*, solving (1) in each time slot would typically involve global (combinatorial) optimization, not allowing a distributed solution. In contrast, a QRA algorithm chooses access probabilities which

$$\text{maximize} \quad \sum_i w_i(t) \log \mu_i, \tag{2}$$

and this *can* be done in a distributed fashion. Thus, although QRA algorithms guarantee stability within a typically smaller region of input rates, they are much more easily implementable in practice.

We would like to mention one more line of previous research (see [3] and references therein), which aims at characterizing *stability region $S$ under static random access strategies* (i.e. those with constant access probabilities) *for the classical Slotted Aloha system* (where all links interfere with each other). It is easy to see that the closure of $S$ contains the saturation throughput region $M$, that is $M \subseteq \bar{S}$. It has long been conjectured that for Slotted Aloha in fact $M = \bar{S}$. This conjecture was proved in brilliant work [3], but only for the case when exogenous arrival processes are dependent in a special way.

The rest of the paper is organized as follows. Section 2 introduces basic notation. Our network model is described in Section 3, and the queueing stability problem for the case of exogenous arrival processes is defined in Section 4. Section 5 defines saturation throughput region. The QRA algorithms and our main stability result (Theorem 1) are presented in Section 6. Section 7 contains the proof of Theorem 1; the key part of this proof (Theorem 2) is then proved separately for QRA-I and QRA-II in Sections 8 and 9.

## 2   Basic Notation and Conventions

We use the notations $R$, $R_+$ and $R_{++}$ for the sets of real, real non-negative and real positive numbers, respectively. Corresponding $I$-times product spaces are denoted $R^I$, $R^I_+$, and $R^I_{++}$. The space $R^I$ is viewed as a standard vector-space, with elements $x \in R^I$ being row-vectors $x = (x_1, \ldots, x_I)$; $x \cdot y$ is scalar product, $\|x\| = (x \cdot x)^{1/2}$ is Eucledian norm, inducing standard metric.

Vector equalities and inequalities are understood componentwise. Sometimes, where it cannot cause confusion, we slightly abusive notation by applying log and exp componentwise,

$$\log x = (\log x_1, \ldots, \log x_I), \quad \exp x = (\exp x_1, \ldots, \exp x_I),$$

and by writing

$$\gamma \times q = (\gamma_1 q_1, \ldots, \gamma_I q_I)$$

for the componentwise product of vectors.

4

For a scalar function $F(x)$, $x \in R^I$, and subset $M \in R^I$, $x^* \in \arg\max_{x \in M} F(x)$ means that $x^*$ maximizes $F(x)$ within $M$. Abbreviation *u.o.c.* means *uniform on compact sets* convergence of functions. Usually, we will consider this convergence for functions (or vector-functions) of the time $t \in [0, \infty)$, in which case u.o.c. convergence means uniform convergence in $[0, b]$ for any $b \geq 0$. We denote $[z]^- \doteq \min\{z, 0\}$, $[z]^+ \doteq \max\{z, 0\}$, and for a non-negative $a$ define $[z]_a^+$ as $z$ if $a > 0$ and $[z]^+$ if $a = 0$.

# 3  Basic Model

We consider the basic model of [14], which is a generalized version of the model of [15, 24]. The system consists of a finite set $\mathcal{N} = \{1, 2, \ldots, N\}$ of *nodes*, and operates in discrete time, with time slots indexed by $t = 0, 1, 2, \ldots$. Let $\mathcal{D}_n \subseteq \mathcal{N} \setminus n$ denote the subset of nodes to which node $n$ has data to send. A node $n$ at any time $t$ may attempt transmission of one unit of data, called packet, or *customer*, to one of the nodes $m \in \mathcal{D}_n$. When this happens, we say that node $n$ makes transmission attempt on the *link* $(n, m)$. We will denote by

$$\mathcal{I} \doteq \{(n, m) \mid n \in \mathcal{N}, \ m \in \mathcal{D}_n\}$$

the set of all system links, and by $I$ its cardinality (i.e., the total number of links). Throughout the paper, for brevity, we often (but not always) denote links $(n, m) \in \mathcal{I}$ by a single index $i$.

We make the following additional model assumptions:
(a1) A node cannot simultaneously (i.e., within the same slot) transmit on two or more different links.
(a2) If a node transmits in a slot, any simultaneous attempt to transmit to this node will fail.
(a3) If there are two or more simultaneous transmissions to a node, they all collide and fail.
(a4) Any transmission attempt by node $n$ will interfere with and "erase" any attempt to receive a packet at any of the nodes within some subset of $\mathcal{N}$, denoted by $\mathcal{N}_n$. (The model of [15, 24] additionally assumes that $m \in \mathcal{N}_n$ implies $n \in \mathcal{N}_m$.)

Note that assumptions (a3) and (a4) imply that $\mathcal{D}_n \subseteq \mathcal{N}_n$. (In other words, a transmission attempt by node $n$ may interfere with receiving at more nodes than it actually sends traffic to.) Also, by assumption (a2), $n \in \mathcal{N}_n$ for all $n$.

We consider a class of *random access* ("Slotted Aloha-type") transmission schemes, defined as follows. In each time slot $t$, each node $n$ chooses the set of *link access probabilities* $p_{nm}$, $m \in \mathcal{D}_n$, such that

$$p_n \doteq \sum_{m \in \mathcal{D}_n} p_{nm} \leq 1. \tag{3}$$

Node $n$ attempts a transmission of one packet (customer) in slot $t$ with probability $p_n$ (which can be called node access probability) independently of other nodes, and when it does

transmit it chooses a particular link $m \in \mathcal{D}_n$ to transmit on, also randomly, with probabilities $p_{nm}/p_n$. Resulting link transmission success probabilities (or, average link *throughputs*) in the slot are given by

$$\mu_{nm} = p_{nm} \prod_{k:\ m \in \mathcal{N}_k,\ k \neq n} (1 - p_k). \qquad (4)$$

We emphasize that *access probabilities $p_{nm}$ may depend on time, past history, and be mutually dependent across nodes and links. However, once $p_{nm}$'s are chosen for a given slot, the node transmission attempts are independent.*

If random access scheme is *static*, that is probabilities $p_{nm}$ stay constant in time, $\mu_{nm}$'s are time-average link throughputs.

The dependence of the set (vector) of link throughputs $\mu = (\mu_i,\ i \in \mathcal{I}) \in R_+^I$ on the set (vector) of link access probabilities

$$p \in \mathcal{P} \doteq \{(p_i,\ i \in \mathcal{I}) \in [0,1]^I \mid (3) \text{ holds}\},$$

given by (4), will be denoted by $\mu(p)$. Clearly, function $\mu(p)$ is continuous.

# 4   Stability Problem

Suppose there is an exogenous arrival process of packets (customers) of average rate $\lambda_i$, to be transmitted on link $i$. To simplify exposition, we will assume that the arrival processes for different links are mutually independent, and the process for link $i$ is given by an i.i.d. sequence $A_i(t),\ t = 0, 1, 2, \ldots$, of non-negative integer random variables, where $A_i(t)$ is the number of packets (customers) arriving in slot $t$. (Obviously, $\lambda_i = EA_i(t)$.) Also for simplicity, we assume that $P\{A_i(t) = 0\} > 0$. (These assumptions on the arrival processes can be replaced by much more general Markov assumptions, e.g. those in [2]. Essentially, all we will need is that the underlying stochastic process describing evolution of the system under the strategy we will introduce later in Section 6 is Markov.)

Customers waiting for transmission form queues, one queue per each link. Queue length (number of waiting customers) at the link $i$ at time $t$ is denoted by $Q_i(t)$.

The problem is to find a dynamic random access strategy, i.e. a dynamic rule for choosing access probabilities $p_i$, such that the link queues remain stochastically stable. For a random access strategy such that the queueing process $Q(t) = (Q_i(t),\ i \in \mathcal{I}),\ t \geq 0$, is a Markov chain (as will be the case for the QRA strategies we introduce later in Section 6), stochastic stability is understood as ergodicity of this Markov chain.

Obviously, stability cannot be expected for arbitrary input rates $\lambda = (\lambda_i,\ i \in \mathcal{I})$. The QRA strategies of Section 6 is such that they ensure stability, as long as $\lambda$ lies within the saturation throughput region, which we introduce next.

# 5 Saturation Throughput Region

We define the system *saturation throughput region $M$*, which we will often call simply *throughput region*, as the set of all non-negative vectors, which can be majorized by vectors of the form $\mu(p)$, namely,

$$M \doteq \{\mu' \in [0,1]^I \mid \mu' \leq \mu(p) \text{ for some } p \in \mathcal{P}\}. \tag{5}$$

Region $M$ has the following simple interpretation. Suppose, each link is "saturated," that is there is unlimited number of packets available for transmission. A vector $\mu'$ is within $M$ if there exists a static random access strategy, with some constant vector $p$ of access probabilities, that provides average throughput of at least $\mu'_i$ on each link $i$. It is easy to observe that region $M$ is a compact subset of the positive quadrant $R_+^I$.

It is shown in [14] that the subset of maximal elements (or, Pareto boundary) of $M$, i.e., set

$$M^* \doteq \{\mu^* \in M \mid \mu^* \leq \mu' \in M \text{ implies } \mu' = \mu^*\}, \tag{6}$$

is essentially equal to the set of points $\mu \in M$ maximizing ("weighted proportional fairness") objective $\sum_{i \in \mathcal{I}} w_i \log \mu_i$, for different non-negative weights $w_i$. The choice of access probabilities $p$ such that $\mu(p)$ maximizes $\sum w_i \log \mu_i$ over $M$ is given by Proposition 1 below, which is Theorem 1 of [14], which in turn is a generalization of the corresponding result in [15].

For each $n \in \mathcal{N}$, let us denote by

$$\mathcal{S}_n \doteq \{(\ell, k) \mid k \in \mathcal{D}_\ell, \ k \in \mathcal{N}_n\}$$

the set of all links $(\ell, k)$ which either originate at $n$ or are such that a transmission by node $n$ interferes with that on $(\ell, k)$.

**Proposition 1** *For arbitrary set of positive weights $w = (w_{nm}, \ (n,m) \in \mathcal{I}) \in R_{++}^I$, there exists a unique set of access probabilities $p \in \mathcal{P}$ which maximizes the function*

$$\sum_{(n,m) \in \mathcal{I}} w_{nm} \log \mu_{nm}(p). \tag{7}$$

*The optimal $p$ is given by:*

$$p_{nm} = \frac{w_{nm}}{\sum_{(\ell,k) \in \mathcal{S}_n} w_{\ell k}}. \tag{8}$$

The dependence of the set (vector) of access probabilities $p \in \mathcal{P}$ on the set (vector) of positive link weights $w \in R_{++}^I$, given by (8), will be denoted by $p(w)$. We will extend the domain of $p(w)$ for all $w \in R_+^I$, using the convention that $p_{nm} = 0$ when $w_{nm} = 0$. (Expression (8) is well defined when $w_{nm} > 0$, but may not be when $w_{nm} = 0$.) Clearly, $p(w)$ is continuous at any point $w \in R_{++}^I$. It is not necessarily continuous at points $w$ with some 0 components. However, some continuity properties, namely those in Proposition 2(iii) below, do hold and will suffice for our purposes.

**Proposition 2** *(i) Function $p(w)$, $w \in R_+^I$, is invariant with respect to scaling of $w$ by a positive constant.*

*(ii) Suppose $w' \in R_+^I$. Then,*

$$\mu' \in \arg\max_{x \in M} \sum w'_{nm} \log x_{nm} \quad \text{if and only if} \quad \mu'_{nm} = \mu_{nm}(p(w')) \text{ for all } (n,m) \text{ with } w'_{nm} > 0. \tag{9}$$

*(We use convention that $0(-\infty) = 0$.)*

*(iii) Suppose $w \to w' \in R_+^I$. Then,*
*(iii.1) If $w'_{nm} = 0$ and $w'_{\ell k} > 0$ for at least one $(\ell, k) \in \mathcal{S}_n$, then $p_{nm}(w) \to 0 = p_{nm}(w')$. In particular, if $w'_{nm} = 0$ and $w'_{nk} > 0$, then $p_{nm}(w) \to 0$.*
*(iii.2) If $w'_{nm} > 0$, then*

$$p_{nm}(w) \to p_{nm}(w') \text{ and } \mu_{nm}(p(w)) \to \mu_{nm}(p(w')). \tag{10}$$

The proof is straightforward - we omit it.

**Remark.** Expression (8) can be equivalently rewritten as

$$p_{nm} = \frac{w_{nm}}{\sum_{k \in \mathcal{N}_n} W_k^{in}}, \tag{11}$$

where

$$W_k^{in} \doteq \sum_{\ell:\ k \in \mathcal{D}_\ell} w_{\ell k} \tag{12}$$

is the sum of the weights of all links "incoming" to node $k$. As explained in [14], given the set of weights $w$, the calculation of access probabilities $p_{nm}$ according to expression (11) can be done in "distributed way," namely, nodes will only need to "know" weights of their own links and exchange a minimum amount of information with their "neighboring" nodes.

# 6  Queue Length Based Dynamic Strategies. Stability Results

Consider the model described in Sections 3-4. Without loss of generality - rather with a *gain* of generality - we assume that a transmission attempt (with non-zero probability) is *allowed* on any link $i$ in any slot $t$, regardless of whether or not there are customers in the queue available for transmission in that slot. (Any transmission attempt interferes with transmissions on the other links in the usual way.) In particular, it is possible to have a "successful transmission attempt" on link $i$ that transmits *no* customer from the corresponding queue, because none was available. We also adopt a convention that the $A_i(t)$

8

customers arriving at link $i$ in slot $t$ are *not* counted into the queue length $Q_i(t)$ at time $t$, but *are* immediately available for transmission at time $t$. (This convention is non-essential, made just to make expressions "cleaner.") Given our conventions, we obviously have the following recurrence for the queue length:

$$Q_i(t+1) = [Q_i(t) + A_i(t) - h_i(t)]^+, \ t = 0, 1, 2, \ldots, \tag{13}$$

where $h_i(t) = 1$ if there was a *successful* transmission attempt on link $i$ at time $t$, and $h_i(t) = 0$ otherwise.

We are going to study the following two dynamic strategies, which will be collectively referred to as QRA algorithms (Queue length based Random Access):

*Each node $n$ maintains* dynamic weights $w_i(t)$, $t = 0, 1, 2, \ldots$, of "its" outgoing links, de-*pending on the corresponding queue lengths. Each node $n$ sets its access probabilities in slot $t$ according to the expression (11). In other words, the set of access probabilities in the net-work at time $t$ is given by $p(w(t))$.* **QRA-I** *algorithm uses weights* $w_i(t) = \alpha_i + \gamma_i Q_i(t)^\beta$, *with parameters* $\alpha_i \geq 0$, $\gamma_i > 0$, $\beta > 0$. **QRA-II** *uses weights* $w_i(t) = \alpha_i \exp[\gamma_i Q_i(t)]^\kappa$, *with parameters* $\alpha_i > 0$, $\gamma_i > 0$, *and* $\kappa \in (0, 1)$.

Given our assumptions on the arrival processes (in Section 4), it is clear that, under both QRA algorithms, the queue length (vector-) process $Q(t)$, $t = 0, 1, 2, \ldots$, is an irreducible (and aperiodic) Markov chain with countable state space.

**Theorem 1** *Suppose vector of input rates $\lambda$ is such that*

$$\lambda < \mu^* \text{ for some } \mu^* \in M. \tag{14}$$

*Then,*
*(i) Markov chain $Q = \{Q(t), \ t = 0, 1, 2, \ldots\}$ is ergodic under the QRA-I policy;*
*(ii) Markov chain $Q$ is ergodic under the QRA-II policy, and the additional (exponential moment) assumption on the input flows:*

$$Ee^{a_1 A_i(0)} < \infty, \text{ for some } a_1 > 0 \text{ and all } i. \tag{15}$$

**Remark 1.** If input processes are not i.i.d., condition (15) can be relaxed to (36) - the latter condition is the one actually used in the proof.

**Remark 2.** In essence, QRA-I algorithm was introduced in [14], in a somewhat different context, where the arrival processes and queue lengths are "virtual," and used to enforce minimum desired throughputs on the links. However, the queueing stability problem was *not* addressed in [14].

**Remark 3.** As already discussed in the introduction, QRA-II algorithm is the "random-access-model" counterpart of EXP algorithm [19, 20]. (The latter algorithm applies to a different model.) The queue weights $w_i(t)$ used in the original EXP algorithm have the form

$$w_i(t) = \alpha_i \exp \frac{[\gamma_i Q_i(t)]}{1 + [\overline{\gamma Q}(t)]^{1-\kappa}},$$

where $\overline{\gamma Q}(t) = (1/I)\sum_i \gamma_i Q_i(t)$. The form $w_i(t) = \alpha_i \exp[\gamma_i Q_i(t)]^\kappa$ was later used in [10] (again, for a different model). These two forms result in equivalent behavior of fluid sample paths (FSP) and local fluid sample paths (LFSP), which are used in our stability proofs, and thus, in principle, either form can be used in QRA-II. The latter form is better suited for a distributed implementation. We emphasize that previous stability analyses in [19, 20, 10] are for a different model, and do not apply to the model of this paper.

# 7 Proof of Theorem 1

We prove stability using fluid limit technique [18, 6, 5, 21, 8]. (For the application of this technique in a discrete-time setting, close to that of this paper, cf. [2].)

Let $Q^{(r)} = (Q^{(r)}(t), \; t = 0, 1, 2, \dots)$, denote a queue length process $Q$ with a fixed initial condition such that $\|Q^{(r)}(0)\| = r$, $r > 0$. In the analysis to follow, all variables associated with a process $Q^{(r)}$ will be supplied with the upper index $(r)$. It will be convenient to extend the definition of the process $Q^{(r)}$ to *continuous time* $t \geq 0$ by adopting the convention that $Q(t) = Q(\lfloor t \rfloor)$.

The following result follows from the state dependent Lyapunov-type stability criterion for countable Markov chains, obtained first by Malyshev and Menshikov [16]. (In the specific form (16) a Markov chain ergodicity criterion was introduced in [18].)

**Proposition 3** *Suppose there exist constants $\epsilon > 0$ and $T > 0$ such that for any sequence of processes $\{Q^{(r)}, r \uparrow \infty\}$, we have*

$$\limsup_{r \to \infty} E[\frac{1}{r}\|Q^{(r)}(rT)\|] \leq 1 - \epsilon \; . \tag{16}$$

*Then the (original, discrete time) Markov chain $Q$ is ergodic.*

To verify condition (16) of Proposition 3, fluid limit technique introduces the sequence of fluid-scaled processes,

$$q^{(r)}(t) \doteq \frac{1}{r}Q^{(r)}(rt), \; t \geq 0. \tag{17}$$

Note that it suffices to verify (16) under the additional condition that the sequence of rescaled initial states $(1/r)Q^{(r)}(0)$ converges, which is equivalent to

$$q^{(r)}(0) \to q(0) \text{ as } r \to \infty, \text{ where } \|q(0)\| = 1. \tag{18}$$

Then, we establish the following result.

10

**Theorem 2** *There exist constants $\epsilon > 0$ and $T > 0$, for which the following holds. Consider a fixed sequence of rescaled processes $\{q^{(r)}, r \uparrow \infty\}$, satisfying condition (18). Then, all processes of the sequence can be constructed on a common probability space, such that the following holds with probability 1. Any subsequence of the sequence of realizations of $\{q^{(r)}, r \uparrow \infty\}$, has in turn a further subsequence u.o.c. converging to a trajectory $q = (q(t), \ t \geq 0)$, which we call a fluid sample path (FSP), and which, in particular, satisfies the following properties:*

$$\|q(0)\| = 1, \tag{19}$$

$$function \ q(t), \ t \geq 0, \ is \ Lipschitz \ continuous, \tag{20}$$

$$\|q(T)\| \leq 1 - \epsilon. \tag{21}$$

Theorem 2 will be proved separately for QRA-I (under assumptions of Theorem 2(i)) and QRA-II (under assumptions of Theorem 2(ii)), in Sections 8 and 9, respectively. Fluid sample paths will be defined differently for QRA-I and QRA-II. We emphasize that not only the equations describing their dynamics will be different, but their *definitions* are in fact somewhat different as well.

Once Theorem 2 is established (for both QRA-I and QRA-II), this completes the proof of Theorem 1. Indeed, for any fixed $T > 0$ the uniform integrability of the family of random variables $(1/r)\|Q^{(r)}(rT)\|$ (indexed by $r$, as in Proposition 3) is easily established, using majorization of the queue lengths by the cumulative arrival processes [18, 6]. This fact and Theorem 2 verify the assertion of Proposition 3.

# 8   Proof of Theorem 2 for QRA-I

## 8.1   Probability space construction and other preliminaries

Let us denote by

$$F_i^{(r)}(t) \doteq \sum_{s=0}^{t-1} A_i(s), \ t = 0, 1, 2, \ldots, \quad i \in \mathcal{I},$$

the total number of customer arrivals at link $i$ by (and excluding) integer time $t$. We also denote by

$$\hat{F}_i^{(r)}(t) = \sum_{s=0}^{t-1} h_i(s), \ t = 0, 1, 2, \ldots, \quad i \in \mathcal{I},$$

the total number of successful transmissions on link $i$ by (and excluding) time $t$.

Without loss of generality, we will assume that random transmission attempt decision by node $n$ at time $t$, given its "current" (depending on $t$) set of access probabilities $p_{nm}, \ m \in$

$\mathcal{D}_n$, is determined by the random variable $y_n(t)$, uniformly distributed in $[0, 1]$. (Random variables $y_n(t)$ are mutually independent across all $n$ and $t$.) Namely, node $n$ assumes some fixed ordering of the nodes in $\mathcal{D}_n$: $m_1, \ldots, m_\ell$. Then, if $y_n(t) \in \phi_{n,m_1} \doteq [0, p_{n,m_1}]$, node $n$ attempts transmission on link $(n, m_1)$; if $y_n(t) \in \phi_{n,m_2} \doteq (p_{n,m_1}, p_{n,m_1} + p_{n,m_2}]$, it attempts on link $(n, m_2)$; and so on. If $y_n(t) \in (p_n, 1]$, node $n$ does not attempt a transmission. We denote

$$Y^{(r)}(t, \xi) = \sum_{s=0}^{t-1} I\{y_n(s) \le \xi_n, \ n \in \mathcal{N}\}, \ t = 0, 1, 2, \ldots, \ \xi = (\xi_1, \ldots, \xi_N) \in [0, 1]^N, \quad (22)$$

where $I\{\cdot\}$ is an event indicator (not to be confused with $I$ as the number of links).

We will use vector notations $F^{(r)}(t) = (F_i^{(r)}(t), \ i \in \mathcal{I})$, $\hat{F}^{(r)}(t) = (\hat{F}_i^{(r)}(t), \ i \in \mathcal{I})$. Finally, we extend the time domain of functions $F^{(r)}(t)$, $\hat{F}^{(r)}(t)$, and $Y^{(r)}(t, \xi)$ to all real $t \ge 0$ by adopting the convention (as we already did for $Q^{(r)}(t)$) that they are constant within each time slot $[\ell, \ell + 1)$ for all integer $\ell \ge 0$.

By our definitions and conventions, $F^{(r)}(0) = 0$, $\hat{F}^{(r)}(0) = 0$, and $Y^{(r)}(0, \xi) = 0$ for every $\xi \in [0, 1]^N$. We also have the following "integral form" of the recurrence (13):

$$Q_i^{(r)}(t) = Q_i^{(r)}(0) + F_i^{(r)}(t) - \hat{F}_i^{(r)}(t) - \left[ \min_{s \in [0,t]} \{Q_i^{(r)}(0) + F_i^{(r)}(s) - \hat{F}_i^{(r)}(s)\} \right]^-, \ t \ge 0, \ i \in \mathcal{I}. \quad (23)$$

Without loss of generality we can assume that processes $F^{(r)}(\cdot)$ and $Y^{(r)}(\cdot, \cdot)$, although carry index $(r)$, do not depend on $r$. (For every $r$ they are constructed from the same underlying sequences i.i.d. random variables.) Then, along any fixed sequence of $r \uparrow \infty$, the following functional strong law of large numbers (FSLLN) properties hold: with probability 1

$$(1/r)F_i^{(r)}(rt) \to \lambda_i t, \ \text{u.o.c.}, \ \forall i \in \mathcal{I}, \quad (24)$$

$$(1/r)Y^{(r)}(rt, \xi) \to [\prod_{n \in \mathcal{N}} \xi_n]t, \ \text{u.o.c.} \quad (25)$$

(In (24) and (25), u.o.c. means that the convergence is uniform on $t$ and $(t, \xi)$, respectively, within any bounded subset.)

## 8.2   Fluid Sample Paths: Definition

Consider a sequence of $r \uparrow \infty$. For each $r$, let $(Q^{(r)}(\cdot), F^{(r)}(\cdot), \hat{F}^{(r)}(\cdot), Y^{(r)}(\cdot, \cdot))$ be a realization (that is, a fixed sample path) of the corresponding random process, with some fixed initial condition $Q^{(r)}(0)$, $\|Q^{(r)}(0)\| = r$. The entire realization is uniquely determined by $Q^{(r)}(0)$, $F^{(r)}(\cdot)$, and $Y^{(r)}(\cdot, \cdot))$. Assume this sequence of *realizations* satisfies conditions (24) and (25).

Consider the following rescaled trajectory for each $r$:

$$(q^{(r)} = (q^{(r)}(t), \ t \geq 0), \ f^{(r)} = (f^{(r)}(t), \ t \geq 0), \ \hat{f}^{(r)} = (\hat{f}^{(r)}(t), \ t \geq 0)),$$

where $f^{(r)}(t) = (1/r)F^{(r)}(rt)$, $\hat{f}^{(r)}(t) = (1/r)\hat{F}^{(r)}(rt)$, and (recall) $q^{(r)}(t) = (1/r)Q^{(r)}(rt)$.

*A triple of vector-functions $(q = (q(t), \ t \geq 0), \ f = (f(t), \ t \geq 0), \ \hat{f} = (\hat{f}(t), \ t \geq 0))$ is called a* fluid sample path *(FSP), if the uniform on compact sets (u.o.c.) convergence*

$$(q^{(r)}, f^{(r)}, \hat{f}^{(r)}) \rightarrow (q, f, \hat{f}) \tag{26}$$

*holds for at least one sequence $(q^{(r)}, f^{(r)}, \hat{f}^{(r)})$ of scaled trajectories (with $r \uparrow \infty$), such that (24) and (25) hold.*

We note that the definition of the FSPs given here, as well as their dynamic properties derived in next Section 8.3, do *not* require condition (14).

## 8.3   Fluid Sample Paths: Basic Dynamics

**Lemma 1** *Any fluid sample path satisfies the following conditions: (19),*

$$\text{functions } q(\cdot), \ f(\cdot), \ \hat{f}(\cdot), \ \text{are Lipschitz continuous, with } \ f_i(t) = \lambda_i t, \tag{27}$$

$$q_i'(t) = [\lambda_i - v_i(t)]_{q_i(t)}^+, \ \forall i \in \mathcal{I}, \ \text{for almost all } t \geq 0, \tag{28}$$

$$v(t) \in \arg\max_{x \in M}[\gamma \times q(t)^\beta] \cdot \log x. \tag{29}$$

**Proof.** Consider a fixed FSP, and any fixed sequence of rescaled paths $(q^{(r)}, f^{(r)}, \hat{f}^{(r)})$ which "defines" this FSP (via convergence (26) and the other conditions). Properties (19) and (27) are obvious, given the FSP construction. (Property $f_i(t) = \lambda_i t$ follows from (24).) Therefore almost all points $t \geq 0$ are such that proper derivatives of component functions of the FSP exist - such time points will be called *regular*. It will suffice to show (28)-(29) for a given regular point $t > 0$. Indeed, switching to rescaled paths in (23) and taking the limit on $r$, we obtain

$$q_i(t) = q_i(0) + f_i(t) - \hat{f}_i(t) - \left[\min_{s \in [0,t]}\{q_i(0) + f_i(s) - \hat{f}_i(s)\}\right]^-, \ t \geq 0, \ i \in \mathcal{I}. \tag{30}$$

Denote $v_i(t) = \hat{f}_i'(t)$ and recall that $f_i'(t) = \lambda_i$. Given (30), we see that the equation in (28) holds trivially if $q_i(t) > 0$. If $q_i(t) = 0$ we must have $q_i'(t) = 0$ (by regularity), and then $\lambda_i - v_i(t) \leq 0$ because otherwise (30) would imply that the right derivative $(d^+/dt)q_i(t)$ exists and is equal to $\lambda_i - v_i(t) > 0$. This proves (28).

13

Let us prove (29). Now, consider the FSP, and the corresponding rescaled trajectories in a small interval $[t, t + \epsilon]$. By continuity of $q(\cdot)$ and by(26), all values of $q^{(r)}(s) \equiv Q^{(r)}(rs)$, $s \in [t, t + \epsilon]$, for all sufficiently large $r$, are close to $q(t)$, as long as $\epsilon > 0$ is small. For each $r$, consider the corresponding unscaled trajectory of $Q^{(r)}(\cdot)$ in the (corresponding) time interval $[rt, rt + r\epsilon]$. It will be convenient to assume that the trajectory with index $r$ uses rescaled dynamic weights $w_i^{(r)}(s) = w_i(s)/r = \alpha_i/r + \gamma_i[Q_i^{(r)}(s)]^\beta/r$. (This does not change anything, by Proposition 2(i).) We see that, in $[rt, rt + r\epsilon]$, $w_i^{(r)}(s)$ is close to $\gamma_i q_i(t)^\beta$. Then, it is easy to see from Proposition 2(iii) that, for any link $i = (n, m)$ with $q_i(t) > 0$, the corresponding interval $\phi_{n,m}$ (see the construction governing transmission attempts) is close to such interval corresponding to the access probabilities $p(\gamma \times q(t))$. (Again, this is for the unscaled trajectory in the interval $[rt, rt + r\epsilon]$, with small $\epsilon$ and for large values of $r$.) Given the definition (22), condition (25), and again Proposition 2(iii), this implies that the time-average rate of successful attempts in the interval $[rt, rt + r\epsilon]$ is close to $\mu_i(p(\gamma \times q(t)^\beta)$. This shows (we omit $\epsilon$-$\delta$ formalities) that, for all links $i$ with $q_i(t) > 0$,

$$v_i(t) = \mu_i(p(\gamma \times q(t)^\beta)),$$

which, by Proposition 2(ii), implies (29). ∎

## 8.4 Stability of Fluid Sample Paths. Conclusion of Proof of Theorem 2

All statements of Theorem 2, except (21), easily follow from the fact that (24) and (25) hold w.p.1, the FSP construction, and FSP property (27). It remains to show that FSPs satisfy (21), which is done in the following

**Theorem 3** *Consider function*

$$\Psi(y) = \frac{1}{\beta + 1} \sum_i \frac{\gamma_i}{\mu_i^*} y_i^{\beta+1}, \ y \in R_+^I. \tag{31}$$

*(It depends on $\mu^*$ as a parameter.) Suppose $\lambda$ satisfies (14), and consider the family of Lipschitz continuous trajectories $(q(t), \ t \geq 0)$, satisfying (28) and (29). Then, for any $\epsilon_1 > 0$ there exists $\epsilon_2 > 0$, such that*

$$\Psi(q(t)) \geq \epsilon_1 \ implies \ \frac{d}{dt}\Psi(q(t)) \leq -\epsilon_2. \tag{32}$$

*As a corollary, there exist $\epsilon > 0$ and $T > 0$ such that (21) holds uniformly for all FSPs.*

**Remark.** Lipschitz continuous trajectories satisfying (28) and (29) also arise in a completely different setting, namely, in the "session level" stability analysis of communication networks with "concave-utility-based" allocation of service rates to different sessions; see [9, 4, 25]. (We want to emphasize that the *derivation* of properties (28) and (29) in our setting is

completely different.) In all the previous work cited above the region $M$ (which has a different meaning from ours) is convex. In our model, region $M$ is not necessarily convex, and in fact non-convex in most cases. However, *for the purposes of establishing trajectory stability* in Theorem 3, convexity of $M$ is *not* important. Since this last point may not be immediately clear from the previous work (some of which does use convexity of $M$, even though it does not have to), and for completeness, we present the proof of Theorem 3.

**Proof.** Denote $\log \mu^* = u^*$ and $\log v(t) = u(t)$. Then, we have

$$
\begin{aligned}
\frac{d}{dt}\Psi(q(t)) &= \sum_i \frac{\gamma_i}{\mu_i^*} q_i(t)^\beta q_i'(t) \\
&= \sum_i \frac{\gamma_i}{\mu_i^*} q_i(t)^\beta [\lambda_i - e^{u_i(t)}] \\
&\leq \sum_i \frac{\gamma_i}{\mu_i^*} q_i(t)^\beta [\lambda_i - e^{u_i^*} - e^{u_i^*}(u_i(t) - u_i^*)] \qquad (33)\\
&= \sum_i \gamma_i q_i(t)^\beta \left[\frac{\lambda_i}{\mu_i^*} - 1\right] - \left[\sum_i \gamma_i q_i(t)^\beta u_i(t) - \sum_i \gamma_i q_i(t)^\beta u_i^*\right] \qquad (34)\\
&\leq \sum_i \gamma_i q_i(t)^\beta \left[\frac{\lambda_i}{\mu_i^*} - 1\right]. \qquad (35)
\end{aligned}
$$

The inequality (33) uses convexity of the exponent function, and inequality (34) $\leq$ (35) is because $\sum_i \gamma_i q_i(t)^\beta u_i(t) - \sum_i \gamma_i q_i(t)^\beta u_i^* \geq 0$ by condition (29). The RHS of (35) is negative and bounded away from 0 as long as $\Psi(q(t))$ is positive and bounded away from 0. ∎

# 9  Proof of Theorem 2 for QRA-II

## 9.1  Preliminaries

Recall that we are now in the conditions of Theorem 2(ii), and therefore under additional assumption (15), which implies large deviations (Cramer's) bound for the input processes. For any $i \in \mathcal{I}$ and any $\nu > 0$, there exists a constant $a = a(\nu) > 0$ such that, for all sufficiently large $n$, uniformly on $k \geq 1$,

$$
\Pr\{|\frac{1}{n}\sum_{t=k}^{k+n-1} A_i(t) - \lambda_i| \geq \nu\} < e^{-an} . \qquad (36)
$$

Let arbitrary $\nu > 0$ and $L > 0$ be fixed. Let us also pick any $\zeta > 0$ such that $\zeta < \eta \equiv 1 - \kappa$. For each $n$, let us cover the interval $[0, rL]$ with $P_L^r \doteq \lfloor rL/r^\zeta \rfloor + 1$ equal non-overlapping $r^\zeta$-long intervals $[(j-1)r^\zeta, jr^\zeta)$, $1 \leq j \leq P_L^r$. Define for each $i \in \mathcal{I}$, and each $\xi \in [0, 1]$,

$F_{i,j}^{(r)} \doteq F_i^{(r)}(jr^\zeta) - F_i^{(r)}((j-1)r^\zeta)$, the number of arrivals of flow $i$ in the time interval $[(j-1)r^\zeta, jr^\zeta)$,

15

$$Y_{i,j,\xi}^{(r)} \doteq Y_i^{(r)}(jr^\zeta, \xi) - Y_i^{(r)}((j-1)r^\zeta, \xi).$$

Let us denote

$$
\begin{aligned}
E_i^r(L, \nu) &= \bigcup_{1 \le j \le P_L^r} \left\{ \left| \frac{F_{i,j}^{(r)}}{r^\zeta} - \lambda_i \right| > \nu \right\}, \\
G_i^r(L, \nu, \xi) &= \bigcup_{1 \le j \le P_L^r} \left\{ \left| \frac{Y_{i,j,\xi}^{(r)}}{r^\zeta} - \xi \right| > \nu \right\}.
\end{aligned}
$$

**Lemma 2** *The following properties hold (with $\mathcal{Q}_+$ being the set of strictly positive rational numbers):*

$$\Pr\left( \bigcup_{\nu, L \in \mathcal{Q}_+} \bigcap_{k=1}^\infty \bigcup_{r=k}^\infty E_i^r(L, \nu) \right) = 0, \quad \forall i \in \mathcal{I}, \tag{37}$$

$$\Pr\left( \bigcup_{\nu, L \in \mathcal{Q}_+} \bigcup_{\xi \in [0,1], \; \xi \in \mathcal{Q}_+} \bigcap_{k=1}^\infty \bigcup_{r=k}^\infty G_i^r(L, \nu, \xi) \right) = 0, \quad \forall i \in \mathcal{I}. \tag{38}$$

*Equivalently, with probability 1, for any rational numbers $L > 0$, $\nu > 0$ and $\xi \in [0, 1]$, there exists finite $k$ such that for all $r > k$,*

$$\max_{i \in \mathcal{I}, \; 1 \le j \le P_L^r} \left| \frac{E_i^r(L, \nu)}{r^\zeta} - \lambda_i \right| \le \nu, \tag{39}$$

$$\max_{i \in \mathcal{I}, \; 1 \le j \le P_L^r} \left| \frac{G_i^r(L, \nu, \xi)}{r^\zeta} - \xi \right| \le \nu. \tag{40}$$

**Proof** uses Cramer's bound (36) and Borel-Cantelli lemma. (See [19], Lemma 4.3, for the proof of (37), and (38) is proved analogously.) ∎

## 9.2 FSP definition and stability. Proof of Theorem 2

Consider the system under the QRA-II algorithm. For this algorithm we define fluid sample paths (FSP) in exactly same way as for QRA-I (in Section 8.2), except we require that, in addition to (or, rather, instead of) (24) and (25), a defining sequence $(q^{(r)}, f^{(r)}, \hat{f}^{(r)})$ of scaled sample paths satisfies stronger conditions (39) and (40). (Such FSPs clearly satisfy the initial condition (19) and Lipschitz condition (27), but certainly *not* the conditions (28) and (29) - the FSP dynamics under QRA-II is very different.) With this FSP definition, the proof of all statements of Theorem 2, except (21), is almost automatic. To establish (21), we prove the following key property of the FSPs under QRA-II algorithm.

**Theorem 4** *Suppose a point $\mu^* \in M \cap R_{++}^I$ is such that for some real $c$, $\gamma_i(\mu_i^* - \lambda) = c$ for all $i$. [We do not necessarily assume (14), and thus a positive $c$ may or may not exist.] Then, in addition to (19) and (20), every FSP satisfies the following condition at every regular point $t \geq 0$ where $\max_i \gamma_i q_i(t) > 0$:*

$$\frac{d}{dt} \max_i \gamma_i q_i(t) \leq -c. \tag{41}$$

The above definition of the FSPs for QRA-II, as well as that for QRA-I, does *not* require condition (14). And Theorem 4 holds regardless of (14). If condition (14) does hold, then $\mu^*$ can be chosen to be the point where the ray in the direction $(1/\gamma_1, \ldots, 1/\gamma_I)$, starting at $\lambda$, hits the boundary of $M$, in which case $c > 0$, and $\frac{d}{dt} \max_i \gamma_i q_i(t) \leq -c < 0$, which of course proves (21).

**Proof of Theorem 4** Properties (19) and (20) are proved the same way as for QRA-I algorithm (and for fluid limits in general). To prove (41), consider a fixed FSP and a sequence of rescaled pre-limit trajectories, defining it. We will use notation: $z(t) \doteq \max_i \gamma_i q_i(t)$, $z^{(r)}(t) \doteq \max_i \gamma_i q_i^{(r)}(t)$, $Z^{(r)}(rt) \doteq \max_i \gamma_i Q_i^{(r)}(rt)$. Suppose (41) does not hold. Then, there exists a regular point $t$ such that $z(t) > 0$ and $z'(t) > -c_1 > -c$. We will show that this leads to a contradiction. We can choose constants $\delta > 0$, $\delta_1 > 0$, and $c_2 \in (c_1, c)$, such that

$$z(s) > \delta_1 , \quad \forall s \in [t, t + \delta] ,$$

and

$$\frac{z(t + \delta) - z(t)}{\delta} > -c_2 .$$

For each $r$, let us now divide the interval $[t, t + \delta]$ into $r^\kappa \delta / \ell$ intervals, each of length $\frac{\ell r^\eta}{r}$, where $\eta = 1 - \kappa$, and $\ell > 0$ is an arbitrary fixed constant. (Since $r^\kappa \delta / \ell$ may not be an integer, we should divide into, say, $\lceil r^\kappa \delta / \ell \rceil$ intervals. To avoid trivial complications and heavy notation, we assume that $r^\kappa \delta / \ell$ *is* integer. It will be clear that we do not lose the correctness of the argument.) Note that in the "unscaled time", each subinterval is of length $\ell r^\eta$.

From the Dirichlet principle, for all sufficiently large $r$, in at least one of the subintervals (of length $\frac{\ell r^\eta}{r}$), the average rate of change of $z^{(r)}(.)$ is greater than or equal to $(-c_2)$. We pick such a subinterval $[s^{(r)}, s^{(r)} + \frac{\ell r^\eta}{r}]$ for each $r$. Let us choose a further subsequence of the sequence of indices $\{r\}$ (which we will still denote by $\{r\}$), such that $s^{(r)} \to s$, for some fixed $s \in [t, t + \delta]$. Obviously, the right end-point $s^{(r)} + \frac{\ell r^\eta}{r}$ of the subinterval also converges to $s$.

From the subsequence $\{r\}$, we choose a further subsequence such that the order of values of $\gamma_i q_i^{(r)}(s^{(r)}), i \in \mathcal{I}$, remains same. For example, without loss of generality, we can assume that

$$z^{(r)}(s^{(r)}) = \gamma_1 q_1^{(r)}(s^{(r)}) \geq \ldots \geq \gamma_I q_I^{(r)}(s^{(r)}) .$$

Finally, for each $i \in \mathcal{I}$, consider the following trajectories:

$$y_i^{(r)}(\tau) \doteq r^\kappa [\gamma_i q_i^{(r)}(s^{(r)} + \tau/r^\kappa) - z^{(r)}(s^{(r)})], \ \tau \in [0, \ell] ,$$

17

and choose a subsequence such that for each $i$,

$$y_i^{(r)}(0) \to y_i(0) \ ,$$

where $\max_i y_i(0) = y_1(0) = 0$ (by our construction), and each other $y_i(0)$ is either finite non-positive or $-\infty$. Let us consider only the case when all $y_i(0)$ are finite. (If not, it is easy to observe that, in the (unscaled) time interval $[rs^{(r)}, rs^{(r)} + r^\eta \ell]$, the queues with $y_i(0) = -\infty$ have asymptotically vanishing impact on the service of queues with $y_i(0) > -\infty$. So, essentially same argument, restricted to the subset of queues with finite $y_i(0)$ applies.)

We notice that the trajectory $y_i^{(r)}(\cdot)$ is obtained from the trajectory $Q_i^{(r)}(\cdot)$ by the time "speedup" of $r^\eta$ and the "space" scaling by the factor $1/r^\eta$ (in addition to the time shift).

In the next step, we can and do choose a further subsequence of $\{r\}$, such that the sequence $\{y_i^{(r)}(\cdot), \ i \in \mathcal{I}\}$, converges u.o.c. in $[0, \ell]$ to a Lipschitz continuous trajectory $\{y_i(\cdot), \ i \in \mathcal{I}\}$, which we call a *local fluid sample path* (LFSP). In this step we use properties (39) and (40) which guarantee that, roughly speaking, the functional law of large numberes holds not only over any (unscaled) interval of length $\delta r$, but also uniformly over the set of $\ell r^\eta$-long subintervals of it.

It is also not hard to see that the LFSP trajectory satisfies the following conditions (42) and (43) at every regular point $\tau \in [0, \ell]$

$$y_i'(\tau) = \gamma_i[\lambda_i - \mu_i(\tau)], \ \forall i \in \mathcal{I}, \tag{42}$$

$$\mu_i(\tau) \in \arg\max_{x \in M}[\alpha \times e^{a_2 y(\tau)}] \cdot \log x, \tag{43}$$

where $a_2 = \kappa/z(s)^\eta$. Indeed, note that $z^{(r)}(s^{(r)})/z(s) \to 1$ and $\gamma_i q_i^{(r)}(s^{(r)} + \tau/r^\kappa)/z(s) \to 1$ as $r \to \infty$, uniformly in $i$ and in $\tau \in [0, \ell]$; and recall that $\gamma_i q_i^{(r)}(s^{(r)} + \tau/r^\kappa) - z^{(r)}(s^{(r)}) = r^{-\kappa} y_i^{(r)}(\tau)$. We have

$$[\gamma_i Q_i^{(r)}(s^{(r)} r + r^\eta \tau)]^\kappa - Z^{(r)}(s^{(r)} r)^\kappa = r^\kappa \left[ [\gamma_i q_i^{(r)}(s^{(r)} + \tau/r^\kappa)]^\kappa - z^{(r)}(s^{(r)})^\kappa \right] =$$

$$= r^\kappa [\kappa \hat{z}^{\kappa-1}][r^{-\kappa} y_i^{(r)}(\tau)],$$

where $\hat{z}$ is a number "between" $z^{(r)}(s^{(r)})$ and $\gamma_i q_i^{(r)}(s^{(r)} + \tau/r^\kappa)$. (Here we simply applied the mean value theorem for an increment of function $x^\kappa$.) But, $\hat{z} \to z(s)$ as $r \to \infty$, uniformly in $i$ and $\tau$, and therefore we have uniform convergence

$$[\gamma_i Q_i^{(r)}(s^{(r)} r + r^\eta \tau)]^\kappa - Z^{(r)}(s^{(r)} r)^\kappa \to \frac{\kappa}{z(s)^\eta} y_i(\tau).$$

It remains to consider the behavior of the sample path in the unscaled time interval $[s^{(r)} r + r^\eta \tau, s^{(r)} r + r^\eta \tau + r^\eta \Delta\tau]$, with small fixed $\Delta\tau$, as $r$ becomes large, and use the form of the QRA-II scheduling rule. (The argument here is analogous to the one used to prove (29) for QRA-I rule. We omit details.)

18

By definition of LFSP, $\max_i y_i(0) = 0$, and by our construction,

$$\max_i y_i(\ell) \geq -c_2\ell > -c\ell. \tag{44}$$

We will use the following Lyapunov function:

$$\Psi(\tau) = \sum_i \frac{1}{\mu_i^* \gamma_i} \alpha_i e^{a_2[y_i(\tau)+c\tau]},$$

where, recall, $\mu^* \in M$ is such that $\gamma_i(\mu_i^* - \lambda_i) = c$ for all $i$. We have

$$\Psi(0) \leq \sum_i \frac{\alpha_i}{\mu_i^* \gamma_i}.$$

We denote $u_i(\tau) = \log \mu_i(\tau)$, $u_i^* = \log \mu_i^*$. Then,

$$
\begin{aligned}
\frac{d}{dt}\Psi(y(\tau)) &= a_2 \sum_i \frac{\alpha_i}{\mu_i^* \gamma_i} e^{a_2[y_i(\tau)+c\tau]}[y_i'(\tau) + c] \\
&= a_2 \sum_i \frac{\alpha_i}{\mu_i^*} e^{a_2[y_i(\tau)+c\tau]}[\mu_i^* - \mu_i(\tau)] \\
&\leq a_2 \sum_i \frac{\alpha_i}{\mu_i^*} e^{a_2[y_i(\tau)+c\tau]}[\mu_i^* - (\mu_i^* + \mu_i^*(u_i(\tau) - u_i^*))] \tag{45} \\
&= a_2 e^{a_2 c\tau}\left[\sum_i \alpha_i e^{a_2 y_i(\tau)} u_i^* - \sum_i \alpha_i e^{a_2 y_i(\tau)} u_i(\tau)\right] \leq 0. \tag{46}
\end{aligned}
$$

Since $\Psi(\tau)$ is non-increasing and $\Psi(0)$ is bounded, we see that for some constant $K$, depending only on the system parameters and on $c$,

$$\max_i y_i(\tau) \leq Kz(s)^\eta - c\tau, \quad \tau \in [0,\ell]. \tag{47}$$

Recall that function $z(\cdot)$ is defined on the time scale of FSP (not LFSP), and it is Lipschitz and $z(0) \leq \max \gamma_i$. Then, $Kz(s)^\eta$ is uniformly bounded within any finite interval (of FSP time scale), and in particular for $s \in [t, t+\delta]$, with $t$ and $\delta$ chosen at the beginning of this proof. Therefore, since $\ell$ could be chosen arbitrarily large, (47) contradicts (44). ∎

# References

[1] N. Abramson. The ALOHA system – Another alternative for computer communications. *Proc. AFIPS Conf.*, Fall Joint Computer Conference, Vol. 37, (1970), pp. 281-285.

[2] M. Andrews, K. Kumaran, K. Ramanan, A. L. Stolyar, R. Vijayakumar, P. Whiting. Scheduling in a Queueing System with Asynchronously Varying Service Rates. *Probability in Engineering and Informational Sciences*, 2004, Vol. 18, pp. 191-217.

[3] V. Anantharam. The stability region of the finite-user slotted Aloha protocol. *IEEE Trans. Inform. Theory,* Vol. 37, (1991), No. 3, pp. 535–540.

[4] T. Bonald and L. Massoulie. Impact of Fairness on Internet Performance. *Proceedings of ACM SIGMETRICS*, 2001.

[5] H. Chen. Fluid Approximations and Stability of Multiclass Queueing Networks: Work-conserving Disciplines. *Annals of Applied Probability*, Vol. 5, (1995), pp. 637-665.

[6] J. G. Dai. On the Positive Harris Recurrence for Open Multiclass Queueing Networks: A Unified Approach Via Fluid Limit Models. *Annals of Applied Probability*, Vol. 5, (1995), pp. 49-77.

[7] J.G.Dai and W.Lin. Maximum Pressure Policies in Stochastic Processing Networks. *Operations Research*, Vol. 53, (2005), pp. 197-218.

[8] J. G. Dai and S. P. Meyn. Stability and Convergence of Moments for Open Multiclass Queueing Networks Via Fluid Limit Models. *IEEE Transactions on Automatic Control*, Vol. 40, (1995), pp. 1889-1904.

[9] G. de Veciana, T. J. Lee, and T. Konstantopoulos. Stability and Performance Analysis of Networks Supporting Elastic Services. *IEEE/ACM Transactions on Networking*, Vol.9, 2001, pp. 2-14.

[10] A. Eryilmaz, R. Srikant, J. Perkins. Stable Scheduling Policies for Fading Wireless Channels. *IEEE/ACM Transactions on Networking*, vol. 13, 2005, pp. 411-424.

[11] S. N. Ethier and T. G. Kurtz. *Markov Processes: Characterization and Convergence.* John Wiley and Sons, New York, 1986.

[12] P. Gupta, Y. Sankarasubramaniam, A. L. Stolyar. Random-Access Scheduling with Service Differentiation in Wireless Networks. *Proceeding of INFOCOM'2005*, Miami, March 13-17, 2005.

[13] P. Gupta, A. L. Stolyar. Throughput Region of Random Access Networks of General Topology. 2005, submitted.

[14] P. Gupta, A. L. Stolyar. Optimal Throughput Allocation in General Random-Access Networks. *Proceeding of CISS'2006*, Princeton, March 2006.

[15] K. Kar, S. Sarkar, L. Tassiulas. Achieving Proportionally Fair Rates using Local Information in Aloha Networks. *IEEE Trans. Autom. Control,* Vol. 49, (2004), No. 10, pp. 1858–1862.

[16] V.A. Malyshev and M.V. Menshikov. Ergodicity, Continuity, and Analyticity of Countable Markov Chains. *Transactions of Moscow Mathematical Society*, Vol. 39, (1979), pp. 3-48.

[17] J. Massey and P. Mathys. The collision channel without feedback. *IEEE Trans. Inform. Theory,* Vol. IT-31, (1985), no. 2, pp. 192–204.

[18] A.N. Rybko and A.L. Stolyar. Ergodicity of stochastic processes describing the operation of open queueing networks. *Problems of Information Transmission*, vol. 28, 1992, pp. 199-220. (Translated from *Problemy Peredachi Informatsii*, vol. 28, no. 3, pp. 3-26, 1992.)

[19] S. Shakkottai and A. L. Stolyar. Scheduling for Multiple Flows Sharing a Time-Varying Channel: The Exponential Rule. *Analytic Methods in Applied Probability. In Memory of Fridrih Karpelevich. Yu. M. Suhov, Editor.* American Mathematical Society Translations, Series 2, Volume 207, pp. 185-202. American Mathematical Society, Providence, RI, 2002.

[20] S. Shakkottai, R. Srikant, and A. L. Stolyar. Pathwise Optimality of the Exponential Scheduling Rule for Wireless Channels. *Advances in Applied Probability*, 2004, Vol. 36, No. 4, pp. 1021-1045.

[21] A.L. Stolyar. On the Stability of Multiclass Queueing Networks: A Relaxed Sufficient Condition via Limiting Fluid Processes. *Markov Processes and Related Fields*, 1(4), 1995, pp. 491-512.

[22] A.L. Stolyar. MaxWeight Scheduling in a Generalized Switch: State Space Collapse and Workload Minimization in Heavy Traffic. *Annals of Applied Probability*, 2004, Vol.14, No.1, pp.1-53.

[23] L.Tassiulas, A.Ephremides. Stability Properties of Constrained Queueing Systems and Scheduling Policies for Maximum Throughput in Multihop Radio Networks. *IEEE Transactions on Automatic Control*, Vol. 37, 1992, pp. 1936-1948.

[24] X. Wang and K. Kar. Distributed Algorithms for Max-min Fair Rate Allocation in Aloha Networks. *Proceedings of the 42nd Annual Allerton Conference*, Urbana-Champaign, 2004.

[25] H. Q. Ye. Stability of Data Networks Under An Optimization-Based Bandwidth Allocation. *IEEE Transactions on Automatic Control*, Vol. 48, 2003, No. 7, pp. 1238-1242.