

# Can Multipath Mitigate Power Law Delays? – Effects of Parallelism on Tail Performance

Jian Tan<sup>1</sup>, Wei Wei<sup>2</sup>, Bo Jiang<sup>2</sup>, Ness Shroff<sup>1</sup>, Don Towsley<sup>2</sup>

<sup>1</sup> Departments of Electrical and Computer Engineering & Computer Science and Engineering  
The Ohio State University, Columbus, OH 43210

<sup>2</sup>Department of Computer Science, University of Massachusetts, Amherst, MA 01003

**Abstract**—Parallelism has often been used to improve the reliability and efficiency of a variety of different engineering systems. In this paper, we quantify the efficiency of parallelism in systems that are prone to failures and exhibit power law processing durations. We focus on the context of transmitting a data unit in communication networks, where parallelism can be achieved by multipath transmission (e.g., multipath routing). We investigate two types of transmission schemes: redundant and split transmission techniques. We find that the power-law transmission delay phenomenon still persists with multipath transmission. In particular, we show that when the transmission delays of each path are characterized by the same power law, redundant multipath transmission can only result in a constant factor performance gain, while order gains are possible when the delays are light tailed. We further compare the performance of redundant transmission and split transmission, and show that there is no clear winner. Depending on the packet size distribution properties and the manner in which splitting is performed, one scheme results in greater performance over the other. Specifically, split transmission is effective in mitigating power law delays if the absolute value of the logarithm of the packet size probability tail is regularly varying with positive index, and becomes ineffective if the above quantity is slowly varying. Based on our analysis, we develop an optimal split transmission strategy, and show that this strategy always outperforms redundant transmission.

## I. INTRODUCTION

Parallelism is a common approach to improve reliability and efficiency in practice. For instance, in peer to peer systems, a file is downloaded in parallel from multiple peers; in grid computing, a job is allocated to multiple machines to be computed simultaneously; and in computer communication networks, multipath routing can be used to improve the efficiency and reliability of data transfer. In one type of parallelism, a file/job is fetched/computed in its entirety, and hence the completion time is the minimum of the completion times from/at the multiple locations. In another type of parallelism, a file/job is split into multiple pieces, fetched/computed independently, and hence the completion time is the maximum of the completion times of all the pieces. In practice, more complicated strategies can be developed by appropriately combining these two types of parallelism. In both cases, we expect better efficiency from using parallelism since the delay is either the minimum one or because a smaller job needs to be completed.

In this paper, we quantify the efficiency of parallelism in mitigating power law tails, which have been shown to be present when a job needs to be retransmitted after a failure occurs. For example, in wireless communication networks, re-

cent studies [7]–[10] show that, contrary to traditional wisdom, when the probability of packet errors is a function of the packet length, retransmission-based protocols may cause power law transmission durations and possibly even zero throughput. Similar results have been reported in other contexts [2], [13]. A natural question to ask is whether and, if so, how, using parallelism can mitigate power law delays, which is the focus of our study.

To focus our discussion, let us consider the notion of parallelism in the context of communication networks, where a data unit can be transmitted using multiple paths (also known as multipath routing or more generally multipath transmission). A data unit can be a file or packet (which are used interchangeably, henceforward), and the transmission needs to restart after a failure (i.e., there is no check point in the transmission). We consider two multipath transmission strategies, *redundant* and *split transmission*, that correspond respectively to the two aforementioned types of parallelism. More specifically, redundant transmission replicates a packet and sends each copy over a different path (we use the terms channel and path interchangeably for the rest of the paper) and therefore, the transmission is successful once the first of the packets arrives at the destination; split transmission, on the other hand, breaks the data unit into several pieces and dispatches each piece along a different path, which completes the transmission when all the pieces arrive at the destination successfully.

We aim to answer the following three questions: (I) Can redundant or split transmission eliminate power laws in transmission delays, and how can the performance gain from multipath transmission be characterized? (II) Is split transmission or redundant transmission more beneficial in mitigating power law delays? and (III) What is the optimal strategy to split packets and dispatch those fragmented pieces to the appropriate paths.

To address the above questions, we generalize the single *channel* model introduced in [9] to a multipath channel model. First note that a channel can be viewed as a medium over which faults can occur causing jobs to be interrupted and retransmitted. In the context of communication networks, this corresponds to a wireless communication channel as in [9], in the context of grid computing the channel may correspond to the processor over which the computations are completed, etc. Henceforth, we will focus on communication networks and

consider the notion of a channel in that context. Specifically, consider a communication network where there are  $K$  paths between a source and destination. The channel dynamics of path  $j$ ,  $1 \leq j \leq K$ , are modeled as an on-off process  $\{(A_i^j, U_i^j)\}_{i \geq 1}$  that alternates between available period,  $A_i^j$ , and unavailable period,  $U_i^j$ . Only in each time period  $A_i^j$  when the channel becomes available, can a packet start its transmission over the path. If the length of  $A_i^j$  is longer than the length of the packet, the transmission is considered successful over path  $j$ ; otherwise, we wait until the beginning of the next available period  $A_{i+1}^j$  and retransmit the packet from the beginning. The above model can be viewed as a first order approximation to channels that may fail. Channel failures can happen due to many reasons. For instance, in a wireless network environment, failures occur due to channel fading, interference and contention with other nodes, multipath effects, obstructions, and node mobility [12]. As a consequence, the signal to noise ratio (SINR) may vary in different time scales. The on periods  $\{A_i^j\}$  in our model correspond to the situation when SINR is high, while the off periods  $\{U_i^j\}$  correspond to the situation when SINR is low.

Our main contributions in this paper can be summarized as follows:

- We show that, when all packets are of the same size, redundant transmission can greatly reduce the transmission delay in the sense that the ratio of the delay distribution tail with and without redundant transmission tends to zero (see Proposition III.1). However, in reality, packet sizes are usually variable due to many other considerations, e.g., reducing communication costs and extra overhead induced from encapsulation. We prove that, when packet sizes are random variables that satisfy  $\log \mathbb{P}[L > x] \approx \alpha^* \log \mathbb{P}[A^j > x]$ , redundant transmission does not change the order of the probability tail of the transmission delays (see Theorem 2), and can only improve the system performance by a constant factor (see Theorem 3).
- We show that split transmission is effective in mitigating power delays if the absolute value of the logarithm of the packet size probability tail is regularly varying with positive index, and becomes ineffective if the above quantity is slowly varying (see Theorems 4 and 5). To illustrate the point, we calculate the effectiveness of split transmission for different packet size distributions. Furthermore, we provide a solution for optimal split when we have heterogeneous paths, and show that this optimal strategy always outperforms redundant transmission (see Theorem 6). To refine the result, we also derive an exact asymptotic for packet delivery time under optimal split transmission (see Theorem 7).

In terms of related work, it was observed in [13] that power law processing times can arise in a system where jobs need to restart once a failure occurs. This observation was rigorously addressed in [2], [8]–[10] for a single channel model. The result reveals that, when the probability of packet

errors is a function of the packet length, retransmission-based protocols could cause heavy-tailed (specifically, power law) transmission durations, even when the data units and channel characteristics are light-tailed. Our study generalizes the single channel model to the one with multiple paths. Multipath transmissions have also been studied in [1] using Extreme Value theory, but only when the number of paths goes to infinity. In this work, we focus on the context of multipath transmissions in computer networks with a fixed (possibly small) number of paths, where multipath transmission has long been used to improve reliability and efficiency (e.g., [4], [5], [11]). Here we want to emphasize that, the packet size distribution has been assumed to have an infinite support in this study, which contradicts the reality that all packet networks (from the Internet to wireless LANs) impose the maximum packet sizes at the different layers of the protocol stack. It can be easily proved that eventually the transmission delay distribution will be light-tailed under this condition. However, as has been shown in [14], [15], this light-tailed behavior occurs with a power law main body of the delay distribution, and this power law behavior may have dominating effects on the system performance since it spans over a time interval that increases very fast with respect to the length of the longest packet. Thus, our assumption on the infinite support of the packet size distribution allows us to study the main body of the transmission delay distribution. While, similar to [14], [15], we can extend our results to the case with packets having finite support, we feel that this would distract from the main insights gained from the paper.

Note that the specific investigation conducted in this paper has been in the context of data transmission in wireless communication networks, especially for lower-power sensor networks where using complicated coding schemes is difficult and often simple operations are preferred to recover failed data. However, the mathematical setting described in Section 2 is quite general, and the results can be extended to many other applications that involve parallelism and job failures, such as computing jobs in grid computing, file downloading in peer to peer networks, parallel experiment planning, and parallel scheduling.

The rest of the paper is organized as follows. Section 2 presents the model description and some results on single path transmission. Redundant transmission and split transmission are investigated in Sections 3 and 4, respectively. Finally, Section 5 concludes the paper.

## II. MODEL DESCRIPTION AND PRELIMINARY RESULTS

Let  $L$  be a random variable that denotes the length of a packet. Assume that there are  $K \geq 1$  paths between the source and destination, as shown in Figure 1. The channel dynamics of path  $j$ ,  $1 \leq j \leq K$  are modeled as an on-off process  $\{(A_i^j, U_i^j)\}_{i \geq 1}$  that alternates between available  $A_i^j$  and unavailable  $U_i^j$  periods, respectively.

Packet transmission can only be initiated at the start of an available period. For a packet transmission started at the beginning of  $A_i^j$ , if  $A_i^j > L$ , the transmission is considered

successful over path  $j$ ; otherwise, we wait until the beginning of the next available period  $A_{i+1}^j$  and retransmit the packet from the beginning.

We study two multipath transmission schemes, namely, redundant transmission and split transmission. Under redundant transmission, the same packet is transmitted over all  $K$  paths, and the transmission is successful as soon as one of the  $K$  duplicates arrives at the destination. Split transmission represents the strategy where a packet is split into  $K$  pieces and each piece is sent over a different path. The transmission is complete once all the  $K$  pieces arrive at the destination successfully.

**Definition II.1** *The number of (re)transmissions of a packet of length  $L_j$  over path  $j$ ,  $1 \leq j \leq K$ , is defined as*

$$N_j \triangleq \inf\{i : A_i^j > L_j\},$$

and, the corresponding transmission time over this path is defined as

$$T_j \triangleq \sum_{i=1}^{N_j-1} (A_i^j + U_i^j) + L_j.$$

- Redundant transmission ( $L_j \equiv L$ ): the transmission completes when the first packet is successfully transmitted over one of the  $K$  paths. Therefore, the total transmission time  $T_r$  for this scheme satisfies

$$T_r \triangleq \min_{1 \leq j \leq K} T_j.$$

- Split transmission ( $\sum_{j=1}^K L_j = L$ ): the transmission completes when all  $K$  pieces of the packet are successfully transmitted. Therefore, the total transmission time  $T_s$  for this scheme satisfies

$$T_s \triangleq \max_{1 \leq j \leq K} T_j,$$

and the total number of retransmissions over  $K$  paths is

$$N \triangleq \sum_{j=1}^K N_j.$$

In this paper, we assume that  $\{U^j, U_i^j\}_{j \geq 1}$  and  $\{A^j, A_i^j\}_{j \geq 1}$ ,  $1 \leq j \leq K$  are mutually independent i.i.d. sequences of random variables, which are also independent of the packet size  $L$ . A sketch of the model depicting the system is shown in Figure 1.

We use the following notation to denote the complementary cumulative distribution functions for  $A^j$ ,  $1 \leq j \leq K$  and  $L$ ,

$$\bar{G}_j(x) \triangleq \mathbb{P}[A^j > x],$$

and

$$\bar{F}(x) \triangleq \mathbb{P}[L > x],$$

with  $\bar{F}(x)$  being continuous eventually. We say  $K$  paths are *homogeneous* if  $A^j \stackrel{d}{=} A$  and  $U^j \stackrel{d}{=} U$  for  $1 \leq j \leq K$ , where “ $\stackrel{d}{=}$ ” denotes equal in distribution. Accordingly, we use  $\bar{G}(x) \triangleq \mathbb{P}[A > x]$ . In general,  $\{A^j\}_{1 \leq j \leq K}$  (and  $\{U^j\}_{1 \leq j \leq K}$ )

need not be identically distributed, which represents the case of *heterogenous* paths.

Throughout this paper, a positive measurable function  $f$  is called regularly varying (at infinity) with index  $\rho$  if

$$\lim_{x \rightarrow \infty} f(\lambda x)/f(x) = \lambda^\rho$$

for all  $\lambda > 0$ . It is called slowly varying if  $\rho = 0$  [3]. Additionally, for any two real functions  $f(t)$  and  $g(t)$ , we use  $f(t) \sim g(t)$  to denote  $\lim_{t \rightarrow \infty} f(t)/g(t) = 1$ . Similarly, we say that  $f(t) \gtrsim g(t)$  if  $\liminf_{t \rightarrow \infty} f(t)/g(t) \geq 1$  and  $f(t) \lesssim g(t)$  if  $\limsup_{t \rightarrow \infty} f(t)/g(t) \leq 1$ . Furthermore, we say that  $f(t) = o(g(t))$  if  $\lim_{t \rightarrow \infty} f(t)/g(t) = 0$  and  $f(t) = O(g(t))$  if  $\limsup_{t \rightarrow \infty} f(t)/g(t) < \infty$ . Also, we use the standard definition of an inverse function  $f^{\leftarrow}(x) \triangleq \inf\{y : f(y) > x\}$  for a non-decreasing function  $f(x)$ ; note that the notation  $f(x)^{-1}$  represents  $1/f(x)$ . We use  $\vee$  to denote max, i.e.,  $x \vee y \equiv \max\{x, y\}$ , and  $\wedge$  to denote min, i.e.,  $x \wedge y \equiv \min\{x, y\}$ .

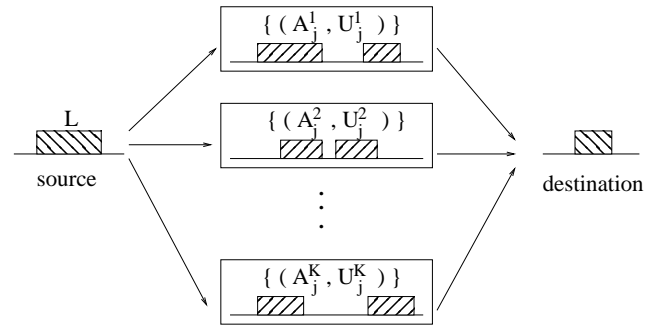


Fig. 1. Multipath transmission over  $K$  channels with failures

#### A. Single path transmission

For the case  $K = 1$ , there is only a single transmission path in the system, hence we let  $A^1 \equiv A$ . The total number of transmissions  $N$  and transmission time  $T = T_r = T_s$  has been studied in [2], [9], [10].

Below we quote Propositions II.1 and II.2 from [9], [10], which show that both  $N$  and  $T$  can follow power law distributions regardless of how heavy or light the tails of  $A$  and  $L$  might be.

**Proposition II.1** *If there exists  $\alpha > 0$  such that*

$$\lim_{x \rightarrow \infty} \frac{\log \mathbb{P}[L > x]}{\log \mathbb{P}[A > x]} = \alpha,$$

then,

$$\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}[N > n]}{\log n} = -\alpha. \quad (1)$$

Additionally, if  $\mathbb{E}[U^{(\alpha \vee 1) + \theta}] < \infty$ ,  $\mathbb{E}[A^{1+\theta}] < \infty$  and  $\mathbb{E}[L^{\alpha+\theta}] < \infty$  for some  $\theta > 0$ , then,

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T > t]}{\log t} = -\alpha. \quad (2)$$

**Proposition II.2** *If*

$$\mathbb{P}[L > x]^{-1} \sim \Phi(\mathbb{P}[A > x]^{-1})$$

where  $\Phi(\cdot)$  is regularly varying with index  $\alpha > 0$ , then, as  $n \rightarrow \infty$ ,

$$\mathbb{P}[N > n] \sim \frac{\Gamma(\alpha + 1)}{\Phi(n)}, \quad (3)$$

and, under the same conditions as in Proposition II.1, as  $t \rightarrow \infty$ ,

$$\mathbb{P}[T > t] \sim \frac{\Gamma(\alpha + 1)(\mathbb{E}[U + A])^\alpha}{\Phi(t)}. \quad (4)$$

**Remark II.1** Proposition II.2 provides more refined results than Proposition II.1 under more restrictive conditions. One can easily check that (3) and (4) imply (1) and (2) by taking logarithms.

**Remark II.2** As mentioned in the introduction, note that the results in the preceding two propositions as well as the ones in the rest of the paper can be readily extended to include packets with bounded sizes using similar techniques as in [14], [15].

### III. REDUNDANT TRANSMISSION

In this section we study redundant transmissions. We begin with  $K$  homogeneous paths, which is followed by the study of the general case of heterogenous paths. We investigate whether sending packets over  $K$  paths can mitigate the power law suffered from single path transmission.

#### A. Homogeneous paths

In this part, we present results for homogeneous paths. We first consider packets of the same size, and then study the more realistic case where packet sizes can be variable.

**Proposition III.1** *If all packets are of constant size  $L \equiv l$  and  $U \equiv 0$ , then,*

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_r > t]}{t} = -K\gamma,$$

where  $\gamma$  is the solution of  $\int_0^1 e^{\gamma x} d\mathbb{P}[A \leq x] = 1$ .

This result can be easily derived using Corollary 3.2 in [2]. From this result, we see that using redundant transmission for equal size packets greatly improves performance, since the decay rate of the delay distribution increases as  $K$  increases, and thus in this case we obtain order improvements in delay performance when using redundant routing. In reality, however, packets are not of equal size. We next present a theorem for the case where the packet size is a random variable.

**Theorem 1** *If*

$$\lim_{x \rightarrow \infty} \frac{\log \bar{F}(x)}{\log \bar{G}(x)} = \alpha,$$

$\mathbb{E}[L^{\alpha+\theta}] < \infty$ ,  $\mathbb{E}[U^{(1 \vee \alpha)+\theta}] < \infty$  and  $\mathbb{E}[A^{1+\theta}] < \infty$  for some  $\theta > 0$ , then,

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_r > t]}{\log t} = -\alpha.$$

**Remark III.1** Comparing the above theorem and Proposition II.1, we observe that, the power law exponent of the total transmission time under redundant transmission is the same as that under single path transmission. Informally speaking, this is because  $T_1, T_2, \dots, T_K$  are not independent, since packets sent over these paths are of the same size.

This theorem is a direct consequence of Theorem 2, which investigates a more general scenario.

#### B. Heterogenous paths

For heterogenous paths, we have the following result when using redundant transmission.

**Theorem 2** *If*

$$\lim_{x \rightarrow \infty} \frac{\log \bar{F}(x)}{\log \bar{G}_j(x)} = \alpha_j \quad (5)$$

for  $1 \leq j \leq K$ , and  $\alpha^* \triangleq \max_{1 \leq j \leq K} \alpha_j > 0$ , then, under the following three conditions I)-III), for some  $\theta > 0$ ,

I)  $\mathbb{E}[L^{\alpha+\theta}] < \infty$ ,

II)  $\max_{1 \leq j \leq K} \mathbb{E}[(U^j)^{(1 \vee \alpha)+\theta}] < \infty$ , and

III)  $\max_{1 \leq j \leq K} \mathbb{E}[(A^j)^{1+\theta}] < \infty$ ,

we have

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_r > t]}{\log t} = -\alpha^*. \quad (6)$$

**Remark III.2** The above theorem implies that the tail behavior of the delay distribution under redundant transmission is determined by the best paths (i.e., the paths with the largest  $\alpha_j$ ).

*Proof of Theorem 2:* First, we establish a lower bound by constructing a new system that has longer available periods than those found on all of the  $K$  paths. The construction is as follows. The new system has an on-off channel characterized by alternating i.i.d. sequences  $\{\bar{A}_i\}$  and  $\{\bar{U}_i\}$ , where

$$\bar{A}_i = \max_{1 \leq j \leq K} A_i^j$$

and  $\bar{U}_i = 0$ . Denote by  $\underline{N}$  the number of transmissions of a packet of length  $L$  over this newly constructed channel.

Now, since  $A_i^j, 1 \leq j \leq K$  are independent, we obtain

$$\mathbb{P}[\bar{A}_i > x] = 1 - \prod_{j=1}^K \mathbb{P}[A_i^j \leq x].$$

Therefore,

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[\bar{A}_i > x]}{\sum_{i=1}^K \bar{G}_j(x)} = 1,$$

coupled with (5), yields

$$\lim_{x \rightarrow \infty} \frac{\log \mathbb{P}[L > x]}{\log \mathbb{P}[A_i > x]} = \alpha^*,$$

which, by Proposition II.1, yields

$$\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}[\underline{N} > n]}{\log n} = -\alpha^*. \quad (7)$$

Define  $\underline{A}_i = \min_{1 \leq j \leq K} A_i^j$  and  $X_i \triangleq \underline{A}_i \mathbf{1}(x_1 < \underline{A}_i < x_2)$ . Choosing  $x_1, x_2$  such that  $\mathbb{E}[X_i] > 0$ , we obtain

$$T_r \geq \sum_{i=1}^{\underline{N}-1} X_i + L. \quad (8)$$

Therefore,

$$\begin{aligned} \mathbb{P}\left[T_r > \frac{t}{\log t}\right] &\geq \mathbb{P}\left[\sum_{i=1}^{\underline{N}-1} X_i > \frac{t}{\log t}\right] \\ &\geq \mathbb{P}\left[\sum_{i=1}^{\underline{N}-1} X_i > \frac{t}{\log t}, \underline{N} > t\right] \\ &\geq \mathbb{P}[\underline{N} > t] - \mathbb{P}\left[\sum_{i=1}^{\underline{N}-1} X_i \leq \frac{t}{\log t}, \underline{N} > t\right] \\ &\geq \mathbb{P}[\underline{N} > t] - \mathbb{P}\left[\sum_{i=1}^{\lfloor t \rfloor} X_i \leq \frac{t}{\log t}\right]. \end{aligned} \quad (9)$$

Since  $\mathbb{E}[e^{\theta X_i}] < \infty$  for some  $\theta > 0$ , we obtain, by a Chernoff bound, for some  $\eta > 0$ ,

$$\begin{aligned} \mathbb{P}\left[\sum_{i=1}^{\lfloor t \rfloor} X_i \leq t/\log t\right] \\ \leq \mathbb{P}\left[\sum_{i=1}^{\lfloor t \rfloor} (\mathbb{E}[X_i] - X_i) \geq \left(\mathbb{E}[X_i] - \frac{1}{\log t}\right)t\right] \\ \leq O(e^{-\eta t}), \end{aligned} \quad (10)$$

which, in combination with (7) and (9), implies

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_r > t]}{\log t} \geq -\alpha^*. \quad (11)$$

Next, we prove the upper bound. Since  $\alpha^* \triangleq \max_{1 \leq j \leq K} \alpha_j > 0$ , there exists  $1 \leq j \leq K$  such that  $\alpha_j = \alpha^*$ . For the  $j$ th path, we have  $T_r \leq T_j$  since  $T_r = \min\{T_1, T_2, \dots, T_K\}$ . Using Proposition II.1, we obtain

$$\overline{\lim}_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_r > t]}{\log t} \leq \lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_j > t]}{\log t} = -\alpha^*. \quad (12)$$

By combining (11) and (12), we complete the proof. ■

Our preceding result characterizes the performance in terms of the ‘‘logarithmic asymptotics’’. Basically, it only contains information about the power law exponent, but yields no information about the pre-factor before the power law term. As a consequence, this result cannot distinguish between redundant transmission and single path transmission. In order

to investigate the performance improvement for redundant transmission, we need a more refined asymptotic result. For a set of regularly varying functions  $\Phi_j(\cdot)$ ,  $1 \leq j \leq K$ , we can compute the exact asymptotic tail of the distribution of  $T_r$ .

**Theorem 3** *If  $\bar{F}(x)^{-1} \sim \Phi_j(\bar{G}_j(x)^{-1})$  and*

$$\lim_{x \rightarrow \infty} \frac{\Phi_j(x)}{\Phi(x)} = \zeta_j > 0, \quad (13)$$

where  $\Phi(\cdot)$  is regularly varying with index  $\alpha > 0$ , then, under the conditions I)-III) in Theorem 2, as  $t \rightarrow \infty$ ,

$$\mathbb{P}[T_r > t] \sim \frac{\Gamma(\alpha + 1)}{\left(\sum_{j=1}^K (\mathbb{E}[A^j + U^j])^{-1} \zeta_j^{1/\alpha}\right)^\alpha} \frac{1}{\Phi(t)}. \quad (14)$$

**Remark III.3** From the preceding result, we see that, redundant transmission improves the system performance by reducing the tail of the distribution by a constant factor. If these  $K$  channels are i.i.d., this constant is equal to  $K^\alpha$ .

In order to prove the theorem, we need the following lemma.

**Lemma 1** *For  $\eta_j > 0$ ,  $1 \leq j \leq K$ ,*

$$\begin{aligned} \mathbb{P}[N_1 > \eta_1 t, N_2 > \eta_2 t, \dots, N_K > \eta_K t] \\ \sim \frac{\Gamma(\alpha + 1)}{\left(\sum_{j=1}^K \eta_j \zeta_j^{1/\alpha}\right)^\alpha} \frac{1}{\Phi(t)}. \end{aligned} \quad (15)$$

The proof of this lemma can be found in the Technical Report [16].

*Proof of Theorem 3:* Due to limited space, we only present the proof of the upper bound. The proof of the lower bound is similar to the upper bound and can be found in [16].

For  $0 < \epsilon < 1$  and  $\eta_j = 1/\mathbb{E}[A^j + U^j]$ , we obtain,

$$\begin{aligned} \mathbb{P}[T_r > (1 + 2\epsilon)t] &= \mathbb{P}\left[\bigcap_{j=1}^K \{T_j > (1 + 2\epsilon)t\}\right] \\ &= \mathbb{P}\left[\bigcap_{j=1}^K \left\{\sum_{i=1}^{N_j-1} (A_i^j + U_i^j) + L > (1 + 2\epsilon)t\right\}\right] \\ &\leq \mathbb{P}\left[\bigcap_{j=1}^K \left\{\sum_{i=1}^{N_j} (A_i^j + \mathbb{E}[U^j]) > t\right\}\right] \\ &\quad + \mathbb{P}\left[\bigcup_{j=1}^K \left\{\sum_{i=1}^{N_j} (U_i^j - \mathbb{E}[U^j]) > \epsilon t\right\}\right] \\ &\quad + \mathbb{P}[L > \epsilon t]. \end{aligned} \quad (16)$$

Then, using union bound, we derive

$$\mathbb{P}[T_r > (1 + 2\epsilon)t] \leq \mathbb{P}\left[\bigcap_{j=1}^K \{N_j > (1 - \epsilon)\eta_j t\}\right]$$

$$\begin{aligned}
& + \sum_{j=1}^K \mathbb{P} \left[ \sum_{i=1}^{(1-\epsilon)\eta_j t} \left( A_i^j \wedge L + \mathbb{E}[U^j] \right) > t \right] \\
& + \sum_{j=1}^K \mathbb{P} \left[ \left\{ \sum_{i=1}^{N_j} \left( U_i^j - \mathbb{E}[U^j] \right) > \epsilon t \right\} \right] \\
& + \mathbb{P}[L > \epsilon t] \\
& \triangleq I_1 + I_2 + I_3 + I_4. \tag{17}
\end{aligned}$$

Using the result (4.20) in [10], we know  $I_2 + I_3 + I_4 = o(1/\Phi(t))$ , which, in view of Lemma 1, yields

$$\mathbb{P}[T_r > t] \lesssim \frac{\Gamma(\alpha + 1)}{\left( \sum_{j=1}^K (\mathbb{E}[A^j + U^j])^{-1} \zeta_j^{1/\alpha} \right)^\alpha} \frac{1}{\Phi(t)}. \tag{18}$$

#### IV. SPLIT TRANSMISSION

Next, we study the case when a packet is split into several pieces and sent over  $K$  independent paths. Using the derived results, we will determine which of the two strategies, split transmission or redundant transmission, results in a lighter distribution tail.

We begin with homogeneous paths, and then investigate heterogenous paths. A fraction  $\gamma_j$  of the packet  $L$  is sent over path  $j$ ,  $\sum_{j=1}^K \gamma_j = 1$ ,  $0 \leq \gamma_j \leq 1$ ,  $1 \leq j \leq K$ . We derive the optimal splitting strategy that minimizes the exponent of the transmission time tail.

##### A. Homogeneous paths

We have the following theorem for split transmission over homogenous paths, where each packet is evenly split into  $K$  pieces. Its proof is a special case of that for heterogeneous paths (see Theorem 5), and hence is omitted.

**Theorem 4** *Under the same conditions in Theorem 1, if there exists  $\beta > 0$ , such that*

$$\lim_{x \rightarrow \infty} \frac{\log \bar{F}(Kx)}{\log \bar{F}(x)} = \beta, \tag{19}$$

then,

$$\lim_{t \rightarrow \infty} \frac{\log P(T_s > t)}{\log t} = -\beta\alpha.$$

**Remark IV.1** Since  $\beta \geq 1$ , comparing the results in Proposition II.1 and Theorem 1, we see that, for homogeneous paths, split transmission is no worse than redundant transmission when packets are split evenly. Split transmission is not beneficial when  $\beta = 1$ , e.g., when  $\log \bar{F}(x)$  is a slowly varying function.

##### B. Heterogenous paths

For heterogenous paths, a packet of size  $L$  is split into  $K$  smaller fragments of sizes  $\gamma_1 L, \gamma_2 L, \dots, \gamma_K L$ , respectively, where  $\sum_{j=1}^K \gamma_j = 1$ ,  $0 \leq \gamma_j \leq 1$ ,  $1 \leq j \leq K$ . We have the following result on packet transmission delay.

**Theorem 5** *If there exist  $\alpha_j, \beta_j$ ,  $j = 1, 2, \dots, K$  such that*

$$\lim_{x \rightarrow \infty} \frac{\log \bar{F}(x)}{\log \bar{G}_j(x)} = \alpha_j, \tag{20}$$

$$\lim_{x \rightarrow \infty} \frac{\log \bar{F}(x)}{\log \bar{F}(\gamma_j x)} = \beta_j, \tag{21}$$

with  $\alpha^\circ \triangleq \min_{1 \leq j \leq K} \beta_j \alpha_j > 0$ , then,

$$\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}[N > n]}{\log n} = -\alpha^\circ,$$

and, under the conditions I)-III) in Theorem 2,

$$\lim_{t \rightarrow \infty} \frac{\log \mathbb{P}[T_s > t]}{\log t} = -\alpha^\circ.$$

**Remark IV.2** When paths are heterogeneous, the packet transmission delay is determined by the best paths under redundant transmission and by the worst paths under split transmission. On the other hand, split transmission only sends a fraction of the packet on each path. Comparing this to Theorem 2, we see that, if  $\min_{1 \leq j \leq K} \beta_j \alpha_j > \max_{1 \leq j \leq K} \alpha_j$ , split transmission is more beneficial than redundant transmission in minimizing the tail behavior; otherwise, redundant transmission is more beneficial. We will show later that, by carefully choosing the way to split packets, split transmission can always result in tail performance that is no worse than redundant transmission.

*Proof of Theorem 5:* We begin with proving the result for  $T_s$ . Since

$$T_s = \max_{1 \leq j \leq K} T_j,$$

we obtain, using a union bound,

$$\max_{1 \leq j \leq K} \mathbb{P}[T_j > t] \leq \mathbb{P}[T_s > t] \leq \sum_{j=1}^K \mathbb{P}[T_j > t], \tag{22}$$

Next, using (20) and (21), we derive

$$\lim_{x \rightarrow \infty} \frac{\log P(\gamma_j L > x)}{\log P(A_j > x)} = \lim_{x \rightarrow \infty} \frac{\beta_k \log \bar{F}(x)}{\log \bar{G}_k(x)} = \beta_j \alpha_j,$$

which, by Proposition II.1, yields

$$\lim_{t \rightarrow \infty} \frac{\log P(T_j > t)}{\log t} = -\beta_j \alpha_j.$$

Thus, for  $\epsilon > 0$ , there exists  $t_0 > 0$  such that for all  $t > t_0$ ,

$$-\beta_j \alpha_j - \epsilon < \frac{\log \mathbb{P}[T_j > t]}{\log t} < -\beta_j \alpha_j + \epsilon.$$

Hence, for  $t > t_0$ , we have

$$\max_{1 \leq j \leq K} \mathbb{P}[T_j > t] > t^{-\alpha^\circ - \epsilon}$$

and

$$\sum_{j=1}^K \mathbb{P}[T_j > t] < K t^{-\alpha^\circ + \epsilon},$$

which, combined with (22) and passing  $\epsilon \rightarrow 0$ , yields

$$\lim_{t \rightarrow \infty} \frac{\log P(T_s > t)}{\log t} = - \min_{1 \leq j \leq K} \{\beta_j \alpha_j\} = -\alpha^\circ.$$

Now, we derive the result for  $N$ . Since  $N_s = \sum_{j=1}^K N_j$ , we have

$$\max_{1 \leq j \leq K} \mathbb{P}[N_j > n] \leq \mathbb{P}[N > n] \leq \sum_{j=1}^K \mathbb{P}\left[N_j > \frac{n}{K}\right]. \quad (23)$$

Proposition II.1 implies

$$\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}[N_j > n/K]}{\log n} = -\beta_j \alpha_j,$$

which, combined with (23) and using a similar argument as in proving the result for  $T_s$ , yields  $\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}[N > n]}{\log n} = -\alpha^\circ$ . ■

1) *Optimal split transmission:* From Theorem 5, we can see that in order to optimize the power law delay tail, we need to choose  $\gamma_1, \gamma_2, \dots, \gamma_K$  so that  $\min_{1 \leq j \leq K} \beta_j \alpha_j$  is maximized. To achieve this, we may speculate that we need to choose  $\gamma_1, \gamma_2, \dots, \gamma_K$  so that  $\beta_1 \alpha_1 = \beta_2 \alpha_2 = \dots = \beta_K \alpha_K$ . The following theorem confirms that this is true when  $\log(1/\bar{F}(x))$  is not slowly varying. [6] is a related work on optimal file split under a different problem setting.

**Theorem 6** Suppose we use split transmission over  $K$  heterogeneous paths, each satisfying (20). If the limit

$$\beta(\gamma) = \lim_{x \rightarrow \infty} \frac{\log \bar{F}(x)}{\log \bar{F}(\gamma x)}$$

exists for all  $0 < \gamma < 1$ , then (i) there exists a unique constant  $\rho \geq 0$  such  $\beta(\gamma) = \gamma^{-\rho}$ ; and (ii) the optimal splitting scheme that minimizes the power law exponent of  $\mathbb{P}[T_s > t]$  satisfies:

a) If  $\rho > 0$ , then

$$\gamma_j^* = \frac{\alpha_j^{1/\rho}}{\sum_{i=1}^K \alpha_i^{1/\rho}}. \quad (24)$$

b) If  $\rho = 0$ , then let  $\gamma_j = 0$  for  $\alpha_j \neq \max_{1 \leq j \leq K} \alpha_j$  and the other  $\gamma_j$  can take arbitrary values.

The corresponding optimal power law exponent for  $\mathbb{P}[T_s > t]$  is  $-\alpha_\rho$ , where

$$\alpha_\rho = \begin{cases} \left( \sum_{i=1}^K \alpha_i^{1/\rho} \right)^\rho, & \rho > 0, \\ \max_{1 \leq j \leq K} \alpha_j, & \rho = 0. \end{cases} \quad (25)$$

**Remark IV.3** In the preceding result, we only minimize the power law exponent. When  $\rho = 0$ , we have  $\beta(\gamma) = 1$ , and  $\log(1/\bar{F}(x))$  is a slowly varying function. In this case, we should only use the best paths, and the scheme in (24) is to split arbitrarily among the best paths. For this case, we need a more refined asymptotic result that accounts for not only the power law exponent but also the exact pre-factors to derive the optimal split strategy. Due to limited space, we do not study this problem. When  $\rho > 0$ , all the channels are utilized, and

the optimal fraction on each path is specified by (24). In this case, one can easily check that the optimal tail exponent is indeed achieved when  $\beta_1 \alpha_1 = \beta_2 \alpha_2 = \dots = \beta_K \alpha_K$ .

**Remark IV.4** Note that  $\alpha_\rho = \left( \sum_{i=1}^K \alpha_i^{1/\rho} \right)^\rho \geq \alpha^*$  with equality if and only if  $\rho = 0$ , where  $\alpha^* = \max_{1 \leq j \leq K} \alpha_j > 0$ , as defined in Theorem 2. Thus, under the assumption of Theorem 5, split transmission achieves a better exponent than redundant transmission if  $\rho > 0$ .

*Proof of Theorem 6:* (i) Note that  $\beta(\gamma) \geq 1$  on  $(0, 1)$ . If  $\beta(\gamma) = 1$  for all  $\gamma \in (0, 1)$ , then  $\beta(\gamma) = \gamma^{-\rho}$  for  $\rho = 0$ . Now assume  $\beta_0 = \beta(\gamma_0) > 1$  for some  $\gamma_0 \in (0, 1)$ . Observe that  $\beta(\gamma_1 \gamma_2) = \beta(\gamma_1) \beta(\gamma_2)$  for any  $\gamma_1, \gamma_2 \in (0, 1)$ . Thus, for any positive integer  $m, n$ ,

$$\beta(\gamma_0^{m/n}) = \left( \beta(\gamma_0^{1/n}) \right)^{n \times m/n} = \left( \beta \left( (\gamma_0^{1/n})^n \right) \right)^{m/n} = \beta_0^{m/n}.$$

Since  $\beta$  is monotonically decreasing and the positive rationals are dense in  $\mathbb{R}^+$ ,

$$\beta(\gamma_0^r) = \beta_0^r, \quad r \in \mathbb{R}^+$$

or, equivalently,

$$\beta(\gamma) = \gamma^{\log \beta_0 / \log \gamma_0} = \gamma^{-\rho}, \quad \gamma \in (0, 1)$$

where  $\rho = -\log \beta_0 / \log \gamma_0 > 0$ . It is clear that  $\rho$  is unique.

(ii) Let  $\{\gamma_j^*\}$  be an optimal split scheme and  $-\alpha_\rho$  the corresponding optimal exponent. By Theorem 5,

$$\alpha_\rho = \min_{j: \gamma_j > 0} \alpha_j (\gamma_j^*)^{-\rho}. \quad (26)$$

If  $\rho = 0$ , then

$$\alpha_\rho = \min_{j: \gamma_j > 0} \alpha_j \leq \max_{1 \leq j \leq K} \alpha_j = \alpha^*$$

with equality if and only if  $\gamma_j = 0$  whenever  $\alpha_j \neq \alpha^*$ .

If  $\rho > 0$ , then (26) gives

$$\gamma_j^* (\alpha_\rho)^{1/\rho} \leq \alpha_j^{1/\rho}, \quad j = 1, 2, \dots, K.$$

Summing over  $j$  and noting  $\sum_j \gamma_j^* = 1$ , we have  $(\alpha^*)^{1/\rho} \leq \sum_{j=1}^K \alpha_j^{1/\rho}$  with equality if  $\gamma_j^*$  is given by (24). ■

2) *Optimal split transmission examples:* To illustrate the results obtained in the preceding section, we compute the optimal split transmission scheme for some typical distributions.

• Weibull distribution. If

$$\begin{aligned} \bar{F}(x) &= P(L > x) = e^{-(\lambda x)^b}, \\ \bar{G}_j(x) &= P(A^j > x) = e^{-(\mu_j x)^b}, \end{aligned}$$

where  $\lambda > 0, \mu_j > 0$ , and  $b > 0$ , then,

$$\alpha_j = \frac{\log \bar{F}(x)}{\log \bar{G}_j(x)} = \frac{-(\lambda x)^b}{-(\mu_j x)^b} = \left( \frac{\lambda}{\mu_j} \right)^b,$$

$$\beta(\gamma) = \frac{\log \bar{F}(x)}{\log \bar{F}(\gamma x)} = \frac{1}{\gamma^b},$$

and

$$\rho = -\log \beta(\gamma) / \log \gamma = b.$$

Therefore, the optimal split is

$$\gamma_j = \frac{\left(\frac{\lambda}{\mu_j}\right)^{1/b}}{\sum_{i=1}^K \left(\frac{\lambda}{\mu_i}\right)^{1/b}} = \frac{\mu_j^{-1/b}}{\sum_{i=1}^K \mu_i^{-1/b}}, \quad j = 1, \dots, K.$$

- Pareto distribution. Consider the case where the size of the packet,  $L$ , and the available time period on path  $j$ ,  $A^j$ , follow Pareto distributions. In this case, we have  $\beta(\gamma) = 1$ . The optimal split transmission strategy is to split among the best paths.

### 3) Exact asymptotic result for optimal split transmission:

Our proposed optimal split transmission minimizes the power law exponent of  $\mathbb{P}[T_s > t]$ . In other words, Theorem 6 only characterizes the tail behavior in the logarithmic scale. Next, to refine the result, we present a theorem on the exact asymptotic result for optimal split transmission. The proof is presented in the Technical Report [16].

**Theorem 7** *If  $\log(\bar{F}(x)^{-1}) = x^\rho l(x)$  where  $\rho > 0$  and  $l(x)$  is slowly varying with*

$$e^{l(x)} \sim e^{l(\gamma x)}$$

for  $\gamma > 0$ , and

$$\bar{F}(x)^{-1} \sim \zeta_j \left(\Phi(\bar{G}_j(x)^{-1})\right)^{\alpha_j/\alpha},$$

where  $\alpha_j, \zeta_j > 0$  and  $\Phi(\cdot)$  is regularly varying with index  $\alpha > 0$ , then, under the conditions I-III) in Theorem 2, as  $t \rightarrow \infty$ ,

$$\mathbb{P}[T_s > t] \sim \sum_{l=1}^K (-1)^{l+1} \sum_{\{j_1, \dots, j_l\} \subseteq \{1, \dots, K\}} \frac{\Gamma(\alpha_\rho + 1)}{\left(\sum_{s=1}^l \eta_{j_s} \zeta_{j_s}^{\frac{1}{\alpha_{j_s}}}\right)^{\alpha_\rho}} \frac{1}{\Phi(t)^{\frac{\alpha_\rho}{\alpha}}},$$

where  $\eta_j \triangleq 1/\mathbb{E}[A^j + U^j]$ ,  $\alpha_\rho \triangleq \left(\sum_{j=1}^K \alpha_j^{1/\rho}\right)^\rho$ .

## V. CONCLUSION

Parallelism is a common approach to improve reliability and efficiency in practice. In this paper, we investigate whether and how parallelism can be used to improve network performance. Specifically, we study whether and how multipath transmission can mitigate power law delays. We show that, when all packets are of the same size, redundant transmission can greatly reduce the transmission delay in the sense that the ratio of the delay distribution tail with and without redundant transmission tends to zero. However, when packet sizes are random variables such that  $\log \mathbb{P}[L > x] \approx \alpha^* \log \mathbb{P}[A^j > x]$ , we prove that, maybe counter intuitively, redundant transmission cannot change the order of the probability tail of the transmission delays, and can only improve the system performance by a constant factor. We also show that split transmission is

effective in mitigating power delays if the absolute value of the logarithm of the packet size probability tail is regularly varying with positive index, and becomes ineffective if the above quantity is slowly varying. Last, we provide an optimal split transmission strategy when the paths are heterogeneous, and further derive an exact asymptotic result for packet delivery time under this scheme. Our results can be extended to many other applications that involve parallelism and job failures, such as computing jobs in grid computing, file downloading in peer to peer networks, parallel experiment planning, and parallel scheduling.

## REFERENCES

- [1] L. N. Andersen and S. Asmussen. Parallel computing, failure recovery and extreme values. *J. Statist. Theory Appl.*, (2):279–292, 2008.
- [2] S. Asmussen, P. Fiorini, L. Lipsky, T. Rolski, and R. Sheahan. Asymptotic behavior of total times for jobs that must start over if a failure occurs. *Mathematics of Operations Research*, 33(4):932–944, November 2008.
- [3] N. H. Bingham, C. M. Goldie, and J. L. Teugels. *Regular Variation*, volume 27. Cambridge University Press, 1987.
- [4] D. G. Cantor and M. Gerla. Optimal transmission in a packet-switched computer network. *IEEE Transactions on Computers*, 23(10), October 1974.
- [5] R. Gallager. A minimum delay transmission algorithm using distributed computation. *IEEE Transactions on Communications*, 25(1), 1977.
- [6] G. Hoekstra, R. Mei, Y. Nazarathy, and B. Zwart. Optimal file splitting for wireless networks with concurrent access. In *NET-COOP '09: Proceedings of the 3rd Euro-NF Conference on Network Control and Optimization*, pages 189–203, Berlin, Heidelberg, 2009. Springer-Verlag.
- [7] P. R. Jelenković and J. Tan. Is ALOHA causing power law delays? In *Proceedings of the 20th International Teletraffic Congress*, Ottawa, Canada, June 2007; *Lecture Notes in Computer Science*, No 4516, pp. 1149–1160, Springer-Verlag, 2007.
- [8] P. R. Jelenković and J. Tan. Are end-to-end acknowledgements causing power law delays in large multi-hop networks? In *14th Inform Applied Probability Conference*, Eindhoven, July 9–11 2007.
- [9] P. R. Jelenković and J. Tan. Can retransmissions of superexponential documents cause subexponential delays? In *Proceedings of IEEE INFOCOM'07*, pages 892–900, Anchorage, Alaska, USA, (submitted August 2006, accepted November 2006), May 2007.
- [10] P. R. Jelenković and J. Tan. Characterizing heavy-tailed distributions induced by retransmissions. Technical report, Department of Electrical Engineering, Columbia University, EE2007-09-07, September 2007. Eprint arXiv: 0709.1138v2.
- [11] L. Kleinrock. *Communication Nets: Stochastic Message Flow and Delay*. New York: McGraw-Hill Book Company, 1964.
- [12] T. S. Rappaport. *Wireless communications: principles and practise*. Prentice Hall, 2 edition, January 2002.
- [13] R. Sheahan, L. Lipsky, P. Fiorini, and S. Asmussen. On the completion time distribution for tasks that must restart from the beginning if a failure occurs. *MAMA 2006 Workshop, Saint-Malo, France*, June 2006.
- [14] J. Tan and N. Shroff. Transition from heavy to light tails in retransmission durations. In *Proceedings of IEEE INFOCOM'2010*, San Diego, California, March 2010.
- [15] J. Tan and N. Shroff. Transition from heavy to light tails in retransmission durations. Technical report, Departments of Electrical and Computer Engineering & Computer Science and Engineering, The Ohio State University, Columbus, OH, 2010.
- [16] J. Tan, W. Wei, B. Jiang, N. Shroff, and D. Towsley. Can multipath routing mitigate power law delays? - effects of parallelism on tail performance. Technical Report UM-CS-2009-055, Department of Computer Science, University of Massachusetts, Amherst, MA, 2009.