# Reliability Equations for Cloud Storage Systems with Proactive Fault Tolerance

Jing Li<sup>10</sup>, Peng Li, Rebecca J. Stones<sup>10</sup>, Gang Wang, *Member, IEEE*, Zhongwei Li, and Xiaoguang Liu

Abstract—As cloud storage systems increase in scale, hard drive failures are becoming more frequent, which raises reliability issues. In addition to traditional reactive fault tolerance, proactive fault tolerance is used to improve a system's reliability. However, there are few studies which analyze the reliability of proactive cloud storage systems, and they typically assume an exponential distribution for drive failures. This paper presents closed-form equations for estimating the number of data-loss events in proactive cloud storage systems using RAID-5, RAID-6, 2-way replication, and 3-way replication mechanisms, within a given time period. The equations model the impact of proactive fault tolerance, operational failures, failure restorations, latent block defects, and drive scrubbing on the systems reliability, and use time-based Weibull distributions to represent processes (instead of homogeneous Poisson processes). We also design a Monte-Carlo simulation method to simulate the running of proactive cloud storage systems. The proposed equations closely match time-consuming Monte-Carlo simulations, using parameters obtained from the analysis of field data. These equations allow designers to efficiently estimate system reliability under varying parameters, facilitating cloud storage system design.

Index Terms—proactive fault tolerance, cloud storage systems, reliability, time-variant failure rates, latent block defects

#### 15 **1** INTRODUCTION

3

5

6 7

8

g

10

11

12

13

14

ODERN-DAY data centers usually host hundreds of 16 L thousands of servers, using hard drives as the primary 17 data storage device. Many challenges are faced for such large 18 data center management [1], [2]. The failure of an individual 19 hard drive might be rare, but a system with thousands of 20 21 drives will regularly experience failures [3], [4], [5], [6]. Drive failure can result in service unavailability, hurting the user 22 experience, and even permanent data loss. Therefore, high 23 reliability is one of the biggest concerns in such systems. 24

25 Traditional cloud storage systems adopt redundancy, e.g., 26 erasure codes and replication, to reconstruct data when drive failure occurs, which is known as *reactive* fault tolerance. To 27 provide satisfactory reliability in large-scale cloud storage 28 systems with high failure frequency, multi-erasure codes (or 29 multiple replicas) must be used, which brings high construc-30 tion and maintenance cost and heavy read/write overhead. 31 Thus, reactive fault tolerance alone cannot meet the demands 32 of the high reliability and service quality in modern data cen-33 ters. Proactive fault tolerance [7], [8], [9], [10], [11], [12], [13], 34 [14], [15], [16], [17] instead predicts drive failures and handles 35 them in advance; with sufficient prediction accuracy and 36

Manuscript received 18 May 2017; revised 13 Aug. 2018; accepted 15 Nov. 2018. Date of publication 0 . 0000; date of current version 0 . 0000. (Corresponding author: Xiaoguang Liu.) Recommended for acceptance by S. Kundu.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TDSC.2018.2882512 effective warning handling, it can significantly enhance the 37 system reliability and reduce costs. 38

When designing proactive cloud storage systems and 39 tweaking parameters to optimize performance, designers 40 must consider factors such as coding redundancy, failure pre-41 diction, the amount of bandwidth used to reconstruct or 42 migrate data after a drive fails or is predicted to fail, and drive 43 scrubbing for eliminating block defects. As such, designers 44 could benefit from an accurate and easy-to-use way to assess 45 the effects of these factors on overall reliability. 46

So far there are only a few relevant studies which analyze 47 the reliability of proactive cloud storage systems [11], [18], 48 [19]. There are some drawbacks to the current research: (a) 49 Inaccurate failure distribution models—the reliability esti- 50 mates are based on the assumption that both hard drive 51 failures and their repairs follow a homogeneous Poisson 52 process with constant failure and restoration rates, which 53 was contested in [20], [21], [22], [23]. (b) Incomplete consid- 54 eration of failures—focusing on whole-drive failures, with-55 out incorporating latent block level failure mode, which, 56 with increases in single-drive and whole-system capacity, 57 can not be ignored [24]. (c) Unrealistic reliability metric— 58 mean time to data loss (MTTDL) is excessive relative to the 59 actual run time of cloud storage systems and does not ade- 60 quately reflect reliability [25].

One could simulate the system to assess reliability more 62 accurately, but at the expense of usability, which would 63 require specialized code, extensive computation, and is 64 time consuming. It is preferable to have reliability equations 65 which are easy to use, applicable to different configurations 66 and executed quickly, especially when frequently tweaking 67 the parameters of a cloud storage system where the executed workload and rates of failures, warnings, etc., fluctu- 69 ate over time. 70

J. Li is with the College of Computer Science and Technology, Civil Aviation University of China, Tianjin, China. E-mail: lijing@nbjl.nankai.edu.cn.

<sup>•</sup> P. Li, R.J. Stones, G. Wang, Z. Li, and X. Liu are with Nankai-Baidu Joint Lab, College of Computer and Control Engineering, Nankai University, Tianjin, China.

E-mail: {lipeng, becky, wgzwp, lizhongwei, liuxg}@nbjl.nankai.edu.cn.



TABLE 1 Mathematical Models for the Expected Number of Data Loss Events in RAID-5, RAID-6, 2-way Replication, and 3-way Replication Systems, Each With Proactive Fault Tolerance

Elerath et al. defined two reliability equations for reac-71 tive RAID-5 [22] and RAID-6 [23] groups without disk 72 failure prediction. In this paper, we extend the work to 73 proactive RAID systems and proactive replication systems 74 with disk failure prediction; in these settings the reliabil-75 ity analysis is more intricate due to the need to factor in 76 77 failure prediction and replica dispersal. We make two main contributions: (a) to incorporate drive failure predic-78 79 tion, we modify the calculation for the cumulative hazard 80 rate and drive availability; and (b) to incorporate replica 81 dispersal, we mathematically derive the probability of 82 data loss.

Specifically, this work generalizes the equations by 83 Elerath et al. [22], [23] for assessing the reliability of proac-84 tive RAID-5 and RAID-6 systems, and we further propose 85 two new equations for assessing the reliability of proactive 86 2-way replication, and 3-way replication systems; these four 87 systems are the most commonly used methods in current 88 data centers. The equations allow for expressions of time-89 dependent failure and restoration rates, and incorporate 90 block defects, media scrubbing processes, and proactive 91 fault tolerance on reliability. We also design an event-driven 92 Monte-Carlo-based simulation method to simulate cloud 93 storage systems with proactive fault tolerance. The pro-94 posed equations and simulations are consistent with one 95 another. Using these equations, system designers can easily 96 97 assess trade-offs, compare schemes, and understand the effects of the parameters on the overall cloud storage system 98 99 reliability, allowing them to better design and optimize infrastructures. 100

The rest of the paper is organized as follows: Section 2 describes the relevant background knowledge. The equations we present are listed in Table 1, and their derivation is given in Section 3. We evaluate their accuracy versus simulations in Section 4, and explore these systems' sensitivity to varying parameters. Section 5 describes some limitations and the efficacy of the models.

## 2 BACKGROUND

#### 2.1 Related Work

Reliability, the focus of this paper, is one of the most important aspects of storage systems and has been studied extensively, especially for RAID systems.

108

109

Gibson et al. [26] found that hard drive failure rates typically followed a "bathtub curve". However they still 114 considered the exponential distribution a useful simplifying 115 assumption for modeling drive failure events. Recently, 116 researchers found that drive failure events were not ade-117 quately modeled by homogeneous Poisson processes [20], 118 [21], [22], [23], and some other work focus on device 119 failure [27]. 120

In this paper, we use Weibull distributions for modeling 121 drive failure events. This is motivated by Schroeder and 122 Gibson [20], who found that hard drive failure rates were 123 not constant with age, and recommended the Weibull distribution for modeling drive failure, which can account for 125 both "infant mortality" and aging drives. 126

Elerath et al. defined two equations for assessing the reliability of RAID-5 [22] and RAID-6 [23] groups, and these 128 papers form the basis of the equations proposed in this 129 paper. Greenan et al. [28] argued that MTTDL was a bad 130 reliability metric, and Elerath et al. went so far as to say it 131 should be "put to rest". However, Iliadis and Venkate-132 san [29] offered a rebuttal. In this paper, we use the 133 expected number of data loss events to measure reliability, 134 consistent with Elerath et al. Elerath et al.'s equations use 135 time-dependent failure and repair rates and included the 136 contributions of both sector defects and data scrubbing. 137 They considered RAID groups without disk failure predic-138 tion. In this work, we extend this work to include (a) 139 proactive fault tolerance and (b) replication systems. 140

For proactive cloud storage systems, there are only a few 141 studies focusing on their reliability. Eckart et al. [18] used 142 Markov models to demonstrate the effect of failure prediction 143 has on a system's MTTDL. They devised models for a single
hard drive, RAID-1, and RAID-5 with proactive fault tolerance. Li et al. extended this study to RAID-6 groups [11] and
replication systems [19]. However, there are some drawbacks
in those studies on reliability of proactive systems (as mentioned in the introduction), which we overcome in this paper.

#### 150 2.2 Drive Failure Modes

*Intermittent failure* is the most common mode of failure. In
this case, a drive or some sectors cannot be accessed temporarily, and are often restored by retrying several times.

Latent block defects are commonly caused by latent sector 154 errors, a permanent inability to access data from certain sec-155 tors (possibly due to physical defects e.g., a scratch), and 156 data corruption, where data stored in a block are incorrect. 157 Latent sector errors are not reported by the drive until the 158 particular sector is accessed. Data corruption can not be 159 reported by the drive even when a defective block is read; it 160 is silent and could have greater impact than other errors. 161

To detect and protect data against the block defects, 162 cloud storage systems usually perform drive scrubbing 163 during idle periods. Drive scrubbing is a background pro-164 cess that proactively reads and checks data from all drive 165 166 blocks. If a defective block is detected, the system recon-167 structs the corrupted content from the available data. The time required to scrub an entire drive varies with the drive 168 169 capacity and the drive scrubbing rate.

A serious type of failure is an *operational failure*, where a whole drive is permanently no longer accessible. Such failure can be repaired only by replacing the drive. Proactive fault tolerance only protects against data loss caused by operational failures.

When an operational failure occurs, the cloud storage 175 system initiates a rebuild process during which it restores 176 the missing data using the accessible surviving data. The 177 rebuild time depends on the amount of data that is trans-178 ferred during the process, and the data transfer rate. More-179 over, to maintain the quality of user service, usually only a 180 fraction of the total bandwidth available is used for a 181 182 rebuild process. As such, the rebuild time will also be influ-183 enced by the foreground activity.

Only the last two failure types—latent block defects and operational failures—impact a storage system's reliability, so we focus on them in this paper.

#### 187 2.3 Proactive Fault Tolerance

Self-Monitoring, Analysis, and Reporting Technology 188 (SMART) is implemented within modern hard drives [30]. 189 SMART monitors and compares drive attributes with pre-190 191 set thresholds, and issues warnings when attributes exceed the thresholds. As a result, systems can act in advance of 192 drive failures, such as by migrating data. This typifies 193 proactive fault tolerance, which fundamentally improves 194 195 system reliability.

Moreover, to improve prediction accuracy, statistical and machine learning methods have been proposed to build hard drive failure prediction models based on SMART attributes [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], some of which achieve good prediction performance, which is reasonable to use in practice. Existing drive failure prediction models are mostly only 202 able to predict operational failures in advance; block-defect 203 prediction models are at an early stage in development and 204 do not have suitable practical performance. For example, 205 recent work by Mahdisoltani et al. [31] does not predict the 206 location of block defects, but only if a disk will incur a 207 block defect. Currently, drive scrubbing and reactive fault 208 tolerance are used to cope with block defects. Thus, this 209 paper only considers the prediction of operational failures 210 (and not the prediction of block defects). 211

In a proactive cloud storage system, a drive failure pre- 212 diction model runs in the background and monitors the 213 drives in real time (with a minor resource cost [14], [16]), 214 periodically outputting their health states (such as once an 215 hour). When an alarm is raised by the model, the data on an 216 at-risk drive is ordinarily migrated (or backed up) to 217 healthy drives immediately. 218

There are two prediction deployment schemes: (a) 219 intra-drive prediction, in which the backups are dealt 220 with locally by the drive; and (b) intra-system prediction, 221 in which backups are dealt with by the system. The latter 222 has greater flexibility, so we study proactive cloud stor- 223 age systems with intra-system prediction. Since failure 224 prediction eliminates most drive failures, it reduces the 225 rate of data loss events. 226

In this paper, we use a simple model for proactive fault 227 tolerance, where a proportion of drives that are about to fail 228 are predicted to fail in advance; the proportion is called the 229 *failure detection rate* (FDR). Rebuild processes for drives that 230 are predicted to fail are assumed to be completed before 231 failure actually occurs. 232

#### 2.4 Reactive Fault Tolerance

In a real-world setting, some drives will inevitably incur 234 operational failure in proactive cloud storage systems since 235 the FDR will not be 100 percent and it takes time to 236 migrate data. It is therefore necessary for proactive cloud 237 storage systems to also use reactive fault tolerance to ensure 238 reliability. In this paper, we include four different reactive 239 fault tolerance methods: RAID-5, RAID-6, 2-way replication, 240 and 3-way replication. 241

In the presence of reactive fault tolerance, data loss 242 requires simultaneous operational failures and/or block 243 defects. Moreover, these failures will not result in data 244 loss if rebuilding has finished and/or scrubbing has taken 245 place. 246

RAID-5 and RAID-6 are popular RAID schemes, in 247 which hard drives are arranged in *RAID groups* of *g* drives. 248 Data *stripes* are spread across multiple drives and are 249 accessed in parallel. Each stripe of a RAID-5 group can 250 tolerate a single failure—an operational failure or block 251 defect—but data loss may occur as the result of simulta-252 neous failures. Each stripe of a RAID-6 group can tolerate 253 any two failures. 254

In a replication cloud storage system, each data block has 255 a certain number of replicas, and the replicas are dispersed 256 over different nodes to improve the probability of blocks 257 available when multiple nodes fail concurrently. Replica- 258 tion systems have two storage properties: no two replicas of 259 a data block are stored on the same node; and replicas of a 260 data block must be found on at least two racks. 261



Fig. 1. Depicting the typical "bathtub curve" behavior of hard-drive failure rates, deduced in [34] from drives in the field.

#### 262 2.5 Weibull Distribution

In a real-world setting, drive failure rates typically follow a 263 "bathtub curve" with high failure rates at the beginning and 264 265 the end of a drive's life-cycle [32]. Fig. 1 depicts the failure rate for a hard drive's life-cycle [32], [33], [34]: after the 266 initial infant mortality, the failure rate enters in low-risk 267 state and starts to wear out after 5 to 7 years. Schroeder and 268 Gibson [20, Fig. 8] subsequently found the Weibull distribu-269 tion is a suitable fit for the empirical cumulative distribution 270 function of time between drive replacements observed in a 271 field-gathered dataset, while the exponential distribution 272 provides a poorer fit. 273

The time required for restoring data after an operational 274 failure (or scrubbing an entire drive) depends on the drive's 275 capacity and the data transfer rate of the drive (or the rate of 276 media scrubbing). To improve the reliability of cloud stor-277 age systems while maintaining a quality user service, the 278 data transfer rates for repairing (and scrubbing rates) are 279 adjusted according to the foreground activity, with high 280 rates when idle and low rates when busy [23]. This was veri-281 fied by Elerath and Schindler [23, Figs. 3, 4, 5, Tab. III], who 282 283 found the Weibull distribution is also a suitable match for time-to-failure, time-to-repair, and scrubbing-time distribu-284 tions derived from field data. 285

Therefore, we use the 2-parameter Weibull distribution to model occurrences of failures, rebuilds, and scrubbing in a cloud storage system. For parameters  $\alpha$  and  $\beta$  in the Weibull distribution, the probability density function f, cumulative density function F, hazard rate h, and cumulative hazard rate H are given as follows:

**293** 294

296

$$f(t) = \beta \frac{t^{\beta-1}}{\alpha^{\beta}} \exp(-(t/\alpha)^{\beta}), \qquad (1)$$

$$F(t) = 1 - \exp(-(t/\alpha)^{\mu}),$$
  

$$h(t) = \beta \frac{t^{\beta-1}}{\alpha^{\beta}}, \text{ and} \qquad (2)$$
  

$$H(t) = \frac{t^{\beta}}{\alpha^{\beta}}.$$

(See e.g., [35, Sec. 3.1.1].) The parameter  $\alpha$  is the *scale parameter* denoting the characteristic life, and  $\beta$  is the *shape parameter* controlling the shape of the distribution.

Weibull distributions are used to express various timedependent distributions with increasing, decreasing, or constant occurrence rates. If  $\beta > 1$ , the hazard rate *h* increases over time, i.e., the probability of an operational failure increases, simulating an aging cloud storage system. With 304  $0 < \beta < 1$ , we model a system with "infant mortality", and 305 with  $\beta = 1$  we have the traditional exponential distribution. 306 The parameter  $\alpha$  gives the characteristic life of drives. Youn- 307 ger drives may have a decreasing failure rate, while the 308 older drives may have increasing failure rates. 309

We compare models using the cumulative hazard rate, 310 i.e., the expected number of failure events from time 0. For a 311 proactive cloud storage system, however, we assume that 312 all drives that are classified as at risk of failure have sufficient *time in advance*, i.e., the time between warning and 314 actual failure. Thus, for proactive cloud storage systems, we 315 scale the cumulative hazard rate accordingly: 316

$$\hat{H}(t) := (1 - \text{FDR}) \frac{t^{\beta}}{\alpha^{\beta}}, \qquad (3)$$

where FDR is the *failure detection rate*. The function  $\hat{H}$  <sup>319</sup> instead gives the expected number of failure events from <sup>320</sup> time 0 which are not predicted in advance. <sup>321</sup>

A drive's *availability* is the proportion of time it can pro- 322 vide service to its users. The availability of a drive in the pres- 323 ence of operational failures, can be estimated using [23], [36]: 324

$$A_{\rm op}(t) := \frac{\alpha_p(t)}{\alpha_p(t) + \rm MTTR},$$
326

where MTTR denotes the *mean time to repair* a drive after an 327 operational failure, and  $\alpha_p$  is the drive *pseudo-characteristic* 328 *life*, modeled by 329

$$(t) = \frac{\alpha^{\beta}}{t^{\beta-1}},$$

331

348

where  $\alpha$  and  $\beta$  are the parameters of the Weibull distribution for operational failures, and *t* is time. 333

 $\alpha_p$ 

For proactive cloud storage systems, we adjust this to 334 account for operational failure prediction: 335

$$\hat{A}_{\rm OP}(t) := \frac{\alpha_p(t)}{\alpha_p(t) + (1 - {\rm FDR})\,{\rm MTTR}},\tag{4}$$

as only unpredicted operational failures need to be repaired. 338 Since the system has intra-system prediction, proactive fault 339 tolerance does not affect the pseudo-characteristic life of 340 drives. 341

We also use a 2-parameter Weibull distribution to model 342 occurrences of block defects, restorations, and scrubbing in 343 a cloud storage systems (block defects are not predicted in 344 advance). The steady-state availability of drives in the presence of block defects is modeled in [36] by 346

$$A_{\text{def}} := \frac{\text{MTTB}}{\text{MTTB} + \text{MTTS}},$$
(5)

where MTTB denotes the *mean time to block defect*, and 349 MTTS denotes the *mean time to scrubbing*, i.e., the mean time 350 between drive scrubbings. 351

For the Weibull distribution, Eq. (1), we have 352

$$MTTR = \alpha \Gamma(1 + 1/\beta), \qquad (6) _{354}$$

TABLE 2 Parameters of the Weibull Distributions [23]

		Drive A	Drive B	Drive C
		SATA	SATA	FC/SCSI
Operational failure	$lpha_f \ eta_f$	$302,016 \\ 1.13$	$4,833,522 \\ 0.576$	$1,058,364 \\ 0.721$
Latent block defect	$lpha_l \ eta_l$	$12,\!325$ 1	42,857 1	$50,\!254$ $1$
Rebuild time	$lpha_r \ eta_r$	22.7 1.65	20.25 1.15	6.75 1.4
Scrubbing time	$lpha_s \ eta_s$	186 1	160 0.97	124 2.1

<sup>355</sup> where  $\Gamma$  is the gamma function, so

357

386

393

$$\Gamma(1+1/\beta) = \int_0^\infty x^{1/\beta} e^{-x} \, dx$$

The values of MTTB and MTTS are also calculated as per Eq. (6), with appropriate parameters (listed in Table 2 for our experiments).

## 361 **3 MATHEMATICAL MODELS**

In this section, we derive equations for the expected number
of data loss events within a period of time, for RAID-5,
RAID-6, 2-way replication, and 3-way replication systems,
each with proactive fault tolerance. The equations are listed
in Table 1.

## 367 3.1 Proactive RAID Systems

If an operational failure or block defect occurs for a RAID
group, the group is said to be in *degraded mode*. Degraded
mode triggered by an operational failure initiates a rebuild
process to repair the failure, whereas degraded mode as a
result of a block defect will be resolved in the course of
scrubbing.

Consistent with [22], [23], we consider a sequence of fail-374 ures as negligible if (a) it requires two or more block defects, 375 as it is improbable that two block defects simultaneously 376 occur which affect the same data on two separate drives (with 377 the number of blocks per drive typically in the tens of thou-378 sands), and (b) if the only surviving copy of a block is lost 379 through a block defect, as the time between block defects 380 (usually thousands or even tens of thousands of hours) is 381 much larger than the rebuild time (typically tens of hours). 382

The expected number of data loss events in time period tin a stripe is thus modeled by:

$$\Pr\left(\begin{array}{c} \text{being at risk of data} \\ \text{loss by op. failure} \end{array}\right) \times \left(\begin{array}{c} \text{no. ways in which an} \\ \text{op. failure could result} \\ \text{in data loss} \end{array}\right) \times \hat{H}(t),$$

where  $\hat{H}(t)$  is the adjusted cumulative hazard rate given in Eq. (3).

Modifying [22, Eq. 5], we estimate the expected number of data loss events in a *g*-drive RAID-5 group within a time period t as

$$N_{\rm R5}(t) := (R_{\rm op} + R_{\rm def}) (g-1) \hat{H}(t), \tag{7}$$

where  $R_{\rm OP} = 1 - \hat{A}_{\rm OP}^g$  is the probability of the group having 394 at least one operational failure and  $R_{\rm def} = 1 - A_{\rm def}^g$  is the 395 probability of the group having at least one block defect. 396 Here,  $\hat{A}_{\rm OP}$  is the adjusted availability of a drive in the pres- 397 ence of operational failures given in Eq. (4), and  $\hat{A}_{\rm OP}^g$  is the 398 probability of all the *g* drives having no operational failure at 399 a given time. After a drive fails, the remaining g - 1 drives are 400 subject to operational failure, thereby giving Eq. (7). 401

Again modifying [23, Eq. 5], we estimate the expected 402 number of data loss events in a *g*-drive RAID-6 group 403 within a time period *t* as

$$N_{\rm R6}(t) := (R_{\rm op-op} + R_{\rm op-def}) (g-2) \hat{H}(t), \qquad (8)$$

where  $R_{\rm OP-OP}$  is the probability of the group having at least  $_{407}$  two operational failures, so  $_{408}$ 

$$R_{\text{op-op}} = \Pr(\text{at least 2 op. fails})$$
  
= 1 - Pr(no op. fail) - Pr(exactly 1 op. fail)  
= 1 -  $\hat{A}_{\text{op}}^{g} - g \, \hat{A}_{\text{op}}^{g-1} (1 - \hat{A}_{\text{op}})$ 

and  $R_{\rm op-def}$  is the probability of the group having at least 411 one operational failure and one block defect on two distinct 412 drives, so 413

$$R_{\rm op-def} = \Pr\left(\begin{array}{l} \text{at least 1 op. fail and 1 block} \\ \text{defect on two distinct drives} \end{array}\right)$$
$$\simeq 1 - \Pr(\text{no op. fail}) - \Pr(\text{no block defect}) \\+ \Pr(\text{no op. fail and no block defect}) \\= 1 - \hat{A}_{\rm op}^g - A_{\rm def}^g + (\hat{A}_{\rm op} A_{\rm def})^g.$$

For proactive RAID systems, we use the adjusted drive 417 availability  $\hat{A}_{\rm OP}$  and cumulative hazard rate  $\hat{H}(t)$  to account 418 for operational failure prediction in Eqs. (7) and (8), rather 419 than  $A_{\rm OP}$  and H(t) in [22, Eq. 5] and [23, Eq. 5], respectively. 420 When the failure detection rate FDR = 0, Eqs. (7) and (8) 421 become [22, Eq. 5] and [23, Eq. 5], respectively, as in reactive 422 RAID systems. 423

## 3.2 Proactive Replication Systems

In a replication system, we use r to denote the number of racks, 425 n to denote the number of nodes in each rack, d to denote the 426 number of drives in each node, and b to denote the number of 427 blocks on each drive. There are thus rnd drives in total.

After an operational failure or block defect, we say the 429 system enters *degraded mode* during the rebuild or scrubbing 430 process. Note, this definition of "degraded mode" is distinct 431 from that for RAID groups: for replication, we consider all 432 of the *rnd* drives in the system (whereas for RAID, we consider those within a group separately). As such, we no lon-434 ger consider *if* the system has block defects, but *how many* 435 block defects it has.

## 3.2.1 Proactive 2-Way Replication

In a 2-way replication system, every data block has two cop-438 ies, and the two copies must be stored on two separate 439 racks. We call such an arrangement a *replica pair*. Given a 440 drive, we can choose any of the (r - 1)nd drives on different racks to extend it to a replica pair. 442

424

r racks (*n* nodes per rack; *d* drives per node)

Fig. 2. 2-way replication. We have an operational failure for some drive D (i.e., degraded mode due to an operational failure). We lose data with probability  $P_{2\text{-way}}$  if an operational failure occurs on any drive on another rack to D.

Suppose both drives in a replica pair incur operational
failures, and drive *D* is one of them. The probability that a
data block *x* on *D* is also stored on the other failed drive is

$$p = \frac{1}{(r-1) nd},$$

assuming duplicate blocks are stored in random replica pairs. The probability that no block on *D* is lost (i.e., not stored on the other failed drive) is  $(1 - p)^b$ , so the probability of data loss caused by these two concurrent failures is

> $P_{2-\text{way}} = 1 - (1 - p)^{b},$ =  $1 - \left(1 - \frac{1}{(r-1)nd}\right)^{b}.$

453

459

465

447

Data loss can only occur when at least two racks simultaneously have corrupted data due to operational failures or
block defects. As with RAID systems, we consider the case
of two block defects erasing shared data to be negligible.

*Case (a):* Two operational failures.

Assuming the system is in degraded mode due to an operational failure of some drive on one rack, data loss occurs with probability  $P_{2-\text{Way}}$  when there is a second operational failure on one of the (r-1)nd drives on the other racks. This situation is illustrated in Fig. 2.

*Case (b):* One operational failure and one block defect.

466 'Given the system has a block defect for block x, data loss 467 will occur when an operational failure occurs for the unique 468 second drive containing x. By definition, a drive has a block 469 defect with probability  $1 - A_{def}$ , so the expected number of 470 drives with block defects is  $(1 - A_{def})rnd$ . This situation is 471 illustrated in Fig. 3.

472 Combining cases (a) and (b), the expected number of data473 loss events within a time period *t* is thus

$$N_{2-\text{way}}(t) = \left(P_{2-\text{way}}(r-1)ndD_{\text{op}} + rnd(1-A_{\text{def}})\right)\hat{H}(t),$$

476 where

475

478

$$D_{\rm op} = 1 - \hat{A}_{\rm op}^{rnd} \tag{11}$$

is the probability of being in degraded mode due to an oper-ational failure.



Fig. 3. 2-way replication. We have a block defect for block x. We will lose x if an operational failure occurs to the unique second drive containing x.

## 3.2.2 Proactive 3-Way Replication

In 3-way replication, each data block is stored 3 times, on 3 482 distinct nodes, out of which 2 nodes are on the same rack 483 and 1 node is on another rack. We call such an arrangement 484 a *replica set*.

Lemma 1. The number of replica sets is

(9)

(10)

$$r(r-1)\binom{n}{2}nd^3.$$
 488

481

486

490

508

The number of replica sets containing a given drive is

$$\frac{3}{2}(r-1)n(n-1)d^2.$$
492
493

The number of replica sets containing two given drives on 494 different nodes in the same rack is (r-1)nd. 495

The number of replica sets containing two given drives on 496 different racks is 2(n-1)d. 497

**Proof.** To count the number of replica sets, we pick one of 498 the *r* racks, from which we choose two nodes (in one of 499 n2 ways), and one of the r - 1 other racks, from which we 500 choose one of the *n* nodes, and from each of the three 501 nodes, we choose one drive (in one of  $d^3$  ways). This gives 502 the first claim in the lemma statement. 503

Given a drive D, it can belong to two types of replica 504 sets; the number of replica sets containing D and another 505 drive in the same rack as D (but different node to D) is 506

$$\underbrace{ \begin{array}{c} \text{choose drive on} \\ \text{same rack as } {}_{D} \\ \overbrace{(n-1)d}^{} \times \overbrace{(r-1)nd}^{} \end{array} }_{\text{choose drive on} \\ \overbrace{(r-1)nd}^{} \times \overbrace{(r-1)nd}^{} \end{array}$$

and there are  $(r-1)n2d^2$  replica sets containing D with 509 the other two drives in a different rack to D. Summing 510 these simplifies to give the second claim. 511

For the third claim, given two drives on different 512 nodes in the same rack, we can choose any of the 513 (r-1)nd drives on different racks to extend them to a 514 replica set. 515

For the fourth claim, given two drives on different 516 racks, we can choose any of the 2(n-1)d drives on dif- 517 ferent nodes in those two racks to extend them to a rep- 518 lica set.  $\Box$  519

520 Suppose all the three drives in a replica set incur operational failures, and drive D is one of them. By the second 521 claim in Lemma 1, the probability of a block on D is also 522 stored on the other two failed drives is 523

$$p_{\rm dup} = \frac{2}{3(r-1)\,n(n-1)d^2}$$

Then the probability [37] of at least one data block being 527 lost caused by the concurrent failures is 528

$$P_{3-\text{way}} = 1 - \left(1 - p_{\text{dup}}\right)^{b},$$
  
=  $1 - \left(1 - \frac{2}{3(r-1)n(n-1)d^{2}}\right)^{b}.$  (12)

There are four situations in which data loss occurs as the 532 result of an additional operational failure, Cases I to IV in 533 the following. 534

Case I: Two operational failures in one rack. 535

For a given *d*-drive node, the probability of no drive in 536 that node incurring an operational failure is  $A_{OD}^d$ . Conse-537 quently, for a given *n*-node rack, the probability of opera-538 tional failures occurring on at least two drives on different 539 nodes in that rack is 540

$$F_{\text{rack}} := 1 - \Pr(\text{no op. fails}) - \Pr\begin{pmatrix} \text{op. fails only} \\ \text{on one node} \end{pmatrix}$$
$$= 1 - (\hat{A}_{\text{op}}^d)^n - n(\hat{A}_{\text{op}}^d)^{n-1}(1 - \hat{A}_{\text{op}}^d)$$

and the probability  $D_{\text{Op-Op}}^{(1)}$  that at least one of the *r* racks 543 incurs operational failures on at least two drives on differ-544 ent nodes in that rack is given by 545

$$D_{\text{op-op}}^{(1)} = 1 - (1 - F_{\text{rack}})^n$$
  
= 1 -  $((\hat{A}_{\text{op}}^d)^n + n (\hat{A}_{\text{op}}^d)^{n-1} (1 - \hat{A}_{\text{op}}^d))^n$ 

547 548

553

558

559

542

525 526

530

531

$$= 1 - ((A_{\text{Op}}) + n(A_{\text{Op}}) - (1 - A_{\text{Op}})).$$
  
Given two operational failures occurring on two dr

549 on different nodes in the same rack, an operational failure 550 on any one of the (r-1)nd drives in the other racks causes 551 data loss with probability  $P_{3-way}$ . 552

Case II: Two operational failures on different racks.

554 For a given rack, the probability of at least one drive incurring an operational failure is  $1 - \hat{A}_{op}^{nd}$ . The probability 555 of at least two racks incurring an operational failure is thus 556

$$D_{\text{op-op}}^{(2)} = 1 - \Pr(\text{no op. fails}) - \Pr\left(\begin{array}{c}\text{exactly one rack}\\\text{with op. fails}\end{array}\right)$$
$$= 1 - \hat{A}_{\text{op}}^{rnd} - r \left(\hat{A}_{\text{op}}^{nd}\right)^{r-1} \left(1 - \hat{A}_{\text{op}}^{nd}\right).$$

Given two operational failures occurring in two different 560 racks, an operational failure on any one of the 2(n-1)d561 562 drives in those racks but in distinct nodes causes data loss with probability  $P_{3-wav}$ . 563

Case III: An operational failure and a block defect on the 564 same rack. 565

This case is illustrated in Fig. 4. Given an operational fail-566 ure, the expected number of drives on different nodes in the 567 same rack with block defects is 568

Pr(block def.) × (no. such drives) = 
$$(1 - A_{def})(n - 1)d$$
.

operational failure: drive D probability  $D_{op}$  of occurring probability 2/(3(n-1)d) of containing x block defect: block x expected number:  $(1 - A_{def})(n-1)d$ Fig. 4. Simultaneously, we have an operational failure for drive D and a

block defect for block x (on a drive on a different node in the same rack to D). If D happens to contain a copy of x, then x is lost if an operational failure occurs to the unique third drive containing x.

Given an operational failure for some drive D and a drive X 571 on a different node in the same rack with defective block  $x_{1}$  572

$$\Pr\left(\begin{array}{c} \text{a copy of } x\\ \text{occurs in } D\end{array}\right) = \frac{\text{no. replica sets containing } D \text{ and } X}{\text{no. replica sets containing } X}$$
$$= \frac{2}{3(n-1)d}$$

by Lemma 1. Thus, given an operational failure for some drive 575 *D*, the expected number of drives on a different node in the 576 same rack with defective block x which also occurs on D is 577

$$\frac{2}{3}(1-A_{\text{def}}).$$
 579

Finally, given an operational failure for some drive D and a 581 drive on a different node in the same rack with defective block 582 x which also occurs on D, there is a unique drive on another 583 rack which must incur an operational failure to lose block x. 584

Case IV: An operational failure and a block defect on the 585 different racks. 586

This case is illustrated in Fig. 5. Given an operational fail-587 ure, the expected number of drives in different racks with 588 block defects is 589

$$(1 - A_{def})(r - 1)nd.$$
 591

Given an operational failure for some drive D and a drive X 592 on a different rack with defective block *x*, 593

$$\Pr\left(\begin{array}{c} \text{a copy of } x\\ \text{occurs in } D\end{array}\right) = \frac{\text{no. replica sets containing } D \text{ and } X}{\text{no. replica sets containing}}$$
$$= \frac{4}{3(r-1)nd}$$

by Lemma 1. Thus, given an operational failure for some 596 drive D, the expected number of drives on a different rack 597 with defective block x which also occurs on D is

$$\frac{4}{3}(1 - A_{\text{def}}).$$
 600





Fig. 5. Simultaneously, we have an operational failure for drive D and a block defect for block x (on a drive on a different rack as D). If D happens to contain a copy of x, then x is lost if an operational failure occurs to the unique third drive containing x.

Finally, given an operational failure for some drive D and a drive on a different rack with defective block x which also occurs on D, there is a unique drive on another rack which must be erased to lose block x.

Combining Cases I-IV: The expected number of data loss events within a time period t is thus

$$N_{3-\text{way}}(t) \simeq \left(P_{3-\text{way}}\left((r-1)nd D_{\text{op-op}}^{(1)} + 2(n-1)d D_{\text{op-op}}^{(2)}\right) + 2 D_{\text{op}}(1-A_{\text{def}})\right) \hat{H}(t)$$
(13)

609 where  $D_{\text{OD}}$  is given by Eq. (11).

#### 610 4 EVALUATION

608

In this section, we compare the predictions of the equations
in Table 1 to simulations and give an analysis of the sensitivity of the equations.

#### 614 **4.1 Equation Verification**

#### 615 4.1.1 Experimental Setup

To test the accuracy of the equations in Table 1, we compare them against Monte-Carlo simulations. We design eventdriven Monte-Carlo simulations, in which there are six types of events that drive the virtual time forward:

- (a) Weibull distributed *operational failure events*, with parameters  $\alpha_f$  and  $\beta_f$ , potentially occurring on each drive;
- (b) Weibull distributed *block defect events*, with parameters  $\alpha_l$  and  $\beta_l$ , potentially occurring on each drive;
- (c) Weibull distributed *failure-rebuild complete events*, with  $\alpha_r$  and  $\beta_r$  denoting recovery time, occurring after an operational failure;
- (d) Weibull distributed *scrubbing complete events*, with  $\alpha_s$ and  $\beta_s$  denoting scrubbing time, periodically occurring on each drive to eliminate simulated block defects on it;

- (e) *warning events*, occurring 300 hours before opera- 632 tional failure for a proportion of drives determined 633 by the failure detection rate FDR; and 634
- (f) Weibull distributed *warning-rebuild complete events*, 635 with  $\alpha_r$  and  $\beta_r$  denoting pre-warning recovery time 636 (chosen to be the same as for failure-rebuild complete 637 events, for simplicity), occurring after a warning event. 638

Events (e) and (f) simulate proactive fault tolerance. Events 639 (c) and (f) trigger the introduction of a new drive (with future 640 failure and block defects). After (b) completes, a future block 641 defect event is added to the drive. When events of type (a) 642 and/or (b) simultaneously occur, under appropriate condi-643 tions (depending on the system) we incur a data loss event, 644 and new data is added to maintain system scale. 645

We enumerate the number of data loss events over a 646 5-year time period, and compare the results to the equations' predictions. A 5-year simulation is repeated until 648 the total number of data loss events is 10 or more, and 649 we average the results. 650

Our choice of time in advance TIA = 300 hours is motivated by [11], where the classification tree prediction model 652 predicted over 95 percent of failures with the TIA around 653 360 hours on a real-world dataset. For the four events (a)- 654 (d), we use the parameters in [23], listed in Table 2, from 655 three representative drive models, which had been in the 656 field for several years. 657

For RAID systems, we set the number of drives in each 658 group g = 15 for RAID-5 and g = 16 for RAID-6, so each 659 RAID group has 14 data drives. We choose 400 RAID 660 groups to model a deployment typical for a single RAID 661 system (5,600 data drives). For 2-way replication, we set the 662 number of racks r = 200, the number of nodes in each rack 663 n = 14, the number of drives in each node d = 4, and the 664 number of blocks in each drive  $b = 10^7$ , and for 3-way replication we set  $(r, n, d, b) = (300, 14, 4, 10^7)$ . In this way, all 666 four systems store the same amount of user data.

#### 4.1.2 Accuracy with Failure Prediction

Fig. 6 plots the expected number of data loss events pre- $_{669}$  dicted by the equations in Table 1 and enumerated by simu- $_{670}$  lation as the failure detection rate varies from 0 to 0.95. We  $_{671}$  include proactive RAID and replication systems for drives  $_{672}$  *A*, *B*, and *C*. For all the drive models, the equation-based  $_{673}$  results of every system closely match the simulation-based  $_{674}$  values; in the average case, they disagree by around  $_{675}$  10 percent (Elerath and Schindler [23] found the difference  $_{676}$  between equation and simulation of around 20 percent).

668

A reliability analysis is often used to assess trade-offs, to 678 compare schemes, and to estimate the effect of several param- 679 eters on storage system reliability. In this setting, a 10 percent 680 error would not significantly influence the analysis, especially 681 considering the Monte-Carlo simulations themselves are also 682 approximations. In a situation where Monte-Carlo simulations are required, the given equations could be used e.g., to 684 quickly reject inferior parameter combinations, after which 685 we can use Monte-Carlo simulations on the remaining cases. 686

Ordinarily, the equations and simulations agree closely,  $^{687}$  but for high FDR, they begin to disagree in some cases, e.g.,  $^{688}$  by a factor of 1.9 for the 3-way replication system with  $^{689}$  FDR = 0.95 on drive C. However, this occurs when the  $^{690}$  expected number of data loss events is low.  $^{691}$ 



Fig. 6. Number of data loss events predicted by the mathematical model and enumerated during simulation as the failure detection rate (FDR) varies, for drives A, B, and C, respectively.

#### 4.1.3 Accuracy as Time Varies 692

Fig. 7 plots the equational and simulated number of data 693 loss events over a t year period, as t varies from 1 year to 694 10 years. Here we use drive A and set FDR = 0.8. We find 695 that the equational results closely match the simulated 696 results. In the worst case, the equation and simulation dis-697 agree by around 30 percent. In the average case, they dis-698 agree by around 10 percent. 699

#### 4.1.4 Accuracy as System Scale Varies 700

Fig. 8 plots the equational and simulated number of data 701 loss events over a 5-year period, as the effective storage 702 space (i.e., the amount of user data, excluding redundancy) 703 varies. We vary the storage space via the number of groups 704 in RAID systems and the number of racks in replication sys-705 tems. Here we continue to use drive A and set FDR = 0.8. 706

We again find that the equational results closely match 707 the simulated results. In the worst case, the equation and 708 simulation disagree by around 30 percent. In the average 709 case, they disagree by around 10 percent. The experiment 710



Fig. 7. The expected number of data loss events over t years, for drive A and FDR = 0.8.

results verify the effectiveness of the reliability equations on 711 cloud storage systems with various scales. 712

Moreover, the equations yield comparable results much 713 faster than the simulations. On a standard PC desktop, with 714 e.g., MATLAB, we can quickly calculate the equational 715 results (in approximately 1ms), while the simulations usu-716 ally take between tens of seconds and tens of hours (even 717 hundreds of hours for a system with high FDR, where the 718 expected number of data loss events is very low) to produce 719 the results for a single set of inputs. 720

If the failure statistics for solid-state storage systems 721 could be obtained (i.e., the inputs in Table 1), the proposed 722 equations could be used for solid-state storage systems. 723

#### Sensitivity Analysis 4.2

In this section, we illustrate how the equations can be used 725 to analyze system sensitivity to varying system parameters. 726 We compare RAID-6 and 3-way replication, which are the 727 most common redundancy schemes. Unless otherwise 728 stated, we use drive A's parameters and set t = 5 years. 729

#### 4.2.1 Sensitivity to Drive Model

While both drives A and B are near-line SATA models, they 731 have the different failure distributions. In particular, drive 732 A has higher operational failure and block defect rates than 733 drive B. Fig. 9 plots the expected number of data loss events 734 for these two drives. 735



Fig. 8. The expected number of data loss events as the effective storage space varies, for drive A and FDR = 0.8. The storage capacity of a drive is denoted m.

724



Fig. 9. The expected number of data loss events of drives A and B, with RAID-6 or 3-way replication (rebuild parameters:  $\alpha_r = 24$ ,  $\beta_r = 2$ ; scrubbing parameters:  $\alpha_s = 240$ ,  $\beta_s = 1$ ).

In terms of data loss events, a designer could learn from 736 this which drive performs better (drive B, in this case) and 737 which storage scheme performs better (RAID-6, in this 738 case). Fig. 9 also shows the change in performance as FDR 739 changes, and the significant difference between proactive 740 and purely reactive fault tolerance. We see a stronger sensi-741 tivity to FDR in 3-way replication and RAID-6 systems, 742 than in 2-way replication and RAID-5 systems. 743

#### 744 4.2.2 Sensitivity to Weibull Parameters

The reliability of a system is affected by the rebuild time, as
concurrent operational failures are more likely to happen
with a longer rebuilding process.

Fig. 10 plots the expected number of data loss events for proactive RAID-6 and 3-way replication systems as the Weibull parameters vary. We vary the rebuild time via  $\alpha_r$ , the time for latent block defects via  $\alpha_l$ , and the scrubbing time via  $\alpha_s$ . All other parameters are those for drive A in Table 2.

We see that rebuild time plays a significant role in the 753 expected number of data loss events, and that 3-way repli-754 cation is more sensitive to the rebuild time than RAID-6. 755 756 This arises as RAID-6 is far more sensitive to block defects (and hence  $A_{def}$ ) than 3-way replication, which is apparent 757 from the model: for RAID-6, the expected number of data 758 loss events scales linearly with  $A_{def'}^{g}$  whereas for 3-way 759 replication, the expected number of data loss events scales 760 linearly with  $A_{def}$ , and with MTTB in the order of years 761 and MTTS in the order of hours, we have  $A_{def}$  close to 1. 762

Fig. 10 shows MTTB and MTTS has a negligible role for 3-way replication, but not for RAID-6. This is consistent with the models: in 3-way replication, two operational failures can result in the loss of a defective block x in only one way, but in RAID-6, two operational failures can result in the loss of block x in  $\binom{g-1}{2}$  ways.

#### 769 4.2.3 Sensitivity to System Scale Parameters

There are several other sensitivity properties that concern system designers and reliability engineers, such as the effects of system scale parameters on the overall cloud storage system reliability. In this subsection, we investigate the sensitivity of a system's reliability to some system scale parameters (including the number of data drives per RAID group, rack



Fig. 10. The expected number of data loss events for drive A as the mean time to rebuild, mean time to block defect, and mean time to scrubbing changes. We include proactive fault tolerance with FDR = 0.8 and reactive fault tolerance with FDR = 0.8

size, node size, and drive capacity). Unless otherwise stated, 776 we use the system parameters described in Section 4.1.1. 777

Fig. 11 plots the expected number of data loss events for 778 RAID-6 and 3-way replication systems, as the number of 779 data drives per group (i.e., g - 2), and the number of nodes 780 per rack (i.e., rack size n) varies. We see that the number of 781 data drives in RAID-6 group and the rack size have a signifi-782 cant impact on RAID-6 and 3-way replication systems reli-783 ability, respectively. 784

Fig. 12 plots the expected number of data loss events for 785 3-way replication systems as the number of drives per node 786 *d* changes, and as the number of blocks per drive *b* changes. 787



Fig. 11. The expected number of data loss events of drive A, as the number of data drives in RAID-6 group, and the rack sizes in 3-way replication system vary. We include proactive fault tolerance with FDR = 0.8 and reactive fault tolerance with FDR = 0.8



Fig. 12. The expected number of data loss events for 3-way replication systems with drive A as the node size (top) and drive size (bottom) change. We include proactive fault tolerance with FDR = 0.8 and reactive fault tolerance with FDR = 0.

We see that the number of drives per node (i.e., the node
size *d*) has a significant impact on 3-way replication systems
reliability. We see that the number of data loss events
increases as *b* increases.

As per the analysis in Section 3.2.2, provided each replica set shares at least one data block, data loss occurs when the three drives in a replica set simultaneously fail, and otherwise occurs with probability  $P_{3-Way}$  given in Eq. (12). Consequently, the larger the value of *b*, the higher the probability of a data loss event through simultaneous failures, but when the  $^{797}$  *b* is greater than a certain value (depending on system size),  $^{798}$  the probability is approximately 1, consistent with Eq. (12).  $^{799}$ 

#### 5 CONCLUSION AND FUTURE WORK

In this paper, we present four equations for proactive RAID 801 and replication cloud storage systems, by which one can 802 predict the overall reliability of systems in the presence of 803 operational drive failures and latent block defects. The 804 equations also incorporate time-variant failure rates and 805 media scrubbing processes. We use Weibull distributions 806 for failure rates, which is now considered more realistic 807 than the simpler exponential distribution. 808

We indicate the usefulness of these equations by investigating the impact of proactive fault tolerance and the system parameters on reliability. While simulations need specialized code and take much longer, they give comparable results to the equations. As such, the equations can help designers to readily explore the design space for their system: 814

- (a) To ensure availability, a designer desires to minimize 815 rebuild and warning migration bandwidth, but this 816 influences the rebuild time, which impacts success- 817 fully protecting at-risk data, which will negatively 818 affect the reliability against data loss. The equations 819 can aid the designer in optimizing this trade-off. 820
- (b) In general, a high failure detection rate incurs a high s21 false alarm rate (FAR), resulting in unnecessary s22 processing costs. The equations can help designers s23 choose a failure predictor with an FDR to achieve a s24 specific level of reliability, while minimizing FAR. 825
- (c) The reliability of systems is also significantly affected <sup>826</sup> by drive model (see Fig. 9). The equations can thus <sup>827</sup> help e.g., to decide whether to construct a system <sup>828</sup> with less reliable, but cheaper drives.
- (d) The intermediate results in the mathematical model 830 are meaningful, which may assist an operator pin- 831 point how data loss occurs in a storage system.
   832

There are some limitations to the mathematical models 833 we present: 834

- The models have a static mean rebuild time, which 835 may not accurately reflect fluctuation due to system 836 usage (which might vary according to the time of the 837 day) and system utilization (how much data is stored 838 on each drive).
- The assumptions break down when FDR approaches 840 1. With FDR = 1, there are no operational failures, 841 and the models will predict no data loss events. Real- 842 istically, even with FDR = 1, data loss events will 843 still occur in cases where there is insufficient time in 844 advance, and as a result of concurrent block defects. 845 We see this behavior in Fig. 6. 846

In practice, cloud storage systems at petabyte or exabyte 847 scales are both dynamic and heterogeneous, since new 848 drives will continuously enter the system as old ones leave 849 due to failure or age. Moreover, correlated failures (e.g., 850 node failures, or simultaneous block defects as a result of a 851 scratch) will also occur in system, which may influence the 852 system's reliability. Therefore, in future work, we plan to: 853 (a) modify the equations to suit other cloud storage systems, 854

such as deduplication and heterogeneous systems, to model
their reliability; and (b) extend the models to incorporate
correlated failures.

#### 858 **ACKNOWLEDGMENTS**

This work is partially supported by the NSF of China (grant number: 61702521), NSF of Tianjin (grant number: 4117JCYBJC15300), the Scientific Research Foundation of Civil Aviation University of China (grant number: 2017QD03S), and the Fundamental Research Funds for the Central Universities. Stones also supported by the Thousand Youth Talents Plan in Tianjin.

#### 866 **REFERENCES**

885

886

893

894

895

896

897

898

- 867 [1] I. Corderí, T. Schwarz, A. Amer, and D. D. E. Long, "Self-adjusting
  868 two-failure tolerant disk arrays," in *Proc. Petascale Data Storage*869 Workshop, 2012, pp. 1–5.
- [2] T. Schwarz, A. Amer, T. Kroeger, E. Miller, D. Long, and
  J. F. Paris, "RESAR: Reliable storage at exabyte scale," in *Proc. Model. Anal. Simul. Comput. Telecommun. Syst.*, 2016, pp. 211–220.
- [3] Q. Xin, E. L. Miller, S. J. Schwarz, and J. E. Thomas, "Evaluation of distributed recovery in large-scale storage systems," in *Proc. High Perform. Distrib. Comput.*, 2004, pp. 172–181.
- [4] S. Mitra, R. K. Panta, M.-R. Ra, and S. Bagchi, "Partial-parallel-repair (PPR): A distributed technique for repairing erasure coded storage," in *Proc. EuroSys*, 2016, pp. 1–14.
- [5] V. Gramoli, G. Jourjon, and O. Mehani, "Disaster-tolerant storage with SDN," in *Proc. Int. Conf. Networked Syst.*, 2015, pp. 278–292.
- [6] V. Estrada-Galiñanes, J. F. Pâris, and P. Felber, "Simple data entanglement layouts with high reliability," in *Proc. Perform. Comput. Commun. Conf.*, 2017, pp. 1–8.
  [7] J. F. Murray, G. F. Hughes, and K. Kreutz-Delgado, "Machine learn-
  - [7] J. F. Murray, G. F. Hughes, and K. Kreutz-Delgado, "Machine learning methods for predicting failures in hard drives: A multipleinstance application," J. Mach. Learn. Res., vol. 6, pp. 783–816, 2005.
- [8] G. F. Hughes, J. F. Murray, K. Kreutz-Delgado, and C. Elkan,
  "Improved disk-drive failure warnings," *IEEE Trans. Rel.*, vol. 51,
  no. 3, pp. 350–357, Sep. 2002.
- [9] J. F. Murray, G. F. Hughes, and K. KreutzDelgado, "Hard drive failure prediction using non-parametric statistical methods," in *Proc. Artif. Neural Netw.*, 2003, pp. 1–4.
  - [10] Y. Zhao, X. Liu, S. Gan, and W. Zheng, "Predicting disk failures with HMM- and HSMM-based approaches," in *Proc. Adv. Data Mining: Appl. Theoretical Aspects*, 2010, pp. 390–404.
  - [11] J. Li, X. Ji, Y. Jia, B. Zhu, G. Wang, Z. Li, and X. Liu, "Hard drive failure prediction using classification and regression trees," in *Proc. Dependable Syst. Netw.*, 2014, pp. 383–394.
- [12] A. Ma, F. Douglis, G. Lu, D. Sawyer, S. Chandra, and W. Hsu,
  "RAIDShield: Characterizing, monitoring, and proactively protecting against disk failures," in *Proc. USENIX File Storage Technol.*,
  2015, pp. 241–256.
- [13] S. Wu, H. Jiang, and B. Mao, "Proactive data migration for improved storage availability in large-scale data centers," *IEEE Trans. Comput.*, vol. 64, no. 8, pp. 2637–2651, Sep. 2015.
  [14] C. Xu, G. Wang, Z. Li, and X. Liu, "Health status and failure
- [14] C. Xu, G. Wang, Z. Li, and X. Liu, "Health status and failure prediction for hard drives with recurrent neural networks," *IEEE Trans. Comput.*, vol. 65, no. 11, pp. 3502–3508, Nov. 2016
- S. Pang, Y. Jia, R. J. Stones, X. Liu, and G. Wang, "A combined Bayesian network method for predicting drive failure times from SMART attributes," in *Proc. Int. Joint Conf. Neural Netw.*, 2016, pp. 4850–4856.
- [16] J. Li, R. J. Stones, G. Wang, Z. Li, X. Liu, and X. Kang, "Being accurate is not enough: New metrics for disk failure prediction," in *Proc. Symp. Reliable Distrib. Syst.*, 2016, pp. 71–80.
- in Proc. Symp. Reliable Distrib. Syst., 2016, pp. 71–80.
  [17] J. Li, R. J. Stones, G. Wang, X. Liu, Z. Li, and X. Ming, "Hard drive failure prediction using decision trees," *Rel. Eng. Syst. Safety*, vol. 164, pp. 55–65, 2017.
- [18] B. Eckart, X. Chen, X. He, and S. L. Scott, "Failure prediction models for proactive fault tolerance within storage systems," in *Proc. Model. Anal. Simul. Comput. Telecommun. Syst.*, 2008, pp. 1–8.
- J. Li, M. Li, G. Wang, X. Liu, Z. Li, and H. Tang, "Global reliability
   evaluation for cloud storage systems with proactive fault tolerance,"
   in Proc. Algorithms Archit. Parallel Process., 2015, pp. 189–203.

- B. Schroeder and G. A. Gibson, "Disk failures in the real world: 925 What does an MTTF of 1,000,000 hours mean to you?" in *Proc.* 926 USENIX File Storage Technol., 2007, pp. 1–16.
- [21] E. Pinheiro, W.-D. Weber, and L. A. Barroso, "Failure trends in a 928 large disk drive population," in *Proc. USENIX File Storage Technol.*, 929 2007, pp. 17–29. 930
- [22] J. G. Élerath, "A simple equation for estimating reliability of an 931 N+1 redundant array of independent disks (RAID)," in Proc. 932 Dependable Syst. Netw., 2009, pp. 484–493.
- [23] J. G. Elerath and J. Schindler, "Beyond MTTDL: A closed-form 934 RAID-6 reliability equation," ACM Trans. Storage, vol. 10, 935 no. 2, 2014, Art. no. 7.
- [24] A. Dholakia, E. Eleftheriou, X. Y. Hu, I. Iliadis, J. Menon, and 937
   K. K. Rao, "A new intra-disk redundancy scheme for high-938 reliability RAID storage systems in the presence of unrecover-939 able errors," ACM Trans. Storage, vol. 4, pp. 1–42, 2008.
- [25] K. M. Greenan, E. L. Miller, and J. J. Wylie, "Reliability of flat 941 XOR-based erasure codes on heterogeneous devices," in *Proc.* 942 *Dependable Syst. Netw.*, 2008, pp. 147–156.
  [26] G. A. Gibson, *Redundant Disk Arrays: Reliable, Parallel Secondary* 944
- [26] G. A. Gibson, Redundant Disk Arrays: Reliable, Parallel Secondary 944 Storage, vol. 368. Cambridge, MA, USA: MIT Press, 1992. 945
- [27] W. Dweik, M. A. Majeed, and M. Annavaram, "Warped-Shield: 946 Tolerating hard faults in GPGPUs," in *Proc. Dependable Syst.* 947 *Netw.*, 2014, pp. 431–442.
- [28] K. M. Greenan, J. S. Plank, and J. J. Wylie, "Mean time to meaningless: MTTDL, Markov models, and storage system reliability," 950 in Proc. USENIX Hot Topics Storage File Syst., 2010, pp. 1–5.
- [29] I. Iliadis and V. Venkatesan, "Rebuttal to "Beyond MTTDL: A 952 closed-form RAID-6 reliability equation"," ACM Trans. Storage, 953 vol. 11, no. 2, 2015, Art. no. 10. 954
- [30] B. Allen, "Monitoring hard disks with SMART," Linux J., vol. 1, 955 no. 117, pp. 74–77, 2004. 956
- [31] F. Mahdisoltani, I. Stefanovici, and B. Schroeder, "Proactive error 957 prediction to improve storage system reliability," in *Proc. USENIX* 958 *Annu. Tech. Conf.*, 2017, pp. 391–402. 959
- [32] J. G. Elerath, "Specifying reliability in the disk drive industry: No 960 more mtbf's," in *Proc. Rel. Maintainability Symp.*, 2000, pp. 194–199. 961
- [33] J. G. Elerath, "AFR: Problems of definition, calculation and 962 measurement in a commercial environment," in *Proc. Rel. Main-* 963 *tainability Symp.*, 2000, pp. 71–76.
- [34] J. Yang and F. B. Sun, "A comprehensive review of hard-965 disk drive reliability," in *Proc. Rel. Maintainability Symp.*, 1999, 966 pp. 403–409. 967
- [35] J. I. McCool, Using the Weibull Distribution: Reliability, Modeling and 968 Inference. Hoboken, NJ, USA: Wiley, 2012. 969
- [36] K. Trivedi, "Availability modeling," 2016. [Online]. Available: 970 http://amod.ee.duke.edu/, Accessed on: Oct. 5, 2016. 971
- [37] P. Li, J. Li, R. J. Stones, G. Wang, Z. Li, and X. Liu, "ProCode: A 972 proactive erasure coding scheme for cloud storage systems," in 973 Proc. IEEE 35th Symp. Reliable Distrib. Syst., 2016, pp. 219–228. 974



Jing Li received the BSc and MSc degrees in 975 computer science and technology from Shan-976 dong University, Jinan, China, in 2004 and 2007, 977 respectively and the PhD degree in computer sci-978 ence from Nankai University, Tianjin, China, in 979 2016. She is currently a teacher with the College 980 of Computer Science and Technology, Civil Avia-981 tion University of China. Her research interests 982 include mass data storage and machine learning. 983



**Peng Li** received the BSc degree in computer sci-984 ence and technology from Tianjin Normal Univer-985 sity, Tianjin, China, in 2012. Now he is working 986 toward the PhD degree in the College of Computer 987 and Control Engineering, Nankai University, Tian-988 jin, China. His research interests include cloud 989 storage, distributed systems and erasure coding. 990

#### LI ET AL.: RELIABILITY EQUATIONS FOR CLOUD STORAGE SYSTEMS WITH PROACTIVE FAULT TOLERANCE



**Rebecca J. Stones** received the PhD degree in pure mathematics from Monash University, in 2010. She now has diverse research interests, including combinatorics and graph theory, codes, search engines and data storage, phylogenetics, and quantitative psychology.



Zhongwei Li received the PhD degree in computer science and technology from Harbin Engineering University, Harbin, China, in 2006. He is1006neering University, Harbin, China, in 2006. He is1007currently an associate professor with the College1008of Software, Nankai University, Tianjin, China.1009His research interests include machine learning1010and mass data storage.1011



Gang Wang received the BSc, MSc, and PhD degrees in computer science from Nankai University, Tianjin, China, in 1996, 1999, and 2002, respectively. He is currently a professor with the College of Computer and Control Engineering, Nankai University. His research interests include storage systems and parallel computing. He is a member of the IEEE.



Xiaoguang Liu received the BSc, MSc, and PhD 1012 degrees in computer science from Nankai University, Tianjin, China, in 1996, 1999, and 2002, 1014 respectively. He is currently a professor in computer science with Nankai University, Tianjin, 1016 China. His research interests include parallel 1017 computing and storage system. 1018

▷ For more information on this or any other computing topic, 1019 please visit our Digital Library at www.computer.org/publications/dlib. 1020