

Performance of Checksums and CRCs over Real Data

Jonathan Stone, Stanford University*
Michael Greenwald, Stanford University*
Craig Partridge, BBN Technologies†
Jim Hughes, Network Systems Corporation

Abstract

Checksum and CRC algorithms have historically been studied under the assumption that the data fed to the algorithms was uniformly distributed. This paper examines the behavior of checksums and CRCs over real data from various UNIX file systems. We show that, when given real data in small to modest pieces (e.g., 48 bytes), all the checksum algorithms have skewed distributions. In one dramatic case, 0.01% of the check values appeared nearly 15% of the time. These results have implications for CRCs and checksums when applied to real data. They also can cause a spectacular failure rate for both the TCP and ones-complement Fletcher checksums when trying to detect certain types of packet splices. When measured over several large file-systems, the 16 bit TCP checksum performed about as well as a 10 bit CRC.

We show that for fragmentation-and-reassembly error models, the checksum contribution of each fragment are, in effect, coloured by the fragment's offset in the splice. This coloring explains the performance of Fletcher's sum on non-uniform data, and shows that placing checksum fields in a packet trailer is theoretically no worse than a header checksum field. In practice, TCP trailer sums outperform even Fletcher header sums.

1 Introduction

The behavior of checksum and cyclic redundancy check (CRC) algorithms have historically been studied under the assumption that the data fed to the algorithms was uniformly distributed. (See, for instance, the work on Fletcher's checksum [2] and the AAL5 CRC[12, 4]). If one assumes random data drawn from a uniform distribution one can show a number of nice error detection properties for various checksums and CRCs. But in the real world,

*Jonathan and Michael's work was supported, in part, by ARPA under Army Contract DABT63-91-K-0001.

†BBN Technologies is a division of GTE Corporation. Craig's work was supported, in part, by the U.S. Department of Defense.

communications data is rarely random. Much of the data is character data, which has distinct skewing towards certain values (for instance, the character 'e' in English). Binary data has similarly non-random distribution of values, such as a propensity to contain zeros.

This paper reports on experiments with running various checksums and CRCs over real data from UNIX file systems. We show that the highly non-uniform distribution of values and the strong local correlation in real data causes extremely irregular distributions of checksum and CRC values. In some tests, less than 0.01% of the possible checksum values occurred over 15% of the time. We particularly examine the effects of this phenomenon when applied to the Internet checksum used for IP, TCP, and UDP [9, 1] and compare it to two variations of Fletcher's checksum. We also report on an experiment with placing the standard TCP checksum in a packet trailer. A trailer checksum noticeably increases the checksum's effectiveness, and we prove why this is so.

2 CRCs vs. Checksums

Before examining the behavior of different algorithms, it is worth briefly discussing the CRC and checksum algorithms we used.

CRCs are based on polynomial arithmetic, base 2. CRC-32 [5] is a 32-bit polynomial with several useful error detection properties. It will detect all errors that span less than 32 contiguous bits within a packet and all 2-bit errors less than 2048 bits apart. It will also detect all cases where there are an odd number of errors. For other types of errors, if they occur in data which has uniformly distributed values, the chance of not detecting an error is 1 in 2^{32} .

The concept of a checksum is less well defined. For the purposes of data communication, the goal of a checksum algorithm is to balance the effectiveness at detecting errors against the cost of computing the check values. Furthermore, it is expected that a checksum will work in conjunction with other, stronger, data checks such as a CRC. For example, MAC layers are expected to use a CRC to check that data was not corrupted during transmission on the local media, and checksums are used by higher layers to ensure that data

was not corrupted in intermediate routers or by the sending or receiving host.

The fact that checksums are typically the secondary level of protection has often led to suggestions that checksums are superfluous. Hard won experience, however, has shown that checksums are necessary. Software errors (such as buffer mismanagement) and even hardware errors (such as network adapters with poor DMA hardware that sometimes fail to fully DMA data) are surprisingly common and checksums have been very useful in protecting against such errors.

The two most popular checksums are the Internet checksum used for IP, TCP, and UDP [9, 1] and Fletcher's checksum [2]. They represent different balances between performance cost and error detection.

The TCP checksum is a 16-bit ones-complement sum of the data. This sum will catch any burst error of 15 bits or less[8], and all 16-bit burst errors except for those which replace one 1's complement zero with another (i.e., 16 adjacent 1 bits replaced by 16 zero bits, or vice-versa). Over uniformly distributed data, it is expected to detect other types of errors at a rate proportional to 1 in 2^{16} . The checksum also has a major limitation: the sum of a set of 16-bit values is the same, regardless of the order in which the values appear. The checksum was chosen by the Internet community in the late 1970s after experimentation on the ARPANET suggested the checksum was good enough and could be implemented efficiently.

Fletcher's checksum is designed to be a more robust error detecting code. The checksum keeps two sums. One sum, A , is a running sum of the data in 8-bit chunks. The other sum, B , is a running sum of each byte multiplied by its position from the end of the packet. This multiplication incorporates positional information into the checksum to protect against movement or transposition of data within the packet. The two 8-bit sums are concatenated to generate a 16-bit checksum. Fletcher also defined a 32-bit version, where 16-bit sums are kept. The algorithm was defined for both ones and twos-complement arithmetic. The version used for the TP4 checksum and in this paper uses 8-bit chunks. When performed in twos-complement, this 16-bit checksum detects all single bit errors, a single error of less than 16 bits in length, and all double bit errors separated by 16 bits or less. Though TP4 uses only the twos-complement version, we investigated both ones- and twos-complement Fletcher sums.

Historically, the TCP checksum and Fletcher's checksum have been viewed as offering a sharp tradeoff between performance and error detection capabilities. The TCP checksum requires one or two additions per machine word of data (assuming the machine word is a multiple of 16 bits long), while Fletcher's sum requires two additions per byte (even if the computation is done in word-sized chunks). As a result, measurements have typically shown the TCP checksum to be two to four times faster [6, 11]. However,

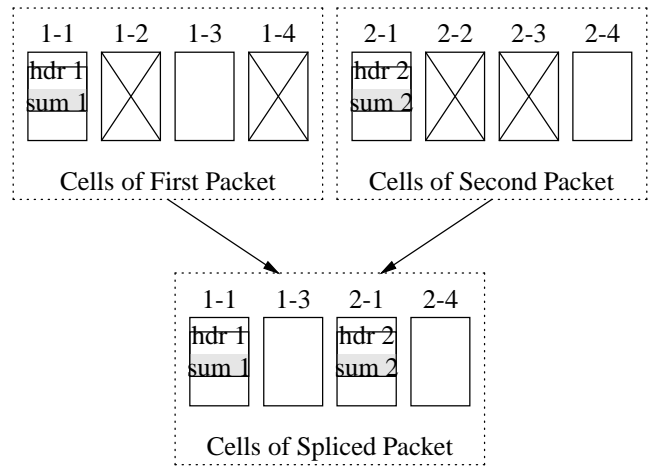


Figure 1: Example AAL5 Splice

that difference may be declining on newer processors, where the memory access time dominates any computational cost.

3 Work with AAL5

This study began as a study of the error scenarios for packet splices in Asynchronous Transfer Mode (ATM) Adaptation Layer 5 (AAL5). The AAL5 work helps motivate the rest of the paper and so is explained briefly here.

3.1 What is a Packet Splice?

AAL5 sends packets as a series of ATM cells, with the last cell specially marked using a bit in the ATM header. A *packet splice* occurs when the right number of cells are dropped such that pieces of two adjacent packets are combined so that they appear to represent one AAL5 packet. Figure 1 illustrates a splice: Two four-cell packets suffer a loss of four cells, such that the first and third cell of the first packet and the first and last cells of the second packet are spliced together to look like a single four-cell packet. It should be noted that ATM does not allow cells to be re-ordered, thus the number of possible splices is limited to those that merely drop, and do not reorder, cells.

Several conditions must be met for a splice to be valid. First, AAL5 stores the length of the packet in the last cell, so the size of the splice must be consistent with the AAL5 length in the last cell. Second, because AAL5 specially marks the last cell of every packet, the last cell of the first packet cannot be part of the splice. Third, the first 40 bytes of the first cell must be a valid TCP/IP header (i.e., have a length consistent with the packet length and certain bits must be set). Unless all three of these requirements are met, the splice will be easily detected without checking the CRC or checksum.

If the three requirements are met, then the splice has to be detected by either the AAL5 CRC (CRC-32) or the higher

layer protocol's checksum (such as the TCP or Fletcher's checksum).

In 1993, an informal study by Bill Marshall and Chuck Kalmanek at AT&T Bell Labs simulated file transfers from a UNIX filesystem (using real data from the filesystem) and examined the performance of the AAL5 CRC. They found a surprising number of cases where the packet splice passed the AAL5 CRC, leading them to wonder if the AAL5 CRC was strong enough. With Marshall's and Kalmanek's assistance, the authors set out to do a more complete set of tests. Those results were reported in an earlier version of this paper, presented at SIGCOMM '95 [7]. Some open questions and surprising results led us to perform a new and more comprehensive series of tests to resolve these issues.

3.2 Testing Splices

Our test program simulated a file transfer with the File Transfer Protocol (FTP) of all files on a file system (or selected directories of a file system) via TCP/IP using AAL5 over ATM. All IP and TCP header fields were filled in as if the file transfer were being done over the loopback interface (127.0.0.1). For each packet, the TCP sequence number was incremented by the data length, and the IP ID was incremented by one. The program then examined all possible splices of two adjacent TCP segments and checked to see if either the TCP checksum or AAL5 CRC failed to detect the splice. The program did not concern itself with splices whose data exactly matched a valid packet, nor with those splices that were detected by IP, TCP, or AAL5 header/trailer checks.

The test program was run over file systems at Network Systems Corporation (NSC), the Swedish Institute of Computer Science (SICS), and Stanford University. The TCP segment sizes examined were 256 bytes long, except for runt packets at the end of files. The first row in Tables 1 through 3 counts the total number of splices inspected. The next row counts how many invalid splices were detected by simple header checks, and so did not need to check the checksum. The row labeled "Identical data" records how many splices resulted in packets that were identical to one of the original packets, and hence would not result in corrupted data (the checksum, of course, was identical). The "Remaining" packets were all incorrect and depended on the checksum and the CRC to detect the corruption. All percentages listed are computed as percent of "Remaining splices". The rows following "Remaining" list the splices missed by the CRC test and the TCP checksum test. There were no splices missed by both CRC and the TCP checksum. The data from each site are broken down by file system. The total number of splices is greater than 2^{32} .

We would expect that the CRC of a splice would match the CRC of the original AAL5 packet at a rate of 1 in 2^{32} (or 0.000000232% of the time). Similarly, we would expect that the TCP checksum would fail to catch bad splices at a

Table 1: CRC and TCP Checksum Results
(256 Byte packets on systems at NSC)

system	code	% remaining	splices
<i>nsc05</i>	Total		7186841747
46411 files	Caught by Header		3593444113
4856193 pkts (98-05-04)	Identical data		17498067
	Remaining splices		3575899567
	Missed by CRC	0.0000000000	0
	Missed by TCP	0.0459554853	1643322
<i>nsc11</i>	Total		6306945748
45627 files	Caught by Header		3152782063
6896637 pkts (98-05-04)	Identical data		22324135
	Remaining splices		3131839550
	Missed by CRC	0.0000000319	1
	Missed by TCP	0.0610412816	1911715
<i>nsc23</i>	Total		4920441461
29444 files	Caught by Header		2459789331
4372688 pkts (98-05-04)	Identical data		50703652
	Remaining splices		2409948478
	Missed by CRC	0.0000000830	2
	Missed by TCP	0.0568444518	1369922
<i>nsc25</i>	Total		8748322301
38187 files	Caught by Header		4372322214
9531889 pkts (98-05-04)	Identical data		65900443
	Remaining splices		4310099644
	Missed by CRC	0.0000000464	2
	Missed by TCP	0.1103037608	4754202
<i>nsc27</i>	Total		5012189213
22319 files	Caught by Header		2505005350
5461908 pkts (98-05-04)	Identical data		16574413
	Remaining splices		2490609450
	Missed by CRC	0.0000000402	1
	Missed by TCP	0.0439271199	1094053
<i>nsc29</i>	Total		5756622285
57299 files	Caught by Header		2878637775
6314509 pkts (98-05-04)	Identical data		19999951
	Remaining splices		2857984559
	Missed by CRC	0.0000000350	1
	Missed by TCP	0.0552609704	1579350
<i>nsc49</i>	Total		5696462431
17663 files	Caught by Header		2846361632
6196298 pkts (98-05-04)	Identical data		16371605
	Remaining splices		2833729194
	Missed by CRC	0.0000000000	0
	Missed by TCP	0.0766246826	2171336
<i>nsc51</i>	Total		4584391161
16864 files	Caught by Header		2290882985
4990431 pkts (98-05-04)	Identical data		14136325
	Remaining splices		2279371851
	Missed by CRC	0.0000000000	0
	Missed by TCP	0.0693654262	1581096
<i>nsc52</i>	Total		8309068498
58132 files	Caught by Header		4153260212
9082777 pkts (98-05-04)	Identical data		40561081
	Remaining splices		4115247205
	Missed by CRC	0.0000000000	0
	Missed by TCP	0.1726656418	7105618

Table 2: CRC and TCP Checksum Results
(256 Byte packets on systems at SICS)

system	code	% remaining	splices
<i>sics.se</i>	Total		3183838883
/src1 48,817 files 3,520,967 pkts (11-24-97)	Caught by Header		1594737950
	Identical data		11000914
	Remaining splices		1578100019
	CRC	0.0000000000	0
	TCP	0.0411719151	649734
<i>sics.se</i>	Total		2902904306
/src2 11,492 files 3,162,423 pkts (11-24-97)	Caught by Header		1450715240
	Identical data		12039586
	Remaining splices		1440149480
	CRC	0.0000000000	0
	TCP	0.0344980161	496823
<i>sics.se</i>	Total		12074080447
/src3 7,845 files 13,097,058 pkts (12-17-97)	Caught by Header		6031140841
	Identical data		12062020
	Remaining splices		6030877586
	CRC	0.0000000000	0
	TCP	0.0088341538	532777
<i>sics.se</i>	Total		5025946678
/src4 33,912 files 5,496,043 pkts (12-17-97)	Caught by Header		2512845921
	Identical data		22171407
	Remaining splices		2490929350
	CRC	0.0000000000	0
	TCP	0.0198888017	495416
<i>sics.se</i>	Total		21107489268
/iss1 204,601 files 23,178,376 pkts (12-17-97)	Caught by Header		10557354562
	Identical data		126239615
	Remaining splices		10423895091
	CRC	0.0000000192	2
	TCP	0.2238580377	23334727
<i>sics.se</i>	Total		6560349785
/opt 141,453 files 7,312,235 pkts (11-24-97) 0.2% executables	Caught by Header		3286741967
	Identical data		152672075
	Remaining splices		3120935743
	CRC	0.0000000320	1
	TCP	0.1703438788	5316323
<i>sics.se</i>	Total		8630623470
/solaris 98,211 files 9,502,013 pkts (12-17-97)	Caught by Header		4318348898
	Identical data		92736322
	Remaining splices		4219538250
	CRC	0.0000000474	2
	TCP	0.1068534691	4508723
<i>sics.se</i>	Total		33661656216
/cna 248,611 files 36,859,417 pkts (12-17-97)	Caught by Header		16832727499
	Identical data		196026754
	Remaining splices		16632901963
	CRC	0.0000000180	3
	TCP	0.1866982627	31053339

Table 3: CRC and TCP Checksum Results
(256 Byte packets on systems at Stanford)

system	code	% remaining	splices
<i>smeg.stanford.edu</i>	Total		8863295657
/u1 198,352 files 9,901,213 pkts (8-20-97)	Caught by Header		4442709123
	Identical data		25715994
	Remaining splices		4394870540
	CRC	0.0000000228	1
	TCP	0.0707199443	3108050
<i>pompano.stanford.edu</i>	Total		1197495954
/usr/local 11,468 files 1,314,390 pkts (11-26-97)	Caught by Header		599005787
	Identical data		6024593
	Remaining splices		592465574
	CRC	0.0000000000	0
	TCP	0.0269563342	159707

rate of 1 in 2^{16} (or 0.001526% of the time). Observe that for the CRC, the CRC must match the CRC of the second AAL5 packet, while for TCP, the checksum over the entire splice must equal zero.

The tables show that for real data, the CRC failure rate is almost perfectly consistent with the expected failure rate for random data, and is therefore not the subject of much further investigation in this paper¹. For TCP, however, the story is different. Between 0.008% and 0.22% of the bad splices passed by the header checks passed the checksum. This is between a factor of 10 and 100 worse than expected, and requires some explanation.

4 Explaining The TCP Checksum Failures

Why does the TCP checksum fail to detect so many splices? The reasons have to do with the distribution of data values and how data from one packet can be mixed with data from another packet.

4.1 Failure Scenarios

We can compute the TCP checksum in pieces and then add the pieces to get the complete packet sum. So, we can think of the TCP checksum of a packet broken into ATM cells as being the sum of the individual checksums of each 48-byte cell.

The usual requirement for a splice to pass the TCP checksum is that the checksum of the splice add up to the checksum of the entire first packet contributing to the splice. Because the splice contains cells of the first and second packets, this requirement can also be expressed as a requirement that the checksum of the cells from the first

¹ The difference between our results and those of Marshall and Kalmanek are the "Identical Data" entries. Given that the payloads were identical, it is not a failure if the CRC doesn't detect these splices as no data-corruption occurs. Their tests did not distinguish the cases of splices with identical data from splices with different data but congruent checksums.

packet not included in the splice must equal the checksum from the cells of the second packet that are included in the splice. If just one cell from the second packet is included in the splice, this requirement reduces to the requirement that the checksum of the cell from the second packet have the same sum as the cell it replaces. In multicell replacements, the sum of the mixes of cells must be equal.

4.2 Distributions of the TCP Checksum

Given random data, a good checksum or CRC should uniformly scatter the checksum values over the entire checksum space. Obviously a checksum algorithm that does not uniformly distribute checksum values (i.e., has hotspots) will be more likely to have multiple cells with the same checksum. Theorem 6 in Appendix A proves that, over uniformly distributed data, the TCP checksum algorithm gives a uniform distribution of checksum values². Thus, any hotspots in the distribution of checksum values are due to non-uniformity of the data, and are not inherent in the TCP checksum algorithm.

4.3 The Distribution of Checksum Values over Single Cells

If the distribution of 16-bit words is completely uniform, the chance of an arbitrary sequence of data in the first packet having the same checksum as an equal-sized arbitrary sequence of data in the second packet is $1/M$, where $M = 2^{16}$. However, the distribution of values over real data is not uniform.

Figure 2 shows three plots summarizing the distribution of checksum values on the filesystem `/u1` on `smeg.dsg.stanford.edu`. The x-axis represents different checksum values, sorted by frequency to better show the distribution. In the PDF graphs, the y-axis is the probability that the given checksum value occurred. Fig. 2(a) shows the entire PDF, and (b) shows a blowup of the most frequent 65 values (0.1%) The CDF (fig. 2c) shows the same 65 values, but here the y-value for a given x represents the cumulative probability that any of the most common x values occurred. If the distribution were uniform then the PDF should simply be a horizontal line at $1/M$, and the CDF a straight line with slope $1/M$.

This data shows that the TCP checksum on real data has hotspots. In the file system shown in the figure (`smeg:/u1`), the top 0.1% of the checksum values occurred 2.5% of the time. If one examines this distributional data over many filesystems, one discovers two things. First, that the single most common checksum value (usually zero) occurs between 0.01% and 1% of the time. Second, that for 48 byte cells the 65 next most frequently occurring checksum values

²The actual requirements are weaker: as long as the values of even one word in the packet is uniformly distributed over all 2^{16} possible values, then the checksum of the entire packet is uniformly distributed over all possible values.

(0.1% of the checksum space) account for between 1% and 5% of the checksum values seen.

4.4 Checksum Distribution over Larger Blocks of Data

Although for uniformly distributed data values the probability distribution of the checksum is uniform independent of the length of the block of data, this is not true for non-uniform data. In that case, the expected probability distribution of the checksum may be computed by

$$P_k[i] = \sum_{j=0}^M (P_{k-1}[j]P_1[i-j])$$

where $P_k[i]$ is the probability that the checksum over a block of length k is equal to i , and where $i - j$ is taken mod M . The dotted line in Figure 2 labeled "Predict $k = 2$ " shows the expected distribution of checksums over blocks 2 cells long, given the checksum distribution over one cell given by $k = 1$.

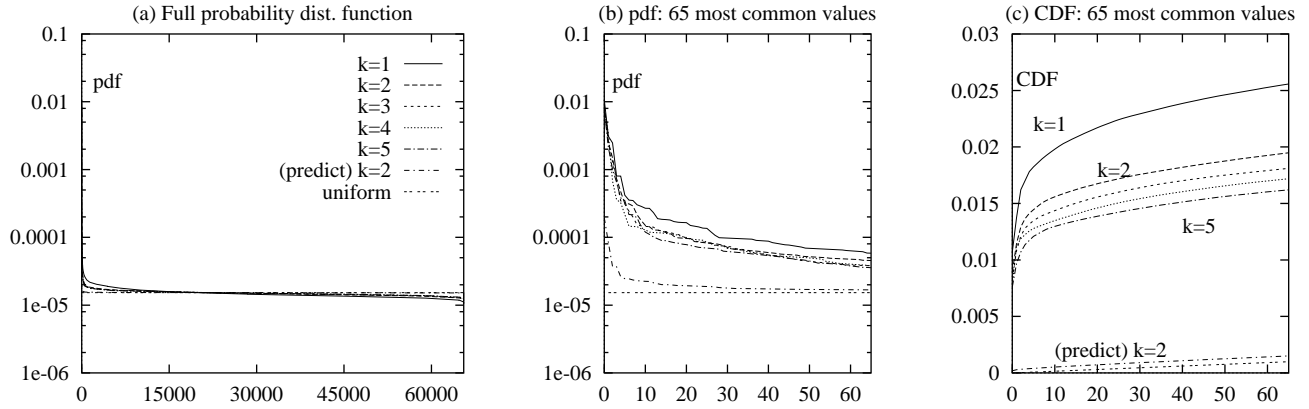
So, if the non-uniformity is uniform – that is, that every cell of data is drawn from the same probability distribution, and that the sum is the sum of *independent* samples – then we would expect the distribution of the sums to conform closely to the dotted line in our graphs. The predicted value for $k = 2$ is already close to uniform for all but the 20 most common values, even though $k = 1$ is decidedly non-uniform. Corollary 3 and Theorem 4 in the appendix show that, regardless of the original distribution, the distribution should get more uniform as k increases.

However, our measurements show that the non-uniformity extends to larger chunks than single words or cells, and that the checksum of one cell is correlated with the checksums of the neighboring cells. The lines labeled $k = 2$, $k = 3$, etc. show the measured distribution of checksums over samples of blocks of length k cells over real data in our file system. The data does get more uniform (seen most clearly in the CDF), but nowhere as quickly as it should if the cells were roughly independent. We believe the samples should be somewhat representative even of non-contiguous blocks. Once again, the checksum values are sorted in decreasing order of probability, to give a clearer picture of the distribution. Note that even over the larger block sizes, although the probability of a match decreases slightly, the distribution is still significantly more non-uniform than expected.

4.5 Filesystem-level Non-Uniformity Is Not The Answer

Given the non-uniform distribution, what, then, is the expected failure rate of the IP/TCP checksum in detecting splices for a given distribution, P , of checksum values? As discussed above, it is simply the probability that the checksum over the cells missing from the first packet is

Figure 2: Distribution of TCP checksum over blocks of k cells in *smeg.dsg.stanford.edu:/u1*.



equal to the checksum over the cells present from the second packet. For a given probability distribution P this probability is

$$P(\text{failure}) = \sum_{i=0}^M P[i]^2$$

Table 4: Probability (as %) of checksum match for substitutions of length k cells.

Length	Uniform	Predicted	Measured
1	0.001526	0.02126770	0.02126770
2	0.001526	0.00153019	0.01494399
3	0.001526	0.00152590	0.01348366
4	0.001526	0.00152590	0.01416288
5	0.001526	0.00152590	0.01108446

Table 4 computes this probability using the measurements of the Stanford file system from Figure 2. It lists the probability that the checksums of two blocks, each k cells long, drawn from anywhere in the same filesystem, will be equal. For each block of length k cells, the first column shows the expected probability given uniform distribution. The second column shows the probability predicted assuming each cell is drawn from the identical, non-uniform distribution. (The particular distribution is the one actually measured for single cells over the *smeg:/u1* file system.) This corresponds to the predicted distribution depicted by the dotted line in Figure 2. The last column lists the probability actually measured for each block size over the entire file system. We can see that even the milder non-uniformity of packet-sized chunks noticeably affects the probability of checksum failure, and that the failure probability does not tail off with larger block sizes as it should if each cell were independent.

Clearly, there is clustering and non-uniformity at a scale larger than single cells. Yet even aggregating the data over chunks of 1, 2, ... and 5 cells is not sufficient to

accurately predict the actual non-uniformity and failure rate, which is still more than 10 times higher than this simple model predicts. There are two issues our initial computation ignores. First, we have measured the probability distribution over the entire file system for chunks of k cells, but we know that distribution of data values is heavily dependent on file type (binary vs. character, executable vs. GIF, even Shakespeare vs. Joyce). Splices come from adjacent packets, which usually come from the same file. Thus real failure rates could be higher than the averaged global distribution would suggest.

For example, consider an extreme (and extremely hypothetical) case in which a file system consists of half binary and half textual data. Imagine that 90% of the cell-sized chunks of binary data had a checksum of 0x0000, and that 90% of the cell-sized chunks of textual data had a checksum of 0x1F00. Considered globally, we'd find 0x0000 45% of the time and 0x1F00 45% of the time, so $\sum p^2$ would be approximately 32% and we'd predict about 32% of the packet splices would incorrectly pass the checksum. However, in reality, for any given file the local distribution would find the most common checksum 90% of the time, and thus the failure rate would be about 81%. Therefore, the global distribution of checksums (measured across an entire filesystem) is not sufficient to accurately predict checksum failure rate: a more localized distribution of checksums is needed.

Even this is not the whole story. If two cells have congruent checksums because the data was identical, then replacement of one cell by the other is not a checksum failure – the packet is unaltered and no corruption will occur. To accurately predict meaningful checksum failures, then, we need to subtract "both congruent and equal" cells from the probability of a match. In a system with uniformly distributed data the odds of finding two 48 byte cells with identical data is 1 in 2^{384} , which is so unlikely as to be utterly negligible. However, in practice it occurs far more frequently. Our actual measurements show that the *most* common reason for checksum congruence is identical data

– identical splices occur 20 to 40 times more frequently than congruent-but-unequal splices. This is another example of non-uniform distribution of the data, but, in this case, a benign one.

4.6 Localized Non-Uniformity of Data

Table 5 shows how the probability changes when we restrict the comparisons to only look at local data. The first column (identical to the last column in Table 4) displays the probability of taking two blocks of data, each k cells long, from anywhere in the entire file-system, and finding that their IP checksums were congruent to each other. The column labeled "Locally congruent" shows the same probability if we limit the search to be within 2 packet lengths (512 bytes). (In order to increase the sample size for the local comparisons, we did not restrict ourselves to contiguous blocks). The final column shows how the probability decreases when we exclude checksum matches for a pair of blocks that contained identical data, as such a substitution would not result in any data corruption. It is still significantly higher than the global rate. (Recall, that if the data were uniformly distributed then every entry in this table should be 0.001526%). If the checksum failures are purely a result of non-uniform distribution, then these sample probabilities should track the measured TCP checksum failure rates.

Table 5: Probability (as %) of checksum match for substitutions of length k cells based on local data.

Length (k)	Globally Congruent	Locally Congruent	Excluding Identical
1	0.02126770	1.58305972	0.20704272
2	0.01494399	1.30267681	0.17226800
3	0.01348366	1.21236431	0.16614066
4	0.01416288	1.15970577	0.16316988

Table 6 compares this distribution data for several file systems with the actual rate of checksum failures for comparable-length substitutions. It is important to note that the sample data only deals in full-size cells, while the measured data deals in 8 byte trailers, too. Thus the byte-length for the sample data is simply $48k$, while the byte-length for the actual data is $48k - 40$. While the exact results vary for each system, there are three things all share. First, they are all in sharp contrast to the expected rate of 0.001526%. Second, the local non-uniformity is significantly worse than the global non-uniformity, and extends over packet-size blocks. Third, the distribution samples correspond roughly to the actual failure rate. But the correspondence is only rough. A small part of this is explained by the difference in byte-length (mainly for $k = 1$). Since "Actual" decreases non-linearly, we have not yet fully explained what is going on. Section 5.4 will return to this and explain the remaining discrepancy. To convert these probabilities to a total failure rate depends

Table 6: Checksum failures on real data Probability (as %) of checksum congruence for blocks of length k cells

<i>smeg.dsg.stanford.edu/u1</i>				
Predicted	0.0212677	0.0015302	0.0015259	0.0015259
Measured Global	0.0212677	0.0149440	0.0134837	0.0141629
Local Congruence	1.5830597	1.3026768	1.2123643	1.1597058
Exclude Identical	0.2070427	0.1722680	0.1661407	0.1631699
Actual	0.1026797	0.1581733	0.0907984	0.0568881
<i>sics.se:/opt</i>				
Predicted	1.1422436	0.0150023	0.0016907	0.0015280
Measured Global	1.1422436	0.9493377	0.8852883	0.8291802
Local Congruence	10.7766645	9.6723695	9.3490614	9.0170788
Exclude Identical	0.3872216	0.4732675	0.6897936	0.6086173
Actual	0.1085216	0.5551069	0.2130342	0.1183174
<i>sics.se:/src1</i>				
Predicted	0.0320218	0.0015407	0.0015259	0.0015259
Measured Global	0.0320218	0.0182235	0.0163506	0.0169735
Local Congruence	1.7774385	1.5562402	1.4326226	1.4595770
Exclude Identical	0.2537653	0.1938629	0.1418823	0.2196530
Actual	0.1037989	0.1143002	0.0499086	0.0283392
<i>sics.se:/src2</i>				
Predicted	0.0204720	0.0015312	0.0015259	0.0015259
Measured Global	0.0204720	0.0154869	0.0146764	0.0140808
Local Congruence	2.1467761	1.9016190	1.8495718	1.7812772
Exclude Identical	0.1094778	0.0995097	0.1345251	0.1108649
Actual	0.1747045	0.1339748	0.0429196	0.0210642

on the likelihood of substitutions of each given length. The odds that a substitution of a given length occurs depends on the type of errors one expects. In our simulation, we can exactly characterize the probability of a k -cell substitution for n -cell packets. (Our typical packets of 256 bytes contain 7 cells). For a splice to be valid the trailer cell of the first packet must be dropped and the trailer cell of the second packet must be kept. Further, to pass the header checks, the header cell of the first packet is usually kept. Therefore, we have $2n - 3$ cells to choose from, with the leading and trailing cells already specified. So there are

$$\binom{2n-3}{n-2}$$

total splices (462 for 7 cell packets) that might pass the header checks. There are

$$\binom{n-2}{k-1} \binom{n-1}{k-1}$$

splices of length k . Our simulation treats every possible substitution as equally likely. This clearly might not be true in all situations. However, the breakdown by substitution length in Table 5 is enough to show that the failure rate will be worse than expected for *all* substitutions, regardless of length. The only question is exactly how bad.

5 Reducing Checksum Failures

In this section we look at various ways to reduce the checksum failure rates.

5.1 Regaining a uniform distribution: compression

We claim that the TCP checksum’s failure to detect many splices is due to the non-uniform distribution of the data being summed. One obvious way to deal with non-uniform data patterns is to compress the data. As an experiment to verify that our diagnosis was correct, we compressed all the files in the file system at SICS that gave the TCP checksum the most trouble (`/opt` on `fafner.sics.se`) and ran our tests on the compressed files. (The compression was Lempel-Ziv, and was performed using the UNIX `compress` command.) The results are shown in Table 7. The interesting result is that the number of splices that passed the checksum is approximately 0.0021%, which is close to the expected rate on uniform data of 0.0015%. This result is a hundred-fold improvement over the .17% miss-rate before compression. So compression clearly helps.

Table 7: CRC and TCP Checksum Results, Compressed Data

(256 Byte packets on systems at SICS)

system	code	% remaining	splices
<i>fafner.sics.se</i>	Total		1549869756
compressed <code>/opt</code> 1,679,166 pkts (5-9-95)	Header		773945117
	Identical		51902
	Remaining		775872737
	CRC	0.0000000000	0
	TCP	0.0021002156	16295

5.2 Alternative checksums: Fletcher

It is not always possible or desirable to compress the data. Another obvious question to ask is whether, without data compression, another checksum algorithm would perform better than TCP’s. An obvious candidate checksum is Fletcher’s checksum[13].³ With our error model, where cells are dropped but no random data is inserted, we might expect the positional B term to improve error detection.

As with TCP, we can compute and analyze Fletcher’s checksum over individual cells rather than entire packets. Recall that the B term of the Fletcher checksum is computed by multiplying each byte by its offset from the end of the packet. We can also compute a local Fletcher checksum over one cell i as A_i , and B_i . To compute the contribution of an individual cell to the total Fletcher sum for the packet, we add A_i to A_{packet} and add R_i to B_{packet} , where $R_i =$

³Unlike [13], our Fletcher’s results perform a sum-to-zero inversion on the transmitted checksum. See Sec. 6.3.

$(B_i + A_i L)$ and L is the offset of the end of the cell from the end of the packet. It should be noted that since all the shifts of data are by a multiple of the cell size (48 bytes), the contribution of the B term for each cell to detect motion is limited to 1 from, at most, $M/GCD(M, 48)$ values (85 and 16 for 1 and 2’s complement, respectively). Both 85 and 16 are considerably smaller than M (255 or 256, respectively).

Table 8 shows the actual results for both 1’s complement (mod 255) and 2’s complement (mod 256) Fletcher’s checksum over several filesystems. The results of the TCP checksum on those filesystems is included for comparison.

We see that Fletcher’s, in general, out-performs the TCP checksum, and in some cases comes within a factor of 2 to a 1 in 2^{16} miss rate. This performance is curious given our results so far. First, Corollary 8 in the Appendix shows that, for uniformly distributed data and replacements larger than single words, Fletcher should not be any stronger than IP/TCP. Second, two empirical measures show that both TCP and Fletcher have a similar non-uniform distribution over individual cells. When looking at plots of checksums over 48-byte cells (Figure 3), the Fletcher’s checksum looks to have a non-uniform curve similar to that of TCP. And when we look at the probability of the checksum that two randomly chosen cells⁴ in the file system match each other, we find a probability of 0.016% for Fletcher 255, 0.013% for Fletcher 256, and 0.011% for IP/TCP.

Table 8: Fletcher’s Checksum Results (256 Byte packets on systems)

System	Missed by	%	splices
<i>sics.se</i> <code>/opt</code>	TCP	0.1703438788	5316323
	F255	0.0044358811	138441
	F256	0.0091286724	284900
<i>smeg.stanford.edu</i> <code>/ul</code>	TCP	0.0707199443	3108050
	F255	0.0862324604	3789805
	F256	0.0046759739	205503
<i>pompano.stanford.edu</i> <code>/usr/local</code>	TCP	0.0269563342	159707
	F255	0.0022121117	13106
	F256	0.0029058228	17216
<i>sics.se</i> <code>/src1</code>	TCP	0.0411719151	649734
	F255	0.0067998225	107308
	F256	0.0054134085	85429
<i>sics.se</i> <code>/src2</code>	TCP	0.0344980161	496823
	F255	0.0023053857	33201
	F256	0.0039193848	56445

Why then does Fletcher perform better than the TCP checksum? The most obvious effect is that the positional dependence of Fletcher’s checksum effectively increases the

⁴This includes all cells, including the short cell at the end of each packet, so the number does not match the “Measured Global” for $k = 1$, given earlier.

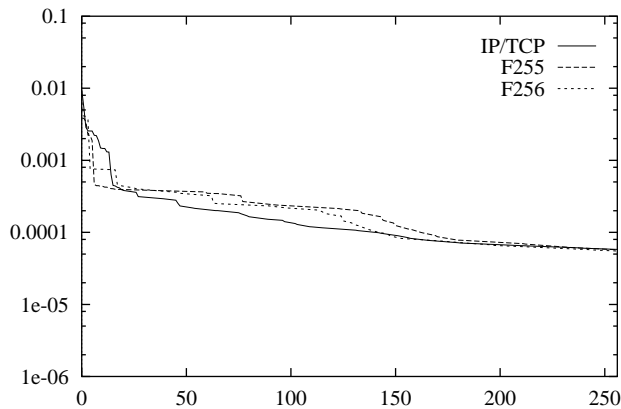


Figure 3: PDF of TCP checksum, F255, and F256 over 48 byte cells in smeg.dsg.stanford.edu:/u1. Most common 256 values.

number of cells changed in a splice. The vast majority of splices which pass IP and TCP header checks include the header cell from the first packet, and therefore the checksum field from the first packet. Each cell from the first packet not included in the splice moves *all* the subsequent cells from the first packet closer to the start of the splice - thus increasing the L_i s component of their contribution to the B field of the splice's checksum, when compared to their L_i contribution to the first packet's checksum. And even if the inserted cells from the second packet are identical to the dropped cells, their L_i s for the dropped cells is *different* than the L_i s of the inserted cells, as they appear later in the splice than in the first packet. The positionality of Fletcher's checksum means that the effective size of the splice is not just the total number of cells replaced, but includes any intervening, "reshuffled" cells from the first packet which lie between the first drop and the last replacement. (Note that this result has no effect on splices that join a prefix of the first packet to a suffix of the second.)

We know that the larger the number of cells, the more uniform the distribution, and thus, the lower the failure rate. However the increased substitution size is not sufficient to explain Fletcher's improvement over IP/TCP. If this hypothesis were true, the Fletcher miss rate should correspond, at best, to the "Actual" rate for $k = 4$ in Table 6. Instead, it's 10 times better. The reshuffling effect is real, but merely increasing the effective number of cells in a splice is only a small part of the story.

The real cause is more subtle. Recall that the condition for checksum failure is that the sum of the 8-bit A_i s be congruent and that the sum of the $(B_i + A_i L_i)$ s also be congruent. The condition on the A 's is identical to the condition for IP checksums. Since the data cells are drawn from the same highly localized non-uniform distribution, their 8-bit A_i terms have a fairly good chance of being congruent - at least 256 times more than the standard 16-

bit TCP sum. But for the B term, each of the R_i terms for individual cells are multiplied by L_i . This permutes the entire distribution. Thus, a given highly probable r drawn from one R_i is unlikely to be drawn from R_j , the distribution of R terms for a nearby cell drawn from the same distribution. In effect, the contribution of each cell to the B term of its packet is *colored* by its offset from the end of the packet (think of coloring the cells by their L_i number). This coloring, and the non-uniformity, combine to make undetected splices less likely. It is well known (we provide a proof in Lemma 9) that the probability of drawing two identical values from a non-uniform distribution is always higher than the probability of drawing two values that differ by any fixed amount. (This is discussed further in the Appendix). Since the data is non-uniform, some terms are more likely than others. The coloring effect of the B term means that the overall B sums of a splice are less likely to be congruent to the original checksum than if the data was uniformly distributed.

The end result is that the standard TCP checksum fails if two observations drawn from the same distribution are equal, while Fletcher fails if two observations drawn from the same distribution differ by a particular amount (where the exact amount varies from splice to splice). Thus, non-uniformity of the data actually strengthens the B field of the checksum. (This was probably *not* an intentional benefit planned by Fletcher.)

Ones complement Fletcher, however, has a weakness that sometimes offsets its probabilistic advantage: since bytes containing either 255 or 0 are considered identical by the checksum, certain common pathological cases cause total failure of Fletcher-255. This is discussed in more detail in Section 5.5.

5.3 Trailer Checksums: Making non-uniformity work for us

Fletcher-256 succeeds in detecting more splices than TCP by taking advantage of the non-uniformity of the data distributions, but it still has drawbacks. It is more expensive to compute, and the non-uniformity can only strengthen 8-bits of the checksum. It turns out that we can use a similar trick to exploit non-uniformity for the standard Internet checksum, with no computational cost. Further, we can strengthen the *entire* 16 bit sum, giving us (for some distributions) 16 bit checksums that are even stronger than 1 in 2^{16} .

The key observation is that with header checksums, the packet header and the packet checksum are located in the same cell of a packet. Thus, either both the header and the checksum covering it are present in a given splice, or neither are. The IP header check and syntactic TCP header checks ensure that almost all splices which are actually checksummed include the header from the first packet. The resulting splices will have the first packet's TCP header,

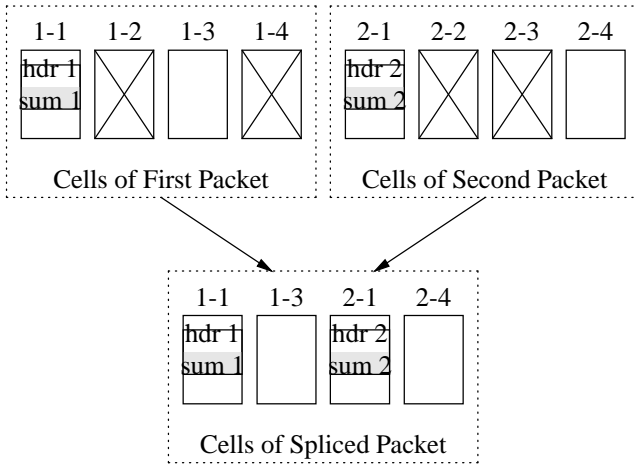


Figure 4: Header checksum fate

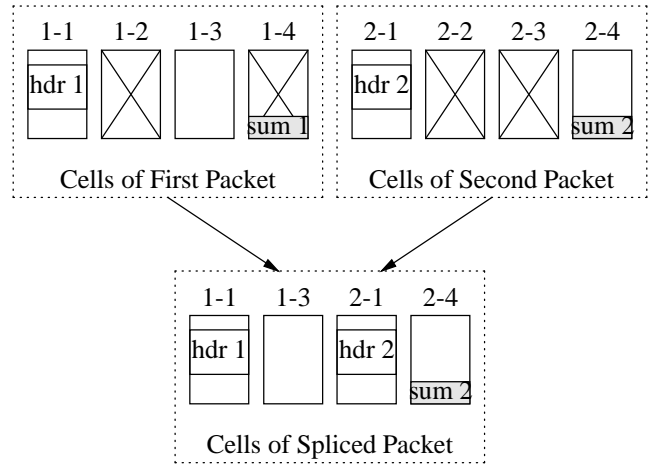


Figure 5: Trailer checksum fate

including TCP sequence number, ACK field, and checksum. As long as the replacement cells in the splice have the same overall checksum as the original packets, the TCP checksum will not detect the splice. Figure 4 shows one such splice diagrammatically. If the TCP checksum was at the end of the TCP packet, instead of in the header, the TCP checksum value would not share fate with the TCP pseudo-header which it covers. Figure 5 shows the same splice as Fig. 4, but with the TCP checksum located in a packet trailer instead of the packet header. Here, the resulting splice has the TCP header from the first packet and the checksum from the second packet. (It also has the header from the second packet, but only half the splices will do so.)

The data cells are, as with Fletcher’s sum, all drawn from the same localized non-uniform distribution and are more likely than 1 in 2^{16} to have congruent sums. But compare the two headers. The only field that changes between adjacent TCP packets in a given flow is the TCP sequence number. The difference between the checksums of the header cells of adjacent packets in a single flow is therefore strongly clustered around the size of the payload. In other words there are actually three different distributions of cells in a packet pair: the payload data, the first header, and the second header. If we separate the checksum value away from the header that it sums and put it in a trailer, we can ensure that there are always three different colors in any given splice – even for splices that only make color-preserving substitutions (e.g., data cell for data cell). Again by Lemma 9, this higher degree of coloring leads to a higher probability of detecting a splice than the standard header checksum, as we show below by case analysis⁵.

What is the probability of a trailer checksum failing? It’s

⁵Our study of trailer checksums was originally motivated by noticing that the AAL5 trailer checksums avoided this fate-sharing, and conjecturing that the predictable header differences would make TCP checksums better. It performed surprisingly well, leading us to the preceding re-analysis of the Fletcher checksum results.

simply that the checksum of the cells inserted from the first packet equal the checksum of the cells dropped from the second packet. (Note that we take the second packet - the source of the trailer sum - as the original, and counting cells from the first packet as insertions,) However, the inserted cells always include a header cell. The inserted cells from the first packet thus have a sum drawn from a distribution that consists of 1 header cell and k data cells. If the second header is dropped, then we again have k data cells and 1 header cell. However, in half of the splices the dropped cells are all data cells, in which case their sum consists of $k + 1$ data cells. The resulting probability will be lower than the probability of an exact match between two checksums drawn from the same distribution (as shown in Lemma 9). This would seem to reduce the failure probability by at most a factor of two.

In the remaining half of the splices, the second header is also dropped. Here the distribution of checksums of the header cell of the second packet does not match the distribution of the first header cell. There are two causes for this result. First, we treat the first header as a header but the second as data, which means we checksum the IP header of the second cell, but not of the first. Second, the header is mostly constant between packets except for an increase in the IP ID field and the TCP sequence number. (Note that in this scenario there is no checksum in the header: the field is left zero; though a practical trailer implementation might perhaps choose to swap the checksum value and the last two bytes of the packet.)

How much lower will the probability of failure be? We conducted an experiment to measure the effectiveness of trailer checksums. We changed the simulator to model a protocol identical to TCP, except that the TCP header checksum is left zero, and the checksum value is appended to the end of the TCP data. The results are shown in Table 9. The failure rate of trailer checksums were significantly better than those of TCP and Fletcher. We note that the failure rate

was actually below 2^{-16} for significant fractions of some file systems. In most cases we noted a failure rate 20 to 50 times lower than for header checksums. We can further

Table 9: **Trailer Checksum Results**
(256 Byte packets on systems)

Filesystem	TCP Misses	Trailer Misses
Uniform	0.001526	0.001526
<i>sics.se:/opt</i>	0.170344	0.004105
<i>smeg.stanford.edu:/u1</i>	0.070720	0.001735
<i>pompano.stanford:/usr/local</i>	0.026956	0.001604
<i>sics.se:/src1</i>	0.041172	0.002351
<i>sics.se:/src2</i>	0.034498	0.002100

test the distribution-coloring analysis by making predictions about the standard header TCP checksum. The number of splices which do not include the header of the first packet are negligible, so there are only two cases: the first header cell followed by all-data cells, and the header from the first packet followed by a mix of data cells and the header from the second packet. In the latter case the splice has replaced a data cell with the header cell from the second packet, and thus should be much less likely to match than the first case. When we went back and examined the data, this prediction was correct. Although roughly half of the splices surviving the header check have the second header included, only 1 in 2^{16} of those passed the TCP checksum. The TCP header checksum was 100-200 times more effective against splices that contained the second header. This result both supports our explanation of the good performance of trailer TCP checksums, and further confirms the utility of our distribution-coloring analysis of checksums.

5.4 Adding cell-coloring to our model

We are now ready to return to the discrepancies between our "Exclude Identical" probabilities in Table 6 of Section 4.6 and the actual measured failure rate. Recall that our sample probabilities predicted total failure rates very accurately for small k (the number of cells in a block), but by the time k increased to 4, the model over-predicted the measured failures by a factor of 3 or 4.

The piece that was missing from our model was the cell-coloring. The sample probabilities in our model were computed using only pure data cells, and thus missed the header effect. In our actual splice simulation, some substitutions of k cells replace a data cell with a header cell. The failure rate for the substitutions with headers should be 1 in 2^{16} , which is ignorable.

What is the probability that a substitution of length k replaced a data cell with a header cell? This is easy to compute. All k cells dropped from the first packet will be

data cells. There are

$$\binom{6}{k-1}$$

possible choices of k cell insertions from the second packet (recall that we must insert the trailing cell of the second packet in the splice). Of these, only

$$\binom{5}{k-1}$$

do not contain the header cell of the second packet. Therefore, to predict the actual failure rate of a k -cell substitution from our "Exclude Identical" samples, we must reduce the sample probability by a factor of

$$\frac{\binom{5}{k-1}}{\binom{6}{k-1}}$$

which equals $(7 - k)/6$. Our sample probabilities now closely match the actual measured failure probabilities, and we are reasonably confident that we have explained the behavior we have observed. Further, the improved performance due to trailer checksums in our packet-splice model seems to be real.

In the past, protocol designers have proposed trailer checksums for various engineering reasons. As far as we know, the argument about improved checksum behavior was not advanced. We conclude that protocol designers should reconsider placing checksums in packet trailers rather than headers, as has been standard practice in Internet protocols to date.

Trailer checksums suffer one apparent drawback. They may unnecessarily reject splices that are identical to an original packet. Consider the scenario where a burst of cell loss splices the front of one packet onto the tail of the following packet, as in Figure 4. If the payload of the splice is identical to the payload of the original packet, then the header checksum should match (since the header of the splice is the header of the first packet), and the packet is accepted. But with trailer checksums, (as in Figure 5, when the payload is identical to the first packet the checksum cannot match: it was computed with the sequence number of the second packet, not the first. So if the contents are identical the checksums will match only if the difference between the inserted and dropped cells is congruent to the difference in sequence number (the payload) between the two packets. By Lemma 9, this is very unlikely. Thus the splice will be rejected even when the contents is correct. The corresponding case (header of the second packet, payload of the first) never comes up, since our error model requires cells to remain ordered. In summary, trailer checksums have a very good chance of detecting a splice even if the resulting packet is a "good" packet.

Table 10 demonstrates this effect on the filesystem */u1* at *smeg.dsg.stanford.edu*. The number of identical splices

rejected by trailer checksums is larger than the number of bad splices they detect that the TCP checksum missed.

The two numbers, however, are not comparable. TCP missed checksums represent undetected data corruption. Spurious rejection by the trailer checksum represents (at worst) a possible performance penalty, it does not cause any data corruption.

Table 10: Header vs. Trailer Checksum Failure Rates

False Positive/Negative	header	trailer
Fails checksum, data identical	0	25,348,910
Passes checksum, data changed	3,108,050	76,270
Fails checksum, data identical	0.0%	0.57%
Passes checksum, data changed	0.07%	0.002%

Comparing missed splices, the trailer checksum misses less than 3% as often as the standard sum, but at the cost of reporting checksum failure on splices that accidentally resulted in a valid packet.

However, a splice in a real network always means that at least one packet has been lost, even if the splice is identical to one of the original packets. So a TCP retransmission will be necessary regardless. Thus the incremental performance impact of triggering retransmission one packet earlier when an identical splice is discarded is not clear.

5.5 Locality of Failure: Pathological Data Patterns

Non-uniformity in the distribution of checksums comes from two causes: non-uniformity of the underlying data, and weakness of a given checksum algorithm for certain patterns of data. That files on a computer system are highly structured is no surprise. We did not expect, however, to discover so many examples of files that were particularly vulnerable to splice-errors.

Though the Fletcher checksums consistently show a lower rate of failure than the standard Internet checksum, they also show a very high degree of locality. Sampling the checksum statistics incrementally during each whole-filesystem run showed sharp spikes in the rate of undetected splices, at the level of individual directories or even files. Manual examination of these files shows that, for each of the checksum algorithms, real data contains pathological data patterns which cause extremely localized rates of high failure.

The most dramatic case is the mod-255 Fletcher sum. This sum has two zeros, 0 and 255. Both these values contribute zero to the cell checksum. Thus, Fletcher mod 255 is susceptible to splice failure on long runs of mixed 0 and 255 bytes. The most dramatic example of this effect is one directory from the Stanford filesystem containing several 8-bit pbm graphs of Internet-backbone RTT measurements.

These graphs were plotted as black-and-white, and thus each byte is either 0 or 255. On these data, combinatorially, 1 in 2 of all permutations are caught by header checks and 1 in 2 of the remainder include a header cell. None of the remaining 25% of all possible permutations are caught by the mod-255 Fletcher. This one directory of files caused so many Fletcher mod-255 failures that on this filesystem, mod-255 Fletcher performs worse than the IP/TCP checksum.

Similarly spectacular mod-255 failures occurred in the Stanford filesystem with a file from a popular PC word processor. This file contained runs of approximately 200 all-zero bytes, followed by a similar number of all-one bytes, between each section of a document.

Fletcher's mod-256 sum behaves slightly differently. It has only one zero, and is not subject to the same dramatic failure as the mod-255 sum. Pathological data patterns for mod-256 do occur, but less frequently. One case we have isolated is hex-encoded PostScript bitmaps which contain identical segments of horizontal lines (e.g. bitmaps containing solid blocks of color, or bitmaps containing parallel lines. Font definitions appear to be a particularly common case). Many common bitmaps appear to have a width, W , that is a power of two. Thus, each ASCII-encoded binary line commonly consists of many "FF"s, and a small number of other two byte values (e.g. "F7") that repeat precisely $W + 1$ apart (The extra byte is due to an ASCII newline.) Though not immediately obvious on inspection, these just happen to combine in such a way that the contribution of 48-byte cells allows splices. We observed a similar effect in BinHex-encoded Macintosh documents stored on our Unix filesystem: very similar lines of 64 bytes followed by an ASCII newline.

Though the overall rate of TCP sum failures is higher than the other sums, and appears to be noisier, we have also isolated a few pathological cases for the standard Internet checksum. One example is Unix gmon.out profiling data. These files often consist mostly of zero entries, with a scattering of a small number of nonzero entries. The non-zero values are often identical. Packetizing this data results in a very small number of checksums. A very large number of splices pass the checksum, resulting in what appears to be scrambled files. A second example is the PostScript bitmap data file mentioned above, which showed pathological behavior for the Internet checksum as well.

Our central point is that the existence of pathological patterns for a given sum is not just theoretical; these patterns occur surprisingly frequently in real filesystem data.

6 Conjectures

In the course of our research we investigated several plausible conjectures that might have explained the TCP checksum failures. We briefly describe several of these blind-alleys.

6.1 The Role of Zero Data

The frequency of the zero checksum led us to study the effects of zeroed data on the checksum. It is no surprise that there are a lot of zeros in filesystem data (the UNIX filesystem has long been optimized such that completely zero blocks did not need to be saved on disk). However, knowing that arbitrarily long zero blocks do not change the IP checksum (zero is the additive identity), we wondered whether this property significantly affected the failure rate independent of the simple fact of their high frequency. In other words: Is there something special about zero? If we replaced all the zeros in the file-system with different values, would the failure rate change?

An approximate first answer is "no, zero is not special because it is the additive identity". If we add one to every word in the file system then the sum of every cell would increase by 24 (48 bytes divided by 2). Similarly, it is easy to demonstrate that the distribution of the sum of any number of cells will contain the same set of values and frequencies, although their mapping will be permuted. So the rate of checksum failure would be unchanged.⁶

It is, however, true that if any single value shows up a disproportionate amount of the time then the failure rate will increase. However, the reason that zero in particular is so common is that several totally independent formats all "happen" to choose zero as a common element. Further, it is likely that this will continue to be the case. Fortunately, although zero checksums do show up very frequently, it is often the result of cells consisting entirely of zeros. A substitution of one all-zero cell for another causes no harm. The problem, therefore, is the frequency of non-zero cells whose checksum is zero, in proximity to all-zero cells or to each other.

6.2 Zero congruent IP/TCP header cells

The TCP checksum is computed over a pseudo-header that covers all but eight bytes from the IP header. In our original simulations, those eight bytes of IP header – including the IP checksum – were not filled in. The TCP checksum is then inverted before it is stored in the header. This causes the checksum of an error-free TCP datagram, (including the TCP header), to be zero.

A full IP header also contains an inverted ones-complement checksum, which meant that the sum of the IP header was also zero. Since all but 8 bytes from the IP header are also covered by the TCP checksum, the checksum over the entire

⁶Zero is special, as we showed in the section on pathological cases, but not because it is the additive identity and does not affect the checksum. Zero's specialness comes from the fact that it is represented by both 0x0000 and 0xFFFF. In reality, adding 1 to every word in the file system would change the distribution of checksums, and might reduce the probability of the most probable value. Cells containing 0xFFFF's would be shifted by less than 24. Whether this would increase or decrease the most probable value depends on the distribution of values in each filesystem.

cell headers is not zero, but rather the checksum of the overlap: IP source and destination addresses, the length, and the TCP protocol ID.

Our earlier results ([7]) were based on simulations that left those eight bytes unfilled. Consider, however, two packets consisting of data that is all zero. This causes the header cell (when considered as data) to have a checksum of zero. The checksum will only be the sum of the header. When the checksum is inverted and stored into the header, we are left with a non-zero cell with a checksum of zero. In our earlier work, these cells were a major source of non-zero cells with a checksum of zero. What is worse, these cell show up precisely in the case when all the cells around it are zero cells (or at least zero-congruent). Thus replacement was common and a major source of splice failures. Filling in the IP header reduced the error rate by three orders of magnitude.

We had conjectured that filling in the IP header would not have much of an effect, because the length, IP addresses, and protocol type do not change between packets during the file transfer, and so the checksums of the header cell remain constant. However, even a constant, non-zero, value is sufficient to distinguish between header cells and zero filled data cells. This simulator deficiency also led us to give undue emphasis to the role of zero-congruent data (as mentioned above).

6.3 Inverted Checksums

Under the TCP and IP specification, the inverse of the checksum is placed in the packet header. This implies that the checksum of a valid segment will be zero. In [7] we cautioned implementors against this approach, since for mostly-zero packets the header cell, too, would be zero. This still is reasonable advice for packet formats as it reduces the frequency of zero congruent cells. However, it is not relevant to TCP and IP because of the overlap of the headers we noted above. To test this conjecture, we ran our tests with a modified version of the TCP checksum that did not invert the checksum before storing it into the packet. The results with the non-inverted checksum were almost identical to the results with an inverted checksum.

6.4 Corrections to SIGCOMM 95 version

As noted in Sec. 6.2 the data in our earlier paper[7] is not accurate. Completely filling in the IP header reduces the overall rate of errors by a factor of from 200 to 1000. In addition, the Fletcher checksum code was mis-implemented as a mixture of mod-255 and mod-256 arithmetic, which led to the Fletcher splice failure rate being higher than the standard TCP checksum. We retract that result; it was an artifact of the buggy Fletcher implementation. That bug was also the motivation for our current investigation of both mod-255 and mod-256 Fletcher sums. The artificially-high Fletcher failure rates also inspired the original work on trailer checksums.

The previous results also suffered from a number of other minor bugs, whose effect was insignificant compared to the two problems above. They are detailed in the Appendix.

7 Observations and Recommendations

The results of the previous sections lead to a number of interesting observations.

First, a non-uniform distribution of data makes failure of the TCP checksum far more likely than one would naively expect. The undetected splice rate in our data for the 16-bit TCP checksum over real data is comparable to uniform data with a 10-bit checksum.

Second, checksum distributions on modest amounts of real data are substantially different from the distributions one would anticipate for uniformly distributed data. This skewed distribution *does* result in significantly higher failure rates of the TCP checksum. In particular, if a router or host has a buffering problem that causes adjacent packets to be merged, the TCP checksum might fail .1% of the time rather than the 0.0015% of the time that purely random data distribution would suggest.

While these scenarios may seem worrisome, there are three pieces of good news.

First, it is important to keep in mind that these error scenarios are all quite rare. This work was initially motivated by studying extremely uncommon AAL5 error scenarios – an error model derived from ATM cell drop splicing two packets into one. In practice, such cell loss can occur due to either congestion or corruption. However, dropping ATM calls independently of each other is now known to cause goodput problems [10]. ATM switch vendors are addressing this problem by dropping all subsequent cells from a packet, once a single cell is dropped. This reduces the probability that a splice will be legal since a trailer will only be delivered if *all* preceding cells have been delivered. The cells from the partial preceding packets will result in a detectably incorrect packet length. This means we can effectively ignore congestion as a source of valid packet splices. Cell loss due to corruption is often estimated at 1 in 10^8 or less. The ATM CRC will fail to detect a splice approximately at a rate of 1 in 2^{32} . Therefore, the chance of the TCP checksum being called upon to detect a splice is much less than 1 in $10^{-8} * 2^{-32}$ or less than one chance in 10^{17} . Moreover, if ATM switch vendors institute Early Packet Discard not just for congestion but for all causes of cell drop, then valid packet splices should never appear.

Second, the packet splice model is, in some sense, a worst-case error model because the substitutions tend to be similar to the data that they replace. This is possibly also true of buffer-management errors, or errors in fragment reassembly. However, in the alternative error models where data is replaced by garbage, while the non-uniformity of the data may still reduce the effectiveness of checksums,

it will only reduce it to the extent that the distribution of the replacement data matches the distribution of the original data. Here, the frequency of long runs of 0's or 1's in the payload may make us slightly more vulnerable to hardware errors that produce similar runs of data. However, hardware failure that produces random bits are unlikely to produce runs of data that look a lot like English prose.

Third, and finally, remedies exist to improve the ability of checksums to work on non-uniform data.

- Compressing data clearly improves the performance of checksums. Since compression also typically reduces file transfer times and saves disk space, there's a strong motivation for FTP archives to compress their files.
- In the future, in the absence of compression, protocol designers should consider avoiding the practice of placing checksums in a protocol header, but instead append them as a trailer to the data being checksummed.
- In general, the checksums are rarely placed in a situation where it is the primary method of failure detection. (We are aware of one exception to this rule. The TCP checksum is the primary method of error detection over SLIP and Compressed SLIP links. That's probably not wise).

What this work simply shows is that checksums are an even less effective error detection method than first thought, because real data often has interesting distributions, and those distributions increase the likelihood of checksum failure.

Acknowledgments

The authors would like to acknowledge the help of Chuck Kalmanek and Bill Marshall of Bell Labs, who discussed issues of study design. We also gratefully acknowledge the help of David Feldmeier of Bellcore and Lansing Sloan of Lawrence Livermore, who helped us with substantially faster CRC computation algorithms, and the Swedish Institute of Computer Science, which allowed us to use its filesystems and one of its fast multiprocessors for some of the test runs.

References

- [1] BRADEN, R., BORMAN, D., AND PARTRIDGE, C. Computing the Internet Checksum. Internet Request For Comments RFC 1071, ISI, September 1988. (Updated by RFCs 1141 and 1624).
- [2] FLETCHER, J. An Arithmetic Checksum for Serial Transmissions. *IEEE Transactions on Communication* 30(1) (January 1983).
- [3] GRAHAM, R., KNUTH, D., AND PATASHNIK, O. *Concrete Mathematics: A Foundation for Computer Science*. Addison Wesley, 1989.

- [4] GREENE, D., AND LYLES, B. Reliability of Adaptation Layers. In *Protocols for High-Speed Networks III* (1992), B. Pehrson, P. Gunningberg, and S. Pink, Eds., Proc. IFIP 6.1/6.4 Workshop.
- [5] JOSEPH L. HAMMOND, JR, E. A. Development of a Transmission Error Model and an Error Control Model. Tech. rep., Georgia Institute of Technology, May 1975. prepared for Rome Air Development Center.
- [6] NAKASSIS, A. Fletcher’s Error Detection Algorithm: How to implement it efficiently and how to avoid the most common pitfalls. *Computer Communication Review* (October 1988), 63–88.
- [7] PARTRIDGE, C., HUGHES, J., AND STONE, J. Performance of Checksums and CRCs Over Real Data. In *Proc. SIGCOMM 1995* (Boston, 1996), vol. 25(4) of *ACM Computer Communications Review*, pp. 68–76.
- [8] PLUMMER, W. W. TCP Checksum Function Design. Internet Engineering Note 45, BBN, 1978. Reprinted in [1].
- [9] POSTEL, J. Transmission Control Protocol. Internet request for comments, ISI, September 1981. 3.
- [10] ROMANOW, A., AND FLOYD, S. Dynamics of TCP traffic over ATM networks. *IEEE JSAC* (May 1995), 633–641. An earlier version of this paper appeared in SIGCOMM ’94, pp. 79-88.
- [11] SKLOWER, K. Improving the Efficiency of the OSI Checksum Calculation. *Computer Communication Review* (October 1989), 44–55.
- [12] WANG, Z., AND CROWCROFT, J. SEAL Detects Cell Misordering. *IEEE Network Magazine* 6(4) (July 1992), 8–19.
- [13] ZWEIG, J., AND PARTRIDGE, C. TCP Alternate Checksum Options. Internet RFC RFC 1143, February 1990.

8 Appendix

This paper contains assertions which depend upon statements that are easily proven, yet not immediately obvious. Detailed explanations in the body of the paper would detract from the main argument. For those interested in the formal justification of some of our statements, we present more detail in this appendix.

8.1 Distributions of checksums

We use the notation $P + R$ to denote the distribution which arises by applying *any* commutative, total function $+$ with a unique inverse on a pair of values drawn from distributions P and R respectively. (In all of our cases we are interested

in the usual arithmetic addition operator). Call $PMax(P)$ the probability of the most likely value in the discrete distribution, P . (We define $PMin(P)$ similarly.) And define $P[i]$ is the probability of selecting i from P .

Lemma 1 $PMax(P + R) \leq \min(PMax(P), PMax(R))$

Proof: For any given x , the probability that the value drawn from $P + R = x$ is given by $S[x] = \sum_i P[i]R[x - i]$. Assume x is the most probable element of S . Without loss of generality, assume that $PMax(P) \leq PMax(R)$. $S[x] = \sum_i P[i]R[x - i] \leq Pmax(P) \sum_i R[x - i] \leq Pmax(P)$ (since $\sum_j R[j] = 1$). Equality would only hold if P were uniformly distributed and if $R[j] \neq 0 \Rightarrow P[x - j] \neq 0$. \square

Lemma 2 *If $\forall i, (P[i] = 0) \Rightarrow (R[x - i] = 0)$, then $PMin(P + R) \geq \max(PMin(P), PMin(R))$*

Proof: Consider the previous proof. Given the non-zero condition on $R[j]$, we are guaranteed that every value in R appears, and so $\sum_j R[j] = 1$, thus $S[x] = \sum_i P[i]R[x - i] \geq Pmin(P) \sum_i R[x - i] \geq Pmin(P)$. \square

This is unremarkable for unbounded discrete distributions. For the maximum, as the number of possible values grows, the probability of any single value must decrease. The conditions on the min require that $|P| \geq |R|$, and that $|S| = |P|$, so it is also unsurprising that the minimum doesn’t decrease. However, for bounded distributions, e.g distributions over the integers mod M , this leads to the following more interesting results.

Corollary 3 *Consider a probability distribution P over the integers mod M . The distribution of the sum, mod M , of j integers drawn from P gets “more uniform” as j increases, in the sense that the minimum probability of any number gets larger and the max probability gets smaller.* \square

Computation: If we have a random variable which can take on M values, with a known distribution of values, then the probability ($P_j[s = k]$) of the sum, s , of j values drawn from this distribution is equal to k , is:

$$\sum_{i=0}^M (P_{j-1}[s = (k - i) \bmod M] \times P_1[s = i]) \quad (1)$$

\square

Corollary 3 shows that each time we add another number to the sum mod M and look at the probability distribution, we increase $PMin(P)$ and decrease $PMax(P)$. We can prove another useful result: for large enough j , $PMin(P)$ and $PMax(P)$ both approach $1/M$ and the distribution approaches uniform.

If P has some zero probability values, then some values in the sum of P might also have zero probability, unless the gcd of M and the entries occurring with non-zero probability is 1. The following theorem applies even if a sum of a

distribution only has M' values with non-zero probability in the following sense: all non-zero values will tend to be equal to $1/M'$.

Theorem 4 (a Central Limit Theorem) *The sum, mod M , of a large number of independent observations from any distribution P tends to have a uniform distribution.*

Proof: We will show that for any given $\epsilon > 0$, there is some j such that $\text{PMax}(P_j) \leq 1/M + \epsilon$. Since $\text{PMax}(P_j)$ is non-increasing as j grows, we know this also holds for all $k > j$. Use the notation \max_j to mean $\text{PMax}(P_j)$, and \min_j to mean $\text{PMin}(P_j)$, when the meaning is clear.

Assume there is a distribution, P , where $\max_j > 1/M + \epsilon$ for all values of j . We can compute a strict upper bound for \max_{j+1} based on \max_j . The largest possible value of \max_{j+1} will arise when the most probable terms from P_0 match the most probable terms from P_j (c.f. exercise in Concrete Mathematics [3], at the bottom of page 38). Assume the probability for the $M - 1$ most common values in P_j are all \max_j , and there is 1 value whose probability is $\leq 1/M$. For P there is at least one value with probability \max_0 , one with probability \min_0 , and $M - 2$ values whose probability sums to $1 - \max_0 - \min_0$.

$$\begin{aligned} \max_{j+1} &\leq \max_0 \max_j + \\ &\quad (1 - \max_0 - \min_0) \max_j + \\ &\quad \min_0 \frac{1}{M} \\ \max_{j+1} &\leq \max_j - \min_0 \times (\max_j - \frac{1}{M}) \\ \max_{j+1} &\leq \max_j - \min_0 \times \epsilon \end{aligned}$$

But after adding $j = \max_0 / (\min_0 \epsilon)$ times, \max_j would be less than 0, given our assumption that \max_j is always greater than $1/M + \epsilon$. So, our assumption must be false.

Thus, for any distribution P and for any ϵ , there is some number j of additions, such that $\text{PMax}(P_j) < 1/M + \epsilon$, so the distribution of P_j tends to the uniform distribution as j gets larger. \square

8.2 Distributions of some checksums over uniformly distributed data

Most existing evaluations of competing checksum algorithms have assumed that single bit errors were common. It is now frequently true that there are CRCs in the data-link layer to protect the integrity of cells on the wire, and ECC to correct memory errors while packets sit in buffers on routers. Thus, the errors that the TCP checksum must protect against are no longer single or double bit errors (which will be detected or corrected by other means), but rather substitution of longer runs of “good” data by (possibly different length) runs of “other” data. How do the IP checksum and Fletcher compare under this substitution model?

This section discusses what their expected behavior would be under substitution errors if the data were, in fact, uniformly distributed⁷.

If we assume all packets are equally likely, then if we look at any unit smaller than the size of the substitution, we can assume that an error consists of replacements drawn uniformly from all strings.

Lemma 5 *The sum S mod M of N numbers, will be uniformly distributed among all M values assuming there is at least one term, U , in the sum which takes on values uniformly distributed mod M*

Proof: Assume $S - U \bmod M$ has an arbitrarily skewed distribution. $\text{PMax}(U) = 1/M$, and $\text{PMin}(U) = 1/M$. By lemma’s 1 and 2, $1/M \geq \text{PMax}(S) \geq \text{PMin}(S) \geq 1/M$. Thus, the probability that $S = x$ for any given x will be precisely $1/M$, so the probabilities are all equal and the distribution is uniform. \square

Theorem 6 *Given uniformly distributed data and the substitution model above, the IP checksum of the modified packet is uniformly distributed over all possible values*

Proof: We assume that errors are replacements drawn from the uniform distribution. Then (assuming replacements larger than a single 16 bit word) every word within the replaced chunk will be uniformly distributed mod M . Therefore, by Lemma 5, the IP checksum will be uniformly distributed under the assumed substitutions, since it is the sum of uniformly distributed words. That is, the checksum will only fail to detect errors (by the replacement string contributing an identical sum to the checksum as the original string) with a probability of 1 out of $2^{16} - 1$. \square

Theorem 7 *Given uniformly distributed data and the substitution model above, the Fletcher checksum of the modified packet is uniformly distributed over all possible values.*

Proof: The same reasoning can be applied to the Fletcher checksum over a chunk of data of size N . The Fletcher checksum consists of two sums. The first is the sum, mod M , of all the bytes in the chunk. The second is the sum mod M of each byte weighted by its offset, L , from the end of the chunk. Call these two sums, respectively, $K1$ and $K2$. The contribution of this chunk (assuming it is L_C from the end of

⁷It is worth noting that one point of the preceding paper is that data values are *not* distributed uniformly and *are* correlated with nearby values, and that, therefore, errors, under the substitution model, are also not distributed uniformly and checksums do not perform as well as expected. This work on uniformly distributed data is still interesting on three counts. First, statements in the main body of the paper depend on results presented here. Second, it provides us with a benchmark against which to measure the actual measured error rate (i.e. what is due to the substitution model and what is due to non-uniform data). Third, encryption and compression are both becoming more common and both tend to produce uniformly distributed data.

the packet) to the Fletcher checksum of the entire packet is straightforward. $K1$ is added, mod M , to the mod M sum of the rest of the packet. $L_C \times K1 + K2$ is added, mod M , to the weighted sum of the rest of the packet. If $K1$ for each chunk is uniformly distributed, then so will $\sum_C K1$. If each $K2$ is uniformly distributed, then so will $\sum_C (L_C \times K1 + K2)$, since by Lemma 5 we only need one uniformly distributed term (and $K2$ is, although $L_C \times K1$ might not be).

That $K1$ is uniformly distributed follows directly from the lemma. $K2$ is only slightly more complicated. As long as the chunks are large enough so that there is a byte B with offset L_0 from the end of the chunk, such that L_0 is relatively prime to M (i.e $\gcd(L_0, M) = 1$), then B 's contribution to $K2$ is uniformly distributed among all M values, and therefore $K2$ itself is also uniformly distributed. Since $L_0 = 1$ is relatively prime to M , as long as the chunk is at least $2 \log_2 M$ bits long, we can apply lemma 5.

We must also show that $K2$ is independent of $K1$, else $\{K1, K2\}$ will not be uniformly distributed. Suppose the last two bytes of the chunk are B_1 and B_0 . Under the assumption of uniform distribution of the data, B_1 and B_0 are both independent and uniformly distributed. B_0 does not affect $K2$ since it is multiplied by 0. As we show the uniform distribution of $K2$ by varying B_1 (as we did in the lemma above), for each B_1 we can choose any value for B_0 to allow $K1$ to take on all values equally, without affecting $K2$. So, for each value K that $K2$ might take on, $K1$ is independent and uniformly distributed. \square

One last complication arises with the Fletcher checksum. Like IP, Fletcher defines the values inserted into the checksum field to be the *negation* of the checksum of the rest of the packet, so that the packet sums to 0. With Fletcher this requires the two bytes of the checksum to be the solution to a system of simultaneous equations. We must show that these two *specific* bytes are independent, since we can no longer magically choose offsets 0 and 1.

Assume the Fletcher checksum $\{F1, F2\}$, is stored in adjacent bytes with offsets L_1 and $L_2 = L_1 - 1$ from the end of the packet. $0 = F1 + F2 + K1 \text{ mod } M$, and $0 = K2 + F1 \times L_1 + F2 \times (L_1 - 1) \text{ mod } M$.

$$\begin{aligned}
F1 &= M - K1 - F2 \\
0 &= K2 + (M - K1 - F2)L_1 + F2(L_1 - 1) \\
0 &= K2 - L_1 K1 - F2L_1 + F2L_1 - F2 \\
0 &= K2 - L_1 K1 - F2 \\
F2 &= K2 - L_1 K1 \\
F1 &= M - K1 - K2 + L_1 K1 \\
&= M + K1(L_1 - 1) - K2
\end{aligned}$$

Since $K2$ is uniformly distributed mod M , so are both $F2$ and $F1$. Since $F1 = F2 - K1 \text{ mod } M$, then $F1$ is still uniformly distributed even if we hold $F2$ fixed (since we can

vary $K2$ internal to $F2$). Therefore, $F1$ is independent of $F2$. \square

Note that $L_1 K1$ will not, in general, be uniformly distributed mod M , since we can't assume that $\gcd(L_1, M) = 1$ (in fact, in our example, L_1 was always equal to 260. $\gcd(260, 255) = 5$ and $\gcd(260, 256) = 4$).

As a curiosity, further note that if $L_2 - L_1$ were not relatively prime to M , then $F1$ and $F2$ would not have been independent or uniformly distributed. (In fact, the equations would not have always had solutions).

Corollary 8 *Given uniformly distributed data, and the substitution model described above, IP and Fletcher checksums are equivalently powerful* \square

8.3 Header checksums vs. trailer checksums

The body of the paper claims that under our splice error model, trailer checksums are stronger than header checksums for nonuniformly distribute data and, no worse for uniformly distributed data. Here we prove that claim.

Lemma 9 *Consider drawing 2 samples, X_0 and X_1 , from any discrete distribution. The probability that $X_0 = X_1$ is greater than or equal to the probability that $X_0 = (X_1 + d) \text{ mod } M$ for any given d .*

Proof: To see this, note that the probability of the former (identical match) is simply $\sum_{i=0}^M P[i]^2$. The probability of the latter (d greater than the first) is $\sum_{i=0}^M P[i]P[i+d]$, where $i + d$ is taken mod M . Double both sums and rearrange terms. Since $(P[i]^2 + P[i+d]^2) \geq 2P[i]P[i+d]$, the former sum is greater than the latter sum. \square

Consider our error model: we substitute j cells from the first packet with j other cells from the second packet. We keep the header cell of the first packet and we keep the trailer cell of the second packet. For a header checksum to fail, the sum s_1 and s_2 of each collection of j cell partial checksums must be equal. For a trailer checksum to fail, the sum s_1 of the j cells missing from the first packet must be d less than the s_2 , assuming that the checksum of the header cell of packet 1 is d less than the checksum of the header cell of packet 2. We distinguish d , the difference between the header cells, since the header cells are drawn from a very different distribution than the data cells, and further, the distribution of the difference of two consecutive header cells is strongly clustered around $d = 256$. Thus, we have:

Theorem 10 *Under our error model of splicing, a trailer checksum will always be at least as powerful as a header checksum.*

Proof: For any given splice we have substituted j cells. Equation 1 on Page 15 gives us the probability distribution of the sum of j cells. The probability that the header checksum fails is the probability that two samples drawn from P_j are

equal. As discussed above, for trailer checksums there is a fixed d , usually 256 in our simulation, computable by looking at the 2 header cells. The probability that the trailer fails is the probability that two samples from P_j differ by d . Lemma 9 above shows that the former is more likely than the latter, thus header checksums are weaker than trailer checksums. \square

Note, that in fact, this depends only on the property that the probability of the checksums over the header cells of two adjacent packets be congruent is lower than the probability that 2 data cells from the same packet be congruent. For computing the actual probability of trailer checksum failure it is useful to be able to model d as a constant 256, but this is not required for the proof.

9 Retractions from the SIGCOMM '95 paper

An earlier version of this paper appeared in SIGCOMM '95 [7].

The central point of that paper still holds: non-uniform distribution of data results in the IP checksum being weaker than expected. Several conjectures expressed in [7] have been resolved and were addressed in the main body of this paper.

However, several minor points and computational details were not correct and we retract them.

First, we expressed surprise (as well we should have) that the Fletcher checksum performed *worse* than the IP checksum. Performance tuning of the Fletcher checksum code used in that paper resulted in an incorrect implementation. The Fletcher code also used a mixture of mod-256 and mod-255 arithmetic and was not computing an accurate Fletcher checksum for either mod-255 or mod-256 Fletcher.

The numbers reported for the Fletcher checksum in that paper were, therefore, not accurate. The corrected numbers reported in this version of the paper show the expected result — Fletcher's detects more splices than TCP. However, the bugs in [7] and its anomalously poor results motivated us to investigate both mod-255 and mod-256 Fletcher, uncovering the pathological cases for mod-255 Fletcher reported here.

The SIGCOMM '95 paper reports numbers where the IP header fields not covered by the TCP checksum were left as zero.

Though covered in the body of this paper, it is important to emphasize it again here: filling in the header significantly reduced the number of matches for zero-congruent cells, and therefore reduced the total number of misses (by three orders of magnitude in some cases). By filling in the IP header in [7] we over-stated the significance of splices including zero-congruent cells and focused too closely on misses involving zero-filled or zero-congruent cells.

Several additional, but relatively minor bugs in the simulator compromised the accuracy of the numbers of all check-

sum algorithms in [7] (but only to a small factor).

First, we used the AAL5 length from the second packet, rather than the apparent IP length from the first cell, for checksum computation. This miscomputed checksums by including data from the last cell beyond the end of the IP payload in the checksum.

Second, this same error arose when testing whether packets were "identical" in payload. This resulted in counting certain splices as checksum failures, when in fact they were simply identical to the original packet, or where the first packet was a prefix of the splice.

Third, we miscomputed the checksum for short packets — that is, packets where the apparent IP header length made the entire TCP packet fit into the first cell and the AAL5 trailer in the second cell. It's well-known that a TCP packet with any user data fills at least two ATM cells. But for packets with 1 to 8 bytes of TCP payload, the entire IP/TCP datagram fits in only one cell and the second cell contains only an AAL5 trailer. Knowing that TCP data packets always take two cells, the simulation in [7] erroneously added a partial checksum for the second cell.

These erroneous calculations did not change the larger picture of TCP checksum performance, but did require us to recompute all data for this version of the paper.

Finally, our code and raw data are available from the URL <http://www.dsg.stanford.edu/>