

TCP Checksum Function Design

William W. Plummer

Bolt Beranek and Newman, Inc.
50 Moulton Street
Cambridge MA 02138

5 June 1978

1. Introduction

Checksums are included in packets in order that errors encountered during transmission may be detected. For Internet protocols such as TCP [1,9] this is especially important because packets may have to cross wireless networks such as the Packet Radio Network [2] and Atlantic Satellite Network [3] where packets may be corrupted. Internet protocols (e.g., those for real time speech transmission) can tolerate a certain level of transmission errors and forward error correction techniques or possibly no checksum at all might be better. The focus in this paper is on checksum functions for protocols such as TCP where the required reliable delivery is achieved by retransmission.

Even if the checksum appears good on a message which has been received, the message may still contain an undetected error. The probability of this is bounded by 2^{-C} where C is the number of checksum bits. Errors can arise from hardware (and software) malfunctions as well as transmission errors. Hardware induced errors are usually manifested in certain well known ways and it is desirable to account for this in the design of the checksum function. Ideally no error of the "common hardware failure" type would go undetected.

An example of a failure that the current checksum function handles successfully is picking up a bit in the network interface (or I/O buss, memory channel, etc.). This will always render the checksum bad. For an example of how the current function is inadequate, assume that a control signal stops functioning in the network interface and the interface stores zeros in place of the real data. These "all zero" messages appear to have valid checksums. Noise on the "There's Your Bit" line of the ARPANET Interface [4] may go undetected because the extra bits input may cause the checksum to be perturbed (i.e., shifted) in the same way as the data was.

Although messages containing undetected errors will occasionally be passed to higher levels of protocol, it is likely that they will not make sense at that level. In the case of TCP most such messages will be ignored, but some could cause a connection to be aborted. Garbled data could be viewed as a problem for a layer of protocol above TCP which itself may have a checksumming scheme.

This paper is the first step in design of a new checksum function for TCP and some other Internet protocols. Several useful properties of the current function are identified. If possible

these should be retained in any new function. A number of plausible checksum schemes are investigated. Of these only the "product code" seems to be simple enough for consideration.

2. The Current TCP Checksum Function

The current function is oriented towards sixteen-bit machines such as the PDP-11 but can be computed easily on other machines (e.g., PDP-10). A packet is thought of as a string of 16-bit bytes and the checksum function is the one's complement sum (add with end-around carry) of those bytes. It is the one's complement of this sum which is stored in the checksum field of the TCP header. Before computing the checksum value, the sender places a zero in the checksum field of the packet. If the checksum value computed by a receiver of the packet is zero, the packet is assumed to be valid. This is a consequence of the "negative" number in the checksum field exactly cancelling the contribution of the rest of the packet.

Ignoring the difficulty of actually evaluating the checksum function for a given packet, the way of using the checksum described above is quite simple, but it assumes some properties of the checksum operator (one's complement addition, "+" in what follows):

- (P1) + is commutative. Thus, the order in which the 16-bit bytes are "added" together is unimportant.
- (P2) + has at least one identity element (The current function has two: +0 and -0). This allows the sender to compute the checksum function by placing a zero in the packet checksum field before computing the value.
- (P3) + has an inverse. Thus, the receiver may evaluate the checksum function and expect a zero.
- (P4) + is associative, allowing the checksum field to be anywhere in the packet and the 16-bit bytes to be scanned sequentially.

Mathematically, these properties of the binary operation "+" over the set of 16-bit numbers forms an Abelian group [5]. Of course, there are many Abelian groups but not all would be satisfactory for use as checksum operators. (Another operator readily

available in the PDP-11 instruction set that has all of these properties is exclusive-OR, but XOR is unsatisfactory for other reasons.)

Albeit imprecise, another property which must be preserved in any future checksum scheme is:

(P5) + is fast to compute on a variety of machines with limited storage requirements.

The current function is quite good in this respect. On the PDP-11 the inner loop looks like:

```
LOOP:  ADD (R1)+,R0    ; Add the next 16-bit byte
        ADC R0        ; Make carry be end-around
        SOB R2,LOOP   ; Loop over entire packet.
```

(4 memory cycles per 16-bit byte)

On the PDP-10 properties P1-4 are exploited further and two 16-bit bytes per loop are processed:

```
LOOP:  ILOB THIS, PTR ; Get 2 16-bit bytes
        ADD SUM, THIS ; Add into current sum
        JUMPGE SUM, CHKSUZ ; Jump if fewer than 8 carries
        LDB THIS, [POINT 20, SUM, 19] ; Get left 16 and carries
        ANDI SUM, 177777 ; Save just low 16 here
        ADD SUM, THIS ; Fold in carries
CHKSUZ: SOJG COUNT, LOOP ; Loop over entire packet
```

(3.1 memory cycles per 16-bit byte)

The "extra" instruction in the loops above are required to convert the two's complement ADD instruction(s) into a one's complement add by making the carries be end-around. One's complement arithmetic is better than two's complement because it is equally sensitive to errors in all bit positions. If two's complement addition were used, an even number of 1's could be dropped (or picked up) in the most significant bit channel without affecting the value of the checksum. It is just this property that makes some sort of addition preferable to a simple exclusive-OR which is frequently used but permits an even number of drops (pick ups) in any bit channel. RIM10B paper tape format used on PDP-10s [10] uses two's complement add because space for the loader program is extremely limited.

Another property of the current checksum scheme is:

- (P6) Adding the checksum to a packet does not change the information bytes. Peterson [6] calls this a "systematic" code.

This property allows intermediate computers such as gateway machines to act on fields (i.e., the Internet Destination Address) without having to first decode the packet. Cyclical Redundancy Checks used for error correction are not systematic either. However, most applications of CRCs tend to emphasize error detection rather than correction and consequently can send the message unchanged, with the CRC check bits being appended to the end. The 24-bit CRC used by ARPANET IMPs and Very Distant Host Interfaces [4] and the ANSI standards for 800 and 6250 bits per inch magnetic tapes (described in [11]) use this mode.

Note that the operation of higher level protocols are not (by design) affected by anything that may be done by a gateway acting on possibly invalid packets. It is permissible for gateways to validate the checksum on incoming packets, but in general gateways will not know how to do this if the checksum is a protocol-specific feature.

A final property of the current checksum scheme which is actually a consequence of P1 and P4 is:

- (P7) The checksum may be incrementally modified.

This property permits an intermediate gateway to add information to a packet, for instance a timestamp, and "add" an appropriate change to the checksum field of the packet. Note that the checksum will still be end-to-end since it was not fully recomputed.

3. Product Codes

Certain "product codes" are potentially useful for checksumming purposes. The following is a brief description of product codes in the context of TCP. More general treatment can be found in Avizienis [7] and probably other more recent works.

The basic concept of this coding is that the message (packet) to be sent is formed by transforming the original source message and adding some "check" bits. By reading this and applying a (possibly different) transformation, a receiver can reconstruct

the original message and determine if it has been corrupted during transmission.

Mo		Ms		Mr
----- A B C -----	code ==>	----- 7 1 4 ... 2 -----	decode ==>	----- A B C ----- plus "valid" flag
Original		Sent		Reconstructed

With product codes the transformation is $M_s = K * M_o$. That is, the message sent is simply the product of the original message M_o and some well known constant K . To decode, the received M_s is divided by K which will yield M_r as the quotient and θ as the remainder if M_r is to be considered the same as M_o .

The first problem is selecting a "good" value for K , the "check factor". K must be relatively prime to the base chosen to express the message. (Example: Binary messages with K incorrectly chosen to be 8. This means that M_s looks exactly like M_o except that three zeros have been appended. The only way the message could look bad to a receiver dividing by 8 is if the error occurred in one of those three bits.)

For TCP the base R will be chosen to be 2^{16} . That is, every 16-bit byte (word on the PDP-11) will be considered as a digit of a big number and that number is the message. Thus,

$$M_o = \sum_{i=0}^N [B_i * (R^{i-1})] \quad , \quad B_i \text{ is } i\text{-th byte}$$

$$M_s = K * M_o$$

Corrupting a single digit of M_s will yield $M_s' = M_s + or - C * (R^{i-j})$ for some radix position j . The receiver will compute $M_s' / K = M_o + or - C * (R^{i-j}) / K$. Since R and K are relatively prime, $C * (R^{i-j})$ cannot be any exact multiple of K . Therefore, the division will result in a non-zero remainder which indicates that M_s' is a corrupted version of M_s . As will be seen, a good choice for K is $(R^{i-b} - 1)$, for some b which is the "check length" which controls the degree of detection to be had for

burst errors which affect a string of digits (i.e., 16-bit bytes) in the message. In fact b will be chosen to be 1, so K will be $2^{16} - 1$ so that arithmetic operations will be simple. This means that all bursts of 15 or fewer bits will be detected. According to [7] this choice for b results in the following expression for the fraction of undetected weight 2 errors:

$$f = 16(k-1) / [32(16k-3) + (6/k)] \quad \text{where } k \text{ is the message length.}$$

For large messages f approaches 3.125 per cent as k goes to infinity.

Multiple precision multiplication and division are normally quite complex operations, especially on small machines which typically lack even single precision multiply and divide operations. The exception to this is exactly the case being dealt with here -- the factor is $2^{16} - 1$ on machines with a word length of 16 bits. The reason for this is due to the following identity:

$$Q \times (R \times j) = Q, \text{ mod } (R-1) \quad 0 \leq Q < R$$

That is, any digit Q in the selected radix ($0, 1, \dots, R-1$) multiplied by any power of the radix will have a remainder of Q when divided by the radix minus 1.

Example: In decimal $R = 10$. Pick $Q = 6$.

$$\begin{aligned} 6 &= 0 \times 9 + 6 = 6, \text{ mod } 9 \\ 60 &= 6 \times 9 + 6 = 6, \text{ mod } 9 \\ 600 &= 66 \times 9 + 6 = 6, \text{ mod } 9 \quad \text{etc.} \end{aligned}$$

More to the point,

$$\begin{aligned} \text{rem}(31415/9) &= \text{rem}((30000+1000+400+10+5)/9) \\ &= (3 \text{ mod } 9) + (1 \text{ mod } 9) + (4 \text{ mod } 9) + (1 \text{ mod } 9) + (5 \text{ mod } 9) \\ &= (3+1+4+1+5) \text{ mod } 9 \\ &= 14 \text{ mod } 9 \\ &= 5 \end{aligned}$$

So, the remainder of a number divided by the radix minus one can be found by simply summing the digits of the number. Since the radix in the TCP case has been chosen to be 2^{16} and the check factor is $2^{16} - 1$, a message can quickly be checked by summing all of the 16-bit words (on a PDP-11), with carries being end-around. If zero is the result, the message can be considered valid. Thus, checking a product coded message is exactly the same complexity as with the current TCP checksum!

In order to form M_s , the sender must multiply the multiple precision "number" M_o by $2^{16} - 1$. Or, $M_s = (2^{16}M_o) - M_o$. This is performed by shifting M_o one whole word's worth of precision and subtracting M_o . Since carries must propagate between digits, but it is only the current digit which is of interest, one's complement arithmetic is used.

$$\begin{array}{r}
 (2^{16})M_o = M_o0 + M_o1 + M_o2 + \dots + M_oX + 0 \\
 - M_o = \quad - (M_o0 + M_o1 + \dots + M_oX) \\
 \hline
 M_s = M_s0 + M_s1 + \dots - M_oX
 \end{array}$$

A loop which implements this function on a PDP-11 might look like:

```

LOOP:  MOV -2(R2),R0 ; Next byte of (216)Mo
        SBC R0      ; Propagate carries from last SUB
        SUB (R2)+,R0 ; Subtract byte of Mo
        MOV R0,(R3)+ ; Store in Ms
        SOB R1,LOOP ; Loop over entire message
  
```

(8 memory cycles per 16-bit byte)

Note that the coding procedure is not done in-place since it is not systematic. In general the original copy, M_o , will have to be retained by the sender for retransmission purposes and therefore must remain readable. Thus the `MOV R0,(R3)+` is required which accounts for 2 of the 8 memory cycles per loop.

The coding procedure will add exactly one 16-bit word to the message since $M_s < (2^{16})M_o$. This additional 16 bits will be at the tail of the message, but may be moved into the defined location in the TCP header immediately before transmission. The receiver will have to undo this to put M_s back into standard format before decoding the message.

The code in the receiver for fully decoding the message may be inferred by observing that any word in M_s contains the difference between two successive words of M_o minus the carries from the previous word, and the low order word contains minus the low word of M_o . So the low order (i.e., rightmost) word of M_r is just the negative of the low order byte of M_s . The next word of M_r is the next word of M_s plus the just computed word of M_r plus the carry from that previous computation.

A slight refinement of the procedure is required in order to protect against an all-zero message passing to the destination. This will appear to have a valid checksum because $M_s'/K = 0/K$

= 0 with 0 remainder. The refinement is to make the coding be $M_s = K \times M_o + C$ where C is some arbitrary, well-known constant. Adding this constant requires a second pass over the message, but this will typically be very short since it can stop as soon as carries stop propagating. Choosing $C = 1$ is sufficient in most cases.

The product code checksum must be evaluated in terms of the desired properties P1 - P7. It has been shown that a factor of two more machine cycles are consumed in computing or verifying a product code checksum (P5 satisfied?).

Although the code is not systematic, the checksum can be verified quickly without decoding the message. If the Internet Destination Address is located at the least significant end of the packet (where the product code computation begins) then it is possible for a gateway to decode only enough of the message to see this field without having to decode the entire message. Thus, P6 is at least partially satisfied. The algebraic properties P1 through P4 are not satisfied, but only a small amount of computation is needed to account for this -- the message needs to be reformatted as previously mentioned.

P7 is satisfied since the product code checksum can be incrementally updated to account for an added word, although the procedure is somewhat involved. Imagine that the original message has two halves, H1 and H2. Thus, $M_o = H1 \times (R^{j-1}) + H2$. The timestamp word is to be inserted between these halves to form a modified $M_o' = H1 \times (R^{j-1}) + T \times (R^{j-1}) + H2$. Since K has been chosen to be $R-1$, the transmitted message $M_s' = M_o' (R-1)$. Then,

$$M_s' = M_s \times R + T \times (R^{j-1}) + P2 \times ((R-1) \times R^{j-2})$$

$$= M_s \times R + T \times (R^{j-1}) + T \times (R^{j-1}) + P2 \times (R^{j-2}) - 2 \times P2 \times R - P2$$

Recalling that R is 2^{16} , the word size on the PDP-11, multiplying by R means copying down one word in memory. So, the first term of M_s' is simply the unmodified message copied down one word. The next term is the new data T added into the M_s' being formed beginning at the $(j+1)$ th word. The addition is fairly easy here since after adding in T all that is left is propagating the carry, and that can stop as soon as no carry is produced. The other terms can be handle similarly.

4. More Complicated Codes

There exists a wealth of theory on error detecting and correcting codes. Peterson [6] is an excellent reference. Most of these "CRC" schemes are designed to be implemented using a shift register with a feedback network composed of exclusive-ORs. Simulating such a logic circuit with a program would be too slow to be useful unless some programming trick is discovered.

One such trick has been proposed by Kirstein [8]. Basically, a few bits (four or eight) of the current shift register state are combined with bits from the input stream (from M_0) and the result is used as an index to a table which yields the new shift register state and, if the code is not systematic, bits for the output stream (M_s). A trial coding of an especially "good" CRC function using four-bit bytes showed this technique to be about four times as slow as the current checksum function. This was true for both the PDP-10 and PDP-11 machines. Of the desirable properties listed above, CRC schemes satisfy only P3 (It has an inverse.), and P6 (It is systematic.). Placement of the checksum field in the packet is critical and the CRC cannot be incrementally modified.

Although the bulk of coding theory deals with binary codes, most of the theory works if the alphabet contains q symbols, where q is a power of a prime number. For instance q taken as 2^{16} should make a great deal of the theory useful on a word-by-word basis.

5. Outboard Processing

When a function such as computing an involved checksum requires extensive processing, one solution is to put that processing into an outboard processor. In this way "encode message" and "decode message" become single instructions which do not tax the main host processor. The Digital Equipment Corporation VAX/780 computer is equipped with special hardware for generating and checking CRCs [13]. In general this is not a very good solution since such a processor must be constructed for every different host machine which uses TCP messages.

It is conceivable that the gateway functions for a large host may be performed entirely in an "Internet Frontend Machine". This machine would be responsible for forwarding packets received

either from the network(s) or from the Internet protocol modules in the connected host, and for reassembling Internet fragments into segments and passing these to the host. Another capability of this machine would be to check the checksum so that the segments given to the host are known to be valid at the time they leave the frontend. Since computer cycles are assumed to be both inexpensive and available in the frontend, this seems reasonable.

The problem with attempting to validate checksums in the frontend is that it destroys the end-to-end character of the checksum. If anything, this is the most powerful feature of the TCP checksum! There is a way to make the host-to-frontend link be covered by the end-to-end checksum. A separate, small protocol must be developed to cover this link. After having validated an incoming packet from the network, the frontend would pass it to the host saying "here is an Internet segment for you. Call it #123". The host would save this segment, and send a copy back to the frontend saying, "Here is what you gave me as #123. Is it OK?". The frontend would then do a word-by-word comparison with the first transmission, and tell the host either "Here is #123 again", or "You did indeed receive #123 properly. Release it to the appropriate module for further processing."

The headers on the messages crossing the host-frontend link would most likely be covered by a fairly strong checksum so that information like which function is being performed and the message reference numbers are reliable. These headers would be quite short, maybe only sixteen bits, so the checksum could be quite strong. The bulk of the message would not be checksummed of course.

The reason this scheme reduces the computing burden on the host is that all that is required in order to validate the message using the end-to-end checksum is to send it back to the frontend machine. In the case of the PDP-10, this requires only 0.5 memory cycles per 16-bit byte of Internet message, and only a few processor cycles to setup the required transfers.

6. Conclusions

There is an ordering of checksum functions: first and simplest is none at all which provides no error detection or correction. Second, is sending a constant which is checked by the receiver. This also is extremely weak. Third, the exclusive-OR of the data may be sent. XOR takes the minimal amount of computer time to generate and check, but is not a good checksum. A two's complement sum of the data is somewhat better and takes no more

computer time to compute. Fifth, is the one's complement sum which is what is currently used by TCP. It is slightly more expensive in terms of computer time. The next step is a product code. The product code is strongly related to one's complement sum, takes still more computer time to use, provides a bit more protection against common hardware failures, but has some objectionable properties. Next is a genuine CRC polynomial code, used for checking purposes only. This is very expensive for a program to implement. Finally, a full CRC error correcting and detecting scheme may be used.

For TCP and Internet applications the product code scheme is viable. It suffers mainly in that messages must be (at least partially) decoded by intermediate gateways in order that they can be forwarded. Should product codes not be chosen as an improved checksum, some slight modification to the existing scheme might be possible. For instance the "add and rotate" function used for paper tape by the PDP-6/10 group at the Artificial Intelligence Laboratory at M.I.T. Project MAC [12] could be useful if it can be proved that it is better than the current scheme and that it can be computed efficiently on a variety of machines.

References

- [1] Cerf, V.G. and Kahn, Robert E., "A Protocol for Packet Network Communications," IEEE Transactions on Communications, vol. COM-22, No. 5, May 1974.
- [2] Kahn, Robert E., "The Organization of Computer Resources into a Packet Radio Network", IEEE Transactions on Communications, vol. COM-25, no. 1, pp. 169-178, January 1977.
- [3] Jacobs, Irwin, et al., "CPODA - A Demand Assignment Protocol for SatNet", Fifth Data Communications Symposium, September 27-9, 1977, Snowbird, Utah
- [4] Bolt Beranek and Newman, Inc. "Specifications for the Interconnection of a Host and an IMP", Report 1822, January 1976 edition.
- [5] Dean, Richard A., "Elements of Abstract Algebra", John Wiley and Sons, Inc., 1966
- [6] Peterson, W. Wesley, "Error Correcting Codes", M.I.T. Press Cambridge MA, 4th edition, 1968.
- [7] Avizienis, Algirdas, "A Study of the Effectiveness of Fault-Detecting Codes for Binary Arithmetic", Jet Propulsion Laboratory Technical Report No. 32-711, September 1, 1965.
- [8] Kirstein, Peter, private communication
- [9] Cerf, V. G. and Postel, Jonathan B., "Specification of Internetwork Transmission Control Program Version 3", University of Southern California Information Sciences Institute, January 1978.
- [10] Digital Equipment Corporation, "PDP-10 Reference Handbook", 1970, pp. 114-5.
- [11] Swanson, Robert, "Understanding Cyclic Redundancy Codes", Computer Design, November, 1975, pp. 93-99.
- [12] Clements, Robert C., private communication.
- [13] Conklin, Peter F., and Rodgers, David P., "Advanced Minicomputer Designed by Team Evaluation of Hardware/Software Tradeoffs", Computer Design, April 1978, pp. 136-7.