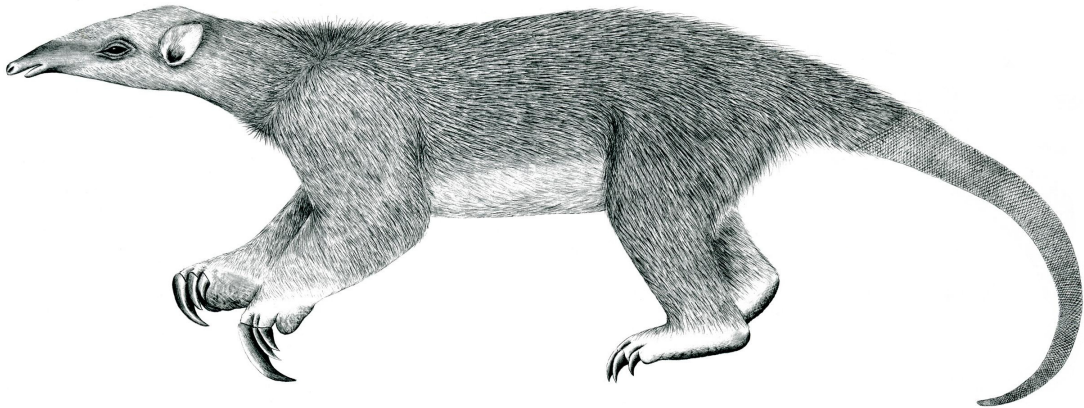

Tranalyzer2

tawk



Awk for Tranalyzer Flow Files



Tranalyzer Development Team

Contents

1	tawk	1
1.1	Description	1
1.2	Dependencies	1
1.3	Installation	1
1.4	Usage	1
1.5	-s and -N Options	2
1.6	Related Utilities	3
1.7	Functions	3
1.8	Examples	9
1.9	t2nfdump	10
1.10	t2custom	11
1.11	Writing a tawk Function	11
1.12	Using tawk Within Scripts	12
1.13	Using tawk With Non-Tranalyzer Files	12
1.14	Awk Cheat Sheet	13
1.15	Awk Templates	14
1.16	Examples	16
1.17	FAQ	17

1 tawk

1.1 Description

This document describes tawk and its functionalities. tawk works just like awk, but provides access to the columns via their names. In addition, it provides access to helper functions, such as `host()` or `port()`. Custom functions can be added in the folder named `t2custom` where they will be automatically loaded.

1.2 Dependencies

gawk version 4.1 is required.

Ubuntu:	sudo apt-get install	gawk
Arch:	sudo pacman -S	gawk
Gentoo:	sudo emerge	gawk
openSUSE:	sudo zypper install	gawk
Red Hat/Fedora¹:	sudo dnf install	gawk
macOS²:	brew install	gawk

1.3 Installation

The recommended way to install tawk is to install `t2_aliases` as documented in `README.md`:

- Append the following line to `~/.bashrc` (make sure to replace `$T2HOME` with the actual path, e.g., `$HOME/tranalyzer2-0.9.0`):

```
if [ -f "$T2HOME/scripts/t2_aliases" ]; then
    . $T2HOME/scripts/t2_aliases          # Note the leading `.'
fi
```

1.3.1 Man Pages

The man pages for tawk and `t2nfdump` can be installed by running: `./install.sh man`. Once installed, they can be consulted by running `man tawk` and `man t2nfdump` respectively.

1.4 Usage

- To list the column numbers and names: `tawk -l file_flows.txt`
- To list the column numbers and names as 3 columns: `tawk -l=3 file_flows.txt`
- To list the available functions: `tawk -g file_flows.txt`
- To list the available functions as 3 columns: `tawk -g=3 file_flows.txt`
- To save the original filename and filter used: `tawk -c 'FILTER' file_flows.txt > file.txt`

¹If the `dnf` command could not be found, try with `yum` instead

²Brew is a packet manager for macOS that can be found here: <https://brew.sh>

- To extract all ICMP flows and the header: `tawk 'hdr() || $l4Proto == 1' file_flows.txt > icmp.txt`
- To extract all ICMP flows without the header: `tawk -H 'icmp()' file_flows.txt > icmp.txt`
- To extract the flow with index 1234: `tawk '$flowInd == 1234' file_flows.txt`
- To extract all DNS flows and the header: `tawk 'hdr() || strtonum($dnsStat)' file_flows.txt`
- To consult the documentation for the function 'func': `tawk -d func`
- To consult the documentation for the functions 'min' and 'max': `tawk -d min,max`
- To consult the documentation for all the available functions: `tawk -d all`
- To consult the documentation for the variable 'var': `tawk -V var`
- To consult the documentation for the variable 'var' with value 0x8a: `tawk -V var=0x8a`
- To decode all variables from tranalyzer2 log file: `tawk -L out_log.txt`
- To decode all variables from tranalyzer2 log file (stdout): `t2 -r file.pcap | tawk -L`
- To convert the output to JSON: `tawk 'json($flowStat "\t" tuple5())' file_flows.txt`
- To convert the output to JSON: `tawk 'aggr(tuple2())' file_flows.txt | tawk 'json()'`
- To create a PCAP with all packets from flow 42: `tawk -x flow42.pcap '$flowInd == 42' file_flows.txt`
- To create a PCAP with packets 4-10: `tawk -P -x pkts-4_to_10.pcap 'packet("4-10")' file_flows.txt`
- To see all ICMP packets in Wireshark: `tawk -k 'icmp()' file_flows.txt`
- To see packet 4, 10 and 42 in Wireshark: `tawk -P -k 'packet("4;10;42")' file_flows.txt`

For a complete list of options, use the `-h` option.

Note that an option not recognized by `tawk` is internally passed to `awk/gawk`. One of the most useful is the `-v` option to set the value of a variable:

- Changing the output field separator:
`tawk -v OFS=',' '{ print $col1, $col2 }' file.txt`
- Passing a variable to `tawk`:
`tawk -v myvar=myvalue '{ print $col1, myvar }' file.txt`

For a complete list of options, run `awk -h`.

1.5 -s and -N Options

The `-s` option can be used to specify the starting character(s) of the row containing the column names (default: ``%'`). If several rows start with the specified character(s), then the last one is used as column names. To change this behavior, the line number can be specified as well with the help of the `-N` option. For example, if rows 1 to 5 start with ``#'` and row 3 contains the column names, specify the separator as follows: `tawk -s '`#' -N 3` If the row with column names does not start with a special character, use `-s '`'`.

1.6 Related Utilities

1.6.1 awkf

Configure `awk` to use tabs, i.e., `'\t'` as input and output separator (prevent issue with repetitive values), e.g.,
`awkf '{ print $4 }' file_flows.txt`

1.6.2 lsx

Display columns with fixed width (default: 40), e.g., `lsx file_flows.txt` or `lsx 45 file_flows.txt`

1.6.3 sortu

Sort rows and count the number of times a given row appears, then sort by the most occurring rows. (Alias for `sort | uniq -c | sort -rn`). Useful, e.g., to analyze the most occurring user-agents: `tawk '{ print $httpUsrAg }' FILE_flows.txt | sortu`

sortup

Same as `sortu`, but display the relative percentage instead of the absolute count. e.g., to analyze the most occurring user-agents: `tawk '{ print $httpUsrAg }' FILE_flows.txt | sortup`

1.6.4 tcol

Display columns with minimum width, e.g., `tcol file_flows.txt`.

1.7 Functions

Collection of functions for `tawk`:

- Parameters between brackets are optional,
- IPs can be given as string ("1.2.3.4"), hexadecimal (0xffffffff) or int (4294967295),
- Network masks can be given as string ("255.255.255.0"), hexadecimal (0xfffff00) or CIDR notation (24),
- Networks can be given as string, hexadecimal or int, e.g., "1.2.3.4/24" or "0x01020304/255.255.255.0",
- String functions can be made case insensitive by adding the suffix `i`, e.g., `streq` → `streqi`,
- Some examples are provided below,
- More details and examples can be found for every function by running `tawk -d funcname`.

Function	Description
<code>hdr()</code>	Use this function in your tests to keep the header (column names).
<code>tuple2()</code>	Return the 2 tuple (source IP and destination IP).
<code>tuple3()</code>	Return the 3 tuple (source IP, destination IP and port).
<code>tuple4()</code>	Return the 4 tuple (source IP and port, destination IP and port).
<code>tuple5()</code>	Return the 5 tuple (source IP and port, destination IP and port, protocol).

Function	Description
<code>tuple6()</code>	Return the 6 tuple (source IP and port, dest. IP and port, proto, VLANID).
<code>host([ip net])</code>	Return true if the source or destination IP is equal to <code>ip</code> or belongs to <code>net</code> . If <code>ip</code> is omitted, return the source and destination IP.
<code>shost([ip net])</code>	Return true if the source IP is equal to <code>ip</code> or belongs to <code>net</code> . If <code>ip</code> is omitted, return the source IP.
<code>dhost([ip net])</code>	Return true if the destination IP is equal to <code>ip</code> or belongs to <code>net</code> . If <code>ip</code> is omitted, return the destination IP.
<code>net([ip net])</code>	Alias for <code>host([ip net])</code> .
<code>snet([ip net])</code>	Alias for <code>shost([ip net])</code> .
<code>dnet([ip net])</code>	Alias for <code>dhost([ip net])</code> .
<code>loopback(ip)</code>	Return true if <code>ip</code> is a loopback address.
<code>mcast(ip)</code>	Return true if <code>ip</code> is a multicast address.
<code>privip(ip)</code>	Return true if <code>ip</code> is a private IP.
<code>port([p])</code>	Return true if the source or destination port appears in <code>p</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>port("1-3")</code> . If <code>p</code> is omitted, return the source and destination port.
<code>dport([p])</code>	Return true if the destination port appears in <code>p</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>dport("1-3")</code> . If <code>p</code> is omitted, return the destination port.
<code>sport([p])</code>	Return true if the source port appears in <code>p</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>sport("1-3")</code> . If <code>p</code> is omitted, return the source port.
<code>ip()</code>	Return true if the flow contains IPv4 or IPv6 traffic.
<code>ipv4()</code>	Return true if the flow contains IPv4 traffic.
<code>ipv6()</code>	Return true if the flow contains IPv6 traffic.
<code>proto([p])</code>	Return true if the protocol number appears in <code>p</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>proto("1-3")</code> . If <code>p</code> is omitted, return the protocol number.
<code>proto2str([p])</code>	Return the string representation of the protocol number <code>p</code> . If <code>p</code> is omitted, return the string representation of the protocol.
<code>icmp([p])</code>	Return true if the protocol is equal to 1 (ICMP).
<code>igmp([p])</code>	Return true if the protocol is equal to 2 (IGMP).
<code>tcp([p])</code>	Return true if the protocol is equal to 6 (TCP).
<code>udp([p])</code>	Return true if the protocol is equal to 17 (UDP).
<code>rsvp([p])</code>	Return true if the protocol is equal to 46 (RSVP).
<code>gre([p])</code>	Return true if the protocol is equal to 47 (GRE).
<code>esp([p])</code>	Return true if the protocol is equal to 50 (ESP).
<code>ah([p])</code>	Return true if the protocol is equal to 51 (AH).
<code>icmp6([p])</code>	Return true if the protocol is equal to 58 (ICMPv6).

Function	Description
<code>sctp([p])</code>	Return true if the protocol is equal to 132 (SCTP).
<code>dhcp()</code>	Return true if the flow contains DHCP traffic.
<code>dns()</code>	Return true if the flow contains DNS traffic.
<code>http()</code>	Return true if the flow contains HTTP traffic.
<code>tcpflags([val])</code>	If <code>val</code> is specified, return true if the specified flags are set. If <code>val</code> is omitted, return a string representation of the TCP flags.
<code>ip2num(ip)</code>	Convert an IP address to a number.
<code>ip2hex(ip)</code>	Convert an IPv4 address to hex.
<code>ip2str(ip)</code>	Convert an IPv4 address to string.
<code>ip62str(ip)</code>	Convert an IPv6 address to string.
<code>ip6compress(ip)</code>	Compress an IPv6 address.
<code>ip6expand(ip[,trim])</code>	Expand an IPv6 address. If <code>trim</code> is different from 0, remove leading zeros.
<code>ip2mask(ip)</code>	Convert an IP address to a network mask (int).
<code>mask2ip(m)</code>	Convert a network mask (int) to an IPv4 address (int).
<code>mask2ipstr(m)</code>	Convert a network mask (int) to an IPv4 address (string).
<code>mask2ip6(m)</code>	Convert a network mask (int) to an IPv6 address (int).
<code>mask2ip6str(m)</code>	Convert a network mask (int) to an IPv6 address (string).
<code>ipinnet(ip,net[,mask])</code>	Test whether an IP address belongs to a given network.
<code>ipinrange(ip,low,high)</code>	Test whether an IP address lies between two addresses.
<code>localtime(t)</code>	Convert UNIX timestamp to string (localtime).
<code>utc(t)</code>	Convert UNIX timestamp to string (UTC).
<code>timestamp(t)</code>	Convert date to UNIX timestamp.
<code>t2split(val,sep [,num[,osep]])</code>	Split values according to <code>sep</code> . If <code>num</code> is omitted or 0, <code>val</code> is split into <code>osep</code> separated columns. If <code>num > 0</code> , return the <code>num</code> repetition. If <code>num < 0</code> , return the <code>num</code> repetition from the end, e.g., -1 for last element. Multiple <code>num</code> can be specified, e.g., "1;-1;2". Output separator <code>osep</code> , defaults to OFS.
<code>splitc(val[,num[,osep]])</code>	Split compound values. Alias for <code>t2split(val, "_", num, osep)</code> .
<code>splitr(val[,num[,osep]])</code>	Split repetitive values. Alias for <code>t2split(val, ";", num, osep)</code> .
<code>valcontains(val,sep,item)</code>	Return true if one item of <code>val</code> split by <code>sep</code> is equal to <code>item</code> .
<code>cvalcontains(val,item)</code>	Alias for <code>valcontains(val, "_", item)</code> .
<code>rvalcontains(val,item)</code>	Alias for <code>valcontains(val, ";", item)</code> .
<code>strisempty(val)</code>	Return true if <code>val</code> is an empty string.
<code>streq(val1,val2)</code>	Return true if <code>val1</code> is equal to <code>val2</code> .

Function	Description
<code>strneq(val1, val2)</code>	Return true if <code>val1</code> and <code>val2</code> are not equal.
<code>hasprefix(val, pre)</code>	Return true if <code>val</code> begins with the prefix <code>pre</code> .
<code>hassuffix(val, suf)</code>	Return true if <code>val</code> finished with the suffix <code>suf</code> .
<code>contains(val, txt)</code>	Return true if <code>val</code> contains the substring <code>txt</code> .
<code>not(q)</code>	Return the logical negation of a query <code>q</code> . This function must be used to keep the header when negating a query.
<code>bfeq(val1, val2)</code>	Return true if the hexadecimal numbers <code>val1</code> and <code>val2</code> are equal.
<code>bitsallset(val, mask)</code>	Return true if all the bits set in <code>mask</code> are also set in <code>val</code> .
<code>bitsanyset(val, mask)</code>	Return true if one of the bits set in <code>mask</code> is also set in <code>val</code> .
<code>isfloat(v)</code>	Return true if <code>v</code> is a floating point number.
<code>isint(v)</code>	Return true if <code>v</code> is an integer.
<code>isnum(v)</code>	Return true if <code>v</code> is a number (signed, unsigned or floating point).
<code>isuint(v)</code>	Return true if <code>v</code> is an unsigned integer.
<code>isip(v)</code>	Return true if <code>v</code> is an IPv4 address in hexadecimal, numerical or dotted decimal notation.
<code>isip6(v)</code>	Return true if <code>v</code> is an IPv6 address.
<code>isiphex(v)</code>	Return true if <code>v</code> is an IPv4 address in hexadecimal notation.
<code>isipnum(v)</code>	Return true if <code>v</code> is an IPv4 address in numerical (int) notation.
<code>isipstr(v)</code>	Return true if <code>v</code> is an IPv4 address in dotted decimal notation.
<code>join(a, s)</code>	Convert an array to string, separating each value with <code>s</code> .
<code>unquote(s)</code>	Remove leading and trailing quotes from a string.
<code>chomp(s)</code>	Remove leading and trailing spaces from a string.
<code>strip(s)</code>	Remove leading and trailing spaces from a string.
<code>lstrip(s)</code>	Remove leading spaces from a string.
<code>rstrip(s)</code>	Remove trailing spaces from a string.
<code>mean(c)</code>	Compute the mean value of a column <code>c</code> . The result can be accessed with <code>get_mean(c)</code> or printed with <code>print_mean([c])</code> .
<code>min(c)</code>	Keep track of the min value of a column <code>c</code> . The result can be accessed with <code>get_min(c)</code> or printed with <code>print_min([c])</code> .
<code>max(c)</code>	Keep track of the max value of a column <code>c</code> . The result can be accessed with <code>get_max(c)</code> or printed with <code>print_max([c])</code> .
<code>abs(v)</code>	Return the absolute value of <code>v</code> .
<code>min2(a, b)</code>	Return the minimum value between <code>a</code> and <code>b</code> .
<code>min3(a, b, c)</code>	Return the minimum value between <code>a</code> , <code>b</code> and <code>c</code> .
<code>max2(a, b)</code>	Return the maximum value between <code>a</code> and <code>b</code> .
<code>max3(a, b, c)</code>	Return the maximum value between <code>a</code> , <code>b</code> and <code>c</code> .

Function	Description
<code>aggr(fields[,val[,num]])</code>	<p>Perform aggregation of <code>fields</code> and store the sum of <code>val</code>. <code>fields</code> and <code>val</code> can be tab separated lists of fields, e.g., <code>\$srcIP4\t\$dstIP4</code>. Results are sorted according to the first value of <code>val</code>. If <code>val</code> is omitted, the empty string or equal to "flows" or "packets" (case insensitive), count the number of records (flows or packets). If <code>num</code> is omitted or 0, return the full list. If <code>num > 0</code> return the top <code>num</code> results. If <code>num < 0</code> return the bottom <code>num</code> results.</p>
<code>aggrrep(fields[,val[,num[,ign_e[,sep]]])</code>	<p>Perform aggregation of the repetitive <code>fields</code> and store the sum of <code>val</code>. <code>val</code> can be a tab separated lists of fields, e.g., <code>\$numBytesSnt\t\$numPktsSnt</code>. Results are sorted according to the first value of <code>val</code>. If <code>val</code> is omitted, the empty string or equal to "flows" or "packets" (case insensitive), count the number of records (flows or packets). If <code>num</code> is omitted or 0, return the full list. If <code>num > 0</code> return the top <code>num</code> results. If <code>num < 0</code> return the bottom <code>num</code> results. If <code>ign_e</code> is omitted or 0, consider all values, otherwise ignore empty values. <code>sep</code> can be used to change the separator character (default: ";").</p>
<code>t2rsort(col[,num[,type]])</code>	<p>Sort the file in reverse order according to <code>col</code>. (Multiple column numbers can be specified by using ";" as separator, e.g., <code>1 ";" 2</code>) If <code>num</code> is omitted or 0, return the full list. If <code>num > 0</code> return the top <code>num</code> results. If <code>num < 0</code> return the bottom <code>num</code> results. <code>type</code> can be used to specify the type of data to sort: "ip", "num" or "str" (default is based on the first matching record).</p>
<code>t2sort(col[,num[,type[,rev]])</code>	<p>Sort the file according to <code>col</code>. (Multiple column numbers can be specified by using ";" as separator, e.g., <code>1 ";" 2</code>) If <code>num</code> is omitted or 0, return the full list. If <code>num > 0</code> return the top <code>num</code> results. If <code>num < 0</code> return the bottom <code>num</code> results. <code>type</code> can be used to specify the type of data to sort: "ip", "num" or "str" (default is based on the first matching record). If <code>rev > 0</code>, sort in reverse order (alternatively, use the <code>t2rsort()</code> function).</p>
<code>t2whois(ip[,o_opt])</code>	<p>Wrapper to call <code>t2whois</code> from <code>tawk</code>. <code>ip</code> must be a valid IPv4/6 address. <code>o_opt</code> is passed verbatim to <code>t2whois -o option</code> (run <code>t2whois -L</code> for more details).</p>
<code>wildcard(expr)</code>	<p>Print all columns whose name matches the regular expression <code>expr</code>.</p>

Function	Description
	If <code>expr</code> is preceded by an exclamation mark, return all columns whose name does NOT match <code>expr</code> .
<code>hnum(num[,mode[,suffix]])</code>	Convert the number <code>num</code> to its human readable form.
<code>json([s])</code>	Convert the string <code>s</code> to JSON. The first record is used as column names. If <code>s</code> is omitted, convert the entire row.
<code>texscape(s)</code>	Escape the string <code>s</code> to make it LaTeX compatible.
<code>bitshift(n,t[,d[,b]])</code>	Shift a byte or of a list of bytes <code>n</code> to the left or right by a given number of bits <code>t</code> . To shift to the left, set <code>d</code> to 0 (default), to shift to the right set <code>d</code> \neq 0. Set <code>b</code> to 16 to force interpretation as hexadecimal, e.g., interpret 45 as 69 (0x45) instead of 45.
<code>nibble_swap(n[,b])</code>	Swap the nibbles of a byte or of a list of bytes <code>n</code> . Set <code>b</code> to 16 to force interpretation as hexadecimal, e.g., interpret 45 as 69 (0x45) instead of 45.
<code>tobits(u, [b])</code>	Convert the unsigned integer <code>u</code> to its binary representation. Set <code>b</code> to 16 to force interpretation as hexadecimal, e.g., interpret 45 as 69 (0x45) instead of 45.
<code>base64(s)</code>	Encode a string <code>s</code> as base64.
<code>base64d(s)</code>	Decode a base64 encoded string <code>s</code> .
<code>urldecode(url)</code>	Decode the encoded URL <code>url</code> .
<code>printerr(s)</code>	Print the string <code>s</code> in red with an added newline.
<code>diff(file[,mode])</code>	Compare two files (<code>file</code> and the input), and print the name and number of the columns which differ. The <code>mode</code> parameter can be used to control the format of the output.
<code>ffsplit([s[,k[,h]])</code>	Split the input file into smaller more manageable files. The files to create can be specified as argument to the function (one comma separated string). If no argument is specified, create one file per column whose name ends with <code>Stat</code> , e.g., <code>dnsStat</code> , and one for <code>pwXType(pw)</code> . If <code>k > 0</code> , then only print relevant fields and those controlled by <code>h</code> , a comma separated list of fields to keep in each file, e.g., <code>"srcIP,dstIP"</code> .
<code>flow([f])</code>	Return all flows whose index appears in <code>f</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>flow("1-3")</code> . If <code>f</code> is omitted, return the flow index.
<code>packet([p])</code>	Return all packets whose number appears in <code>p</code> (comma or semicolon separated). Ranges may also be specified using a dash, e.g., <code>packet("1-3")</code> . If <code>p</code> is omitted, return the packet number.
<code>follow_stream(f[,of[,d[,pf[,r[,nc]]]])</code>	Return the payload of the flow with index <code>f</code> .

Function	Description
	<p><code>of</code> can be used to change the output format [default: 0]:</p> <ul style="list-style-type: none"> 0: Payload only, 1: prefix each payload with packet/flow info, 2: JSON, 3: Reconstruct (pipe the output to <code>xxd -p -r</code> to reproduce the binary file). <p><code>d</code> can be used to only extract a specific direction ("A" or "B") [default: "" (A and B)].</p> <p><code>pf</code> can be used to change the payload format [default: 0]:</p> <ul style="list-style-type: none"> 0: ASCII, 1: Hexdump, 2: Raw/Binary, 3: Base64. <p><code>r</code> can be used to prevent the analysis of TCP sequence numbers (no TCP reassembly and reordering).</p> <p><code>nc</code> can be used to print the data without colors.</p>
<code>shark(q)</code>	Query flow files according to Wireshark's syntax.

1.8 Examples

Collection of examples using `tawk` functions:

Function	Description
<code>dnsZT()</code>	Return all flows where a DNS zone transfer was performed.
<code>exeDL([n])</code>	Return the top N EXE downloads.
<code>httpHostsURL([f])</code>	Return all HTTP hosts and a list of the files hosted (sorted alphabetically). If <code>f > 0</code> , print the number of times a URL was requested.
<code>nonstdports()</code>	Return all flows running protocols over non-standard ports.
<code>passivedns()</code>	Extract all DNS server replies from a flow file. The following information is reported for each reply: FirstSeen, LastSeen, Type (A or AAAA), TTL, Query, Answer, Organization, Country, AS number
<code>passwords([val[, num]])</code>	Return information about hosts sending authentication in cleartext. If <code>val</code> is omitted or equal to "flows", count the number of flows. Otherwise, sum up the values of <code>val</code> . If <code>num</code> is omitted or 0, returns the full list. If <code>num > 0</code> return the top <code>num</code> results. If <code>num < 0</code> return the bottom <code>num</code> results.
<code>postQryStr([n])</code>	Return the top N POST requests with query strings.

Function	Description
ssh()	Return the SSH connections.
topDnsA([n])	Return the top N DNS answers.
topDnsIp4([n])	Return the top N DNS answers IPv4 addresses.
topDnsIp6([n])	Return the top N DNS answers IPv6 addresses.
topDnsQ([n])	Return the top N DNS queries.
topHttpMimesST([n])	Return the top HTTP content-type (type/subtype).
topHttpMimesT([n])	Return the top HTTP content-type (type only).
topSLD([n])	Return the top N second-level domains queried (google.com, yahoo.com, ...).
topTLD([n])	Return the top N top-level domains (TLD) queried (.com, .net, ...).

1.9 t2nfdump

Collection of functions for `tawk` allowing access to specific fields using a syntax similar as `nfdump`.

Function	Description
ts()	Start Time — first seen
te()	End Time — last seen
td()	Duration
pr()	Protocol
sa()	Source Address
da()	Destination Address
sap()	Source Address:Port
dap()	Destination Address:Port
sp()	Source Port
dp()	Destination Port
pkt()	Packets — default input
ipkt()	Input Packets
opkt()	Output Packets
byt()	Bytes — default input
ibyt()	Input Bytes
obyt()	Output Bytes
flg()	TCP Flags
mpls1()	MPLS label 1
mpls2()	MPLS label 2
mpls3()	MPLS label 3
mpls4()	MPLS label 4
mpls5()	MPLS label 5
mpls6()	MPLS label 6
mpls7()	MPLS label 7
mpls8()	MPLS label 8
mpls9()	MPLS label 9

Function	Description
<code>mpls10()</code>	MPLS label 10
<code>mpls()</code>	MPLS labels 1–10
<code>bps()</code>	Bits per second
<code>pps()</code>	Packets per second
<code>bpp()</code>	Bytes per package
<code>oline()</code>	nfdump line output format (<code>-o line</code>)
<code>olong()</code>	nfdump long output format (<code>-o long</code>)
<code>oextended()</code>	nfdump extended output format (<code>-o extended</code>)

1.10 t2custom

Copy your own functions in this folder. Refer to Section 1.11 for more details on how to write a tawk function. To have your functions automatically loaded, include them in the file `t2custom/t2custom.load`.

1.11 Writing a tawk Function

- Ideally one function per file (where the filename is the name of the function)
- Private functions are prefixed with an underscore
- Always declare local variables 8 spaces after the function arguments
- Local variables are prefixed with an underscore
- Use uppercase letters and two leading and two trailing underscores for global variables
- Include all referenced functions
- Files should be structured as follows:

```
#!/usr/bin/env awk
#
# Function description
#
# Parameters:
# - arg1: description
# - arg2: description (optional)
#
# Dependencies:
# - plugin1
# - plugin2 (optional)
#
# Examples:
# - tawk `funcname()` file.txt
# - tawk `{ print funcname() }` file.txt

@include "hdr"
```

```
@include "_validate_col"

function funcname(arg1, arg2, [8 spaces] _locvar1, _locvar2) {
    _locvar1 = _validate_col("colname1;altcolname1", _my_colname1)
    _validate_col("colname2")

    if (hdr()) {
        if (__PRIHDR__) print "header"
    } else {
        print "something", $_locvar1, $colname2
    }
}
```

1.12 Using tawk Within Scripts

To use tawk from within a script:

1. Create a TAWK variable pointing to the script: `TAWK="$T2HOME/scripts/tawk/tawk"`
2. Call tawk as follows: `$TAWK `dport(80)` file.txt`

1.13 Using tawk With Non-Tranalyzer Files

tawk can also be used with files which were not produced by Tranalyzer.

- The input field separator can be specified with the `-F` option, e.g., `tawk -F ',' 'program' file.csv`
- The row listing the column names, can start with any character specified with the `-s` option, e.g., `tawk -s '#' 'program' file.txt`
- All the column names must not be equal to a function name
- Valid column names must start with a letter (a-z, A-Z) and can be followed by any number of alphanumeric characters or underscores
- If no column names are present, use the `-t` option to prevent tawk from trying to validate the column names.
- If the column names are different from those used by Tranalyzer, refer to Section 1.13.1.

1.13.1 Mapping External Column Names to Tranalyzer Column Names

If the column names are different from those used by Tranalyzer, a mapping between the different names can be made in the file `my_vars`. The format of the file is as follows:

```
BEGIN {
    _my_srcIP = non_t2_name_for_srcIP
    _my_dstIP = non_t2_name_for_dstIP
    ...
}
```

Once edited, run tawk with the `-i $T2HOME/scripts/tawk/my_vars` option and the external column names will be automatically used by tawk functions, such as `tuple2()`. For more details, refer to the `my_vars` file.

1.13.2 Using tawk with Bro/Zeek Files

To use tawk with Bro/Zeek log files, use one of `--bro` or `--zeek` option:

```
tawk -bro '{ program }' file.log tawk -zeek '{ program }' file.log
```

1.14 Awk Cheat Sheet

- Tranalyzer flow files default field separator is `'\t'`:
 - **Always** use `awk -F'\t'` (or `awkf/tawk`) when working with flow files.
- Load libraries, e.g., tawk functions, with `-i`: `awk -i file.awk 'program' file.txt`
- Always use `strtonum` with hex numbers (bitfields)
- Awk indices start at 1
- Using tawk is recommended.

1.14.1 Useful Variables

- `$0`: entire line
- `$1, $2, ..., $NF`: column 1, 2, ...
- `FS`: field separator
- `OFS`: output field separator
- `ORS`: output record separator
- `NF`: number of fields (columns)
- `NR`: record (line) number
- `FNR`: record (line) number relative to the current file
- `FILENAME`: name of current file
- To use external variables, use the `-v` option, e.g., `awk -v name="value" '{ print name }' file.txt.`

1.14.2 Awk Program Structure

```
awk -F'\t' -i min -v OFS='\t' -v h="$(hostname)" `
  BEGIN { a = 0; b = 0; }           # Called once at the beginning
  /^A/ { a++ }                     # Called for every row starting with char A
  /^B/ { b++ }                     # Called for every row starting with char B
      { c++ }                       # Called for every row
  END { print h, min(a, b), c }    # Called once at the end
' file.txt
```

1.15 Awk Templates

- Print the whole line:

```
- tawk '{ print }' file.txt
- tawk '{ print $0 }' file.txt
- tawk 'FILTER' file.txt
- tawk 'FILTER { print }' file.txt
- tawk 'FILTER { print $0 }' file.txt
```

- Print selected columns only:

```
- tawk '{ print $srcIP4, $dstIP4 }' file.txt
- tawk '{ print $1, $2 }' file.txt
- tawk '{ print $4 "\t" $6 }' file.txt
- tawk '{
  for (i = 6; i < NF; i++) {
    printf "%s\t", $i
  }
  printf "%s\n", $NF
}' file.txt
```

- Keep the column names:

```
- tawk 'hdr() || FILTER' file.txt
- awkf 'NR == 1 || FILTER' file.txt
- awkf '/^%/ || FILTER' file.txt
- awkf '/^[[:space:]]*[[[:alpha:]]][[:alnum:]]_*/ || FILTER' file.txt
```

- Skip the column names:

```
- tawk '!hdr() && FILTER' file.txt
- awkf 'NR > 1 && FILTER' file.txt
- awkf '!/^%/ && FILTER' file.txt
- awkf '!/^[[:space:]]*[[[:alpha:]]][[:alnum:]]_*/ && FILTER' file.txt
```

- Bitfields and hexadecimal numbers:

```
- tawk 'bfeq($3,0)' file.txt
- awkf 'strtonum($3) == 0' file.txt
- tawk 'bitsanyset($3,1)' file.txt
- tawk 'bitsallset($3,0x81)' file.txt
- awkf 'and(strtonum($3), 0x1)' file.txt
```

- Split compound values:

```
- tawk '{ print splitc($16, 1) }' file.txt # first element
```



```

- tawk '{ print splitc($16, -1) }' file.txt # last element
- awkf '{ split($16, A, "_"); print A[1] }' file.txt
- awkf '{ n = split($16, A, "_"); print A[n] }' file.txt # last element
- tawk '{ print splitc($16) }' file.txt
- awkf '{ split($16, A, "_"); for (i=1;i<=length(A);i++) print A[i] }' file.txt

```

- Split repetitive values:

```

- tawk '{ print splitr($16, 3) }' file.txt # third repetition
- tawk '{ print splitr($16, -2) }' file.txt # second to last repetition
- awkf '{ split($16, A, ";"); print A[3] }' file.txt
- awkf '{ n = split($16, A, ";"); print A[n] }' file.txt # last repetition
- tawk '{ print splitr($16) }' file.txt
- awkf '{ split($16, A, ";"); for (i=1;i<=length(A);i++) print A[i] }' file.txt

```

- Filter out empty strings:

```

- tawk '!strisempty($4)' file.txt
- awkf '!(length($4) == 0 || $4 == "\"\"")' file.txt

```

- Compare strings (case sensitive):

```

- tawk 'streq($3,$4)' file.txt
- awkf '$3 == $4' file.txt
- awkf '$3 == \"text\"' file.txt

```

- Compare strings (case insensitive):

```

- tawk 'streqi($3,$4)' file.txt
- awkf 'tolower($3) == tolower($4)' file.txt

```

- Use regular expressions on specific columns:

```

- awkf '$8 ~ /^192.168.1.[0-9]{1,3}$/' file.txt # print matching rows
- awkf '$8 !~ /^192.168.1.[0-9]{1,3}$/' file.txt # print non-matching rows

```

- Use column names in awk:

```

- tawk '{ print $srcIP4, $dstIP4 }' file.txt
- awkf `
NR == 1 {
  for (i = 1; i <= NF; i++) {
    if ($i == "srcIP4") srcIP4 = i
    else if ($i == "dstIP4") dstIP4 = i
  }
  if (srcIP4 == 0 || dstIP4 == 0) {
    print "No column with name srcIP4 and/or dstIP4"
    exit
  }
}
`

```

```

    }
  }
  NR > 1 {
    print $srcIP4, $dstIP4
  }
' file.txt
- awkf `
  NR == 1 {
    for (i = 1; i <= NF; i++) {
      col[$i] = i
    }
  }
  NR > 1 {
    print $col["srcIP4"], $col["dstIP4"];
  }
' file.txt

```

1.16 Examples

1. Pivoting (variant 1):

- (a) First extract an attribute of interest, e.g., an unresolved IP address in the `Host :` field of the HTTP header:

```
tawk 'aggr($httpHosts)' FILE_flows.txt | tawk '{ print unquote($1); exit }'
```

- (b) Then, put the result of the last command in the `badguy` variable and use it to extract flows involving this IP:

```
tawk -v badguy="$(!)" 'host(badguy)' FILE_flows.txt
```

2. Pivoting (variant 2):

- (a) First extract an attribute of interest, e.g., an unresolved IP address in the `Host :` field of the HTTP header, and store it into a `badip` variable:

```
badip="$(tawk 'aggr($httpHosts)' FILE_flows.txt | tawk '{ print unquote($1);exit }')"
```

- (b) Then, use the `badip` variable to extract flows involving this IP:

```
tawk -v badguy="$badip" 'host(badguy)' FILE_flows.txt
```

3. Aggregate the number of bytes sent between source and destination addresses (independent of the protocol and port) and output the top 10 results:

```
tawk 'aggr($srcIP4 "\t" $dstIP4, $numBytesSnt, 10)' FILE_flows.txt
```

```
tawk 'aggr(tuple2(), $numBytesSnt "\t" "Flows", 10)' FILE_flows.txt
```

4. Sort the flow file according to the duration (longest flows first) and output the top 5 results:

```
tawk 't2sort(duration, 5)' FILE_flows.txt
```

5. Extract all TCP flows while keeping the header (column names):

```
tawk `hdr() || tcp()` FILE_flows.txt
```

6. Extract all flows whose destination port is between 6000 and 6008 (included):

```
tawk `dport("6000-6008")` FILE_flows.txt
```

7. Extract all flows whose destination port is 53, 80 or 8080:

```
tawk `dport("53;80;8080")` FILE_flows.txt
```

8. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using host or net):

```
tawk `shost("192.168.1.0/24")` FILE_flows.txt
```

```
tawk `snet("192.168.1.0/24")` FILE_flows.txt
```

9. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using ipinrange):

```
tawk `ipinrange($srcIP4, "192.168.1.0", "192.168.1.255")` FILE_flows.txt
```

10. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using ipinnet):

```
tawk `ipinnet($srcIP4, "192.168.1.0", "255.255.255.0")` FILE_flows.txt
```

11. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using ipinnet and a hex mask):

```
tawk `ipinnet($srcIP4, "192.168.1.0", 0xffffffff00)` FILE_flows.txt
```

12. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using ipinnet and the CIDR notation):

```
tawk `ipinnet($srcIP4, "192.168.1.0/24")` FILE_flows.txt
```

13. Extract all flows whose source IP is in subnet 192.168.1.0/24 (using ipinnet and a CIDR mask):

```
tawk `ipinnet($srcIP4, "192.168.1.0", 24)` FILE_flows.txt
```

For more examples, refer to `tawk -d` option, e.g., `tawk -d aggr`, where every function is documented and comes with a set of examples. The complete documentation can be consulted by running `tawk -d all`.

1.17 FAQ

1.17.1 Can I use tawk with non Tranalyzer files?

Yes, refer to Section 1.13.

1.17.2 Can I use tawk functions with non Tranalyzer column names?

Yes, edit the `my_vars` file and load it using `-i $T2HOME/scripts/tawk/my_vars` option. Refer to Section 1.13.1 for more details.

1.17.3 Can I use tawk with files without column names?

Yes, use the `-t` option to prevent tawk from trying to validate the column names.

1.17.4 The row listing the column names start with a '#' instead of a '%'. Can I still use tawk?

Yes, use the `-s` option to specify the first character, e.g., `tawk -s '#' 'program'`

1.17.5 Can I process Bro/Zeek log files with tawk?

Yes, use the `--zeek` option.

1.17.6 Can I process a CSV (Comma Separated Value) file with tawk?

The simplest way to process CSV files is to use the `--csv` option. This sets the input and output separators to a comma and considers the first row to be the column names.

```
tawk --csv 'program' file.csv
```

Alternatively, the input field separator can be changed with the `-F` option and the output separator with `-O ','` or `-v OFS=','`. Note that tawk expects the column names to be the last row starting with a '%'. This can be changed with the `-s` and `-N` options (Section 1.5).

```
tawk -F ',' -v OFS=',' -s "" -N 1 'program' file.csv
```

1.17.7 Can I produce a CSV (Comma Separated Value) file from tawk?

The output field separator (OFS) can be changed with the `-O 'fs'` or `-v OFS='fs'` option. To produce a CSV file, run tawk as follows: `tawk -O ',' 'program' file.txt` or `tawk -v OFS=',' 'program' file.txt`

1.17.8 Can I write my tawk programs in a file instead of the command line?

Yes, copy the program (without the single quotes) in a file, e.g., `prog.txt` and run it as follows:

```
tawk -f prog.txt file.txt
```

1.17.9 Can I still use column names if I pipe data into tawk?

Yes, you can specify a file containing the column names with the `-I` option as follows:

```
cat file.txt | tawk -I colnames.txt 'program'
```

1.17.10 Can I use tawk if the row with the column names does not start with a special character?

Yes, you can specify the empty character with `-s ""`. Refer to Section 1.5 for more details.

1.17.11 I get a list of syntax errors from gawk... What is the problem?

The name of the columns is used to create variable names. If it contains forbidden characters, then an error similar to the following is reported.

```
gawk: /tmp/fileBndhdf:3: col-name = 3
gawk: /tmp/fileBndhdf:3:      ^ syntax error
```

Although tawk will try to replace forbidden characters with underscore, the best practice is to use only alphanumeric characters (A-Z, a-z, 0-9) and underscore as column names. Note that a column name **MUST NOT** start with a number.

1.17.12 Tawk cannot find the column names... What is the problem?

First, make sure the comment char (`-s` option) is correctly set for your file (the default is ``%'`). Second, make sure the column names do not contain forbidden characters, i.e., use only alphanumeric and underscore and do not start with a number. If the row with column names is not the last one to start with the separator character, then specify the line number with the `-N` option as follows: `tawk -N 3'` or `tawk -s '`#' -N 2`. Refer to Section 1.5 for more details.

1.17.13 How to make tawk faster?

Tawk tries to validate the column names by ensuring that no column names is equal to a function name and that all column names used in the program exist. This verification process is quite slow and can easily be disabled by using the `-t` option.

1.17.14 Wireshark refuses to open PCAP files generated with tawk `-k` option...

If Wireshark displays the message `Couldn't run /usr/bin/dumpcap in child process: Permission Denied.`, then this means that your user does not belong to the `wireshark` group. To fix this issue, simply run the following command `sudo gpasswd -a YOUR_USERNAME wireshark` (you will then need to log off and on again).

1.17.15 Tawk reports errors similar to `free(): double free detected in tcache 2`

Tawk uses `gawk -M` option to handle IPv6 addresses. For some reasons, this option is regularly affected by bugs... If you do not need IPv6 support, you can simply comment out line 653 in `tawk`:

```
OPTS=(
  #-M -v PREC=256          # <-- Add the leading sharp ('#') here
  -v __PRIHDR__=$PRIHDR
  -v __UNAME__="$ (uname) "
)
```