dnsop Internet-Draft Intended status: Standards Track Expires: June 9, 2017

J. Dickinson J. Haque S. Dickinson Sinodun IT T. Manderson J. Bond ICANN December 6, 2016

C-DNS: A DNS Packet Capture Format draft-ietf-dnsop-dns-capture-format-00

Abstract

This document describes a data representation for collections of DNS messages. The format is designed for efficient storage and transmission of large packet captures of DNS traffic; it attempts to minimize the size of such packet capture files but retain the full DNS message contents along with the most useful transport meta data. It is intended to assist with the development of DNS traffic monitoring applications.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at http://datatracker.ietf.org/drafts/current/.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on June 9, 2017.

Copyright Notice

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (http://trustee.ietf.org/license-info) in effect on the date of

Dickinson, et al. Expires June 9, 2017

[Page 1]

publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

$\underline{1}$. Introduction	<u>3</u>
2. Requirements Terminology	<u>4</u>
$\underline{3}$. Data Collection Use Cases	<u>4</u>
$\underline{4}$. Design Considerations	<u>6</u>
5. C-DNS conceptual overview	<u>8</u>
<u>6</u> . Choice of CBOR	<u>8</u>
$\underline{7}$. C-DNS CBOR format	<u>8</u>
<u>7.1</u> . CDDL definition	<u>8</u>
<u>7.2</u> . Format overview	<u>9</u>
<u>7.3</u> . File header contents	<u>9</u>
<u>7.4</u> . File preamble contents	<u>10</u>
7.5. Configuration contents	<u>10</u>
<u>7.6</u> . Block contents	<u>11</u>
<u>7.7</u> . Block preamble map	<u>12</u>
<u>7.8</u> . Block statistics	<u>12</u>
<u>7.9</u> . Block table map	<u>13</u>
<u>7.10</u> . IP address table	<u>13</u>
<u>7.11</u> . Class/Type table	<u>13</u>
<u>7.12</u> . Name/RDATA table	<u>14</u>
<u>7.13</u> . Query Signature table	<u>14</u>
<u>7.14</u> . Question table	<u>16</u>
7.15. Resource Record (RR) table	<u>17</u>
<u>7.16</u> . Question list table	<u>17</u>
7.17. Resource Record list table	<u>17</u>
7.18. Query/Response data	<u>18</u>
7.19. Address Event counts	<u>20</u>
<u>8</u> . C-DNS to PCAP	<u>21</u>
<u>8.1</u> . Name Compression	<u>22</u>
9. Data Collection	<u>23</u>
<u>9.1</u> . Matching algorithm	<u>23</u>
<u>9.2</u> . Message identifiers	<u>24</u>
<u>9.2.1</u> . Primary ID (required)	<u>24</u>
<u>9.2.2</u> . Secondary ID (optional)	<u>24</u>
<u>9.3</u> . Algorithm Parameters	<u>24</u>
<u>9.4</u> . Algorithm Requirements	<u>24</u>
<u>9.5</u> . Algorithm Limitations	<u>25</u>
<u>9.6</u> . Workspace	<u>25</u>
<u>9.7</u> . Output	<u>25</u>
9.8. Post Processing	<u>25</u>

Dickinson, et al. Expires June 9, 2017 [Page 2]

<u>10</u> . IANA Considerations \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots	. <u>25</u>
<u>11</u> . Security Considerations	. <u>25</u>
<u>12</u> . Acknowledgements	. <u>25</u>
<u>13</u> . Changelog	. <u>26</u>
<u>14</u> . References	. <u>26</u>
<u>14.1</u> . Normative References	. <u>26</u>
<u>14.2</u> . Informative References	. <u>27</u>
<u>14.3</u> . URIS	. <u>28</u>
Appendix A. CDDL	. <u>28</u>
Appendix B. DNS Name compression example	. <u>33</u>
<u>B.1</u> . NSD compression algorithm	. <u>34</u>
B.2. Knot Authoritative compression algorithm	. <u>35</u>
B.3. Observed differences	. <u>35</u>
Appendix C. Comparison of Binary Formats	. <u>35</u>
Appendix D. Sample data on the C-DNS format	. <u>36</u>
D.1. Comparison to full PCAPS	. <u>36</u>
D.2. Block size choices	. <u>36</u>
D.3. Blocking vs more simple output	. <u>36</u>
Authors' Addresses	. <u>37</u>

1. Introduction

There has long been a need to collect DNS queries and responses on authoritative and recursive name servers for monitoring and analysis. This data is used in a number of ways including traffic monitoring, analyzing network attacks and DITL [dit1].

A wide variety of tools already exist to facilitate the collection of DNS traffic data. DSC [dsc], packetq [packetq], dnscap [dnscap] and dnstap [dnstap]. However, there is no standard exchange format for large DNS packet captures and PCAP [pcap] or PCAP-NG [pcapng] are typically used in practice. Such file formats can contain much additional information not directly pertinent to DNS traffic analysis which unnecessarily increases the capture file size.

There has also been work on using other text based formats to describe DNS packets [<u>I-D.daley-dnsxml</u>], [<u>I-D.hoffman-dns-in-json</u>] but these are largely aimed at producing convenient representations of single messages.

Many DNS operators may receive 100's of thousands of queries per second on a single name server instance so a mechanism to minimize the storage size (and therefore upload overhead) of the data collected is highly desirable.

This documents focusses on the problem of capturing and storing large packet capture files of DNS traffic. with the following goals in mind:

Dickinson, et al. Expires June 9, 2017 [Page 3]

- o Minimize the file size for storage and transmission
- o Minimizing the overhead of producing the packet capture file and the cost of any further (general purpose) compression of the file to minimise the size

This document contains

- o A discussion of the some common use cases in which such DNS data is collected. See <u>Section 3</u>.
- o A discussion of the major design considerations in developing an efficient data representation for collections of DNS messages. See Section 4.
- o A definition of a CBOR [RFC7049] representation of a collection of DNS messages. This will be referred to as the C-DNS format (Compacted-DNS). See Section 7.
- o Notes on converting C-DNS back to PCAP format. See Section 8.
- o Some high level implementation considerations for applications designed to produce C-DNS, e.g. a query response matching algorithm. See Section 9.

2. Requirements Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC2119].

3. Data Collection Use Cases

In an ideal world it would be optimal to collect full packet captures of all packets going in or out of a name server. However, there are several design choices or other limitations that are common to many DNS installations and operators.

- o Servers are hosted in a variety of situations
 - * Operator self hosted servers
 - * Third party hosting (including multiple third parties)
 - * Third party hardware (including multiple third parties)
- o Data is collected under different conditions

Dickinson, et al. Expires June 9, 2017 [Page 4]

- * On well provisioned servers running in a steady state.
- * On heavily loaded servers
- * On virtualized servers
- * On servers that are under attack
- * On servers that are unwitting intermediaries in attacks
- o Traffic can be collected via a variety of mechanisms
 - * On the same hardware as the name server itself
 - * Using a network tap to listen in from another server
 - * Using port mirroring to listen in from another server
- o The capabilities of data collection (and upload) networks vary
 - * Out-of-band networks with the same capacity as the in-band network
 - * Out-of-band networks with less capacity than the in-band network
 - * Everything on the in-band network

Clearly, there is a wide range of use cases from very limited data collection environments (third party hardware, servers that are under attack, packet capture on the name server itself and no out-of-band network) to 'limitless' environments (self hosted, well provisioned servers, using a network tap or port mirroring with an out-of-band networks with the same capacity as the in-band network). In the former, it is unfeasible to reliably collect full PCAPS especially if the server is under attack. In the latter case, collection of full PCAPs may be reasonable.

As a result of these restrictions the data format discussed below was designed with the most limited use case in mind such that

- o Data collection will occur on the same hardware as the name server itself.
- o Collected data will be stored on the same hardware as the name server itself, at least temporarily.

Dickinson, et al. Expires June 9, 2017 [Page 5]

- o Collected data being returned to some central analysis system will use the same network interface as the DNS queries and responses.
- o There are multiple third party servers involved.

and therefore minimal storage size of the capture files is a major factor.

Another consideration for any application that records DNS traffic is that the running of the name server software and the transmission of DNS queries and responses is the most important job of a name server. Any data collection system co-located with the name server will need to be intelligent enough to carefully manage its CPU, disk, memory and network utilization. Hence this use case benefits from a format that has a relatively low overhead to produce and minimizes the requirement for further potentially costly compression.

However, it was also essential that interoperability with less restricted infrastructure was maintained. In particular it is highly desirable that the resulting collection format should facilitate the re-creation of common formats (such as PCAPs) that are as close to the original as is realistic given the restrictions above.

<u>4</u>. Design Considerations

This section presents some of the major design considerations used in the development of the C-DNS format.

o The basic unit of data is a combined DNS Query and the associated Response (a 'Q/R data item'). The same structure will be used for unmatched queries and responses. Queries without responses will be captured omitting the Response data. Responses without gueries will be captured omitting the Query data (but using the Query section from the Response, if present, as an identifying QNAME).

Rationale: A Query and Response represents the basic level of a clients interaction with the server. Also, combining the Query and Response into one item lowers storage requirements due to commonality in the data in most cases.

o Each Q/R data item will comprise a default Q/R data description and a set of optional sections. Inclusion of optional sections shall be configurable.

Rationale: Different users will have different requirements for data to be available for analysis. Users with minimal requirements should not have to pay the cost of recording full data, however this will limit the ability to reconstruct PCAPS. For example omitting the

Dickinson, et al. Expires June 9, 2017 [Page 6]

Resource Records from a Response will reduce the files size, and in principle responses can be synthesized if there is enough context.

 Multiple Q/R items will be collected into blocks in the format.
 Common data in a block will be abstracted and referenced from individual Q/R items by indexing. The maximum number of Q/R items in a block will be configurable.

Rationale: This blocking and indexing provides a significant reduction in the volume of file data generated. Whilst introducing complexity it provides compression of the data that makes use of knowledge of the DNS packet structure.

[TODO: Further discussion on commonality between DNS packets e.g.

- o common query signatures
- o for the authoritative case there are a finite set of valid responses and much commonality in NXDOMAIN responses]

It is anticipated that the files produced will be subject to further compression using general purpose compression tools. Measurements show that blocking significantly reduces the CPU required to perform such strong compression. See <u>Appendix D</u>.

 Meta-data about other packets received should also be included in each block. For example counts of malformed DNS packets and non-DNS packets (e.g. ICMP, TCP resets) sent to the server are of interest.

It should be noted that any structured capture format that does not capture the DNS payload byte for byte will likely be limited to some extent in that it cannot represent 'malformed' DNS packets. Only those packets that can be transformed reasonably into the structured format can be represented by it. So if a query is malformed this will lead to the (well formed) DNS responses with error code FORMERR appearing as 'unmatched'.

[TODO: Need further discussion of well-formed vs malformed packets and how name servers view this definition.]

[TODO: Need to develop optional representation of malformed packets within CBOR and what this means for packet matching. This may influence which fields are optional in the rest of the representation.]

Packets such as those described above can be separately recorded in a PCAP file for later analysis.

Dickinson, et al. Expires June 9, 2017 [Page 7]

5. C-DNS conceptual overview

The following figures show purely schematic representations of the C-DNS format to convey the high-level structure of the C-DNS format. Section 7 provides a detailed discussion of the CBOR representation and individual elements.

Figure showing the C-DNS format (PNG) [1] Figure showing the C-DNS format (SVG) [2] Figure showing the Q/R data item and Block tables format (PNG) [3] Figure showing the Q/R data item and Block tables format (SVG) [4]

6. Choice of CBOR

This document presents a detailed format description using CBOR, the Concise Binary Object Representation defined in [RFC7049].

The choice of CBOR was made taking a number of factors into account.

- o CBOR is a binary representation, and so economical in storage space.
- o Other similar representations were investigated, and whilst all had attractive features, none had a significant advantage over CBOR. See Appendix C and Appendix D - for some discussion of this.
- o CBOR is an IETF Standard and familiar to IETF participants, and being based on the successful JSON text format, requires very little familiarization for those in the wider industry.
- o CBOR can also be easily converted to JSON for debugging and other human inspection requirements.
- o CBOR data schemas can be described using CDDL [I-D.greevenbosch-appsawg-cbor-cddl].

7. C-DNS CBOR format

7.1. CDDL definition

The CDDL definition for the C-DNS format is given in Appendix A.

Dickinson, et al. Expires June 9, 2017 [Page 8]

7.2. Format overview

A C-DNS file begins with a file header containing a file type identifier and preamble. The preamble contains information on the collection settings.

This is followed by a series of data blocks.

A block consists of a block header, containing various tables of common data, and some statistics for the traffic received over the block. The block header is then followed by a list of the Q/R pairs detailing the queries and responses received during the block. The list of Q/R pairs is in turn followed by a list of per-client counts of particular IP events that occurred during collection of the block data.

The exact nature of the DNS data will affect what block size is the best fit, however sample data for a root server indicated that block sizes in the low 1000's give good results. See Appendix D.2 for more details.

If no field type is specified then the field is unsigned.

In the following

- o For all quantities that contain bit flags, bit 0 indicates the least significant bit.
- o Items described as indexes are the index of the data item in the referenced table. Indexes are 1-based. An index value of 0 is reserved to mean not present.

7.3. File header contents

The file header contains the following:

+-----+ | Field | Type | Description +-----+ | File type | Text string | String identifying the file type | | ID File| Map of items| Collection information for the| preamble|| whole file. | File Blocks | Array of | The data blocks | | Blocks | +-----+

Dickinson, et al. Expires June 9, 2017 [Page 9]

Internet-Draft C-DNS: A DNS Packet Capture Format December 2016

<u>7.4</u>. File preamble contents

The file preamble contains the following:

+		+	•+
 +	Field	Туре	Description
	Format version	Unsigned 	Indicates version of format used in file.
	Configuration 	Map of items	The collection configuration.
	Generator ID	Text string 	String identifying the collection program. Optional.
 	Host ID	Text string 	String identifying the collecting host. Blank if converting an existing PCAP file. Optional.

<u>7.5</u>. Configuration contents

The collection configuration contains the following items. All are optional.

ш.			
	Field	Туре	Description
	Query timeout 	Unsigned 	To be matched with a query, a response must arrive within this number of seconds.
	Skew timeout 	Unsigned 	The network stack may report a response before the corresponding query. A response is not considered to be missing a query until after this many micro- seconds.
	Snap length 	Unsigned 	Collect up to this many bytes per packet.
	Promiscuous mode	Unsigned 	1 if promiscuous mode was enabled on the interface, 0 otherwise.
	Interfaces 	Array of text strings	Identifiers of the interfaces used for collection.

Dickinson, et al. Expires June 9, 2017 [Page 10]

	VLAN IDS	Array of unsigned	Identifiers of VLANs selected for collection.
	Filter	Text string	"tcpdump" [<u>pcap</u>] style filter for input.
	Query collection options	Unsigned	Bit flags indicating sections in Query packets to be collected.
			Bit 0. Collect second and subsequent question sections. Bit 1. Collect Answer sections. Bit 2. Collect Authority sections. Bit 3. Collection Additional sections.
	Response collection options	Unsigned	Bit flags indicating sections in Response packets to be collected.
			Bit 0. Collect second and subsequent question sections. Bit 1. Collect Answer sections.
			Bit 2. Collect Authority sections. Bit 3. Collection Additional sections.
	Accept RR types	Array of text strings	A set of RR type names [<u>rrtypes</u>]. If not empty, only the nominated RR types are collected.
	Ignore RR types	Array of text strings	A set of RR type names [<u>rrtypes</u>]. If not empty, all RR types are collected except those listed. If present, this item must be empty if a non-empty list of Accept RR types is present.

<u>7.6</u>. Block contents

Each block contains the following:

Dickinson, et al. Expires June 9, 2017 [Page 11]

+-----+ | Field | Type | Description | Block | Map of items | Overall information for the | preamble | | block. | Statistics about the block. | Block | Map of | statistics | statistics | Block | Map of tables | The tables containing data | tables | referenced by individual Q/R | entries. | Array of Q/Rs | Details of individual Q/R pairs. | | Q/Rs | Address| Array of Address| Per client counts of ICMP| Event| Event counts| messages and TCP resets. | Counts +-----+

7.7. Block preamble map

The block preamble map contains overall information for the block.

+----+ | Field | Type | Description +-----+ | Timestamp | Array of | A timestamp for the earliest record in the | | unsigned | block. The timestamp is specified as a | CBOR array with two elements as in Posix | struct timeval. The first element is an | unsigned integer time_t and the second is | | an unsigned integer number of | microseconds. The latter is always a value | L | between 0 and 999,999.

[TODO: Extend to support pico/nano. Also do this for Time offset and Response delay]

7.8. Block statistics

[TODO: Add block statistics]

Dickinson, et al. Expires June 9, 2017 [Page 12]

7.9. Block table map

The block table map contains the block tables. Each element, or table, is an array. The following tables detail the contents of each block table.

The Present column in the following tables indicates the circumstances when an optional field will be present. A Q/R pair may be:

- o A Query plus a Response.
- o A Query without a Response.
- o A Response without a Query.

Also:

- o A Query and/or a Response may contain an OPT section.
- o A Question may or may not be present. If the Query is available, the Question section of the Query is used. If no Query is available, the Question section of the Response is used. Unless otherwise noted, a Question refers to the first Question in the Question section.

So, for example, a field listed with a Present value of QUERY is present whenever the Q/R pair contains a Query. If the pair contains a Response only, the field will not be present.

7.10. IP address table

This table holds all client and server IP addresses in the block. Each item in the table is a single IP address.

7.11. Class/Type table

This table holds pairs of RR CLASS and TYPE values. Each item in the table is a CBOR map.

Dickinson, et al. Expires June 9, 2017 [Page 13]

```
+----+
| Field | Description |
+----+
| Class | CLASS value. |
| Type | TYPE value. |
+----+
```

[TODO: Can this be optimized? Should a class of IN be inferred if not present?]

7.12. Name/RDATA table

This table holds the contents of all NAME or RDATA items in the block. Each item in the table is the content of a single NAME or RDATA.

+ Field	+ Ty	/pe	+ Description +	+ -
Data 	B) st 	/te ring	The NAME or RDATA contents. NAMEs, and labels within RDATA contents, are in uncompressed label format. +	

7.13. Query Signature table

This table holds elements of the Q/R data that are often common to between different individual Q/R records. Each item in the table is a CBOR map. Each item in the map has an unsigned value and an unsigned key.

The following abbreviations are used in the Present (P) column

- O Q = QUERY
- o A = Always
- o QT = QUESTION
- o Q0 = QUERY, OPT
- o QR = QUERY & RESPONSE
- O R = RESPONSE

+-----| Field | P | Description

Dickinson, et al. Expires June 9, 2017 [Page 14]

+	+	++
Server address	A 	The index in the IP address table of the server IP address.
Server port	A 	The server port.
 Transport flags 	 A 	Bit flags describing the protocol used to service the query. Bit 0 is the least significant bit. Bit 0. Transport type. 0 = UDP, 1 = TCP. Bit 1. IP type. 0 = IPv4, 1 = IPv6.
Q/R signature flags 	 A 	<pre>Bit flags indicating information present in this Q/R pair. Bit 0 is the least significant bit. Bit 0. 1 if a Query is present. Bit 1. 1 if a Response is present. Bit 2. 1 if one or more Question is present. Bit 3. 1 if a Query is present and it has an OPT Resource Record. Bit 4. 1 if a Response is present and it has an OPT Resource Record. Bit 5. 1 if a Response is present but has no Question.</pre>
 Query OPCODE	 Q 	Query OPCODE.
Q/R DNS flags 	A 	<pre> Bit flags with values from the Query and Response DNS flags. Bit 0 is the least significant bit. Flag values are 0 if the Query or Response is not present. Bit 0. Query Checking Disabled (CD) flag. Bit 1. Query Authenticated Data (AD) flag. Bit 2. Query reserved (Z) flag. Bit 3. Query Recursion Available (RA) flag. Bit 4. Query Authoritative Answer (AA) flag. Bit 7. Query DNSSEC answer OK (D0) flag. Bit 9. Response Authenticated Data (AD) flag. Bit 10. Response Recursion Available (RA) flag. Bit 11. Response Recursion Available (RA) flag. Bit 12. Response Recursion Available (RA) flag. Bit 13. Response TrunCation (TC) flag. Bit 14. Response Authoritative Answer (AA) </pre>

Dickinson, et al. Expires June 9, 2017 [Page 15]

1		flag.
 Query RCODE	 Q 	 Query RCODE. If the Query contains OPT, this value incorporates any EXTENDED_RCODE_VALUE.
 Question Class/Type	 QT 	 The index in the Class/Type table of the CLASS and TYPE of the first Question.
 Question QDCOUNT	 QT 	 The QDCOUNT in the Query, or Response if no Query present.
Query ANCOUNT	Q Q	Query ANCOUNT.
Query ARCOUNT	Q Q	Query ARCOUNT.
 Query NSCOUNT	Q Q	Query NSCOUNT.
 Query EDNS version	 QO 	 The Query EDNS version.
 EDNS UDP size	QO 	 The Query EDNS sender's UDO payload size
 Query OPT RDATA	 QO 	 The index in the NAME/RDATA table of the OPT RDATA.
Response RCODE 	R R 	 Response RCODE. If the Response contains OPT, this value incorporates any EXTENDED_RCODE_VALUE.

7.14. Question table

This table holds details on individual Questions in a Question section. Each item in the table is a CBOR map containing a single Question. Each item in the map has an unsigned value and an unsigned key. This data is optionally collected.

Dickinson, et al. Expires June 9, 2017 [Page 16]

+---+
| Field | Description |
+---+
QNAME	The index in the NAME/RDATA table of the QNAME.			
Class/Type	The index in the Class/Type table of the CLASS and			
		TYPE of the Question.		
+---+

7.15. Resource Record (RR) table

This table holds details on individual Resource Records in RR sections. Each item in the table is a CBOR map containing a single Resource Record. This data is optionally collected.

7.16. Question list table

This table holds a list of second and subsequent individual Questions in a Question section. Each item in the table is a CBOR unsigned. This data is optionally collected.

+----+ | Field | Description | +----+ | Question | The index in the Question table of the individual | | | Question. |

7.17. Resource Record list table

This table holds a list of individual Resource Records in a Answer, Authority or Additional section. Each item in the table is a CBOR unsigned. This data is optionally collected.

Dickinson, et al. Expires June 9, 2017 [Page 17]

+	·+· +·	Description	-+- -+
RR 		The index in the Resource Record table of the individual Resource Record.	

7.18. Query/Response data

The block Q/R data is a CBOR array of individual Q/R items. Each item in the array is a CBOR map containing details on the individual Q/R pair.

Note that there is no requirement that the elements of the Q/R array are presented in strict chronological order.

The following abbreviations are used in the Present (P) column

- o Q = QUERY
- o A = Always
- o QT = QUESTION
- O QO = QUERY, OPT
- O QR = QUERY & RESPONSE
- o R = RESPONSE

Each item in the map has an unsigned value (with the exception of those listed below) and an unsigned key.

- o Query extended information and Response extended information which are of Type Extended Information.
- Response delay which is an integer (This can be negative if the network stack/capture library returns them out of order.)

Dickinson, et al. Expires June 9, 2017 [Page 18]

+------+ | Field | P | Description | Time offset | A | Q/R timestamp as an offset in microseconds | from the Block pre-amble Timestamp. The | timestamp is the timestamp of the Query, or | the Response if there is no Query. | A | The index in the IP address table of the | Client | | client IP address. l address | Client port | A | The client port. | Transaction | A | DNS transaction identifier. ID | Query | A | The index of the more information on the Q/R | signature | | in the Query Signature table. | Client | Q | The IPv4 TTL or IPv6 Hoplimit from the Query | hoplimit | | packet. | QR | The time different between Query and Response, | Response delay | | in microseconds. | Question | QT | The index in the NAME/RDATA table of the QNAME NAME | | for the first Question. | Response | R | The size of the DNS message (not the packet | | containing the message, just the DNS message) size | that forms the Response. | Q | Extended Query information. This item is only | Query | extended | present if collection of extra Query | information is configured. | information | | R | Extended Response information. This item is Response | extended | | only present if collection of extra Query | information | | information is configured. +---------+

The collector always collects basic Q/R information. It may be configured to collect details on Question, Answer, Authority and Additional sections of the Query, the Response or both. Note that only the second and subsequent Questions of any Question section are collected (the details of the first are in the basic information), and that OPT Records are not collected in the Additional section.
Dickinson, et al. Expires June 9, 2017 [Page 19]

The Extended information is a CBOR map as follows. Each item in the map is present only of collection of the relevant details is configured. Each item in the map has an unsigned value and an unsigned key.

+	
Field	Description
Question	The index in the Questions list table of the entry
	listing the second and subsequent Question sections
	for the Query or Response.
Answer	The index in the RR list table of the entry listing
	the Answer Resource Record sections for the Query or
	Response.
Authority	The index in the RR list table of the entry listing
	the Authority Resource Record sections for the Query
	or Response.
Additional 	The index in the RR list table of the entry listing the Additional Resource Record sections for the Query or Response.

7.19. Address Event counts

This table holds counts of various IP related events relating to traffic with individual client addresses.

Dickinson, et al. Expires June 9, 2017 [Page 20]

Field	 Type	Description
Event type 	Unsigned 	The type of event. The following events types are currently defined: 0. TCP reset. 1. ICMP time exceeded. 2. ICMP destination unreachable. 3. ICMPv6 time exceeded. 4. ICMPv6 destination unreachable. 5. ICMPv6 packet too big.
Event code 	Unsigned 	A code relating to the event. Optional.
Address index 	Unsigned 	The index in the IP address table of the client address.
Count +	Unsigned 	The number of occurrences of this event during the block collection period.

8. C-DNS to PCAP

It is possible to re-construct PCAP files from the C-DNS format. However this is a lossy process and some of the issues with reconstructing both the DNS payload and the full packet stream are outlined here.

Firstly the reconstruction depends on whether or not all the optional sections of both the query and response were captured in the C-DNS file. Clearly if they were not all captured the reconstruction is imperfect.

Secondly, even if all sections of the response were captured name compression presents a challenge in reconstructing the DNS response payload byte for byte. <u>Section 8.1</u> discusses this is more detail.

Thirdly, not all transport information is captured in the C-DNS format. For example, the following aspects of the original packet stream cannot be re-constructed from the C-DNS format:

- o IP Fragmentation
- o TCP stream information:
 - * Multiple DNS messages may have been sent in a single TCP segment

Dickinson, et al. Expires June 9, 2017 [Page 21]

- * A DNS payload may have be split across multiple TCP segments
- * Multiple DNS messages may have be sent on a single TCP session

o Malformed DNS messages and non-DNS packets

Simple assumptions can be made on the reconstruction - fragmented and DNS-over-TCP messages can be reconstructed into 'single' packets and a single TCP session can be constructed for each TCP packet.

Additionally if the malformed and non-DNS packets are captured separately into PCAPs they can be merged with PCAPs reconstructed from C-DNS to produce a more complete packet stream.

8.1. Name Compression

All the names stored in the C-DNS format are full domain names; no DNS style name compression is used on the individual names within the format. Therefore when reconstructing a packet name compression must be used in order to re-produce the on the wire representation of the packet.

[RFC1035] name compression works by substituting trailing sections of a name with a reference back to the occurrence of those sections earlier in the packet. Not all name server software uses the same algorithm when compressing domain names within the responses. Some attempt maximum recompression at the expense of runtime resources, others use heuristics to balance compression and speed and others use different rules for what is a valid compression target.

This means that responses to the same question from different name server software which match in terms of DNS payload content (header, counts, RRs with name compression removed) do not necessarily match byte for byte on the wire.

From the C-DNS format it is not possible to ensure that the DNS response payload is reconstructed byte for byte. However it can at least, in principle, be reconstructed to have the correct payload length (since the original response length is captured) if there is enough knowledge of the commonly implemented name compression algorithms. For example, a simplistic approach would be to try each algorithm in turn to see if it reproduces the original length, stopping at the first match. This would not guarantee the correct algorithm has been used as it is possible to match the length whilst still not matching the on the wire bytes but without further information added to the C-BOR this is the best that can be achieved.

Dickinson, et al. Expires June 9, 2017 [Page 22]

<u>Appendix B</u> presents an example of two differing compression algorithms used by well known name server software.

9. Data Collection

This section describes a non-normative proposed algorithm for the processing of a captured stream of DNS queries and responses and matching queries/responses where possible.

For the purposes of this discussion, it is assumed that the input has been pre-processed such that:

- All IP fragmentation reassembly, TCP stream reassembly etc. has already been performed
- Each message is associated with transport meta-data required to generate the Primary ID (see below)
- Each message has a well-formed DNS header of 12 bytes and (if present) the first RR in the query section can be parsed to generate the Secondary ID (see below).
 - * As noted earlier, this requirement can result in a malformed query being removed in the pre-processing stage, but the correctly formed response with RCODE of FORMERR being present

DNS messages are processed in the order they are delivered to the application.

o It should be noted that packet capture libraries do not necessary provide packets in strict chronological order.

[TODO: Discuss the corner cases resulting from this in more detail.]

<u>9.1</u>. Matching algorithm

A schematic representation of the algorithm for matching Q/R pairs is shown in the following diagram:

Figure showing the packet matching algorithm format (PNG) [5]

Figure showing the packet matching algorithm format (SVG) [6]

and further details of the algorithm are given in the following sections.

Dickinson, et al. Expires June 9, 2017 [Page 23]

Internet-Draft C-DNS: A DNS Packet Capture Format December 2016

9.2. Message identifiers

9.2.1. Primary ID (required)

A Primary ID can be constructed for each message which is composed of the following data:

- 1. Source IP Address
- 2. Destination IP Address
- 3. Source Port
- 4. Destination Port
- 5. Transport
- 6. DNS Message ID

9.2.2. Secondary ID (optional)

If present, the first question in the Question section is used as a secondary ID for each message. Note that there may be well formed DNS queries that have a QDCOUNT of 0, and some responses may have a QDCOUNT of 0 (for example, RCODE=FORMERR or NOTIMP)

9.3. Algorithm Parameters

1. Configurable timeout

9.4. Algorithm Requirements

The algorithm is designed to handle the following input data:

- 1. Multiple queries with the same Primary ID (but different Secondary ID) arriving before any responses for these queries are seen.
- 2. Multiple queries with the same Primary and Secondary ID arriving before any responses for these queries are seen.
- 3. Queries for which no later response can be found within the specified timeout.
- 4. Responses for which no previous query can be found within the specified timeout.

Dickinson, et al. Expires June 9, 2017 [Page 24]

Internet-Draft C-DNS: A DNS Packet Capture Format December 2016

9.5. Algorithm Limitations

For cases 1 and 2 listed in the above requirements, it is not possible to unambiguously match queries with responses. The solution to this employed in this algorithm is to match to the earliest query with the correct Primary and Secondary ID.

9.6. Workspace

A FIFO structure is used to hold the Q/R items during processing.

<u>9.7</u>. Output

The output is a list of Q/R data items. Both the Query and Response elements are optional in these items, therefore Q/R items have one of three types of content:

- 1. Paired Q/R messages
- 2. A query message (no response)
- 3. A response message (no query)

The timestamp of a list item is that of the query for cases 1 and 2 and that of the response for case 3.

<u>9.8</u>. Post Processing

When ending capture, all remaining entries in the Q/R FIFO should be treated as timed out queries.

10. IANA Considerations

None

11. Security Considerations

Any control interface MUST perform authentication and encryption.

Any data upload MUST be authenticated and encrypted.

12. Acknowledgements

The authors wish to thank CZ.NIC, in particular Tomas Gavenciak, for many useful discussions on binary formats, compression and packet matching. Also Jan Vcelak and Wouter Wijngaards for discussions on name compression.

Dickinson, et al. Expires June 9, 2017 [Page 25]

Thanks to Robert Edmonds, Paul Hoffman and Jerry Lundstroem for review.

Also, Miek Gieben for mmark [7]

<u>13</u>. Changelog

draft-ietf-dnsop-dns-capture-format-00

- o Changed dnstap.io to dnstap.info
- o qr_data_format.png was cut off at the bottom
- o Update authors address
- o Improve wording in Abstract
- o Changed DNS-STAT to C-DNS in CDDL
- o Set the format version in the CDDL
- o Added a TODO: Add block statistics
- o Added a TODO: Add extend to support pico/nano. Also do this for Time offset and Response delay
- Added a TODO: Need to develop optional representation of malformed packets within CBOR and what this means for packet matching. This may influence which fields are optional in the rest of the representation.
- o Added section on design goals to Introduction
- o Added a TODO: Can Class be optimised? Should a class of IN be inferred if not present?

draft-dickinson-dnsop-dns-capture-format-00

o Initial commit

14. References

<u>**14.1</u>**. Normative References</u>

[RFC1035] Mockapetris, P., "Domain names - implementation and specification", STD 13, <u>RFC 1035</u>, DOI 10.17487/RFC1035, November 1987, <<u>http://www.rfc-editor.org/info/rfc1035</u>>.

Dickinson, et al. Expires June 9, 2017 [Page 26]

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, DOI 10.17487/RFC2119, March 1997, <<u>http://www.rfc-editor.org/info/rfc2119</u>>.
- [RFC7049] Bormann, C. and P. Hoffman, "Concise Binary Object Representation (CBOR)", <u>RFC 7049</u>, DOI 10.17487/RFC7049, October 2013, <<u>http://www.rfc-editor.org/info/rfc7049</u>>.

<u>14.2</u>. Informative References

- [ditl] DNS-OARC, "DITL", 2016, <<u>https://www.dnsoarc.net/oarc/data/ditl</u>>.
- [dnscap] DNS-OARC, "DNSCAP", 2016, <<u>https://www.dns-oarc.net/tools/</u> dnscap>.
- [dnstap] dnstap.info, "dnstap", 2016, <<u>http://dnstap.info/</u>>.
- [I-D.daley-dnsxml]

Daley, J., Morris, S., and J. Dickinson, "dnsxml - A standard XML representation of DNS data", <u>draft-daley-dnsxml-00</u> (work in progress), July 2013.

[I-D.greevenbosch-appsawg-cbor-cddl]

Vigano, C. and H. Birkholz, "CBOR data definition language (CDDL): a notational convention to express CBOR data structures", <u>draft-greevenbosch-appsawg-cbor-cddl-09</u> (work in progress), September 2016.

- [I-D.hoffman-dns-in-json] Hoffman, P., "Representing DNS Messages in JSON", drafthoffman-dns-in-json-09 (work in progress), October 2016.
- [packetq] .SE The Internet Infrastructure Foundation, "PacketQ", 2014, <<u>https://github.com/dotse/PacketQ</u>>.
- [pcap] tcpdump.org, "PCAP", 2016, <<u>http://www.tcpdump.org/</u>>.
- [pcapng] Tuexen, M., Risso, F., Bongertz, J., Combs, G., and G. Harris, "pcap-ng", 2016, <<u>https://github.com/pcapng/</u> pcapng>.
- [rrtypes] IANA, "RR types", 2016, <<u>http://www.iana.org/assignments/</u> dns-parameters/dns-parameters.xhtml#dns-parameters-4>.

Dickinson, et al. Expires June 9, 2017 [Page 27]

14.3. URIS

- [1] https://github.com/dns-stats/draft-dns-capture- format/blob/master/cdns_format.png
- [2] <u>https://github.com/dns-stats/draft-dns-capture-</u> format/blob/master/cdns_format.svg
- [3] https://github.com/dns-stats/draft-dns-captureformat/blob/master/gr_data_format.png
- [4] https://github.com/dns-stats/draft-dns-captureformat/blob/master/qr_data_format.svg
- [5] https://github.com/dns-stats/draft-dns-capture- format/blob/master/packet_matching.png
- [6] https://github.com/dns-stats/draft-dns-capture- format/blob/master/packet_matching.svg
- [7] <u>https://github.com/miekg/mmark</u>
- [8] <u>https://www.nlnetlabs.nl/projects/nsd/</u>
- [9] https://www.knot-dns.cz/

Appendix A. CDDL

```
; CDDL specification of the file format for C-DNS,
; which describes a collection of DNS Query/Response pairs.
File = [
   file-type-id : tstr,
                                  ; "C-DNS"
   file-preamble : FilePreamble,
    file-blocks : [* Block],
]
FilePreamble = {
    format-version => uint, ; 1
   ? configuration => Configuration,
   ? generator-id => tstr,
   ? host-id => tstr,
}
format-version = 0
configuration = 1
generator-id = 2
host-id
         = 3
```

Dickinson, et al. Expires June 9, 2017 [Page 28]

```
Configuration = {
   ? query-timeout => uint,
   ? skew-timeout
                   => uint,
   ? snaplen
                    => uint,
   ? promisc
                   => uint,
                  => [* tstr],
   ? interfaces
                   => [* uint],
   ? vlan-ids
   ? filter
                    => tstr,
   ? query-options => uint,
                               ; See below
   ? response-options => uint,
   ? accept-rr-types => [* tstr],
   ? ignore-rr-types => [* tstr],
}
; query-options and response-options are bitmasks. A bit set adds in the
; specified sections.
; second & subsequent question sections = 1
; answer sections = 2
; authority sections = 4
; additional sections = 8
query-timeout
               = 0
skew-timeout = 1
snaplen
               = 2
promisc
              = 3
interfaces
              = 4
vlan-ids
              = 5
filter
               = 6
query-options = 7
response-options = 8
accept-rr-types = 9;
ignore-rr-types = 10;
Block = {
   preamble
                      => BlockPreamble,
   ? statistics
                     => BlockStatistics, ; Much of this could be derived
   tables
                      => BlockTables,
   queries
                      => [* QueryResponse],
   address-event-counts => [* AddressEventCount],
}
preamble
                   = 0
statistics
                   = 1
tables
                   = 2
queries
                   = 3
address-event-counts = 4
```

Dickinson, et al. Expires June 9, 2017 [Page 29]

```
December 2016
```

```
BlockPreamble = {
   start-time => Timeval
}
start-time = 1
Timeval = [
   seconds
                : uint,
   microseconds : uint,
]
BlockStatistics = {
   ? total-packets
                          => uint,
   ? total-pairs
                          => uint,
   ? unmatched_queries
                          => uint,
   ? unmatched_responses => uint,
   ? malformed-packets
                          => uint,
   ? non-dns-packets
                          => uint,
   ? out-of-order-packets => uint,
   ? missing-pairs
                          => uint,
   ? missing-packets
                          => uint,
   ? missing-non-dns
                          => uint,
}
total-packets
                    = 0
total-pairs
                    = 1
unmatched_queries
                    = 2
unmatched_responses = 3
malformed-packets
                    = 4
                    = 5
non-dns-packets
out-of-order-packets = 6
missing-pairs
                    = 7
missing-packets
                    = 8
missing-non-dns
                    = 9
BlockTables = {
    ip-address => [* bstr],
    classtype => [* ClassType],
   name-rdata => [* bstr],
                                      ; Holds both Name RDATA and RDATA
   query_sig => [* QuerySignature]
   ? qlist => [* QuestionList],
   ? qrr => [* Question],
   ? rrlist => [* RRList],
   ? rr => [* RR],
}
ip-address = 0
classtype = 1
```

Dickinson, et al. Expires June 9, 2017 [Page 30]

```
Internet-Draft
                  C-DNS: A DNS Packet Capture Format
                                                         December 2016
name-rdata = 2
query_sig = 3
qlist
          = 4
qrr
          = 5
         = 6
rrlist
rr
          = 7
QueryResponse = {
    time-useconds
                                     ; Time offset from earliest record
                       => uint,
    client-address-index => uint,
    client-port
                       => uint,
   transaction-id
                     => uint,
    query-signature-index => uint,
    ? client-hoplimit => uint,
    ? delay-useconds
                         => int,
                                      ; Times may be -ve at capture
    ? query-name-index => uint,
                      => uint, ; und die
=> QueryResponseExtended,
   ? response-size
                                        ; DNS size of response
   ? query-extended
    ? response-extended => QueryResponseExtended,
}
time-useconds
                     = 0
client-address-index = 1
client-port
                     = 2
transaction-id
                     = 3
query-signature-index = 4
client-hoplimit
                     = 5
delay-useconds
                     = 6
query-name-index
                     = 7
response-size
                     = 8
query-extended
                     = 9
response-extended = 10
ClassType = {
    type => uint,
   class => uint,
}
type = 0
class = 1
QuerySignature = {
    server-address-index
                           => uint,
    server-port
                          => uint,
    transport-flags
                           => uint,
    qr-sig-flags
                         => uint,
    ? query-opcode
                           => uint,
    gr-dns-flags
                           => uint,
```

Dickinson, et al. Expires June 9, 2017 [Page 31]

```
? query-rcode
                           => uint,
   ? query-classtype-index => uint,
   ? query-qd-count
                         => uint,
   ? query-an-count
                           => uint,
   ? query-ar-count
                          => uint,
   ? query-ns-count
                          => uint,
   ? edns-version
                           => uint,
   ? udp-buf-size
                           => uint,
   ? opt-rdata-index
                           => uint,
   ? response-rcode
                           => uint,
}
server-address-index = 0
server-port
                     = 1
transport-flags
                     = 2
gr-sig-flags
                     = 3
query-opcode
                     = 4
                     = 5
qr-dns-flags
query-rcode
                     = 6
query-classtype-index = 7
query-qd-count
                     = 8
query-an-count
                     = 9
query-ar-count
                     = 10
query-ns-count
                     = 11
edns-version
                     = 12
udp-buf-size
                    = 13
opt-rdata-index
                    = 14
response-rcode
                     = 15
QuestionList = [
    * uint,
                                     ; Index of Question
]
Question = {
                                     ; Second and subsequent questions
    name-index
                   => uint,
                                     ; Index to a name in the name-rdata table
   classtype-index => uint,
}
name-index
           = 0
classtype-index = 1
RRList = [
    * uint,
                                     ; Index of RR
1
RR = {
                                     ; Index to a name in the name-rdata table
   name-index
                   => uint,
    classtype-index => uint,
```

Dickinson, et al. Expires June 9, 2017 [Page 32]

```
Internet-Draft
                   C-DNS: A DNS Packet Capture Format
                                                            December 2016
    tt1
                    => uint,
    rdata-index
                    => uint,
                                      ; Index to RDATA in the name-rdata table
}
ttl
            = 2
rdata-index = 3
QueryResponseExtended = {
    ? question-index => uint,
                                      ; Index of QuestionList
                                       ; Index of RRList
   ? answer-index
                       => uint,
   ? authority-index => uint,
   ? additional-index => uint,
}
question-index
                 = 0
answer-index
                 = 1
authority-index = 2
additional-index = 3
AddressEventCount = {
                     => &AddressEventType,
    ae-type
    ? ae-code
                     => uint,
   ae-address-index => uint,
    ae-count
                     => uint,
}
                 = 0
ae-type
ae-code
                 = 1
ae-address-index = 2
ae-count
                = 3
AddressEventType = (
    tcp-reset: 0,
    icmp-time-exceeded
                           : 1,
    icmp-dest-unreachable : 2,
    icmpv6-time-exceeded
                           : 3,
    icmpv6-dest-unreachable: 4,
    icmpv6-packet-too-big : 5,
)
```

Appendix B. DNS Name compression example

The basic algorithm which follows the guidance in [RFC1035] is simply to collect each name, and the offset in the packet at which it starts, during packet construction. As each name is added, it is offered to each of the collected names in order of collection, starting from the first name. If labels at the end of the name can be replaced with a reference back to part (or all) of the earlier

Dickinson, et al. Expires June 9, 2017 [Page 33]

name, and if the uncompressed part of the name is shorter than any compression already found, the earlier name is noted as the compression target for the name.

The following tables illustrate the process. In an example packet, the first name is example.com.

++	+	+	+
N Name	Uncompressed	Compression Target	
1 example.com	+		
++	+	+	• +

The next name added is bar.com. This is matched against example.com. The com part of this can be used as a compression target, with the remaining uncompressed part of the name being bar.

++-	+	+	ŀ
N Name	Uncompressed	Compression Target	
++ 1 example.com 2 bar.com ++	+ bar +	+ 1 + offset to com +	+

The third name added is www.bar.com. This is first matched against example.com, and as before this is recorded as a compression target, with the remaining uncompressed part of the name being www.bar. It is then matched against the second name, which again can be a compression target. Because the remaining uncompressed part of the name is www, this is an improved compression, and so it is adopted.

+---+ | N | Name | Uncompressed | Compression Target | +---+ 3 www.bar.com www 2 +---+----+

As an optimization, if a name is already perfectly compressed - in other words, the uncompressed part of the name is empty - no further names will be considered for compression.

B.1. NSD compression algorithm

Using the above basic algorithm the packet lengths of responses generated by NSD $[\underline{8}]$ can be matched almost exactly. At the time of

Dickinson, et al. Expires June 9, 2017 [Page 34]

writing, a tiny number (<.01%) of the reconstructed packets had incorrect lengths.

B.2. Knot Authoritative compression algorithm

The Knot Authoritative [9] name server uses different compression behavior, which is the result of internal optimization designed to balance runtime speed with compression size gains. In brief, and omitting complications, Knot Authoritative will only consider the ONAME and names in the immediately preceding RR section in an RRSET as compression targets.

A set of smart heuristics as described below can be implemented to mimic this and while not perfect it produces output nearly, but not quite, as good a match as with NSD. The heuristics are:

- 1. A match is only perfect if the name is completely compressed AND the TYPE of the section in which the name occurs matches the TYPE of the name used as the compression target.
- 2. If the name occurs in RDATA:
 - a If the compression target name is in a query, then only the first RR in an RRSET can use that name as a compression target.
 - b The compression target name MUST be in RDATA.
 - c The name section TYPE must match the compression target name section TYPE.
 - d The compression target name MUST be in the immediately preceding RR in the RRSET.

Using this algorithm less than 0.1% of the reconstructed packets had incorrect lengths.

B.3. Observed differences

In sample traffic collected on a root name server around 2-4% of responses generated by Knot had different packet lengths to those produced by NSD.

Appendix C. Comparison of Binary Formats

Several binary representations were considered in particular CBOR, Apache Avro and Protocol Buffers.

Dickinson, et al. Expires June 9, 2017 [Page 35]

Protocol Buffers and Avro both require a data schema, and validate data being stored against that schema.

[TODO: Finish pros and cons of CBOR vs Avro vs Protocol buffers - tools, schema, adoption, etc.]

The difference in file sizes were mostly minimal See Appendix D.3.

Appendix D. Sample data on the C-DNS format

This section presents some example figures for the output size of capture files when using different block sizes, data representations and binary formats. The data is sample data for a root instance.

[TODO: This section needs more work..]

D.1. Comparison to full PCAPS

As can be seen in more detail below for this sample data the compressed C-DNS files are around 30% the size of the full compressed PCAPs. It should also be noted that experiments showed that compression of the C-DNS format required very roughly an order of magnitude less CPU resources than compression of full PCAPSs when using one core from a 3.5GHz i7 processor.

D.2. Block size choices

[TODO: Discuss trade-off of file block size vs memory consumption.]

[TODO: Add graph that demonstrates block size of 5000 is optimal for the sample data used.]

D.3. Blocking vs more simple output

Some experiments were conducted producing output in a very simple format involving a single record per Q/R data item (akin to a .csv representation). The aim here was to examine whether the blocking mechanism (using a block size of 5000) was worth the complexity, particularly after compression of the output file using several general purpose compression tools. The original PCAP file was 325.79M and compressed using xz to 24.3Mb.

Dickinson, et al. Expires June 9, 2017 [Page 36]

It might be expected that blocking is exploiting commonality that a general purpose compression engine could also exploit, and the figures do indeed bear this out. The more powerful (and resource-consuming) the compression, the closer the compressed simple file size gets to the compressed chunk file size. With no compression, the blocked output size is typically half that of the simple output, but as greater degrees of compression are applied the gap shrinks. However, even with the stronger compressor, the chunked output remains roughly 5-10% smaller than the simple output. This, and the higher gains at lower compression, might be significant, depending on the target environment.

[TODO: Add data on reduction in CPU overhead of compressing blocked output vs simple output.]

This was repeated using some other binary representations:

Format Output size lz4 gzip xz +	+	+	++	+ +
json-simple 189.85M 25.59M 16.03M 9.74M avro-simple 43.31M 16.07M 11.92M 7.99M	Format	Output size	1z4 gz	zip xz
avro-block 17.44M 12.94M 10.08M 7.18M protobuf-simple 46.02M 15.79M 11.59M 7.94M protobuf-block 22.08M 15.43M 10.91M 7.40M	<pre> json-simple avro-simple avro-block protobuf-simple protobuf-block +</pre>	189.85M 43.31M 17.44M 46.02M 22.08M	25.59M 16.0 16.07M 11.9 12.94M 10.0 15.79M 11.9 15.43M 10.9	03M 9.74M 92M 7.99M 08M 7.18M 59M 7.94M 91M 7.40M

There's not a lot to choose between the three contenders with simple output. Avro produces the smaller output, CBOR the next and Protocol Buffers the largest, but the different is under 10%. However, with blocking, while CBOR and Protocol Buffers are again within a few percentage points of each other (though Protocol Buffers now has a slight advantage), Avro produces files in the region of 20% smaller, and holds a diminishing advantage through increased compression.

Authors' Addresses

John Dickinson Sinodun IT Magdalen Centre Oxford Science Park Oxford OX4 4GA

Email: jad@sinodun.com
Dickinson, et al. Expires June 9, 2017 [Page 37]

Jim Hague Sinodun IT Magdalen Centre Oxford Science Park Oxford OX4 4GA

Email: jim@sinodun.com

Sara Dickinson Sinodun IT Magdalen Centre Oxford Science Park Oxford OX4 4GA

Email: sara@sinodun.com

Terry Manderson ICANN 12025 Waterfront Drive Suite 300 Los Angeles CA 90094-2536

Email: terry.manderson@icann.org

John Bond ICANN 12025 Waterfront Drive Suite 300 Los Angeles CA 90094-2536

Email: john.bond@icann.org

Dickinson, et al. Expires June 9, 2017 [Page 38]