

# Data Formats

## HashSets.com Databases

**Please Note: Within this document we have removed some older and insignificant table fields to help you in simplifying table importation and subsequent table data analysis. Please note the changes further below and implement into your own new or preexisting database tables.**

Updated: 12 October 2023

### 1. Introduction

This document describes the formats of various database tables provided by HashSets.com to customers who subscribe to our 'Database Exports' (also known as 'database dumps'). Throughout this document we will cover only trusted, known-good and/or non-threatening hash set tables from our database.

Since November 2003, Whitehat Computer Forensics, LLC, has been performing hash file calculations of common electronic files found within various types of computer operating systems, workstations and servers.

Our largest hash set consist of operating system files from more than 698+ operating system versions we have installed and forensically analyzed since 2003. The process involves installing the operating systems onto sterile hard drives and/or virtual disks, performing analysis using commercially available computer forensic software to discover all files, file attributes and meta-data and most importantly the MD5, SHA1, SHA-256 and SHA-512 hash values.

All of the aforementioned findings are then imported into individual operating system tables which subsequently form our HashSets.com database and online search engine. We concurrently offer to our Gold and Platinum subscription members the database tables exported into large tab delimited text files which may be reimported into their own database or third-party software products. Each table also provides a header row for simple field name explanation and identification (see further below for details).

### 2. Database Overview of the 'Primary Database Export' Table

The HashSets Database consists of fifteen very large database tables. From these tables there is one master table named '**Primary Database Export**' table.

The '**Primary Database Export**' table contains file information that never changes including the individual file's corresponding MD5, SHA1, SHA-256 and SHA-512 hash values, the file's first 32 bytes in hexadecimal, the first 128 ASCII characters (similar to using Unix or Linux 'strings' command), logical byte file size, and so on.

The '**Primary Database Export**' table contains hash values of safe, non-threatening, known-good computer files only. There are also NO duplicate hash values in this table. This is considered the master table that would be associated with the remaining fourteen tables discussed further below.

## 'Primary Database Export' Table Schema

Field Name	Type	Length	Description
<b>MD5</b>	char	32	128-bit Message Digest 5 (hash value) of a specific file.
<b>SHA_1</b>	char	40	160-bit Secure Hash Algorithm message digest (hash value) of a specific file.
<b>SHA_256</b>	char	64	256-bit Secure Hash Algorithm message digest (hash value) of a specific file.
<b>Header_HEX</b>	varchar	64	First 32 bytes in Hexadecimal format of the file.
<b>128_Bytes_ASCII</b>	varchar	128	First detectable 128 bytes in ASCII format (Similar to performing Linux or UNIX STRINGS of a file).
<b>Signature</b>	varchar	255	Potential file header signature(s), if known.
<b>Logical_Size</b>	bigint	17	File size in byte format.
<b>NSRL</b>	char	3	File known to be found within the National Software Reference Library (NSRL) Dataset released by the US Government. Will be marked with either 'Yes' or 'No'.
<b>key_field</b>	int	11	Table key that uniquely defines a record, if provided.
<b>SHA_512</b>	char	128	512-bit Secure Hash Algorithm message digest (hash value) of a specific file, if available.
<b>CRC32</b>	char	8	32-bit Checksum of a specific file, if known.
<b>Fuzzy_SSdeep</b>	varchar	255	Context Triggered Piecewise Hash values (CTPH). Also called "fuzzy" hash values. (Future Use Only/Work in Progress).

<b>Ignorable</b>	char	7	<p>Intended to identify potentially ignorable or non-intrinsic hash values. For example, common operating system log files (.log, .evt, .evtx) which may have different contents and hash values if found later on other installations of operating systems, different computers or devices.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>
<b>Operating_System_File</b>	char	7	<p>Hash values found within an Operating System.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>
<b>Other_Computer_File</b>	char	7	<p>Hash values found outside of operating system installations. For example, third-party software, applications, drivers, utilities, etc, that we downloaded directly from manufacturer websites, etc.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>
<b>MS_Windows_OS</b>	char	7	<p>Hash values found within Microsoft Windows operating system installations and updates.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>
<b>Linux_OS</b>	char	7	<p>Hash values found within Linux operating system installations and updates.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>
<b>BSD_OS</b>	char	7	<p>Hash values found within BSD (UNIX like operating system) installations and updates.</p> <p>Will be marked with either 'Yes' or 'Unknown'.</p>

<b>macOS</b>	char	7	Hash values found within Apple's macOS (formerly OS X) operating system installations and updates.  Will be marked with either 'Yes' or 'Unknown'.
<b>Solaris_OS</b>	char	7	Hash values found within Oracle's (formerly Sun) Solaris operating system installations and updates.  Will be marked with either 'Yes' or 'Unknown'.
<b>Gold_Disks</b>	char	7	Hash values found within Oracle's (formerly Sun) Solaris operating system installations and updates.  Will be marked with either 'Yes' or 'Unknown'.
<b>Installation_Discs</b>	Char	7	Hash values found within operating system installation media (CD, DVD and .ISO) files.  Will be marked with either 'Yes' or 'Unknown'.

### 3. All Other Remaining Tables Described

The remaining fourteen tables are divided into operating systems (MS Windows, Linux, BSD, macOS and Solaris) and non-operating system groups (MS Windows App Store, Mac App Store, Gold Disks, Installation Media, third-party manufacturer's Software & Hardware Applications, etc).

Below is a description of each of the fourteen remaining tables. It should be mentioned briefly that over the course of many years (20+ years) the size of MS Windows and macOS operating system tables grew too large to contain within one table each. Therefore, in July 2021, we decided to divide those tables into smaller tables for convenience and ease of importation. The MS Windows tables were ultimately divided into three smaller tables (Windows North America, Windows Europe, Windows Asia and Windows Middle-East). For macOS we divided the previous single table into two tables ('macOS versions 8 thru 10' and 'macOS versions 11 and above').

- **Windows\_North\_America\_Database\_Export** – This database table contains various Microsoft Windows operating system versions for the United States, Canada (French Canadian), Mexico (Spanish Mexican) and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating

system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.

- **Windows\_Europe\_Database\_Export** – This database table contains various Microsoft Windows operating system versions for Europe (Western and Eastern) and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **Windows\_Asia\_Database\_Export** – This database table contains various Microsoft Windows operating system versions for Asia (Japan, China and so on) and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **Windows\_Middle\_East\_Database\_Export** – This database table contains various Microsoft Windows operating system versions for the Middle East (Hebrew, etc) and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **Linux\_Database\_Export** – This database table contains various Linux operating system distributions and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **macOS\_8\_thru\_10\_Database\_Export** – This database table contains various Apple macOS 8 and 9 (legacy) and macOS 10 (OS X) operating system versions and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **macOS\_11\_and\_above\_Database\_Export** – This database table contains Apple macOS 11 and above operating system versions and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.

- **BSD\_Database\_Export** – This database table contains various BSD (UNIX like) operating system distributions and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **Solaris\_Database\_Export** – This database table contains various Solaris operating system versions and the individual file or folder names, file name extensions, last Modified/Accessed/Created file dates, and so on. Basically, the individual file or folder information that may not remain the same throughout the purpose of a MS Windows operating system regardless of individual file hash values (MD5, SHA1, SHA-256 and SHA-512) that generally remain the same until the file itself is ultimately changed or modified.
- **Applications\_and\_Hardware\_Software\_Drivers\_Database\_Export** – This database table contains the details of files gathered from third-party utilities and software applications as well as software drivers from common hardware manufacturer websites. Due to the large number of files within this category the installation process (installing onto any ordinary running operating system) before hashing could not be performed. As an alternative, the analysis consisted of unpacking executable files and compressed files, when possible, then subsequently gathering hash value calculations and file meta-data details.
- **Windows\_App\_Store\_Database\_Export** – This database table contains MS Windows 8 and 10 Applications commonly found within the MS Windows App Store. Specifically, downloadable business, game, education, and other apps which were installed, analyzed and then gathered into MD5, SHA-1 and SHA-256 hash sets.
- **Mac\_App\_Store\_Database\_Export** – This database table contains macOS applications commonly found within the macOS App Store. Specifically, downloadable business, game, education and other apps which were installed, analyzed, hashed and then gathered into MD5, SHA-1 and SHA-256 hash sets.
- **Gold\_Disks\_Windows\_Database\_Export** – This database table contains installations of various Microsoft Windows operating systems with the added installations of commonly used programs and software which we have analyzed and hashed. For example, Microsoft Visual Studio, Microsoft Office Suites, common web browsers, etc. Efforts are made periodically to update the Gold Disks installations, but it is best efforts only currently.
- **Microsoft\_Windows\_Installation\_Media\_Database\_Export** – This database table contains Microsoft Windows installation media (CD, DVD and ISO) which we have gathered, analyzed, hashed and then subsequently placed into MD5, SHA-1, SHA-256 hash sets.

Other Optional Table(s):

- **US\_Government\_Database\_Export** – This optional database table contains common non-threatening known hash values consisting of US Government (federal, state, local and military) publicly accessible website images, logos, multimedia files, office documents (.doc, .pdf, .xls, .ppt, etc). Please note that this table may be removed completely sometime in the future as these formerly common hash values may no longer be of use or popular to computer forensic or

computer security professionals. For the time being we have furnished this table for optional download and use.

- **File\_Extensions** – This optional table contains general information or descriptions pertaining various file name extensions (.exe, .pdf, .dll, etc).

For example, the below file extension .EXE within this table could be described as one of the following from historical popular use:

EXECUTABLE FILE ::: SELF-DISPLAYING IMAGE ::: SELF-EXTRACTING ARCHIVE  
 ::: SETTLERS 4 SAVE FILE ::: PDP-10 PAGE-MAPPED EXECUTABLE BINARY FILE :::  
 PLAYSTATION EXECUTABLE FILE ::: OUT-OF-PROCESS CODE COMPONENT FILE :::  
 MICROSOFT LINKER EXE INPUT FILE EXTENSION ::: DATAFLEX RUNTIME FILE  
 EXTENSION ::: SELF-EXTRACTING ARCHIVE ::: MIME: APPLICATION/OCTET-  
 STREAM FILE EXTENSION ::: MIME: APPLICATION/X-MSDOWNLOAD

This information above and within the table is used only as one of many starting points when analyzing computer files. It is not to be used as a guarantee that a particular file with a specific file extension is truly associated with any particular software, program, third-party utility, hardware device, etc.

## 4. Operating Systems and Non-operating System Table Schema

The following depicts the data elements for the previously mentioned fourteen tables. To associate the ‘Primary Database Export’ table with any of the above fourteen tables you would use the MD5 hash value from both tables as the linking “Key”.

Note: As mentioned at the beginning of this document the below table schema no longer includes older fields labeled ‘Record\_Date’, ‘Key\_Field’, ‘From\_Recycle\_Bin’, ‘From\_Free\_Space’, ‘Moved\_From\_Location’, ‘Graphic’ and ‘Inode\_Number’.

Note: Additionally, we have lengthened the below fields labeled ‘Extraneous’ and ‘Log’ from 3 to 7.

Field Name	Type	Length	Description
<b>MD5</b>	char	32	128-bit Message Digest 5 (hash value) of a specific file.
<b>Name</b>	varchar	255	Names of files and folders.
<b>File_Ext</b>	varchar	255	The file name’s extension, if applicable.
<b>Description</b>	varchar	75	A general description of the file or folder.

<b>Last_Accessed</b>	datetime	0	<p>Date and Time that the file was last accessed.</p> <p>Uses the following Date and Time format (DD/MM/YYYY Hour:Minute:Seconds). An example would be (20/11/2020 14:06:49)</p>
<b>File_Created</b>	datetime	0	<p>Date and Time that the file was created.</p> <p>Uses the following Date and Time format:</p> <p>DD/MM/YYYY Hour:Minute:Seconds</p> <p>An example would be (20/11/2020 14:06:49)</p>
<b>Last_Written</b>	datetime	0	<p>Date and Time that the file was last modified or written.</p> <p>Uses the following Date and Time format:</p> <p>DD/MM/YYYY Hour:Minute:Seconds</p> <p>An example would be (20/11/2020 14:06:49)</p>
<b>Full_Path</b>	text	0	The full path to the file or folder.
<b>Quick_Category</b>	varchar	75	<p>Used for quick identification of a group of files and folders that were examined together:</p> <p>Example 'Quick Categories':</p> <ul style="list-style-type: none"> <li>• Windows 11 Professional (64bit) – French;</li> </ul>



			<ul style="list-style-type: none"> <li>• Windows 11 Professional (64bit) – Russian;</li> <li>• Windows 10 Enterprise (64bit) - Chinese Simplified;</li> <li>• FreeBSD 4.6 (32bit);</li> <li>• Web Application;</li> <li>• Mac App Store;</li> <li>• Windows App Store;</li> </ul>
<b>File_Notes</b>	varchar	255	Internal Use Only: Analysis notes we mentioned internally during our analysis.
<b>Major</b>	varchar	75	<p>Main grouping label of a file or folder.</p> <p>Some Example of ‘Majors’:</p> <ul style="list-style-type: none"> <li>• Operating Systems;</li> <li>• Applications;</li> <li>• Etc.</li> </ul>
<b>Minor</b>	varchar	75	<p>Secondary or subgrouping from the file or folder’s ‘Major’.</p> <p>Some Example of ‘Minors’:</p> <ul style="list-style-type: none"> <li>• Installation;</li> <li>• Software Updates &amp; Fixes;</li> <li>• Decompressed and/or Extracted;</li> <li>• Etc.</li> </ul>
<b>Operating_System</b>	varchar	75	<p>Name of the affiliated operating system, if applicable.</p> <p>Some examples:</p>

			Windows Linux macOS
<b>Manufacturer</b>	varchar	150	Name of the manufacturer, if known or applicable.
<b>Version</b>	varchar	50	Version name, if known or applicable.
<b>Is_Deleted</b>	char	7	If the file or folder was found to be deleted by the host operating system, if applicable.  Will be marked with either 'True' or 'False'.
<b>Website</b>	varchar	75	Website source, if applicable.
<b>Geographic_Location</b>	varchar	50	Geographic location of the manufacturer's Headquarters. If not known or not applicable then manufacturer's intended audience.  Will be marked with either: <ul style="list-style-type: none"> <li>• North America</li> <li>• South America</li> <li>• Europe</li> <li>• Asia</li> <li>• Asia/Pacific</li> <li>• Middle East</li> <li>• North America/Europe/Asia</li> <li>• Worldwide</li> <li>• Undetermined</li> </ul> Potentially other regions or regional groups will be added in the future.
<b>Extraneous</b>	char	7	(Work in Progress) Intended to identify ignorable files that are unique to one particular situation or instance

			For example system log files, registry files, etc.  Will be marked with either 'Yes' or 'Unknown'.
<b>Log</b>	char	7	(Work in Progress) Intended to identify potential log files.  Will be marked with either 'Yes' or 'Unknown'.
<b>Bad_Extension</b>	varchar	255	File found to have a suspected bad file extension in comparison to the file's signature/header information.  Will be marked with either 'True', 'False' or 'Unknown'.
<b>actual_file</b>	varchar	255	True if an actual file. False if derived or generated from an actual file (e.g. data broken out from compound files, EXIF data from graphic images, file metadata, and so on.  Will be marked with either 'True', 'False' or 'Unknown'.
<b>file_class</b>	varchar	255	Class of file such as Regular File, Symbolic Link, etc, if known.
<b>folder</b>	varchar	255	Identified if a Folder.  Will be marked with either 'True' or 'False'.
<b>category</b>	varchar	255	Specific type of file, if available, such as EXE, Text, Unicode, 7 bit text, etc.
<b>compressed</b>	varchar	255	Found as a compressed file, if applicable.
<b>compressed_file_size</b>	bigint	15	The compressed size of the file

			in bytes (compressed files only).
<b>Compression_Method</b>	varchar	255	Compressed file's compression method (Zip files only) such as Deflated, Stored, etc.
<b>Extract_Version</b>	varchar	255	Compressed file's extraction version (Zip files only), if identifiable.
<b>permissions</b>	varchar	255	Primarily UNIX like read/write/execute permissions, if identifiable.
<b>UID</b>	int	255	Primarily UNIX like User ID, if identifiable.
<b>Group_Name_UNIX</b>	varchar	255	UNIX like Group Name, if identifiable.
<b>GID</b>	int	255	UNIX like Group ID, if identifiable.
<b>Username</b>	varchar	255	Username of the file (Unix file systems only), if identifiable.
<b>Container</b>	varchar	255	Does the file, disk, partition, etc, have any children.  Will be marked with either 'True' or 'False'.
<b>Encrypted</b>	varchar	255	Found to have Encryption file like qualities.
<b>Deleted_Date</b>	datetime	0	The Date deleted (Unix file systems only).
<b>checksum</b>	varchar	255	Checksum computed on compressed files (Zip files only).
<b>Hash_Search_Engine_Record_Date</b>	date	0	The date we added the file or folder's data into our databases.

## 5. In Summary

If, for whatever reason, you run into any issues or problems understanding the database table structures, data or field types then please feel free to reach out to us via our website [HashSets.com](https://HashSets.com). We will make every effort to provide you with a reasonable amount of additional information to help you better understand the database tables.