



Confidence in a connected world.



## Clean Data Profiling

Bartek Uscilowski

Julie Weber

October 2008

# Introduction & outline

- The paradigm change
- Clean Data sets
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



File Name	Size	Type	Date
...	11 KB	Application Extension	23/08/2014
...	108 KB	Application Extension	14/04/2014
...	487 KB	Application Extension	14/04/2014
...	63 KB	Application	14/04/2014
...	71 KB	Compiled HTML Help...	29/08/2014
...	76 KB	Application Extension	14/04/2014
...	20 KB	Application	14/04/2014
...	24 KB	RLL File	13/04/2014
...	68 KB	Application	19/03/2014
...	101 KB	Application	14/04/2014
...	33 KB	Application	14/04/2014
...	49 KB	Application	28/02/2014
...	57 KB	Application Extension	14/04/2014
...	16 KB	Application Extension	14/04/2014
...	380 KB	Application	14/04/2014
...	336 KB	Application Extension	14/04/2014
...	1 KB	SRG File	07/05/1996
...	25 KB	Application	14/04/2014
...	40 KB	Windows Script Co...	29/08/2014
...	60 KB	Help File	29/08/2014
...	39 KB	Application	14/04/2014
...	1 KB	RAM File	29/08/2014
...	14 KB	Application Extension	29/08/2014
...	181 KB	Application Extension	14/04/2014
...	13 KB	Application Extension	14/04/2014

# Outline

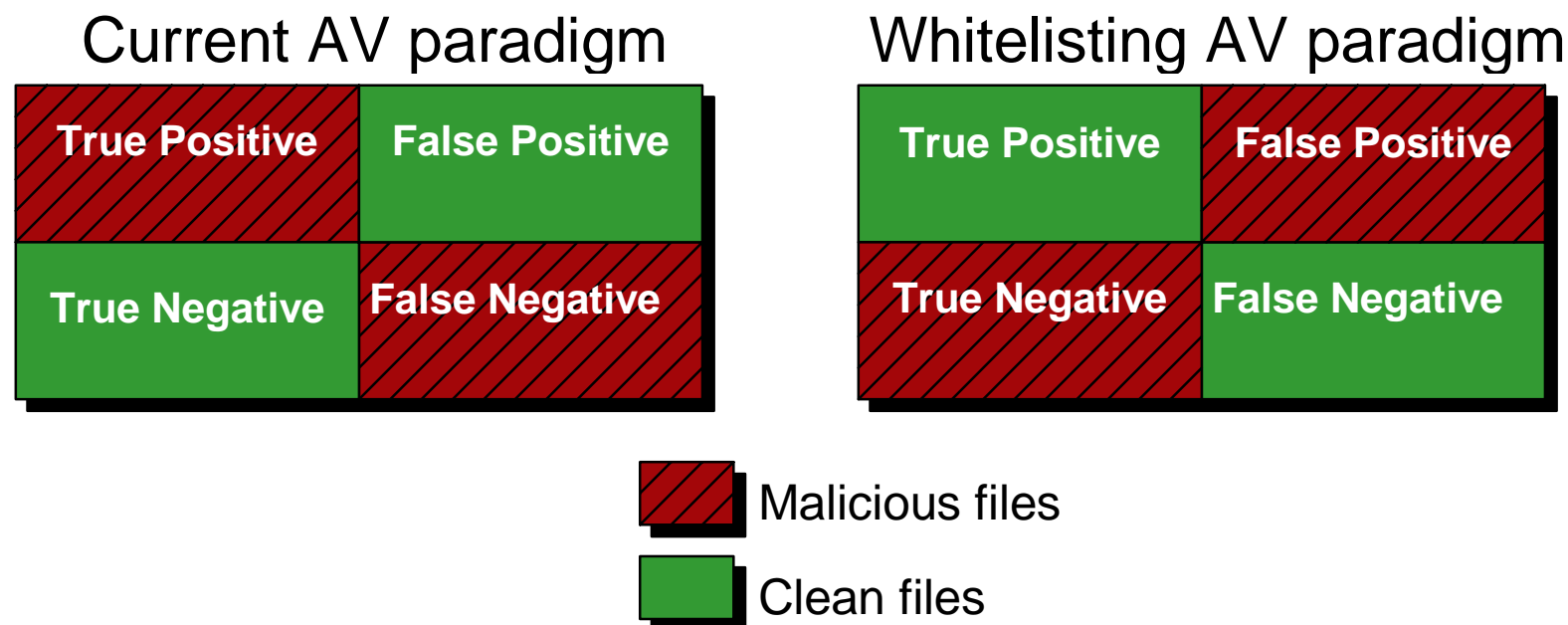
- The paradigm change
- Clean Data sets
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



File Name	Size	Type	Date
...	11 KB	Application Extension	23/08/2014
...	108 KB	Application Extension	14/04/2014
...	487 KB	Application Extension	14/04/2014
...	63 KB	Application	14/04/2014
...	71 KB	Compiled HTML Help...	29/08/2014
...	76 KB	Application Extension	14/04/2014
...	20 KB	Application	14/04/2014
...	24 KB	RLL File	13/04/2014
...	68 KB	Application	19/03/2014
...	101 KB	Application	14/04/2014
...	33 KB	Application	14/04/2014
...	49 KB	Application	28/02/2014
...	57 KB	Application Extension	14/04/2014
...	16 KB	Application Extension	14/04/2014
...	380 KB	Application	14/04/2014
...	336 KB	Application Extension	14/04/2014
...	1 KB	SRG File	07/05/1980
...	25 KB	Application	14/04/2014
...	40 KB	Windows Script Co...	29/08/2014
...	60 KB	Help File	29/08/2014
...	39 KB	Application	14/04/2014
...	1 KB	RAM File	29/08/2014
...	14 KB	Application Extension	29/08/2014
...	181 KB	Application Extension	14/04/2014
...	13 KB	Application Extension	14/04/2014

# The paradigm change

- Exponential growth of malware
- Current AV paradigm possibly insufficient
- Increasing importance of clean data



# Outline

- The paradigm change
- Clean Data sets
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



File Name	Size	Type	Date
...	11 KB	Application Extension	23/08/2014
...	108 KB	Application Extension	14/04/2014
...	487 KB	Application Extension	14/04/2014
...	63 KB	Application	14/04/2014
...	71 KB	Compiled HTML Help...	29/08/2014
...	76 KB	Application Extension	14/04/2014
...	20 KB	Application	14/04/2014
...	24 KB	RLL File	13/04/2014
...	68 KB	Application	19/03/2014
...	101 KB	Application	14/04/2014
...	33 KB	Application	14/04/2014
...	49 KB	Application	28/02/2014
...	57 KB	Application Extension	14/04/2014
...	16 KB	Application Extension	14/04/2014
...	380 KB	Application	14/04/2014
...	336 KB	Application Extension	14/04/2014
...	1 KB	SRG File	07/05/1980
...	25 KB	Application	14/04/2014
...	40 KB	Windows Script Co...	29/08/2014
...	60 KB	Help File	29/08/2014
...	39 KB	Application	14/04/2014
...	1 KB	RAM File	29/08/2014
...	14 KB	Application Extension	29/08/2014
...	181 KB	Application Extension	14/04/2014
...	13 KB	Application Extension	14/04/2014

# Clean Data Sets

- Set of data derived from legitimate software
- Targeted
  - Specific purpose (FP ?)
  - Clean data is associated with software that is known and identified in advance
- Non targeted
  - Degree of randomness
  - May have specific purpose
  - Web crawler, user images...

# Outline

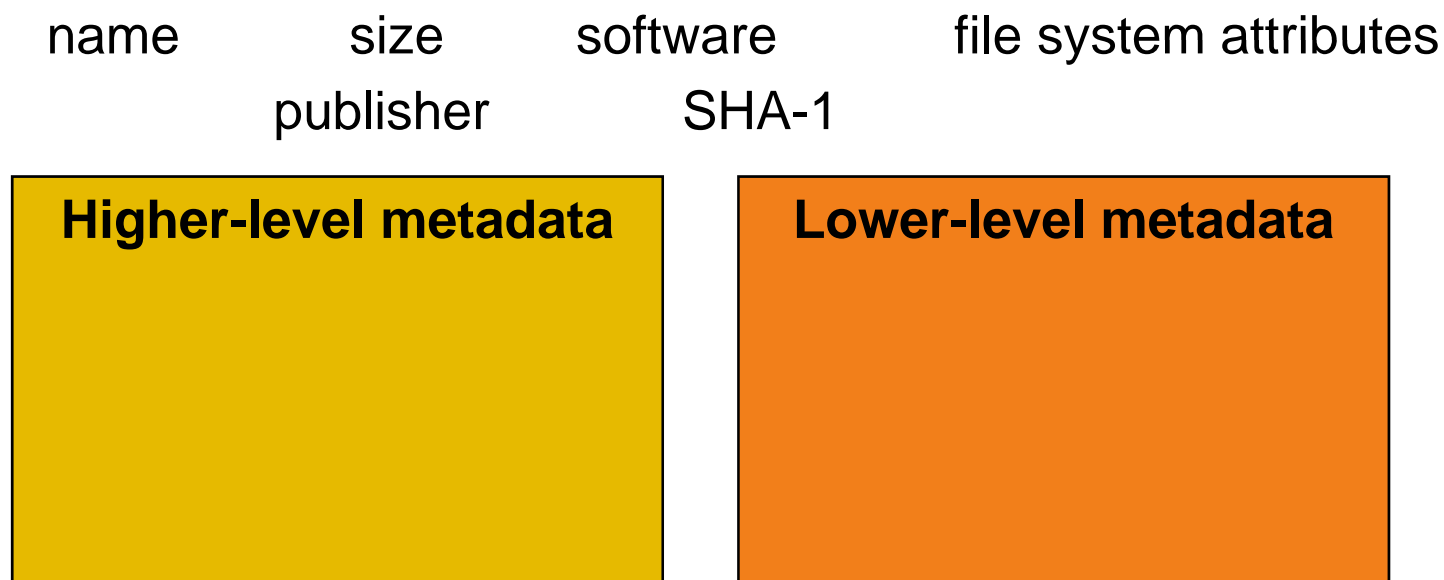
- The paradigm change
- Clean Data sets
- **Metadata**
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



File Name	Size	Type	Date Modified
...	11 KB	Application Extension	23/08/2014
...	108 KB	Application Extension	14/04/2014
...	487 KB	Application Extension	14/04/2014
...	63 KB	Application	14/04/2014
...	71 KB	Compiled HTML Help...	29/08/2014
...	76 KB	Application Extension	14/04/2014
...	20 KB	Application	14/04/2014
...	24 KB	RLL File	13/04/2014
...	68 KB	Application	19/03/2014
...	101 KB	Application	14/04/2014
...	33 KB	Application	14/04/2014
...	49 KB	Application	28/02/2014
...	57 KB	Application Extension	14/04/2014
...	16 KB	Application Extension	14/04/2014
...	380 KB	Application	14/04/2014
...	336 KB	Application Extension	14/04/2014
...	1 KB	SRG File	07/05/1980
...	25 KB	Application	14/04/2014
...	40 KB	Windows Script Co...	29/08/2014
...	60 KB	Help File	29/08/2014
...	39 KB	Application	14/04/2014
...	1 KB	RAM File	29/08/2014
...	14 KB	Application Extension	29/08/2014
...	181 KB	Application Extension	14/04/2014
...	13 KB	Application Extension	14/04/2014

# Metadata

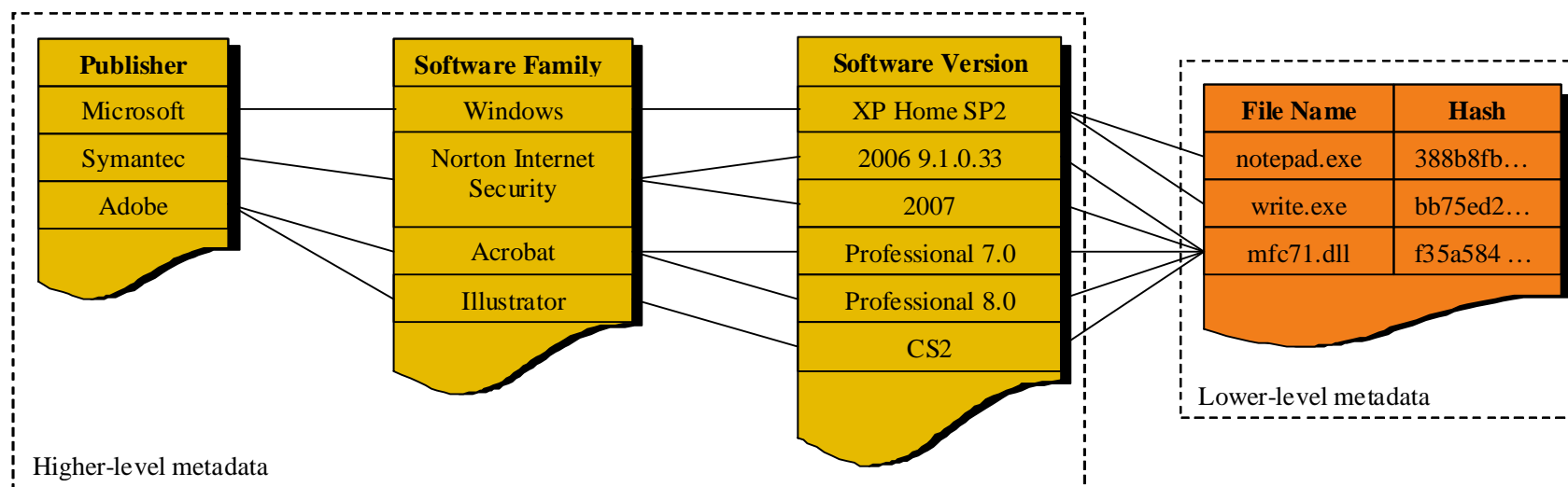
- Describes the features of the set
- Information relative to a file
  - Example:





# Metadata

- Relationship between metadata levels



- Files vs associated software & publisher
- Lower-level – extracted from files
- Higher-level – annotated or profiled

# Outline

- The paradigm change
- Clean Data sets
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary

File Name	Size	Type	Date
...	...	...	...
...	11 KB	Application Extension	23/08/2014
...	108 KB	Application Extension	14/04/2014
...	487 KB	Application Extension	14/04/2014
...	63 KB	Application	14/04/2014
...	71 KB	Compiled HTML Help...	29/08/2014
...	76 KB	Application Extension	14/04/2014
...	20 KB	Application	14/04/2014
...	24 KB	RLL File	13/04/2014
...	68 KB	Application	19/03/2014
...	101 KB	Application	14/04/2014
...	33 KB	Application	14/04/2014
...	49 KB	Application	28/02/2014
...	57 KB	Application Extension	14/04/2014
...	16 KB	Application Extension	14/04/2014
...	380 KB	Application	14/04/2014
...	336 KB	Application Extension	14/04/2014
...	1 KB	SRG File	07/05/1980
...	25 KB	Application	14/04/2014
...	40 KB	Windows Script Co...	29/08/2014
...	60 KB	Help File	29/08/2014
...	39 KB	Application	14/04/2014
...	1 KB	RAM File	29/08/2014
...	14 KB	Application Extension	29/08/2014
...	181 KB	Application Extension	14/04/2014
...	13 KB	Application Extension	14/04/2014

# Lower-level metadata

- Hashes
  - MD5, SHA-1, SHA-2 family
  - cryptographic libraries

- Size, attributes

- **MIME and Magic** **Example for mfc71.dll**

```
MD5: f35A584E947A5B401FEB0FE01DB4A0D7
SHA-1: 664dc99e78261a43d876311931694b6ef87cc8b9
SHA-256: 23A1570B8518459265dbd54b9bbc4247c9d4d45bb99b3b4baec589973d
```

```
C:\>file -i C:\windows\notepad.exe
C:\windows\notepad.exe: application/x-dosexec

C:\>file C:\windows\notepad.exe
C:\windows\notepad.exe: PE executable for MS Windows (GUI) Intel 80386 32-bit
```

# Lower-level metadata

```

sigcheck v1.53 = sigcheck
Copyright (c) 2004-2008 Mark Russinovich
sysinternals www.sysinternals.com
C:\windows\notepad.exe:
PE version information
  Signers:
    → Version information embedded in file
    → Sigcheck Microsoft Windows Verification Intermediate PCA
      Microsoft Root Authority
Digital signatures
C:\>file C:\windows\notepad.exe
Publisher: Microsoft Corporation
C:\windows\notepad.exe: MS executable for MS Windows (GUI) Intel 80386 32-bit
  → Signers chain, dates
  → Sigcheck: Signtool Microsoft« Windows« Operating System
    Version: 5.1.2600.2180
    File version: 5.1.2600.2180 (xpsp_sp2_rtm.040803-2158)
    Original Name: NOTEPAD.EXE
    Internal Name: Notepad
    Copyright: © Microsoft Corporation. All rights reserved.
    Comments: n/a
  
```

# Lower-level metadata

```

C:\>signtool verify /timestamped \WINDOWS\2008\pad1a49
Timestamp Verified by:
Verifying Co\WINDOWS\SoftwareAuthority
FileIssuedBy Microsoft Root\WINDOWS\system32\CatRoot\{F750E6C3-38EE-11D1-85E5-
00C04F8937FE}\sp1/cat/2020 08:00:00
SHA1Hash Microsoft Root Authority
Issued to: Microsoft Root Authority
IssuedBy Microsoft Root Authority
ExpiresBy: Microsoft Root Authority
SHA1Hash A43489169A20090D93D032CCAF37E7FE20A8B419
SHA1 hash: 3EA99A60058275E0ED83B892A909449F8C33B245
Issued to: Microsoft Windows Verification Intermediate PCA
IssuedBy Microsoft Root Authority
ExpiresBy: Microsoft Root Authority
SHA1Hash A3245CA9520DD6C96880E292DD85E2671CAE9E
SHA1 hash: A2D57D63CF331B177BE147088FEABEC7388BE01D
Issued to: Microsoft Windows Component Publisher
Successfully verified: Microsoft Windows Verification Intermediate PCA
Expires: 10/06/2009 23:07:51
Number of files successfully verified: 271F375A249EE9DE2D8E1AA363
Number of warnings: 0
Number of errors: 0

```

# Lower-level metadata

- Caveats

- PE Version information

- Detached signatures

- Where are .cat files?


- How to use .cat files?

- File type based on filename extension

Symantec
Symantec Corp
Symantec Corp.
Symantec Corporation
Symantec, Inc.
Symantec Ltd.
Symantec New Zealand Limited
Symantec Corporationn
Symantec Corporation <a href="http://www.symantec.com">http://www.symantec.com</a>

# Outline

- The paradigm change
- Clean Data sets
- Clean data generation
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



...res.dll
...seq.dll
...mgr.exe
...chm
...fg.dll
...fg.exe
...fg.ril
...ntDiag.exe
...brd.exe
...srv.exe
...pack.exe
...api.dll
...ncfg32.dll
...mid.exe
...ndial32.dll
...MDIALOG.SRG
...nd32.exe
...ndlib.wsc
...mgr32.hlp
...mon32.exe
...mos.ram
...mb32.dll
...mcs.dll
...mcs2.dll



... Application Extension	23/08/2014
108 KB Application Extension	14/04/2014
487 KB Application Extension	14/04/2014
63 KB Application	14/04/2014
71 KB Compiled HTML Help...	29/08/2014
76 KB Application Extension	14/04/2014
20 KB Application	14/04/2014
24 KB RLL File	13/04/2014
68 KB Application	19/03/2014
101 KB Application	14/04/2014
33 KB Application	14/04/2014
49 KB Application	28/02/2014
57 KB Application Extension	14/04/2014
16 KB Application Extension	14/04/2014
380 KB Application	14/04/2014
336 KB Application Extension	14/04/2014
1 KB SRG File	07/05/1980
25 KB Application	14/04/2014
40 KB Windows Script Co...	29/08/2014
60 KB Help File	29/08/2014
39 KB Application	14/04/2014
1 KB RAM File	29/08/2014
14 KB Application Extension	29/08/2014
181 KB Application Extension	14/04/2014
13 KB Application Extension	14/04/2014

# Higher-level metadata

- Nature

- Publisher
- Software name
- Software version
- Locale
- Source of software/file

Symantec
Norton Internet Security
2006 9.1.0.33
German
Read-only physical media (e.g. CD/DVD)

- Sourcing

- Targeted datasets
  - known in advance
  - tracked
- Non-targeted datasets
  - Profiling techniques



# Higher-level metadata (profiling)

## 1. Cross-reference

- Pros
  - Accurate
  - Consistent
- Cons
  - Requires annotated set
  - Quality of the reference set

# Higher-level metadata (profiling)

## 2. PE Version information

- Pros
  - Easily automated
  - Extracted directly from files
  - Could be interpolated on non-PE files
- Cons
  - Extracted only from PE files
  - Not accurate. A file may be shipped with various applications
  - Inconsistent / missing

# Higher-level metadata (profiling)

## 3. Web

- **P** [??0x80070422:???????,????????????????????? ...](#)  
DLL - MFCDLL Shared Library - Retail Version - **f35a584e947a5b401feb0fe01db4a0d7** O40 - Explorer.EXE - Microsoft Corporation - C:\Program Files\Common ...  
[bbs.winos.cn/viewthread.php?tid=35163](#) - 48k - Cached - Similar pages - Note this
- [MFC71.dll - File Information \[Is mine authentic?\]](#)  
File Hash: **F35A584E947A5B401FEB0FE01DB4A0D7**. File HashType: MD5. File Size: 1.01 MB (1060864 Bytes). File Modified Time: 8/4/2006 4:00:00 PM. Date Added: ...  
[www.programchecker.com/file/8466.aspx](#) - 76k - Cached - Similar pages - Note this
- **C** [AOFroobs.com :: View topic - Problem With Loader](#)  
27 Aug 2007 ... **MFC71.dll:f35a584e947a5b401feb0fe01db4a0d7**  
msvcp71.dll:561fa2abb31dfa8fab762145f81667c2  
— msvcr71.dll:86f1895ae8c5e8b17d99ece768a70732 ...  
— [www.aofroobs.com/forum/viewtopic.php?t=5309](#) - 75k - Cached - Similar pages - Note this
- [Spyware Browser Antispyware, adware and malware removal.](#)  
**f35a584e947a5b401feb0fe01db4a0d7**, MFCDLL Shared Library - Retail Version. MFC71.DLL  
Microsoft Corporation. Classified as: Good application / MFC ...  
[www.antispyware.tv/reports/-/231/?cmd=1](#) - 185k - Cached - Similar pages - Note this

# Higher-level metadata (profiling)

## 4. Source i.e. URL

- Pros

FTP source:  
— Easy for files gathered by crawling the internet  
`ftp://gamefiles.blueyonder.co.uk/beyondergames/unrealtournament/modifications/defencealliance/patches/win32/dabeta1.exe`

This would indicate that the file is a patch for the game Unreal Tournament published by Epic Games FTP

- Does not establish all metadata, mostly either publisher or software
- Difficult for automation but not impossible.

FTP source:  
`ftp://ftp.dell.com/sas-non-raid/BR123319.exe`

This would indicate that the file is a Dell driver

# Outline

- The paradigm change
- Clean Data sets
- Clean data generation
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- **Benefits of profiling**
- Summary



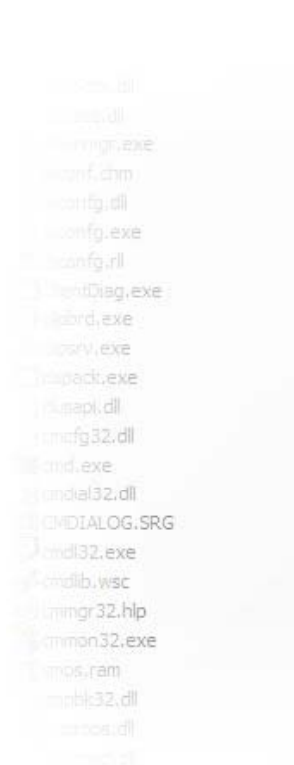
Size	File Name	File Type	Date
33 KB	Application Extension	Application Extension	29/08/2014
108 KB	Application Extension	Application Extension	14/04/2014
487 KB	Application Extension	Application Extension	14/04/2014
63 KB	Application	Application	14/04/2014
71 KB	Compiled HTML Help...	Compiled HTML Help...	29/08/2014
76 KB	Application Extension	Application Extension	14/04/2014
20 KB	Application	Application	14/04/2014
24 KB	RLL File	RLL File	13/04/2014
68 KB	Application	Application	19/03/2014
101 KB	Application	Application	14/04/2014
33 KB	Application	Application	14/04/2014
49 KB	Application	Application	28/02/2014
57 KB	Application Extension	Application Extension	14/04/2014
16 KB	Application Extension	Application Extension	14/04/2014
380 KB	Application	Application	14/04/2014
336 KB	Application Extension	Application Extension	14/04/2014
1 KB	SRG File	SRG File	07/05/1999
25 KB	Application	Application	14/04/2014
40 KB	Windows Script Co...	Windows Script Co...	29/08/2014
60 KB	Help File	Help File	29/08/2014
39 KB	Application	Application	14/04/2014
1 KB	RAM File	RAM File	29/08/2014
14 KB	Application Extension	Application Extension	29/08/2014
181 KB	Application Extension	Application Extension	14/04/2014
13 KB	Application Extension	Application Extension	14/04/2014

## Benefits of profiling

- Analysis & fitness for business (Quality Assurance)
- Making purpose built sets
- Improvement / extension of existing data sets
- Easy access to samples (hash lookup)
- Detailed information for engineers and customers

# Outline

- The paradigm change
- Clean Data sets
- Clean data generation
- Metadata
  - Lower-level metadata
  - Higher-level metadata
- Benefits of profiling
- Summary



Size	Type	Date
33 KB	Application Extension	29/08/2014
108 KB	Application Extension	14/04/2014
487 KB	Application Extension	14/04/2014
63 KB	Application	14/04/2014
71 KB	Compiled HTML Help...	29/08/2014
76 KB	Application Extension	14/04/2014
20 KB	Application	14/04/2014
24 KB	RLL File	13/04/2014
68 KB	Application	19/03/2014
101 KB	Application	14/04/2014
33 KB	Application	14/04/2014
49 KB	Application	28/02/2014
57 KB	Application Extension	14/04/2014
16 KB	Application Extension	14/04/2014
380 KB	Application	14/04/2014
336 KB	Application Extension	14/04/2014
1 KB	SRG File	07/05/1999
25 KB	Application	14/04/2014
40 KB	Windows Script Co...	29/08/2014
60 KB	Help File	29/08/2014
39 KB	Application	14/04/2014
1 KB	RAM File	29/08/2014
14 KB	Application Extension	29/08/2014
181 KB	Application Extension	14/04/2014
33 KB	Application Extension	14/04/2014

# Summary

- Profiling :
- What is the metadata?
- How to obtain the metadata?
  - Lower level is in the file
  - Higher level:
    - known for targeted sets
    - 4 techniques for non-targeted sets:
      - Cross Reference
      - PE Version Information
      - Web Search
      - Source



# Summary

- Clean Data Sets: what are they for?
  - Currently: FP prevention
  - Future: potential paradigm change?

# Questions



Clean data is becoming more and more important

Metadata is vital

Several profiling techniques exist



Confidence in a connected world.



# Thank You!

[Bartlomiej\\_uscilowski@symantec.com](mailto:Bartlomiej_uscilowski@symantec.com)

[Julie\\_weber@symantec.com](mailto:Julie_weber@symantec.com)