

# Empirical Analysis of Factors Affecting Malware URL Detection

**Marie Vasek** & Tyler Moore  
SMU, Dallas, TX  
mvasek@smu.edu

eCrime Researchers Summit

September 17, 2013

# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

Conclusion

# So Much Malware



# So Much Malware



# So Many AntiVirus Products



# Research Questions

**Our hypothesis:** At least part of the difference in performance across blacklists can be explained by characteristics of the malware URL or the blacklisting service.

1. Impact of type of malware URL on blacklisting?
  - Payload type (exploit kit, Fake AV, etc.)
  - URL features (IP address, path)
2. Impact of type of blacklisting service?
  - Blocking or advisory
  - Cost

# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

Conclusion

# Malware Domain List

- Publicly Available
- Low Volume
  - 722 URLs in our collection
  - *BUT* it's free, open, & no data-sharing issues

M A L W A R E   D O M A I N   L I

[Homepage](#) | [Malware Domain List](#) | [Forums](#) | [RSS update feed](#) | [Contact us](#)

WARNING: All domains on this website should be considered dangerous. If you not know what you are doing here, it is recommended you leave right away. This website is a resource for security professionals and enthusiasts.

Date (UTC)	Domain	IP	Reverse Lookup	Description
2013/07/11_12:54	www.keepsaketributes.com/calc/images/ac0094cbbe/? ==wMw1mLulWYt9VbzxXO5IDMxIzN3QD O5UDO2w3LlJmYjRTOWAz Yh9ycldWYtL2LjxWYj9S bvNmLzVGd1JWayRXZrF2 cwVWZr5yd3d3LvoDc0RHa8NnZ	69.16.206.133	host2.jbatkins.com.	exploit kit
2013/07/11_12:54	www.lowes-pianos-and-organs.com/images/d521c2a038/? zAXbu4Wah12XtNHf0YTNwgTNwIjM0 ATMzIDfvGzMwEmMjFjM1	67.222.109.112	d15.altserver.com.	ycuF2Zy9WLk Waw1ycld3bs Dc0RHa8NnZ



# VirusTotal

- Publicly Available
- API accessible
- Google-Owned
- Checks each URL against 38 blacklists



# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

Conclusion

# What We Do

1. Check MDL for new URLs.
2. Scan new URLs through VT.
3. Scan new URL every other hour for 48 hours after.
4. Afterwards, scan every day for 14 days.

# Table of Contents

Intro

Data Sources

Methodology

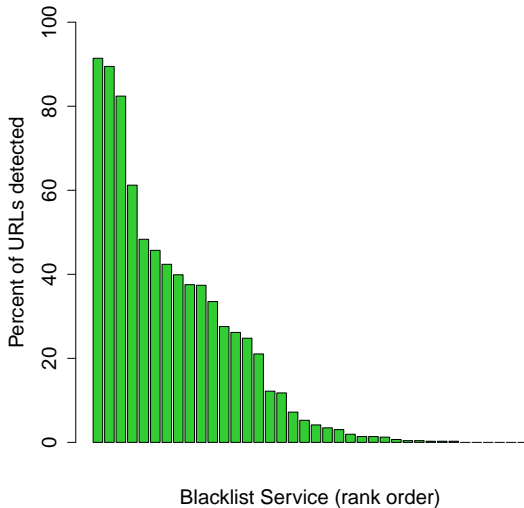
**Data**

Factors

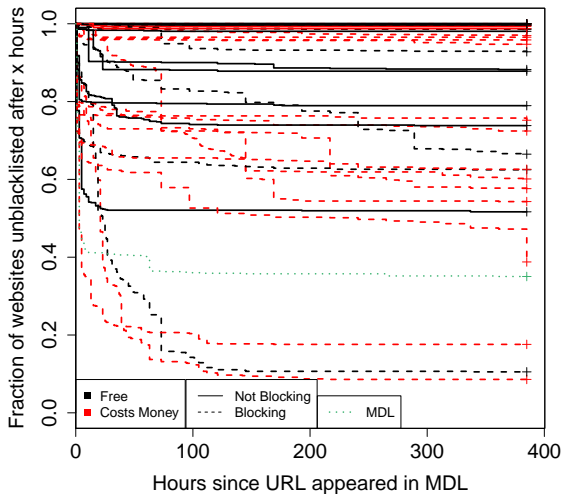
Results

Conclusion

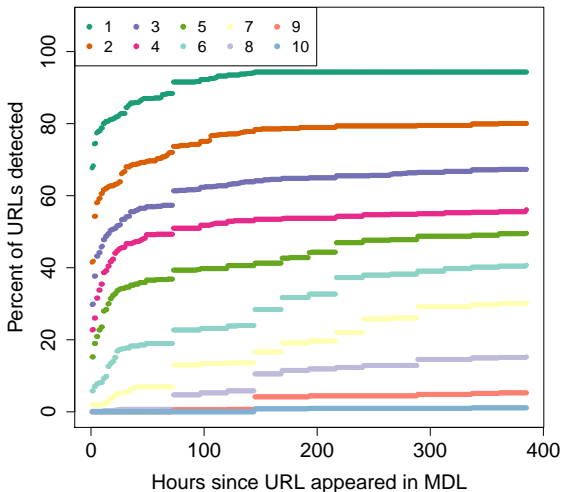
# How Blacklisting Detection Rates Vary



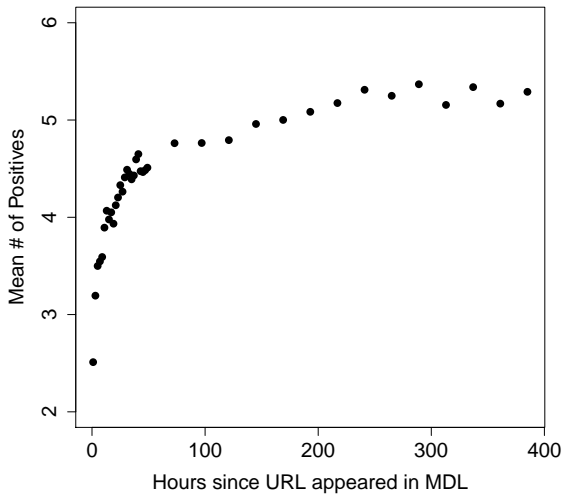
# How Blacklisting Detection Rates Vary



# Percentage of URLs Detected by at least $n$ Services



# Blacklisting over Time





# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

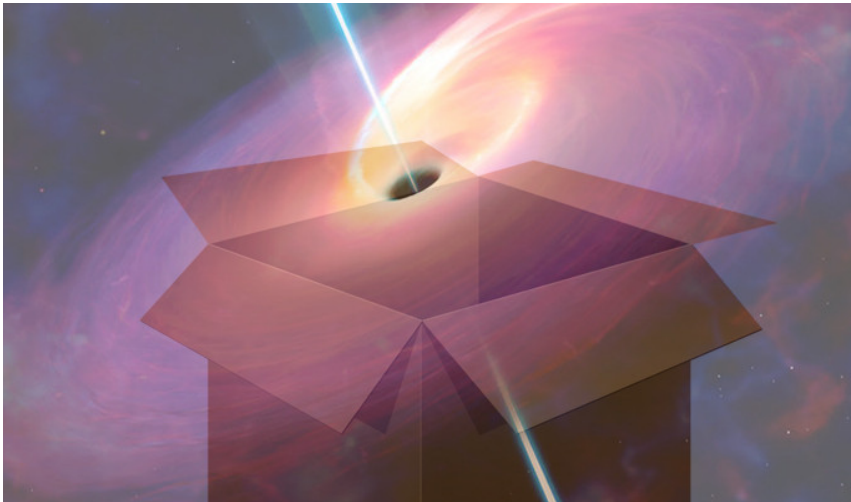
Conclusion

# Research Questions

**Our hypothesis:** At least part of the difference in performance across blacklists can be explained by characteristics of the malware URL or the blacklisting service.

1. Impact of type of malware URL on blacklisting?
  - Payload type (exploit kit, Fake AV, etc.)
  - URL features (IP address, path)
2. Impact of type of blacklisting service?
  - Blocking or advisory
  - Cost

# Malware Type



IP/Domain

127.0.0.1/bad.html

icanhazcats.com/milk.php

Path?

fluffybunnies.org/

goldphish.net/bubbles.php

# Type of Malware

Malware Type	#	%	IP/Domain	#	%	Path?	#	%
Executable	175	24%	IP Address	124	17%	Has Path	675	93%
Fake AV	65	9%	Domain	598	83%	No Path	47	7%
Styx	51	7%						
Blackhole Lnd.	149	21%						
Other	282	39%						

# Blocks?

## Web Page Blocked

Access to the web page you were trying to visit has been blocked in accordance with company policy. Please contact your system administrator if you believe this is in error.

**User:** 129.119.9.166

**URL:** [www.malwareconference.org/index.php?option=com\\_docman\\_task=cat\\_view\\_gid=79\\_Itemid=51](http://www.malwareconference.org/index.php?option=com_docman_task=cat_view_gid=79_Itemid=51)

**Category:** malware-sites

# Costs Money





## Blacklist Types

Blocks?	Costs?	Description
✓	✓	Paid AV
✓	✗	Free AV & Web Browser Blacklists
✗	✓	Malware Site Checkers
✗	✗	Malware URL lists

## Blacklist Types

Blocks?	#	%	Costs?	#	%
Blocks Users	22	58%	Costs Money	17	45%
Doesn't Block	16	42%	Free	21	55%

# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

Conclusion

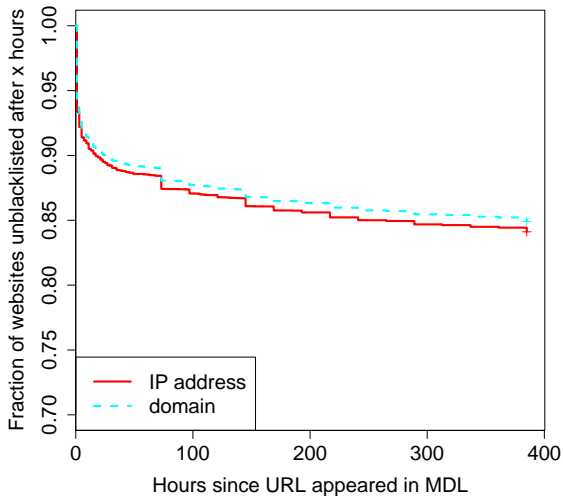
# Logistic Regression: URL ever Blacklisted by Service?

	Coefficient	Odds Ratio	p value
<b>Intercept</b>	-2.304	<b>0.100</b>	0.000
<i>URL features</i>			
IP address	-0.004	0.996	0.945
<b>Has a Path?</b>	-0.196	<b>0.822</b>	0.031
<b>Executable</b>	1.017	<b>2.765</b>	0.000
<b>Fake AV</b>	-0.838	<b>0.433</b>	0.000
<b>Styx</b>	0.746	<b>2.109</b>	0.000
<b>Blackhole Landing Page</b>	0.196	<b>1.217</b>	0.000
<i>Malware Blacklist Features</i>			
<b>Blocks Users?</b>	0.611	<b>1.843</b>	0.000
<b>Costs Money</b>	0.298	<b>1.347</b>	0.000
$\chi^2 = 1144.501, p \text{ value} = 0.000$			

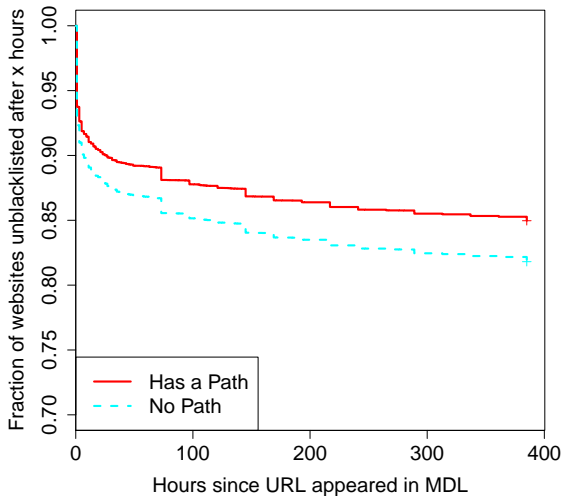
## Survival Regression: Blacklist Timing

	Coefficient	Odds Ratio	<i>p</i> value
<i>URL features</i>			
IP address	0.056	1.058	0.210
<b>Has a Path?</b>	-0.207	<b>0.811</b>	0.012
<b>Executable</b>	0.896	<b>2.449</b>	0.000
<b>Fake AV</b>	-0.814	<b>0.443</b>	0.000
<b>Styx</b>	0.750	<b>2.118</b>	0.000
<b>Blackhole Landing Page</b>	0.179	<b>1.196</b>	0.000
<i>Malware Blacklist Features</i>			
<b>Blocks Users?</b>	0.538	<b>1.713</b>	0.000
<b>Costs Money</b>	0.300	<b>1.351</b>	0.000
$R^2 = 0.055$			

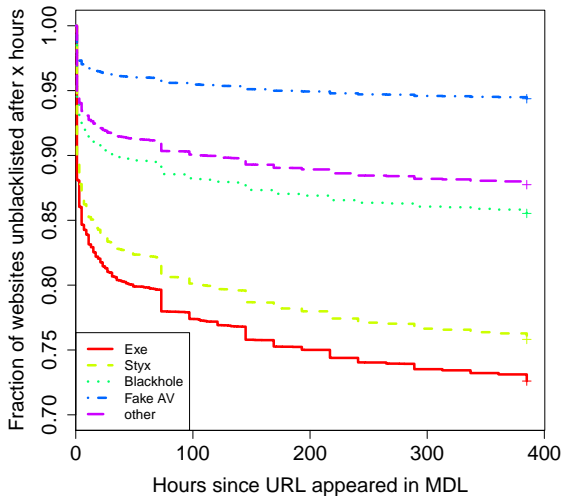
# Survival Regression



# Survival Regression

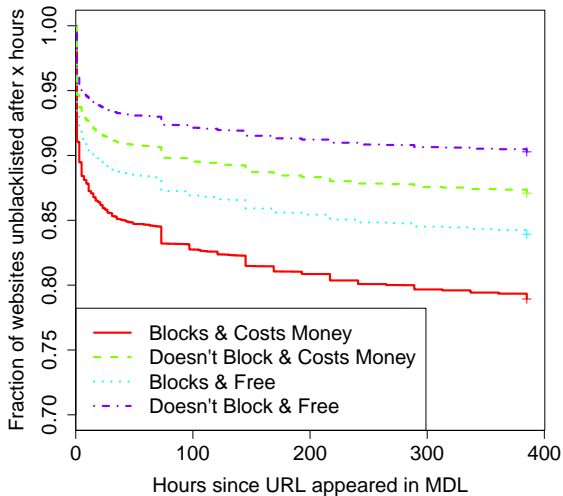


# Survival Regression





# Survival Regression



# Table of Contents

Intro

Data Sources

Methodology

Data

Factors

Results

Conclusion

# Summary

- Known Exploit kits more likely to be blacklisted
  - ... but Fake AV is not ..?
- Executables more likely to be blacklisted
- Blocking users → more likely to blacklist
- Costing money → more likely to blacklist

# Limitations

- MDL → low volume
  - other exploit kits
  - AS / DNS features
  - Malicious vs. Hacked?
- MDL → used as input to blacklists
  - Blacklisted b/c found on VT or found organically?
  - Are things on MDL still bad on day  $n$ ?
- VT → not representative sample of malware blacklists
  - Blacklists not on VT?
  - Blacklists on VT that didn't blacklist anything?

Questions?

