



# Introduction of Cybersecurity AI dataset In Korea

2022.10.20

Dr. Lee Jeong Min

# Contents

- I Cybersecurity AI Dataset Introduction
- II Best Practices
- III Conclusion

|

Cybersecurity  
AI Dataset  
**Introduction**

# 1. AI Dataset for Cybersecurity



## 01 AlphaGO(2016)

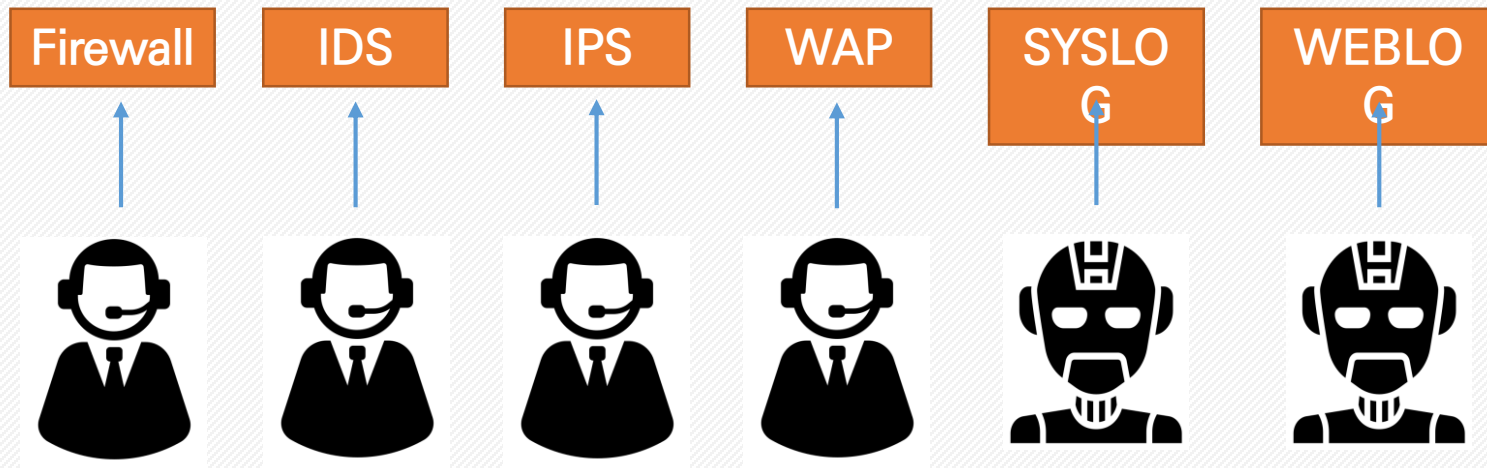


# 1. AI Dataset for Cybersecurity



## 02 Need for AI in Cybersecurity

- ➔ Expansion of Security data(Firewall, IDS, IPS , WAP, System Log, Web Log, etc.)



# 1. AI Dataset for Cybersecurity



## 03 Require high-quality cybersecurity AI learning datasets

Artificial intelligence will contribute to saving money and time by autonomously identifying or responding to potential cyber attacks (WEF, 2020)



👉 The results of the AI model determine the quality of the learning dataset

# 1. AI Dataset for Cybersecurity



## 04 Result of Cybersecurity AI Dataset(1/2)



### Improving the code vs. the data

	Steel defect detection	Solar panel	Surface inspection
Baseline	76.2%	75.68%	85.05%
Model-centric	+0% (76.2%)	+0.04% (75.72%)	+0.00% (85.05%)
Data-centric	+16.9% (93.1%)	+3.06% (78.74%)	+0.4% (85.45%)

Source : Andrew Ng(2021)

# 1. AI Dataset for Cybersecurity



## 04 Result of Cybersecurity AI Dataset(1/2)

Category		Dataset based Diagnostic Name of Malware	Dataset based Attribute Name of Malware	Dataset based Malware of resent Social Issues
Malware	Number Of Data	3 Hundred Million	1 Hundred Million	120,000
	File Type	More than 24 file types Include EXE	More than 30 file types Include PDF	More than 10 file types Include EXE,
	Labelling	More than 11,800 Family Dataset	More than 3,717 Similar/Variant Dataset	26 Keyword of recent Incidents



# 1. AI Dataset for Cybersecurity



## 04 Result of Cybersecurity AI Dataset(2/2)

Category		Dataset based Detection of Cybersecurity Incident	Dataset based Cyber Attack Tactics	Dataset based Reenact of recent Cybersecurity Incident
Incident	Number Of Data	2 Hundred Million	1 Hundred Million	1.2 Hundred Million
	Environment	6 Heterogeneous Equipment Operating Environment (FW, IDS, IPS, WAF, etc)	Utilize an automated malware analysis platform	6 Heterogeneous Equipment Operating Environment (FW, IDS, IPS, WAF, etc)
	Labelling	More than 17 normal and aggressive acts	More than 230 types based on attack techniques	15 Cybersecurity Incident Scenarios



# Best Practices



### 01 Best Practices of Cybersecurity AI Dataset



## Sharing 8 Best Practices

- ▶ Cooperate with Private/Public Cybersecurity Organizations
- ▶ Spam Filtering with Teleco, Malware Detection with Game Publisher, Intrusion Detection with Monitoring Org.(CERT), Etc.

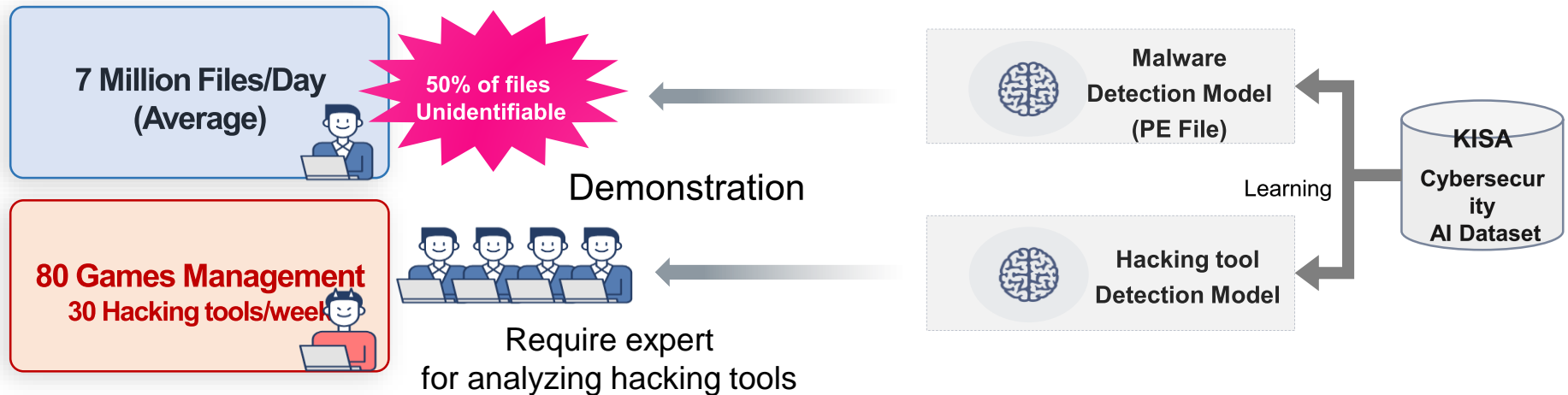
👉 <https://www.youtube.com/watch?v=nVijOJD3Efk>

# 2. Best Practices



## 1 Best Practices of AI Dataset – ① Nexon | Game publisher

- **[Demonstration]** Development of AI model for detections/classification of malicious files and hack-tools
- **[Result]** Reducing manual analysis time more than 70%, and Increasing business efficiency more than 30%
  - The effect of preventing large-scale infection for PCs and mobile devices of hundreds of millions of users by supplementing the undetected areas missed by pattern-based anti-virus.



# 2. Best Practices

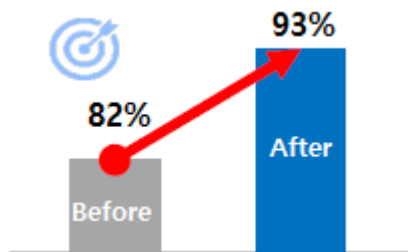


## 2 Best Practices of AI Dataset – ② KLID\* | Public SOC

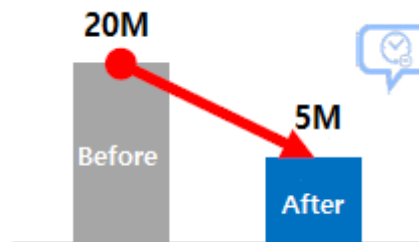
\*KOREA LOCAL INFORMATION RESEARCH & DEVELOPMENT INSTITUTE

- **[Demonstration]** Applying AI Dataset for to security monitoring model for 17 local governments in Korea
- **[Result]** Increase AI model detection performance more than 5 to 30% by learning the latest intrusion scenario dataset
  - The effect of proactive defend against latest threats such as spear phishing, penetration into the internal network and etc. that former AI models have not detected.

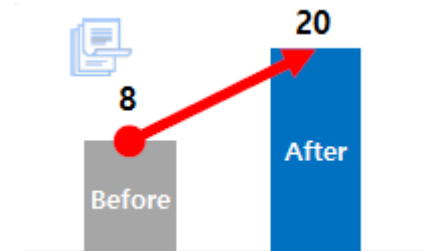
### Increase Accuracy



### Reduce Response Time



### Increase Detection



## 2. Best Practices



### 3 Best Practices of AI Dataset – ③ KT | Telecommunication Corp.

- **[Demonstration]** Improving AI model accuracy for detecting malwares attached to an e-mail
- **[Result]** Strengthen learning capabilities by applying additional datasets to malware detection model operated by KT
  - As a result of re-learning AI with additional datasets, the detection rate is increased (83% → 92.6%), and the count of error detection files is decreased (224 → 68)

	TPR	TP	FP	TN	FN
<b>Original Model</b>	83%	632	6	224	491
<b>New Model (KISA Dataset)</b>	92.6%	788	32	68	465

True Positive : Detect actual malicious code as malicious code

False Positive : Detects actual normal code as malicious code

True Negative : Detects actual normal code as normal code

False Negative : Detect actual malicious code as normal code



# Conclusion



## 01 2022 Cybersecurity AI Datasets

6 Hundred Million

Application

Active  
Monitoring

Threat  
Profiling



# 3. Conclusion



Q

Our Cybersecurity AI Dataset is good in South Korea Environment,  
But how about in other countries?



Make standard procedure and standard format  
about Cybersecurity AI dataset.



# Q/A

[jmlee@kisa.or.kr](mailto:jmlee@kisa.or.kr)