

MADAR: Efficient Continual Learning for Malware Analysis with Distribution-Aware Replay

Saidur Rahman^{1,3}, Scott Coull², Qi Yu³, Matthew Wright³

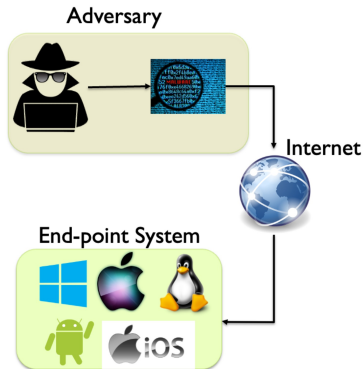
¹University of Texas at El Paso (UTEP), ²Google

³Rochester Institute of Technology

The Conference on Applied Machine Learning in Information Security (CAMLIS)
October 23, 2025



Cybercrime and Malware



200 Million to 1.2 Billion in 10 years Growth
600%

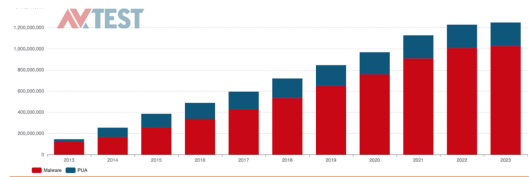
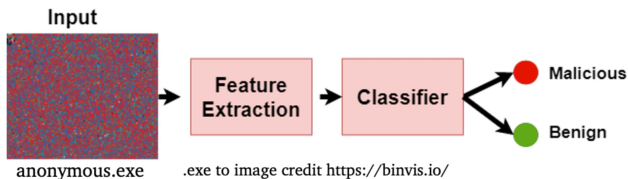


Figure: Growth of Malware and Potential Unwanted Applications (PUA)¹

¹ <https://www.av-test.org/en/statistics/malware/>

Malware Analysis and Machine Learning

- Supervised Machine Learning (ML)
- Static malware analysis
 - Computational efficiency
 - Easy-to-Scale
 - Existing expert knowledge
- Significant performance
 - LightGBM on EMBER²
 - ROC AUC 0.996



²H. S. Anderson and P. Roth, "EMBER: an open dataset for training static pe malware machine learning models," arXiv, 2018. ▶ ◀ ≡ 🔍 ↺

Ever Evolving Growth of Malware

- AV-TEST \Rightarrow 450K *new* malware and PUA *each day*¹
- VirusTotal \Rightarrow 1.8M *unique* software samples *each day*³



¹ <https://www.av-test.org/en/statistics/malware/>

³ VirusTotal, <https://www.virustotal.com/gui/intelligence-overview>

Ever Evolving Growth of Malware

- AV-TEST \Rightarrow 450K *new* malware and PUA *each day*¹
- VirusTotal \Rightarrow 1.8M *unique* software samples *each day*³

Huge data volumes drive up costs and training times



¹ <https://www.av-test.org/en/statistics/malware/>

³ VirusTotal, <https://www.virustotal.com/gui/intelligence-overview>

Less than Ideal Solutions

Expanding Training Effort

expend tremendous effort to frequently retrain over all the data



Less than Ideal Solutions

Expanding Training Effort

expend tremendous effort to frequently retrain over all the data

Remove Older Samples

allows attackers to revive older malware instead of writing new ones



Figure: from ⁴

⁴ <http://www.martybucella.com/E199.gif>

Less than Ideal Solutions

Expanding Training Effort

expend tremendous effort to frequently retrain over all the data

Remove Older Samples

allows attackers to revive older malware instead of writing new ones

Expanding Training Effort

at the cost of not adjusting to changes in the distribution



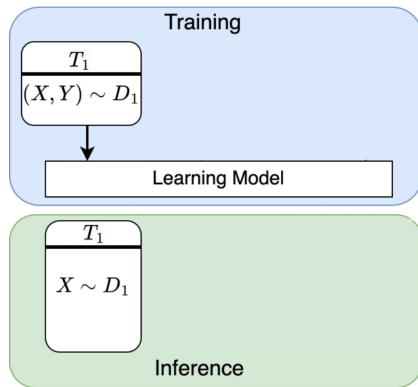
Figure: from ⁴



⁴ <http://www.martybucella.com/E199.gif>

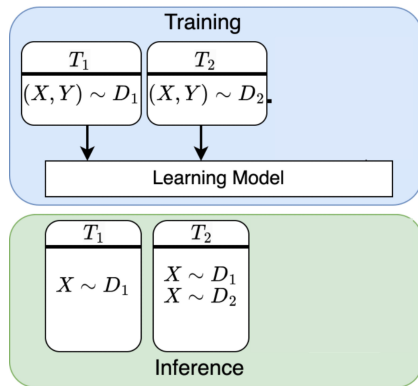
Continual Learning

- Acknowledges
 - Continuous distributional shift
- Non-stationary data
 - Observed periodically
(T_1, T_2, \dots, T_N)
 - Different data distribution in each period (D_1, D_2, \dots, D_N)
 - Data from each period is referred to as task
 - $task_N \in (T_N, D_N)$
 - New class/ new samples/
new objective



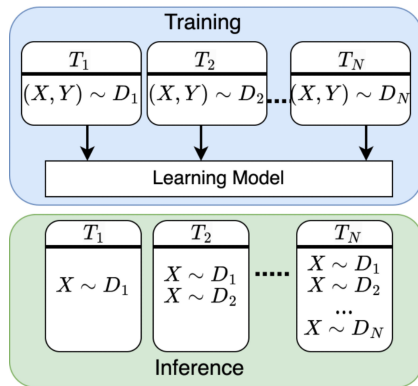
Continual Learning

- Acknowledges
 - Continuous distributional shift
- Non-stationary data
 - Observed periodically (T_1, T_2, \dots, T_N)
 - Different data distribution in each period (D_1, D_2, \dots, D_N)
 - Data from each period is referred to as task
 - $task_N \in (T_N, D_N)$
 - New class/ new samples/ new objective



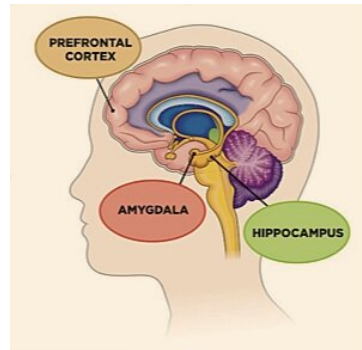
Continual Learning

- Acknowledges
 - Continuous distributional shift
- Non-stationary data
 - Observed periodically (T_1, T_2, \dots, T_N)
 - Different data distribution in each period (D_1, D_2, \dots, D_N)
 - Data from each period is referred to as task
 - $task_N \in (T_N, D_N)$
 - New class/ new samples/ new objective



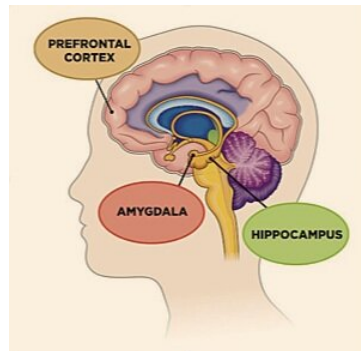
Continual Learning

- Inspired by human learning process
 - Continuous learning
 - Observe and learn
 - Storage → abstract representation in the hippocampus
- Relax the need to store all the data
 - Reduce storage cost
- Reduce computational cost



Continual Learning

- Inspired by human learning process
 - Continuous learning
 - Observe and learn
 - Storage → abstract representation in the hippocampus
- Relax the need to store all the data
 - Reduce storage cost
- Reduce computational cost



Challenge → Catastrophic Forgetting

Forgetting would reintroduce vulnerabilities

Catastrophic Forgetting (CF)

Neural Networks suffer from catastrophic forgetting⁵

- Forget the old tasks, unlikely to happen in human learning

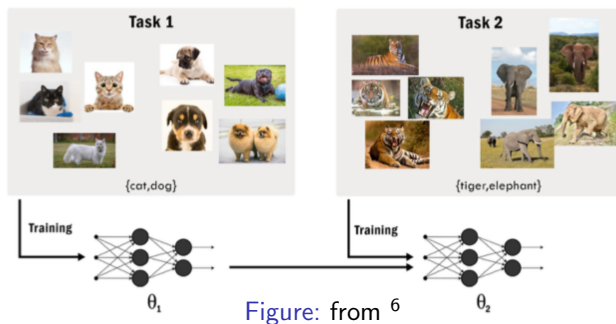


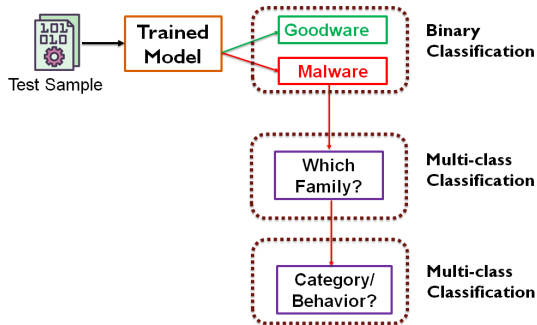
Figure: from ⁶

⁵ McCloskey and Cohen, Catastrophic interference in connectionist networks: The sequential learning problem, Psychology of learning and motivation, 1989.

⁶ https://mrifkikurniawan.github.io/blog/2021/Catastrophic_Forgetting_in_Neural_Networks_Explained

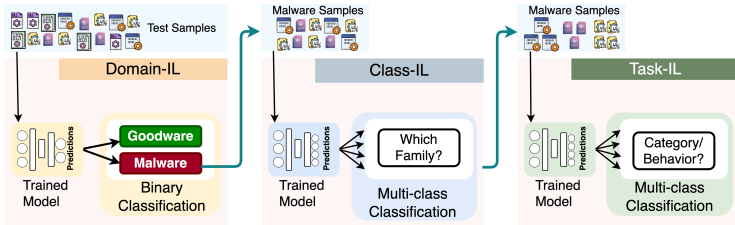
Malware Classification Pipeline

- Family
 - Citadel
 - Observe and learn
 - Gameover
 - Cthonic, and so on
- Category/Behavior
 - Adware
 - Ransomware
 - Banking Trojan
 - Backdoor, and so on



CL in Malware Classification Pipeline

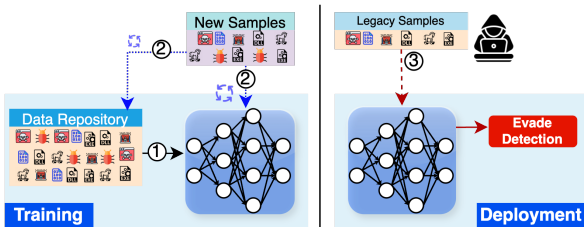
- Domain Incremental Learning (**Domain-IL**)
 - Distribution shift
 - Emergence of new malware
- Class Incremental Learning (**Class-IL**)
 - New malware family
- Task Incremental Learning (**Task-IL**)
 - New malware category



Threat Model

Retrograde Malware Attack (RMA)

- ① Initial Training and Deployment
- ② Incremental Updates and Forgetting
- ③ Attack Phase

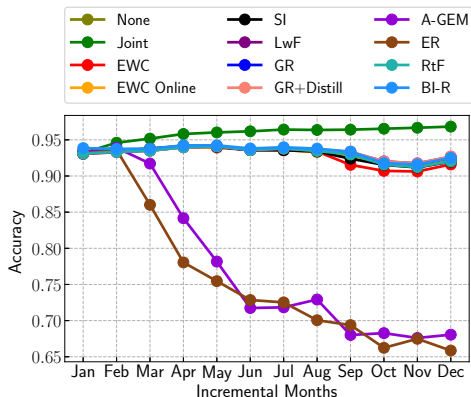


EMBER Domain-IL

- Benchmarks

- None → No CL techniques applied
- Joint → Static training (training over accumulated data)

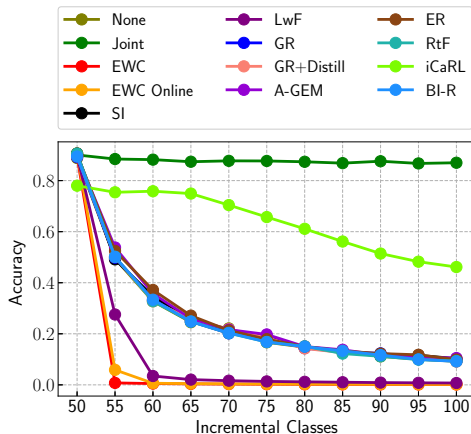
None of the CL techniques are effective in the Domain-IL setting



⁷ Rahman, M.S., Coull, S. and Wright, M., On the limitations of continual learning for malware classification. In Conference on Lifelong Learning Agents, 2022.

EMBER Class-IL

- 10 of the 11 methods performed poorly
- Only iCaRL performing marginally better against the Joint replay baseline



⁷ Rahman, M.S., Coull, S. and Wright, M., On the limitations of continual learning for malware classification. In Conference on Lifelong Learning Agents, 2022.

Malware samples for each task

- Belong to multiple families
- Indicating sub-distributions within malware distribution

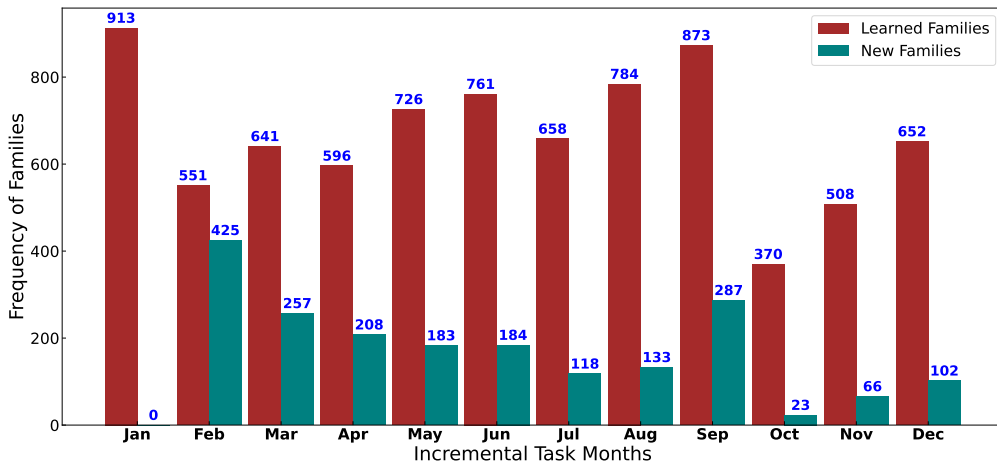
Task	#of Goodware	#of Malware	#of Malware Families
January	29423	32491	913
February	22915	31222	976
March	21373	20152	898
April	25190	26892	804
May	23719	22193	909
June	23285	25116	945
July	24799	26622	776
August	23634	21791	917
September	26707	37062	1160
October	29955	56459	393
November	50000	50000	574
December	50000	50000	754

Malware samples for each task

- Belong to multiple families
- Indicating sub-distributions within malware distribution

Task	#of Goodware	#of Malware	#of Malware Families
January	29423	32491	913
February	22915	31222	976
March	21373	20152	898
April	25190	26892	804
May	23719	22193	909
June	23285	25116	945
July	24799	26622	776
August	23634	21791	917
September	26707	37062	1160
October	29955	56459	393
November	50000	50000	574
December	50000	50000	754

Emergence of New Families in each Task



Summary of Exploratory Analysis

- Malware distribution in each task
 - Contains multiple sub-distributions
 - On an average around 800 families
- Lot of new novel families emerge
 - Old families observed infrequently
- Substantial #of samples wo/ AV class labels
- Priorities change over time
 - Prominent families do not remain prominent

MADAR: Malware Analysis with Distribution-Aware Replay

- CL technique should capture both representative and discriminative samples^{8,9}
- Heterogeneity among the replay samples
 - Family based sample selection
 - To accommodate varying families
 - Representative samples
 - Samples closer to the cluster mean
 - Discriminative (outlier) samples
 - Samples farther away from the mean

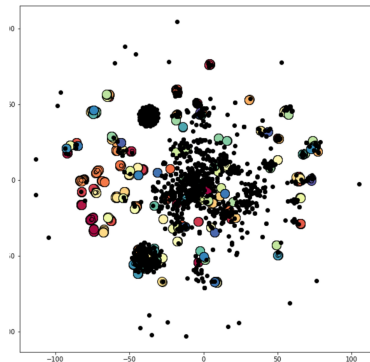
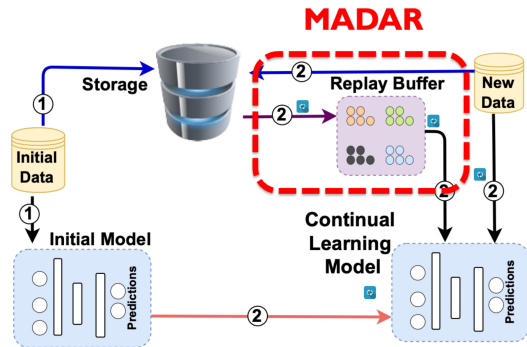


Figure: t-SNE projection of EMBER malware from January 2018

⁸ Aljundi, Rahaf, et al. "Gradient based sample selection for online continual learning." NeurIPS 2019.

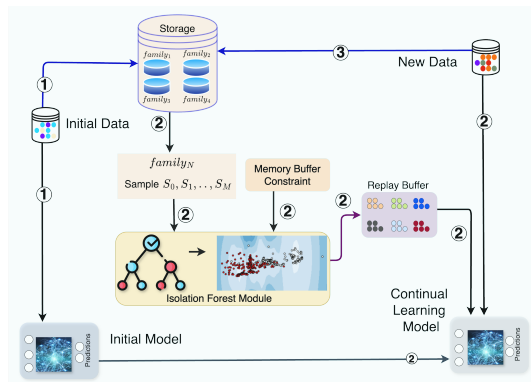
Replay-based CL for Malware Classification

- ① Initial Phase
 - Initialize model w/ available data
 - Store the available data
- ② CL Phase
 - Initialize model → CL Model
 - Replay some old data from the storage
 - Use (some/all) new data



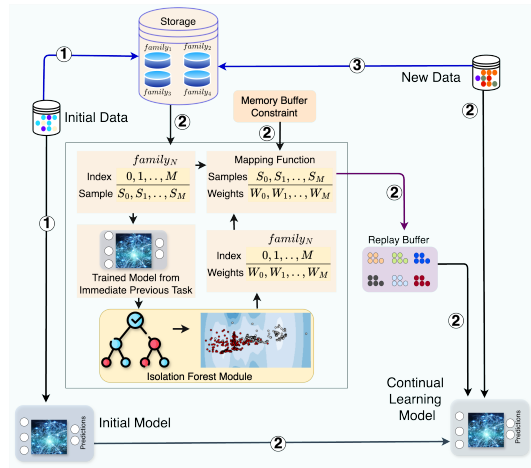
MADAR → Isolation Forest based Sampling

- ① Initial Phase
 - Initialize model w/ available data
 - Store the available data
- ② CL Phase
 - Initialize model → CL Model
 - IFS Module
 - Replay Buffer



MADAR^θ → Anomalous Weights based Sampling

- Hidden representation
 - Weights of the model
 - Anomalous and Similar weights
 - Backtrack to raw feature space
- Low dimension
 - Faster to process than raw feature space
 - i.e., 2381 → 256 (for EMBER)



Evaluation → MADAR and MADAR^θ in Domain-IL

Group	Method	EMBER Budget				AZ Budget			
		1K	100K	200K	400K	1K	100K	200K	400K
Baselines	Joint	96.4±0.3				97.3±0.1			
	None	93.1±0.1				94.4±0.1			
	GRS	93.6±0.3	95.3±0.7	95.9±0.1	96.0±0.3	95.3±0.1	97.1±0.1	97.1±0.1	97.2±0.1
Prior Work	ER	80.6±0.1	69.9±0.1	70.0±0.1	70.0±0.1	40.4±0.1	42.6±0.1	44.0±0.1	48.6±1.1
	AGEM	80.5±0.1	70.0±0.1	70.0±0.2	70.0±0.1	45.4±0.1	53.7±0.6	54.2±0.3	56.7±0.3
	GR	93.1±0.2				93.3±0.4			
	RtF	93.2±0.2				93.4±0.2			
	BI-R	93.4±0.1				93.5±0.1			
Ours	MADAR-R	93.7±0.1	95.3±0.6	96.0±0.1	96.1±0.1	95.8±0.1	97.0±0.1	97.0±0.1	97.0±0.1
	MADAR-U	93.6±0.2	95.3±0.1	95.5±0.1	95.8±0.1	95.7±0.1	95.2±0.1	95.4±0.1	96.3±0.2
	MADAR ^θ -R	93.6±0.1	95.8±0.1	96.1±0.1	96.1±0.1	95.8±0.2	96.9±0.1	97.1±0.1	97.2±0.1
	MADAR ^θ -U	93.5±0.2	95.2±0.2	95.6±0.1	95.7±0.1	95.6±0.1	96.8±0.1	97.0±0.1	97.1±0.1

Evaluation → MADAR and MADAR^θ in Class-IL

Group	Method	EMBER				AZ			
		Budget				Budget			
		100	1K	10K	20K	100	1K	10K	20K
Baselines	Joint	86.5±0.4				94.2±0.1			
	None	26.5±0.2				26.4±0.2			
	GRS	51.9±0.4	75.4±0.7	83.5±0.1	84.6±0.2	43.8±0.7	70.2±0.4	86.4±0.2	89.1±0.2
Prior Work	TAMiL	32.2±0.3	35.3±0.2	38.2±0.3	38.8±0.2	53.4±0.3	57.6±0.3	63.5±0.1	67.7±0.3
	iCaRL	53.9±0.7	60.0±1.0	64.6±0.8	66.8±1.1	43.6±1.2	61.7±0.7	81.5±0.6	84.6±0.5
	ER	27.5±0.1	28.0±0.1	28.0±0.1	28.2±0.1	50.8±0.7	58.9±0.2	62.9±0.7	64.2±0.4
	AGEM	27.3±0.1	27.7±0.1	28.2±0.1	28.2±0.1	27.3±0.7	27.1±0.3	28.2±1.0	28.0±0.8
	GR	26.8±0.2				22.7±0.3			
	RtF	26.5±0.1				22.9±0.3			
	BI-R	26.9±0.1				23.4±0.2			
	MalCL	54.5±0.3				59.8±0.4			
Ours	MADAR-R	68.0±0.4	76.0±0.3	83.2±0.2	84.0±0.2	59.4±0.6	71.9±0.5	86.3±0.1	89.1±0.1
	MADAR-U	66.4±0.4	79.4±0.4	84.8±0.1	85.8±0.3	57.3±0.5	76.2±0.2	89.8±0.1	91.5±0.1
	MADAR ^θ -R	67.9±0.3	72.7±0.5	83.2±0.1	84.5±0.2	58.8±0.3	71.0±0.7	85.1±0.2	88.1±0.1
	MADAR ^θ -U	67.5±0.3	78.5±0.4	85.3±0.1	86.2±0.2	58.5±0.7	74.7±0.2	88.7±0.1	90.7±0.1

Evaluation → MADAR and MADAR^θ in Task-IL

Group	Method	EMBER Budget				AZ Budget			
		100	1K	10K	20K	100	1K	10K	20K
Baselines	Joint	97.0±0.3				98.8±0.2			
	None	74.6±0.7				74.5±0.2			
Prior Work	GRS	86.9±0.3	93.6±0.3	94.7±0.3	95.0±0.1	85.2±0.1	90.8±0.1	93.5±0.1	95.2±0.1
	TAMiL	72.8±0.1	86.9±0.2	90.3±0.1	94.2±0.7	80.5±0.4	91.5±0.2	93.5±0.1	94.8±0.2
	ER	67.4±0.3	89.5±0.5	94.8±0.2	95.4±0.1	83.6±0.2	92.3±0.3	96.2±0.1	97.5±0.2
	AGEM	79.6±0.2	83.8±0.4	86.1±0.2	89.3±0.1	76.7±0.5	85.3±0.1	86.7±0.2	91.3±0.3
	GR	79.8±0.3				75.6±0.2			
	RtF	77.8±0.2				74.2±0.3			
	BI-R	87.2±0.3				85.4±0.2			
Ours	MADAR-R	92.1±0.2	93.8±0.2	94.8±0.2	95.6±0.1	86.0±0.3	92.4±0.1	96.7±0.1	97.9±0.2
	MADAR-U	93.4±0.2	93.9±0.3	95.6±0.1	95.8±0.2	88.1±0.3	94.5±0.3	98.1±0.1	98.7±0.1
	MADAR ^θ -R	93.1±0.2	93.6±0.1	94.6±0.2	94.7±0.3	87.3±0.3	93.2±0.2	95.9±0.1	96.9±0.1
	MADAR ^θ -U	93.2±0.1	93.8±0.2	94.8±0.1	95.5±0.3	87.9±0.2	93.6±0.1	97.2±0.2	98.1±0.1

Summary of the Key Findings

- Prior CL techniques → do not work well for malware tasks
 - Due to the complexity of the data and unique non-stationary nature
- Malware distribution represents heterogeneity among and within families
- MADAR: Distribution Aware Replay Technique
 - State-of-the-art performance
 - Domain-IL → Ratio variants (MADAR-R and MADAR^θ-R)
 - Class-IL and Task-IL → Uniform variants (MADAR-U and MADAR^θ-U)

Takeaways

- ① Evolving growth of malware is a challenging problem
 - Require an ever evolving and intelligent system for effective malware classification and detection
- ② Continual Learning (CL) is an ideal candidate
 - CV based CL systems fall short to mitigate catastrophic forgetting in malware domain
- ③ CL for malware domain →
 - Must consider the heterogeneous nature and complexities of malware data distribution
 - Lots of open research questions
- ④ MADAR achieves state-of-the-art performance in several configurations



Thank You Question?

