

ANDROMAK:

State of the art Android ML based detection systems at your fingertips

Matthieu Simonin (SED Rennes)
David Bromberg (WIDE Team)

Context 1/4

Interested in the experimentation practices in various contexts

- Distributed systems¹
 - Organizer / participant of meetings around experimentations
 - ▶ XUG meetings: <https://xug.gitlabpages.inria.fr/meetings/>
(Check out the upcoming sessions)
 - ▶ AskTheSED: coming soon
-

- Recently interested in Android malware classifications methodology
 - ▶ Static analysis with scale requirements
 - ▶ Environment local machine, Grid'5000

¹Cherrueau et al. EnosLib: A Library for Experiment-Driven Research in Distributed Computing.
10.1109/TPDS.2021.3111159

$$APK : x \xrightarrow{C} C(x) \in \{0, 1\}$$

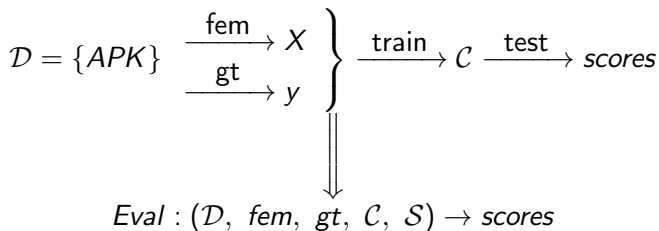
- Some contributions focus on building C
 - ▶ *"Here is a new C and I'll show you how good it is"*
- Some contributions focus on evading an existing C
 - ▶ Black/white box adversarial / poisoning
 - ▶ *"This dataset evades this C , here's how we made it."*
- Some (less) contributions focus on the overall methodology
 - ▶ Reproducibility
 - ▶ Bias identification and evaluation
 - ▶ Brings nuances: *"yes your C is better in this setting BUT..."*

The three kinds of contributions are complementary²

²BlackBurn et al.: The Truth, The Whole Truth, and Nothing But the Truth: A Pragmatic Guide to Assessing Empirical Evaluations <https://doi.org/10.1145/2983574>

Context 3/4

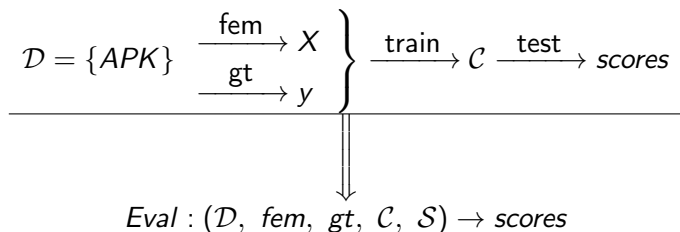
In a supervised context:



Some baseline:

- Drebin paper:
 - ▶ $\text{Eval}(120K + 5K(\text{drebin}), \text{fem}_{\text{drebin}}, \text{VirusTotal}^*, \text{SVC}, 66/33)$
 - ▶ comparison with other related fem on the same dataset
- Mamadroid paper:
 - ▶ \mathcal{D} : product of GW (oldbenign/new benign) + MW(drebin/ V.Share)
 - ▶ \mathcal{C} : Random Forest / k-NN
 - ▶ S : 10-fold CV + mean f1 scores

Context 4/4



You need **Variations**³ into the pipeline == parameter change

- \mathcal{D} : dataset
- fem: feature extraction method (core of many contributions)
- gt: ground truth (there isn't one true ground-truth)
- \mathcal{C} : model / hyperparameters (including DNN models...)
- selection procedure (5/10 fold-CV, time-aware split ...)
- and the associated statistical analysis (hypothesis testing ...)

³Feitelson, D.: From repeatability to reproducibility and corroboration. *ACMSIGOPS Oper. Syst. Rev.*49, 3–11 (2015)

Question

Are we there in the Android Malware classification context ?

Not quite⁵

Eval: (\mathcal{D} , fem , gt , \mathcal{C} , \mathcal{S}) \rightarrow scores

hard to vary

- No reference datasets
- Unavailability / deprecated
- Tight coupling between \mathcal{D} and fem/gt

easy to vary

- Nice APIs (sklearn / tensorflow ...)

⁵Daoudi et al. Lessons Learnt on Reproducibility in Machine Learning Based Android Malware Detection. Empir Software Eng 26, 74 (2021). <https://doi.org/10.1007/s10664-021-09955-7>

Walkthrough Andromak

A library to play with: " $Eval : (\mathcal{D}, fem, gt, \mathcal{C}, \mathcal{S}) \rightarrow scores$ " in the Android context

- Goal 1: Allow experimenters to build easily various evaluations (write their own *Eval* functions)
- Goal 2: Knowledge base around classical techniques around Android malware classification

Design principles:

- Provide an implementation compatible for minimalistic infrastructure (single machine or cluster of machines with a minimal shared storage)
- Favor(abuse of?) fonctionnal aspects of the language

Knowledge base (datasets)

The library has built-in support for

Dataset	#sample	MW%	Ref
CICMalAnal17	4K	50%	10.1016/j.cose.2011.12.012
GM19	10K	50%	hal-02288116
Maldroid20	17K	75%*	10.1109/DASC-PICom-...
Malscan19	30K	50%	10.1109/ASE.2019.00023
Tess19	120K	10%	usenix/Pendlebury19

How ? Store abstraction

- Generic API for getting/storing APK, metadata, any transformation of APK/metadata
- Can be populated from an existing dataset (on disk or dl from androzo)
- Can be backed in a FS (or a DB – theoretically)

Knowledge base (fem)

Various *fem* (decoupled from the datasets):

fem	Réf
Drebin	ndss/ArpSHGR14
Mamadroid	10.1145/3313391
Malscan	10.1109/ASE.2019.00023
Revealdroid	10.1145/3162625
Hindroid	0.1145/3097983.3098026

How ? Rely on the definition of

- *extract* : $APK \rightarrow Feature$
 - ▶ can be full python or wrap a docker execution
 - ▶ **serializable \Rightarrow implicit parallelism/distribution on a cluster**
- *dump* : $Feature \rightarrow bytes$
 - ▶ compression / sparsification
 - ▶ storage optimization

Example: massively parallel fem execution

extract serializable \Rightarrow implicit parallelism/distribution on a cluster

- screenshot of *fem_{drebin}* execution on 1M apk
- Andromak uses Dask⁶ transparently (here on Grid5000 using 50 nodes / 250 workers)



⁶<https://docs.dask.org/en/latest/>

Example: dataset / fem decoupling

FEM		MalAnal17	GM19	MalDroid20	Malscan17	Tess19
reference	apps	2,126	10,000	17,187	30,715	129,728
	mw	426	5,000	13,150	15,430	12,737
	gw	1,700	5,000	4,037	15,285	116,991
common	apps	1,940	9,158	15,781	28,856	124,979
	mw	407	4,323	12,172	14,108	12,027
	gw	1,533	4,835	3,609	14,748	112,952
drebin	apps	2,116	10,000	17,186	30,715	129,728
	mw	426	5,000	13,149	15,430	12,737
	gw	1,690	5,000	4,037	15,285	116,991
hindroid	apps	2,111	9,930	16,661	30,690	129,658
	mw	425	4,941	12,630	15,430	12,706
	gw	1,686	4,989	4,031	15,260	116,952
malscan	apps	2,121	9,947	16,702	30,715	129,344
	mw	425	4,947	12,671	15,430	12,679
	gw	1,696	5,000	4,031	15,285	116,665
mamadrogard	apps	2,121	9,945	16,661	30,715	129,677
	mw	425	4,945	12,630	15,430	12,706
	gw	1,696	5,000	4,031	15,285	116,971
mamadroid	apps	1,978	9,288	15,881	29,333	125,660
	mw	408	4,430	12,211	14,527	12,145
	gw	1,570	4,858	3,670	14,806	113,515
revealdroid	apps	2,100	9,764	16,325	30,142	129,258
	mw	420	4,781	12,355	14,897	12,633
	gw	1,680	4,983	3,970	15,245	116,625

Table 3. Extraction effectiveness. number of applications (resp. malware, resp. goodware) that are initially in the various datasets (reference), or successfully extracted for each extraction method. Applications that are successfully extracted by all the extraction techniques are accounted in the common bloc.

Knowledge base (Domain specific ML)

- Some default \mathcal{C} (for reproducibility purpose)
 - ▶ Related classifiers (best ones according to the related papers)
 - ★ e.g. $\mathcal{C}_{drebin} = SVC$
 - ▶ Multilayer Perceptron ⁷
- Specific selection method / Hypothesis testing
 - ▶ Time aware splits (usenix/Pendlebury19)
 - ▶ Wilcoxon, Ranks ...

⁷Grosse et al. (2017) Adversarial Examples for Malware Detection. ESORICS 2017.

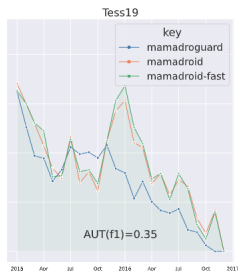
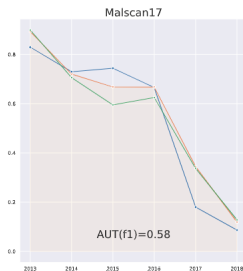
Example: comparison of 3 variants

Eval: (\mathcal{D} , *fem*, *gt*, \mathcal{C} , \mathcal{S}) \rightarrow scores with 3 different *fem* and 5 different \mathcal{C}

- \mathcal{S} : 5x2-fold CV (mean) f1-scores

	MalAnal17		GM19		MalDroid20		Malscan17		Tess19	
	coupled	free	coupled	free	coupled	free	coupled	free	coupled	free
mamadroguard	0.82	0.82	0.92	0.92	0.97	0.97	0.95	0.95	0.74	0.75
mamadroid	0.82	0.81	0.92	0.92	0.97	0.97	0.96	0.96	0.80	0.80
mamadroid-fast	0.82	0.80	0.93	0.93	0.96	0.97	0.96	0.96	0.81	0.81

- \mathcal{S} : time-aware model selection with AUT(f1) score



Conclusion

- Andromak
 - ▶ Let you introduce in your evaluation pipeline
 - ★ variations
 - ★ scale
 - ▶ Future:
 - ★ research report / paper + release
 - ★ sandboxed analysis, new datasets, visualisations ...
- Talk to SED people about your experimentations / problems
 - ▶ XUG / AskTheSED events