























<code>javax/crypto/spec/desKeySpec;.&lt;init&gt;:([B)V</code>	0.057388
<code>android/app/Application;.&lt;init&gt;:()V</code>	0.056203
<code>java/lang/reflect/Field;.getType():Ljava/lang/class;</code>	0.037152
<code>java/lang/class;.forName:(Ljava/lang/string;)Ljava/lang/class;</code>	0.027028
<code>org/json/JSONArray;.toString():Ljava/lang/string;</code>	0.014487
<code>javax/crypto/Cipher;.getInstance:(Ljava/lang/string;)javax/crypto/Cipher;</code>	0.011853
<code>java/io/ByteArrayOutputStream;.flush():V</code>	0.008844
<code>java/io/RandomAccessFile;.writeLong:(J)V</code>	0.007683
<code>android/widget/Toast;.makeText:(Landroid/content/Context;Ljava/lang/CharSequence;)I</code>	0.005666
<code>java/lang/reflect/Array;.getInt:(Ljava/lang/object;I)I</code>	0.003103

Figure 8. API-calls and attention weights for an Android malware app

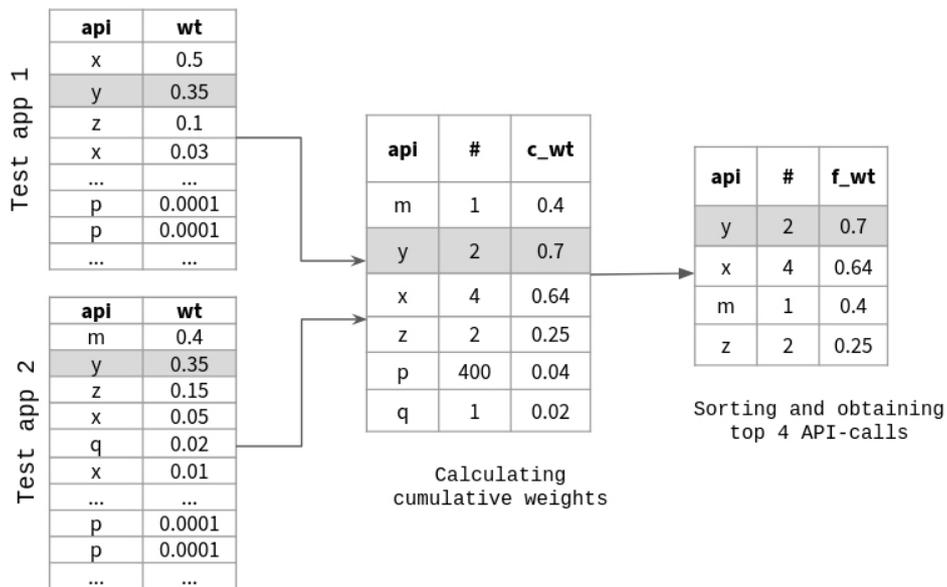


Figure 9. Top feature selection strategy: For the sake of readability, here we illustrate the selection process with a tiny set of two apps whereas the goal is to identify only top 4 features.

In the current work, we finally extracted 200 top features of malicious apps. In order to prove that these top features are valid, we plan to compare them with the manually extracted features in a future work.

Someone may wonder that instead of considering so many individual API calls as in the current work, whether it is better to merge API calls into groups such as permission related, string related, File/Network IO related, and so on. Although the features will then become more coarse-grained, it is important to explore this route in the future.

## 8. Conclusions

We leveraged attention-based deep learning approaches for security vetting of Android apps. API-calls extracted

from the Android apps were used as the artifacts for deep-learning models. Such API-calls are quasi-sequential in nature. We experimented with two attention-based models: Bi-LSTM Attention and Self-Attention. Both models gave competitive results. Additionally, after analyzing the attention weights from these two models, we identified top 200 API-calls that reflect the maliciousness of an Android app.

Our experiments show that deep learning models can be implemented for large scale Android app security vetting. This can save human time and effort from manually handpicking the malicious features. The top features identified from the models can be further studied by the research community to potentially discover new malware signature.

**Table 5.** Common API-calls from the two models

Common API-calls
Ljava/io/file;mkdir:()z
Ljava/io/file;exists:()z
Ljava/io/fileinputstream;close:()v
Ljava/lang/reflect/method;setaccessible:(z)v
landroid/content/res/resources;getassets:()landroid/content/res/assetmanager;
Ljava/io/file;delete:()z
Ljava/io/filenotfoundexception;printstacktrace:()v
landroid/telephony/telephonymanager;getdeviceid:()Ljava/lang/string;
landroid/os/asynctask;onPostExecute:(Ljava/lang/object;)v
Ljava/io/file;listfiles:()Ljava/io/file;
landroid/content/res/resources;getboolean:(i)z
landroid/util/base64;decode:([bi]b
landroid/app/application;onConfigurationChanged:(landroid/content/res/configuration;)v
Ljava/lang/stringbuilder;insert:(iljava/lang/string;)Ljava/lang/stringbuilder;
lorg/apache/http/statusline;getStatusCode:()i
landroid/content/sharedpreferences;edit:()landroid/content/sharedpreferences\$editor;
Ljava/io/objectoutputstream;writeObject:(Ljava/lang/object;)v
landroid/app/application;onTerminate:()v
Ljava/io/inputstream;mark:(i)v
Ljava/lang/process;destroy:()v
Ljava/util/zip/zipinputstream;closeEntry:()v

**Table 6.** Accuracy of standard ML algorithms

ML Model	App Class	Precision	Recall	F1-Score	Area Under PR Curve
Bernoulli Naive Bayes	Malicious	0.6209	0.6105	0.6156	0.6845
	Benign	0.8709	0.8757	0.8733	
K-Nearest Neighbor (K=5)	Malicious	0.9071	0.8881	0.8975	0.9471
	Benign	0.9630	0.9697	0.9663	
Support Vector Machine	Malicious	0.9133	0.8940	0.9035	0.9547
	Benign	0.9649	0.9717	0.9683	

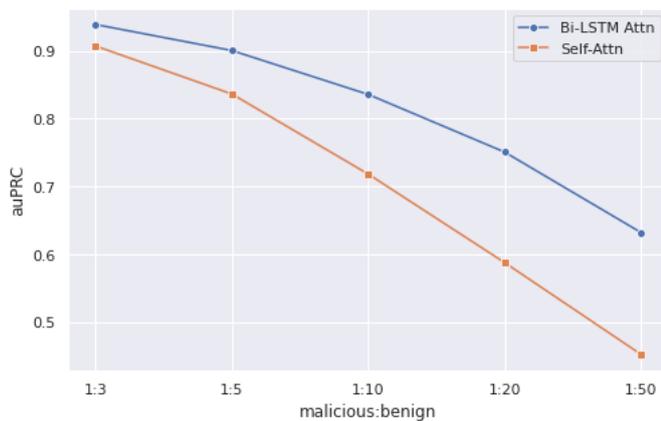
**Table 7.** Deep learning models' accuracy with high quality malware

Model	Review Type	Precision	Recall	F1-Score	Area Under PR Curve
Bi-LSTM Attn.	Benign	0.9554	0.9772	0.9662	0.9548
	Malicious	0.9266	0.8632	0.8938	
Self Attention	Benign	0.9283	0.9586	0.9432	0.9291
	Malicious	0.9084	0.7949	0.8479	

**Acknowledgement.** This work has been partially supported by the U.S. National Science Foundation (NSF) under grant no. 1718214 and 1717871. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF.

## References

- [1] STATCOUNTER (2021), Android OS Market Share. URL <https://gs.statcounter.com/os-market-share/mobile/worldwide>.
- [2] CHAULAGAIN, D., POUDEL, P., PATHAK, P., ROY, S., CARAGEA, D., OU, X. and LIU, G. (2020) Hybrid analysis of android apps for security vetting using deep learning. In *IEEE conference on communications and network security (CNS)*.
- [3] PASCANU, R., MIKOLOV, T. and BENGIO, Y. (2012) Understanding the Exploding Gradient Problem. *arXiv e-print* : arXiv:1211.5063.
- [4] BENGIO, Y., SIMARD, P. and FRASCONI, P. (1994) Learning Long-term Dependencies with Gradient Descent is Difficult. *IEEE Transactions on Neural Networks* : 157–166.



**Figure 10.** Experimenting with varying malware apps to benign apps ratio in the test set

- [5] GERS, F.A., SCHMIDHUBER, J. and CUMMINS, F. (1999) Learning to forget: continual prediction with lstm. In *1999 Ninth International Conference on Artificial Neural Networks ICANN 99. (Conf. Publ. No. 470)*: 850–855.
- [6] BAHDANAU, D., CHO, K. and BENGIO, Y. (2016), Neural machine translation by jointly learning to align and translate. [1409.0473](#).
- [7] LUONG, M.T., PHAM, H. and MANNING, C.D. (2015), Effective approaches to attention-based neural machine translation. [1508.04025](#).
- [8] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A.N., KAISER, L. et al. (2017), Attention is all you need. [1706.03762](#).
- [9] ARZT, S., RASTHOFER, S., FRITZ, C., BODDEN, E., BARTEL, A., KLEIN, J., LE TRAON, Y. et al. (2014) FlowDroid: Precise Context, Flow, Field, Object-sensitive and Lifecycle-aware Taint Analysis for Android Apps. *SIGPLAN* : 259–269.
- [10] WEI, F., ROY, S., OU, X. and , R. (2018) Amandroid: A Precise and General Inter-component Data Flow Analysis Framework for Security Vetting of Android Apps. *ACM Transactions on Privacy and Security* : 1–32.
- [11] ARP, D., SPREITZENBARTH, M., HÜBNER, M., GASCON, H. and RIECK, K. (2014) DREBIN: Effective and Explainable Detection of Android Malware in Your Pocket. In *Symposium on Network and Distributed System Security (NDSS)*: 23–26.
- [12] ROY, S., DELOACH, J., LI, Y., HERNDON, N., CARAGEA, D., OU, X., RANGANATH, V.P. et al. (2015) Experimental study with real-world data for android app security analysis using machine learning. In *Proceedings of the 31st Annual Computer Security Applications Conference (ACSAC)*: 81–90.
- [13] AAFER, Y., DU, W. and YIN, H. (2013) DroidAPIMiner: Mining API-Level Features for Robust Malware Detection in Android. In *International conference on security and privacy in communication systems*: 86–103.
- [14] ONWUZURIKE, L., MARICONTI, E., ANDRIOTIS, P., DE CRISTOFARO, E., ROSS, G. and STRINGHINI, G. (2017) MaMaDroid: Detecting Android Malware by Building Markov Chains of Behavioral Models. *arXiv e-prints* : arXiv:1612.04433.
- [15] GONG, L., LI, Z., QIAN, F., ZHANG, Z., CHEN, Q., QIAN, Z., LIN, H. et al. (2020) Experiences of landing machine learning onto market-scale mobile malware detection. *Proceedings of the Fifteenth European Conference on Computer Systems* .
- [16] KE, X., LI, Y., DENG, R.H. and CHEN, K. (2018) DeepRefiner: Multi-layer Android Malware Detection System Applying Deep Neural Networks. In *2018 IEEE European Symposium on Security and Privacy (EuroS P)*: 473–487.
- [17] HOU, S., SAAS, A., CHEN, L. and YE, Y. (2016) Deep4MalDroid: A Deep Learning Framework for Android Malware Detection Based on Linux Kernel System Call Graphs. In *2016 IEEE/WIC/ACM International Conference on Web Intelligence Workshops (WIW)*: 104–111.
- [18] KARBAB, E.B., DEBBABI, M., DERHAB, A. and MOUHEB, D. (2018) MalDozer: Automatic Framework for Android Malware Detection using Deep Learning. *Digital Investigation* : S48–S59.
- [19] (2018), Genymotion Android Emulator. URL <https://www.genymotion.com/>.
- [20] SUTSKEVER, I., VINYALS, O. and LE, Q.V. (2014), Sequence to sequence learning with neural networks. [1409.3215](#).
- [21] CHO, K., VAN MERRIENBOER, B., GULCEHRE, C., BAHDANAU, D., BOUGARES, F., SCHWENK, H. and BENGIO, Y. (2014), Learning phrase representations using rnn encoder-decoder for statistical machine translation. [1406.1078](#).
- [22] WU, B., CHEN, S., GAO, C., FAN, L., LIU, Y., WEN, W. and LYU, M.R. (2021) Why an android app is classified as malware: Toward malware classification interpretation. *ACM Trans. Softw. Eng. Methodol.* **30**(2).
- [23] ZHOU, P., SHI, W., TIAN, J., QI, Z., LI, B., HAO, H. and XU, B. (2016) Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*: 207–212.
- [24] HOCHREITER, S. and SCHMIDHUBER, J. (1997) Long Short-Term Memory. *Neural Computation* : 1735–1780.
- [25] SCHUSTER, M. and PALIWAL, K. (1997) Bidirectional Recurrent Neural Networks. *IEEE Transactions on Signal Processing* : 2673–2681.
- [26] CHO, K., VAN MERRIENBOER, B., BAHDANAU, D. and BENGIO, Y. (2014), On the properties of neural machine translation: Encoder-decoder approaches. [1409.1259](#).
- [27] ALLIX, K., BISSYANDÉ, T.F., KLEIN, J. and LE TRAON, Y. (2016) AndroZoo: Collecting Millions of Android Apps for the Research Community. In *Proceedings of the 13th International Conference on Mining Software Repositories*: 468–471.
- [28] PENDLEBURY, F., PIERAZZI, E., JORDANEY, R., KINDER, J. and CAVALLARO, L. (2019) TESSERACT: Eliminating experimental bias in malware classification across space and time. In *28th USENIX Security Symposium*: 729–746.