

## Abstract

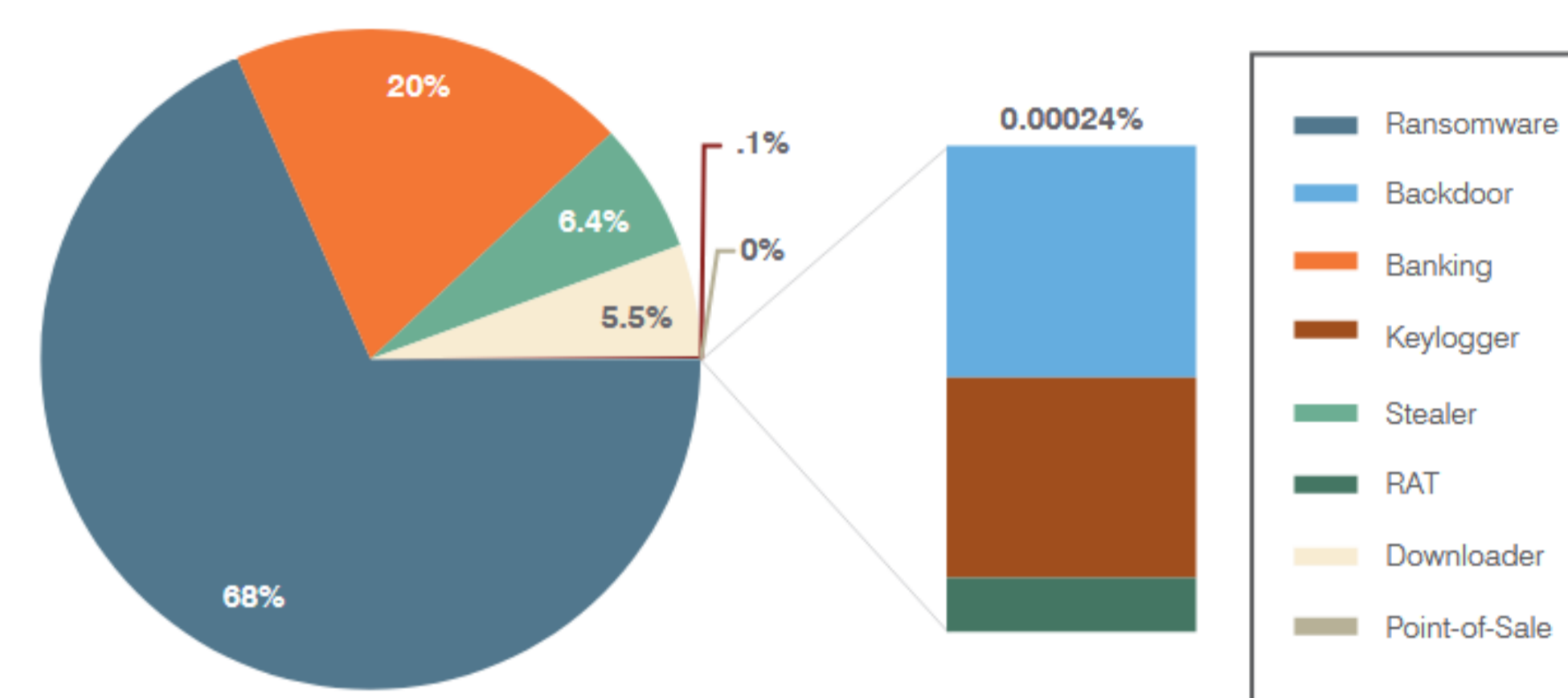
Owning the largest share of Smartphone platform tendency, Android OS has changed to an undeniable target for many malware authors. According to real world scenario, most of Android malware applications don't benefit from having labels and applying supervised learning wouldn't be the right solution to pick. That's because of the major drawback of supervised learning, i. e., requiring lots of labeled data that is so expensive to collect. Addressing this issue, in this work, we employ a semi-supervised deep learning architecture, called, ladder networks that minimizes the sum of supervised and unsupervised cost functions by backpropagation at the same time. We craft two different types of feature vectors, namely, API Call frequency and API Call sequences to feed into ladder networks. The experimental results show a reasonable gap of 7-8% in total accuracy for 100 used labels and all used labels. At last, we fine-tune the learning parameters of ladder networks to improve its generalization ability and present corresponding graphs and diagrams

## Previous Research

Year	Authors	Publication Title	Focus Area	Analysis Technique
2013	Aafer et. Al.	Droidapiminer: Mining api-level features for robust malware detection in android [SecureComm]	Android Malware Detection	Lightweight approach based on semantic information inside bytecode of the applications such as, invoked critical API calls, package level information as well as some important parameters
2014	Zhang et. al.	Semantics-aware android malware classification using weighted contextual api dependency graphs [ACM SIGSAC]	Android Malware Detection	Semantic-based approach that classifies Android malware via a weighted contextual API dependency graph
2017	Mariconti et. al.	Mamadroid: Detecting android malware by building markov chains of behavioral models [NDSS]	Android Malware Detection	Static Android malware detection system based on the sequences of abstracted API Calls which are used to build a Markov chain
2016	Hou et. al.	Deep4maldroid: A deep learning framework for android malware detection based on linux kernel system call graphs [WIW]	Android Malware Detection using Deep Learning	Dynamic Android malware detection based on a weighted graph of Linux kernel system calls using Stacked Auto-Encoders

## Android Malware Categories

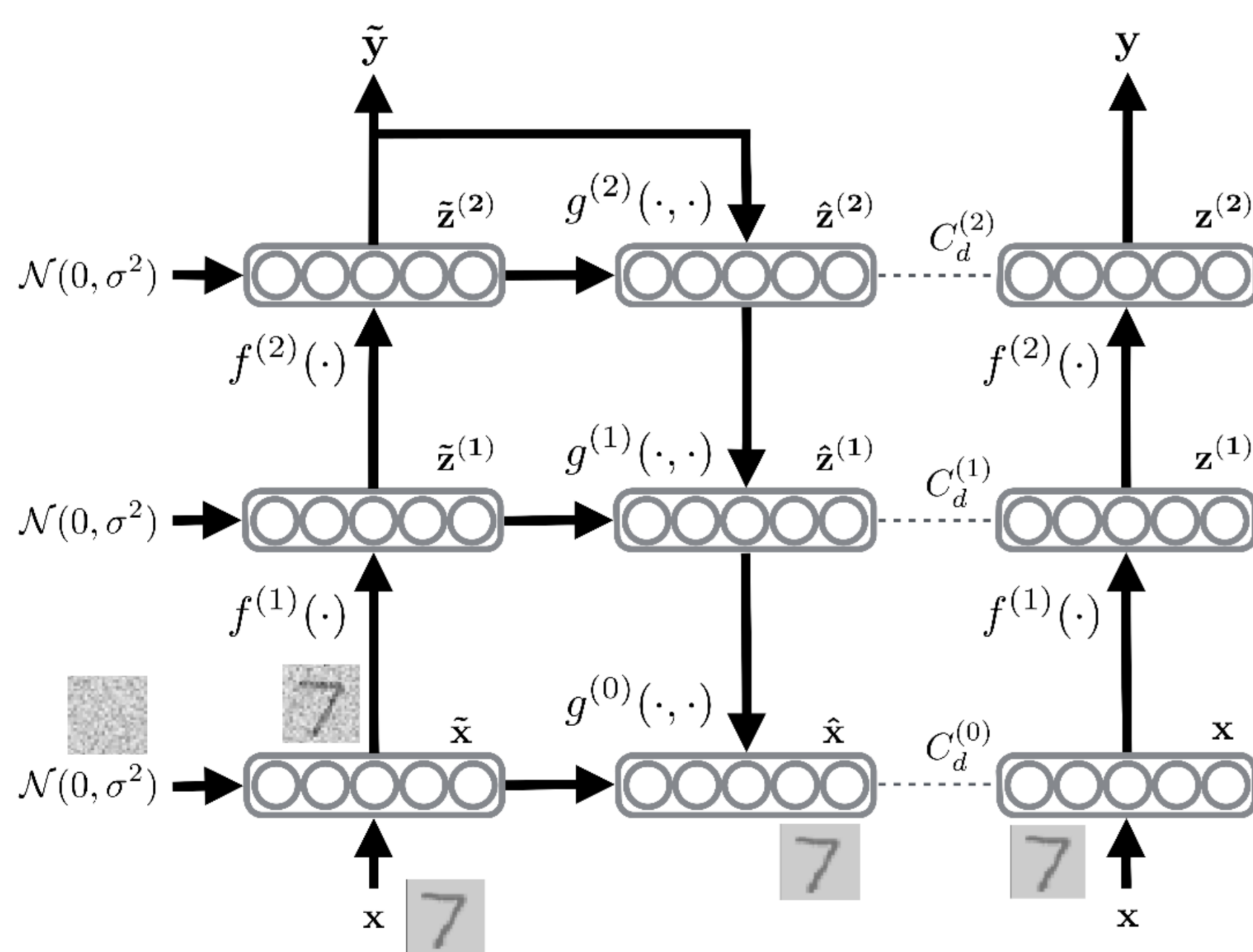
Malware by Category, Q2 2017



Comparison of quarterly overall message volume

## Ladder Network

- Auto-encoder with lateral shortcut connections from the encoder to the decoder at each level of the hierarchy
- Allows the higher levels of the network to discard details and focus on representing more abstract invariant features



## Experiments

```

1- {}
2- "0": {
3   "inputs": [],
4   "package": "com.android.music",
5   "outputs": {
6     "ReturnType": "java.lang.Boolean",
7     "ReturnValue": "true"
8   },
9   "class": "java.io.File",
10  "timestamp": "1477336771863",
11  "method": "exists"
12 },
13- "1": {
14  "inputs": [],
15  "package": "com.android.music",
16  "outputs": {
17    "ReturnType": "java.lang.Boolean",
18    "ReturnValue": "true"
19  },
20  "class": "java.io.File",
21  "timestamp": "1477336771863",
22  "method": "exists"
23 },
24- "2": {
25  "inputs": [],
26  "package": "com.android.music",
27  "outputs": {
28    "ReturnType": "java.lang.Boolean",
29    "ReturnValue": "true"
30  },
31  "class": "java.io.File",
32  "timestamp": "1477336831841",
33  "method": "exists"
34 },

```

TABLE I  
PARAMETER FINE-TUNING

Parameter Values	
layers	[4761, 1000, 500, 250, 250, 250, 5]
learning rate	0.02
Decay-after	5
Batch-size	20
Epoch	20

Algorithm 1 Calculation of the output y and cost function C of the Ladder network

**Require:**  $x(n)$

# Corrupted encoder and classifier  
 $\tilde{h}^{(0)} \leftarrow \tilde{z}^{(0)} \leftarrow x(n) + \text{noise}$   
 # Decoder and denoising  
 for  $l = 1$  to  $L$  do  
    $\tilde{z}^{(l)} \leftarrow \text{batchnorm}(W^{(l)}\tilde{h}^{(l-1)} + \text{noise})$   
    $\tilde{h}^{(l)} \leftarrow \text{activation}(\gamma^{(l)} \odot (\tilde{z}^{(l)} + \beta^{(l)}))$   
 end for  
 $P(\tilde{y} | x) \leftarrow \tilde{h}^{(L)}$   
 # Clean encoder (for denoising targets)  
 $h^{(0)} \leftarrow z^{(0)} \leftarrow x(n)$   
 for  $l = 1$  to  $L$  do  
    $z_{\text{pre}}^{(l)} \leftarrow W^{(l)}h^{(l-1)}$   
    $\mu^{(l)} \leftarrow \text{batchmean}(z_{\text{pre}}^{(l)})$   
    $\sigma^{(l)} \leftarrow \text{batchstd}(z_{\text{pre}}^{(l)})$   
    $z^{(l)} \leftarrow \text{batchnorm}(z_{\text{pre}}^{(l)})$   
    $h^{(l)} \leftarrow \text{activation}(\gamma^{(l)} \odot (z^{(l)} + \beta^{(l)}))$   
 end for

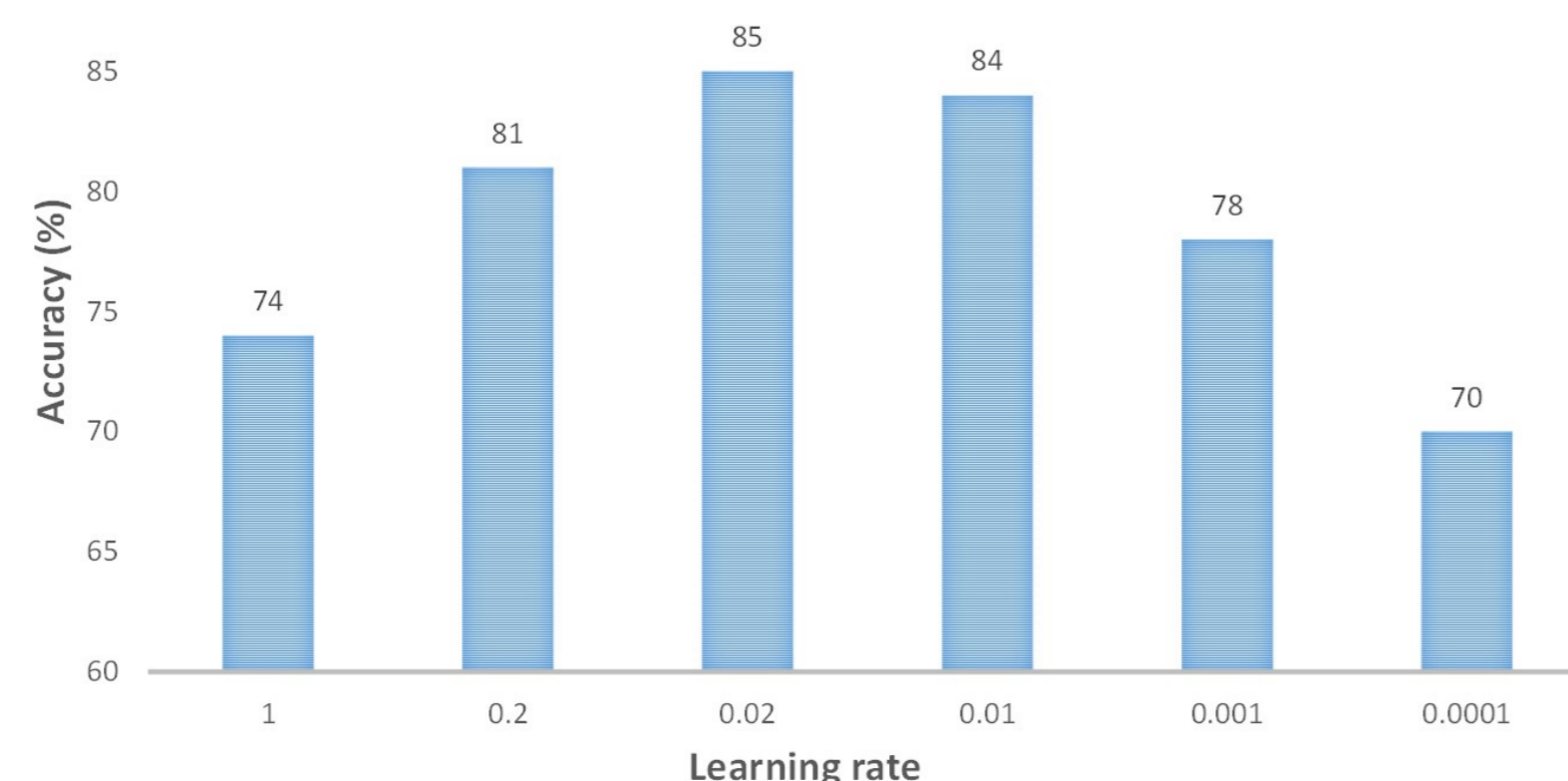
# Final classification:  
 $P(y | x) \leftarrow h^{(L)}$   
 # Decoder and denoising  
 for  $l = L$  to  $0$  do  
   if  $l = L$  then  
    $u^{(L)} \leftarrow \text{batchnorm}(\tilde{h}^{(L)})$   
   else  
    $u^{(l)} \leftarrow \text{batchnorm}(V^{(l+1)}\tilde{z}^{(l+1)})$   
   end if  
    $\forall i: \tilde{z}_i^{(l)} \leftarrow g(\tilde{z}_i^{(l)}, u_i^{(l)})$  # Eq. (2)  
    $\forall i: \tilde{z}_{i, \text{BN}}^{(l)} \leftarrow \frac{\tilde{z}_i^{(l)} - \mu_i^{(l)}}{\sigma_i^{(l)}}$   
 end for  
 # Cost function C for training:  
 $C \leftarrow 0$   
 if  $t(n)$  then  
    $C \leftarrow -\log P(\tilde{y} = t(n) | x(n))$   
 end if  
 $C \leftarrow C + \sum_{l=0}^L \lambda_l \|z^{(l)} - \tilde{z}_{\text{BN}}^{(l)}\|^2$  # Eq. (3)

## Dataset and Results

TABLE II  
ACCURACY BASED ON NUMBER OF USED LABELS

Accuracy % with # of used labels	100	All
API Call Frequency	59%	66%
API Call Sequence	77%	85%

Type	Number
Banking	250
Ransomware	250
Adware	250
Botnet	250
Benign	250



## Future Work

### Practical

- We have to add more data to our dataset since doing so will improve the results and considering that we can get our hands on more data, this is definitely going to make a noticeable improvement in accuracy
- In the near future, there is going to be more works on semi-supervised android malware detection and we need to compare our method to the other algorithms
- There are more datasets that we can trust and as a result, perform our model on them. This makes our method more comparable to others and also may open some new aspects to the problem.

### Theoretical

- The analysis of feature is rather incomplete, as we can extract more and better features from API call sequences and frequencies. Much remains to be done in this regard