

Performance Evaluation of a CNN for Malware Detection (Maling Dataset)

Problem Statement: The paper addresses automated **malware detection and classification** by treating malware binaries as images. It converts executables (PE files) into grayscale images and trains a CNN to recognize malware families ¹. The motivation is that traditional signature-based methods are failing, so deep CNNs can “accurately and effectively detect malware” ². In summary, the problem is a multi-class classification of malware vs. benign samples using image-based representations.

Dataset: The experiments use the **Maling dataset**, a standard malware image corpus. Maling contains **9,339 malware samples** across **25 families** ³ ⁴. The paper notes that all images are grayscale and initially sized 32×32, so the data are loaded and reshaped accordingly ⁵. (In code the authors load Maling via `maling.load_data()` and reshape to shape (32,32,1) ⁶, indicating 32×32 gray images.) The data are split (via `maling.load_data()`) into training and test sets.

CNN Architecture: A simple CNN is defined in Keras. The **baseline model** has a feature-extraction “front end” of convolution and pooling layers, followed by a classification “back end” of dense layers ⁷. For example, the baseline uses one Conv2D layer (32 filters, 3×3, ReLU), a 2×2 max-pool, flattening, then two Dense layers (100-unit ReLU and 10-unit softmax). After establishing this baseline, an **improved model** is created by “increasing the depth” in a VGG-like fashion: adding a second Conv layer (also 64 filters) before pooling ⁸. Thus the improved CNN has two 64-filter conv layers (with max-pooling) instead of one. (The paper describes this in text; code for the improved model is not shown but is said to follow the double-conv pattern ⁸.)

Evaluation Metrics & Protocol: The model is evaluated primarily using **classification accuracy**. The authors compile the model with categorical cross-entropy loss and monitor accuracy ⁹. They perform **5-fold cross-validation**: “evaluate the model using five-fold cross-validation (k=5)” ¹⁰. They plot train/test loss and accuracy curves to check fit (reporting no obvious overfitting) ¹¹. The final reported metric is the mean accuracy (with standard deviation) over folds. No other metrics (precision, recall, etc.) are reported, so accuracy is the main performance measure.

Results: The **baseline CNN** achieves about **98% accuracy**, and the **improved (deeper) CNN** achieves **~99.18% accuracy** ² ¹². In the cross-validation folds shown, the improved model reached 99.183% in one fold ². The paper reports a final mean accuracy $\approx 99.18\%$ ($\pm 0.5\%$) ¹², indicating very high performance. Table 3 in the paper compares this to prior work: e.g. Nataraj et al. (2011) 98%, VASAN et al. 98.82%, Moussa et al. 99.13% etc., showing the proposed model is at or above state-of-the-art ² ¹². The authors also mention using the final model to predict on *new* malware samples (hold-out test), though no separate test-set results are explicitly given beyond cross-val accuracy ².

Critical Analysis

Methodology: Converting malware binaries to images (by interpreting byte values as pixels) is a known approach in malware classification ¹. The authors follow this paradigm but do not contribute a novel representation – the novelty claimed is in model “optimization”. Methodologically, the paper constructs a baseline CNN (simple conv+pool) and then deepens it. The idea of increasing depth (following a VGG-like pattern) is reasonable for capacity, and the authors correctly use cross-entropy loss and accuracy metric.

However, the **method description has gaps**. It is unclear how class imbalance is handled: Maling’s 25 families have unequal sample counts, but the paper reports only overall accuracy. They use 10 output units in the code despite 25 classes (the last Dense is size 10), suggesting they may have restricted the task to 10 classes, but this is never explained. This discrepancy (25 families vs. 10-class output) raises questions about what exactly was done. Also, only accuracy is used; for an imbalanced multi-class problem, precision/recall per class or confusion matrices would be informative but are absent. The methodology section is also interwoven with tutorial-like code (listing data loading and model definition), which is clear but not very rigorous.

Experimental Setup: The paper uses 5-fold cross-validation on the Maling dataset ¹⁰, which is good for robust estimation on a small dataset. Learning curves in Fig.7 (not shown here) indicate training/test accuracy curves match, suggesting the model is not overfitting ¹¹. The optimizer (SGD lr=0.01, momentum 0.9) and loss function are specified. However, details like number of epochs, batch size, or random seed are not given. The baseline vs. improved model experiment is straightforward: the deeper model does slightly better (from ~98% to ~99.18%). There is no validation on truly novel data beyond CV, except a mention of “predictions on new samples” without results.

Validity of Conclusions: The main claim is that the proposed CNN “outperforms most CNN models in the literature” ². The reported accuracy (~99.18%) is indeed on par with or slightly above prior published results (as seen in Table 3 ² ¹²). So the conclusion of high performance is supported by the numbers. However, the **scope is limited**: results hold on the Maling dataset only. That dataset is static and relatively small (9k images), so a 99% accuracy is less surprising. In practice, new malware families and larger datasets exist, so it’s unclear how well this model generalizes. The paper does not test on other datasets or real-world streams. Also, since only accuracy is reported, it’s not known if certain classes dominate the performance. In summary, the conclusions are valid **for this dataset**, but the paper’s framing (“secure model outperforms others”) should be tempered by dataset limitations.

Key Contributions and Takeaways

- **New CNN for Malware:** The paper proposes a CNN “from scratch” for malware image classification, achieving high accuracy ¹³.
- **Data Science Approach:** It highlights an “expert data science” process: establish a baseline CNN, then systematically improve it (e.g. by adding layers) ¹³ ⁸.
- **Performance Evaluation:** The study rigorously evaluates using k-fold cross-validation, obtaining ~98% (baseline) to ~99.18% (improved) accuracy ² ¹³.
- **Comparison with Literature:** It compares its results to past work, showing its accuracy is competitive or better (Table 3) ² ¹².

- **Practical ML Workflow:** As a takeaway, the paper emphasizes the importance of methodical model tuning (changing optimizers, network depth) to boost performance ⁸.

Reproducibility and Implementation

- **Data Availability:** The Malimg dataset is publicly available (e.g. via the original authors or Kaggle) ³ ⁴. Reproducing experiments would require downloading Malimg, which contains 9,339 grayscale images of malware from 25 families.
- **Model Clarity:** The authors provide Keras code snippets for loading data and defining the CNN ⁵ ⁷. The baseline architecture (one Conv2D+pool + two Dense layers) is clearly shown. The improvements (double conv with 64 filters) are described in words ⁸. Thus the CNN design and training procedure (SGD optimizer, lr=0.01, momentum=0.9, cross-entropy loss) are given.
- **Implementation Feasibility:** Implementing the model in PyTorch or TensorFlow should be straightforward given these details. The code is standard (use of `model = Sequential()`, conv/pooling, etc.). The evaluation (five-fold cross-validation) is explicitly mentioned ¹⁰. One caveat is that the code shown uses `input_shape=(28, 28, 1)` but text says images are 32×32 ⁵; the implementer must resolve this minor inconsistency (perhaps cropping or resizing to 28×28).
- **Missing Details:** Some aspects are unspecified (e.g. exact train/test split method beyond “load_data”, number of epochs, batch size, random seed). The output layer size suggests 10 classes, but Malimg has 25 – so it’s unclear if only a subset of classes was used. These gaps would need to be clarified when reproducing.
- **Overall:** Reproduction is feasible with moderate effort. The necessary components (dataset, network architecture, training loops) are described. The paper does not release code, but the pseudocode is explicit enough that an experienced student could rebuild the model and obtain similar accuracy.

Strengths and Weaknesses

Strengths (for reproducibility/learning): - The paper includes code excerpts and parameter settings, which is excellent for a learning project. Seeing a real example of Keras model definition and cross-validation setup is educational ⁵ ¹⁰. - It systematically steps through baseline vs. improved models, illustrating the process of model iteration. - The use of a known dataset (Malimg) and comparisons to prior work adds context. - The writing is clear on the *intended* approach (though not always technically precise), and the contributions are explicitly listed ¹³.

Weaknesses: - **Limited novelty:** The idea (CNN on malware images) is not new, so the paper reads more like an application report than a research breakthrough. - **Lack of detail/clarity:** Key details (e.g. number of classes, data splitting) are ambiguous. The modeling choices (why 10 output units, why 32×32 vs 28×28) are not explained. - **Evaluation scope:** Only one dataset is used. An academic reproducibility project might criticize the narrow evaluation and suggest testing additional malware corpora. - **No error analysis:** The paper gives raw accuracy but no insight into which malware families are confused or how robust the model is to variations. This limits learning from the results. - **Claims vs. Reality:** Phrases like “expert data science approach” are vague. A critical student would note the contributions are mostly incremental.

Overall, for a student project focused on CNN practice, this paper is **useful**: it shows a full pipeline from data loading to model tuning. But academically, it is **weak**: it doesn’t explore the problem deeply or provide strong validation beyond high accuracy on a fixed dataset. A reproducibility effort would be straightforward but might also reveal the paper’s oversights (e.g. clarify data issues).

Sources: The above analysis is based on the paper's content ¹ ⁵ ¹³ ⁸ ¹⁰ ¹² and relevant external dataset info ⁴ .

¹ ³ ⁵ **arxiv.org**
⁶ ⁷ ⁸ <https://arxiv.org/pdf/2301.11161>
⁹ ¹⁰ ¹¹
¹² ¹³

² **[2301.11161] New Approach to Malware Detection Using Optimized Convolutional Neural Network**
<https://ar5iv.org/pdf/2301.11161>

⁴ **Malimg Dataset | Papers With Code**
<https://paperswithcode.com/dataset/malimg>