

# Tăng Cường Phát Hiện Mã Độc Android Với Bảo Đảm Quyền Riêng Tư Thông Qua Federated Learning

Lê Vũ Tuấn Anh

Trường Đại học Công nghệ Thông tin - ĐHQG TP.HCM

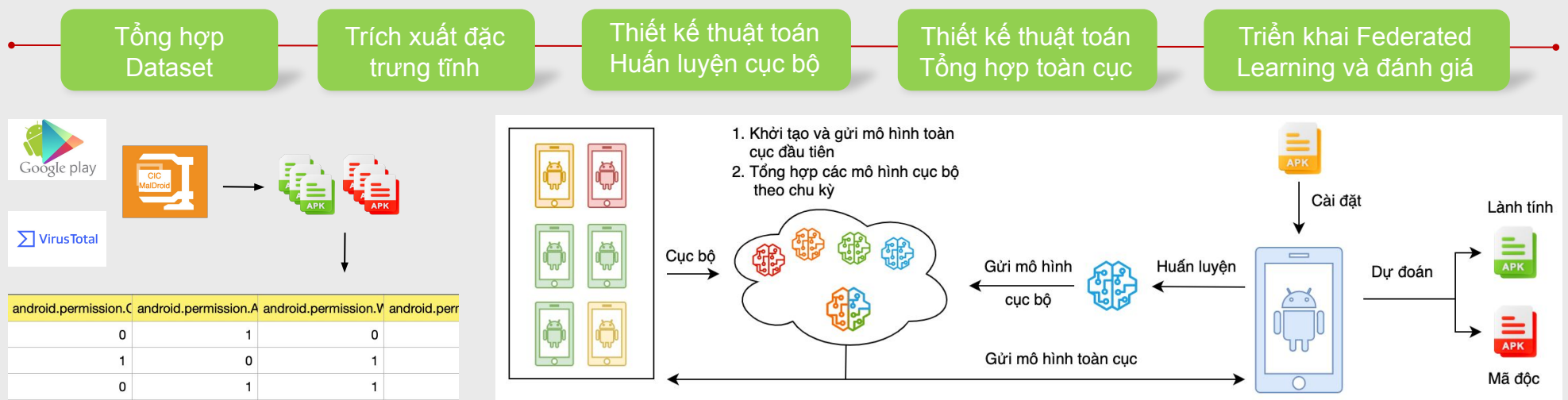
## Mục tiêu ?

- Tổng hợp Dataset các file APK mã độc và lành tính mới nhất và phổ biến từ nhiều nguồn.
- Thiết kế thuật toán huấn luyện cục bộ phát hiện mã độc Android và thuật toán tổng hợp tham số toàn cục cho Federated Learning.
- Đánh giá tính khả thi và hiệu suất hệ thống với nhiều kịch bản khác nhau về số lượng thiết bị, số lượng đặc trưng và số chu kỳ.

## Lý do chọn đề tài ?

- Mô hình phát hiện mã độc Android cần được huấn luyện, cập nhật chủ động với dữ liệu mới từ người dùng.
- Phương pháp truyền thống thu thập thông tin từ thiết bị để huấn luyện tập trung, gây rủi ro rò rỉ dữ liệu cá nhân nhạy cảm.
- Sử dụng Federated Learning là một hướng tiếp cận mới, tiềm năng trong huấn luyện mô hình phát hiện mã độc, đồng thời bảo đảm được tính riêng tư.

## Tổng quan



Hình 1. Tập tin CSV chứa đặc trưng sau khi trích xuất

Hình 2. Mô hình hệ thống đề xuất cho nghiên cứu

## Mô tả

### 1. Tổng hợp Dataset

- Tải xuống các Dataset mã độc Android: CIC, VirusTotal...
- Tải thêm mẫu lành tính từ Google Play Store.
- So MD5 các mẫu để loại bỏ trùng lặp.

### 5. Triển khai Federated Learning và đánh giá

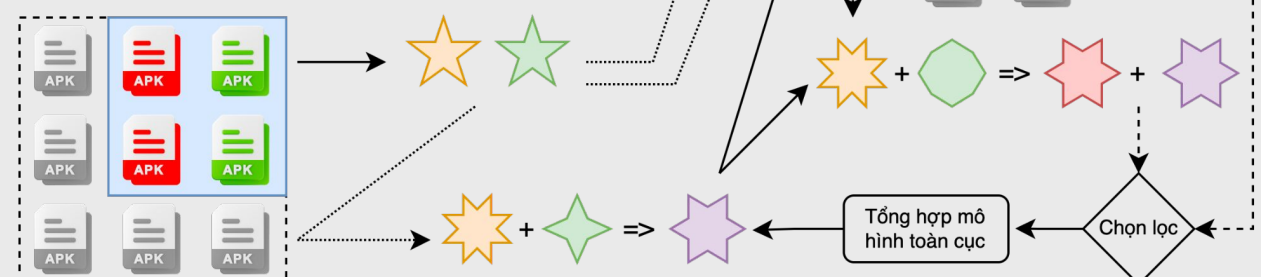
- Chia dữ liệu huấn luyện thành các bộ có nhãn và không nhãn, phân phối cho máy chủ và thiết bị Android.
- Thử nghiệm các thư viện học máy và khung triển khai Federated Learning của TensorFlow trên Android và máy chủ.
- Thử nghiệm sử dụng Docker để tạo môi trường nhiều máy ảo Android.
- Thống kê các chỉ số độ đo để đánh giá mô hình cho từng kịch bản thử nghiệm: thay đổi số lượng đặc trưng, thay đổi thuật toán, thay đổi số lượng thiết bị và chu kỳ. Từ đó, phân tích và kết luận tính khả thi cũng như hiệu suất hệ thống.

### 2. Trích xuất đặc trưng

- Máy chủ thử nghiệm sử dụng Python và Androguard để trích xuất permission, intent-filter... từ file AndroidManifest.xml.
- Thiết bị Android thử nghiệm API của lớp PackageManager để trích xuất đặc trưng tương tự từ các ứng dụng đã cài đặt.

### 4. Thiết kế thuật toán Tổng hợp toàn cục

- Tìm hiểu những nghiên cứu đã có về các yếu tố tác động đến việc tổng hợp mô hình toàn cục.
- Phân tích và lựa chọn thuật toán FedAvg, FedDyn... để tổng hợp mô hình cục bộ sau mỗi chu kỳ.



Hình 3. Sơ đồ đề xuất huấn luyện cục bộ và tổng hợp mô hình toàn cục