

Obfuscation (Làm rối mã) – Kỹ thuật và Ứng dụng trong An toàn Phần mềm

[Tên tác giả]

18 tháng 9, 2025

Tóm tắt nội dung

Làm rối mã (*code obfuscation*) là kỹ thuật biến đổi chương trình nhằm bảo vệ bí mật thuật toán và dữ liệu bằng cách gây khó khăn cho việc đọc và phân tích ngược mã. Tài liệu này cung cấp một khảo sát toàn diện về lĩnh vực làm rối mã trong bối cảnh an toàn phần mềm. Mở đầu, chúng tôi trình bày bối cảnh và động lực, định nghĩa chính thức theo tài liệu học thuật kinh điển [4], và mô hình đe dọa giả định. Tiếp theo, một hệ thống phân loại (taxonomy) chi tiết được đề xuất, phân chia các kỹ thuật làm rối mã thành các lớp và nhóm rõ ràng. Trong từng lớp tiêu biểu, chúng tôi đi sâu phân tích một kỹ thuật đại diện – từ nguồn gốc học thuật, nguyên lý hoạt động, ví dụ mã trước và sau khi biến đổi, cho đến đánh giá an toàn và tác động hiệu năng. Ma trận so sánh đa tiêu chí được đưa ra giúp định lượng mức độ *Resilience*, *Stealth*, chi phí hiệu năng và chi phí triển khai của mỗi kỹ thuật. Tài liệu cũng điểm qua các công cụ làm rối mã phổ biến, các phương pháp tấn công và giải rối (deobfuscation) hiện có, cùng các trường hợp ứng dụng thực tế (bảo vệ ứng dụng di động, DRM, malware) và thảo luận về khía cạnh pháp lý, đạo đức. Cuối cùng, chúng tôi đề xuất các hướng nghiên cứu tương lai và đưa ra khuyến nghị thực tiễn cho việc tích hợp kỹ thuật làm rối mã vào quy trình phát triển phần mềm một cách hiệu quả.

Mục lục

1	Giới thiệu	2
1.1	Bối cảnh và Động lực	2
1.2	Định nghĩa chính thức	2
1.3	Mô hình đe dọa	2
1.4	Mục tiêu và Phạm vi	3
2	Phân loại các Kỹ thuật Obfuscation	3
2.1	Obfuscation Bố cục (Layout Obfuscation)	3
2.2	Obfuscation Dữ liệu (Data Obfuscation)	4
2.3	Obfuscation Luồng điều khiển (Control-Flow Obfuscation)	4
2.4	Obfuscation qua Trừu tượng hóa (Abstraction-Based Obfuscation)	4
2.5	Obfuscation Môi trường (Environment-Based Obfuscation)	5
2.6	Các hướng tiếp cận Lý thuyết (Theoretical Paradigms)	5
3	Phân tích Chi tiết các Kỹ thuật Obfuscation Tiêu biểu	7
3.1	Đổi tên Định danh (Identifier Renaming)	7
3.2	Làm phẳng Luồng điều khiển (Control-Flow Flattening)	8
3.3	Mã hóa Chuỗi ký tự (String Encryption)	9
3.4	Obfuscation dạng Máy ảo (Virtualization-based Obfuscation)	11
3.5	Chống gỡ lỗi (Anti-Debugging)	12
3.6	Mật mã Hộp trắng (White-box Cryptography)	13

4	Ma trận Đánh giá và So sánh	15
5	Công cụ và Hệ sinh thái	15
6	Các phương pháp Tấn công và Giải rối mã (Deobfuscation)	17
7	Ứng dụng Thực tế và Các Vấn đề Liên quan	18
7.1	Bảo vệ Ứng dụng Di động	18
7.2	Bảo vệ Quyền nội dung số (DRM)	18
7.3	Mã độc và Phần mềm gian lận	18
7.4	Rủi ro Pháp lý và Đạo đức	19
8	Hướng dẫn Tích hợp vào Quy trình Kỹ thuật	19
9	Hướng Nghiên cứu Tương lai	20
10	Kết luận	21

1 Giới thiệu

1.1 Bối cảnh và Động lực

Trong kỷ nguyên mà mã nguồn và mã bytecode thường được phân phối rộng rãi, nguy cơ bị dịch ngược (decompile) và phân tích trái phép ngày càng cao. Các nền tảng như Java và .NET cho phép kẻ tấn công dễ dàng lấy được mã gần giống mã nguồn từ file biên dịch. Từ đó, tài sản trí tuệ (thuật toán độc quyền, bí mật kinh doanh) và các thông tin nhạy cảm (như khóa mật mã) có thể bị lộ. Bên cạnh đó, trên không gian mạng, lập trình viên xấu có thể chỉnh sửa hoặc tái phân phối trái phép phần mềm của người khác. Để đối phó, nhiều kỹ thuật bảo vệ phần mềm đã ra đời, trong đó làm rối mã được xem là một phương pháp khả thi và thiết thực nhất để ngăn chặn việc phân tích ngược [4]. Việc làm rối mã đúng cách giúp tăng *độ khó* cho đối thủ khi cố gắng hiểu hoặc tái tạo logic chương trình, từ đó bảo vệ hiệu quả tài sản phần mềm. Tuy nhiên, kỹ thuật này cũng có hai mặt: nó có thể bị lạm dụng bởi malware nhằm qua mặt công cụ phòng chống mã độc. Do đó, nhu cầu nghiên cứu một cách hệ thống về làm rối mã là rất cấp thiết [5].

1.2 Định nghĩa chính thức

Định nghĩa

Làm rối mã là quá trình chuyển đổi một chương trình P thành một chương trình P' tương đương về chức năng, nhưng P' khó hiểu và khó phân tích ngược hơn đáng kể so với P . Nói cách khác, code obfuscation không làm thay đổi đầu ra và hành vi của chương trình, nhưng khiến cho việc xác định cấu trúc và ý nghĩa bên trong chương trình trở nên thách thức đối với con người và công cụ phân tích tự động [4].

1.3 Mô hình đe dọa

Trong tài liệu này, chúng tôi giả định mô hình tấn công kiểu *Man-At-The-End* – tức kẻ tấn công có toàn quyền truy cập vào bản chương trình đã biên dịch (binary hoặc bytecode) của nạn nhân. Đối thủ có thể sử dụng mọi kỹ thuật phân tích tĩnh (như disassembler, decompiler, phân tích dòng dữ liệu) và phân tích động (như debug, instrument code khi chạy) để tìm kiếm thông tin nhạy cảm hoặc điểm yếu. Mục tiêu của kẻ tấn công có thể là đánh cắp thuật toán

độc quyền, trích xuất khóa mật mã nhúng trong ứng dụng, hoặc loại bỏ các cơ chế bảo vệ bản quyền (license check). Chúng tôi cũng giả định đối thủ có kiến thức sâu về hệ thống và có thể thực thi chương trình trong môi trường tùy ý (ví dụ: chạy ứng dụng trong trình giả lập hoặc máy ảo để quan sát hành vi). Tuy nhiên, kẻ tấn công không có sẵn mã nguồn gốc trừ khi họ dịch ngược thành công. Với mô hình đe dọa này, nhiệm vụ của kỹ thuật làm rối mã là khiến chi phí và thời gian cần thiết cho việc tấn công tăng lên mức không hiệu quả so với giá trị thu được.

1.4 Mục tiêu và Phạm vi

Mục tiêu của sách trắng này là cung cấp một cái nhìn toàn diện và chuyên sâu về kỹ thuật làm rối mã dành cho cả đối tượng kỹ sư thực hành và nhà nghiên cứu. Chúng tôi tập trung vào các kỹ thuật làm rối mã ở cấp độ phần mềm (mã nguồn, mã máy) và cách áp dụng chúng trong bối cảnh bảo vệ ứng dụng. Các chủ đề liên quan như đóng gói (packing), watermarking hay bảo vệ bằng phần cứng sẽ chỉ được đề cập thoáng qua nếu có liên quan, vì phạm vi chính ở đây là các biến đổi mã chương trình nhằm ẩn giấu logic.

Tài liệu được cấu trúc như sau: Phần 2 trình bày hệ thống phân loại các kỹ thuật làm rối mã, dựa trên nền tảng lý thuyết từ Collberg et al. [4] và các phân loại cập nhật. Phần 3 đi sâu phân tích chi tiết từng kỹ thuật obfuscation tiêu biểu thuộc các nhóm khác nhau, kèm ví dụ minh họa, đánh giá an toàn và hiệu năng, cũng như lưu ý về xử lý chuỗi tiếng Việt trong mã. Phần 4 đưa ra ma trận so sánh các kỹ thuật theo bốn tiêu chí chính để người đọc tiện đánh giá tổng quan. Phần 5 điểm qua các công cụ (tool) thông dụng hỗ trợ làm rối mã trong thực tế. Phần 6 tổng hợp các phương pháp tấn công và giải rối mã, giúp người đọc hiểu được hạn chế của từng kỹ thuật phòng thủ. Trong phần 7, chúng tôi thảo luận các ứng dụng thực tế của obfuscation (bảo vệ ứng dụng di động, DRM và malware) và xem xét khía cạnh pháp lý, đạo đức của việc triển khai kỹ thuật này. Phần 8 đưa ra một checklist hướng dẫn cách tích hợp obfuscation vào quy trình kỹ thuật phần mềm (CI/CD) sao cho hiệu quả. Phần 9 đề xuất các hướng nghiên cứu tiềm năng trong tương lai, đặc biệt trước những thách thức mới (bảo vệ mô hình AI/ML, chống lại phân tích dựa trên phần cứng). Cuối cùng, phần 10 kết luận với những khuyến nghị thực tiễn rút ra từ toàn bộ nội dung đã thảo luận.

2 Phân loại các Kỹ thuật Obfuscation

Dựa trên các nghiên cứu nền tảng [4] cùng những phát triển mới hơn, chúng tôi phân loại các kỹ thuật làm rối mã thành các lớp chính sau đây, mỗi lớp bao gồm một nhóm các kỹ thuật có mục tiêu và bản chất tương đồng. Sơ đồ cây ở Hình 1 minh họa cấu trúc phân loại gồm 3 tầng: **Lớp** (class) → **Nhóm** (category) → **Kỹ thuật cụ thể** (technique).

2.1 Obfuscation Bố cục (Layout Obfuscation)

Lớp này bao gồm các biến đổi hình thức mã nguồn hoặc mã máy nhằm loại bỏ thông tin trực quan giúp người đọc hiểu chương trình nhưng không ảnh hưởng đến logic thực thi. Những kỹ thuật này chủ yếu đánh vào *tính dễ đọc* của mã mà không làm thay đổi luồng điều khiển hay dữ liệu. Các nhóm tiêu biểu trong lớp bố cục gồm:

Đổi tên và tái cấu trúc: Thay đổi tên biến, tên hàm thành các chuỗi vô nghĩa hoặc khó hiểu; thay đổi thứ tự khai báo và sắp xếp lại mã nguồn mà không ảnh hưởng tới luồng chương trình. Ví dụ, đổi tên biến `totalScore` thành `a_1Z` sẽ gây khó khăn cho người đọc trong việc liên hệ biến với ý nghĩa ban đầu.

Loại bỏ thông tin và định dạng: Loại bỏ hoặc mã hóa các chuỗi ký tự chú thích, ký hiệu gỡ lỗi và metadata; xóa bỏ thụt đầu dòng, xuống hàng để mã trở thành một khối liên tục khó quan sát. Mục tiêu là làm mất đi manh mối từ định dạng và thông tin phi chức năng.

Mã vô hại (Dead code): Chèn thêm các lệnh không ảnh hưởng đến kết quả tính toán, chẳng hạn như những đoạn mã không bao giờ được thực thi (vì nằm sau lệnh `return` hoặc trong điều kiện luôn sai), hoặc các phép tính dư thừa không tác động đến biến đầu ra. Mã “rác” này khiến chương trình phình to và khó theo dõi hơn đối với người phân tích.

2.2 Obfuscation Dữ liệu (Data Obfuscation)

Lớp kỹ thuật này tập trung vào việc biến đổi cách thức lưu trữ và xử lý dữ liệu nội tại của chương trình, nhằm ẩn giấu ý nghĩa thực sự của dữ liệu. Mục tiêu là gây khó khăn cho việc phân tích giá trị biến, hằng số và cấu trúc dữ liệu.

Mã hóa và làm mờ dữ liệu: Mã hóa các hằng số nhạy cảm (đặc biệt là chuỗi văn bản) thành dạng không trực quan, chỉ giải mã tại thời điểm cần sử dụng [3]. Ví dụ, một chuỗi khóa API `"secretKey"` có thể lưu trữ dưới dạng một mảng số đã mã hóa và được giải mã tại runtime.

Biến đổi cấu trúc dữ liệu: Thay đổi cách tổ chức dữ liệu: chia nhỏ một biến thành nhiều biến con (splitting), gộp nhiều biến thành một cấu trúc phức tạp (merging), hay tráo đổi thứ tự phần tử trong mảng kèm theo cập nhật chỉ mục. Điều này làm cho việc suy luận quan hệ giữa các dữ liệu trở nên khó khăn hơn.

Ẩn giấu luồng trong dữ liệu: Đưa một phần logic điều khiển chương trình vào dữ liệu, ví dụ sử dụng bảng tra cứu (lookup table) để quyết định luồng thay vì các câu lệnh điều kiện rõ ràng. Kỹ thuật này làm mờ ranh giới giữa code và data, khiến người phân tích khó phát hiện logic điều khiển thực sự.

2.3 Obfuscation Luồng điều khiển (Control-Flow Obfuscation)

Đây là lớp kỹ thuật làm rối trực tiếp cấu trúc luồng thực thi của chương trình, khiến biểu đồ luồng điều khiển (CFG) của chương trình trở nên phức tạp hoặc khác xa so với cấu trúc ban đầu.

Biến đổi cấu trúc điều khiển: Thay các cấu trúc điều khiển chuẩn (if-else, vòng lặp) bằng các cấu trúc tương đương nhưng phức tạp hơn. Một ví dụ kinh điển là kỹ thuật *làm phẳng luồng điều khiển* (control-flow flattening) [7]: chương trình được bao trong một vòng lặp vô hạn và một câu lệnh `switch`, mọi nhánh logic được chuyển thành các `case` trong switch, luồng đi được điều hướng bằng một biến trạng thái.

Chèn điều kiện opaque: Thêm vào chương trình các điều kiện kiểm tra sử dụng *điều kiện mờ đục* (opaque predicate) – tức những biểu thức logic luôn luôn trả về true hoặc false nhưng người phân tích khó nhận biết tính chất đó [4]. Mục đích là tạo các nhánh giả phức tạp: ví dụ, một `if (X*X - X % 2 == 0)` luôn đúng với mọi số nguyên X nhưng không hiển nhiên khi nhìn lướt qua, khiến mã phân nhánh nhiều đường vô nghĩa để đánh lừa công cụ phân tích tĩnh.

Giả cấu trúc và nhảy không tuần tự: Phá vỡ tính tuần tự của mã bằng cách sử dụng `goto` hoặc các lệnh nhảy không theo cấu trúc lập trình thông thường. Code có thể bị “xé nhỏ” và nhảy qua lại, làm người đọc khó theo dõi luồng. Kết hợp với các kỹ thuật trên, điều này làm cho CFG trở nên rối rắm với nhiều cạnh và nút giả.

2.4 Obfuscation qua Trừu tượng hóa (Abstraction-Based Obfuscation)

Các kỹ thuật lớp này nâng cao mức độ trừu tượng của chương trình, hoặc chuyển đổi chương trình sang một hình thức thực thi khó theo dõi hơn, chẳng hạn như thông qua một máy ảo trung gian.

Máy ảo hóa mã: Chuyển đổi một phần hoặc toàn bộ chương trình ban đầu thành bytecode của một máy ảo tùy biến, sau đó chèn vào chương trình một bộ thông dịch (interpreter) cho bytecode đó [1]. Người tấn công khi phân tích sẽ phải đối mặt với mã máy ảo phức tạp thay vì mã gốc. Ví dụ, thay vì thực thi trực tiếp các phép toán, chương trình sẽ chứa một máy ảo với các opcode riêng và bộ xử lý tương ứng, ẩn đi logic gốc bên trong.

Đa hình và biến hình: Tạo ra nhiều phiên bản mã khác nhau (về mặt nhị phân) từ cùng một mã nguồn, thường bằng cách thay thế các phần của chương trình bằng những đoạn mã tương đương nhưng khác biệt về triển khai. Kỹ thuật này làm giảm khả năng nhận dạng mã dựa trên mẫu bởi mỗi bản phân phối có fingerprint riêng.

Trừu tượng hóa qua ngôn ngữ kịch bản: Đưa mã nhảy cảm sang dạng script hoặc ngôn ngữ khác chạy trong môi trường runtime (như JavaScript, Lua), sau đó nhúng script này vào chương trình chính. Việc này tách biệt logic sang một ngôn ngữ khác khiến quá trình phân tích phức tạp hơn do phải hiểu thêm tầng ngôn ngữ trung gian.

2.5 Obfuscation Môi trường (Environment-Based Obfuscation)

Không chỉ bản thân mã, các kỹ thuật còn có thể khai thác môi trường thực thi để tăng cường bảo vệ.

Chống gỡ lỗi và phân tích động: Tích hợp các cơ chế phát hiện đang bị debug (ví dụ: kiểm tra flag của hệ điều hành, hoặc đo thời gian thực thi bất thường) và thực hiện hành động phản ứng (dừng chương trình, đánh lừa debugger). Tương tự, chống phân tích động bao gồm phát hiện chạy trong máy ảo hoặc sandbox và thay đổi hành vi chương trình khi phát hiện (anti-VM, anti-emulation).

Ràng buộc môi trường: Ràng buộc phần mềm với một môi trường cụ thể: ví dụ, chỉ chạy nếu phát hiện phần cứng/thiết bị hợp lệ (theo một ID duy nhất), hoặc chỉ hoạt động nếu tích hợp với một server xác thực. Mọi nỗ lực tách phần mềm ra khỏi môi trường đó sẽ làm chương trình bị vô hiệu.

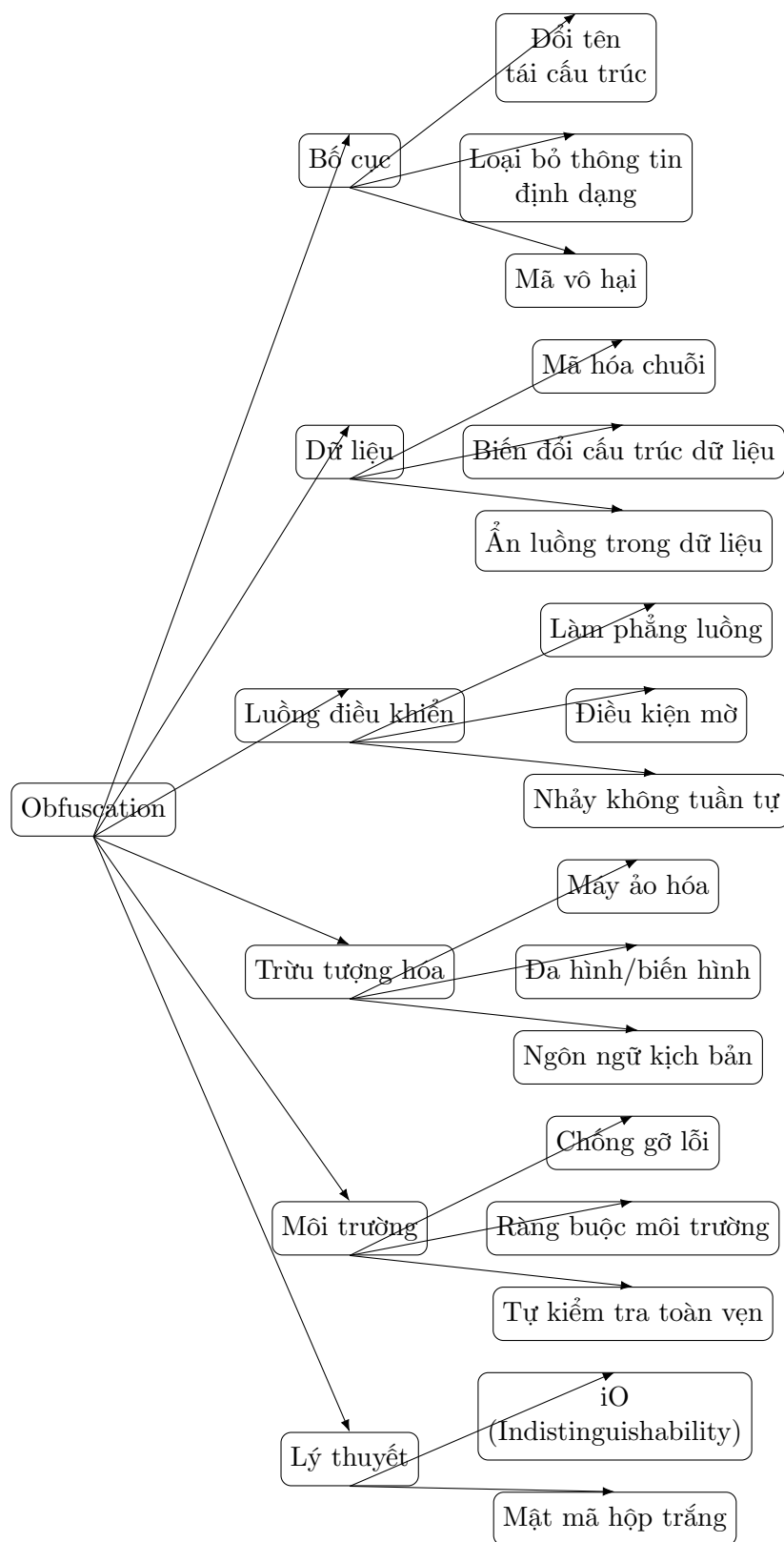
Kiểm tra và làm rối tính toàn vẹn: Chèn các checksum, hash kiểm tra mã code tại runtime để phát hiện xem mã có bị chỉnh sửa (patch) hay không. Đồng thời, có thể tự mã hóa một số phần của chính chương trình, chỉ giải mã khi chạy (self-modifying code). Điều này gây khó cho kẻ tấn công vì nếu sửa đổi mã thì chương trình có thể tự phát hiện và ngừng hoạt động.

2.6 Các hướng tiếp cận Lý thuyết (Theoretical Paradigms)

Bên cạnh các kỹ thuật thực tiễn kể trên, có những hướng tiếp cận mang tính hàn lâm và lý thuyết cao, nhằm đạt được mức bảo mật “hoàn hảo” về mặt toán học.

Làm rối bất khả phân biệt (Indistinguishability Obfuscation): Đây là khái niệm trong mật mã học, theo đó một bộ làm rối mã O được gọi là iO nếu với mọi hai chương trình P_1, P_2 có cùng chức năng, $O(P_1)$ và $O(P_2)$ không thể phân biệt được với nhau về mặt quan sát [2]. Nghiên cứu về iO rất sôi động trong thập kỷ qua, tuy nhiên cho đến nay chưa có giải pháp iO nào đủ hiệu quả để ứng dụng rộng rãi.

Mật mã hộp trắng (White-box Cryptography): Đây là hướng tiếp cận bảo vệ các thuật toán mật mã khi thực thi trong môi trường mà đối thủ có toàn quyền quan sát. Một ví dụ điển hình là mã hóa White-box AES [3], trong đó khóa bí mật được trộn lẫn vào các bảng tra cứu thay vì lưu rõ ràng, khiến kẻ tấn công rất khó trích xuất khóa dù có được mã chương trình. White-box crypto thực chất là sự kết hợp giữa lý thuyết mật mã và làm rối mã ở cấp độ thao tác bit và byte.



Hình 1: Sơ đồ phân loại các kỹ thuật làm rối mã theo lớp, nhóm và kỹ thuật cụ thể.

3 Phân tích Chi tiết các Kỹ thuật Obfuscation Tiêu biểu

Dưới đây, chúng tôi chọn lọc và phân tích sâu một số kỹ thuật làm rối mã tiêu biểu, mỗi kỹ thuật đại diện cho một lớp đã nêu ở trên nhằm đảm bảo tính bao quát. Mỗi tiểu mục cung cấp thông tin về nguồn gốc học thuật, nguyên lý hoạt động, ví dụ mã trước và sau khi áp dụng, phân tích mức độ an toàn, tác động hiệu năng, và lưu ý đặc biệt liên quan đến xử lý Tiếng Việt (nếu có).

3.1 Đổi tên Định danh (Identifier Renaming)

Nguồn gốc học thuật. Kỹ thuật đổi tên biến, hàm và các định danh khác xuất hiện từ những công cụ làm rối mã thương mại đầu tiên dành cho Java (như SourceGuard, Dotfuscator). Về mặt học thuật, Collberg et al. [4] đã mô tả đây là một trong các phép biến đổi bố cục đơn giản nhưng hiệu quả thấp. Hầu hết các bộ công cụ obfuscation (ví dụ: ProGuard cho Java) đều hỗ trợ tính năng này vì tính dễ triển khai và không làm thay đổi logic chương trình.

Nguyên lý hoạt động. Thuật toán rất trực quan: duyệt toàn bộ bảng ký hiệu (symbol table) của chương trình, thay thế mỗi tên định danh (tên biến, tên hàm, lớp, v.v.) bằng một tên mới không gợi ý ý nghĩa. Tên mới có thể được tạo ngẫu nhiên (ví dụ: `a`, `b1`, `Xyz`) hoặc bằng các ký tự khó phân biệt. Một biến thể tinh vi hơn là dùng các ký tự Unicode trông giống chữ cái Latin (homoglyph) để đặt tên biến, khiến người đọc nhầm lẫn (ví dụ: ký tự Cyrillic thay cho `a` Latin). Kỹ thuật này không tác động đến mã máy đã biên dịch trừ phần tên (trong trường hợp ngôn ngữ có thông tin tên trong binary, như .NET, Java bytecode). Với ngôn ngữ biên dịch tĩnh như C/C++, tên biến cục bộ không tồn tại trong binary, do đó đổi tên chỉ hữu ích ở mức mã nguồn (nhưng tên hàm xuất hoặc symbol public vẫn có thể đổi được).

Ví dụ minh họa (Trước & Sau). Đoạn mã dưới đây cho thấy hiệu quả của việc đổi tên trên một hàm đơn giản:

Đoạn mã 1: Mã trước khi đổi tên định danh

```
1 // H m t nh giai t h a
2 int tinhGiaithua(int n) {
3     if (n <= 1) return 1;
4     return n * tinhGiaithua(n - 1);
5 }
```

Đoạn mã 2: Mã sau khi đổi tên định danh thành các chuỗi vô nghĩa

```
1 int f1(int x) {
2     if (x <= 1) return 1;
3     return x * f1(x - 1);
4 }
```

Chú thích: Trong ví dụ trên, tên hàm `tinhGiaithua` đã được đổi thành `f1`, tên biến `n` đổi thành `x`. Các tên mới không mang ý nghĩa, gây khó khăn hơn đôi chút cho người đọc trong việc hiểu mục đích của hàm.

Phân tích An toàn. Xét về *resilience* (khả năng chống chịu phân tích tự động), kỹ thuật này rất yếu. Các công cụ dịch ngược hiện đại hầu như không gặp trở ngại gì với mã đã đổi tên – chúng vẫn tái tạo được cấu trúc điều khiển và dữ liệu, chỉ là các tên biến mất nghĩa [9]. Về mặt gây khó khăn cho con người (*potency*), đổi tên chỉ có tác dụng ở mức bề mặt: lập trình viên sẽ mất thời gian hơn để hiểu vai trò của từng biến/hàm, nhưng với ngữ cảnh đủ lớn, họ vẫn suy luận ra được. Kỹ thuật này không thêm bất kỳ lớp bảo vệ nào chống lại phân tích động.

Tác động Hiệu năng. *Không đáng kể.* Đổi tên định danh hoàn toàn không làm thay đổi mã sinh ra khi thực thi, do đó không có overhead về tốc độ hay dung lượng bộ nhớ. Một ngoại lệ nhỏ: trong ngôn ngữ chạy trên máy ảo (Java, .NET), việc loại bỏ tên có thể giảm nhẹ kích thước file (do tên ngắn hơn và có thể xóa bớt metadata), nhưng ảnh hưởng này không đáng kể.

Lưu ý: Kỹ thuật đổi tên chỉ tác động tới tên định danh, không liên quan đến dữ liệu chuỗi hay ký tự Unicode. Do đó, nó không gây ảnh hưởng trực tiếp đến các chuỗi tiếng Việt trong mã nguồn. Tuy nhiên, nếu mã nguồn ban đầu có sử dụng biến hoặc hàm đặt tên bằng tiếng Việt có dấu, việc đổi tên sẽ thay thế chúng bằng tên mới (thường bằng ASCII), dẫn đến mất hoàn toàn dấu vết tiếng Việt trong mã sau khi làm rồi. Điều này có thể gây chút khó khăn khi gỡ rối lỗi (debug) nếu lập trình viên quen thuộc với các tên biến tiếng Việt, nhưng về chức năng không có vấn đề gì.

3.2 Làm phẳng Luồng điều khiển (Control-Flow Flattening)

Nguồn gốc học thuật. Kỹ thuật làm phẳng luồng điều khiển lần đầu được đề cập bởi D. Low trong luận văn Thạc sĩ về bảo vệ Java bytecode [7]. Sau đó, nhiều nghiên cứu đã mở rộng và phân tích tính hiệu quả của kỹ thuật này [1]. Hiện nay, làm phẳng luồng điều khiển được triển khai trong một số công cụ như Obfuscator-LLVM (OLLVM) [6] và Tigress, cho thấy tính thực tiễn của nó trong việc gây khó khăn cho việc tách và hiểu các khối mã lệnh.

Nguyên lý hoạt động. Thay vì để chương trình thực thi tuần tự qua các cấu trúc phân nhánh tự nhiên, kỹ thuật này đưa mọi khối lệnh (basic block) về cùng một mức trong một vòng lặp lớn. Một biến trạng thái (**state**) được sử dụng để quyết định khối lệnh nào sẽ thực thi tiếp theo. Cụ thể:

- Tạo một vòng lặp vô hạn (hoặc vòng lặp `while(true)`) bao bọc toàn bộ thân hàm.
- Bên trong vòng lặp, dùng một cấu trúc `switch-case` trên biến `state`. Mỗi `case` tương ứng với một khối lệnh gốc của hàm.
- Cuối mỗi khối (cuối mỗi case), cập nhật `state` giá trị mới tương ứng với khối lệnh kế tiếp sẽ thực thi, rồi `break` để quay lại đầu vòng lặp.
- Vòng lặp kết thúc khi `state` nhận giá trị đặc biệt (ví dụ `state = -1`) và có một `case` cho giá trị đó thực hiện lệnh `break` khỏi vòng lặp.

Kết quả là luồng điều khiển nguyên gốc (với nhiều nhánh if/else, vòng lặp lồng nhau) bị “làm phẳng” thành duy nhất một vòng lặp + switch đơn cấp. Người phân tích tĩnh sẽ thấy một vòng lặp vô hạn phức tạp và rất khó để tái tạo lại cấu trúc điều khiển ban đầu.

Ví dụ Demo Trước và Sau: Để minh họa, xem đoạn mã giả C dưới đây trước và sau khi áp dụng làm phẳng luồng điều khiển. Mã gốc có một cấu trúc rẽ nhánh đơn giản:

Đoạn mã 3: Mã gốc với cấu trúc điều kiện đơn giản

```
1 if (x < 0) {
2     foo();
3 } else {
4     bar();
5 }
```

Sau khi làm phẳng, mã được chuyển thành dạng dùng `while` và `switch` với biến trạng thái:

Đoạn mã 4: Mã sau khi áp dụng kỹ thuật làm phẳng luồng điều khiển

```
1 int state = 0;
2 while (1) {
3     switch(state) {
4         case 0:
5             if (x < 0) {
6                 state = 1;
7             } else {
8                 state = 2;
9             }
10    }
```



```

10         break;
11     case 1:
12         foo();
13         state = 3;
14         break;
15     case 2:
16         bar();
17         state = 3;
18         break;
19     case 3:
20         goto EXIT;
21 }
22 }
23 EXIT: ;

```

Chú thích: Ban đầu, `state = 0` vào case 0, tại đây xác định sẽ nhảy sang `state = 1` hoặc 2 tương ứng với nhánh `if/else`. Sau khi thực thi `foo()` (case 1) hoặc `bar()` (case 2), chương trình đặt `state = 3` và `break` để chuyển đến case 3, nơi lệnh `goto EXIT` dùng để thoát khỏi vòng lặp vô hạn.

Phân tích An toàn. Về *resilience*, làm phẳng luồng điều khiển được xem là kỹ thuật khá mạnh. Các công cụ decompiler tĩnh gặp nhiều khó khăn để tái tạo mã nguồn có cấu trúc từ mã máy đã bị làm phẳng: vòng lặp vô hạn và switch-case khiến decompiler thường phải giữ nguyên dưới dạng `goto` và `label` khó hiểu, thay vì `if/else` gốc. Ngay cả các hệ thống phân tích dòng lệnh nâng cao cũng bị tăng độ phức tạp [1]. Tuy nhiên, với phân tích động hoặc symbolic execution, kỹ thuật này không phải không thể vượt qua: kẻ tấn công có thể chạy chương trình và lần lượt khám phá các giá trị của `state` để suy ra cấu trúc ban đầu. *Stealth* của kỹ thuật này khá thấp – cấu trúc vòng lặp vô hạn và switch với nhiều case trông rất bất thường so với mã thông thường, dễ bị phát hiện bằng mắt (dấu hiệu đặc trưng).

Tác động Hiệu năng. *Trung bình.* Mã bị làm phẳng thường chậm hơn một chút do chi phí của vòng lặp và switch trong mỗi lần lặp. Cụ thể, mỗi khối lệnh thay vì nhảy thẳng tới khối kế tiếp, giờ đây phải qua bước `break` và kiểm tra switch, gây thêm một tầng overhead. Ngoài ra, các lệnh `goto` hoặc `switch` có thể ảnh hưởng tới ưu hóa CPU (làm giảm khả năng dự đoán nhánh). Trong nhiều trường hợp kỹ thuật, overhead hiệu năng chỉ ở mức vài phần trăm, nhưng nếu hàm làm phẳng có rất nhiều khối cơ bản, thời gian trong vòng lặp sẽ tăng đáng kể.

Lưu ý: Kỹ thuật làm phẳng luồng điều khiển hoạt động ở mức cấu trúc mã lệnh, không liên quan đến xử lý chuỗi hoặc ký tự Unicode. Do đó, nó không ảnh hưởng trực tiếp đến các chuỗi ký tự tiếng Việt trong chương trình. Ví dụ, nếu chương trình gốc in ra chuỗi tiếng Việt, sau khi làm phẳng thì chuỗi đó vẫn in bình thường, không bị biến đổi. Không có rủi ro đặc biệt nào liên quan đến mã hóa UTF-8 khi áp dụng kỹ thuật này.

3.3 Mã hóa Chuỗi ký tự (String Encryption)

Nguồn gốc học thuật. Việc mã hóa (hay còn gọi là *ẩn* hoặc *xáo trộn*) các chuỗi văn bản bên trong mã nguồn là một kỹ thuật phổ biến, được cả các nhà phát triển phần mềm thương mại lẫn tác giả malware sử dụng. Mặc dù không có một bài báo “kinh điển” riêng về kỹ thuật này, nó thường được nhắc đến trong các khảo sát về obfuscation và bảo mật mã độc [1]. Một số công cụ chuyên dụng như DexGuard (cho Android) và nhiều packer mã độc cung cấp sẵn chức năng mã hóa chuỗi.

Nguyên lý hoạt động. Tất cả các chuỗi văn bản dạng rõ (plain-text) trong mã – chẳng hạn thông báo lỗi, URL, khóa API – sẽ được chuyển thành dạng mã hóa để tránh bị lộ khi kẻ tấn công duyệt qua phần dữ liệu của file nhị phân. Có nhiều phương pháp mã hóa:

- Mã hóa đối xứng đơn giản: Ví dụ, XOR từng byte ký tự với một khóa cố định; hoặc dịch chuyển giá trị ASCII (Caesar cipher). Kết quả là chuỗi trở thành một mảng byte khó đọc.

- Mã hóa bằng thuật toán tiêu chuẩn: Sử dụng AES, DES... với một khóa bí mật nhúng sẵn trong chương trình để mã hóa chuỗi. Khi cần dùng đến chuỗi (ví dụ hiển thị thông báo), chương trình sẽ giải mã tạm thời.
- Biến đổi tùy ý: Xáo trộn thứ tự ký tự, thay thế bằng mã Unicode vô nghĩa, hay thậm chí nén lại (compression) kèm mã giải nén.

Chương trình sau đó phải chứa thêm code giải mã/giải nén các chuỗi này khi cần thiết ở runtime. Phần code này cũng có thể được làm rối để tránh dễ dàng tìm được khóa hoặc thuật toán.

Ví dụ minh họa (Trước & Sau). Xét tình huống đơn giản: chương trình in ra một thông báo chứa tiếng Việt:

Đoạn mã 5: Mã trước khi mã hóa chuỗi

```
1 printf("Xin chào! Khá API của bạn khá hợp lý.");
```

Sau khi áp dụng kỹ thuật mã hóa chuỗi, chuỗi rõ trên sẽ không xuất hiện trực tiếp trong file nhị phân. Thay vào đó, ta có thể thấy một mảng byte mã hóa và đoạn code giải mã:

Đoạn mã 6: Mã sau khi áp dụng mã hóa chuỗi ký tự

```
1 // Mảng byte để mã hóa (vòng lặp XOR với khóa 0x5A)
2 unsigned char s[] = {0x13, 0x3F, 0x3D, ..., 0x00};
3 int len = sizeof(s)/sizeof(s[0]);
4 for(int i=0; i<len-1; i++) {
5     s[i] ^= 0x5A; // Giải mã XOR từng byte
6 }
7 printf("%s", (char*)s);
```

Chú thích: Mảng `s` chứa dữ liệu đã mã hóa (các giá trị hex). Vòng lặp `for` thực hiện giải mã XOR với khóa `0x5A` cho từng byte (bỏ qua byte kết thúc chuỗi `0x00`). Sau vòng lặp, mảng `s` được khôi phục thành chuỗi gốc và in ra. Người đọc tinh sẽ không thấy chuỗi "Xin chào!..." trong mã, thay vào đó chỉ thấy các byte khó hiểu và thuật toán giải mã.

Phân tích An toàn. Kỹ thuật mã hóa chuỗi nâng cao *stealth* đáng kể cho dữ liệu nhạy cảm: các chuỗi rõ không còn hiện diện, tránh bị tìm thấy dễ dàng bằng cách quét binary. Về *resilience*, nó chống lại phân tích tĩnh khá tốt – trừ phi công cụ phân tích có thể tự động nhận diện mẫu mã và thực thi giả lập để giải mã, nếu không chuỗi vẫn ẩn. Tuy nhiên, đối với phân tích động, kỹ thuật này không quá bền: khi chương trình chạy, nó buộc phải giải mã chuỗi (như trong ví dụ trên), do đó kẻ tấn công có thể đặt breakpoint hoặc theo dõi bộ nhớ để lấy được chuỗi gốc. Các công cụ dynamic instrumentation như Frida có thể hook ngay tại hàm `printf` hoặc hàm giải mã để thu thập dữ liệu trước khi in ra. Do đó, mã hóa chuỗi chỉ thực sự ngăn cản những kẻ tấn công không có khả năng phân tích động.

Tác động Hiệu năng. *Thấp.* Việc mã hóa/giải mã chuỗi thường chỉ tốn vài phép toán trên mỗi byte, ít khi ảnh hưởng đáng kể đến tổng thể hiệu năng ứng dụng (trừ khi có một lượng rất lớn chuỗi phải giải mã liên tục). Tuy nhiên, cần lưu ý về chi phí bộ nhớ: nếu nén chuỗi hoặc mã hóa phức tạp, có thể cần bộ nhớ tạm để lưu trữ kết quả giải mã. Trong đa số trường hợp, chi phí này không đáng kể so với khả năng chấp nhận của ứng dụng.

Lưu ý: Đối với chuỗi tiếng Việt chứa ký tự Unicode nhiều byte, kỹ thuật mã hóa chuỗi vẫn hoạt động bình thường nếu được triển khai cẩn thận. Lập trình viên phải chú ý mã hóa và giải mã theo đúng định dạng byte của chuỗi UTF-8. Ví dụ: chữ "á" trong "chào" được biểu diễn bằng 2 byte (`0xC3 0xA1`). Nếu dùng XOR, cần áp dụng XOR trên từng byte; nếu nén hoặc mã hóa phức tạp hơn, cần đảm bảo quá trình giải mã phục hồi đúng chuỗi UTF-8 ban đầu để không bị lỗi hiển thị. Tóm lại, kỹ thuật này không phân biệt nội dung là ASCII hay Unicode, miễn là xử lý theo byte một cách đồng nhất.

3.4 Obfuscation dạng Máy ảo (Virtualization-based Obfuscation)

Nguồn gốc học thuật. Sử dụng máy ảo tùy biến để bảo vệ mã là một kỹ thuật mạnh mẽ xuất phát từ lĩnh vực bảo vệ phần mềm chống sao chép và mã độc. Năm 2008, nghiên cứu của Sharif et al. giới thiệu việc malware dùng máy ảo giả lập CPU để gây khó cho việc phân tích [1]. Trong công nghiệp, VMProtect và Themida là các packer thương mại nổi tiếng áp dụng thủ thuật “ảo hóa” mã. Về học thuật, Banescu et al. (2016) cũng phân tích độ bền vững của mã ảo hóa trước tấn công symbolic [1]. Công cụ Tigress của Collberg cung cấp chế độ chuyển mã C thành bytecode VM và bộ thông dịch, cho phép thực nghiệm rộng rãi kỹ thuật này.

Nguyên lý hoạt động. Ý tưởng cốt lõi là định nghĩa một máy ảo giả lập (có thể hiểu như một CPU đơn giản do ta tự thiết kế), rồi biên dịch đoạn mã cần bảo vệ thành các “lệnh” (bytecode) của máy ảo đó. Chương trình gốc sẽ được thay thế bằng:

- Một mảng hằng số chứa bytecode (các opcode và toán hạng) biểu diễn logic gốc.
- Một hàm thông dịch (interpreter) có cấu trúc vòng lặp `while` đọc từng opcode và thực hiện hành động tương ứng trên một mô phỏng bộ trạng thái (thanh ghi, stack, v.v.).

Do máy ảo là tùy biến, không công khai đặc tả, nên người phân tích phải tốn công “dịch ngược” bộ thông dịch mới hiểu được từng lệnh ảo nghĩa là gì. Có thể hình dung mức độ phức tạp như phải dịch ngược một chương trình viết bằng một ngôn ngữ lạ. Kỹ thuật này có thể kết hợp với các obfuscation khác: ví dụ, mã thông dịch cũng có thể bị làm rối (flatten, opaque predicates), bytecode có thể được mã hóa, v.v.

Ví dụ Demo Trước và Sau: Giả sử ta có hàm đơn giản tính tổng hai số:

Đoạn mã 7: Mã gốc cần bảo vệ bằng ảo hóa

```
1 int add(int a, int b) {
2     return a + b;
3 }
```

Sau khi áp dụng máy ảo hóa, chương trình có thể biến thành:

Đoạn mã 8: Mã sau khi được chuyển sang bytecode và bộ thông dịch máy ảo

```
1 // Định nghĩa các opcode cho VM
2 #define OP_LOAD_A 0x01
3 #define OP_LOAD_B 0x02
4 #define OP_ADD    0x03
5 #define OP_STORE  0x04
6 #define OP_HALT   0xFF
7 // Bytecode cần thực hiện (a + b)
8 unsigned char code[] = { OP_LOAD_A, OP_LOAD_B, OP_ADD, OP_HALT };
9 // Bộ thông dịch VM
10 int run_vm(int a, int b) {
11     int regA = 0, regB = 0, acc = 0;
12     int pc = 0;
13     while (1) {
14         unsigned char op = code[pc++];
15         switch(op) {
16             case OP_LOAD_A: acc = a; break;
17             case OP_LOAD_B: acc = b; break;
18             case OP_ADD:    regA = acc; acc = regA + b; break;
19             case OP_HALT:   return acc;
20         }
21     }
22 }
```

Chú thích: Bytecode `code[]` gồm các lệnh: nạp biến a, nạp biến b, thực hiện cộng, rồi dừng. Hàm `run_vm` thực thi mã ảo này: khi gặp `OP_LOAD_A` thì đưa tham số a vào `acc`, gặp `OP_ADD` thì lấy `acc` (đang chứa a) cộng với b, v.v. Kết quả cuối cùng trả về qua `acc`. Mặc dù ví dụ đơn

giản, nhưng với một chương trình lớn, bytecode có thể rất dài và phức tạp, làm cho ai muốn phân tích phải dò từng case trong **switch** tương ứng với từng opcode.

Phân tích An toàn. Về *resilience*, obfuscation dạng máy ảo hiện được coi là một trong những kỹ thuật mạnh nhất. Reverse engineer sẽ phải tốn công trích xuất được bytecode, hiểu rõ cơ chế VM, rồi viết trình giải mã logic. Nhiều nghiên cứu chỉ ra rằng mã ảo hóa nếu kết hợp với các biến đổi phụ có thể "cầm cự" đáng kể trước tấn công tự động [1]. Dẫu vậy, bản thân kỹ thuật này không phải là bất khả xâm phạm: các công cụ symbolic execution mạnh có thể khôi phục phần nào logic bằng cách coi máy ảo như một hàm toán học và giản ước nó. Ngoài ra, đã có các công cụ chuyên dụng để devirtualize (giải ảo) cho một số VM thông dụng bằng cách phân tích mẫu. Về *stealth*, do máy ảo quá khác thường so với mã thông thường, nó có dấu hiệu đặc trưng (rất nhiều lệnh switch-case, mảng bytecode lớn), do đó dễ nhận biết sự hiện diện của obfuscation này.

Tác động Hiệu năng. *Cao.* Cái giá phải trả cho sự bảo vệ của VM là hiệu năng. Chương trình bị chậm đi rõ rệt do mỗi lệnh gốc giờ thực thi qua nhiều bước thông dịch (fetch opcode, switch, thực hiện hành vi). Thông thường, mã ảo hóa có thể chậm hơn hàng chục lần so với mã gốc nếu ảo hóa những phần tính toán nhiều. Kích thước code cũng tăng (phải chứa bộ thông dịch, bảng bytecode). Việc CPU không thể tối ưu nhánh dự đoán trong VM loop càng làm giảm hiệu quả. Vì lý do đó, kỹ sư thường chỉ áp dụng ảo hóa cho một vài đoạn code quan trọng thay vì toàn bộ chương trình.

Lưu ý: Kỹ thuật máy ảo hóa không ảnh hưởng tới xử lý chuỗi ký tự hay encoding, vì nó thay đổi cách thực thi logic chứ không động đến dữ liệu chuỗi. Nếu trong đoạn code được ảo hóa có xử lý chuỗi tiếng Việt, thì khi chạy máy ảo kết quả cũng tương tự (miễn là VM được viết đúng chức năng). Tuy nhiên, do VM thường hoạt động ở mức rất thấp (toán hạng số nguyên, thanh ghi), nên chuỗi Unicode vẫn sẽ được thao tác ở dạng các byte mã UTF-8 như bình thường. Tóm lại, không có rủi ro nào đặc thù cho dữ liệu tiếng Việt khi dùng kỹ thuật này.

3.5 Chống gỡ lỗi (Anti-Debugging)

Nguồn gốc học thuật. Các kỹ thuật chống gỡ lỗi và phân tích động thường không xuất phát từ một nghiên cứu đơn lẻ mà được tổng hợp qua kinh nghiệm thực tiễn của giới an toàn thông tin. Nhiều tài liệu kỹ thuật (như Microsoft SDK, Linux manpage) liệt kê các API và cơ chế có thể dùng để phát hiện debugger. Về học thuật, Chen et al. (2008) có đề cập anti-debug như một thành phần trong chiến lược bảo vệ phần mềm đa tầng, và nhiều phân tích malware case study cũng liệt kê các chiêu thức anti-debug mà malware sử dụng.

Nguyên lý hoạt động. Ý tưởng chính là lập trình viên thêm vào chương trình các đoạn code kiểm tra xem môi trường runtime có dấu hiệu đang bị debug/hook hay không, từ đó thực hiện hành động đối phó. Một số kỹ thuật tiêu biểu:

- Kiểm tra flag hệ thống: Trên Windows, hàm `IsDebuggerPresent()` trả về true nếu tiến trình hiện tại đang gắn với debugger. Ngoài ra, có thể kiểm tra một số trường trong *Process Environment Block* (PEB) của tiến trình (các cờ `BeingDebugged`, `NtGlobalFlag`).
- Bẫy thời gian: Chèn các lệnh `Sleep()` có thời lượng ngắn trong vòng lặp và đo thời gian thực thi. Nếu có debugger, do bước qua/thiết lập breakpoint tốn thời gian, tổng thời gian sẽ dài hơn ngưỡng cho phép, khi đó chương trình nghỉ ngơi và tự thoát.
- Gây sự kiện đặc biệt: Ví dụ, trên x86, lệnh `INT 3` (breakpoint) sẽ làm chương trình văng nếu không có debugger bắt, ngược lại nếu có debugger thì nó can thiệp làm chương trình tiếp tục. Dựa vào điều này, chương trình có thể thực hiện `INT 3` và bắt ngoại lệ; nếu không thấy ngoại lệ (tức không ai xử lý, nghĩa là có debugger), chương trình sẽ biết.

- Chống tháo gỡ: Sử dụng các thủ thuật như làm rối `Import Table` hoặc gọi API theo cách động (Runtime API resolving) để debugger khó theo dõi; hoán đổi thứ tự các sự kiện khởi tạo để debugger không kịp gắn vào.

Khi phát hiện có debugger, chương trình thường sẽ thực hiện hành động như thoát ngay lập tức, hoặc thậm chí giả vờ chạy sai (cung cấp kết quả sai) nhằm đánh lạc hướng người debug.

Ví dụ minh họa (Trước & Sau). Mã gốc:

Đoạn mã 9: Mã gốc không có chống debug

```
1 int main() {
2     do_sensitive_task();
3     return 0;
4 }
```

Mã sau khi thêm cơ chế chống debug:

Đoạn mã 10: Mã có cơ chế phát hiện debugger và phản ứng

```
1 #include <windows.h>
2 int main() {
3     if (IsDebuggerPresent()) {
4         ExitProcess(1); // Tho t n u đang b debug
5     }
6     do_sensitive_task();
7     return 0;
8 }
```

Chú thích: Ở đây, chương trình gọi API `IsDebuggerPresent()` của Windows. Nếu hàm trả về true (nghĩa là có debugger), chương trình ngay lập tức thoát (mã exit code 1). Nếu không, nó thực hiện công việc nhạy cảm bình thường. Trong thực tế, mã anti-debug có thể phức tạp hơn nhiều và được rải ở nhiều chỗ trong chương trình.

Phân tích An toàn. Về mặt chống lại phân tích thủ công (con người), anti-debug gây phiền toái lớn: người muốn debug phải tìm cách vượt qua các cơ chế này (ví dụ: vá nhị phân để skip `IsDebuggerPresent` hoặc sử dụng các công cụ chuyên dụng chặn API). Tuy vậy, anti-debug hầu như không ảnh hưởng đến phân tích tĩnh tự động (vì các check này chỉ kích hoạt khi chạy). Đối với kẻ tấn công có kinh nghiệm, các kỹ thuật anti-debug quen thuộc có thể bị vô hiệu hóa khá nhanh (vd: dùng plugin để giả lập phản hồi API). Nói cách khác, *resilience* không cao trước một đối thủ cứng, nhưng *stealth* lại tốt: chương trình trông có vẻ chạy bình thường, chỉ khi debug mới xảy ra vấn đề, do đó phân tích viên có thể mất thời gian mới nhận ra đang bị anti-debug "chơi khăm".

Tác động Hiệu năng. *Không đáng kể/Thấp.* Các check anti-debug như đọc flag hay gọi API thường rất nhẹ so với logic chính. Chúng chủ yếu chạy một lần lúc khởi động hoặc định kỳ giãn cách nên khó gây ảnh hưởng hiệu năng đáng kể. Một số chiêu phức tạp hơn (như bẫy thời gian trong vòng lặp chặt chẽ) có thể làm chậm chương trình nếu sử dụng không cẩn thận, nhưng nhìn chung overhead thấp hơn nhiều so với các kỹ thuật obfuscation khác.

Lưu ý: Các kỹ thuật chống debug không liên quan đến xử lý chuỗi ký tự hay dữ liệu Unicode, nên chúng không gây vấn đề gì đặc thù với Tiếng Việt. Chẳng hạn, việc gọi `IsDebuggerPresent()` hay đo thời gian thực thi không ảnh hưởng đến cách chương trình xử lý chuỗi Unicode. Do đó, nếu một ứng dụng có giao diện tiếng Việt hoặc xử lý dữ liệu tiếng Việt, việc thêm anti-debug sẽ không làm thay đổi hay lỗi các chức năng đó.

3.6 Mật mã Hộp trắng (White-box Cryptography)

Nguồn gốc học thuật. Khái niệm mật mã hộp trắng được giới thiệu bởi Chow et al. năm 2002 [3], xuất phát từ bài toán bảo vệ khóa bí mật khi thuật toán mật mã được thực thi trong

môi trường mà đối thủ có toàn quyền quan sát. Đây là hướng kết hợp chặt chẽ giữa lý thuyết mật mã và làm rối mã: mục tiêu là ẩn giấu khóa và quy trình mã hóa trong một “mê cung” các bảng tra cứu và mã hoán vị sao cho kẻ tấn công không thể trích xuất khóa dù có mã chương trình. Kỹ thuật white-box nổi tiếng nhất là triển khai AES hộp trắng của Chow (gồm các bảng lookup thay thế cho mỗi vòng AES). Nhiều cải tiến và tấn công đã được công bố trong hai thập kỷ qua, tuy nhiên white-box crypto vẫn được sử dụng trong một số hệ thống DRM và ứng dụng ví điện tử.

Nguyên lý hoạt động. Thay vì giữ khóa mật mã dưới dạng một mảng rõ trong chương trình, white-box sẽ trộn khóa vào logic thuật toán. Cách làm điển hình:

- Biến đổi mỗi bước của thuật toán mật mã (vốn là các phép toán đại số như XOR, S-box) thành một hoặc nhiều bảng tra cứu (lookup table) đã “điền sẵn” kết quả tính toán khi có khóa. Những bảng này bao gồm cả hiệu ứng của khóa, nên nếu nhìn vào bảng không dễ suy ra khóa.
- Thêm các bước hoán vị ngẫu nhiên giữa các vòng tính toán (so-called *random bijections*) để làm mất mối liên hệ trực tiếp giữa input-output bảng và khóa.
- Code thực thi chỉ là liên tục tra cứu các bảng này và kết hợp kết quả.

Nói cách khác, toàn bộ thuật toán (ví dụ AES) được “unroll” thành các hằng số (các bảng) và mã đan xen phức tạp. Người tấn công khi dịch ngược sẽ thấy hàng loạt bảng số dài và phép toán trên chỉ số, rất khó nhận ra logic AES hay khóa nằm ở đâu.

Ví dụ minh họa (Trước & Sau). Giả sử một bước đơn giản của thuật toán mã hóa sử dụng khóa:

Đoạn mã 11: Mã AES bình thường (trích đoạn một vòng)

```
1 uint8_t state = input_byte ^ roundKey;
2 state = SBox[state];
```

Trong triển khai white-box, hai bước “XOR với khóa” và “tra S-box” có thể gộp thành một bảng duy nhất:

Đoạn mã 12: Mã sau khi chuyển thành bảng tra cứu hộp trắng

```
1 uint8_t T0[256] = { /* 256 phần tử = SBox[i ^ k] với k là khóa */ };
2 uint8_t state = T0[input_byte];
```

Chú thích: Mảng T0 ở đây đóng vai trò như S-box “đã trộn khóa”: mỗi giá trị T0[i] thực chất bằng $S(i \oplus k)$, với k là một byte của khóa round. Khi mã chạy, thay vì XOR rồi tra S-box, nó chỉ việc lấy trực tiếp T0[input_byte]. Kẻ tấn công nhìn vào T0 sẽ rất khó đoán được k vì bảng trông như một hoán vị ngẫu nhiên 256 phần tử. Trong triển khai thật, có nhiều bảng phức tạp hơn và thêm các bước xáo trộn giữa các vòng.

Phân tích An toàn. Mật mã hộp trắng được thiết kế để chống lại *attacker* mạnh có quyền quan sát mọi thứ (thậm chí debug từng bước). Về lý thuyết, nếu được xây dựng hoàn hảo, attacker không thể tách được khóa bí mật từ các bảng. Tuy nhiên, thực tế phần lớn các scheme white-box công bố đều đã bị bẻ gãy bởi các phân tích nâng cao (chẳng hạn kỹ thuật gian lận bố cục - algebraic attack). Mặc dù vậy, đối với attacker thông thường, white-box vẫn gây khó khăn lớn: code rất cồng kềnh và khó dò, cần kỹ năng mật mã để phân tích. *Resilience* nhìn chung cao hơn hẳn so với các obfuscation thông thường (vì được hậu thuẫn bởi lý thuyết mật mã). *Stealth* lại thấp: mã có rất nhiều bảng và phép tính lạ, dung lượng lớn, dễ nhận diện là một thứ đã bị làm rối nặng.

Tác động Hiệu năng. *Cao.* Triển khai white-box thường cực kỳ nặng: ví dụ, AES hộp trắng kích thước code có thể hàng trăm KB so với AES tiêu chuẩn chỉ vài KB, và tốc độ chậm hơn hàng chục lần do tra cứu bảng liên tục thay vì tính toán trực tiếp. Điều này hạn chế phạm vi áp dụng: white-box chủ yếu dùng cho các thành phần đặc thù (như module cấp phép bản quyền) chứ không thể áp dụng tràn lan.

Lưu ý: Kỹ thuật white-box nhằm bảo vệ thuật toán mật mã và khóa bí mật, không tương tác với dữ liệu văn bản hay chuỗi ký tự trong chương trình. Vì vậy, nó không ảnh hưởng đến xử lý tiếng Việt. Nếu chương trình có sử dụng chuỗi tiếng Việt ở đâu đó, phần code white-box (vốn chỉ liên quan đến tính toán bit và byte khóa) sẽ không tác động đến chuỗi đó. Chỉ cần lưu ý rằng white-box làm chương trình nặng nề hơn, nhưng không tạo lỗi hay rủi ro gì liên quan đến encoding Unicode.

4 Ma trận Đánh giá và So sánh

Để đánh giá định lượng các kỹ thuật làm rối mã, chúng tôi sử dụng bốn tiêu chí chính:

- **Resilience:** Khả năng kỹ thuật giữ vững hiệu quả khi đối mặt với các công cụ giải rối tự động (mức chống chịu trước tấn công tự động hóa).
- **Stealth:** Mức độ kỹ thuật ẩn mình, không tạo ra những dấu hiệu bất thường dễ bị phát hiện (cả bởi người phân tích lẫn máy).
- **Performance Overhead:** Mức suy giảm hiệu năng (thời gian chạy, dung lượng bộ nhớ) mà kỹ thuật gây ra cho chương trình.
- **Chi phí triển khai (Implementation Cost):** Độ phức tạp khi hiện thực kỹ thuật trong phát triển phần mềm, bao gồm yêu cầu về thời gian, tài nguyên và ảnh hưởng tới quy trình (ví dụ: khó khăn khi debug, build).

Bảng 1 dưới đây tổng hợp đánh giá của chúng tôi đối với ít nhất 8 kỹ thuật obfuscation tiêu biểu. Mỗi tiêu chí được chấm theo thang: *Cao*, *Trung bình*, *Thấp* hoặc *Không đáng kể*. Cột “*Lý giải*” cung cấp diễn giải ngắn gọn cho điểm số, dựa trên phân tích kỹ thuật đã trình bày ở các phần trước.

Kỹ thuật	Resilience	Stealth	Overhead	Chi phí	
Đổi tên định danh	<i>Thấp</i>	<i>Cao</i>	<i>Không</i>	<i>Thấp</i>	
Chèn điều kiện mờ (Opaque predicates)	<i>Trung bình</i>	<i>Thấp</i>	<i>Thấp</i>	<i>Thấp</i>	
Làm phẳng luồng điều khiển	<i>Cao</i>	<i>Thấp</i>	<i>Trung bình</i>	<i>Trung bình</i>	Biến đổi CFG t
Máy ảo hóa mã	<i>Cao</i>	<i>Thấp</i>	<i>Cao</i>	<i>Cao</i>	Rất khó phân tích v
Mã hóa chuỗi ký tự	<i>Trung bình</i>	<i>Cao</i>	<i>Thấp</i>	<i>Thấp</i>	H
Chống gỡ lỗi	<i>Thấp</i>	<i>Cao</i>	<i>Không</i>	<i>Thấp</i>	Chỉ cản trở đượ
Mật mã hộp trắng	<i>Trung bình</i>	<i>Thấp</i>	<i>Cao</i>	<i>Cao</i>	
Chèn mã vô hại	<i>Thấp</i>	<i>Trung bình</i>	<i>Thấp</i>	<i>Thấp</i>	Mã vô dụng có thể bị các tối ưu

5 Công cụ và Hệ sinh thái

Trong thực tiễn, có nhiều công cụ phần mềm hỗ trợ tự động hoá việc làm rối mã. Dưới đây là một số công cụ tiêu biểu và đặc điểm của chúng:

ProGuard/R8 (Java/Android): ProGuard là trình thu nhỏ và làm rối mã mã nguồn mở dành cho Java và Android, hiện đã được thay thế phần lớn bởi R8 tích hợp trong Android build pipeline. Công cụ này tập trung vào các kỹ thuật mức độ nhẹ như đổi tên lớp/phương thức, loại bỏ mã không dùng, tối ưu hoá bytecode và loại bỏ thông tin debug. ProGuard *không hỗ trợ* các kỹ thuật phức tạp như mã hóa chuỗi hay biến đổi control-flow. Ưu điểm là tốc độ xử lý nhanh và giữ cho kích thước file APK nhỏ, đồng thời dễ tích hợp vào Gradle. Nhược điểm là mức độ bảo vệ hạn chế – kẻ tấn công vẫn có thể dễ dàng dịch ngược logic ứng dụng Android đã qua ProGuard, chỉ gặp khó khăn nhỏ do tên bị làm rối.

Obfuscator-LLVM (OLLVM, C/C++): Đây là một nhánh của trình biên dịch LLVM do Junod et al. phát triển [6], bổ sung các bước biến đổi obfuscation ở mức mã trung gian LLVM. OLLVM là mã nguồn mở và được cộng đồng duy trì qua nhiều fork. Các kỹ thuật chính được hỗ trợ gồm: *Control-Flow Flattening* (làm phẳng luồng điều khiển), *Bogus Control Flow* (chèn nhánh và mã giả), và *Instruction Substitution* (thay thế một câu lệnh bằng chuỗi câu lệnh tương đương phức tạp hơn). OLLVM làm việc ở cấp độ compiler nên tích hợp vào quy trình build C/C++ khá thuận tiện (chỉ thay trình `clang` bằng phiên bản đã chỉnh sửa). Nhờ tích hợp sâu, OLLVM có thể bảo vệ cả mã native ở mức độ thấp. Tuy nhiên, do là mã nguồn mở, các mẫu obfuscation của OLLVM cũng được giới nghiên cứu thu thập và nhận diện, vì vậy cần kết hợp thêm các lớp bảo vệ khác nếu đối thủ đã biết trước việc sử dụng OLLVM.

Tigress (C, nguồn-nguồn): Tigress là một công cụ làm rối mã dạng nguồn sang nguồn dành cho ngôn ngữ C, được phát triển trong giới học thuật (Đại học Arizona) và cung cấp miễn phí cho mục đích nghiên cứu. Khác với OLLVM, Tigress vận hành ở mức mã nguồn C: nó nhận mã C và sinh ra mã C đã được làm rối. Tigress hỗ trợ rất nhiều kỹ thuật nâng cao, từ các biến đổi đã nêu như control-flow flattening, opaque predicates đến *Virtualization* (chuyển một hàm thành bytecode máy ảo và bộ thông dịch tương ứng), *Dynamic Code Generation* (tạo mã tự biến đổi khi chạy), v.v. Nhờ sự đa dạng này, Tigress thường được dùng để tạo bộ mẫu malware hoặc benchmark cho nghiên cứu giải rối mã. Nhược điểm của Tigress là mã đầu ra rất khó đọc và hầu như không thể duy trì thủ công; nó phù hợp cho việc tạo binary bảo vệ hơn là tiếp tục phát triển trên mã nguồn đã làm rối.

JavaScript Obfuscator (JS): Trong lĩnh vực web, một công cụ nổi bật là thư viện `javascript-obfuscator` (có trên npm) cho phép làm rối mã JavaScript. Đây là công cụ mã nguồn mở, được sử dụng rộng rãi để bảo vệ mã nguồn front-end. Nó áp dụng nhiều kỹ thuật: đổi tên biến thành tên rất dài hoặc ký tự Unicode khó phân biệt, loại bỏ khoảng trắng và comment, mã hóa chuỗi (chuỗi được chuyển thành các biểu thức toán học hoặc mã ASCII), flatten một số cấu trúc điều khiển, chèn code vô hại ngẫu nhiên, và thậm chí chèn các cơ chế runtime anti-debug (ví dụ: phát hiện `debugger` statement). Do JavaScript là ngôn ngữ thông dịch động, các công cụ như `javascript-obfuscator` phải đảm bảo mã sau khi làm rối vẫn chạy được trên mọi trình duyệt, do đó chúng thường tránh các kỹ thuật quá phức tạp. Mức độ bảo vệ của obfuscator JS ở mức vừa phải: đủ làm nản lòng người xem mã bình thường, nhưng hacker có kinh nghiệm có thể dùng công cụ như *deobfuscator* hoặc trực tiếp phân tích động trên trình duyệt để rút gọn mã.

Ngoài các công cụ trên, còn nhiều tên tuổi khác trong hệ sinh thái làm rối mã: chẳng hạn, cho .NET có Dotfuscator (thương mại) và ConfuserEx (mã nguồn mở); cho iOS có các plugin obfuscator Swift/Objective-C; cho nhúng (firmware) có các packer như Themida, VMProtect (thường dùng kèm trong bảo vệ game/DRM). Mỗi công cụ thường tập trung vào một nhóm ngôn ngữ hoặc nền tảng nhất định, và không ngừng phát triển để đối phó với các kỹ thuật tấn công mới.

6 Các phương pháp Tấn công và Giải rối mã (Deobfuscation)

Mặc cho những nỗ lực bảo vệ, các kỹ thuật obfuscation đều có thể bị tấn công ở mức độ nào đó. Dưới đây là các phương pháp chính mà chuyên gia phân tích sử dụng để "giải rối" chương trình, cùng ví dụ tiêu biểu:

Phân tích tĩnh: Sử dụng các công cụ đọc và phân tích mã máy (disassembler, decompiler) để cố gắng tái tạo lại cấu trúc chương trình gốc mà không thực thi nó. Những kỹ thuật làm rối như đổi tên, chèn mã vô hại thường bị vô hiệu hóa dễ dàng bởi phân tích tĩnh – ví dụ, trình decompiler vẫn sẽ hiển thị cấu trúc điều khiển if/else rõ ràng dù tên biến là vô nghĩa. Một số nghiên cứu đã phát triển công cụ giải rối tự động, như Udupa et al. (2005) đề xuất thuật toán loại bỏ control-flow phức tạp [9]. Phân tích tĩnh cũng bao gồm việc áp dụng các tối ưu hoá ngược (vd: loại bỏ mã chết, inline hàm) để đơn giản hoá chương trình đã bị làm rối.

Phân tích động: Kê tấn công chạy chương trình trong môi trường kiểm soát (ví dụ: sandbox, máy ảo) và quan sát hành vi thực tế. Cách này đặc biệt hiệu quả để vượt qua các kỹ thuật ẩn dữ liệu: chẳng hạn, chuỗi ký tự mã hoá cuối cùng sẽ phải giải mã để sử dụng, do đó kẻ tấn công có thể đặt breakpoint ngay sau hàm giải mã để lấy được chuỗi gốc. Tương tự, các cơ chế anti-debug có thể bị vô hiệu bằng cách vá bỏ đoạn code kiểm tra debugger hoặc dùng công cụ như Frida để can thiệp API. Phân tích động thường được hỗ trợ bởi các framework như *Pin*, *DynamoRIO* (dành cho binary) hoặc công cụ instrument runtime (đối với bytecode) nhằm thu thập các thông tin luồng chạy, bảng jump, giá trị biến, ... vượt qua lớp rối che phủ bên ngoài.

Thực thi biểu tượng (Symbolic Execution): Đây là kỹ thuật rất mạnh, dùng công cụ như angr [8] để mô phỏng thực thi chương trình với đầu vào là các biểu thức thay vì giá trị cụ thể. Bằng cách cho phép duyệt nhiều nhánh một lúc, symbolic execution có thể tự động tìm các điều kiện đầu vào kích hoạt những phần mã nhất định. Điều này giúp giải các *opaque predicate* (bằng cách chứng minh biểu thức luôn đúng/sai) và bóc tách được luồng điều khiển thật khỏi luồng giả. Nghiên cứu của Banescu et al. cho thấy ngay cả kỹ thuật mạnh như ảo hóa cũng khó chống lại symbolic execution hoàn toàn [1]. Ví dụ, symbolic execution có thể mô hình hoá máy ảo và suy ra quan hệ giữa bytecode và kết quả, từ đó dần phục hồi logic gốc.

Phân tích dựa trên Machine Learning: Gần đây, các phương pháp học máy và học sâu bắt đầu được áp dụng nhằm nhận diện và khôi phục mã làm rối. Ý tưởng chung là sử dụng mạng neural để học đặc trưng của mã đã obfuscate và deobfuscate. Một ví dụ là huấn luyện mô hình dịch ngược (decompiler neuron) để chuyển mã máy đã rối về mã nguồn gần đúng. Ngoài ra, phân tích cấu trúc đồ thị chương trình bằng ML có thể phân loại được những vùng mã nào đã bị làm rối (ví dụ: phát hiện pattern flattening) và từ đó áp dụng biện pháp loại bỏ. Tuy các kỹ thuật này còn trong giai đoạn đầu, chúng hứa hẹn tự động hoá một phần quá trình giải rối mà trước đây phải làm thủ công.

Giải ảo (De-virtualization): Đây là hướng tấn công đặc thù nhắm vào kỹ thuật obfuscation bằng máy ảo. Bằng cách phân tích bytecode và trình thông dịch, chuyên gia có thể cố gắng xây dựng lại chương trình gốc. Một số cách tiếp cận bao gồm: khôi phục bảng chuyển trạng thái của VM, xác định từng opcode ảo tương ứng với chức năng gì (thông qua thực nghiệm), sau đó viết một bộ giải mã để biến code ảo thành code gốc. Công việc này rất phức tạp, nhưng đã có những thành công nhất định trong cộng đồng reverse (ví dụ: giải mã bytecode VMProtect trong các game bị bẻ khóa). Mặt khác, công cụ symbolic execution cũng tỏ ra hữu ích trong de-virtualization khi nó có thể suy luận qua lớp ảo. Chẳng hạn, thay vì phải hiểu từng lệnh VM, symbolic execution coi toàn bộ VM như một hộp đen toán học và trực tiếp liên kết đầu ra chương trình với đầu vào, từ đó loại bỏ sự cần thiết của việc hiểu chi tiết VM.

7 Ứng dụng Thực tế và Các Vấn đề Liên quan

7.1 Bảo vệ Ứng dụng Di động

Vấn nạn phổ biến trên Android và iOS là ứng dụng hợp pháp bị *repackaging* (giải nén, chỉnh sửa và đóng gói lại) hoặc bị ăn cắp mã nguồn/thuật toán bởi đối thủ cạnh tranh. Do đó, các nhà phát triển di động hầu như luôn áp dụng một mức độ obfuscation nhất định. Trên Android, công cụ ProGuard/R8 thường được sử dụng mặc định để thu nhỏ và làm rối mã APK trước khi phát hành. Đối với các ứng dụng nhạy cảm (như ứng dụng ngân hàng, ví điện tử, game online), nhiều lớp bảo vệ nâng cao hơn được triển khai: mã hóa chuỗi để giấu khóa API, kiểm tra tính toàn vẹn để phát hiện app bị sửa đổi, anti-debug để ngăn chặn kẻ tấn công hook khi app chạy, thậm chí có SDK bảo mật chuyên dụng (VD: Guardsquare DexGuard) cung cấp thêm obfuscation nâng cao trên Android. Trên iOS, do đặc thù mã native, các kỹ thuật như OLLVM và các thủ thuật anti-tamper (phát hiện jailbreak, debugger) được tích hợp thẳng vào app. Kết quả là những ứng dụng di động hiện đại, đặc biệt trong lĩnh vực tài chính và nội dung số, trở thành “hộp đen” khó phân tích. Điều này bảo vệ nhà phát triển khỏi việc bị sao chép bất hợp pháp và bảo vệ người dùng (tránh app giả mạo), nhưng cũng gây không ít khó khăn cho các nhà nghiên cứu bảo mật khi cần kiểm tra độc lập ứng dụng.

7.2 Bảo vệ Quyền nội dung số (DRM)

Digital Rights Management (DRM) là lĩnh vực áp dụng obfuscation rất mạnh mẽ. Mục tiêu của DRM là ngăn người dùng hoặc hacker bẻ khóa cơ chế bảo vệ bản quyền của nội dung số (nhạc, phim, phần mềm) để sao chép hay phân phối lậu. Các phần mềm phát nội dung (như trình phát Blu-ray, ứng dụng streaming Netflix/Spotify, game PC) thường tích hợp nhiều lớp bảo vệ:

- Sử dụng **white-box cryptography** để bảo vệ khóa giải mã nội dung: ví dụ, khóa AES dùng để giải mã phim được nhúng theo kiểu hộp trắng, rất khó trích xuất trực tiếp.
- Áp dụng **ảo hóa và làm rối luồng điều khiển** cho các đoạn mã kiểm tra bản quyền: điển hình như Denuvo (cơ chế chống crack game PC) bao gồm những đoạn code ảo hoá nhiều tầng, đòi hỏi hacker phải thực thi hàng triệu lệnh giả mới tới được logic game thật, làm trì hoãn việc crack.
- **Anti-tamper, Anti-debug**: Nếu phát hiện người dùng chỉnh sửa binary (patch) hoặc chạy trong môi trường debugger, chương trình DRM có thể dừng hoạt động hoặc xoá nội dung giải mã khỏi bộ nhớ.

Nhờ những biện pháp này, các hãng nội dung số duy trì được mô hình kinh doanh (bản quyền theo thuê bao, bán lẻ) trước nạn vi phạm bản quyền. Tuy nhiên, mặt trái là đôi khi các biện pháp DRM làm giảm trải nghiệm người dùng hợp pháp (ứng dụng nặng nề, ngốn CPU) và gây tranh cãi về quyền sở hữu của người dùng đối với sản phẩm đã mua. Việc bẻ khóa DRM thường được thực hiện bởi các nhóm hacker có tổ chức, đòi hỏi trình độ cao và thời gian hàng tháng phân tích – nhưng lịch sử cho thấy không có hệ thống DRM nào là bất khả xâm phạm: sớm hay muộn các khóa bảo vệ vẫn bị crack, buộc nhà sản xuất phải cập nhật kỹ thuật mới liên tục.

7.3 Mã độc và Phần mềm gian lận

Các tác giả phần mềm độc hại (*malware*) và công cụ gian lận (cheat, bot) là những “khách hàng” nhiệt thành của công nghệ obfuscation. Mục tiêu của họ là trốn tránh sự phát hiện của chương trình diệt virus và gây khó khăn cho chuyên gia phân tích malware. Hầu hết malware hiện đại khi lan truyền đều được đóng gói (pack) hoặc mã hoá để che giấu nội dung thực sự. Ví

dụ, một virus Trojan có thể được bọc trong nhiều lớp packer (UPX, MPRESS, v.v.), khi thực thi mới giải nén payload cuối. Mã độc còn tận dụng kỹ thuật *polymorphic* và *metamorphic*: mỗi biến thể phát tán ra sẽ tự thay đổi một phần mã, sử dụng các chuỗi lệnh tương đương khác nhau, nhằm tránh bị nhận diện mẫu. Đối với mã độc trên Android, thường thấy việc mã hóa chuỗi (để ẩn URL máy chủ điều khiển) và làm rối control-flow (để gây khó cho việc phân tích hành vi). Bên cạnh malware, các phần mềm gian lận (như hack game, tool auto) cũng dùng obfuscation để trốn tránh kiểm duyệt của nhà phát triển game. Họ có thể chèn anti-debug để ngăn game phát hiện quá trình can thiệp, mã hoá các tương tác với game để giấu hành vi bất thường. Tất cả những điều này khiến cuộc chiến giữa bên phòng thủ (antivirus, anti-cheat) và bên tấn công liên tục leo thang: khi kỹ thuật obfuscation mới xuất hiện, phía phòng thủ lại phát triển giải pháp giải rối hay nhận diện tương ứng, và vòng xoáy tiếp diễn.

7.4 Rủi ro Pháp lý và Đạo đức

Việc sử dụng các kỹ thuật làm rối mã cũng đặt ra một số vấn đề pháp lý và đạo đức. Về pháp lý, nhiều quốc gia (như Hoa Kỳ với DMCA) có luật cấm hành vi *vượt qua các biện pháp bảo vệ kỹ thuật* – trong đó, obfuscation được xem như một biện pháp bảo vệ bản quyền. Điều này có nghĩa là người dùng hoặc nhà nghiên cứu nếu cố gắng giải mã, dịch ngược phần mềm đã bị làm rối có thể vô tình vi phạm pháp luật, trừ các ngoại lệ được cho phép (ví dụ: nghiên cứu bảo mật, tương thích). Mặt khác, về phía nhà phát triển, cần lưu ý không vi phạm quy định khi sử dụng obfuscation: chẳng hạn một số giấy phép phần mềm mã nguồn mở (GPL) có thể coi việc phân phối chương trình ở dạng “đã làm rối khó hiểu” là không phù hợp với yêu cầu cung cấp mã nguồn.

Về đạo đức và trải nghiệm người dùng, lạm dụng obfuscation có thể dẫn đến thiếu minh bạch. Phần mềm thương mại đóng gói kín vốn đã khó kiểm tra, nay còn obfuscation mạnh sẽ càng gây nghi ngờ về những gì ẩn bên trong (liệu có hành vi thu thập dữ liệu người dùng không được công bố?). Đối với cộng đồng nghiên cứu bảo mật, obfuscation tạo ra “vùng cấm” cản trở việc audit độc lập các phần mềm phổ biến – ví dụ, rất khó để phân tích một ứng dụng chat mã hoá đầu cuối xem có lỗ hổng hay cửa hậu không, nếu ứng dụng đó bị làm rối nặng. Từ góc nhìn doanh nghiệp, obfuscation quá mức cũng gây khó khăn cho chính đội ngũ phát triển nội bộ trong việc bảo trì, gỡ lỗi sản phẩm (trừ khi duy trì bản build không obfuscation để debug). Vì vậy, các chuyên gia khuyến cáo chỉ nên áp dụng obfuscation ở mức cần thiết, kết hợp với các biện pháp bảo mật khác (ký mã, kiểm thử bảo mật) thay vì coi obfuscation là giải pháp “một nút bấm” thay thế mọi biện pháp an ninh. Dù sao, obfuscation bản chất là “bảo mật nhờ che giấu” (security through obscurity); nếu lớp che giấu bị xuyên thủng, hệ thống sẽ lộ nguyên trạng. Do đó tính minh bạch và cân bằng lợi ích – rủi ro cần được xem xét khi đưa obfuscation vào một dự án phần mềm.

8 Hướng dẫn Tích hợp vào Quy trình Kỹ thuật

Áp dụng obfuscation trong một dự án phần mềm đòi hỏi lên kế hoạch cẩn thận, đặc biệt khi đưa vào quy trình tích hợp liên tục (CI/CD). Dưới đây là một danh sách kiểm tra (checklist) những thực hành tốt giúp việc tích hợp obfuscation hiệu quả và an toàn:

1. **Tách biệt cấu hình debug và release:** Đảm bảo rằng phiên bản phần mềm dùng cho nội bộ (debug build) không áp dụng obfuscation, để thuận tiện cho việc gỡ lỗi. Chỉ bật obfuscation ở bước build *release* cuối cùng. Điều này giúp đội phát triển tránh phải làm việc trên mã đã bị làm rối, đồng thời giữ cho quá trình phát triển chính xác.
2. **Xây dựng deterministic (quyết định):** Nếu công cụ obfuscator có sử dụng ngẫu nhiên (ví dụ: sinh tên biến ngẫu nhiên, xáo trộn khác nhau mỗi lần), hãy cố gắng cố định seed hoặc cấu hình để build ra kết quả giống nhau giữa các lần chạy cùng một phiên bản mã

nguồn. Build có tính quyết định cao giúp việc so sánh phiên bản (diff) và debug sự cố dễ dàng hơn, đồng thời tránh trường hợp mỗi lần phát hành tạo ra một binary hoàn toàn khác (gây khó khăn cho kiểm thử).

3. **Quản lý file ánh xạ (mapping):** Nhiều obfuscator sinh ra file ánh xạ giữa tên/ cấu trúc gốc và sau khi làm rối (ví dụ: ProGuard tạo file `mapping.txt` ánh xạ tên lớp và phương thức). Các file này rất quan trọng cho việc phân tích crash (stacktrace) và debug sự cố sau khi phát hành. Cần lưu trữ chúng một cách an toàn (ví dụ: trong kho mật, chỉ đội ngũ đáng tin cậy truy cập) vì nếu file mapping bị rò rỉ, kẻ tấn công có thể dùng nó để dễ dàng đảo ngược chương trình.
4. **Kiểm thử hiệu năng tự động:** Tích hợp các bài kiểm thử hiệu năng trong pipeline CI/CD cho bản build đã obfuscate. So sánh thời gian chạy, mức sử dụng CPU, bộ nhớ giữa phiên bản thường và phiên bản làm rối. Nếu phát hiện obfuscation làm chậm hệ thống quá mức chấp nhận (ví dụ: >10% đối với tác vụ quan trọng), cần điều chỉnh – có thể tắt bớt một số biến đổi nặng (như ảo hóa) hoặc tối ưu lại code. Kiểm thử định kỳ giúp đảm bảo việc bảo vệ không đánh đổi quá nhiều hiệu năng.
5. **Kiểm thử hồi quy (regression test):** Luôn chạy đầy đủ bộ test chức năng trên binary sau khi obfuscation. Obfuscation đúng về lý thuyết không làm thay đổi hành vi, nhưng thực tế đôi khi công cụ có lỗi hoặc một số cấu trúc đặc thù có thể bị thay đổi logic (đặc biệt khi kết hợp tối ưu và obfuscate). Một bộ test tự động chạy trên bản obfuscate sẽ đảm bảo không có chức năng nào bị hỏng ngầm.
6. **Giám sát và cập nhật:** Xem obfuscation như một thành phần sống của quy trình bảo mật. Theo dõi các báo cáo mới về lỗ hổng hay khả năng giải rối của kỹ thuật đang dùng. Ví dụ, nếu có nghiên cứu mới chỉ ra phương pháp giải phẳng control-flow flattening hiệu quả, đội ngũ có thể xem xét bổ sung lớp bảo vệ khác. Đồng thời, cập nhật phiên bản công cụ obfuscator để vá các lỗi (nếu có) và cải thiện thuật toán.

Tuân thủ các thực hành trên giúp cân bằng giữa bảo mật và ổn định của phần mềm. Obfuscation nên được coi là một phần của chiến lược bảo mật tổng thể, phối hợp với các quy trình phát triển an toàn (Security Development Lifecycle) chứ không thay thế chúng.

9 Hướng Nghiên cứu Tương lai

Mặc dù kỹ thuật làm rối mã đã có lịch sử hàng chục năm, nhiều thách thức và hướng đi mới vẫn đang mở ra:

Obfuscation cho AI/ML: Sự nổi lên của các mô hình học máy phức tạp (như mạng nơ-ron sâu) đặt ra câu hỏi: làm sao bảo vệ được tài sản trí tuệ bên trong mô hình (kiến trúc và trọng số) khi triển khai trên thiết bị người dùng? Các mô hình on-device có thể bị trích xuất (model extraction) hoặc bị đảo ngược để sao chép. Nghiên cứu gần đây của Zhou et al. đã đề xuất kỹ thuật “Model Obfuscation” cho mạng neural, bao gồm việc đổi tên lớp, ẩn tham số và tiêm các lớp giả để làm mô hình khó hiểu hơn [10]. Đây là lĩnh vực mới mẻ, kết hợp giữa bảo mật phần mềm và học máy. Trong tương lai, chúng ta có thể thấy các framework hỗ trợ làm rối trực tiếp mô hình (ví dụ: obfuscate mạng TensorFlow) trước khi phát hành, tương tự như cách ProGuard làm với mã Java.

Chống phân tích dựa trên phần cứng: Khi các biện pháp bảo vệ phần mềm ngày càng mạnh, đối thủ có thể chuyển sang khai thác cấp phần cứng. Ví dụ, họ dùng thiết bị chuyên dụng để đọc trộm nội dung bộ nhớ, đo đặc thời gian hoặc năng lượng tiêu thụ để suy ra hoạt động nội bộ (tấn công kênh bên). Obfuscation truyền thống hầu như không chống lại được các kiểu tấn công này. Do đó, một hướng nghiên cứu là *kết hợp obfuscation với phần cứng bảo mật*: sử dụng enclave (như Intel SGX) để chạy những đoạn code quan trọng trong môi trường mà

ngay cả debug phần cứng cũng không thấy rõ, hoặc áp dụng kỹ thuật che giấu lưu lượng bộ nhớ (oblivious RAM) để tránh lộ mẫu truy cập. Ngoài ra, việc xáo trộn trình tự lệnh ở mức vi kiến trúc (micro-op obfuscation) cũng là hướng đang được xem xét để gây nhiễu cho các công cụ phân tích ở cấp chip.

Indistinguishability Obfuscation (iO) khả thi: Về mặt lý thuyết, iO được xem là “Chén Thánh” của obfuscation – đảm bảo chương trình bị làm rối an toàn tuyệt đối trừ khi $P = NP$. Sau nhiều năm tưởng chừng bất khả thi [2], các nhà mật mã đã đạt được tiến triển trong việc xây dựng iO (dựa trên các giả định hàm học đa tuyến tính). Tuy các đề xuất hiện tại còn phi thực tế (kích thước mã tăng cực lớn và rất chậm), nhưng trong tương lai, nếu iO khả thi được tối ưu hoá, ta có thể có những obfuscator “một bước” biến mọi chương trình thành không thể đảo ngược với bằng chứng an toàn chặt chẽ. Đây là hướng nghiên cứu mang tính hàn lâm cao nhưng tiềm năng ứng dụng cách mạng cho bảo vệ phần mềm.

Tự động hoá và tối ưu hoá quá trình obfuscation: Hiện tại, việc chọn kỹ thuật nào, áp dụng ở đâu phần lớn dựa trên kinh nghiệm chuyên gia. Tương lai có thể chứng kiến các hệ thống tự động đề xuất chiến lược obfuscation. Chẳng hạn, dùng giải thuật di truyền hoặc học tăng cường để tìm tổ hợp biến đổi tối ưu (đạt độ bảo vệ cao mà overhead thấp) cho một ứng dụng cụ thể. Cũng có thể xuất hiện các công cụ “obfuscator compiler” thông minh hơn: tự phân tích code và quyết định chỗ nào nên áp dụng opaque predicate, chỗ nào nên ảo hóa, v.v. Điều này đòi hỏi định nghĩa được các hàm mục tiêu (như “độ khó đảo ngược” định lượng được) – một thách thức khoa học cần giải quyết.

Nhìn chung, cuộc đấu giữa kỹ thuật làm rối mã và kỹ thuật giải rối sẽ tiếp tục trong tương lai. Mỗi bước tiến của bên này sẽ kéo theo bước tiến của bên kia. Tuy nhiên, với sự kết hợp liên ngành (mật mã học, học máy, phần cứng) và nhận thức ngày càng cao về an toàn phần mềm, chúng ta có cơ sở để hy vọng vào những phương pháp làm rối mã hiệu quả và an toàn hơn, phục vụ cho việc bảo vệ tài sản số trong những thập kỷ tới.

10 Kết luận

Làm rối mã là một lĩnh vực kết hợp giữa nghệ thuật và khoa học, với mục tiêu chung là bảo vệ bí mật phần mềm bằng cách làm cho mã trở nên khó hiểu với kẻ tấn công. Trong tài liệu này, chúng tôi đã trình bày một bức tranh toàn cảnh về obfuscation: từ định nghĩa, phân loại các kỹ thuật, phân tích ưu nhược điểm từng phương pháp, so sánh định lượng, cho đến ứng dụng thực tế và cách tích hợp vào quy trình phát triển. Có thể thấy, không có kỹ thuật làm rối nào là “vạn năng” – mỗi phương pháp đều có điểm yếu có thể bị khai thác (như đã liệt kê trong bảng so sánh). Do đó, để đạt hiệu quả bảo vệ cao, người phát triển nên kết hợp nhiều kỹ thuật một cách hợp lý (đa tầng), đồng thời thường xuyên cập nhật trước các phương pháp tấn công mới.

Về mặt thực tiễn, obfuscation nên được áp dụng một cách có chọn lọc: tập trung vào những thành phần nhạy cảm (thuật toán lõi, khoá bí mật) thay vì rải đều trên toàn bộ mã, nhằm cân bằng giữa độ an toàn và hiệu năng. Các kỹ thuật nặng như ảo hóa hay hộp trắng chỉ nên dùng cho các module thật sự cần bảo vệ tối đa. Luôn duy trì quy trình kiểm thử chặt chẽ để đảm bảo chương trình sau khi làm rối vẫn hoạt động đúng và hiệu năng trong giới hạn cho phép. Đồng thời, tuân thủ các vấn đề pháp lý khi triển khai (đặc biệt trong bối cảnh DRM, bản quyền).

Tóm lại, làm rối mã là một công cụ mạnh trong hộp công cụ bảo vệ phần mềm, nhưng không phải là giải pháp thay thế cho các biện pháp bảo mật khác. Kết hợp obfuscation với thiết kế phần mềm an toàn, quản lý khoá chặt chẽ, và cơ chế bảo mật đa lớp sẽ tạo nên tuyến phòng thủ vững chắc. Chúng tôi hy vọng sách trắng này đã cung cấp cái nhìn sâu sắc và thực tiễn để độc giả hiểu và vận dụng hiệu quả kỹ thuật làm rối mã trong bối cảnh an toàn phần mềm hiện nay.

Tài liệu tham khảo

Tài liệu

- [1] Banescu, Sebastian, et al. *Code Obfuscation against Symbolic Execution Attacks*. Proceedings of the 2016 ACM Workshop on Software PROtection (SPRO'16), 2016.
- [2] Barak, Boaz, et al. *On the (Im)possibility of Obfuscation*. Advances in Cryptology – CRYPTO 2001, LNCS 2139, pp. 1–18, 2001.
- [3] Chow, Sherman S. M., et al. *White-Box Cryptography and an AES Implementation*. Proceedings of the 2002 ACM CCS-9 Workshop on DRM, pp. 1–15, 2002.
- [4] Collberg, Christian, Clark Thomborson, and Douglas Low. *A Taxonomy of Obfuscating Transformations*. Technical Report 148, Department of Computer Science, University of Auckland, 1997.
- [5] Collberg, Christian, and Jasvir Nagra. *Surreptitious Software: Obfuscation, Watermarking, and Tamperproofing for Software Protection*. Addison-Wesley Professional, 2010.
- [6] Junod, Pascal, et al. *Obfuscator-LLVM – Software Protection for the Masses*. Proceedings of the 1st ACM Workshop on Software PROtection (SPRO'15), pp. 15–21, 2015.
- [7] Low, Daniel. *Java Control Flow Obfuscation*. Master's thesis, University of Auckland, 1998.
- [8] Shoshitaishvili, Yan, et al. *angr: A Platform-Agnostic Binary Analysis Framework*. (Tool Demonstration), 2016.
- [9] Udupa, Srinivas, Saumya Debray, and Matias Madou. *Deobfuscation: Reverse Engineering Obfuscated Code*. Proceedings of the 12th Working Conference on Reverse Engineering (WCRE 2005), pp. 45–54, 2005.
- [10] Zhou, Mingyi, et al. *ModelObfuscator: Obfuscating Model Information to Protect Deployed ML-Based Systems*. Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2023), pp. 352–364, 2023.