

個人情報及び個人識別子を含むファイルと通信を検出するための双子の環境

張 世申¹ 三村 賢次郎¹ 新城 靖¹

受付日 xxxx年0月xx日, 採録日 xxxx年0月xx日

概要: インタネットユーザは普段ネットワークサービスを利用する時に, プライバシデータがサービス提供者に収集されることがよくある. 機械学習ブームの原因でユーザデータの重要性がより高くなる. サービスの提供者は自分のアプリケーションの体験を向上させるため, ユーザの個人データを大規模な計算や分析, 販売などを行うことがよくある. 通常のログインを通じて情報漏洩とは違い, ユーザがログインしなくてもユーザトラッキングの手段でユーザのプライバシーが漏れることもよく発生する. これらのユーザの身元を識別できるタグは, 広告やオークションのウェブサービスに利用され, ユーザの使用に悪影響をもたらすこともよく発生する. 本研究では, ユーザプライバシーの保護を目標として, ユーザトラッキングに使われる情報を検出, 保護したい. そこで本研究では Docker コンテナ技術を利用し, 双子の環境というユーザトラッキングを検出できるアプリケーションの仮想実行環境を提案した. 双子の環境を利用する双子のブラウザを実装し, ユーザトラッキングを検出した.

キーワード: ユーザトラッキング, HTTP, コンテナ, ブラウザ

Twins environment of detecting communication and file with privacy data or privacy identifier

ZHANG SHISHEN¹ KENJIRO MIMURA¹ YASUSHI SHINJO¹

Received: xx xx, xxxx, Accepted: xx xx, xxxx

Abstract: This document is a guide to prepare a draft for submitting to IPSJ Journal, and the final camera-ready manuscript of a paper to appear in IPSJ Journal, using L^AT_EX and special style files. Since this document itself is produced with the style files, it will help you to refer its source file which is distributed with the style files.

Keywords:

1. はじめに

PCで動作するアプリケーションは, Webブラウザのように明示的に通信を行うものだけでなく, オフィスツールのように, ユーザの意図しない通信を行うものがある. Webブラウザであっても, ユーザトラッキングのために暗黙的に通信を行うことがある. このような意図しない通信により住所・氏名等の個人情報, および cookie のように個人情報と結び付けられた個人識別子が送信されることが

ある.

本研究では, コンテナという OS 層の仮想実行環境を利用し, 個人情報及び個人識別子が含まれているファイルと通信を検出することを提案する [5]. そして, ファイルや通信から不要な情報を削除するツールを実装する. 個人情報および個人識別子を検出する対象は, 次の2つである.

- ファイル
- ネットワーク通信

たとえば, Webブラウザは, 個人識別子として訪問履歴や, フォームの内容, Cookieなどを保存する. ユーザトラッキングでよく利用される Cookieにはサーバが配る

¹ 筑波大学
University of Tsukuba

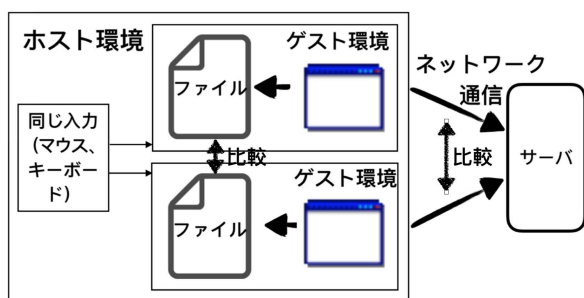


図1 双子の環境

ユニークな ID やセッション ID が含まれている。しかし、現在の Web ブラウザやオフィスツールは、非常に複雑であり、これらの識別子がどのファイルにどのような形式で保存されているかを調べることは容易ではない。本研究では、双子の環境を実装して、個人情報および個人識別子が含まれているファイルや通信を検出する。なお、以下では、個人情報と記載した時にも、個人情報と個人識別子の両方を含むものとする。

2. 双子の環境と双子のブラウザによる個人情報の検出

双子の環境とは、プログラムファイルやデータファイル等の内容がほとんど同じであるような2つの仮想実行環境である(図1)。人間における双子の研究では、異なる双子が異なる環境で育てられた際にそれぞれの医学的、遺伝子的、心理学的性格を調査することで、どのような違いが生まれるかを調査する場合が多い。本研究で提案する双子の環境では、類似の2つの環境で同一のプログラムをそれぞれ実行し、同一の入力を与える。そして、2つのプログラムの動作上の相違点を検出する。

本研究では、双子の環境で動作する Web ブラウザとして、双子のブラウザを開発している[4]。双子のブラウザとは、双子の環境で協調動作する2つのブラウザである。

本研究では、双子のブラウザを用いてサーバによるユーザトラッキングを検出する。ユーザトラッキングの手法としては、Cookie を使う方法や URL にタグを埋め込む方法がある。それ以外に、Flash Cookie や HTML5 の IndexedDB などのストレージを使う方法もある。本研究では双子のブラウザを用いてブラウザが作成する全てのファイルの差分およびブラウザが発信するネットワーク通信の内容の差分を調査する。その差分にユーザトラッキングのための情報が含まれる可能性が高い。その差分を削除、または修正することでユーザトラッキングを阻止することができると思われる。

3. コンテナによる双子の環境の実装

本研究では、環境内のファイルを調査するが、一般的な

仮想マシンではホストはゲスト OS のファイルシステムを直接的アクセスできないという問題がある。そこで本研究では、コンテナという OS 層の仮想マシンを使う。

コンテナは、仮想化技術の一種である。VMware や Xen, KVM などの一般的な仮想マシンとは違い、コンテナはハードウェア層ではなく、OS 層の仮想マシンである。コンテナは Linux カーネルの cgroups と namespaces 機能を利用し、リソース管理と隔離を提供している。

本研究ではコンテナを実装する仕組みとして Docker を使う。Docker では、Overlay File System というファイルシステムが利用可能である。このファイルシステムではゲスト OS のファイルをホスト OS からアクセスできるため、ファイルの差分を取得することが容易である。

ゲスト OS が生成したファイルは Overlay File System の Upper Layer に保存される。したがって、Upper directory のみ読み込むことでファイル内容の変化を得ることができる。

4. ファイルの差分検出

図2に、ファイルの差分検出の仕組みを示す。まず、同じイメージを利用する2つのコンテナを起動し、対象となるプログラム(主に双子のブラウザ)を実行する、ホストからの入力を2つ複製して、2つのコンテナを操作する。それからコンテナが生成したファイルをホスト OS で取得する。そして、ファイルの種類ごとにファイル内容に前処理を行い、その結果を diff コマンドに与える。

4.1 前処理とテキスト化

プログラムでよく利用される保存形式としてはテキスト形式、データベース形式、および、マーシャリング形式の3つがある。マーシャリング形式としてはよく JSON と XML がよく使われる。本研究では Web ブラウザ Firefox と Google Chrome のファイル保存形式を調査した。その結果、テキスト、JavaScript、JSON、XML、SQLite、BerkeleyDB、LevelDB の7種類あることがわかった。

コンテナが出力するファイルがテキストファイルであれば、そのまま diff コマンドで差分を取ることができる。しかしながら、JSON や XML では、そのまま diff コマンドに与えても大量の出力がなされ、目的とする個人情報が埋もれてしまうという問題がある。また、diff コマンドは、バイナリファイルを扱うことができない。

そこで本研究ではプログラムが生成したファイルに対して前処理を行い、diff で差分を取る前に識別しやすい形に変換する(図2)。形式分類モジュールで、ファイルを形式で分類する。関係データベースは dump ツールでテキストを生成し、自動増加の主キー列を消して、属性内容でソートする。JSON ファイルの標準化するために、json.tool を利用した。

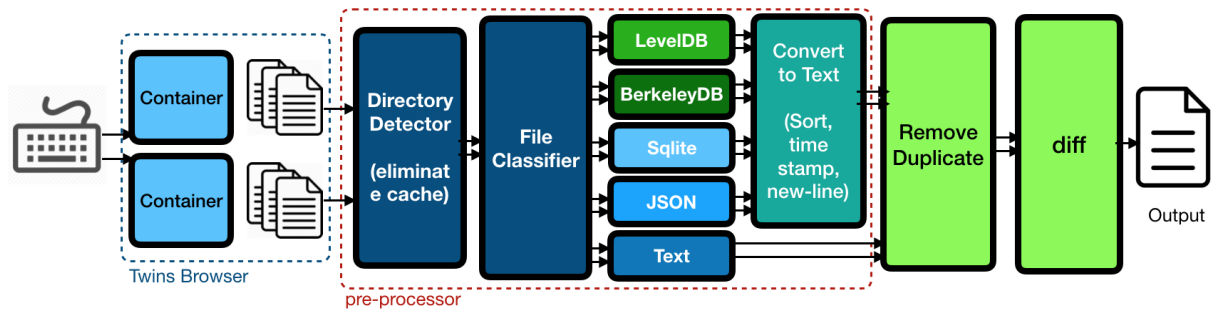


図 2 Docker Overlay Filesystem によるファイル差分検出

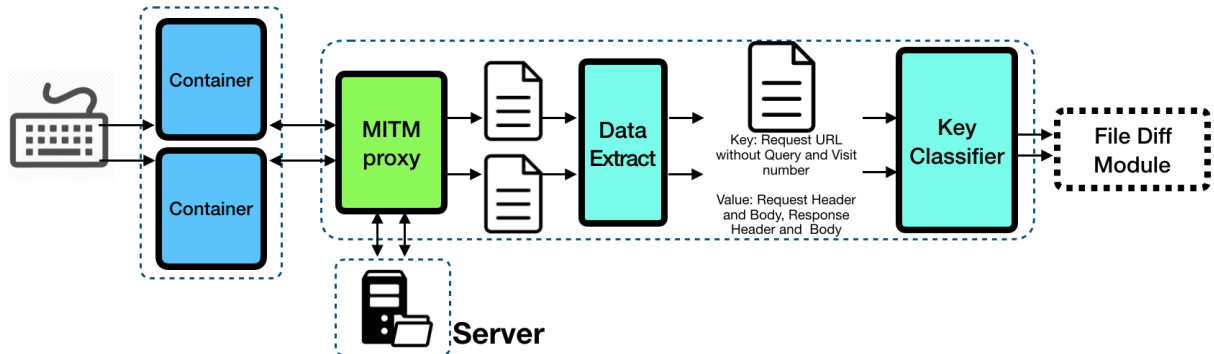


図 3 MITM-proxy を用いた HTTP のメッセージの差分検出

4.2 タイムスタンプの扱い

ネットワーク通信を行うプログラムは、様々なタイムスタンプをファイルに保存する。単純にファイルの差分を得ると、タイムスタンプの差によるものが大量に生成され、重要な差分が埋もれてしまう。

本研究では、タイムスタンプを次のように分類して扱う。

- ・ リモート： リモートの通信相手が指定したもの。例えば、HTTP の Last-Modified:ヘッダに由来するもの。
- ・ ローカル： ローカルの OS からシステム・コールで取得したものに由来するもの。

リモート・タイムスタンプは、外部に発信されることが想定されている。例えば、HTTP の応答メッセージに含まれた Last-Modified. ヘッダの値は、同じコンテンツを再取得する時に、要求メッセージの if-Modified-Since. ヘッダに含まれて発信される。この値は、個人識別子としてユーザトラッキングに使われることがあることが知られている。

一方、ローカル・タイムスタンプは、外部に発信されなければ、個人識別子にはなり得ない。したがって、双子の環境の実装では、ローカル・タイムスタンプの違いを排除したい。本研究では、双子の環境で実行したタイム関連のシステム・コール `gettimeofday()` と `clock_gettime()` をオーバーライドする。そのため、本研究で作成した動的リンクライブラリを LD_PRELOAD で置き換える。置き換えたシステム・コールは、環境変数で指定された固定の日付を返す。

このような処理を行ったとしても、アプリケーションの

内部で独自にシステム・コールで得たタイムスタンプを加工して利用していることがある。例えば、Firefox では、システム・コール `clock_gettime()` で得られたナノ秒単位の時刻に、1 から 4 まで加えた値を利用している。この問題を解決するために、本研究では、テキスト化した後、ローカル・タイムスタンプとそれを加工したものと思われる数字を定数で置き換える。

4.3 ランダム性の排除

プログラムは多く予想出来ない行動を行う。その結果、ファイルの内容に差が生まれ、本来検出したい個人情報を覆い隠してしまう。そこで、本研究では、そのようなランダムな行動を排除する。

まず、乱数によるプログラムのランダム行動を抑止するために、乱数デバイス `/dev/urandom` を置き換える。本研究では擬似乱数生成器デバイスドライバを作り、双子のコンテナのインスタンスに同じシードを与える。

マルチスレッドやマルチプロセスのアプリケーションのスケジュールも実行結果に影響を与える。Chrome ブラウザの Task モデルは、UI と IO スレッド以外にワカスレッドが多数存在する。一つの作業がどのワカスレッドにより実行されるか予想できない。その結果、Chrome ブラウザでの実験中には、TID のような予測できない結果が出力される。そして、PC のリアルタイムパフォーマンスによる動的にワカスレッドの数を調整することもある [3]。そこで本研究では、テキスト化してダンプする時にスレッド

識別子を含まないようにする。

5. ネットワークメッセージの差分

本研究では、まず、HTTPを対象としてネットワークメッセージをキャプチャする。コンテナ内で実行されるプログラムが送受信されているメッセージは、HTTPSにより暗号化されていることがある。そして、多数なHTML5やJavascriptの新機能もよく支持することが要求されている。そこで、本研究では、MITM-Proxy(Man-In-The-Middle-Proxy)*1を使ってHTTPやHTTPSの通信をキャプチャする。MITM-ProxyはHTTPS通信のキャプチャができ、Keep-Aliveなどの持続的接続機能もつける。

ネットワーク通信の差分取る仕組みを図3に示す。HTTP通信の内容をKey-Valueの形式に変える。Keyとしては現在RequestのURLをQuery Stringを排除した部分と訪問の回数をを用いている。Valueは残りのQuery String, Header, body, Response Header, bodyにする。Dumpした内容はHTTP通信のヘッダContent-TypeやContent-Encodingによる分類して、ふさわしい処理を行う。例えば”Content-Encoding: gzip”ヘッダの通信は圧縮した形を展開し, ”Transfer-Encoding: chunked”ヘッダの通信はchunkの属性によって、多数のChunkを組み合わせる。それから、全部のテキストファイル、例えばhtmlやjavascript, cssなどのファイルはテキストファイルとして4.1節で述べたファイル差分検出モジュールに与える。テキストファイルではない通信、例えばpng,faviconなどの通信はバイナリファイルとして扱う。

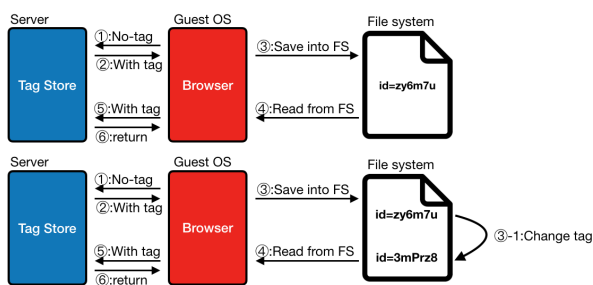


図4 ユーザトラッキングとタグプロテクション。上のイメージはブラウザが一回目と二回目の訪問する時に一部の個人識別子の移動ディレクション。一回目はタグがないからサーバが生成して、サーバのストレージとクライアント同時に保存する。下のイメージは個人識別子を保護する時に、一回目と二回目の訪問する間で個人識別子をなんかの手段で直す。

6. 差分内容の隔離と保護

トラッキングに使われる嫌疑があったタグを二回目サーバを訪問する時に、サーバが正しい認識できないような形に変える必要がある。図4のように、一回目の通信で生成

した新しいタグをそして、タグの隔離効果をテストするため、修正前の元のタグはホストに保存しなければならない。ファイルとネットワーク通信の差分の対応は違う。

6.1 ファイルの差分の置き換え

ファイルの差分内容は全部diffコマンドで見つかったから、出力ファイルは解析し、オリジナルファイルのふさわしい部分を処理する必要がある。ファイルの形によって処理手段は違うから、次のように扱う。

- ・テキストファイル: テキストファイルの対応は、まずdiffの出力ラインから差分タグを検出、差分タグを双子の環境の各ホスト環境のストレージに保存する。それからファイル名でゲストOSのファイルシステムに対応したファイルをawkで差分タグを直す。直したタグは直した前のタグと同じように保存する。タグの変更は、乱数を生成、元のタグと同じ形式に変わって、コンテナFSに書き込む。よく使われたタグの形式はBase64, Hash(MD5,SHA256), UUIDなどの形式がある。それらのタグをふさわしいランダムジェネレータで取って代わり、タグの保護ができる。
- ・データベース: テキストファイルと同じように差分タグを検出、保存する。タグの直すはawkではなく、データベースの管理ツールを利用して簡単で直せる。関係型データベースはSQLのupdateを利用し、No-SQLは直接にKey-Valueを直す。
- ・バイナリファイル: バイナリファイルの差分は直接的にdiffコマンドを使うのができないから、バイナリファイル中のテキストの部分を引き出して差分を取る。

6.2 ネットワーク通信の差分の置き換え

ネットワークの通信の差分タグは通信中にMITM-Proxyによるファイルシステムに保存する。保存の形は節のようにテキストファイルに保存して、diffコマンドで差分を取る。取り除いた差分タグと置き換えたタグはファイルの差分内容と同じようにホストのストレージに保存する。

6.3 ネットワーク通信の保護

ネットワーク通信で検出したタグは、全部コンテナファイルシステムに見つけれられるではない。その原因でネットワーク通信のタグを動的に検出し、置き換える必要がある。ここでは上記の手順でホストに保存した個人情報候補のタグを探し、HTTP-Requestをのちに関連したタグを通信中に直す。MITM-Proxyのたくさんのフィルタより簡単に実装できる。

7. 実験

7.1 双子のブラウザが生成したファイルの差分

本研究室ではFirefoxとGoogle Chromeで双子のブ

*1 <https://mitmproxy.org/>

表 1 差分がないファイルと差分があるファイルの数

	Firefox			Google Chrome		
	個人情報が ない	差分がある		個人情報が ない	差分がある	
		不明なバイナリ	テキストと扱 えるバイナリ (行 数)		不明なバイナリ	テキストと扱 えるバイナリ (行 数)
更新されたファイル	52	0	7(28)	98	5	47(114)
乱数生成器の置き換え	55	0	4(10)	128	5	17(46)
ローカルタイムスタンプの固定	55	0	4(8)	128	5	17(41)
前処理後	55	0	4(8)	136	5	9(15)

表 2 ファイルの差分の例

ファイル名	コンテナ 1	コンテナ 2
cookies.sqlite	value = 132=S1fjWe4xjUJJowuImYQyi...	value = 132=ablTT1GOqBIYlBX-MQ...
places.sqlite	guid = xhIxtPu6zAJ7	guid = IlfETnEs0Dr4
sessionstore.js	"docshellUUID": "{4c4508da-9468-430b-8eac-0484dcc43e5d}"	"docshellUUID": "{bfcff6ba-42c2-4731-a65c-ae1ba7b1cc0e}"
prefs.js	user_pref("browser.slowStartup.averageTime", 14971)	user_pref("browser.slowStartup.averageTime", 18103)

ブラウザを実装している．今回，ユーザトラッキングを行っている代表的な Web サイトとして Google を選択し <http://www.google.com/> を訪問する実験を行った．これは，ログインを必要としない簡単なサイトである．

7.1.1 Firefox

まず Firefox ブラウザを双子のブラウザとして双子の環境で実行した．コンテナの全てのファイルの変化をリストアップしたところ，全部で 59 個のファイルが作成された．その中に 29 個が /.cache にある一時的なファイルのためファイル名によるフィルタにより自動的に除外された．残る 30 個のファイルに対して前処理を行って，テキスト化し diff コマンドで差分を取った結果は全部 28 行であった．4 章で述べた方法に従ってテキスト化した差分を削減した結果，最終的に差分は 8 行だった．識別できないバイナリファイルはなかった．

表 3 に，発見した差分の一部を示す．cookies.sqlite には，HTTP Cookies が保存されていることがわかる．places.sqlite は，訪問履歴を保持するファイルである．その中に，アクセスした URL が保存されている．このように，URL の中に，ユーザトラッキングに利用可能なタグが埋め込まれていることがわかる．

この実験の結果，提案手法により，ファイルを特定し，使用者の識別子を検出することができ，ユーザトラッキングに使われる個人識別情報が含まれることが確認された．

7.1.2 Chrome

2 番目の実験は Google Chrome で行った．結果を表 1 に表す．差分があるファイルは Firefox より多い．識別できないバイナリファイルが 5 個残された．

7.2 ネットワーク通信内容の差分

ブラウザが行うネットワーク通信の内容の差分を調査した．訪問したサイトは，7.1 節と同じである．

全部の通信の数は 21 個あり，テキスト形式のファイル以外に png が 6 個，ico が 1 個あった．20 個のメッセージの内容に差分があり，1 番目の通信だけは差分がなかった．ネットワーク通信の差分はタイムスタンプと重複を除き全部で 22 行あった．応答メッセージの内容に差分があるファイルが 2 個あった．それらの URL は www.google.com と www.google.com/gen_204 であった．また，URL が違う通信の数は 3 個あった．

メッセージの内容には cookies が現れている．これは，cookies.sqlite(表 3) にも含まれている．

7.3 プライバシ保護

上記の実験で検出した個人情報データは，6 節で紹介した方法でタグの引き換え，及び削除でプライバシーを保護したい．保護する結果の評価は前と同じように双子の環境で動き，図 5 のように，双子のブラウザで訪問を二回行って，保護する結果は成功かどうかでテストする．もし二回目の Query の結果に一回目の Query 内容が含まれば，タグの保護がうまく行かなかったと表す．

7.3.1 評価方

実験でテストするサイトは Firefox ブラウザを利用，Google，Bing と Youtube のメインページでテストする．まず Prev-Key と Post-Key を決まる．それから双子のブラウザを利用し，初期化したブラウザアプリケーションで各サイトを訪問する．訪問した時に Prev-Key を Query

表 3 ネットワーク通信の差分の例

URL	コンテナ 1	コンテナ 2
www.google.com	Cookie:1P_JAR=2017-12-12-04,NID=119=BGA6eL4YT9Q...	Cookie:1P_JAR=2017-12-12-04,NID=rzmrga.L3s...
www.google.com/gen_204	ei:pdcsW5TSLIb28AXosLGYBA	ei:pdcsW-vCDMfS8QWq-I34Bg
ssl.gstatic.com/gb/imagess/i1_1967ca6a.png	age:253629	age:253628

表 4 トラッキング保護する結果

		一回目の結果 (コンテナ1)		一回目の結果 (コンテナ2)		二回目の結果 (既存環境)		二回目の結果 (テスト環境)	
Query key		Key "Prev"	Key "Post"	Key "Prev"	Key "Post"	Key "Prev"	Key "Post"	Key "Prev"	Key "Post"
テストサイト	www.google.com								
	www.bing.com								
	www.youtube.com								

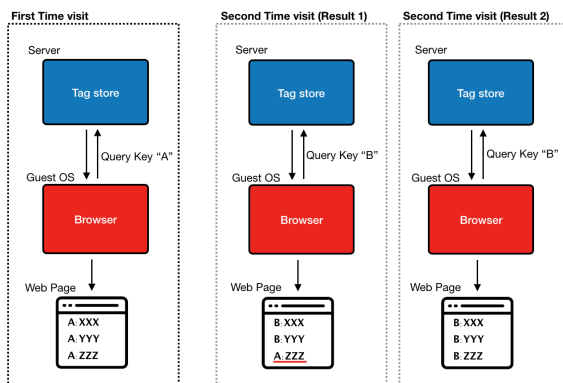


図 5 個人識別子を保護する結果の評価。一回目の訪問する時に、ブラウザがキーワード「A」を入力し、サーバに送信する結果は、すべて A 関連のデータである。二回目の訪問にはキーワード「B」を入力が、もし結果には「A」の結果が出れば個人識別子はよく保護していないとわかる

し、結果のウェブページに Prev-Key と Post-Key の出現回数を記録する。次にコンテナの実行を一時中止、2 つのコンテナに一つのコンテナの個人識別子を置き換える。置き換えたコンテナは、テスト環境となり、置き換えないそのまま実行した環境は既存環境とみなす。2 つのコンテナの実行結果のウェブページに Prev-Key と Post-Key の出現回数を比べて、保護の結果を評価する。それ以外に、取り除いた個人識別子が以後の通信に出現した場合でもある。これが ever-cookie[6] というトラッキングである。

7.3.2 評価

テストの結果は表 4 に示す。

未完成

8. 関連研究

Qubes-OS[1] は高いセキュリティを実現するための OS である。Qubes-OS は仮想計算機モニタ (Xen) を用いて、アプリケーションを隔離された仮想実行環境で実行する。

Qubes-OS には、複数の仮想実行環境のファイルを比較する機能はない。

Blink-Docker[2] はコンテナ中でブラウザを実行することで、Canvas Fingerprint を利用したユーザトラッキングからユーザを保護する。Canvas Fingerprint はハードウェアや OS のわずかな違いによってクライアントを特定する手法である。Blink-Docker はコンテナ技術を利用し、毎回 Fingerprint を変えることができる。本研究では、通信内容を検査し、そのようなユーザトラッキングを検出したいと考えている。

9. まとめ

本研究では人間の双子の研究に参考した双子の環境を提案した。本研究では Docker コンテナを利用し、軽量で類似の仮想環境を作り、ファイルと通信内容に含まれる個人情報を検出する。

本研究で開発された双子のブラウザを双子の環境で実行し、それらが生成するファイルの差分を調査するツールを実装した。また、Man-In-The-Middle Proxy を用いて、通信内容を比較するツールを実装した。実装した双子のブラウザを利用し、個人情報と個人識別子である可能性が高いユニークなタグがいくつ検出した。そして、個人識別子を変更の方で保護する。

未完成

参考文献

- [1] Qubes OS - A reasonably secure operating system: <https://www.qubes-os.org/>, accessed: 2017-11-22.
- [2] P Laperdrix, W Rudametkin, B Baudry: "Blink: A moving-target approach to fingerprint diversification" 37th IEEE Symposium on Security and Privacy, Poster Session, [Online] http://www.ieee-security.org/TC/SP2016/poster-abstracts/59-poster_abstract.pdf, (2016).

- [3] Javier Verd, Juan Jos Costa, Alex Pajuelo: "Dynamic Web worker pool management for highly parallel javascript web applications" *Concurrency and Computation: Practice and Experience*, 28(13):35253539, September 2015.
- [4] 張世申, 新城靖, 三村賢次郎: "個人情報を含むファイルと通信を検出ための双子の環境の提案", 情報処理学会第29回コンピュータシステムシンポジウム ポスターセッション, 2ページ (2017).
- [5] 三村賢次郎, 新城靖, 張世申: "Webサービスごとに隔離されたブラウジング環境の提案", 同上.
- [6] Gunes Acar, Christian Eubank, Steven Englehardt, Marc Juarez, Arvind Narayanan, Claudia Diaz, "The Web Never Forgets: Persistent Tracking Mechanisms in the Wild", *CCS '14 Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security* Pages 674-689.