

ビッグメモリアプリケーション向けの Multi-Variant 監視機構

清水 祐太郎^{†1,a)} 山田 浩史^{†1,b)}

概要：ソフトウェアの脆弱性を利用した攻撃は広く知られている。システムソフトウェアには未だに C/C++ などの型安全でない言語で記述されているものもあり、依然としてソフトウェアは攻撃の脅威に晒されている。近年の OS やコンパイラには、ASLR や SSP のような攻撃に対する緩和機構が実装されている。しかしながら、いずれの機構も情報の漏洩によって容易に回避できる可能性がある。多くの攻撃が特定のセキュリティ機構やそのバージョンに依存することに着目して、アプリケーションの複製を同時に実行する MVEEs と呼ばれる手法が存在する。それぞれの挙動を比較することで、既存の緩和機構を活かしながらセキュリティの向上が期待できる。その一方で、メモリを大量に利用するアプリケーションに対して MVEEs を適用する場合にはメモリの浪費が懸念される。本研究では、内容が一致するページを併合することでビッグメモリアプリケーションに対して効率的に MVEEs を適用する方法を提案する。Linux kernel 4.13.9 に提案手法の実装を行い、物理メモリの消費量を計測したところ、最大でアプリケーションを単体で動作させたのと同程度までメモリの消費を抑制することができた。

1. はじめに

ソフトウェアの脆弱性を利用した攻撃は広く知られている。様々な対策機構が提案、実用化されているにもかかわらず、依然としてソフトウェアは攻撃の脅威に晒されている。Java や C# のような型安全なプログラミング言語の登場によってこれらの問題は緩和されているものの、未だ多くのシステムソフトウェアは C/C++ などの安全でない言語で記述されているものも多い。そこで、近年のオペレーティングシステム (OS) やコンパイラには、脆弱性に対する攻撃を成功しにくくするためのセキュリティ機構が実装されている。Address Space Layout Randomization (ASLR) [1] ではアプリケーションの実行時に仮想アドレス空間にマッピングするバイナリやスタック、ヒープなどのアドレスをランダムにする。そうすることで攻撃の際に必要なデータが存在するアドレスを不明確にする。Stack Smashing Protector (SSP) [2] では、関数のスタックフレームの上位に Canary と呼ばれるランダム値を配置しておく。関数から返る際にこの値をチェックすることでスタック破壊を検知し、アプリケーションを停止させるこ

とができる。しかしながら、いずれの機構も情報の漏洩を伴うバグが存在している場合は、攻撃者は容易に機構を回避することが可能となってしまう。

そこで、多くの攻撃が特定のセキュリティ機構やそのバージョンに依存することに着目して、複数のセキュリティ機構を稼働させる Multi-Variant Execution Environments (MVEEs) [3-9] という手法が存在する。MVEEs では、同一アプリケーションを複数同時に動作させ、それぞれの動作を比較する。具体的には、アプリケーションの複製 (バリエーションと呼ぶ) を別々に実行し、それぞれが発行するシステムコールやプロセスの状態をモニタが監視する。挙動が異なるバリエーションを検知するとシステムは攻撃を受けていると認識し、全てのバリエーションを直ちに終了させることができる。MVEEs と極めて近い動作を行うが、全く同じアプリケーションではなく、バージョンが多少異なるアプリケーションを同時に動作させる N-Version [10, 11] と呼ばれる手法もある。これらの手法では、ASLR や SSP といった既存の緩和機構を活かしながら、よりセキュリティを向上させることができる。同一アプリケーションを実行しても、仮想アドレス空間や Canary 値は異なる。そのため、特定のバリエーションのみを狙った Exploit が他のバリエーションでも攻撃を成功させる確率は非常に低い。

しかしながら、Memcached [12] などの In-memory DB のようなメモリを大量に利用するアプリケーション (ビッグ

^{†1} 現在、東京農工大学
Presently with Tokyo University of Agriculture and Technology

a) simiyu@asg.cs.tuat.ac.jp

b) hiroshiy@asg.cs.tuat.ac.jp

グメモリアプリケーション) に対して MVEEs を適用する場合にはメモリの浪費が懸念される。バリエーション同士は互いに独立した仮想アドレス空間を持っている。そして、それぞれが別々の物理メモリを確保して利用している。たとえば、Amazon DynamoDB Accelerator (DAX) [13] では 488 GiB までメモリを搭載することが可能であり、こうしたアプリケーションに MVEEs を適用すると大量のメモリが必要となる。これにより、バリエーション数を増やすことが困難となる。また、このメモリの逼迫は、近年主流であるサーバ統合の集約率を低下させてしまう。

本研究では、ビッグメモリアプリケーションに対して効率的に MVEEs を適用する方法を提案する。提案機構では、全バリエーションに対して与えられる入力は等しいため、内部に記録されるデータも同じものになることに着目する。提案機構は、各バリエーションのメモリを走査し、同一内容のメモリ領域を検出、共有しながら実行を進める。これにより、ビッグメモリアプリケーションの複数稼働時の使用メモリ量を削減する。

本論文の貢献は次のとおりである。

- バリエーション間でメモリ内容が同一の領域を共有することで、物理メモリの浪費を抑える機構を提案した。提案機構は次の特徴を持つ。提案機構は、ビッグメモリアプリケーションを MVEEs に適用した際にも、物理メモリの使用量を抑えることができる。また、アプリケーションのソースコードを変更する必要はない。加えて、提案機構はマルチスレッドをサポートする。
- 提案機構を実現するためのソフトウェア機構を提示した。同一内容のメモリ領域を検出するために、カーネルレベルによるページ共有機構を利用する。これにより、各バリエーション間のアドレス空間における独立性を保ちながら、メモリ使用量を節約する。また、既存のセキュリティ機構である ASLR の適用時においてページの共有率を高めるために、ヒープオブジェクト内のポインタ値を変換する。
- Linux kernel 4.13.9 に実装を行いマイクロベンチマークを用いて物理メモリの消費量の計測を行った。実験により、最大でアプリケーションを単体で動作させたのと同程度までメモリの消費を抑制することができた。

本論文の構成は次のとおりである。第 2 章で本研究の背景と問題点について触れ、第 3 章でその解決策として本手法を提案する。第 4 章では提案手法を実現するためのシステム的设计、第 5 章では Linux カーネルおよびアプリケーションへの実装方法について述べる。第 6 章において提案手法の評価を行い、第 7 章で本研究のまとめと今後の課題を示す。

2. 背景

2.1 既存のセキュリティ機構

システムプログラムにおいて、メモリ管理やメモリへのアクセス方法のミスに起因する脆弱性が存在する。領域外参照や Use-After-Free [14]、未初期化領域の読み込みによるセンシティブなデータのリークなどが挙げられる。Java や C# のような型安全なプログラミング言語ではこれらの問題は緩和されてきた。しかし、未だ多くのプログラムは C/C++ などの安全でない言語で記述されているため、これらの問題は解決はされていない。そのため、これまでに脆弱性に対しての攻撃を通しにくくすることを目的とした緩和機構が開発されてきた。

2.1.1 Address Space Layout Randomization

Address Space Layout Randomization (ASLR) [1] とは、バイナリやスタック、ヒープなどをランダムな位置に配置するセキュリティ機構である。プロセスを実行する際、ローダはプログラム本体のバイナリのみならず、ライブラリやスタックなど実行に必要なデータを仮想アドレス空間に展開する。攻撃者はライブラリ内に存在しているコードやヒープやスタックに配置されたデータなど、既にメモリ内に存在している情報を利用して Exploit を構築することが多い。これらが無作為に配置することによって、攻撃のために必要なデータの位置を推測されにくくする。

2.1.2 Stack Smashing Protection

Stack Smashing Protection (SSP) [2] は、スタックオーバーフロー [15] による局所変数およびスタックフレームの破壊を緩和・検知する機構である。これは GCC [16] をはじめとするコンパイラによって実現される。局所変数においてオーバーフローを起こした場合、その上位に位置するデータは改竄されてしまう。通常スタックフレームの上位には、前の関数のスタックフレームのベースポインタ、リターンアドレス、引数が配置されている。SSP では、関数の引数や他の局所変数の改竄を難しくする。オーバーフローし得る配列よりも下位に局所変数を配置し、引数も同様に下位に移動して利用することで実現する。また、この機構ではベースポインタとリターンアドレスの改竄を検知することもできる。Canary と呼ばれるランダムな値を局所変数とベースポインタとの間に配置する。関数から返る際にこの値が変化しているかどうかチェックする。変化している場合は、スタックフレームが破壊され、上位のデータが改竄されたことを検出する。

2.2 Multi-Variant Execution Environments

MVEEs では一つのアプリケーションのレプリカとなるプロセスを作成し、それぞれの挙動を比較することでセキュリティの向上を図る。2.1 節で挙げた既存のセキュリ

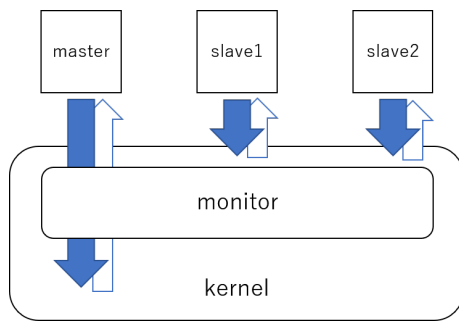


図 1 モニタとバリエント

ティ機構では、メモリアドレスや Canary といったセンシティブな情報が漏洩してしまうと、途端に攻撃に対する防御力を下げる。攻撃者は漏洩した情報を基に Exploit を組むためである。MVEEs では、各レプリカごとに仮想アドレス空間や Canary 値も異なっている。そのため、特定のレプリカのみを狙った攻撃がたとえ成功したとしても、他のレプリカに対しては失敗するため、攻撃の検知及び緩和が可能となる。

MVEEs は図 1 に示すように、モニタとレプリカ（バリエント）から成る。モニタはバリエント全体の動作の指揮を行う役割を担う。モニタが行うバリエントの監視はあくまで挙動の比較であり、メモリの内容や状態の完全一致を狙うものではない。そのため、ASLR や SSP のようなオペレーティングシステムやコンパイラに組み込まれているセキュリティ機構を活かすことができる。

動作の比較を行うために、バリエントは一定の間隔で同期を行う必要がある。一つの命令が実行されるごとにプロセスの状態を取得することも可能だが、これはあまりにもオーバーヘッドが大きいため実用的ではない。最も一般的な同期のタイミングは、プロセスがシステムコールを発行する際である。システムコールはユーザアプリケーションが単体で行うことのできない処理をカーネルに依頼するものである。情報の漏洩や新たなプロセスの立ち上げなど、攻撃者が何かしらの任意の動作をプログラムに行わせようとした際には必ず行われる処理である。

モニタは特定のバリエントをマスターと定め、他のバリエントをスレーブとして扱うマスターバリエントの発行したシステムコールとその結果を記録しておき、各スレーブバリエントが同期ポイントに達するタイミングで動作を比較する。異なる挙動を検出した場合、モニタは攻撃を受けたと判断し、直ちに全バリエントの実行を停止させることができる。

MVEEs のモニタリング方式は、ユーザアプリケーションによる実装とカーネルによる実装の 2 通りに大別される。ユーザアプリケーションによるモニタは `sys_ptrace()` システムコールを用いて実現する。この手法ではカーネルに手を入れる必要がないため容易に導入することができる。しかしその反面、実行時にはコンテキストスイッチが多くな

り、またバリエントのデータをコピーする手間も増えるため、オーバーヘッドが増加する傾向にある。カーネルによる実装では、システムコールを呼び出した際のハンドラとしてモニタの処理を割り込ませる手法を取る。この場合は、ほとんど遅延なくモニタに処理を渡すことができるため、`ptrace` によるようなオーバーヘッドは無い。だが、カーネルへ変更を加える場合では導入の手間や TCB の増加が懸念される。

2.3 関連研究

MVEEs については、モニタの実現方法や高速化、マルチスレッドへの対応など様々な側面から研究がなされている。MvArmor [4] はハードウェアの仮想化支援を受けたユーザアプリケーションレベルの MVX (= MVEEs) モニタである。モニタは Dune [17] を利用して仮想環境上で動作する。そのため、モニタはカーネルに手を加えずに特権命令を発行することができる。モニタは仮想環境においてプロセスの状態に直接アクセスすることができるため、コンテキストスイッチの増加を避けることができる。アプリケーションが発行した `syscall` を受け取ったモニタは、`vmcall` を発行する。この際、VM Exit が行われてホストのカーネルにロードされた Dune module に処理が渡るが、`syscall` ごとに毎回この処理を行うのはコストが大きい。そこで、本手法ではメモリ管理や `pid` の取得など、一部のシステムコールは VM Exit せずにモニタ内で処理を行う。

Orchestra [5] はユーザ空間で動作するモニタである。`sys_ptrace()` システムコールを用いてバリエントを監視することで、OS カーネルを変更せずにモニタの実装を実現する。同期ポイントはバリエントがシステムコールを呼び出す際としており、システムコールの番号や引数の比較を行う。本手法では脅威モデルとしてスタックオーバーフローを据えている。GCC に変更を加えることで、スタックの成長方向が逆 [7] となるプログラムを用意する。このプログラムを通常のプログラムと共に動作させることで、一方ではスタックオーバーフローによるリターンアドレスが成功するが、もう一方では失敗するようになる。

Taming Parallelism [6] では特定のバリエントにおけるイベントの実行順序を補足し、他のバリエントで再現することでバリエント間における動作の多様化を防ぐ手法を提案している。MVEEs にはマルチスレッドへの適用が困難という問題が存在する。バリエントごとでスレッドのスケジューリング順序が異なると、それに伴ってシステムコールや共有メモリへのアクセスなどの順番が変化し、動作が多様化してしまう恐れがある。動作の多様化を検知したモニタは、これが攻撃によるものなのか否か区別することができず、誤検知に繋がる。本手法ではアプリケーションがロックを取る順番を一致させ、結果的に同一順序でメモリにアクセスを行うようにする。ロックの順序を監視す

るためにバイナリレベルでロックに相当する命令を検出し、LLVM [18] を用いて命令の前後に同期用関数を挿入する。この関数内部で専用のシステムコールを発行し、適当なロックの取得を試みているスレッドのみ許可し、他はブロックする。また、ストール時間を減らすために WoC と呼ばれる論理クロックの概念を導入し、依存関係のあるロック操作の順序のみを再現する。

Bunshin [10] は、セキュリティ機構を部分適用させたバリエントを作動させる N-Version システムである。メモリエラーに対する様々な対策手法が提案されているが、それらはオーバーヘッドが大きく、中には実行時間が2倍以上になるようなものも存在する。対策を一部分にのみ適用すると、実行時間の回復は見込めるがその反面保護性能は低下してしまう。また、複数の対策手法を同時に適用することが難しい場合がある。本手法では、複数の対策手法を別々に施したバリエントを用意し、それらを統合して監視する。コストが大きい対策手法については、同一の手法を一部ずつ異なる部分に対して適用したバリエントを用意することでオーバーヘッドの減少を試みる。部分適用の分割部分の決定のために対策手法の適用の有無でプロファイリングを行い、check code を挿入する。対策手法の適用と check code 挿入の位置によってバリエントの実行時間を調整し、それぞれのオーバーヘッドが同じになるようにする。バリエントの同期ポイントはシステムコールの呼び出し時としている。

2.4 問題点

MVEEs ではバリエントを生成した分だけ物理メモリは消費される。バリエントはそれぞれ独立した仮想アドレス空間を持っている。Linux には無駄なページの複製および物理ページの消費を低減するため、Copy-on-Write と呼ばれる機構が存在する。sys_fork() などによってページの複製が必要となった際、一時的に同一物理ページを参照するようにし、書き込みがあった際に初めて複製を行う。MVEEs のバリエントについても同様である。sys_fork() してバリエントを作成した際、マッピングされたページに書き込みを行うまでは物理ページは共有されている。しかし、ひとたび書き込みを行えば新たな物理ページが確保される。バリエントは並列に実行されてはいるが、メモリにアクセスするタイミングは必ずしも一致するということはない。そのため、常にページを共有しておくことは不可能であり、最終的にほとんどのページを個別の物理ページに割り当てることになる。

バリエントを作成した後に sys_execve() によって新たにプログラムを実行した場合でも物理メモリの無駄は生じる。sys_fork() のみを行った場合とは異なり、この場合は初めからほとんどすべての物理ページは共有されていない。しかしながら、MVEEs では全バリエントに対して与えられ

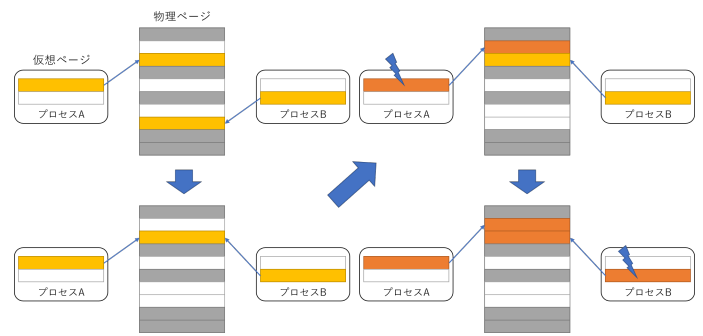


図 2 ページの併合・分裂サイクル

る入力は等しいため、内部に記録されるデータは同じものになる。すなわち、バリエントの数だけ同一内容のページが重複してしまう。

大量にメモリを消費するプログラムに MVEEs を適用する場合、これは無視することのできない無駄である。MVEEs を適用して n 個のバリエントを作成した場合、消費する物理メモリは n 倍近くとなる。数 MiB 程度のメモリを使用するアプリケーションに MVEEs を適用した場合にはそれほど大きな問題にはならない。しかしながら、memcached のようなひとつのサービスで数 GiB のメモリを利用する状況では、たった一つのサービスとそのレプリカで、ほとんどの物理メモリを消費してしまう場合も考えられる。実際に、Amazon DynamoDB Accelerator (DAX) [13] を提供する Amazon Elastic Compute Cloud (Amazon EC2) [19] のインスタンスでは 488GiB まで利用可能としている。物理メモリの逼迫によってスラッシングが発生などすると、MVEEs の監視下でない他のプロセスに対しても性能に悪影響を及ぼす恐れがある。

3. 提案

本研究では、メモリを大量に扱うアプリケーションを対象とした MVEEs を提案する。提案手法では、こうしたアプリケーションを MVEEs 上で効率的に実行する。提案手法は以下のデザインゴールに基づいている。

(1) 物理メモリ使用量を抑える。

MVEEs を利用している状況においても資源の消費を抑制し、他のプロセスへの影響を極力避ける。

(2) ユーザアプリケーションのソースコードを改変しない。利用者がソースコードに手を入れることなく提案手法を利用できることは、適用のハードルを下げることに貢献する。

(3) multi-thread アプリケーションをサポートする。

近年のアプリケーションではワーカを分散させて、比較的空きのあるワーカに処理を割り振ることが行われている。そのため、複数スレッドに対してもモニタリングを行える必要がある。

MVEEs では、全てのバリエントには等しい入力が与え

られるため、ページ内部に保持されるデータもまた等しい。併合を行わない場合は、それぞれのバリエントで個々の物理メモリを確保し利用している。したがって、内容が重複するページが複数存在することになる。本手法では、これら内容が等しいページを全バリエントを通して一つにまとめる。ページの併合・分裂サイクルを図2に示す。併合先となる物理ページを内容が一致する物理ページの中から一つ選び、全てのバリエントの仮想ページからの参照をこれに向ける。結果として、アプリケーションを単体で動作させた状態と同等程度までメモリ消費の削減が期待できる。あるバリエントがこの併合したページに対して書き込みを行おうとした際、直接ページに書き込みが行ってしまうと他のバリエントでは不整合が生じる可能性がある。このような場合には再び新たな物理ページを用意し、参照を指し直すことで分裂させる。

提案手法を実現するためには、定期的に全バリエントの仮想ページを走査する必要がある。しかしながら、MVEEsのモニタがこれを行うのはコストが大きく、バリエントの動作速度を著しく低下させる恐れがある。また、MVEEsに単純にページの併合を適用しようとしても高い併合率は見込めない。なぜなら、ページ内部にメモリアドレスを含む場合は、その値はASLRによってバリエント同士で一致しないためである。本研究ではこれらの課題に対処できるように機構を設計し、コストの低いページ併合とバリエント同士の高いページ併合率を実現する。

4. 設計

4.1 MVEEs モニタ

4.1.1 モニタ・バリエントの生成

MVEEsのモニタを生成する際には、ユーザプログラムが自身を監視対象に追加するためのシステムコールを発行する。このとき、カーネルはモニタを新たなカーネルスレッドとして生成し、対象プロセス固有の管理構造にモニタとして記録する。その後プロセスの複製によりバリエントの生成が行われるが、このとき管理構造も同時に複製されるため、モニタの情報が引き継がれる。

今回実装したMVEEsは、モニタ、マスターバリエント、1個以上のスレーブバリエントで1組の監視構成である。構成を図3に示す。1つのモニタは1組のバリエントのみを監視し、他のバリエントの組の動作には関知しない。監視対象のバリエントがsys_fork()した場合には再びモニタを生成する。バリエントそれぞれの子プロセスの管理構造に、生成したモニタを記録することで新たな1組の監視構成を成す。

4.1.2 システムコールの監視

モニタはリングバッファを用いてアプリケーションが発行したシステムコールを管理する。リングバッファには発行したシステムコールの番号および引数が保存される。モ

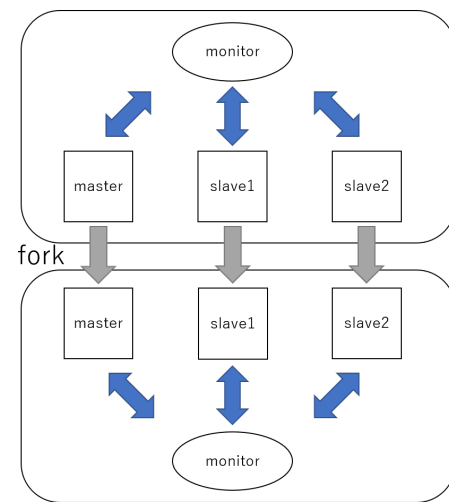


図3 fork時のモニタとバリエント

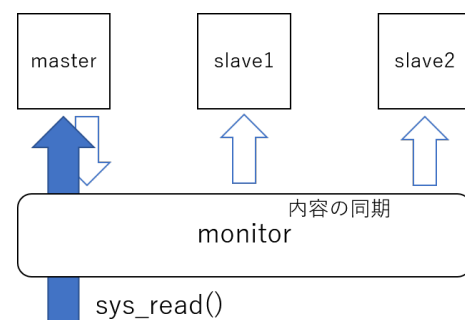


図4 readの同期

ニタはマスターおよびスレーブバリエントの全てがシステムコールを発行したことを確認すると、リングバッファに記録された情報を精査して実行の可否を決定する。

モニタの動作は、バリエントが発行しようとしているシステムコールによって異なる。システムコール番号が一致していることは必須であるが、引数についてはアドレスが含まれる場合があるため、単純に値のみを比較することはできない。実行することでプロセス内部にのみ変更が加わるシステムコールであれば、バリエント全てでそのまま実行することは何ら問題ではない。しかしながら、I/Oのような外部に対して影響を及ぼすシステムコールでは、それぞれが実行してしまうと不整合が生じる恐れがある。これを考慮して、マスターバリエントのみ実行を許可し、スレーブバリエントの実行は不許可にする必要が生じる場合がある。

4.1.3 外部との通信

バリエントが外部と情報をやり取りするためにはモニタの介入が必要となる。sys_read()によってデータを受け取る際には、一旦マスターバリエントのみがシステムコールを実行してメモリ内にデータを格納する。スレーブバリエントに対しては図4のように戻り値および内容を同期し、あたかもシステムコールが実行されたかのように見せかける。

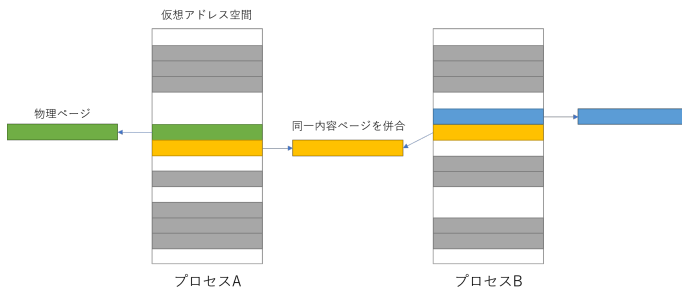


図 5 プロセス間でのページの併合

sys.write() によってデータを受け渡す際には、内容が完全に一致しているかどうかを検証してから実行を行う。read と同様に、この場合でも実際にシステムコールの実行を行うのはマスターバリエーションのみであり、スレーブバリエーションに対しては戻り値を同期する。ここで出力するメモリの内容が完全に一致しているかどうかを検証することにより、センシティブな情報であるメモリアドレスを漏洩させることはできなくなる。ASLR によってページ単位でランダムに配置が行われているため、アドレスが含まれる箇所を write しようとした際に内容が一致しなくなるためである。

4.2 同一内容ページの併合

ページ内容の比較及び併合処理は、モニターやバリエーションとは独立したカーネルスレッドで行う。併合の処理では各プロセスが利用している仮想アドレス空間を走査し、内容が一致するページを見つけ出す必要がある。そのため、アプリケーションや MVEEs モニタが逐次的に行おうとすると、その間に本来の処理を行うことができず、スループットの低下を招く。別スレッドに処理を委ねることで、MVEEs の動作に影響を及ぼすことを避けることができる。

ページの併合は、複数の異なる仮想ページのページテーブルから、図 5 のように同一の物理ページを指すように変更することで実現する。ページテーブルから外された物理ページは参照されないため、解放することができる。しかし、同一内容のページを併合したとしても、各バリエーションは再び該当ページに書き込みを行う場合がある。その際に、直接ページに書き込みが行えてしまうと他のバリエーションでは不整合が生じる可能性がある。そこで、書き込み時には Copy-on-Write [20] によって再び別の物理ページが設けられ、分裂するようにする。それ以降は別に設けたページを利用することで、他のバリエーションには影響を及ぼさずに既に併合していたページに対して書き込みを行うことができる。

書き込み頻度の高いページを併合することはコストが大きい。先述の仕組みにより、併合と分裂を頻繁に繰り返すことはシステム全体の性能劣化を招く恐れがある。内容の一致する全ページを併合するのではなく、書き込み頻度が

低く安定しているページのみを対象とするべきである。

4.3 ポインタのハッシュ化

単純なデータのみを含むページだけではなく、ポインタを含むデータ構造が存在するページに対しても併合は行うことができるのが好ましい。ページ内にアプリケーションのデータのみが含まれる場合はバリエーション間でページ内容は完全に一致する。しかしながら内部にポインタを含む場合は、そのポインタは ASLR によってランダムに配置されたアドレスであるため、バリエーション間で一致する可能性は非常に低い。この場合、単純にページを併合することはできない。これらのページに対しても併合が行えるように、ポインタの表現に手を加えることを考える。

本研究ではポインタをハッシュ値として扱うことで、メモリアドレスを含むページに対してもバリエーション間で併合を可能にする。このポインタのハッシュ値を、以後ハッシュポインタと呼称する。

4.3.1 ハッシュ化とアドレスの解決

ハッシュポインタは、アドレスと一対一で対応する一意な値であるが可逆である必要はない。生成したハッシュポインタと実際の仮想アドレスを対応付けるテーブルを保持することで、ハッシュポインタを利用してメモリにアクセスする際に解決を行う。

アドレスをハッシュに変換する作業、およびその解決はカーネル内で行う。プログラムはメモリアドレスをメモリに書き出す際に、生のアドレスを引数にとって変換用のシステムコール sys.hashptr_gen() を発行する。依頼を受けたカーネルは所定の手順に従ってアドレスをハッシュ値へと変換し、テーブルに対応付けを追加する。ハッシュ値への変換には、プロセス固有のランダム値で排他的論理和を取るなどの処理を挟むため、復元は容易ではない。変換したハッシュ値を受け取ったプログラムは、これをポインタとして扱いメモリに記録する。ハッシュポインタから仮想アドレスを解決する場合も同様である。図 6 に示すように、プログラムは解決したいハッシュポインタを引数に解決用のシステムコール sys.hashptr_res() を発行し、カーネルはハッシュポインタテーブルを走査する。対応付けられた仮想アドレスが存在する場合はこれを返し、プログラムはメモリにアクセスすることができる。

この機構に対応させるために、プログラムのコードには手を加える必要はない。ソースコードが存在する場合、LLVM [18] によってコンパイル時にポインタへのアクセスを捕捉することが可能である。メモリアドレスを書き出す直前には変換用の、読みだした直後には解決用のシステムコールを発行する関数を挿入する。そうすることで、プログラムはメモリのハッシュポインタを読み書きする以外では、これを利用していることを意識せずに動作することができる。

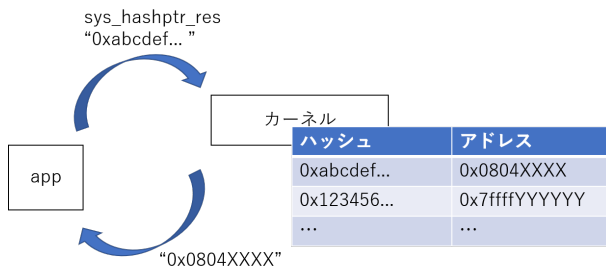


図 6 ハッシュポインタの解決

ソースコードが存在しないバイナリに対しても、ポインタのハッシュ化機構は一部対応可能である。全ての anonymous page を対象とすることはできないが、ヒープ領域に対しては malloc() および free() 等のヒープチャンクを操作する関数をフックすることで対応ができる。memcached [12] のような In-memory DB などは、ヒープを用いてデータを格納することが多い。ユーザプログラムがポインタを含んだデータ構造を内部に持たせなかったとしても、ヒープにはチャンクのフリーリストを構成するためのリスト構造が保持されている。これもまた、ページの併合を阻害する要因の一つである。そこで、glibc [21] に実装されているヒープの操作関数は用いずに、同等の機能を実装したコードを LLVM でコンパイルした共有ライブラリを利用する。これにより、リスト構造のポインタに対してもハッシュ化することができる。

バリエーション間でページを併合するためには、生成されるハッシュポインタが全バリエーションを通して一致している必要がある。マスターバリエーションのアドレスを基に生成したハッシュポインタを、スレーブバリエーションに対しても配布することでこれを実現する。

4.3.2 ハッシュポインタテーブルの管理

ハッシュポインタテーブルは task_struct 構造体から参照することで、プロセスごとに紐づけて保持する。プロセスが sys.fork() した際には、子プロセスにおいてもこれまで変換したハッシュポインタを解決できる必要があるため、親プロセスが保持しているハッシュテーブルを参照できるようにする。一旦はテーブルを共有していても問題はないが、親または子プロセスが新たなハッシュを生成した際に共有している同じテーブルにそれを登録することは好ましくない。したがって、このような場合にはテーブルを丸ごと複製し、新たなテーブルを task_struct 構造体に登録して利用するようにする。

sys.execve() によって仮想アドレス空間が一新される際には、以前のハッシュと仮想アドレスの対応は不要となる。そのため、sys.execve() 時にはこれまで利用していたテーブルのエントリは全て削除する。

4.3.3 ASLR との共存

今回提案するハッシュによるポインタの表現では ASLR

の利点を損なう恐れがある。攻撃者が Exploit を作成する際、漏洩したメモリアドレスを利用することがある。MVEEs では単一のバリエーションを狙った Exploit を送ったとしても他のバリエーションでは正常に動作しないため攻撃を検知することができる。しかしながらハッシュポインタを用いた場合、バリエーション間ではハッシュの値は同期しているため、Exploit に含まれたハッシュポインタを自動で解決し、攻撃が成功してしまう可能性がある。

攻撃者にハッシュポインタを利用させないためには、ハッシュの値そのものを漏洩させないことと推測させないことが重要である。MVEEs では、4.1.3 小節で述べたように write 時に内容を比較しているが、ハッシュ値は一致しているためこのままでは容易に漏洩してしまう。そこで、出力しようとしている領域を 8byte 単位でアライメントしたうえで、その範囲にハッシュポインタが含まれるかどうかを確認する。ハッシュポインタが存在している場合は不正な動作として実行を停止する。しかしながら、この方法でハッシュ値の漏洩を完全に防ぐことはできない。glibc に実装されている printf 系関数の、その使用方法の誤りに起因する Format String Bugs (FSB) [22] と呼ばれるバグが存在する。このバグを利用した書式指定文字列攻撃 [23] など、システムコールを直接利用せずに値を別の形式に変換して出力することでハッシュ値を漏洩させることも考えられる。

だが、漏洩させたハッシュ値から、攻撃者が必要とするハッシュ値を生成することは非常に難解である。4.3.1 小節に示したように、このハッシュ値はメモリアドレスそのものではなく、処理を施した値を基に計算している。メモリアドレス自体とそれを 2byte シフトした値、さらにランダム値の排他的論理和を取ったものである。総当たりで計算することでハッシュ値からこの元の値を復元することは原理的には可能であるが、64bit 空間でこれを試行するのは現実的ではない。仮に元の値が復元できたとしても、ASLR によってメモリアドレスもページサイズ単位でランダムとなっているため任意のオフセットを加えたメモリアドレスのハッシュを求めることは非常に難しい。それに加えて、ハッシュポインタは既に変換したアドレスとそのハッシュの組み合わせしか解決できないという性質もある。以上の事から、ASLR とハッシュポインタを組み合わせることで、攻撃者が任意のアドレスのハッシュを生成して攻撃を成功させる可能性は非常に低く抑えることができるといえる。

5. 実装

本章では、提案機構の実装について述べる。対象とする OS カーネルは Linux 4.13.9 とする。OS カーネルには MVEEs のモニタと、ハッシュポインタの生成および解決の機能を実装し、またページの併合を行う。アプリケーションがハッシュポインタを利用するようにするため

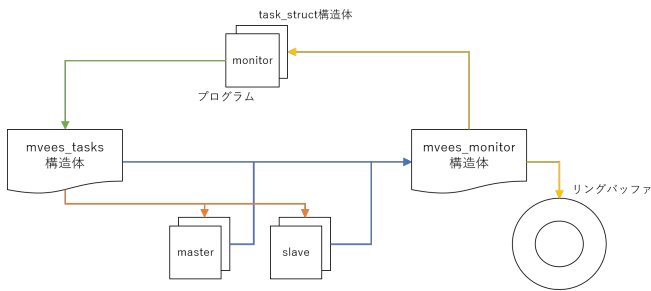


図 7 モニタとバリエーションの管理構造

LLVM の Pass を作成し、ソースコードに手を加えることなくバイナリに変更を加える。また、本提案手法をアプリケーションが容易に利用できるようにするため、API を作成しライブラリを提供する。

5.1 Linux カーネルへの実装

本節では、MVEEs のモニタとハッシュポインタを実現するために OS カーネルに対して行った実装について述べる。ページの併合については、現在 Kernel Samepage Merging (KSM) [24] の機構を利用している。

5.1.1 モニタとバリエーションの管理構造

各バリエーションは、同一の `mvees_monitor` 構造体を共有し、その構造体へのポインタを `task_struct` 構造体内に保持する。`mvees_monitor` 構造体には、5.1.2 小節で後述するリングバッファや、モニタそのものの `task_struct` 構造体への参照が含まれる。

モニタは、`mvees_tasks` 構造体を持っている。この構造体には各バリエーションの `task_struct` 構造体への参照が保持される。これによりマスターバリエーション、および全てスレーブバリエーションが管理されることになる。また、バリエーションが保持している `mvees_monitor` 構造体への参照も持つことでリングバッファへのアクセスが可能になる。図 7 にこの管理構造を示す。

5.1.2 リングバッファを利用したシステムコールの監視

各バリエーションが発行するシステムコールのリクエストは、全て `syscall_req` 構造体のリングバッファに記録される。記録した後、バリエーションはスリープし、モニタがチェックするまで待機する。この構造体に記録する内容はシステムコールの番号および引数、またモニタによってシステムコールの実行が許可されたか否かを示すフラグである。

図 8 に示すように、モニタはリングバッファに記録されたシステムコールのリクエストを順番に取り出し処理を行う。全バリエーションのリクエストが出揃った時点で比較を行う。比較の結果、一致が確認できればシステムコールに応じて `syscall_req` 構造体のフラグに許可を記録し、バリエーションのプロセスを起床させる。システムコールによっては実行の許可ではなくスキップを命じる場合もある。この場合、バリエーションはシステムコールの実行を行わずに戻り値

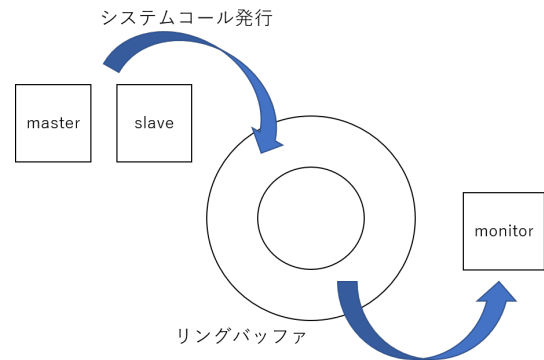


図 8 リングバッファによる管理

同期の処理に飛ぶ。全バリエーションに対するチェックが終了したのち、リングバッファが空の間はモニタはスリープし、無駄なリソースの消費を避ける。

システムコールを実行した後の戻り値の比較も同様である。全バリエーションは `syscall_res` 構造体のリングバッファに戻り値を格納し、モニタがチェックを行う。

5.1.3 ハッシュポインタの利用と管理

カーネルはアプリケーションに対して表 1 に示すシステムコールを提供する。アプリケーションから `sys_hashptr_gen()` の発行を受け取ったカーネルは、引数のメモリアドレスからハッシュ値を生成して返す。`sys_hashptr_gen()` ではその逆で、ハッシュ値から元のメモリアドレスを解決して返す。

ハッシュとメモリアドレスの対応は `hashptr` 構造体で管理されている。この構造体は二分木構造で接続され、ハッシュとアドレスのいずれからでも該当するエンタリにたどり着けるようにしている。ハッシュをキーとした二分木はアドレスの解決に、アドレスをキーとした二分木は生成済みのハッシュを高速に返すために利用される。木の先頭はそれぞれ 256 個用意されており、ハッシュでは最上位バイト、アドレスでは下位 2 バイト目の値で振り分けられる。

ユーザプログラムがハッシュポインタの生成をカーネルに依頼した際、カーネルはアドレスをキーとして構築されている二分木を辿ってエンタリを探す。該当エンタリが存在している場合は即座に記録されているハッシュを返す。存在しない場合は新たに生成したハッシュ値とアドレスを格納したエンタリを生成し、ハッシュとアドレスそれぞれの二分木につなげる。初めからアドレスの木を辿らずに、ハッシュ値の木のみを辿りエンタリの有無確認および繋げる方が動作としては高速だが、アドレスによる二分木を扱っているのはハッシュの衝突を考慮したためである。

表 1 システムコール一覧

syscall	機能
<code>sys_hashptr_gen</code>	ハッシュポインタの生成
<code>sys_hashptr_res</code>	メモリアドレスの解決
<code>sys_write_nohash</code>	内容からハッシュを除いた write


```
void main(int argc, char *argv[], char *envp[]){  
    if(argc<3) return;  
    mvees(atoi(argv[1]));  
    execve(argv[2], &argv[2], envp);  
}
```

図 9 簡易ランチャープログラム

5.2 ユーザアプリケーションの対応

本章では、ユーザアプリケーションを本提案手法に対応させるために実装を行った機構について述べる。

5.2.1 MVEEs の利用

アプリケーションのバリエーションを作成し MVEEs の監視モニタの配下に置くためには、4.1.1 で述べたように `sys_mvees()` システムコールを発行する必要がある。アプリケーションは作成したいスレーブバリエーションの数を引数にとって、`sys_mvees()` を呼び出す。したがって、アプリケーションを監視対象とするためには起動した時点でこのシステムコールを呼び出すようにコードに変更を加える必要がある。しかしながら、それではソースコードが入手できないバイナリに対して適用が難しい。

既存のプログラムを監視する場合には、ランチャーとなるプログラムを介せばよい。図 9 に示すように、`sys_mvees()` を発行して自身を監視対象としてから `sys_execve()` によって既存のプログラムを起動する。バリエーションとモニタの参照関係は `task_struct` 構造体に記録されているため、`execve` によって別のプログラムを起動したとしても MVEEs の動作に必要な情報は引き継がれる。したがって、そのままモニタリングを継続することが可能となる。

5.2.2 ハッシュポインタの利用

アプリケーションのソースコードに変更を加えることなく提案手法を利用するため、LLVM を用いてコンパイルを行う。この際、ポインタに対するアクセスを検出し、指定する関数を呼び出す Pass を作成し用いる。ポインタ変数へ値を代入する処理では、その直前に `Addr2Hash()` 関数を呼び出し、その戻り値を変数へ格納する。ポインタ変数から値を取り出す際には、レジスタへ値を移動した直後に `Hash2Addr()` 関数を呼び出し、以後解決したアドレスを利用する。

5.1.3 小節で示したシステムコールを利用するため、ライブラリで API を提供する。先述の `Addr2Hash()` と `Hash2Addr()` 関数は、それぞれ `sys_hashptr_gen()` と `sys_hashptr_res()` を発行する関数である。

ヒープチャンクのフリーリスト構造もハッシュポインタに対応させるため、glibc のヒープの実装を模倣した独自のライブラリを作成した。このプログラムも上記同様の Pass を用いて LLVM でコンパイルすることでこれを実現する。既存のプログラムを動作させる際には、LD_PRELOAD で指定することで本ライブラリのヒープ操作系関数を利用するようになる。

6. 実験

6.1 実験環境

本提案手法の評価を、表 2 に示すコンピュータで行う。Dell PowerEdge T630 II は 16 コア、125GiB のメモリを搭載している。OS は Ubuntu 16.04.1 LTS、カーネルは提案手法を実装した Linux 4.13.9 である。

6.2 システムのコスト評価

MVEEs を適用した状態でマイクロベンチマークを実行し、処理にかかる時間を測定する。スレーブバリエーションは一つのみとする。MVEEs を利用していない通常の状態でもベンチマークを測定し、MVEEs によるコストを評価する。

また、ハッシュポインタの利用に対するコスト評価も行う。ハッシュポインタを有効にした条件下でマイクロベンチマークを実行し、処理に要する時間を測定する。

6.2.1 実験方法

マイクロベンチマークプログラムで、メモリに対する読み書きのスループットを計測する。確保するメモリサイズは 1GiB とし、各ページに対してシーケンシャルにアクセスを行う。これを 1 セットとし、10 回くりかえして処理に要した時間を測る。

メモリに書き込みアクセスを行う際、数回に一度システムコールを発行する。MVEEs の同期ポイントの出現頻度を変化させ、コストの増加の程度を確認する。本実験では `brk()` システムコールを計測のために発行する。発行頻度は、メモリアクセス 10 万回に対して 0 回、1 回、5 回、10 回の 4 通りとする。これは、1GiB 分のページにアクセスする間にそれぞれ 0 回、2.6 回、13 回、26 回程度システムコールを発行することになる。ハッシュポインタの利用によるコスト評価では、システムコールは発行しない。

比較対象として、MVEEs を用いないネイティブ環境での測定も行う。

6.2.2 実験結果

マイクロベンチマークを実行した際の、スループット結果を図 10 に示す。メモリにアクセスを行っている際にシステムコールの発行を行っていない場合では、MVEEs モニタを利用した場合でもネイティブと比べて 4% ほどスループットが低下している。10 万回のアクセスのうち 1 回システムコールを発行した際には、40% ほどの遅延が発

表 2 実験環境

機器名	Dell PowerEdge T630
CPU	Intel(R) Xeon(R) CPU E5-2623 v4 @ 2.60GHz
キャッシュ	10240KB/コア
コア数	16
RAM	125GiB

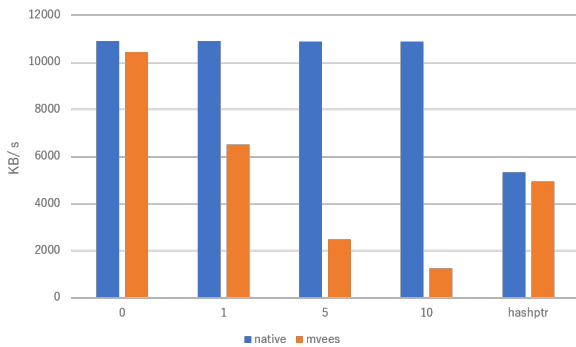


図 10 MVEEs モニタとハッシュポインタのスループット

生している。5回、10回発行した場合にはそれぞれ77%、88%ほど遅延している。同期ポイントが増えるにつれて、スループットが低下が非常に大きくなる傾向にある。このような結果になったのは、同期の度に全バリエントが一旦停止するためであると考えられる。

ハッシュポインタを利用した場合には、MVEEs を用いない場合でも50%程度スループットが低下している。一つのページへのアクセスに対して、複数回アドレスのハッシュ化と解決が行われている。そのため、その度に発行されている `sys.hashptr_gen()` と `sys.hashptr_res()` が大きなボトルネックとなっている。また、同一アドレスのハッシュ化や解決が何度も行われる傾向にあるため、キャッシュなどの利用による変換の効率化が求められる。

6.3 実メモリ使用量評価

MVEEs を適用した状態でマイクロベンチマークを実行し、消費される物理メモリ量を測定する。ハッシュポインタを利用しない条件下でも同様に測定を行い、提案手法による物理メモリ使用量の削減を確認する。

6.3.1 実験方法

6.2節で用いたものと同じベンチマークプログラムを用いて実験を行う。スレーブバリエントの数は一つとする。

本実験で確保するメモリサイズは4GiBとする。始めに全ページを初期化し、あらかじめ必要な物理ページが全て確保された状態から測定を開始する。ページの初期化に用いるデータには、ポインタが存在するデータ構造を想定してアドレスを含むようにする。実験ではこれらのページに対して10セットだけシーケンシャルに書き込みを繰り返す。ただし、1セットごとに10秒の休止を挟む。

また、本実験では偏りのある書き込みのワークロードを想定する。全体の1%、10%、50%の領域に対してのみ書き込みを繰り返し、他の領域の内容に変更は加わえない。このとき、各バリエントにおける物理メモリ使用量を `procs` の `smaps` から1秒ごとに読み取る。

6.3.2 実験結果

ハッシュポインタを有効にした条件の下でマイクロベン

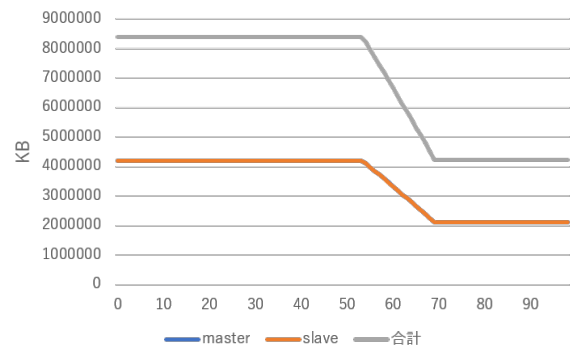


図 11 1%に偏ったワークロードでの物理メモリ使用量

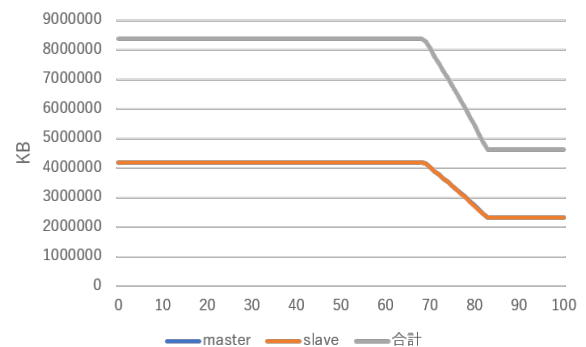


図 12 10%に偏ったワークロードでの物理メモリ使用量

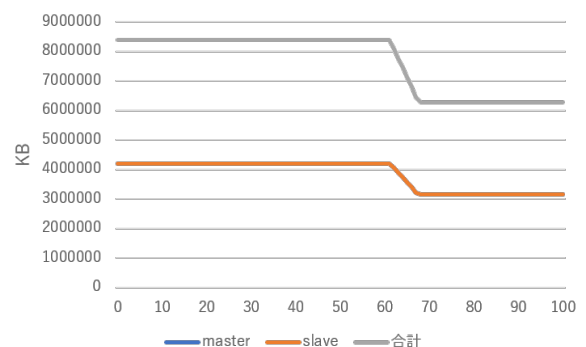


図 13 50%に偏ったワークロードでの物理メモリ使用量

チマークを実行した際の、物理メモリの使用量の遷移を図11, 12, 13に示す。マスターバリエントとスレーブバリエントの物理メモリ使用量は、どの条件下においても常にほぼ等しい。スレーブバリエント数は1つとしているため、計測開始直後の合計の物理メモリ使用量は、単体で動作させた場合の2倍である8GiBとなっている。

開始から1分前後の時点でKSMがページの併合を行っている。頻繁に書き込みが行われるページは併合の対象外となるため、50%に書き込みが偏ったワークロードでは4GiBの半分が重複し、物理メモリ使用量の合計は6GiB程度となっている。それに対して、1%に偏ったワークロードでは、ほとんどアプリケーションを単体で動作させた場合と同等程度まで使用量の合計が減少している。

なお、ハッシュポインタを無効とした状態で書き込みが

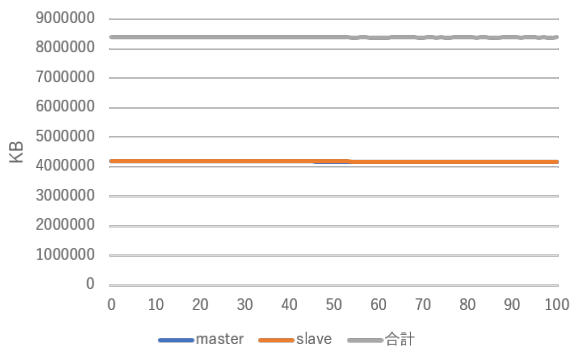


図 14 ハッシュポインタ無効時の物理メモリ使用量

1%に偏ったワークロードをかけた際の、物理メモリの使用量は図 14 に示すとおりである。書き込みが行われるページが非常に偏っていたとしても、ページにはメモリアドレスが含まれるためページが内容が一致せず、併合することはできない。

以上より、提案手法を MVEEs に適用することで、ワークロードによっては最大でアプリケーションを単体で動作させた場合と同等程度まで物理メモリの使用量を削減することができることが確認された。

7. おわりに

本研究では、大量にメモリを消費するアプリケーションに対して MVEEs を適用する際に、同一内容のページを併合することで物理メモリの消費を低減する手法を提案した。しかしながら、ASLR によってプロセスの仮想アドレス空間はランダムに配置されているため、そのままではメモリアドレスが一致せずページを併合することはできない。そこで、アドレスをハッシュ値を用いて扱うようにし、バリエーション同士でこの値が一致するようにしたことでページの併合を可能とした。

提案手法をカーネルに組み込み、またハッシュポインタを扱うように LLVM を用いてアプリケーションをコンパイルした。これらに対して性能を評価し、手法を適用していない場合との比較を行った。その結果、提案手法を適用していない場合ではバリエーションの数だけ物理メモリが消費されているのに対して、適用した場合は最大でアプリケーション単体が消費するのと同程度まで使用量を削減することができた。ただし、削減できる量はワークロードに大きく依存し、書き込みを行うページに偏りが少ない場合ではこれほどの使用量削減は期待しづらい。

今後の課題として、MVEEs モニタと連携したページの併合機構の実装が挙げられる。現在、ページの併合は KSM の実装を利用しているため、モニタとは独立してページのスキャンを行っている。この状態でも十分な物理メモリ使用量の削減は行えるが、併合後に再度 CoW によって分割されたページの再併合が行われないなどの不都合が生じ

るためである。また、より効率的な MVEEs の監視機構の実装が必要である。今回行った MVEEs のモニタによるコスト評価から、システムコールの発行が多いほどオーバーヘッドが非常に大きくなることが判明した。これは、システムコールを発行する度に全バリエーションを同期させ、揃うまで停止をさせているためであると考えられる。マスターバリエーションのみは入出力以外では停止させず後々比較を行うことで、スループットの低下を抑えた監視が期待できる。

参考文献

- [1] Eklectix Inc. Kernel address space layout randomization [LWN.net]. <https://lwn.net/Articles/569635/>.
- [2] Crispin Cowan, Calton Pu, Dave Maier, Heather Hintony, Jonathan Walpole, Peat Bakke, Steve Beattie, Aaron Grier, Perry Wagle, and Qian Zhang. StackGuard: Automatic adaptive detection and prevention of buffer-overflow attacks. *Science*, p. 5, 1998.
- [3] Stijn Volckaert, Bjorn De Sutter, Tim De Baets, and Koen De Bosschere. GHUMVEE-Efficient, Effective And FlexibleReplication. *International Symposium on Foundations and Practice of Security FPS 2012: Foundations and Practice of Security*, pp. 261–277, 2012.
- [4] Koen Koning, Herbert Bos, and Cristiano Giuffrida. Secure and efficient multi-variant execution using hardware-assisted process virtualization. *Proceedings - 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks, DSN 2016*, No. Mvx, pp. 431–442, 2016.
- [5] Babak Salamat, Todd Jackson, Andreas Gal, and Michael Franz. Orchestra: intrusion detection using parallel execution and monitoring of program variants in user-space. *EuroSys '09 Proceedings of the 4th ACM European conference on Computer systems*, p. 14, 2009.
- [6] Stijn Volckaert, Bart Coppens, Bjorn De Sutter, and Michael Franz. Taming Parallelism in a Multi-Variant Execution Environment. *Proceeding EuroSys '17 Proceedings of the Twelfth European Conference on Computer Systems*, pp. 270–285, 2017.
- [7] Babak Salamat, Andreas Gal, and Michael Franz. Reverse stack execution in a multi-variant execution environment. *Workshop on Compiler and Architectural Techniques for Application Reliability and Security*, pp. 1–7, 2008.
- [8] Babak Salamat, Andreas Gal, Jackson Karthikeyan, Todd Manivannan, Gregor Wagner, and Michael Franz. Multi-variant program execution: Using multi-core systems to defuse buffer-overflow vulnerabilities. *Proceedings - CISIS 2008: 2nd International Conference on Complex, Intelligent and Software Intensive Systems*, pp. 843–848, 2008.
- [9] Stijn Volckaert, Bart Coppens, Alexios Voulimeneas, Andrei Homescu, Per Larsen, Bjorn De Sutter, and Michael Franz. Secure and Efficient Application Monitoring and Replication. *2016 USENIX Annual Technical Conference (USENIX ATC 16)*, pp. 167–179, 2016.
- [10] Meng Xu, Kangjie Lu, Taesoo Kim, and Wenke Lee. Bunshin: Compositing Security Mechanisms through Diversification (with Appendix). 2017.
- [11] P Hosek and C Cadar. Varan the Unbelievable: An efficient N-version execution framework. *Proceeding AS-*

- PLOS '15 Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 339–353, 2015.
- [12] Dormando. memcached - a distributed memory object caching system. <http://www.memcached.org/>.
- [13] Amazon Web Services Inc. Amazon DynamoDB Accelerator (DAX). <https://aws.amazon.com/dynamodb/dax/>.
- [14] Byoungyoung Lee, Chengyu Song, Yeongjin Jang, and Tielei Wang. Preventing Use-after-free with Dangling Pointers Nullification. *NDSS Symposium 2015*, No. February, pp. 8–11, 2015.
- [15] Crispin Cowan, Perry Wagle, Calton Pu, Steve Beattie, and Jonathan Walpole. Buffer overflows: Attacks and defenses for the vulnerability of the decade. *Foundations of Intrusion Tolerant Systems, OASIS 2003*, pp. 227–237, 2003.
- [16] GNU Project. GCC, the GNU Compiler Collection. <https://gcc.gnu.org/>.
- [17] Adam Belay, Andrea Bittau, Ali Mashtizadeh, David Terei, David Mazi, and Christos Kozyrakis. Dune : Safe User-level Access to Privileged CPU Features.
- [18] llvm-admin team. The LLVM Compiler Infrastructure. <https://llvm.org/>.
- [19] Amazon Web Services Inc. Amazon Elastic Compute Cloud (Amazon EC2). <https://aws.amazon.com/ec2/>.
- [20] Eklektix Inc. User-space page fault handling [LWN.net]. <https://lwn.net/Articles/550555/>.
- [21] GNU Project. The GNU C Library. <https://www.gnu.org/software/libc/>.
- [22] Microsoft Corporation. Format String Bugs - MSDN. <https://msdn.microsoft.com/en-us/library/ee823826%28v=cs.20%29.aspx?f=255&MSPPError=-2147217396>.
- [23] Fatih Kilic, Thomas Kittel, and Claudia Eckert. Blind format string attacks. *Lecture Notes of the Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering, LNICST*, Vol. 153, pp. 301–314, 2015.
- [24] Andrea Arcangeli, Izik Eidus, and Chris Wright. Increasing memory density by using KSM. *Proceedings of the linux symposium*, pp. 19–28, 2009.