

# Trace format

---

## *Binary file format description, rationale and its uses*

Note: although the binary trace format described in this document is a pretty stable version, it is bound to evolve. As of such, both this document and the underlying file format are subject to change.

- 2013.10 : First stable release
- 2013.11 : Tool suite version 2.1
  - Support for *XED* version 20131119
- 2014.01 : Tool suite version 3.0
  - Support for *XED* version 20140106
  - Slight changes in **Execution** table (see **execRec.Flags**)
- 2014.05 : Slight corrections to documentation only
  - **cStr/uStr** member sizes and offsets
  - some **insRec** member sizes
- 2015.07 : Tool suite version 4.0
  - Support for *XED* version 20150121
  - Introduced *common register name* concept (see **insRec** versions in 7.3.2 and context map registers in 11.1)

## Table of Contents

1	Introduction.....	4
1.1	About traces .....	4
2	Format design rationale .....	5
3	Data types.....	7
3.1	Basic types.....	7
3.2	File “pointers” (offsets) .....	8
3.3	Arrays.....	8
3.4	Structures .....	8
3.5	Files.....	9
4	Trace format overview .....	10
4.1	Components .....	10
4.2	Generic table file format .....	11
5	Common structure types.....	13
6	Execution table.....	14
6.1	Underlying data format.....	14
6.2	File format .....	16
6.3	Versions.....	17
6.4	Execution record format .....	17
6.5	Execution order issues.....	20
7	Instruction table .....	21
7.1	Underlying data format.....	21
7.2	Decoded instruction .....	21
7.3	Instruction table file format .....	24
8	System calls table .....	29
8.1	Name and argument count .....	29
8.2	Argument values .....	30
8.3	Extras .....	30
8.4	File format .....	30
9	Module table .....	32

10	Thread table .....	33
11	Context map .....	34
11.1	Registers .....	34
11.2	File format .....	35
12	Address space usage .....	36
13	Trace.xml .....	37
14	External tables .....	38
14.1	Intel <i>XED</i> .....	38
14.2	System call IDs .....	38

# 1 Introduction

This document describes trace files format and its design rationale. It also explains the underlying data format, its logic, and roughly how a trace analyzer is expected to use trace data. All this extra information is threaded around the file format itself.

The reader is assumed to possess sufficient knowledge on *Microsoft Windows* architecture and on *Intel* x86 and x64 architecture and assembly. For thorough information on the latter, see *Intel® 64 and IA-32 Architectures Software Developer Manuals*:

<http://www.intel.com/content/www/us/en/processors/architectures-software-developer-manuals.html>

## 1.1 About traces

A raw trace is first generated from a running application with a home-made *Pin tool* called *UMTracer*, before being indexed and converted to be more usable. Raw traces come in two parts:

- Instrumentation trace:
  - Akin to static code analysis
  - Records instructions
    - Addresses and code bytes, not their execution.
    - Recording occurs when those instructions are instrumented, i.e. “augmented” to generate the execution trace when executed
  - Special events (load/unload modules)
- Execution trace:
  - Akin to dynamic code analysis
  - Generated by instrumented code as it runs
  - Records instructions execution
    - Registers and referenced memory values
  - Special events
    - Kernel and user mode transitions (system calls, context changes, exceptions)
    - Thread creation and termination
  - Sequential number on each record

A side tool, *UMTracerRemote* allows turning instrumentation on and off during the process. This only affects *instructions* instrumentation since *UMTracer* always record system calls and other special events. For general information about code instrumentation, see *Pin* documentation and papers:

<http://software.intel.com/en-us/articles/pin-a-dynamic-binary-instrumentation-tool>

*UMTIndex*, another home-made tool, converts a raw trace into an *indexed trace* (simply called *trace* in the remainder of this document). It validates the raw trace, recovers what it can if the trace seems incomplete (abrupt process or OS termination for instance), and reorders the execution trace according to sequence numbers. The result is stored into a few binary files.

## 2 Format design rationale

The trace format design was driven by a few use cases that all share a similar base scenario:

- A raw trace is generated while running an application with *UMTracer*
- The raw trace is then verified and indexed
  - Can be done at a later time, on a different machine, in a batch, etc.
- The trace is expected to contain
  - Several billions of small (< 100 bytes) records
  - Millions of widely variable size (10 bytes to 100 MB) records.
- One or many trace analyses then are performed offline
  - The trace is read-only
    - Annotations and other generated data are stored elsewhere.
  - Analyses are typically of sequential nature
    - Records are processed in chronological or control flow wise order

Given the above and the sheer amount of data to process, data fetching cost can be relatively significant to trace analysis intrinsic complexity. Data fetching performance boils down to a trade-off between:

- Easy data access:
  - Fully expanded values (propagated leading zero or sign bit)
  - Pre-computed pointers, offsets, etc.
    - Information redundancy in various forms
  - Structured row-column tables.
- Fast data access:
  - Related data are physically close to each other
    - Disk buffers and memory cache friendly
    - Semi-structured stream-like records (variable fields and size)
  - Minimal file size:
    - Little or no redundancy
    - Data reuse on identical chunks of data across records (multiple file “pointers” to the same location)
    - Disk access can be a bottleneck, unless running on powerful systems (high transfer speed to read trace into memory, and enough memory to fit the whole trace, perhaps 100 GB in size).

Therefore, smaller data collections (instructions, modules, thread descriptions, etc.) favor ease of access, while potentially larger collections (executed instructions) lean towards fast data access

The trace can be stored either in a typical DB or in home-made binary files. The latter was chosen for a few reasons:

- DB engines limitations:

- Existing DB engines can be pretty slow when it comes to requesting billions of records in a near sequential fashion.
  - The mere overhead of generating/decoding DB requests and copying/marshalling data may be 10 to 100 times higher than data processing itself.
- Custom format processing efficiency:
  - Flat binary files with (known) custom structures can be mapped into the trace analyzer process memory, in whole or in part.
  - A trace analyzer implemented in a language that provides full control over memory layout of structured types (C or C++ for instance) allows direct mapping of implementation language structures over trace data without any copy or marshalling.

Moreover, currently envisaged trace analysis use cases do not perform typical DB operations besides requesting records by their index. Therefore, the lack of DB requests versatility does not warrant reinventing the wheel in trace analyzer implementations. Nevertheless, should the need arise for DB-oriented operations, a new front-end could be added to UMTIndex to generate a proper DB from a raw trace.

## 3 Data types

This section describes data type names and conventions used throughout this document.

Some general considerations:

- Constants
  - Unless otherwise noted, constants in this document are in decimal.
  - Hexadecimal constants are prepended with **0x**, such as **0x1234**.
- Signed-ness
  - Types can be explicitly specified as signed or unsigned.
  - When unspecified, sign does not matter and trace reader implementation is free to use whatever representation seems fit. Only the actual binary value (bits) matters.
  - Some constants in this document may be written as signed/unsigned for the sake of readability. For instance, a 32-bit value written as **-1** means **0xffffffff**. Unless otherwise specified, it does not matter if the trace reader implementation uses it as a signed (**-1**) or unsigned (**4294967295**) value.
- Byte ordering
  - Little-endian (Intel x86 and x64)
  - Low order byte is stored first, and high-order byte is stored last. Value **0x12345678** is stored as 4 bytes: **0x78 0x56 0x34 0x12**.
- Size vs count/length
  - Throughout this document, unless otherwise noted,
    - “size” is in bytes,
    - “count” or “length” is a number of elements in an array or a list,
  - **sizeof(type)** denotes the size of a type, most often a composite type (see 3.4 Structures later in this section).

### 3.1 Basic types

- **i8, i16, i32, i64**: integral values (8 to 64 bits) that can safely be implemented with either signed or unsigned integral values. Such an unspecified signedness type is typically used in two cases:
  - The value is raw binary contents or some ID, and is not used in any arithmetic operation.
  - It represents an unsigned quantity, such as a string length, but the highest half of the value’s range (with the highest bit set, which would otherwise be signed) will never occur. For example, a string length of the type **i64** means it is safe to assume that the value will remain within  $[0, 2^{63}-1]$ .
- **si8, si16, si32, si64**: SIGNED integral values (8 to 64 bits)
- **ui8, ui16, ui32, ui64**: UNSIGNED integral values (8 to 64 bits)
- **bool**: one-byte value that is either 0 or 1.

- **ac**: ASCII character, similar to **i8**.
- **uc**: UTF-16 Unicode “code unit” (not “character”, as some characters are composed using more than one code unit).
  - Similar to **i16**.
  - See <http://msdn.microsoft.com/en-us/library/windows/desktop/dd374081.aspx> about code units.

## 3.2 File “pointers” (offsets)

- **Type \***: Signed 64-bit file offset into the same file. It is relative to a specific point in the file, given in that field’s description. The “type” does not appear in the data itself, and complex types can be omitted in the description to improve readability.

## 3.3 Arrays

- **Type [count]**: Array of **count** elements of type **Type**.
  - All elements are packed without extra padding. When **Type** explicitly contains padding (see **pad** below), this padding is kept, even if it occurs at the end of **Type**.
  - Elements are indexed from **0** to **count-1** inclusively, and found at offset **i \* sizeof(Type)** after element **0**.
  - Total size in bytes of the array is always **count \* sizeof(Type)**.
  - If the element count is not specified in the type description, it must come from another data field.
- **pad[size]**: “size” bytes of padding. Those bytes contain garbage and are meant for alignment purpose only (may improve performance or ease structure definitions).
- **buf[size]**: “size” bytes buffer of undefined type. Within a table, this usually holds several variable-size data pointed to (via file pointers) by other table fields.

## 3.4 Structures

As with many programming languages, structures are tuples whose number of elements and their respective type are fixed. However, trace structure elements may also have a variable size when specified as such. In such cases, the position of following elements, relative to the beginning of the structure, is not statically computable. Also, there is no implicit alignment of elements, on 4 or 8-byte boundaries for instance. Any alignment is explicitly included using a “padding” element.

Structures are defined by their name and size, then by the list of their element. Each element has an offset, type, name and description:

- Offset is in bytes, noted in decimal
  - From the beginning of the structure
  - The offset is always “offset of previous element” + **sizeof(previous element type)**.



- Padding (for alignment purposes) when needed is always explicitly included
- Name is only used as a reference in descriptions of other fields, for the sake of clarity. It does not appear anywhere within the data itself.
- Description may use the following abbreviations:
  - EID: Unique execution record ID (64 bits), as found in the execution trace.
    - Record IDs are unsigned numbers in theory, although they will never reach  $2^{63}$  in supported use cases. So signed 64-bit values will do just fine too.
  - FP: File pointer. Relative to a specific point in the file, given in the description.

## 3.5 Files

Content of the various trace files is described like any structure type.

## 4 Trace format overview

### 4.1 Components

**Execution** table is the main table. It contains all executed instructions of the trace, chronologically ordered (to some extent, see 6.5 Execution order issues). Its records directly or indirectly refer to other tables. Content of those tables is summarized below, and Figure 1 depicts basic relation between all tables.

- **System calls**
  - Extra contextual information recorded at system calls entry and exit points.
- **Instructions**
  - All instructions ever instrumented during trace recording
    - Basic information: module it belongs to, address, code bytes
    - Decoded instruction: explicit and implicit operands, constants, text disassembly
    - Data pattern for recorded values in execution trace
  - Generated from the raw instrumentation trace
- **Modules**
  - DLLs loaded at some point into the traced process
  - Information about them (file name, size), not their actual contents
- **Threads**
  - Basic information on all threads of the traced process
  - Does not include control-flow information, which is found in the execution table

Two tables are not referred to by others, but will tell a trace analyser how to process data:

- **Context map**
  - Register values are sometimes dumped all at once in an execution record as a monolithic “context”
  - This map tells where each register within that context is.
- **Address space usage**
  - List of all memory address (4kb page ranges) ever referred to by executed instructions.

Another file, *trace.xml*, also provides information on how to process data (which “enum tables” to use) as well as some meta-data (e.g. traced process command-line).

Finally, a few “enum” tables, external to the trace, may be needed to trace analyzer tools developers. They provide a somewhat meaningful representation (literal names) of numerical IDs used in other tables (system call IDs, instruction mnemonics, etc.).

More details on each table are presented in their respective section of this document.

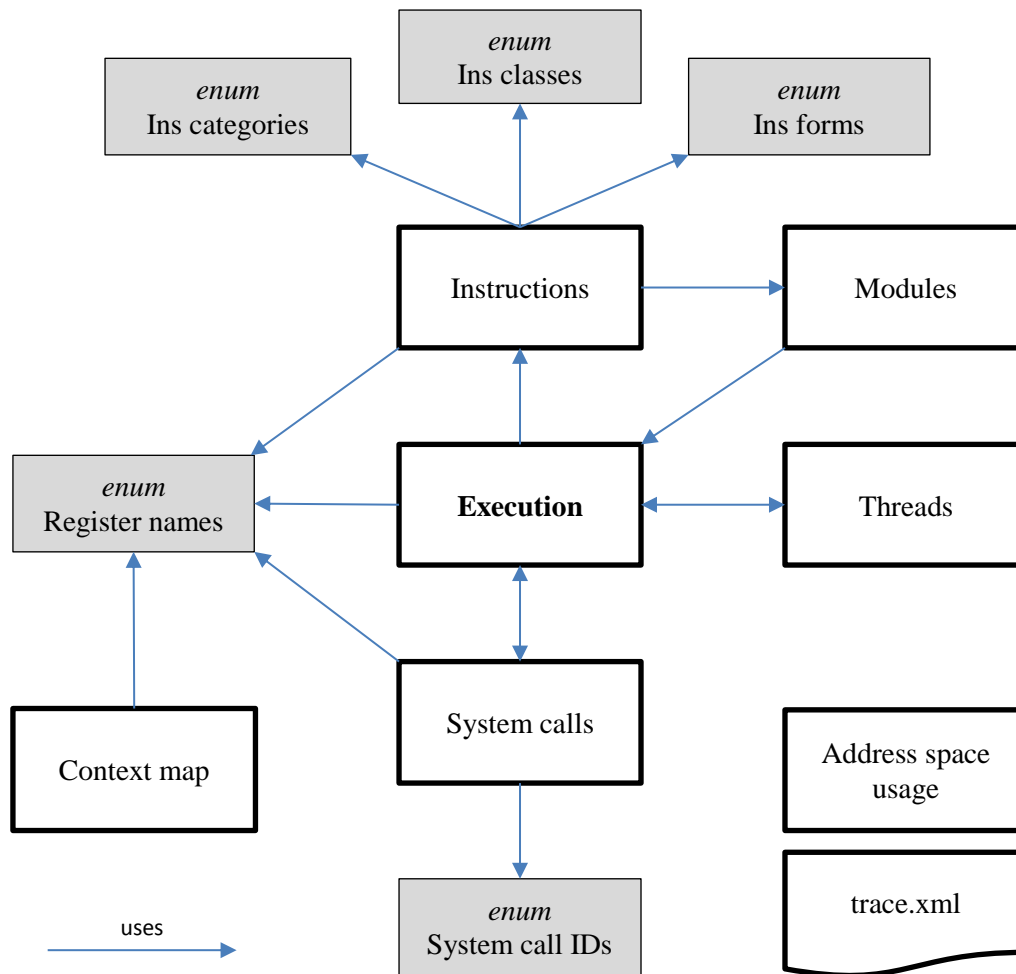


Figure 1 - Trace tables

## 4.2 Generic table file format

Table file format comes in two flavours:

- Self-indexed, fixed-size records
- Variable-size packed records list

### 4.2.1 Fixed-size records

This file format is mostly used by tables that are expected to be small (typically < 1GB) and mapped or loaded all at once in memory for processing. All tables but **Execution** and **System calls** use this format.

The table is merely an array of records, with direct access to any record by its 0-based ID. If variable-size fields are needed, they are located in an extra area at the end of the file and are pointed to from fixed-size fields in the table. Trace tables stored in that format share a common header and the following layout, where **record** is table-specific:

Offset	Type	Name	Description
0	i32	HeaderVer	Header version number. Layout of the remainder of this header depends on its version. Only version 0 is currently supported for all tables.
4	i32	RecordVer	Record layout version number. This refines the <code>record</code> type to use (which type or which version of this type). Most fixed-size record types are still at version 0.
8	i64	RowCount	Number of records in the table
16	i64	RowSize	Size of each record, always <code>sizeof(record)</code> for a given <code>RecordVer</code>
24	record *	TableFP	File position (absolute, from start of file) of the first record (ID = 0) of the table.
32+	...	...	Table-specific additional header information.
TableFP	record [RowCount]	Table	Array of <code>RowCount</code> contiguous records.

#### 4.2.2 Variable-size packed records list

This file format does not have a specific structure. Stream-like records (fields may depend on some header values) are simply collected back to back (ordered or not) in a file. Records are not accessed directly by their ID. A second file provides an offset index where each entry points (absolute file offset) to the corresponding record in the former file.

`Execution` and `System calls` tables use this format.

## 5 Common structure types

A few structure types are not intrinsically related to any table.

**cStr**: null-terminated character (8 bits) string with explicit size

Size: **Length** + 3

Offset	Type	Name	Description
0	<b>i16</b>	<b>Length</b>	String length in characters (bytes), not including trailing null
2	<b>ac[Length+1]</b>		String contents followed by a null character (0)

**uStr**: UTF-16 null-terminated string with explicit size

Size: **Length**\*2 + 4

Offset	Type	Name	Description
0	<b>i16</b>	<b>Length</b>	String length in UTF-16 “code units” (not characters), not including trailing null. The actual number of Unicode characters is always lesser or equal to <b>Length</b> . About “code units”, see <a href="http://msdn.microsoft.com/en-us/library/windows/desktop/dd374081.aspx">http://msdn.microsoft.com/en-us/library/windows/desktop/dd374081.aspx</a>
2	<b>uc[Length+1]</b>		String contents followed by a null character (16-bit 0)

## 6 Execution table

### 6.1 Underlying data format

**Execution** table contains every executed instruction in chronological order (or almost, see 6.5 Execution order issues), i.e. their ID is an execution sequential number. It also contains references to some execution events such as thread creation and system calls. Records therefore have a “type” and some data fields that depend on that type. They also share some common fields.

#### 6.1.1 What is in an execution record

Informally, records contain:

- Thread ID: refers to a thread record in the **Thread** table.
- Thread control-flow: IDs of the previous and next execution records according to the control-flow for the current thread
- Flags: bit flags providing generic information about the execution at that point (32/64-bit, instrumentation turned on/off, etc.)
- Type: instruction, system call entry/exit, etc.
- Type-dependant data.

For most record types, type-dependant data provides a complete “local context”, that is, recorded values for all registers and/or memory locations relevant to that instruction. For example, a system call or context change records provide a full register context dump, and instruction-type execution records contain all operands values, both explicit and implicit. Therefore, in most cases, a trace analyzer does not have to maintain a virtual machine state (registers and memory contents) to analyze a record. In other words, as far as a virtual machine state is concerned, records are pretty much self-contained. This may imply some redundancy and larger table files, but may greatly improve analysis performance and simplicity afterwards.

Still, maintaining a virtual machine state or searching through previous and next records may be required for replay-debug scenarios or thorough system call analysis. In the latter case, in a user-mode trace, system calls are akin to mega-instructions with several operands and probably several side effects as well. Basically, operands are the system call parameters, and those are included in execution records. However, semantically, operands also refer to:

- The various OS-specific structures they might point to in memory (which in turn may contain pointers to other structures, and so on).
- Kernel objects, via handles.

Neither UMTracer nor UMTIndex tools deal with those for now, except for a very few experimental cases that are not described in this document (but for which execution records may contain a few extra values recorded from memory). This would require specific handling for each of the few thousands known

system calls as well as the hundreds (or more) of data structures used to make those calls. Nevertheless, a trace analyzer can implement some workarounds:

- Data structures in memory (system call parameters): if their values are really used by the user-mode calling code, those values will be read or written in instructions preceding or following a call. They will appear in their own record as with any memory access.
- Kernel objects: those are opaque to user-mode code, although in several cases the Win32 native API provides ways to query them. Obviously, this cannot be done offline from a trace. Only the tracer could do it live. The best a trace analyzer could do is to track a handle origin (a call to `NtOpenFile` for instance) and try to figure out what it refers to. Nevertheless, tracking what happens to that object from that point has its limits: handles may refer to objects shared with other processes and internal object state may thus be modified without any hint in the recorded trace.

### 6.1.2 Organizing huge record collections

Records' size varies from one to the other, so they need a "file offset and size" index. Semantically, this looks like the following ("educated random" values):

ID	Offset	Size		Thread	Type	Flags	Prev	Next	Data
0	16	18	→	1	Begin	0x3	none	1	
1	34	24	→	1	Ins	0x3	0	3	
2	58	19	→	2	Begin	0x3	none	4	
3	77	220	→	1	Ins	0x1	1	6	...
4	297	23	→	2	Ins	0x1	2	5	
5	320	17	→	2	Ins	0x1	4	7	
...	...	...		...	...	...	...	...	...

In reality, such an index takes about 25% of the overall table size and Prev/Next together take another 25%, meaning only half of the table contains actual execution data. This would be worse should the ID be stored in the table as well. Since the execution table can be pretty large, a lot of data must be read even though only part of it might be used. The following optimizations were applied to alleviate this issue:

- Remove "Size"
  - Records are stored in order (offset is strictly increasing) and size of a record can be inferred with the offset of the following record.
  - "Offset" column contains an extra offset, pointing right after the last record
- Move "Prev" and "Next" columns into a separate file.
  - Trace analyzers that only rely on sequential execution order have no use of thread control-flow and may thus entirely avoid reading it.

## 6.2 File format

**Execution** table is split into three files, as shown in Figure 2:

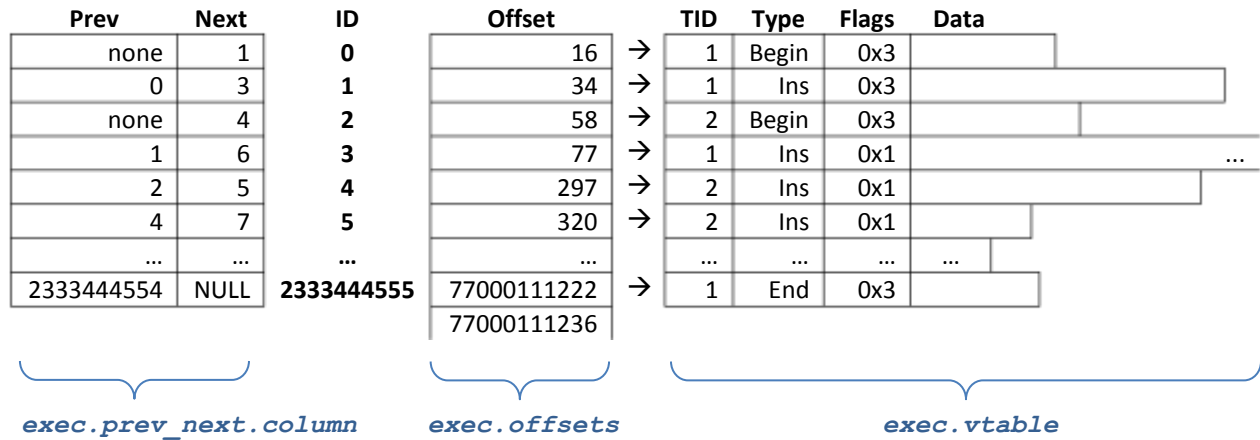


Figure 2 – Execution table files

**exec.vtable**: packed collection of variable-size execution records, less thread CF columns

Size: variable

Offset	Type	Name	Description
0	i32	Version	Version number (see versions in 6.3)
4	pad[4]		
8	i64	RecordCount	Valid IDs for all three files are in [ 0, RecordCount [
Offset[i]	execRec		Variable-length record (see <b>execRec</b> structure below) for each <b>i</b> in [ 0, RecordCount [. Records are sorted by their ID.

**exec.prev\_next.column**: thread CF information, i.e. “Prev” and “Next” columns of the table.

Size: RecordCount \* 16

Offset	Type	Name	Description
0	i64Pair[RecordCount]	PrevNext	Each pair <b>PrevNext[i]</b> holds respectively the IDs of the previous and next records, thread CF-wise, for record <b>i</b> in [ 0, RecordCount [. ID <b>-1</b> is used when there is no previous or next record.

**exec.offsets**: Absolute file offsets into **exec.vtable** for each record.

Size: (RecordCount + 1) \* 8

Offset	Type	Name	Description
0	i64[RecordCount+1]	Offset	Absolute file offsets into <b>exec.vtable</b> for each record, and an extra file offset pointing after the end of the last record. Any record <b>i</b> is always fully contained in <b>exec.vtable</b> within [ Offset[i], Offset[i+1] [.



## 6.3 Versions

- Version 0
  - Initial release version, as described in this section unless otherwise noted
- Version 1
  - A few changes among `execRec.Flags`, see section 6.4.

## 6.4 Execution record format

**execRec**: execution record, combining a common header and variable type-dependent data

Size: variable (at least 6)

Offset	Type	Name	Description
0	i32	ThreadID	Thread ID (in <code>Thread</code> table) this record belongs to.
4	i8	Type	Record type, see type descriptions below
5	i8	Flags	Combination of bit flags, see below
6	?	Data	Type-dependent data, see type descriptions below

`execRec.Flags` can be a combination of the following values:

- **0x01**: 64-bits. Instruction or special event occurred within the context of a 64-bit mode (aka long-mode) code segment. 32-bit otherwise.
- **0x02**: Instruction instrumentation was OFF before execution of this record (see 1 Introduction). When set, context values (register values and perhaps memory references) from previous records cannot be assumed to be correct for the current record. However, context values in current record data are valid for that record, and can be propagated to the next record if it does not have this flag set.
- **0x04**: non-executed conditional instruction. Applies to the following conditional instructions: `cmovcc`, `fcmovcc`, `jcc` (conditional jumps) and `repcc`-prefixed instructions. This flag is set when the condition `cc` was false and the remaining of the instruction was not carried on. All input and (unchanged) output operand values nevertheless appear in the current record data.
- **0x08**: set on the first iteration of a `repcc`-prefixed instruction.
- **0x10 to 0x80**: one or more of those flags are set when the instruction could not be completely traced, either because it caused a low-level exception/fault or the application abruptly terminated (processed killed). Although input and output operands values appear in the current record, they may not be reliable. Meaning of those flags depends on `Execution` record version:
  - **Version 0**
    - **0x10**: exception occurred while recording input operands. Both input and output operands values contain garbage.
    - **0x20**: exception occurred during instruction execution. Input values are properly recorded, but output values are not.

- **0x40**: exception occurred when tracing output memory operands. For some reason, although the instruction executed fine, output operands could not be read back to be recorded.
- **0x80**: internal error, probably a bug in the tracer tool.
- **Version 1**
  - **0x10**: inputs are not reliable. Most often, register values are fine, but memory operand values are probably not. This includes pre-execution value of output-only operands.
  - **0x20**: outputs are not reliable. Can't tell whether output operand values were fully, partially or not at all written to the trace.
  - **0x40**: instruction triggered an exception/fault. Access violation is probably the most common cause.
  - **0x80**: aborted process. Traced application was abruptly terminated while a thread was *around* this instruction (about to execute it, executing it, or just done with it).

Fields in `execRec.Data` depend on `execRec.Type`. The following table list all possible types, their meaning and the fields found in `execRec.Data`. Some fields are structures and those structure types are defined after the table. Multiple fields, if there are, are always packed without any padding.

Type	Description		
	Field type	Field name	Field description
0	<b>Instruction</b>		
	i64	InsID	Corresponding instruction record ID in <a href="#">Instruction</a> table
	buf[?]	OperVals	Values of all instruction operands, both implicit and explicit. Provides the instruction “local context” (see 6.1.1 What is in an execution record). The size and data layout of the buffer is given in the instruction record.
1	<b>Thread begin</b> Exactly once for each thread, before any other record for that thread.		
	rmContext	Context	Registers context at the thread start point, before the first user-mode instruction, as well as the thread environment block (TEB).
2	<b>Thread end</b> Exactly once for each thread that terminated before the traced application exited. This is the last record for a thread.		
	i32	ExitCode	Windows thread exit code
3	<b>Application end</b> The very last record of the trace. Present only if the traced application exited gracefully.		
	i32	ExitCode	Windows process exit code
4	<b>System call entry</b> Thread is about to enter kernel mode through one of the various kind of system calls. Typically follows a <a href="#">syscall/sysenter/int 2e</a> instruction record.		
	i64	SyscallID	Corresponding system call record ID in <a href="#">System calls</a> table.
	rmContext	Context	Value of all registers just before the call. Additional information, such as call parameters on the stack, is in the <a href="#">System calls</a> table.
5	<b>System call exit</b> Thread just exit from kernel mode system call. Does not always immediately follow a system call entry record, and may not even match an entry at all (matching is done according to the stack pointer). Actually, system calls can be nested, involve user-mode callbacks, and can end on an exception or a context change. Also, some system calls, such as <a href="#">NtContinue</a> , can shortcut a few system call exits at once.		
	i64	SyscallID	Corresponding system call record ID in <a href="#">System calls</a> table, or -1 if there is no corresponding system call entry record (thus nothing in <a href="#">System calls</a> table either). The latter case is semantically equivalent to a context change.
	rmContext	Context	Value of all registers right after the call. Additional information may be found in the <a href="#">System calls</a> table.
6	<b>Skipped system call exit</b> A previous system call entry occurred, but starting from this record, the stack was unwound further than it would be on a normal system call exit (a “return” to the user-mode caller). In other words, the thread is not anymore in the context of that system call, and that call will never have a “System call exit” record.		
	i64	SyscallID	Corresponding system call record ID in <a href="#">System calls</a> table.
8	<b>Context change – Asynchronous Procedure Call</b>		
9	<b>Context change – Exception Handling</b>		
10	<b>Context change – Callback</b>		
11	<b>Context change – Unknown reason</b> For one of the above reasons, the thread execution context just changed.		
	rmContext	Context	Value of all registers right after the context change, before carrying execution of the next instruction.

**rmContext**: full register context (all register values) and optionally some relevant values from memory.

Size: variable (at least 16)

Offset	Type	Name	Description
0	i32	MemCount	Number of elements in <b>MemRefs</b>
4	pad[4]		
8	buf[?] *	RegsFP	Points to a buffer that holds all registers value. Size and data layout of this buffer is given by the <b>Context map</b> (see section 11). FP is relative to the start of this structure (its position in the file)
16	contextMem [MemCount]	MemRefs	Array of optional memory references added to the context.

**contextMem**: one memory reference, part of a **rmContext**.

Size: 24

Offset	Type	Name	Description
0	i64	Address	Memory address (high 32 bits are always 0 on 32-bit traces)
8	i64	Size	Size of the memory buffer
16	buf[Size] *	DataFP	Points to memory content. FP is relative to the start of the parent <b>rmContext</b> structure (its position in the file)

**Note:** Data pointed to from those structures is stored right after **rmContext** instance. That way, each record remains a monolithic stream of bytes. Any record **i** is always fully contained in **exec.vtable** within [ **Offset[i]**, **Offset[i+1]** ], where **Offset** is found in **exec.offsets**.

## 6.5 Execution order issues

The actual execution order that took place while tracing the application might slightly differ from recorded order. Instrumenting an instruction in the original application involves “recompiling” together that instruction and additional tracing code. A single instruction then turns into several instructions and in some cases, the original instruction itself is broken apart into a few instructions. A thread may therefore be pre-empted anywhere “during” an instrumented instruction.

Also, the tracer records values from memory before or after the actual instruction reads/writes that memory. Should that memory value be changed by another thread in between, there may be a discrepancy between the recorded value and the value used by the application. This is a rare event that is likely to occur only when tracing thread synchronization code (getting a lock for instance) or if there is some thread synchronization bug in the traced application. While this could be avoided in the tracing tool (*Pin* provides required functionalities), no efforts have been put so far regarding that issue. Should we ever need such accuracy (to track multi-thread bugs in applications for instance) we might improve the tracer accordingly.

Neither of those issues affects the binary file format. They may only produce wrong memory references values in very few execution record data.

## 7 Instruction table

### 7.1 Underlying data format

**Instruction** table is referred to by most execution records and is needed to do pretty much anything useful with those execution records. Instruction records are organized as an array of fixed-size records and two non-indexed collections of variable-size sub-records pointed to by fields from the fixed-size part:

- Fixed-size part: static instruction information
  - Actual instruction: module it is from, address, code bytes, pointers to corresponding data in the two other sub-tables.
  - Recorded by the tracer as instructions are instrumented:
    - Not from some module-wide static code analysis.
    - Every instruction that is potentially about to be executed is instrumented, including self-modifying code.
    - Only instrumented instructions are recorded.
  - Instructions can be instrumented more than once for various reasons (namely the user turning instrumentation on an off during recording) and therefore recorded more than once into the table, on two distinct unrelated entries.
- Variable-size part: decoded instruction
  - Complete description of each instruction
    - Explicit and implicit operands (name and size)
    - Immediate values
    - Semantics revealed through enum values (instruction category, mnemonic and specific form), imported from Intel *XED* (see section 14.1) and subject to version change over time.
  - Data layout of operand values in execution record (the local context, see 6.1.1 What is in an execution record)
  - The same variable-size sub-record can be pointed to by several identical instructions from their respective fixed-size sub-record.
- Variable-size part: disassembly
  - Textual disassembly of each instruction, in typical Intel format.
  - Identical strings are merged, i.e. many instructions may point to the same string.

### 7.2 Decoded instruction

Intel instruction set is large, versatile and inherently complex. So is any structure that aims to uniformly describe all possible instructions. Figure 3 illustrates the basic relations between the above three sub-tables and execution records, with an emphasis on the “decoded instruction” sub-table. For the sake of clarity, some record fields were omitted or informally summarized.

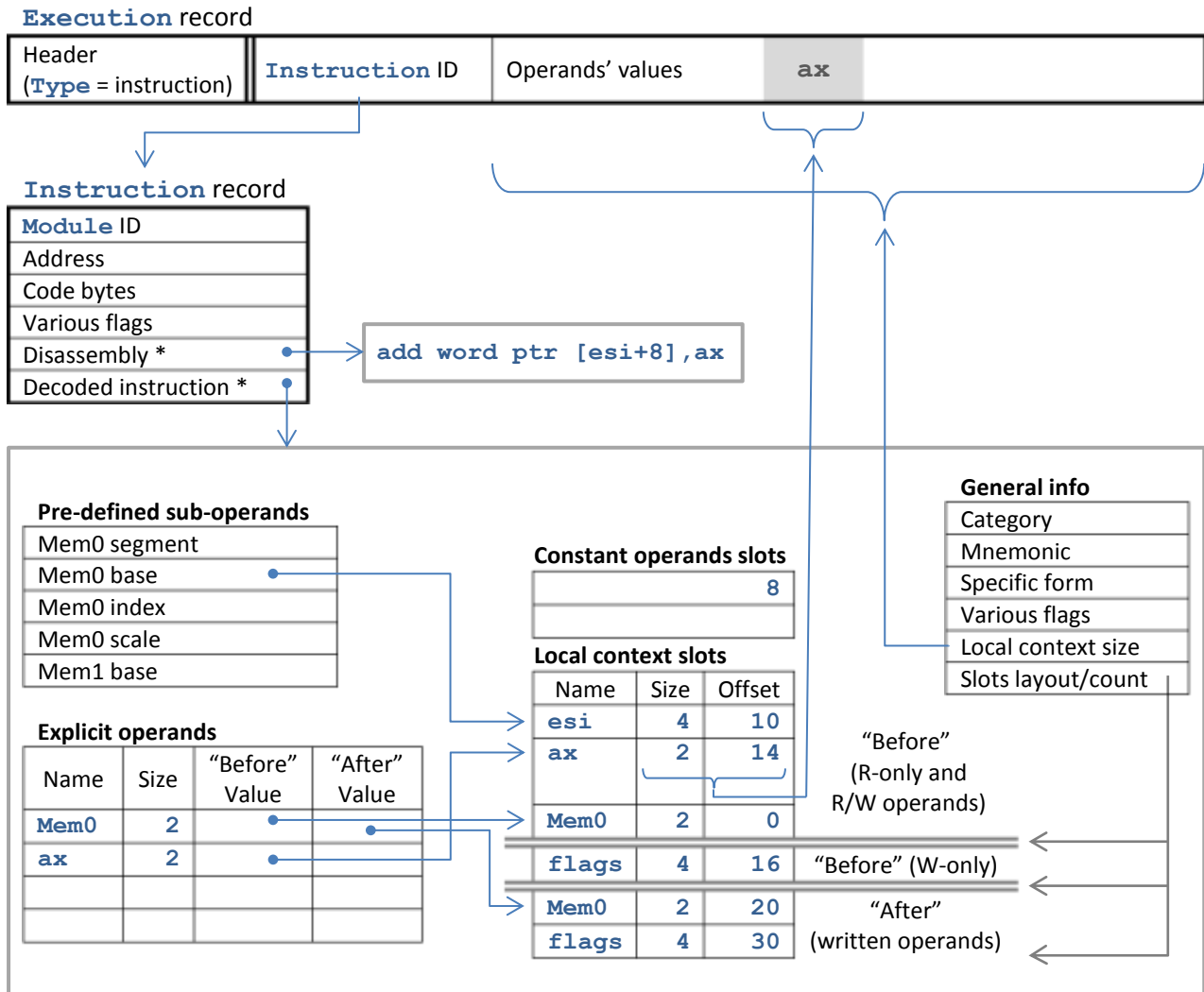


Figure 3 – Simplified decoded instruction example

### 7.2.1 Local context slots

Local context slots are central to a decoded instruction. They list exactly all registers and memory references that are potentially used and/or modified. This also corresponds to the execution record "local context" buffer data layout. Slots are stored in a contiguous variable-size array and are grouped according to whether they refer to values before or after instruction execution. Pre-execution values of write-only operands are provided as a convenience to a trace analysis tool, although they are not actual inputs for that instruction. Slots referring to those values form a third slot group, between the two others. Besides this grouping, there is no particular order among slots.

Overlapping or "contiguous" operands can be merged in the execution record and sometimes in the slots as well (yet duplicates are allowed). For instance, if an instruction reads both `al` and `eax`, only `eax` value will be stored in the execution record and most often, only `eax` will appear in the slot list as well. Therefore, local context slots do not provide a clear picture of explicit instruction operands. This is

instead provided by another level of indirection: explicit operands and pre-defined sub-operands (see the corresponding tables in Figure 3).

Explicit operands are just that: those that appear textually in disassembly or Intel documentation, in the same order. The current instruction set supports up to four explicit operands. Memory operands, whether explicit (`mov [rax], rdx`) or implicit (such as `[rdi]` in `stos` for instance), also have sub-operands in their “memory phrase” `seg:[base + scale*index + offset]`. Only one full memory phrase is supported per instruction, and the second memory operand, when there is, only has a base register. Each operand and sub-operand in use refers to one of the local context slots (including the two “constant” slots).

An explicit operand may refer only to a portion of the value the corresponding slot refers to. For instance, in `mov byte ptr [eax], al`, a single slot will be allocated for `eax` with a value size of 4, while the second explicit operand will point to that slot but will specify a size of 1. Explicit operands may also specify an additional offset into the value referred to by the slot. This is only required when an operand is `ah`, `bh`, `ch` or `dh` register and the corresponding slot refers to a larger enclosing register.

A trace analysis tool operating on instruction semantics would query instructions via its explicit operands and memory phrase sub-operands. But if only the list of input and output operands (or their values) is required, the tool may go directly for the local context slots.

### 7.2.2 Constant operands

Only explicit constants are found in the instruction record. Implicit constants (such as stack pointer offset on `push` or `pop`) are part of the instruction semantics and/or can be inferred from some fields of the instruction record (effective address or operand width, for instance). Constants do not need to be stored in the execution record, and are stored instead in two reserved operand slots.

If one of the explicit operands is a memory reference, the first constant slot is reserved for the memory phrase offset. Explicit constant operands (immediate values and branch offsets) simply point to either slot. In all cases, IP-relative offsets (in relative branches or 64-bit RIP-relative addressing) are converted to absolute constants using the address of the current instruction.

### 7.2.3 Memory operands

Operand name is a simple enum value. Most names relate to real registers (see section 14.1) but a few refer to memory operands: `MEM0`, `MEM1`, `STACKPUSH`, `STACKPOP`. For these operands, the memory address is stored in the execution record along with the value. This alleviates the need to compute it from its constituents (memory phrase, stack pointer, etc.) although their values are all stored in the execution record as well. Note that memory-phrase/stack registers always have their “pre-execution” value (an address) stored in the execution record. Also, addresses in `FS` or `GS` segments are converted to flat linear addresses. Base address of each segment is provided in the local context as a convenience, but is not need to get the memory-based operand actual address.

Another special operand name, **AGEN**, is used solely for the **lea** instruction. It uses the same memory phrase sub-operands than **MEMO**, but does not have an associated memory value. Only the “address” is stored in the execution record, just like register values.

## 7.3 Instruction table file format

### 7.3.1 Generic table format

The **Instruction** table is stored in **ins.itable**. Its header follows the fixed-size table format described in section 4.2.1, where **record** type is **insRec**, described below. It also has the following table-specific fields in its header:

Offset	Type	Name	Description
32	<b>buf[ExtraSize] *</b>	<b>ExtraFP</b>	Absolute file pointer to the start of the two collections of all variable-size sub-records (decoded instruction and disassembly).
40	<b>i64</b>	<b>ExtraSize</b>	Total size of those two collections

Data at **ExtraFP** is not readily usable. It is rather accessed via instruction records, found at **Table[i]**.

### 7.3.2 Versions

**RecordVer**, from the table header, can take the following values:

- Version 0
  - Initial release version, identical to latest version described in this section unless otherwise noted
- Version 1
  - Semantic change only, i.e. binary compatible with version 0. **lcSlot.LargestEnclosing** was changed to **lcSlot.CommonRegName**. See details in structure file format below. While it cannot be used anymore to get the “largest known enclosing register” of a given register, **CommonRegName** remains coherent with **CommonName** table in the context map (see section 11) just like with version 0.

### 7.3.3 Data structures

In the following structures, a “local context slot index” (LC slot index) can be either an index into **decodedIns.LCSlots** array or one of the following special values:

- **SLOT\_IMM0 (0xfd)**: value is found in **decodedIns.ConstSlot[0]**
- **SLOT\_IMM1 (0xfe)**: value is found in **decodedIns.ConstSlot[1]**
- **SLOT\_NA (0xff)**: no corresponding slot, i.e. does not exist in local context



Also, operand names can either be physical register names (enum constants from Intel *XED*, **subject to version change over time**, see section 14.1) and/or one of the following special values:

- **MEM0** (**0xffff**): first memory-based operand of an instruction, referenced by a memory phrase of the form `seg:[base + scale*index + offset]`, where most components are optional.
- **MEM1** (**0xfffe**): second memory-based operand of an instruction
- **AGEN** (**0xfffd**): address generator: integral result of a memory phrase, without reference to memory at that “address”.
- **STACKPUSH** (**XED\_REG\_STACKPUSH**): output operand is on the stack at address `rsp/esp/sp - n`, where `n` depends on the operand size and pre-execution stack pointer value is used. Local context value in execution record provides the computed address like it does for any memory operands.
- **STACKPOP** (**XED\_REG\_STACKPOP**): input operand is on the stack at address `rsp/esp/sp`. Local context value in execution record provides this address like it does for any memory operand.
- **FS\_BASE** (**XED\_REG\_FSBASE**) and **GS\_BASE** (**XED\_REG\_GSBASE**, from latest *XED*): in a user-mode trace, it would be pretty useless to provide `fs/gs` selector values when those are used as segment overrides in memory phrases. Instead, their base linear address is provided.
- **INVALID** (**XED\_REG\_INVALID**, likely to always remain **0x0000**): no associated register (to name a constant or non-existent operand for instance).

**insRec**: fixed-size part of an instruction record

Size: 64

Offset	Type	Name	Description
0	i64	ID	Self-ID in the current table.
8	i8	Type	Record type, provides a hint to the analyzer: <ul style="list-style-type: none"> <li>• 0: normal instruction</li> <li>• 1: software interrupt (except alternate system call)</li> <li>• 2: 32-bit fast system call (<b>syscall/sysenter</b>)</li> <li>• 3: 64-bit fast system call (<b>syscall/sysenter</b>)</li> <li>• 4: alternate system call (<b>int 0x2e</b> on Windows)</li> <li>• 5: WOW64 system call</li> <li>• 6: other unknown system call</li> <li>• 8: NOP <ul style="list-style-type: none"> <li>○ There are several ways with various mnemonics to encode a “no operation”. All of them have this flag set.</li> </ul> </li> <li>• 9: Prefetch <ul style="list-style-type: none"> <li>○ Prefetched values are not recorded by the tracer and corresponding execution record operand values are likely to contain garbage.</li> </ul> </li> </ul>
9	bool	CSis64Bits	1 if the instruction was executed from a “long mode” 64-bit code segment, 0 otherwise.

10	i8	CodeSize	Number of valid bytes in <a href="#">CodeBytes</a> array.
11	pad[1]		
12	i32	ModuleID	Module ID (in <a href="#">Module</a> table) this instruction is from, or <a href="#">-1</a> if this instruction was from floating code in memory.
16	i64	Address	Address where this instruction was found (zero-extended on 32-bit applications). Note that addresses are not unique: <ul style="list-style-type: none"> <li>• Instructions can be instrumented more than once.</li> <li>• Different code can be loaded/unloaded at a given address.</li> </ul>
24	buf[15]	CodeBytes	Actual bytes of the instruction. Only the first <a href="#">CodeSize</a> bytes are valid.
39	pad[1]		
40	cStr *	DisasmFP	Disassembly string of instruction, Intel style. FP is relative to <a href="#">ExtraFP</a> in the file header.
48	decodedIns *	DecodedFP	Fully decoded instruction. FP is relative to <a href="#">ExtraFP</a> in the file header. See <a href="#">decodedIns</a> below.
56	pad[8]		

[decodedIns](#): main (and somewhat monstrous) variable-size part of an instruction record

Size: variable (at least 88)

Offset	Type	Name	Description
0	i16	Category	Instruction category, from <a href="#">xed_category</a> enum table (subject to version change over time). See section 14.1.
2	i16	Mnemonic	Instruction mnemonic, from <a href="#">xed_iclass</a> enum table (subject to version change over time). See section 14.1.
4	i16	Form	Specific instruction form for the given mnemonic, from <a href="#">xed_iform</a> enum table (subject to version change over time). See section 14.1.
6	i8	Flags	Bitwise combination of zero or more of the following: <ul style="list-style-type: none"> <li>• <a href="#">0x01</a>: has <a href="#">rep</a> or <a href="#">repe/repz</a> prefix</li> <li>• <a href="#">0x02</a>: has <a href="#">repne/repnz</a> prefix</li> <li>• <a href="#">0x04</a>: has <a href="#">lock</a> prefix</li> <li>• <a href="#">0x08</a>: instruction executes in 64-bit long mode, but operates on 32-bit register operands, in which case all output operands have their highest 32 bits zeroed</li> <li>• <a href="#">0x10</a>: instruction potentially reads from memory</li> <li>• <a href="#">0x20</a>: instruction potentially writes to memory</li> </ul>
7	i8	ExpCount	Explicit operand count (0 to 4 incl.)
8	expOper[4]	ExpOper	Array of explicit operands. Only the first <a href="#">ExpCount</a> elements are valid. See <a href="#">expOper</a> below.
56	i8	EffOperSize	Effective operand width in bytes (2, 4 or 8). Not the actual operands width, but rather related to instruction encoding and some zero/sign-extend rules. See Intel manual about “operand-size and address-size attributes”.
57	i8	EffAddrSize	Effective address size in bytes (2, 4 or 8). Not the actual size of code addresses, but rather related to instruction encoding and some zero/sign-extend rules. See Intel manual about “operand-size and address-size attributes”.
58	i16	LCSize	Size of local context data in execution record.

60	i8	FirstOutSlot	Index of the first LC slot for “before execution” value of write-only operand (see Figure 3), thus also the number of “before execution” values slots for read-only and read-write operands.
61	i8	FirstAfterSlot	Index of the first LC slot for after-execution operands (see Figure 3), i.e. after before-execution write-only operand slots.
62	i8	LCCount	Total count of LC slots (see Figure 3). Length of <b>LCSlots</b> array.
63	i8	Mem0SegSlot	LC slot index of the memory phrase segment override for any operand named <b>MEM0</b> or <b>AGEN</b> .
64	i8	Mem0BaseSlot	LC slot index of the memory phrase base register for any operand named <b>MEM0</b> or <b>AGEN</b> .
65	i8	Mem0Scale	Scale factor (1, 2, 4 or 8) of the memory phrase base register for any operand named <b>MEM0</b> or <b>AGEN</b> , or 0 if there is not an index register as well.
66	i8	Mem0IndexSlot	LC slot index of the memory phrase index register for any operand named <b>MEM0</b> or <b>AGEN</b> .
65	i8	Mem1BaseSlot	LC slot index of the memory phrase base register for any operand named <b>MEM1</b> .
68	pad[4]		
72	i64[2]	ConstSlot	Two constant slot values. Those are not <b>lcSlot</b> structures but only the actual value of an operand or memory phrase offset. If a <b>MEM0</b> or <b>AGEN</b> operand is present in this instruction, <b>ConstSlot[0]</b> contains the memory phrase offset, and the number of significant bytes (up to 8) is given by <b>lcSlot.AddrWidth</b> of the corresponding operand slot. For explicit operands that directly refer to a constant slot, the number of valid bytes is given by the operand size.
88	lcSlot [LCCount]	LCSlots	Array of all local context slots, see <b>lcSlot</b> below.

**expOper**: description of an instruction explicit operand

Size: 12

Offset	Type	Name	Description
0	i16	RegName	Register name (from Intel <i>XED</i> , see section 14.1) or one of the special values (see top of this section on page 24). May be different than <b>RegName</b> of the corresponding before/after slots since the latter may refer to a larger enclosing register. For instance, in <b>mov al, [ebx+eax]</b> , the first explicit operand may refer to an input slot that contains <b>eax</b> (no need to have <b>al</b> as a separate input slot).
2	i16	Size	Operand value size in bytes. Lesser or equal to <b>ValueSize</b> of the corresponding before/after slots. Even for memory operands, this is the size of the value in memory, not counting the address itself stored in the execution record. When <b>RegName</b> is <b>AGEN</b> , this is the size of the computed address. Note: while <b>Size</b> is most often a register size (up to 64 bytes for ZMM registers, some instructions may have memory operands of nearly 2700 bytes ( <b>xsave/xrstor</b> and relatives for instance).
4	i8	Flags	Combination of at least one of <ul style="list-style-type: none"> <li>• <b>0x01</b>: operand is read</li> <li>• <b>0x02</b>: operand is written</li> </ul>

5	i8	BeforeSlot	LC slot index for this operand “before execution” value. Never <b>SLOT_NA</b> .
6	i8	AfterSlot	LC slot index for this operand “after execution” value. <b>SLOT_NA</b> for read-only operands.
7	i8	BeforeOffset	Byte offset to add to <b>LCSlots[BeforeSlot].Offset</b> to get the actual offset, in the execution record local context, where the operand “before execution” value is stored. For register operands, this is always 0, except for <b>ah, bh, ch dh</b> registers and only if the corresponding slot refers to a larger enclosing register (for instance, operand is <b>ah</b> but slot refers to <b>ax</b> , hence an additional offset of 1). For memory operands, this is always 8 since a 64-bit address is stored before the actual memory value, even on 32-bits traces.
8	i8	AfterOffset	Same as <b>BeforeOffset</b> , for “after execution” operand value.
9	pad[3]		

**lcSlot**: local context slot, defines a portion of the local context data layout in an execution record

Size: 12

Offset	Type	Name	Description
0	i16	RegName	Register name (from Intel <i>XED</i> , see section 14.1) or one of the special values (see top of this section on page 24).
2	i16	CommonRegName  ( <b>insRec</b> version >= 1, see section 7.3.2 for previous version)	<b>RegName</b> can be a subset or a superset of other registers (e.g. <b>al, ah, ax, eax, rax</b> among themselves). A <i>common name</i> is given to all registers of the same <i>group</i> , as defined in section 11.1. This comes handy for a trace analysis tool looking for a specific register in the LC (e.g. looking for <b>al</b> when the LC actually contains <b>eax</b> ). It can search for matching minimum size and common name (given by the <i>context map</i> , see section 11), instead of looking for all possible enclosing registers ( <b>al, ax, eax, rax</b> in the above example).
4	i16	ValueSize	Operand value size. For memory operands, this is the size of the value in memory, not counting the address itself stored in the execution record. When <b>RegName</b> is <b>AGEN</b> , this is the size of the computed address. Note: while <b>Size</b> is most often a register size (up to 64 bytes for ZMM registers), some instructions may have memory operands of nearly 2700 bytes ( <b>xsave/xrstor</b> and relatives for instance).
6	i16	Offset	Offset, relative to the start of the execution record LC data, where the register/memory value ( <b>ValueSize</b> bytes wide) is found. For memory operands, a 64-bit address (zero-extended from original if needed) precedes the value.
8	i8	AddrWidth	0 for non-memory operands. Otherwise, the significant size of the address value in the execution record LC. Addresses in the execution record are always 64-bit (zero-extended if needed), but only the first <b>AddrWidth</b> bytes are significant.
9	pad[3]		

## 8 System calls table

System calls are taken care of in the **Execution** table. There is however some data pertaining to a system call as a whole (or to both corresponding “entry” and “exit” execution records), namely the “name” and arguments (count and values), among others, which are included in the **System calls** table.

### 8.1 Name and argument count

Internally, *Windows* system calls are made through a common entry point, passing a numerical ID as the kernel function identifier, which is then dispatched via a branch table. Branch tables are OS version-specific but corresponding kernel functions are also exported by name. The latter is more consistent across various *Windows* versions since those functions are part of the Native API. Nevertheless, from a version to another, some functions are added, removed or have their argument changed. The argument count is somewhat exported as well, although there is no easy reliable way to automatically know how many arguments have the calls on 64-bit *Windows* with 4 or less arguments.

For the sake of trace analysis, function name and argument count (where possible) from all supported OSes’ branch tables (see *supported OSes.docx*), were automatically extracted and merged into a unified table. Functions with different name or argument count were manually given different unique IDs. The result is found in *syscall\_unified\_id\_v\*.h*, provided with this documentation, where *\** is either 0 or 1 for now, the only two supported versions. Which version to use is given by the **System calls** table header.

Obviously, as newer OS versions are added in the future, this list is likely to change, and new versions will be generated. However, IDs are only added, never removed from a version to the next. Therefore, it is always possible to convert IDs from a previous version to the latest known version.

Most enumerated names are constructed on the following pattern:

**SYSCALL\_ExportedName\_MMmm\_N**

Where **MM** and **mm** are respectively major and minor version of the earliest *Windows* (among supported versions) to provide this function, and **N** is the argument count. For 64-bit OS calls with 4 or less arguments, argument count is assumed to be equal to the 32-bit version argument count, which is known. Should there not be a 32-bit equivalent, **N** is written as **x64**. Note that **N** might be wrong in some cases, the *Native API* being mostly undocumented. Some functions with a different argument count may also have the same exported name, even on the same *Windows* version. Other bits of text were added to a few names to resolve conflicts between OS versions.

A trace analysis tool that does system call specific analysis should rely on that uniform ID in each system call record (and the developer can rely on the enumeration name), although the original internal call ID is also provided as a convenience.

## 8.2 Argument values

Argument values are available from registers and/or the stack, depending on the system call convention used, and found in the few execution records preceding the call. To make things easier for the analysis tool, argument values are properly extracted and copied into a simple array in system call records.

## 8.3 Extras

Several system calls use kernel objects handles as arguments. Those objects are opaque to user-mode code, although in several cases the Win32 native API provides way to query them. Such queries must be done online, at trace time. For now, UMTracer records such extra information for a handful of functions and UMTIndex forwards it to system call records. However, this is experimental code, format for those extras is far from definitive and its description is beyond the scope of this document. See section 6.1.1 for more information and minimal workarounds.

## 8.4 File format

**System call** table is a collection of variable-size records, split into two files.

***syscall.vtable***: packed collection of variable-size system call records.

Size: variable

Offset	Type	Name	Description
0	i32	Version	Format version number (only version 0 is currently supported)
4	i32	UnifiedIDsVer	Uniform ID enum version (only version 0 and 1 are currently supported). See section 8.1 above.
8	i64	RecordCount	Valid IDs for both files are in [ 0, RecordCount [
Offset[i]	syscallRec		Variable-length record (see <b><i>syscallRec</i></b> structure below) for each <i>i</i> in [ 0, RecordCount [. Offset and size are given in <b><i>syscall.offsets</i></b> . Records in the file are <i>not</i> sorted by their ID.

***syscall.offsets***: Absolute file offsets and size for each record in ***syscall.vtable***.

Size: RecordCount \* 16

Offset	Type	Name	Description
0	i64Pair[RecordCount]	(Offset, Size)	Each pair holds the absolute file offset and size of the corresponding record in <b><i>syscall.vtable</i></b> .

**syscallRec**: fixed-size part of an instruction record

Size: 64

Note: All data pointed to (via file pointers) by this record is stored consecutively so that the whole record (including pointed data) fits in a single block of size **Size[i]**.

Offset	Type	Name	Description
0	i16	RawID	OS version-specific function ordinal. Relevant only when combined with the OS version the trace was created on.
2	i16	UniformID	OS version-independent kernel function ID. See section 8.1 above.
4	i8	Convention	System call convention used to make that call: <ul style="list-style-type: none"> <li>• 0: Unknown</li> <li>• 1: 32-bit fast system call (<b>syscall/sysenter</b>)</li> <li>• 2: 64-bit fast system call (<b>syscall/sysenter</b>)</li> <li>• 3: alternate system call (<b>int 0x2e</b> on Windows)</li> <li>• 4: WOW64 system call</li> <li>• 5: software interrupt (<b>int</b> instruction, except alternate system call)</li> </ul>
5	bool	Completed	1 if the system call returned like a normal function ( <b>EndEID</b> is valid), 0 otherwise. In the latter case, either the call never completed before the end of the thread or trace ( <b>EndEID</b> = -1) or its exit was shortcut by any kind of context jump ( <b>EndEID</b> is valid). See system call related execution in section 6.4.
6	i8	ArgCount	Number of arguments in array pointed to by <b>Args</b> . Always 4 for 64-bit calls with an unknown argument count between 0 and 4.
7	i8	ExtraCount	Number of “extra” items in array pointed to by <b>ExtrasFP</b> . Undocumented, see section 8.3 above.
8	i64	BeginEID	System call entry execution record ID.
16	i64	EndEID	System call exit execution record ID. It is -1 if the thread was still in the context of the call when it stopped (or where the tracer stopped recording). Otherwise, it is either a normal (CF-wise) system call execution record, or a skipped system call exit ( <b>Completed</b> = 0, see system call related execution in section 6.4).
24	i64	BeforeSP	Value of stack pointer ( <b>esp/rsp</b> ) before the call. Zero-extended for 32-bit traces.
32	i64	Result	Return value in <b>eax/rax</b> at exit point, zero-extended for 32-bit traces. Invalid (garbage) if <b>Completed</b> = 0.
40	i64[ArgCount] *	ArgsFP	FP on the array of <b>ArgCount</b> argument values. FP is relative to the start of this <b>syscallRec</b> record. Arguments of 32-bit calls are zero-extended to 64 bits. Arguments are in left-to-right order of exported kernel function prototypes.
48	buffer	ExtrasFP	FP on the array of <b>ExtraCount</b> “extras” items. FP is relative to the start of this <b>syscallRec</b> record. Undocumented, see section 8.3 above.

## 9 Module table

Modules (.exe and .dll) content is not stored directly in the trace. It rather appears in the trace as it is used, either through memory references in executed instructions, or instruction code bytes themselves. Optionally, UMTracer can dump in the trace all files mapped into memory (this is how modules are typically imported into the process). Module initial contents will then appear as a simple memory block dumped along with the registers context on system call exit execution record. This is not complete since some modules are already mapped into memory when the tracer starts, and this can be rather inefficient. For example, in a file explorer window, the same few MBs in size DLL can be mapped and unmapped for each visible file.

The **Module** table is stored in `module.itable`. Its header follows the fixed-size table format described in section 4.2.1, where **record** type is `moduleRec`, described below. It also has the following table-specific fields in its header:

Offset	Type	Name	Description
32	<code>buf[ExtraSize] *</code>	<code>ExtraFP</code>	Absolute file pointer to a buffer containing collections of strings pointed to by module records.
40	<code>i64</code>	<code>ExtraSize</code>	Total size of the buffer pointed to by <code>ExtraFP</code> .

**moduleRec**: fixed-size part of an instruction record

Size: 48

Offset	Type	Name	Description
0	<code>i64</code>	<code>LoadEID</code>	Load time, expressed through the approximate ID of the execution record about to be recorded at that point.
8	<code>i64</code>	<code>UnloadEID</code>	Unload time, expressed through the approximate ID of the execution record about to be recorded at that point. Set to <code>-1</code> if the trace recording ends before module is unloaded.
16	<code>i64</code>	<code>StartAddr</code>	Memory address where the module was loaded
24	<code>i64</code>	<code>EndAddr</code>	Memory address of the last module byte (inclusive)
32	<code>uStr *</code>	<code>WinNameFP</code>	FP to full path and file name, in usual Windows format, such as <code>C:\Windows\System32\mf.dll</code> FP is relative to start of file (absolute FP)
40	<code>uStr *</code>	<code>DeviceNameFP</code>	FP to Low-level device path, such as <code>\Device\HarddiskVolume3\Windows\System32\mf.dll</code> FP is relative to start of file (absolute FP)



## 10 Thread table

The **Thread** table is stored in *thread.itable*. Its header follows the fixed-size table format described in section 4.2.1, where **record** type is **threadRec**, described below.

**threadRec:**

Size: 48

Offset	Type	Name	Description
0	i32	ID	Self-ID in the current table.
4	i32	WinTID	Windows thread ID while it was running.
8	i64	TIBAddr	Address of Thread Information Block. TIB content is recorded along with registers context in “thread begin” execution records (see section 6.4).
16	i64	FirstEID	First execution record for that thread, normally of type “thread begin”.
24	i64	LastEID	Last execution record for that thread, normally of type “thread end”, or <b>-1</b> if the tracer was interrupted while the thread was still running.
32	i64	RecCount	Number of execution records that belong to that thread.

## 11 Context map

### 11.1 Registers

Some execution records contain a registers context dump. All register values are stored in a contiguous buffer. The **Context map** defines its data layout, i.e. where each register is given a register name. Although the current version of UMTracer does not generate mixed 32-bit and 64-bit code traces, the **Context map** supports both of them at once. It also supports registers that may not exist on the CPU the trace was made, although those are marked as “not present”.

Incidentally, this map also provides generic static information about registers and relations among themselves. Namely, it gives the *largest enclosing register* for each register as well as a *common register name* for all overlapping registers.

The largest enclosing register for a given register is taken from available registers on the CPU and mode (32/64-bit) the trace was made. For instance, the largest register enclosing **al** is **eax** on 32-bit traces, but **rax** on 64-bit traces. Enclosures are:

- For each general purpose registers, as well as flags register and instruction pointer:
  - 64-bit reg  $\supset$  32-bit reg  $\supset$  16-bit reg  $\supset$  8-bit reg
  - Possible 8-bit reg includes **ah**, **bh**, **ch** and **dh**
- x87 registers (80-bit)  $\supset$  MMX registers (64-bit)
  - In reverse order, by Intel hardware design: **st0**  $\supset$  **mm7**, **st1**  $\supset$  **mm6**, ... **st7**  $\supset$  **mm0**
- ZMM registers (512-bit)  $\supset$  YMM registers (256-bit)  $\supset$  XMM registers (128-bit)
  - Whichever the largest supported on the traced CPU
- For all others, including pseudo-registers, the largest enclosing register is that register itself

Common register names are independent of the traced CPU architecture and 32/64-bit mode. The current mapping is:

- Legacy 16-bit general purpose registers, including flags and instruction pointer:
  - All their relatives map to the 16-bit name, e.g. **al**, **ah**, **ax**, **eax**, **rax**  $\rightarrow$  **ax**
  - This includes registers added in the x64 extension, i.e. **sil**  $\rightarrow$  **si** and **bpl**  $\rightarrow$  **bp**
- GPRs added in the x64 extensions use their 64-bit version name, e.g. **r81**, **r8w**, **r8d**, **r8**  $\rightarrow$  **r8**
- Each MMX registers map to enclosing x87 registers
  - In reverse order, as with above enclosing registers, i.e. **mm7**  $\rightarrow$  **st0**, **mm6**  $\rightarrow$  **st1**, ...
- XMM, YMM and ZMM registers map to corresponding XMM
- For all others, including pseudo-registers, the common register name is that register itself

Previous versions of UMTracer (using **insRec** version 0) used a different mapping: all registers would map to their largest *known-at-that-time* register, which could change over time, as newer tools may *know* more registers. Analysis tools relying on the mapping provided here, instead of hardcoded constants, won't require any changes. But the new mapping may ease implementation since it is stable: it may only be augmented as new registers are introduced, but never otherwise modified, as long as unrelated registers are not merged into larger registers (this never happened so far in x86 history).

## 11.2 File format

The file format has been designed for quick lookups, without the need for bounds checking. Per register information is indexed by register *name*, i.e. its ID from *Intel XED* register enum (see section 14.1). There are about 280 registers names for now (of those about 260 being real physical registers), and it's logical to expect some more to be introduced every few years. Therefore, all arrays have a nice and easy size of 512 elements and unused elements are filled with valid “no register here” data. Also, all register names are represented using 16-bit integers.

**context\_map.i**table: context dump data layout and inclusion relations among registers

Size: 17424

Offset	Type	Name	Description
0	ui16	DumpSize	Size of a context dump buffer. All offsets in <b>RegMap32</b> and <b>RegMap64</b> are smaller than <b>DumpSize</b> .
2	ui16	ElemCount	Count of valid elements in all four arrays below. All Arrays have 512 elements, but only the first <b>ElemCount</b> are meaningful.
4	pad[12]		
16	i16[512]	CommonName	<i>Common register name</i> (see 11.1 above)
1040	ac[16][512]	StrName	Array of 512 “strings”, the textual name of each register. Each string is actually a fixed 16-bytes array and names are null-terminated.
9232	regInfo[512]	RegMap32	Per register contextual information for instructions in 32-bit mode code segment.
13328	regInfo[512]	RegMap64	Per register contextual information for instructions in 64-bit mode (a.k.a. <i>long mode</i> ) code segment.

**regInfo**: contextual information about a register (32/64-bit dependant)

Size: 8

Offset	Type	Name	Description
0	ui16	Offset	Offset of the register in a context dump buffer, relative to the start of the buffer, or <b>0xffff</b> if the current register is not part of this 32/64-bit context.
2	ui8	Size	Size of the register in a context dump buffer, or <b>0</b> if the current register is not part of this 32/64-bit context.
3	pad[1]		
4	i16	EffLargest	<i>Effective largest enclosing register</i> for current register (see 11.1 above)
6	pad[2]		

## 12 Address space usage

This table contains all memory addresses referenced by the trace, either in instruction operands or context dumps. Memory addresses are rounded to 4 kB page granularity, and consecutive pages are grouped into page ranges. Therefore, this table is likely a superset of addresses that are really accessed.

The **Address space usage** table is stored in *address\_space.itable*. Its header follows the fixed-size table format described in section 4.2.1, where **record** type is **addressRange**, described below.

**addressRange**: page-granularity address range

Size: 16

Offset	Type	Name	Description
0	i64	<b>First</b>	First byte of a page range, always a multiple of 4 kB.
8	i64	<b>Last</b>	Last byte (included) of a page range. Always greater than <b>First</b> and one less than a multiple of 4 kB.

## 13 Trace.xml

Trace.xml is meant to contain small amounts of heterogeneous data, as well as trace meta-data. It uses plain XML, no specific format or schema being specified for now. Element path names make them somewhat self-descriptive, and besides a few exceptions below, are not described in this document.

File contents is expected to change in the future, although a few elements are most likely to remain because of their potential use by a trace analyzer tool:

- **/Xed/Version**: *Intel XED* version used to generate the trace. Tells which external *XED* tables to use (see section 14.1). *XED* version is actually derived from the version string supplied by *XED* library, but this could change in the future. Currently, version 20150121 is the latest supported.
  - An analyzer should not try to process a trace made with a newer *XED* version since most constants from *XED* tables are likely to be wrong.
  - Traces made with earlier *XED* versions can have all their IDs converted on the fly to the latest version using simple translation tables. Indeed, for all four enum tables, newer versions are supersets of older versions.
- **/App/ProcessID**: Could be useful to some system calls analysis.

## 14 External tables

### 14.1 Intel XED

UMTracer and UMTIndex both use *Intel XED*, a powerful instruction decoding/encoding library, either directly or through Intel *Pin*. **Instruction** table provides detailed enough instruction decoding so that *XED library* (or linking against any C library at all) should not be required to analyze a trace. Nevertheless, some fields refer to lengthy *XED* enumerations tables (register name, instruction mnemonic, etc.) with self-descriptive value names. Those are required for a trace analyzer tool developer to make sense out of those IDs.

While it would be possible to translate them into home-made enums, it has not been required so far. It might be in the future, should *XED* versioning become problematic (*XED* instruction set support is always up to date with *upcoming* Intel processors). Latest relevant header files from *XED* are found in **XED v20150121** folder along with this document (see section 13 to identify which version is used for a given trace). Readily usable in C/C++, their contents can be easily adapted for other languages such as Java or Python.

- **xed-reg-enum.h**: register “name” enum. Some values are not likely to find their way in a user-mode trace (**cr<sub>x</sub>** or **dr<sub>x</sub>** registers), but most others are used. In particular, **XED\_REG\_INVALID** is used in trace records to mean “no register”, “empty operand slot”, etc.
- **xed-iclass-enum.h**: instruction mnemonics
- **xed-iform-enum.h**: specific form of an instruction, e.g. **add reg,reg** and **add mem,reg** are two forms of the same mnemonic.
- **xed-icategory-enum.h**: instruction categories. May provide easier analysis dispatching than a lengthy switch-case on form or mnemonic.

Although not directly used by trace records, **idata.txt**, found along with other *XED* files, provides the complete mapping between instruction-related enums. For more information about *XED*, see the *Intel Pin* page: <http://www.pintool.org/downloads.html>.

### 14.2 System call IDs

As explained in section 8.1, system call identifiers (their ordinal in kernel function tables) are OS version-specific. The tracer converts them to a uniform ID so that the same ID relates to the same exported kernel function (as long as the function semantics and argument count remains the same across OS versions). The latest version of that enum table can be found in **syscall\_unified\_id\_v1.h** (see section 138.4 to retrieve which version is used for a given trace). Enum names are described in section 8.1.