# Kernel Pool Exploitation on Windows 7

Tarjei Mandt | Black Hat DC 2011

# About Me

- ▶ Security Researcher at Norman
  - ▶ Malware Detection Team (MDT)
- ▶ Interests
  - ▶ Vulnerability research
  - ▶ Operating systems internals
  - ▶ Exploit mitigations
- ▶ Reported some bugs in the Windows kernel
  - ▶ Windows VDM Task Initialization Vulnerability (MS10-098)
  - ▶ Windows Class Data Handling Vulnerability (MS10-073)
- ▶ I have a Twitter account ☺
  - ▶ @kernelpool

▶

# Agenda

▸ Introduction

▸ Kernel Pool Internals

▸ Kernel Pool Attacks

▸ Case Study / Demo

  ▸ MS10-098 (win32k.sys)

  ▸ MS10-058 (tcpip.sys)

▸ Kernel Pool Hardening

▸ Conclusion

# Introduction

Kernel Pool Exploitation
on Windows 7

# Introduction

- Exploit mitigations such as DEP and ASLR do not prevent exploitation in every case
    - JIT spraying, memory leaks, etc.
- Privilege isolation is becoming an important component in confining application vulnerabilities
    - Browsers and office applications employ "sandboxed" render processes
    - Relies on (security) features of the operating system
- In turn, this has motivated attackers to focus their efforts on privilege escalation attacks
    - Arbitrary ring0 code execution $\rightarrow$ OS security undermined

# The Kernel Pool

▶ Resource for dynamically allocating memory

▶ Shared between all kernel modules and drivers

▶ Analogous to the user-mode heap

  ▶ Each pool is defined by its own structure

  ▶ Maintains lists of free pool chunks

▶ Highly optimized for performance

  ▶ No kernel pool cookie or pool header obfuscation

▶ The kernel executive exports dedicated functions for handling pool memory

  ▶ **ExAllocatePool**\* and **ExFreePool**\* (discussed later)

▶

# Kernel Pool Exploitation

- An attacker's ability to leverage pool corruption vulnerabilities to execute arbitrary code in ring 0
    - Similar to traditional heap exploitation
- Kernel pool exploitation requires careful modification of kernel pool structures
    - Access violations are likely to end up with a bug check (BSOD)
- Up until Windows 7, kernel pool overflows could be generically exploited using write-4 techniques
    - SoBeIt[2005]
    - Kortchinsky[2008]

# Previous Work

▸ Primarily focused on XP/2003 platforms

▸ How To Exploit Windows Kernel Memory Pool
  ▸ Presented by SoBeIt at XCON 2005
  ▸ Proposed two write-4 exploit methods for overflows

▸ Real World Kernel Pool Exploitation
  ▸ Presented by Kostya Kortchinsky at SyScan 2008
  ▸ Discussed four write-4 exploitation techniques
  ▸ Demonstrated practical exploitation of MS08-001

▸ All the above exploitation techniques were addressed in Windows 7 (Beck[2009])

# Contributions

▸ Elaborate on the internal structures and changes made to the Windows 7 (and Vista) kernel pool

▸ Identify weaknesses in the Windows 7 kernel pool and show how an attacker may leverage these to exploit pool corruption vulnerabilities

▸ Propose ways to thwart the discussed attacks and further harden the kernel pool

# Kernel Pool Internals

Kernel Pool Exploitation
on Windows 7

# Kernel Pool Fundamentals

- Kernel pools are divided into types
  - Defined in the **POOL_TYPE** enum
  - Non-Paged Pools, Paged Pools, Session Pools, etc.
- Each kernel pool is defined by a *pool descriptor*
  - Defined by the **POOL_DESCRIPTOR** structure
  - Tracks the number of allocs/frees, pages in use, etc.
  - Maintains lists of free pool chunks
- The initial descriptors for paged and non-paged pools are defined in the **nt!PoolVector** array
  - Each index points to an array of one or more descriptors

# Kernel Pool Descriptor (Win7 RTM x86)

- kd> **dt nt!_POOL_DESCRIPTOR**
  - +0x000 PoolType                         : _POOL_TYPE
  - +0x004 PagedLock                        : _KGUARDED_MUTEX
  - +0x004 NonPagedLock       : Uint4B
  - +0x040 RunningAllocs       : Int4B
  - +0x044 RunningDeAllocs   : Int4B
  - +0x048 TotalBigPages       : Int4B
  - +0x04c ThreadsProcessingDeferrals : Int4B
  - +0x050 TotalBytes                        : Uint4B
  - +0x080 PoolIndex                         : Uint4B
  - +0x0c0 TotalPages                        : Int4B
  - +0x100 PendingFrees       : Ptr32 Ptr32 Void
  - +0x104 PendingFreeDepth: Int4B
  - +0x140 ListHeads                         : [512] _LIST_ENTRY

# Non-Uniform Memory Architecture

- In a NUMA system, processors and memory are grouped together in smaller units called *nodes*
  - Faster memory access when local memory is used
- The kernel pool always tries to allocate memory from the ideal node for a process
  - Most desktop systems only have a single node
- Each node is defined by the **KNODE** data structure
  - Pointers to all **KNODE** structures are stored in the **nt!KeNodeBlock** array
  - Multiple processors can be linked to the same node
- We can dump NUMA information in WinDbg
  - kd> **!numa**

# NUMA Node Structure (Win7 RTM x86)

- kd> **dt nt!_KNODE**
  - +0x000 PagedPoolSListHead  : _SLIST_HEADER
  - +0x008 NonPagedPoolSListHead        : [3] _SLIST_HEADER
  - +0x020 Affinity                : _GROUP_AFFINITY
  - +0x02c ProximityId            : Uint4B
  - +0x030 NodeNumber            : Uint2B
  - +0x032 PrimaryNodeNumber  : Uint2B
  - +0x034 MaximumProcessors  : UChar
  - +0x035 Color                    : UChar
  - +0x036 Flags                    : _flags
  - +0x037 NodePad0              : UChar
  - +0x038 Seed                    : Uint4B
  - +0x03c MmShiftedColor      : Uint4B
  - +0x040 FreeCount              : [2] Uint4B
  - +0x048 CachedKernelStacks  : _CACHED_KSTACK_LIST
  - +0x060 ParkLock              : Int4B
  - +0x064 NodePad1              : Uint4B

Array index to associated pool descriptor on NUMA compatible systems

# NUMA on Intel Core i7 820QM

kd> **!numa**

NUMA Summary:

------------

Number of NUMA nodes : 1

Number of Processors : 8

Single node, despite multicore CPU architecture

MmAvailablePages     : 0x00099346

KeActiveProcessors   :

********---------------------------------------------------- (00000000000000ff)

NODE 0 (FFFFF80003412B80):

  Group           : 255 (Assigned, Committed, Assignment Adjustable)

  ProcessorMask   :  (ff)

  ProximityId     : 0

  Capacity        : 8

  Color           : 0x00000000

  […]

# Non-Paged Pool

- Non-pagable system memory
  - Guaranteed to reside in physical memory at all times
- Number of pools stored in **nt!ExpNumberOfNonPagedPools**
- On uniprocessor systems, the first index of the **nt!PoolVector** array points to the non-paged pool descriptor
  - kd> **dt nt!_POOL_DESCRIPTOR poi(nt!PoolVector)**
- On multiprocessor systems, each node has its own non-paged pool descriptor
  - Pointers stored in **nt!ExpNonPagedPoolDescriptor** array

# Paged Pool

- Pageable system memory
  - Can only be accessed at IRQL < DPC/Dispatch level
- Number of paged pools defined by **nt!ExpNumberOfPagedPools**
- On uniprocessor systems, four (4) paged pool descriptors are defined
  - Index 1 through 4 in **nt!ExpPagedPoolDescriptor**
- On multiprocessor systems, one (1) paged pool descriptor is defined per node
- One additional paged pool descriptor is defined for prototype pools / full page allocations
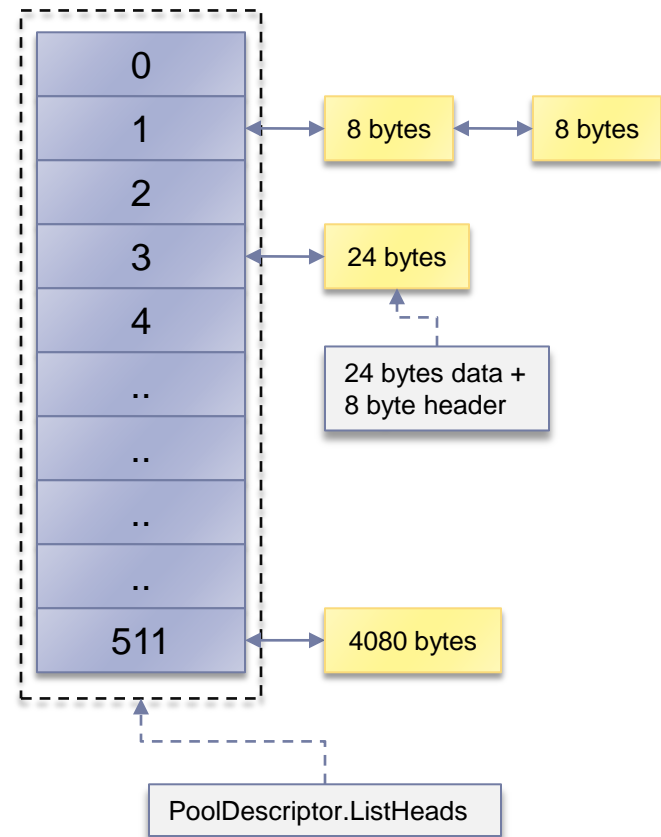  - Index 0 in **nt!ExpPagedPoolDescriptor**

# Session Paged Pool

- Pageable system memory for session space
  - E.g. Unique to each logged in user
- Initialized in **nt!MiInitializeSessionPool**
- On Vista, the pool descriptor pointer is stored in **nt!ExpSessionPoolDescriptor** (session space)
- On Windows 7, a pointer to the pool descriptor from the current thread is used
  - KTHREAD->Process->Session.PagedPool
- Non-paged session allocations use the global non-paged pools

# Pool Descriptor Free Lists (x86)

▸ Each pool descriptor has a *ListHeads* array of 512 doubly-linked lists of free chunks of the same size

  ▸ 8 byte granularity

  ▸ Used for allocations up to 4080 bytes

▸ Free chunks are indexed into the ListHeads array by block size

  ▸ BlockSize: (NumBytes+0xF) >> 3

▸ Each pool chunk is preceded by an 8-byte pool header

# Kernel Pool Header (x86)

- kd> **dt nt!_POOL_HEADER**
  - +0x000 PreviousSize        : Pos 0, 9 Bits
  - +0x000 PoolIndex           : Pos 9, 7 Bits
  - +0x002 BlockSize           : Pos 0, 9 Bits
  - +0x002 PoolType            : Pos 9, 7 Bits
  - +0x004 PoolTag             : Uint4B
- *PreviousSize*: BlockSize of the preceding chunk
- *PoolIndex*: Index into the associated pool descriptor array
- *BlockSize*: (NumberOfBytes+0xF) >> 3
- *PoolType*: Free=0, Allocated=(PoolType|2)
- *PoolTag*: 4 printable characters identifying the code responsible for the allocation
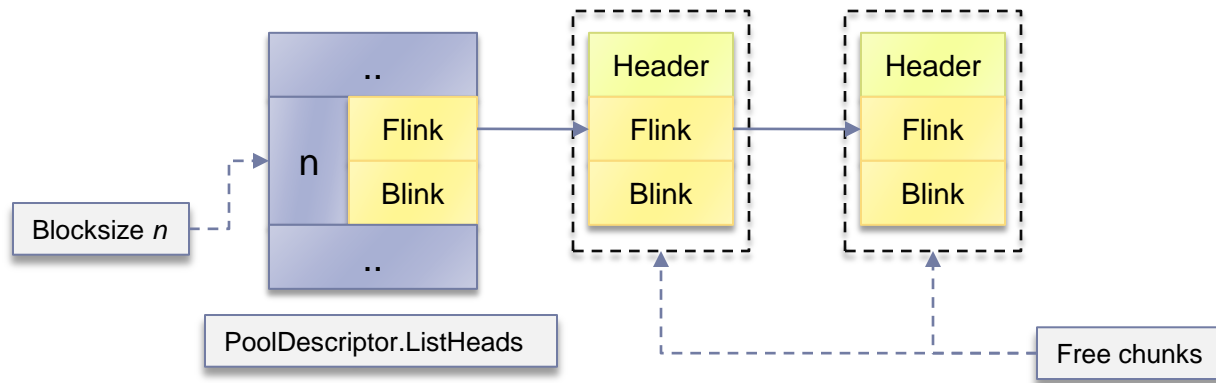
# Kernel Pool Header (x64)

- kd> **dt nt!_POOL_HEADER**
  - +0x000 PreviousSize   : Pos 0, 8 Bits
  - +0x000 PoolIndex       : Pos 8, 8 Bits
  - +0x000 BlockSize        : Pos 16, 8 Bits
  - +0x000 PoolType         : Pos 24, 8 Bits
  - +0x004 PoolTag           : Uint4B
  - +0x008 ProcessBilled    : Ptr64 _EPROCESS
- *BlockSize*: (NumberOfBytes+0x1F) >> 4
  - 256 ListHeads entries due to 16 byte block size
- *ProcessBilled*: Pointer to process object charged for the pool allocation (used in quota management)

# Free Pool Chunks

▸ If a pool chunk is freed to a pool descriptor ListHeads list, the header is followed by a **LINK_ENTRY** structure

- ▸ Pointed to by the ListHeads doubly-linked list

- ▸ kd> **dt nt!_LIST_ENTRY**
  +0x000 Flink        : Ptr32 _LIST_ENTRY
  +0x004 Blink        : Ptr32 _LIST_ENTRY

# Lookaside Lists

▸ Kernel uses *lookaside lists* for faster allocation/deallocation of small pool chunks

  ▸ Singly-linked LIFO lists

  ▸ Optimized for performance – e.g. no checks

▸ Separate per-processor lookaside lists for pagable and non-pagable allocations

  ▸ Defined in the Processor Control Block (KPRCB)

  ▸ Maximum BlockSize being 0x20 (256 bytes)

  ▸ 8 byte granularity, hence 32 lookaside lists per type

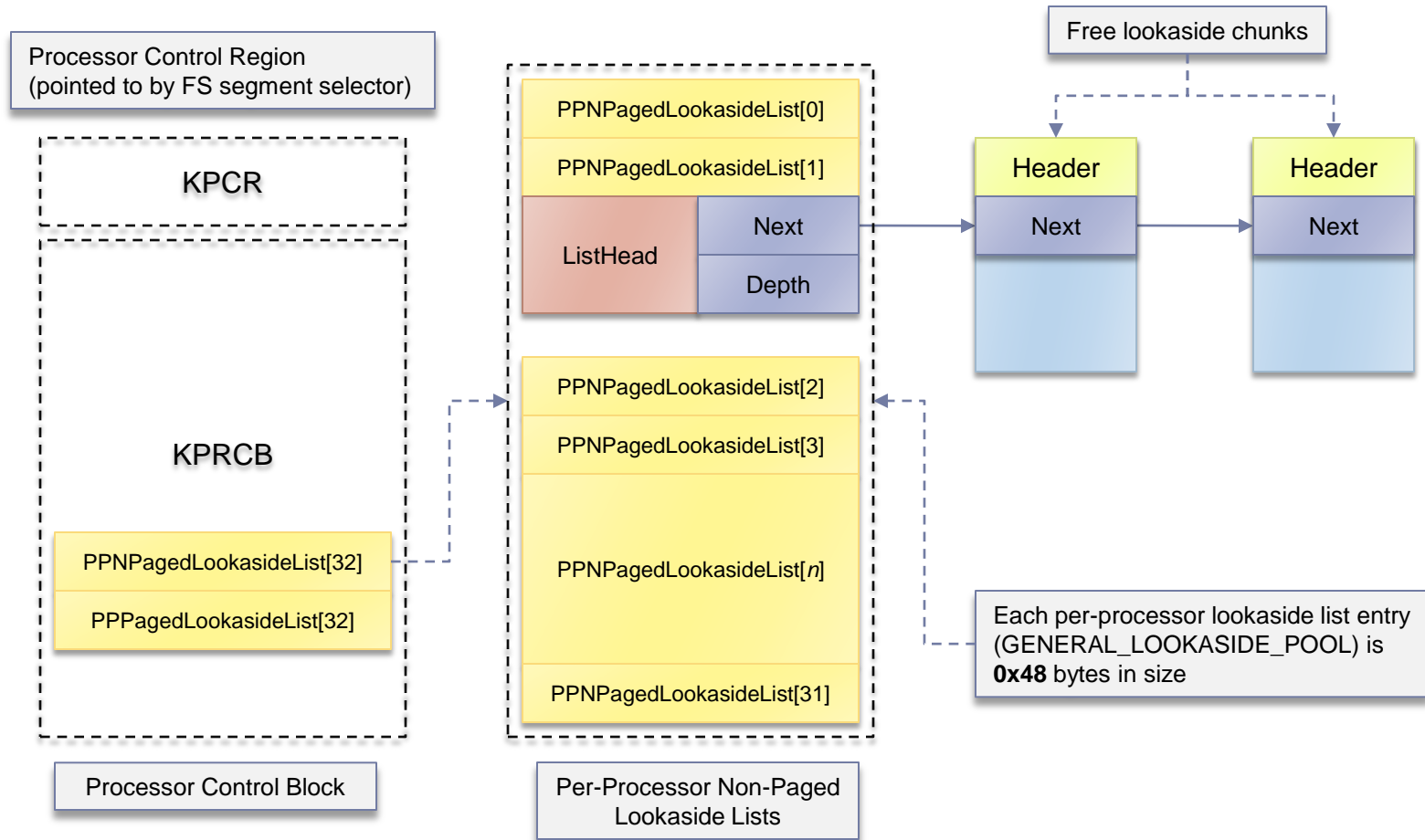▸ Each lookaside list is defined by a **GENERAL_LOOKASIDE_POOL** structure

▸

# General Lookaside (Win7 RTM x86)

- kd> **dt _GENERAL_LOOKASIDE_POOL**
  - +0x000 ListHead              : _SLIST_HEADER
  - +0x000 SingleListHead        : _SINGLE_LIST_ENTRY
  - +0x008 Depth                 : Uint2B
  - +0x00a MaximumDepth          : Uint2B
  - +0x00c TotalAllocates        : Uint4B
  - +0x010 AllocateMisses        : Uint4B
  - +0x010 AllocateHits          : Uint4B
  - +0x014 TotalFrees            : Uint4B
  - +0x018 FreeMisses            : Uint4B
  - +0x018 FreeHits              : Uint4B
  - +0x01c Type                  : _POOL_TYPE
  - +0x020 Tag                   : Uint4B
  - +0x024 Size                  : Uint4B
  - […]

# Lookaside Lists (Per-Processor)



Processor Control Region
(pointed to by FS segment selector)

KPCR

KPRCB

PPNPagedLookasideList[32]

PPPagedLookasideList[32]

Processor Control Block

PPNPagedLookasideList[0]

PPNPagedLookasideList[1]

ListHead    Next    Depth

PPNPagedLookasideList[2]

PPNPagedLookasideList[3]

PPNPagedLookasideList[$n$]

PPNPagedLookasideList[31]

Per-Processor Non-Paged
Lookaside Lists

Free lookaside chunks

Header    Next

Header    Next

Each per-processor lookaside list entry
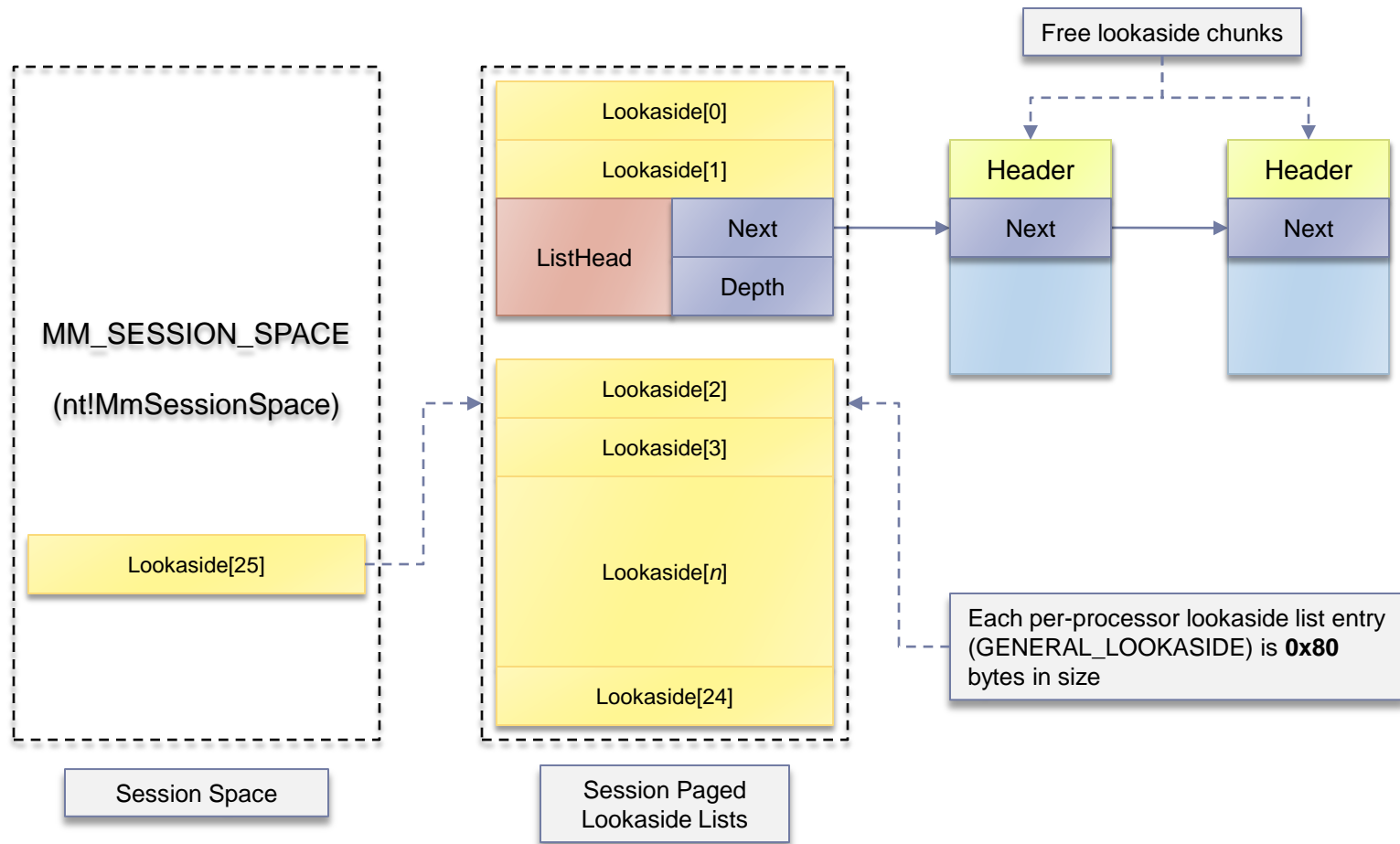(GENERAL_LOOKASIDE_POOL) is
**0x48** bytes in size

# Lookaside Lists (Session)

- Separate per-session lookaside lists for pagable allocations
  - Defined in session space (**nt!ExpSessionPoolLookaside**)
  - Maximum BlockSize being 0x19 (200 bytes)
  - Uses the same structure (with padding) as per-processor lists
  - All processors use the same session lookaside lists
- Non-paged session allocations use the per-processor non-paged lookaside list
- Lookaside lists are disabled if *hot/cold separation* is used
  - **nt!ExpPoolFlags** & 0x100
  - Used during system boot to increase speed and reduce the memory footprint

# Lookaside Lists (Session)

# Dedicated Lookaside Lists

- Frequently allocated buffers (of fixed size) in the NT kernel have dedicated lookaside lists
  - Object create information
  - I/O request packets
  - Memory descriptor lists
- Defined in the processor control block (KPRCB)
  - 16 **PP_LOOKASIDE_LIST** structures, each defining one per-processor and one system-wide list

# Large Pool Allocations

▸ Allocations greater than 0xff0 (4080) bytes

▸ Handled by the function **nt!ExpAllocateBigPool**

  ▸ Internally calls **nt!MiAllocatePoolPages**

    ▸ Requested size is rounded up to the nearest page size

  ▸ Excess bytes are put back at the end of the appropriate pool descriptor ListHeads list

▸ Each node (e.g. processor) has 4 singly-linked lookaside lists for big pool allocations

  ▸ 1 paged for allocations of a single page

  ▸ 3 non-paged for allocations of page count 1, 2, and 3

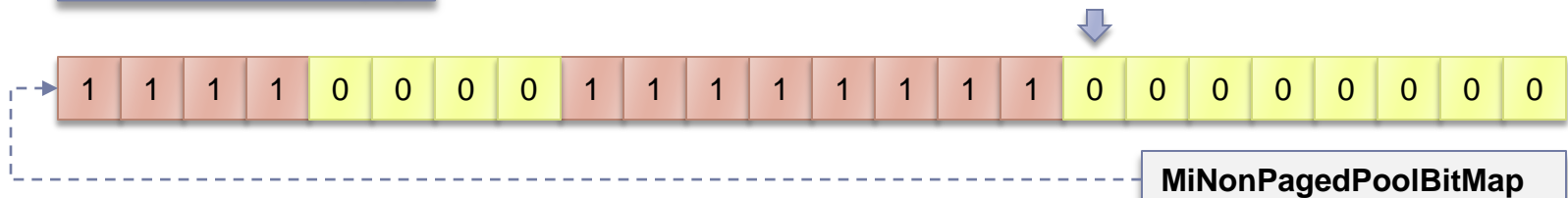  ▸ Defined in **KNODE** (KPCR.PrcbData.ParentNode)

# Large Pool Allocations

- If lookaside lists cannot be used, an *allocation bitmap* is used to obtain the requested pool pages
    - Array of bits that indicate which memory pages are in use
    - Defined by the **RTL_BITMAP** structure
- The bitmap is searched for the first index that holds the requested number of unused pages
- Bitmaps are defined for every major pool type with its own dedicated memory
    - E.g. **nt!MiNonPagedPoolBitMap**
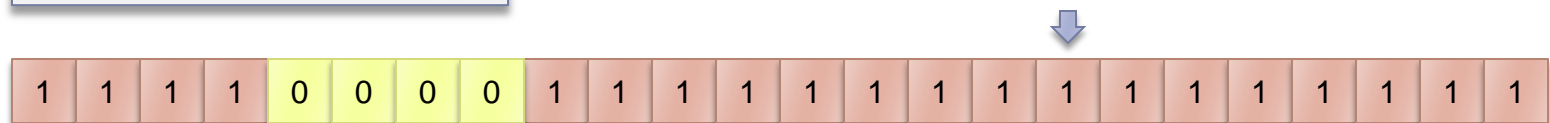- The array of bits is located at the beginning of the pool memory range

# Bitmap Search (Simplified)

1. **MiAllocatePoolPages(NonPagedPool, 0x8000)**

2. **RtlFindClearBits(...)**

| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**MiNonPagedPoolBitMap**

3. **RtlFindAndSetClearBits(...)**

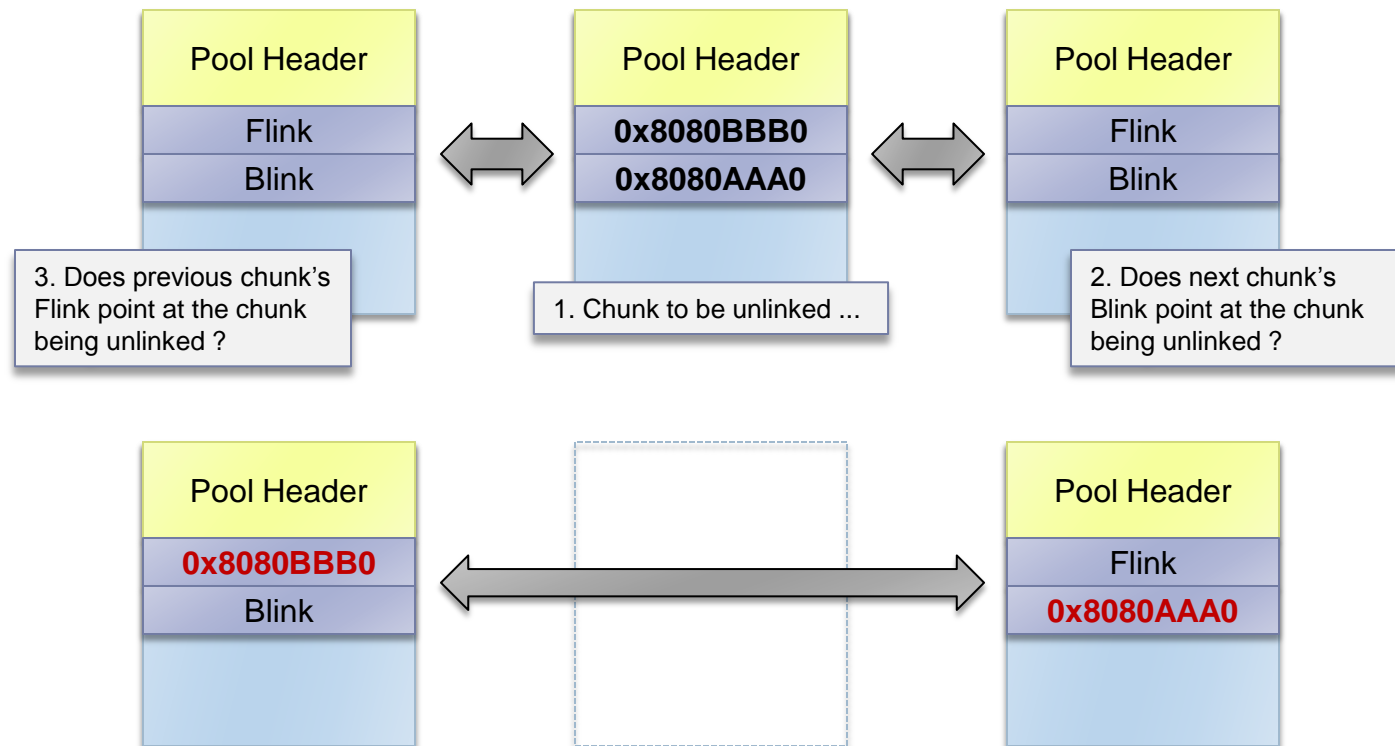| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

4. PageAddress = **MiNonPagedPoolStartAligned** + ( BitOffset << 0xC )

# Allocation Algorithm

▸ The kernel exports several allocation functions for kernel modules and drivers to use

▸ All exported kernel pool allocation routines are essentially wrappers for **ExAllocatePoolWithTag**

▸ The allocation algorithm returns a free chunk by checking with the following (in order)

  ▸ Lookaside list(s)

  ▸ ListHeads list(s)

  ▸ Pool page allocator

▸ Windows 7 performs *safe unlinking* when pulling a chunk from a free list ([Beck[2009]](#))

▸

# Safe Pool Unlinking

# ExAllocatePoolWithTag (1/2)

- **PVOID ExAllocatePoolWithTag(POOL_TYPE PoolType, SIZE_T NumberOfBytes, ULONG Tag)**
- If NumberOfBytes > 0xff0
  - Call **nt!ExpAllocateBigPool**
- If PagedPool requested
  - If (PoolType & SessionPoolMask) and BlockSize <= 0x19
    - □ Try the session paged lookaside list
    - □ Return on success
  - Else If BlockSize <= 0x20
    - □ Try the per-processor paged lookaside list
    - □ Return on success
  - Try and lock paged pool descriptor (round robin)

# ExAllocatePoolWithTag (2/2)

- ### Else (NonPagedPool requested)
  - If BlockSize <= 0x20
    - Try the per-processor non-paged lookaside list
    - Return on success
  - Try and lock non-paged pool descriptor (local node)
- ### Use ListHeads of currently locked pool
  - For n in range(BlockSize,512)
    - If ListHeads[n] is empty, try next BlockSize
    - Safe unlink first entry and split if larger than needed
    - Return on success
  - If failed, expand the pool by adding a page
    - Call **nt!MiAllocatePoolPages**
    - Split entry and return on success

# Splitting Pool Chunks

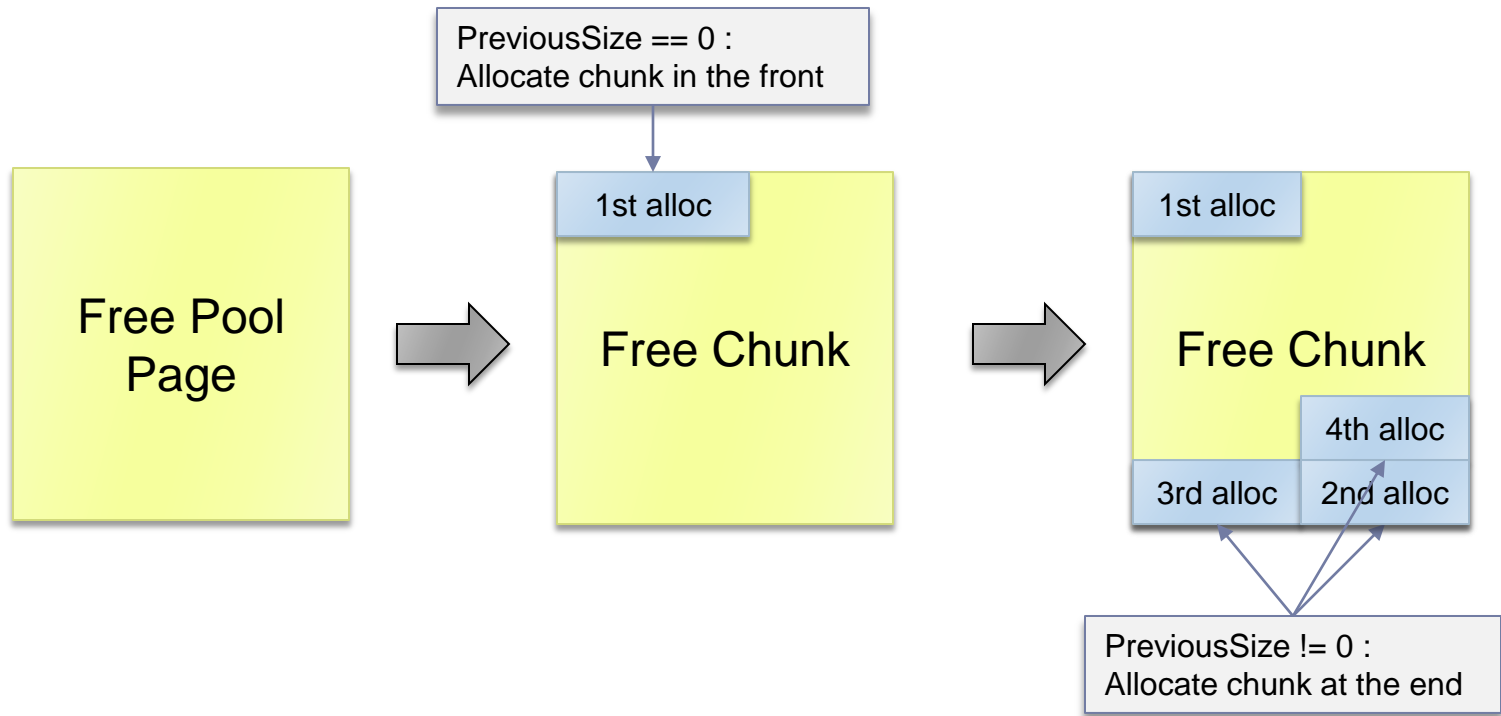▶ If a chunk larger than the size requested is returned from ListHeads[n], the chunk is split

  ▶ If chunk is page aligned, the requested size is allocated from the <u>front of the chunk</u>

  ▶ If chunk is <u>not</u> page aligned, the requested size is allocated at the <u>end of the chunk</u>

▶ The remaining fragment of the split chunk is put at the <u>tail</u> of the proper ListHeads[n] list

# Splitting Pool Chunks

# Free Algorithm

▸ The free algorithm inspects the pool header of the chunk to be freed and frees it to the appropriate list

  ▸ Implemented by **ExFreePoolWithTag**

▸ Bordering free chunks may be merged with the freed chunk to reduce fragmentation

  ▸ Windows 7 uses safe unlinking in the merging process

# ExFreePoolWithTag (1/2)

- **VOID ExFreePoolWithTag(PVOID Address, ULONG Tag)**
- If Address (chunk) is page aligned
  - Call **nt!MiFreePoolPages**
- If Chunk->BlockSize != NextChunk->PreviousSize
  - BugCheckEx(BAD_POOL_HEADER)
- If (PoolType & PagedPoolSession) and BlockSize <= 0x19
  - Put in session pool lookaside list
- Else If BlockSize <= 0x20 and pool is local to processor
  - If (PoolType & PagedPool)
    - Put in per-processor paged lookaside list
  - Else (NonPagedPool)
    - Put in per-processor non-paged lookaside list
- Return on sucess

▷

# ExFreePoolWithTag (2/2)

- ## If the DELAY_FREE pool flag is set
  - If pending frees >= 0x20
    - Call **nt!ExDeferredFreePool**
  - Add to front of pending frees list (singly-linked)
- ## Else
  - If next chunk is free and not page aligned
    - Safe unlink and merge with current chunk
  - If previous chunk is free
    - Safe unlink and merge with current chunk
  - If resulting chunk is a full page
    - Call **nt!MiFreePoolPages**
  - Else
    - Add to front of appropriate ListHeads list

# Merging Pool Chunks

**1**

| Pool Header (free) | | Pool Header (busy) | | Pool Header (free) |

Chunk to be freed

**2**

Next chunk unlinked

| Pool Header (free) | | Pool Header (busy) | | unlinked chunk |

**3**

Previous chunk unlinked

Merge with next

| unlinked chunk | Pool Header (busy) | BlockSize updated |

**4**

Marked as free and returned

Merge with previous

| Pool Header (free) | BlockSize updated |

# Delayed Pool Frees

▸ A performance optimization that frees several pool allocations at once to amortize pool acquisition/release

  ▸ Briefly mentioned in mxatone[2008]

▸ Enabled when **MmNumberOfPhysicalPages** >= 0x1fc00

  ▸ Equivalent to 508 MBs of RAM on IA-32 and AMD64

  ▸ **nt!ExpPoolFlags** & 0x200

▸ Each call to **ExFreePoolWithTag** appends a pool chunk to a singly-linked deferred free list specific to each pool descriptor

  ▸ Current number of entries is given by **PendingFreeDepth**

  ▸ The list is processed by the function **ExDeferredFreePool** if it has 32 or more entries

▸

# ExDeferredFreePool

- **VOID ExDeferredFreePool(PPOOL_DESCRIPTOR PoolDescriptor, BOOLEAN bMultiThreaded)**
- For each entry on pending frees list
  - If next chunk is free and not page aligned
    - Safe unlink and merge with current chunk
  - If previous chunk is free
    - Safe unlink and merge with current chunk
  - If resulting chunk is a full page
    - Add to full page list
  - Else
    - Add to front of appropriate ListHeads list
- For each page in full page list
  - Call **nt!MiFreePoolPages**

# Free Pool Chunk Ordering

- Frees to the lookaside and pool descriptor ListHeads are always put in the front of the appropriate list
  - Exceptions are remaining fragments of split blocks which are put at the tail of the list
  - Blocks are split when the pool allocator returns chunks larger than the requested size
    - Full pages split in **ExpBigPoolAllocation**
    - ListHeads[n] entries split in **ExAllocatePoolWithTag**
- Allocations are always made from the most recently used blocks, from the front of the appropriate list
  - Attempts to use the CPU cache as much as possible

# Kernel Pool Attacks

Kernel Pool Exploitation
on Windows 7

# Overview

- Traditional ListEntry Attacks (< Windows 7)
- ListEntry Flink Overwrite
- Lookaside Pointer Overwrite
- PoolIndex Overwrite
- PendingFrees Pointer Overwrite
- Quota Process Pointer Overwrite

# ListEntry Overwrite (< Windows 7)

▸ All free list (ListHeads) pool chunks are linked together by LIST_ENTRY structures

▸ Vista and former versions do not validate the structures' forward and backward pointers

▸ A ListEntry overwrite may be leveraged to trigger a write-4 in the following situations

  ▸ Unlink in merge with next pool chunk

  ▸ Unlink in merge with previous pool chunk

  ▸ Unlink in allocation from ListHeads[n] free list

▸ Discussed in Kortchinsky[2008] and SoBeIt[2005]

# ListEntry Overwrite (Merge With Next)



**1**

Pool Header (busy)

Pool overflow →

Pool Header | List Entry

**2**

Chunk to be freed

Pool Header (busy)

Pool Header | List Entry

PoolType | Flink | Blink

PoolType set to 0 (free)

**3**

Pool Header (free)

unlinked chunk → write-4

Chunk size is updated to accomodate the merged chunk

When the **overflowing chunk** is freed, the next bordering chunk is merged and unlinked

# ListEntry Overwrite (Merge With Previous)

# ListEntry Flink Overwrite

- Windows 7 uses safe unlinking to validate the LIST_ENTRY pointers of a chunk being unlinked
- In allocating a pool chunk from a ListHeads free list, the kernel fails to properly validate its forward link
  - The algorithm validates the ListHeads[n] LIST_ENTRY structure instead
- Overwriting the <u>forward link</u> of a free chunk may cause the address of ListHeads[n] to be written to an attacker controlled address
  - Target ListHeads[n] list must hold at least two free chunks

# The Not So Safe Unlink

After unlink
- **FakeEntry.Blink** = ListHeads[n]
- **ListHeads[n].Flink** = FakeEntry

Pool Descriptor ListHeads

Chunk to be unlinked

FakeEntry

**ListHeads[n].Flink**
(validated in safe unlink)

Index for BlockSize *n*,
**Flink** points to first
chunk to be allocated

Pool overflow

ListEntry

Flink

Blink

Pool Header

Flink

Blink

Pool Header

Flink

Blink

**ListHeads[n].Blink**
(validated in safe unlink)

**NextEntry.Blink**
(validated in safe unlink)

**PreviousEntry.Flink**
(validated in safe unlink)

# ListEntry Flink Overwrite

▸ In the following output, the address of ListHeads[n] (**esi**) in the pool descriptor is written to an attacker controlled address (**eax**)

▸ Pointers are not sufficiently validated when allocating a pool chunk from the free list

```
eax=80808080 ebx=829848c0 ecx=8cc15768 edx=8cc43298 esi=82984a18 edi=829848c4
eip=8296f067 esp=82974c00 ebp=82974c48 iopl=0        nv up ei pl zr na pe nc
cs=0008  ss=0010  ds=0023  es=0023  fs=0030  gs=0000        efl=00010246

nt!ExAllocatePoolWithTag+0x4b7:
8296f067 897004        mov     dword ptr [eax+4],esi ds:0023:80808084=????????
```
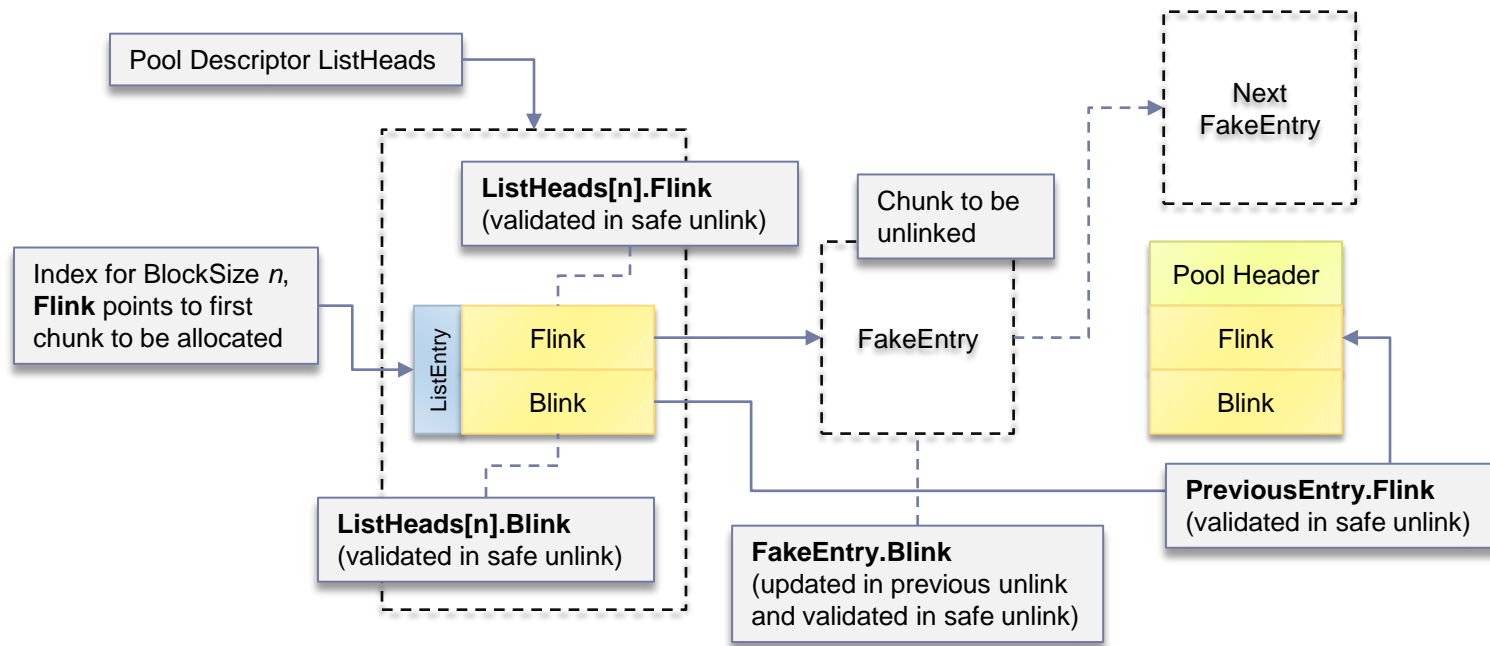
# ListEntry Flink Overwrite

▸ After unlink, the attacker may control the address of the next allocated entry

  ▸ **ListHeads[n].Flink** = FakeEntry

▸ FakeEntry can be safely unlinked as its blink was updated to point back to ListHeads[n]

  ▸ **FakeEntry.Blink** = ListHeads[n]

▸ If a user-mode pointer is used in the overwrite, the attacker could fully control the contents of the next allocation

▸

# ListEntry Flink Overwrite



Pool Descriptor ListHeads

Index for BlockSize *n*, **Flink** points to first chunk to be allocated

**ListHeads[n].Flink**
(validated in safe unlink)

**ListHeads[n].Blink**
(validated in safe unlink)

ListEntry

Flink

Blink

Chunk to be unlinked

FakeEntry

**FakeEntry.Blink**
(updated in previous unlink and validated in safe unlink)

Next FakeEntry

Pool Header

Flink

Blink

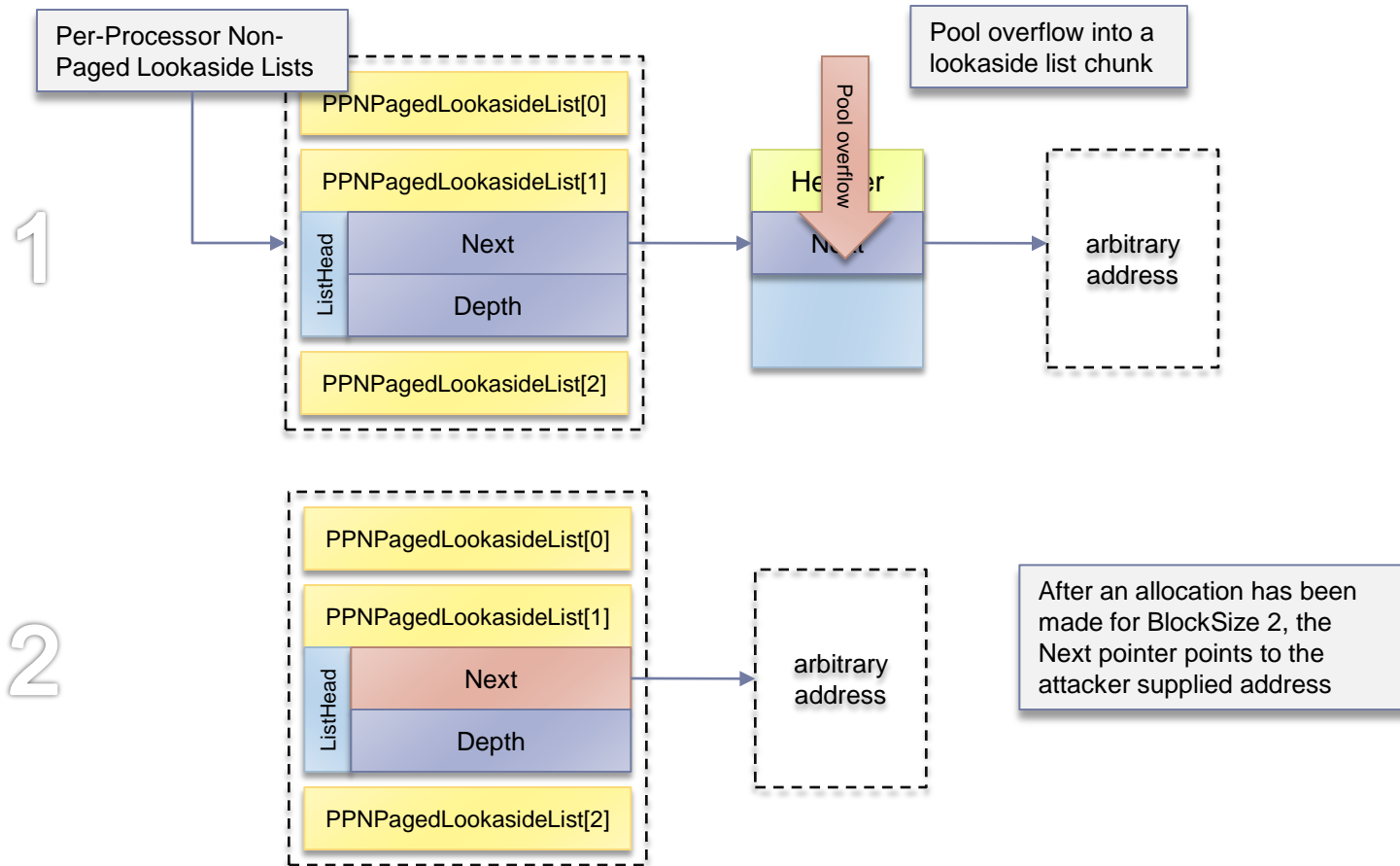**PreviousEntry.Flink**
(validated in safe unlink)

# Lookaside Pointer Overwrite

▸ **Pool chunks and pool pages on lookaside lists are singly-linked**

  ▸ Each entry holds a pointer to the <u>next</u> entry

  ▸ Overwriting a next pointer may cause the kernel pool allocator to return an attacker controlled address

▸ **A pool chunk is freed to a lookaside list if the following hold**

  ▸ BlockSize <= 0x20 for paged/non-paged pool chunks

  ▸ BlockSize <= 0x19 for paged session pool chunks

  ▸ Lookaside list for target BlockSize is not full

  ▸ Hot/cold page separation is not used

▸

# Lookaside Pointer Overwrite (Chunks)

Per-Processor Non-Paged Lookaside Lists

Pool overflow into a lookaside list chunk

**1**

PPNPagedLookasideList[0]

PPNPagedLookasideList[1]

ListHead

Next

Depth

PPNPagedLookasideList[2]

Pool overflow

Header

Next

arbitrary address

**2**

PPNPagedLookasideList[0]

PPNPagedLookasideList[1]

ListHead

Next

Depth

PPNPagedLookasideList[2]

arbitrary address

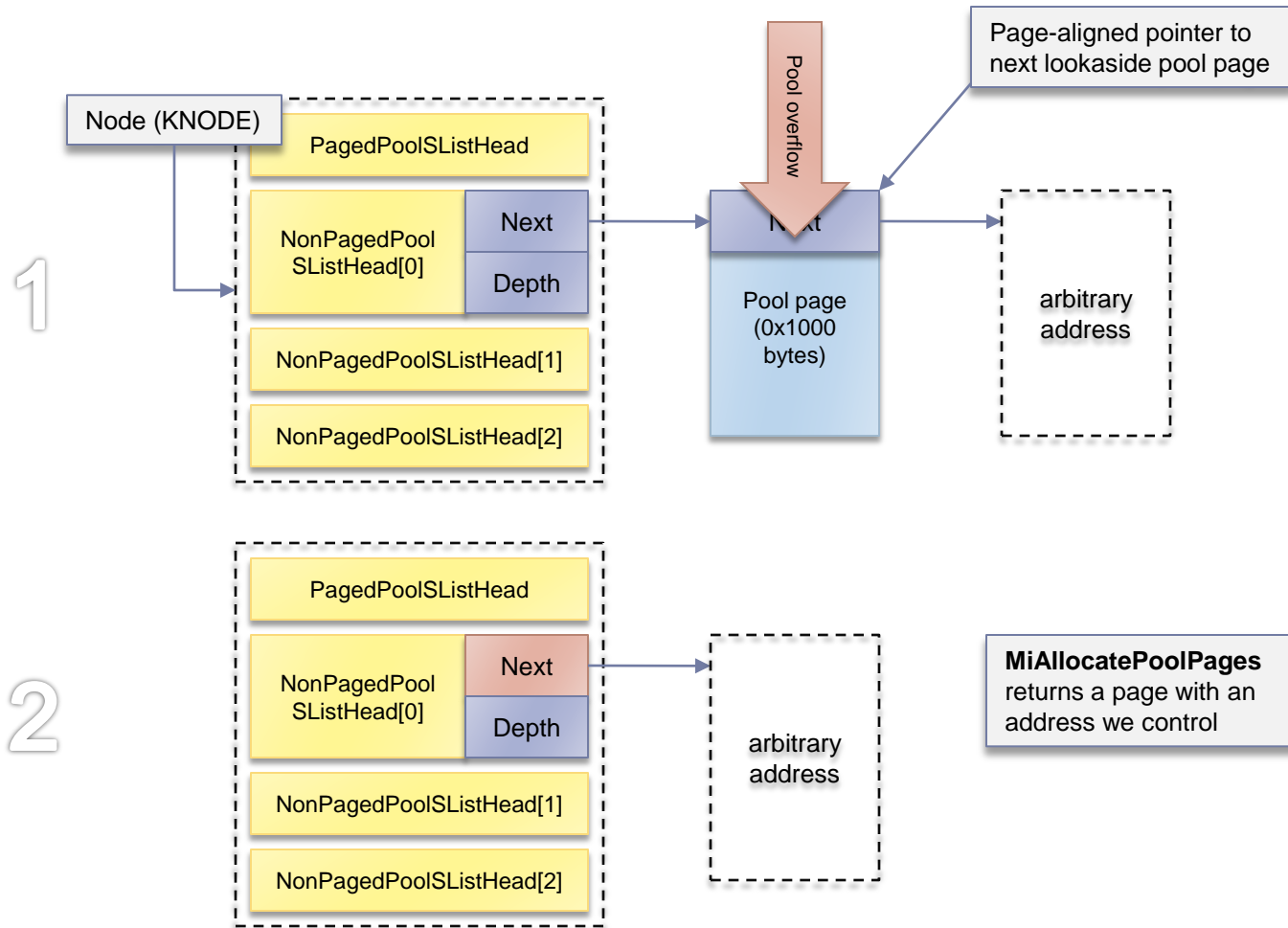After an allocation has been made for BlockSize 2, the Next pointer points to the attacker supplied address

# Lookaside Pointer Overwrite (Pages)

- A pool page is freed to a lookaside list if the following hold
  - NumberOfPages = 1 for paged pool pages
  - NumberOfPages <= 3 for non-paged pool pages
  - Lookaside list for target page count is not full
    - Size limit determined by physical page count in system
- A pointer overwrite of lookaside pages requires at most a pointer-wide overflow
  - No pool headers on free pool pages!
  - Partial pointer overwrites may also be sufficient
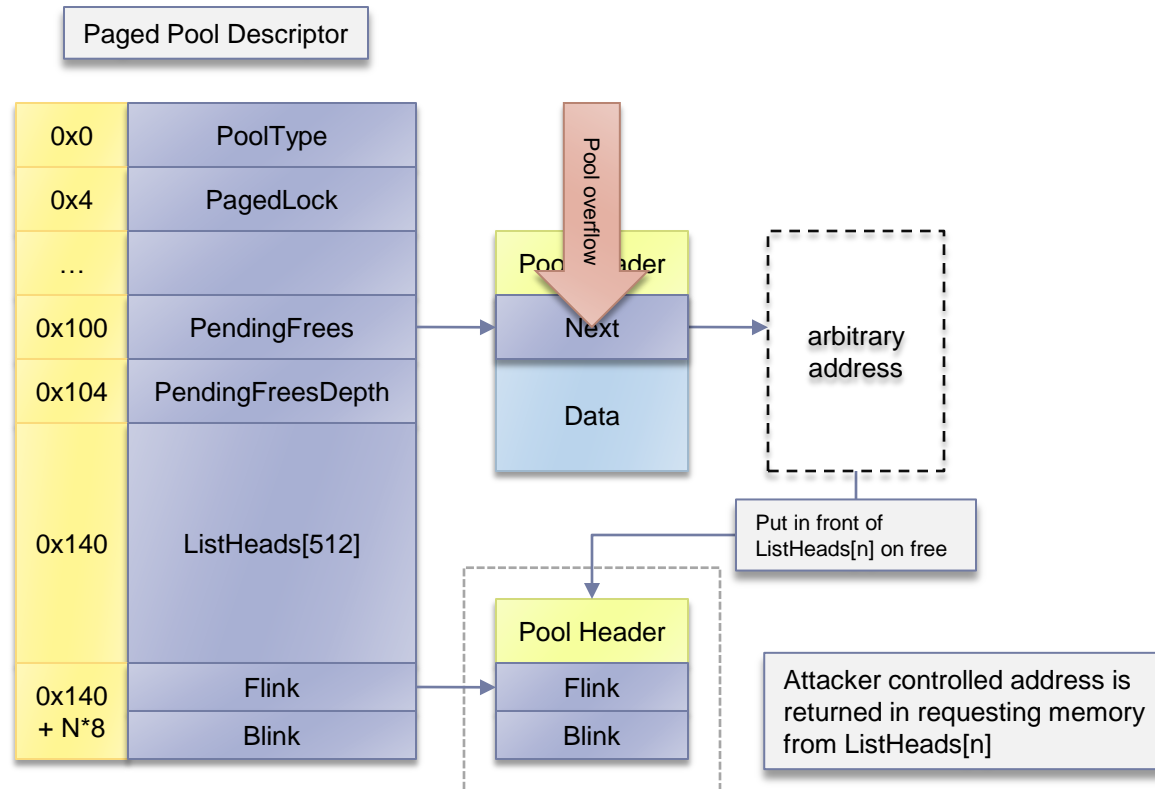
# Lookaside Pointer Overwrite (Pages)

# PendingFrees Pointer Overwrite

- Pool chunks waiting to be freed are stored in the pool descriptor deferred free list
  - Singly-linked (similar to lookaside list)
- Overwriting a chunk's <u>next pointer</u> will cause an arbitrary address to be freed
  - Inserted in the front of ListHeads[n]
  - Next pointer must be NULL to end the linked list
- In freeing a user-mode address, the attacker may control the contents of subsequent allocations
  - Must be made from the same process context

# PendingFrees Pointer Overwrite

Paged Pool Descriptor

| | |
|---|---|
| 0x0 | PoolType |
| 0x4 | PagedLock |
| … | |
| 0x100 | PendingFrees |
| 0x104 | PendingFreesDepth |
| 0x140 | ListHeads[512] |
| 0x140 + N*8 | Flink |
| | Blink |

Pool overflow

| Pool Header |
|---|
| Next |
| Data |

arbitrary address

Put in front of ListHeads[n] on free

| Pool Header |
|---|
| Flink |
| Blink |

Attacker controlled address is returned in requesting memory from ListHeads[n]

# PendingFrees Pointer Overwrite Steps

▸ Free a chunk to the deferred free list

▸ Overwrite the chunk's next pointer

  ▸ Or any of the deferred free list entries (32 in total)

▸ Trigger processing of the deferred free list

  ▸ Attacker controlled pointer freed to designated free list

▸ Force allocation of the controlled list entry

  ▸ Allocator returns user-mode address

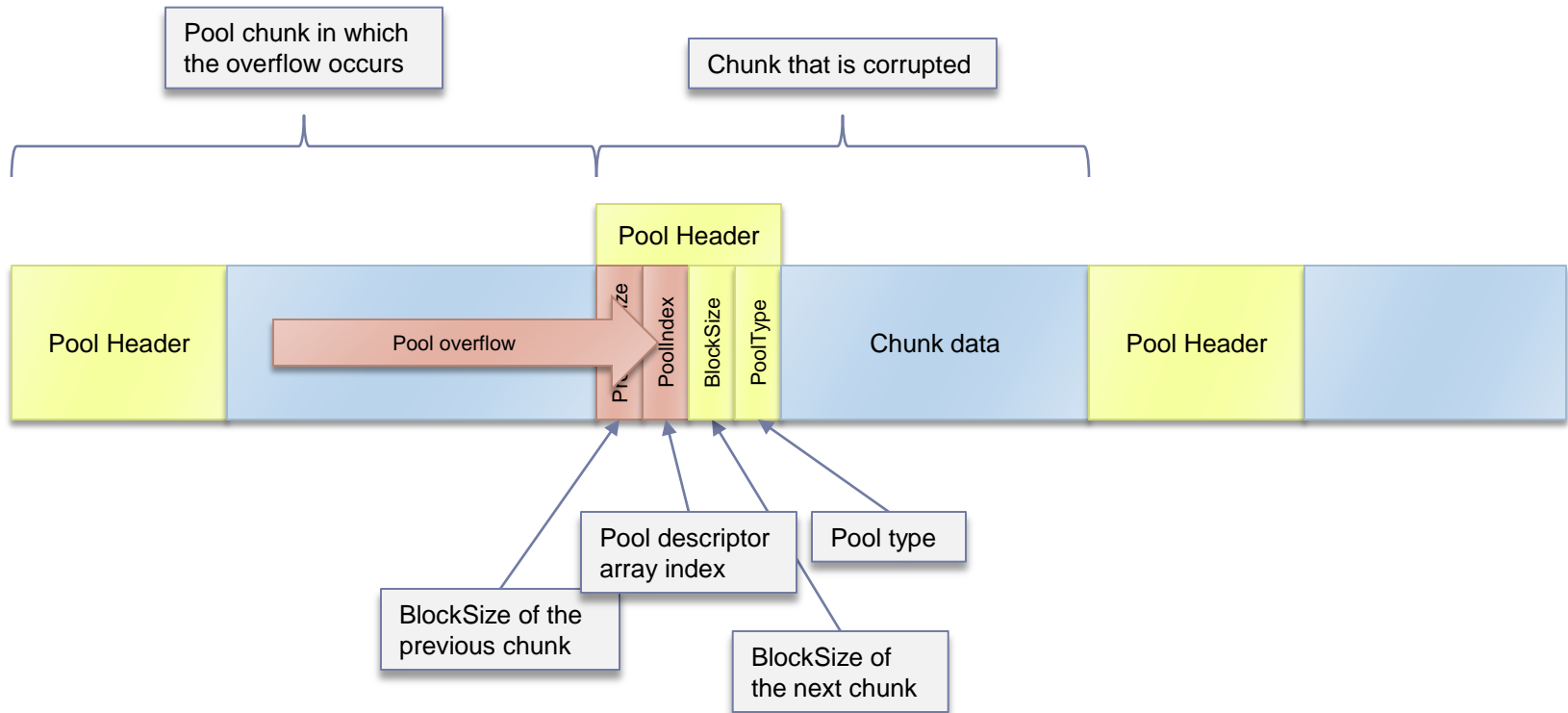▸ Corrupt allocated entry

▸ Trigger use of corrupted entry

▸

# PoolIndex Overwrite

- A pool chunk's <u>PoolIndex</u> denotes an index into the associated pool descriptor array
- For paged pools, PoolIndex always denotes an index into the **nt!ExpPagedPoolDescriptor** array
  - On checked builds, the index value is validated in a compare against **nt!ExpNumberOfPagedPools**
  - On free (retail) builds, the index is <u>not</u> validated
- For non-paged pools, PoolIndex denotes an index into **nt!ExpNonPagedPoolDescriptor** when there are <u>multiple NUMA nodes</u>
  - PoolIndex is <u>not</u> validated on free builds
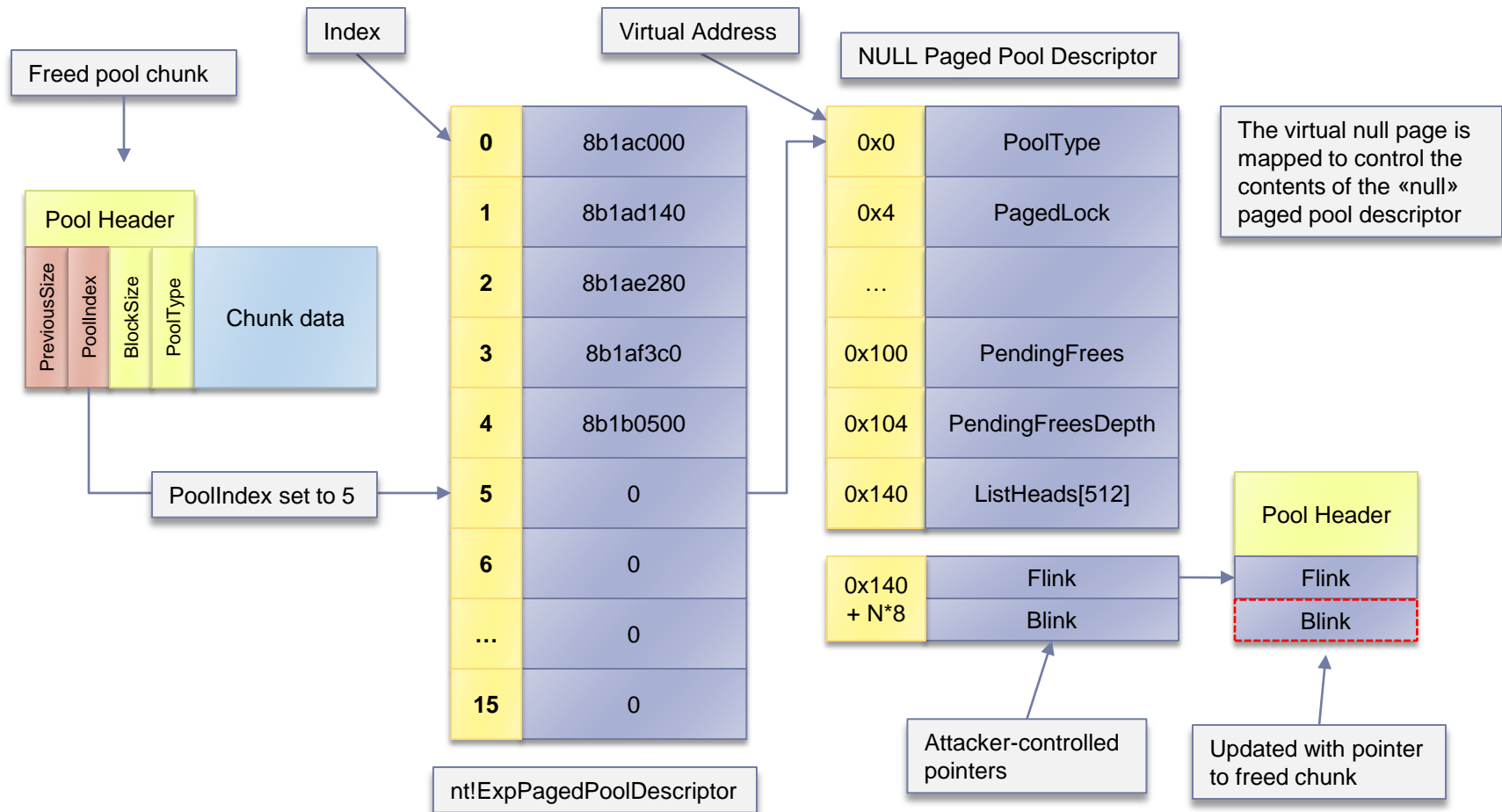
# PoolIndex Overwrite

# PoolIndex Overwrite

▶ A malformed PoolIndex may cause an allocated pool chunk to be <u>freed</u> to a null-pointer pool descriptor

  ▶ Controllable with null page allocation

  ▶ Requires a 2 byte pool overflow

▶ When <u>linking in</u> to a controlled pool descriptor, the attacker can write the address of the freed chunk to an arbitrary location

  ▶ No checks performed when "linking in"

  ▶ All ListHeads entries are fully controlled

  ▶ **ListHeads[n].Flink->Blink** = FreedChunk

▶

# PoolIndex Overwrite



Freed pool chunk

Pool Header

PreviousSize | PoolIndex | BlockSize | PoolType | Chunk data

PoolIndex set to 5

Index

Virtual Address

| 0 | 8b1ac000 |
| 1 | 8b1ad140 |
| 2 | 8b1ae280 |
| 3 | 8b1af3c0 |
| 4 | 8b1b0500 |
| 5 | 0 |
| 6 | 0 |
| … | 0 |
| 15 | 0 |

nt!ExpPagedPoolDescriptor

NULL Paged Pool Descriptor

| 0x0 | PoolType |
| 0x4 | PagedLock |
| … | |
| 0x100 | PendingFrees |
| 0x104 | PendingFreesDepth |
| 0x140 | ListHeads[512] |

| 0x140 + N*8 | Flink |
| | Blink |

The virtual null page is mapped to control the contents of the «null» paged pool descriptor

Pool Header

Flink

Blink

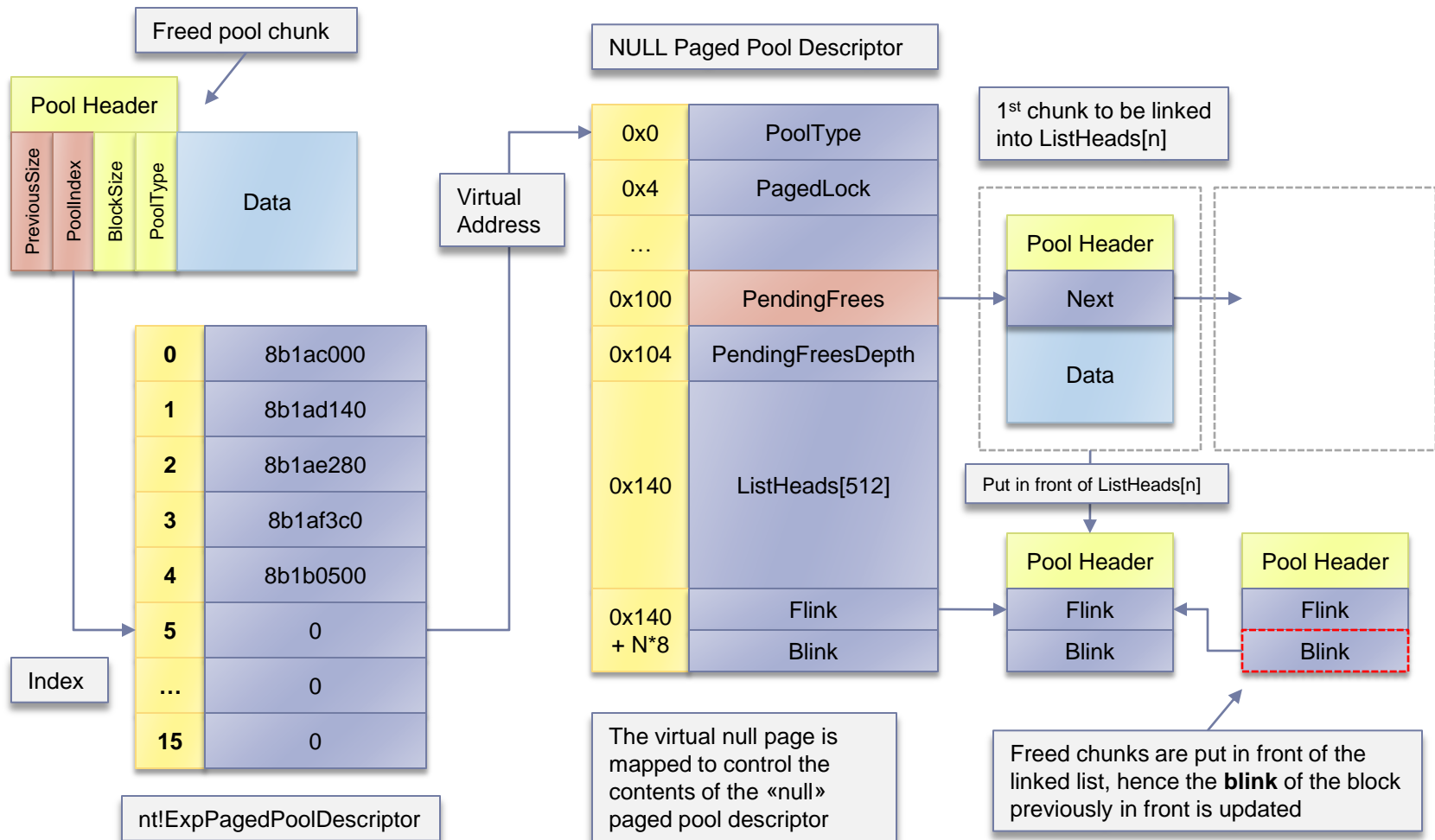Attacker-controlled pointers

Updated with pointer to freed chunk

# PoolIndex Overwrite (Delayed Frees)

- If delayed pool frees is enabled, the same effect can be achieved by creating a fake PendingFrees list
    - First entry should point to a user crafted chunk
- The **PendingFreeDepth** field of the pool descriptor should be >= 0x20 to trigger processing of the PendingFrees list
- The free algorithm of **ExDeferredFreePool** does basic validation on the crafted chunks
    - Coalescing / safe unlinking
    - The freed chunk should have busy bordering chunks

# PoolIndex Overwrite (Delayed Frees)

Freed pool chunk

Pool Header

| PreviousSize | PoolIndex | BlockSize | PoolType | Data |

NULL Paged Pool Descriptor

Virtual Address

nt!ExpPagedPoolDescriptor

Index

| 0 | 8b1ac000 |
|---|---|
| 1 | 8b1ad140 |
| 2 | 8b1ae280 |
| 3 | 8b1af3c0 |
| 4 | 8b1b0500 |
| 5 | 0 |
| … | 0 |
| 15 | 0 |

| 0x0 | PoolType |
|---|---|
| 0x4 | PagedLock |
| … | |
| 0x100 | PendingFrees |
| 0x104 | PendingFreesDepth |
| 0x140 | ListHeads[512] |
| 0x140 + N*8 | Flink |
| | Blink |

The virtual null page is mapped to control the contents of the «null» paged pool descriptor

$1^{st}$ chunk to be linked into ListHeads[n]

Pool Header

| Next |
| Data |

Put in front of ListHeads[n]

Pool Header

| Flink |
| Blink |

Pool Header

| Flink |
| Blink |

Freed chunks are put in front of the linked list, hence the **blink** of the block previously in front is updated

# PoolIndex Overwrite (Example)

▸ In controlling the PendingFrees list, a user-controlled virtual address (**eax**) can be written to an arbitrary destination address (**esi**)

▸ In turn, this can be used to corrupt function pointers used by the kernel to execute arbitrary code

```
eax=20000008 ebx=000001ff ecx=000001ff edx=00000538 esi=80808080 edi=00000000
eip=8293c943 esp=9c05fb20 ebp=9c05fb58 iopl=0        nv up ei pl nz na po nc
cs=0008  ss=0010  ds=0023  es=0023  fs=0030  gs=0000          efl=00010202

nt!ExDeferredFreePool+0x2e3:
8293c943 894604         mov     dword ptr [esi+4],eax ds:0023:80808084=????????
```

# PoolIndex Overwrite (!= PagedPool)

▸ The described technique can be used on any pool type if the chunk's PoolType is overwritten

  ▸ E.g. force a memory block to be part of a paged pool

  ▸ Requires also the BlockSize to be overwritten

▸ The BlockSize value must match the PreviousSize value of the next block

  ▸ FreedBlock->BlockSize = NextBlock->PreviousSize

  ▸ No problem if the size of the next block is known

  ▸ May also create a fake bordering chunk embedded in the corrupted chunk

# PoolIndex Overwrite (!= PagedPool)

**Before overflow**

Pool Header

Pool Header

PreviousSize | PoolIndex | BlockSize | PoolType

Pool Header

**After overflow**

Pool Header

Fake Header

Pool Header

Pool overflow

PreviousSize | PoolIndex | BlockSize | PoolType

PreviousSize | PoolIndex | BlockSize | PoolType

Pool Header

PagedPool

Must overflow BlockSize in order to get to PoolType (should be >= 0x20)

Create an embedded next pool chunk (busy) such that **Entry->BlockSize = NextEntry->PreviousSize**

# Quota Process Pointer Overwrite

▸ Quota charged pool allocations store a pointer to the associated process object

  ▸ **ExAllocatePoolWithQuotaTag(…)**

  ▸ x86: last four bytes of pool body

  ▸ x64: last eight bytes of pool header

▸ Upon freeing a pool chunk, the quota is released and the process object is dereferenced

  ▸ The object's reference count is decremented

▸ Overwriting the process object pointer could allow an attacker to free an in-use process object or corrupt arbitrary memory

# Quota Process Pointer Overwrite

# Quota Process Pointer Overwrite

▸ Quota information is stored in a EPROCESS_QUOTA_BLOCK structure

  ▸ Pointed to by the EPROCESS object

  ▸ Provides information on limits and how much quota is being used

▸ On free, the charged quota is returned by subtracting the size of the allocation from the quota used

  ▸ An attacker controlling the quota block pointer could decrement the value of an arbitrary address

  ▸ More on this later!

# Arbitrary Pointer Decrement

Address of executive process object controlled by the attacker

EPROCESS

EPROCESS_QUOTA_BLOCK

Usage counter decremented on free, for which the address is controlled by the attacker

Quota charged pool allocation (x86)

Pool Header

Pool overflow

Process Pointer

Pool Header

# Summary of Attacks

▸ Corruption of busy pool chunk
  ▸ BlockSize <= 0x20
    ▸ PoolIndex + PoolType/BlockSize Overwrite
    ▸ Quota Process Pointer Overwrite
  ▸ BlockSize > 0x20
    ▸ PoolIndex (+PoolType) Overwrite
    ▸ Quota Process Pointer Overwrite

▸ Corruption of free pool chunk
  ▸ BlockSize <= 0x20
    ▸ Lookaside Pointer Overwrite
  ▸ BlockSize > 0x20
    ▸ ListEntry Flink Overwrite / PendingFrees Pointer Overwrite

▸

# Case Studies

Kernel Pool Exploitation
on Windows 7

# Case Study Agenda

- ▶ Two pool overflow vulnerabilities
    - ▶ Both perceived as difficult to exploit
- ▶ CVE-2010-3939 (MS10-098)
    - ▶ Win32k CreateDIBPalette() Pool Overflow Vulnerability
- ▶ CVE-2010-1893 (MS10-058)
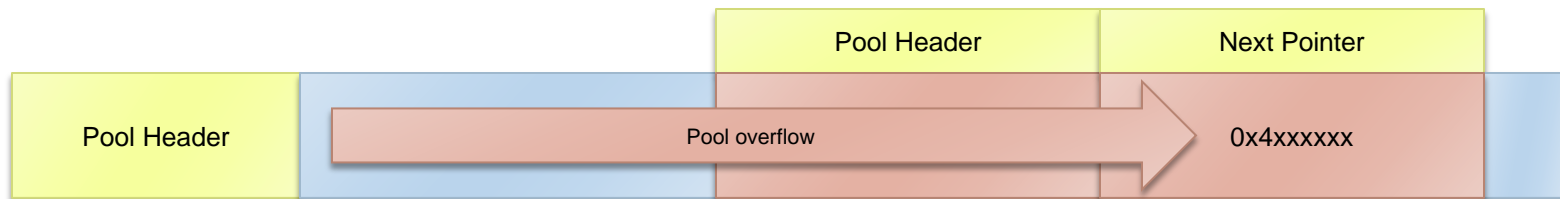    - ▶ Integer Overflow in Windows Networking Vulnerability

# CVE-2010-3939 (MS10-098)

▸ Pool overflow in win32k!CreateDIBPalette()

  ▸ Discovered by Arkon

▸ Function did not validate the number of color entries in the color table used by a bitmap

  ▸ BITMAPINFOHEADER.biClrUsed

▸ Every fourth byte of the overflowing buffer was set to 0x4

  ▸ Can only reference 0x4xxxxxx addresses (user-mode)
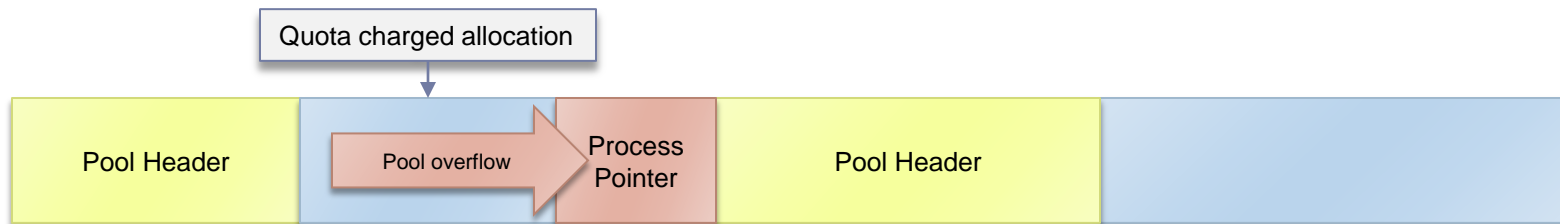
  ▸ PoolType is always set to NonPaged

| Pool Header | | Pool overflow → | Pool Header | | | | |
|---|---|---|---|---|---|---|---|
| | | | X | X | X | 0x2 | |

PoolType = NonPaged | InUse
(0x2 due to bit alignment of field on x86)

# CVE-2010-3939 (MS10-098)

▸ The attacker could coerce the pool allocator to return a user-mode pool chunk

  ▸ ListEntry Flink Overwrite

  ▸ Lookaside Overwrite

▸ Requires the kernel pool to be cleaned up in order for execution to continue safely

  ▸ Repair/remove broken linked lists

| Pool Header | Next Pointer |
|:---:|:---:|
| | |

| Pool Header | Pool overflow | 0x4xxxxxx |
|:---:|:---:|:---:|

# CVE-2010-3939 (MS10-098)

▸ **Vulnerable buffer is also quota charged**

▸ Can overwrite the process object pointer (x86)

▸ No pool chunks are corrupted (clean!)

▸ **Tactic: Decrement the value of a kernel-mode window object procedure pointer**

▸ Trigger the vulnerability n-times until it points to user-mode memory and call the procedure

Quota charged allocation

| Pool Header | Pool overflow | Process Pointer | Pool Header | |

# CVE-2010-3939 (MS10-098)

- ## Quota Process Pointer Overwrite
  - Demo

# CVE-2010-1893 (MS10-058)

- Integer overflow in tcpip!IppSortDestinationAddresses()
  - Discovered by Matthieu Suiche
  - Affected Windows 7/2008 R2 and Vista/2008
- Function did not use safe-int functions consistently
  - Could result in an undersized buffer allocation, subsequently leading to a pool overflow

# IppSortDestinationAddresses()

▸ Sorts a list of IPv6 and IPv4 destination addresses

  ▸ Each address is a SOCKADDR_IN6 record

▸ Reachable from user-mode by calling WSAIoctl()

  ▸ Ioctl: SIO_ADDRESS_LIST_SORT

  ▸ Buffer: SOCKET_ADDRESS_LIST structure

▸ Allocates buffer for the address list

  ▸ iAddressCount * sizeof(SOCKADDR_IN6)

  ▸ No overflow checks

```
typedef struct _SOCKET_ADDRESS_LIST {
  INT               iAddressCount;
  SOCKET_ADDRESS Address[1];
} SOCKET_ADDRESS_LIST, *PSOCKET_ADDRESS_LIST;
```
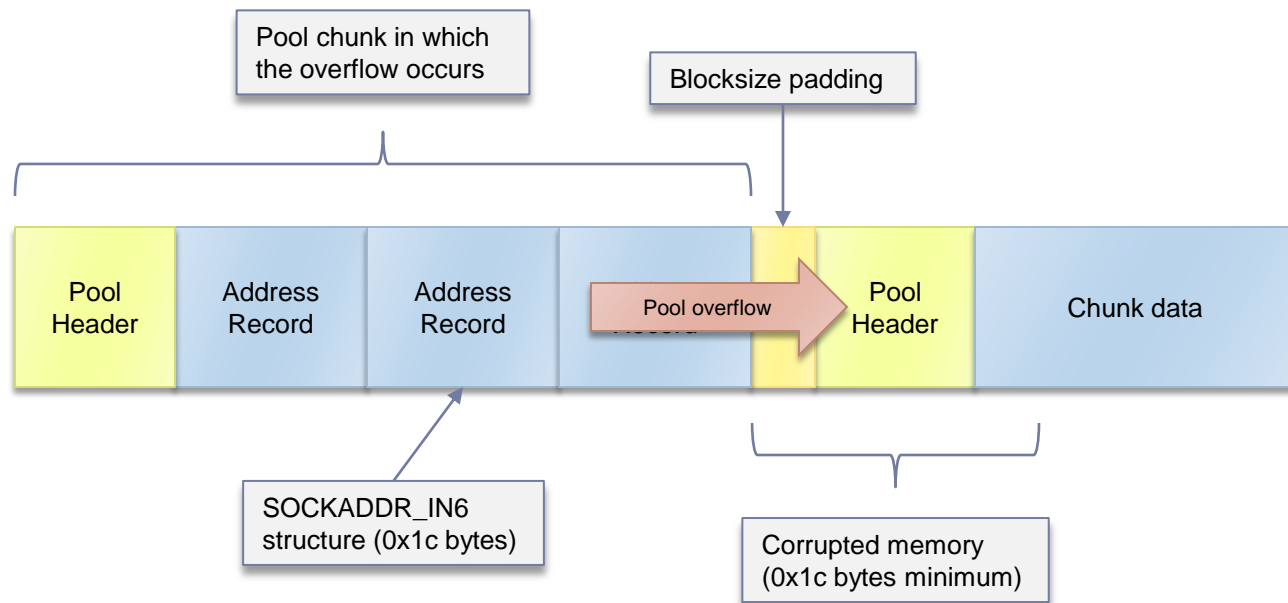
# IppFlattenAddressList()

- Copies the user provided address list to the allocated kernel pool chunk

- An undersized buffer could result in a pool overflow
  - Overflows the next pool chunk with the size of an address structure (0x1c bytes)

- Stops copying records if the size != 0x1c or the protocol family != AF_INET6 (0x17)
  - Possible to avoid trashing the kernel pool completely

- The protocol check is done after the memcpy()
  - We can overflow using any combination of bytes

# Pool Overflow

# Exploitation Tactics

- Can use the PoolIndex attack to extend the pool overflow to an arbitrary memory write
  - Must overwrite a busy chunk
- Overwritten chunk must be freed to ListHeads lists
  - BlockSize > 0x20
  - Or… fill the lookaside list
- To overflow the desired pool chunk, we must defragment and manipulate the kernel pool
  - Allocate chunks of the same size
  - Create "holes" by freeing every other chunk

# Kernel Pool Manipulation (1)

▸ ## What do we use to fill the pool ?

- ▸ Depends on the pool type
- ▸ Should be easy to allocate and free

▸ ## NonPaged

- ▸ Kernel objects introduce low overhead
  - ▸ NtAllocateReserveObject
  - ▸ NtCreateSymbolicLinkObject

▸ ## PagedPool

- ▸ Unicode strings (e.g. object properties)

# Kernel Pool Manipulation (2)

▶ Create holes by freeing every second allocation

  ▶ The vulnerable buffer is later allocated in one of these holes

▶ Freeing the remaining allocations after triggering the vulnerability mounts the PoolIndex attack

```
kd> !pool @eax
 Pool page 976e34c8 region is Nonpaged pool

 976e32e0 size: 60 previous size: 60 (Allocated) IoCo (Protected)
 976e3340 size: 60 previous size: 60 (Free) IoCo
 976e33a0 size: 60 previous size: 60 (Allocated) IoCo (Protected)
 976e3400 size: 60 previous size: 60 (Free) IoCo
 976e3460 size: 60 previous size: 60 (Allocated) IoCo (Protected)
*976e34c0 size: 60 previous size: 60 (Allocated) *Ipas
        Pooltag Ipas : IP Buffers for Address Sort, Binary : tcpip.sys
 976e3520 size: 60 previous size: 60 (Allocated) IoCo (Protected)
 976e3580 size: 60 previous size: 60 (Free) IoCo
 976e35e0 size: 60 previous size: 60 (Allocated) IoCo (Protected)
 976e3640 size: 60 previous size: 60 (Free) IoCo
```

# CVE-2010-1893 (MS10-058)

- Kernel pool manipulation + PoolIndex overwrite
  - Demo

# Kernel Pool Hardening

Kernel Pool Exploitation
on Windows 7

# ListEntry Flink Overwrites

▶ Can be addressed by properly validating the flink and blink of the <u>chunk being unlinked</u>
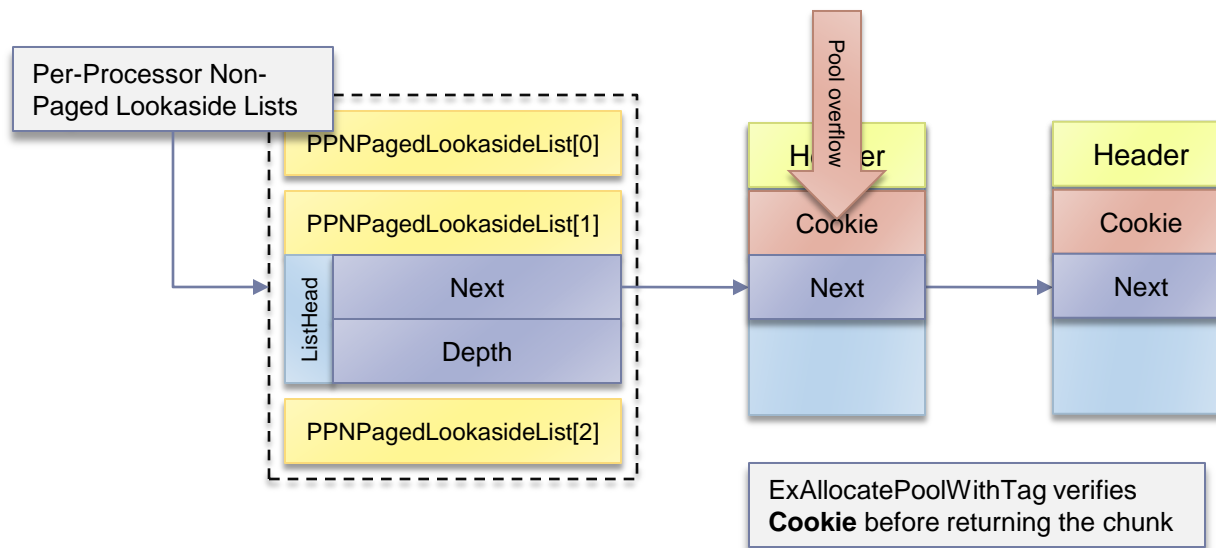
  ▶ Yep, that's it...

# Lookaside Pointer Overwrites

▸ Lookaside lists are inherently insecure

  ▸ Unchecked embedded pointers

▸ All pool chunks must reserve space for at least the size of a LIST_ENTRY structure

  ▸ Two pointers (flink and blink)

▸ Chunks on lookaside lists only store a single pointer

  ▸ Could include a cookie for protecting against pool overflows

▸ Cookies could also be used by PendingFrees list entries

▸

# Lookaside Pool Chunk Cookie

# PoolIndex Overwrites

- Can be addressed by validating the PoolIndex value before freeing a pool chunk
  - E.g. is PoolIndex > nt!ExpNumberOfPagedPools ?
- Also required the NULL-page to be mapped
  - Could deny mapping of this address in non-privileged processes
  - Would probably break some applications (e.g. 16-bit WOW support)

# Quota Process Pointer Overwrites

▸ Can be addressed by encoding or obfuscating the process pointer

  ▸ E.g. XOR'ed with a constant unknown to the attacker

▸ Ideally, no pointers should be embedded in pool chunks

  ▸ Pointers to structures that are written to can easily be leveraged to corrupt arbitrary memory

# Conclusion

Kernel Pool Exploitation
on Windows 7

# Future Work

- Pool content corruption
  - Object function pointers
  - Data structures

- Remote kernel pool exploitation
  - Very situation based
  - Kernel pool manipulation is hard
  - Attacks that rely on null page mapping are infeasible

- Kernel pool manipulation
  - Becomes more important as generic vectors are addressed

# Conclusion

▸ The kernel pool was designed to be fast
- ▸ E.g. no pool header obfuscation

▸ In spite of safe unlinking, there is still a big window of opportunity in attacking pool metadata
- ▸ Kernel pool manipulation is the key to success

▸ Attacks can be addressed by adding simple checks or adopting exploit prevention features from the userland heap
- ▸ Header integrity checks
- ▸ Pointer encoding
- ▸ Cookies

▸

# Questions ?

- Email: kernelpool@gmail.com
- Blog: http://mista.nu/blog
- Twitter: @kernelpool

# References

- **SoBeIt[2005]** – SoBeIt
  How to exploit Windows kernel memory pool,
  X'con 2005

- **Kortchinsky[2008]** – Kostya Kortchinsky
  Real-World Kernel Pool Exploitation,
  SyScan 2008 Hong Kong

- **Mxatone[2008]** – mxatone
  Analyzing Local Privilege Escalations in win32k,
  Uninformed Journal, vol. 10 article 2

- **Beck[2009]** – Peter Beck
  Safe Unlinking in the Kernel Pool,
  Microsoft Security Research & Defense (blog)