# GDC

- Cold, Hard Cache
  Insomniac Games' Cache Simulator

- Andreas Fredriksson
  Lead Engine Programmer, Insomniac Games

# Hi, I'm Andreas

- I lead the Core tools and infrastructure team at Insomniac

# Hi, I'm Andreas

- I lead the Core tools and infrastructure team at Insomniac

- Today's topic: Insomniac Games' CacheSim

  - Custom tooling for measuring cache effectiveness

# Hi, I'm Andreas

- I lead the Core tools and infrastructure team at Insomniac

- Today's topic: Insomniac Games' CacheSim
  - Custom tooling for measuring cache effectiveness
  - Excited to be open sourcing this and sharing with you all today
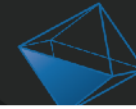
# Cache and memory sizes, visually

# Cache and memory sizes, visually



4GB DRAM

# Cache and memory sizes, visually

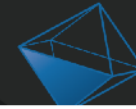4GB DRAM          2 MB L2

# Cache and memory sizes, visually

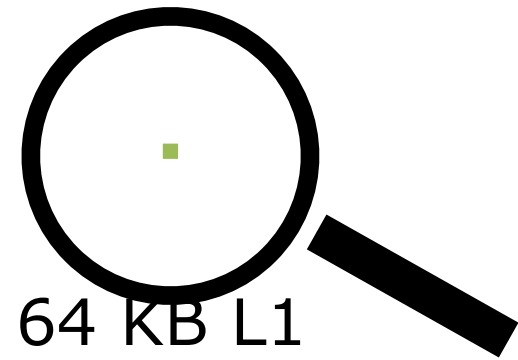4GB DRAM                    2 MB L2                    64 KB L1
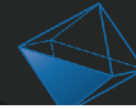
# Cache and memory sizes, visually

4GB DRAM              2 MB L2              64 KB L1

# Cache, memory access speeds in cycles



DRAM    L2    D1

# Background

- Cache orders of magnitude faster and smaller than RAM
  - What you put into them makes a huge difference

# Background

- Cache orders of magnitude faster and smaller than RAM

  - What you put into them makes a huge difference

- Memory operations are extremely easy to add to a program

  - Costs are hidden and non-obvious

# Background

- Cache orders of magnitude faster and smaller than RAM

  - What you put into them makes a huge difference

- Memory operations are extremely easy to add to a program

  - Costs are hidden and non-obvious

- We desperately need actionable data on access patterns

  - Not a wealth of options for the performance-aware programmer

# Sampling profilers

- They've basically won – most profilers are sample based
  - Great for many workflows
  - Leverage CPU features to gather HW stats about cache

# Sampling profilers

- They've basically won – most profilers are sample based
    - Great for many workflows
    - Leverage CPU features to gather HW stats about cache
- Limitation: Only every N instructions are sampled (N is large)

# Sampling profilers

- They've basically won – most profilers are sample based
  - Great for many workflows
  - Leverage CPU features to gather HW stats about cache
- Limitation: Only every N instructions are sampled (N is large)
- Not ideal for smaller, "bursty" workloads
  - Statistical means less reproducible for smaller things
  - Point you in the right direction, but that's about it

# Outside the sampling space

- Cachegrind – part of Valgrind
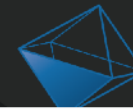  - Simulate a cache based on the program's instruction stream

# Outside the sampling space

- Cachegrind – part of Valgrind
  - Simulate a cache based on the program's instruction stream
- Pros:
  - Extremely thorough – every memory access is simulated

# Outside the sampling space

- Cachegrind – part of Valgrind
  - Simulate a cache based on the program's instruction stream
- Pros:
  - Extremely thorough – every memory access is simulated
- Cons:
  - Linux only
  - All or nothing
  - Extremely slow to get to a point of interest

# Outside the sampling space

- Cachegrind – part of Valgrind
  - Simulate a cache based on the program's instruction stream
- Pros:
  - Extremely thorough – every memory access is simulated
- Cons:
  - Linux only
  - All or nothing
  - Extremely slow to get to a point of interest
- Also – a vendor-specific tool from prevgen was awesome
  - Could do simulated captures on demand for a short time

# Why you want cache simulation tooling

```cpp
void PushBuffer::SetTextureAssets(uint32_t start_slot, uint32_t slot_count,
                                  const TextureAsset** textures_assets, uint32_t slot_mask, uint32_t hq_mask)
{
  // …
  for(uint32_t itex = 0; itex < slot_count; ++itex, tex_unit_test <<= 1)
  {
    Texture const*         tex            = (slot_mask & tex_unit_test ) ? textures_assets[itex]->GetTexture() : NULL;
    // …
    if(tex != NULL)
    {
      view = (textures_assets[itex]->GetFormatFlags() & TextureFormatFlags::kIsCube) ? tex->m_View : tex->m_ViewArray;
      samp = (hq_mask & tex_unit_test) ? textures_assets[itex]->GetAnisoSampler() : tex->m_SamplerState;
    }
    // …
```

# Why you want cache simulation tooling

```cpp
void PushBuffer::SetTextureAssets(uint32_t start_slot, uint32_t slot_count,
                                  const TextureAsset** textures_assets, uint32_t slot_mask, uint32_t hq_mask)
{
  // …
  for(uint32_t itex = 0; itex < slot_count; ++itex, tex_unit_test <<= 1)
  {
    Texture const*          tex            = (slot_mask & tex_unit_test ) ? textures_assets[itex]->GetTexture() : NULL;
    // …
    if(tex != NULL)
    {
      view = (textures_assets[itex]->GetFormatFlags() & TextureFormatFlags::kIsCube) ? tex->m_View : tex->m_ViewArray;
      samp = (hq_mask & tex_unit_test) ? textures_assets[itex]->GetAnisoSampler() : tex->m_SamplerState;
    }
    // …
```

2800 L2 misses in a frame

UBM

# It's the small things

- ## A member had moved cache lines
  - ### Accessing the 16-bit FormatFlags field was now a guaranteed L2 miss

```
//HOT: Cache Line 1
Texture        m_Texture;
TextureAsset*  m_Default;

uint32_t       m_ResidentSize;
uint32_t       m_StreamSize;


uint32_t       m_AnisoSamplerIndex;
... other fields...


// COLD: Cache Line 2

uint16_t       m_FormatFlags;
uint16_t       m_Flags;

int16_t        m_TopHeight;
int16_t        m_TopWidth;
... other fields ...
```

UBM

# It's the small things

- ## A member had moved cache lines
  - ### Accessing the 16-bit FormatFlags field was now a guaranteed L2 miss

```
//HOT: Cache Line 1
Texture        m_Texture;
TextureAsset*  m_Default;

uint32_t       m_ResidentSize;
uint32_t       m_StreamSize;


uint32_t       m_AnisoSamplerIndex;
... other fields...


// COLD: Cache Line 2

uint16_t       m_FormatFlags;
uint16_t       m_Flags;

int16_t        m_TopHeight;
int16_t        m_TopWidth;
... other fields ...
```

# It's the small things

- A member had moved cache lines
  - Accessing the 16-bit FormatFlags field was now a guaranteed L2 miss

- Fix: Literally swap two lines in header

```cpp
//HOT: Cache Line 1
Texture        m_Texture;
TextureAsset*  m_Default;

uint32_t       m_ResidentSize;
uint32_t       m_StreamSize;

uint32_t       m_AnisoSamplerIndex;
... other fields...

// COLD: Cache Line 2

uint16_t       m_FormatFlags;
uint16_t       m_Flags;

int16_t        m_TopHeight;
int16_t        m_TopWidth;
... other fields ...
```

# It's the small things

- A member had moved cache lines
  - Accessing the 16-bit FormatFlags field was now a guaranteed L2 miss
- Fix: Literally swap two lines in header
- 150-250 us savings depending on view

```
//HOT: Cache Line 1
Texture        m_Texture;
TextureAsset*  m_Default;

uint32_t       m_ResidentSize;
uint32_t       m_StreamSize;

uint32_t       m_AnisoSamplerIndex;
... other fields...

// COLD: Cache Line 2

uint16_t       m_FormatFlags;
uint16_t       m_Flags;

int16_t        m_TopHeight;
int16_t        m_TopWidth;
... other fields ...
```

# It's the small things

- A member had moved cache lines
  - Accessing the 16-bit FormatFlags field was now a guaranteed L2 miss

- Fix: Literally swap two lines in header

- 150-250 us savings depending on view

- Time investment: 30 minutes

```
//HOT: Cache Line 1
Texture       m_Texture;
TextureAsset* m_Default;

uint32_t      m_ResidentSize;
uint32_t      m_StreamSize;

uint32_t      m_AnisoSamplerIndex;
... other fields...

// COLD: Cache Line 2

uint16_t      m_FormatFlags;
uint16_t      m_Flags;

int16_t       m_TopHeight;
int16_t       m_TopWidth;
... other fields ...
```

# More reasons you want cache tooling

- Small utilities miss cache a lot sometimes
  - ```
    Vec3 SceneObject::GetPosition() { return m_ObjToWorld.v[3]; }
    ```

# More reasons you want cache tooling

- Small utilities miss cache a lot sometimes

  - `Vec3 SceneObject::GetPosition() { return m_ObjToWorld.v[3]; }`

- Can't optimize a single return statement! /tableflip

  - But you can optimize callers of that function to be less naïve

# More reasons you want cache tooling

- ## Small utilities miss cache a lot sometimes
  - `Vec3 SceneObject::GetPosition() { return m_ObjToWorld.v[3]; }`
- ## Can't optimize a single return statement! /tableflip
  - But you can optimize callers of that function to be less naïve
- ## Need to track call stacks as well, pass blame up the stack

# More reasons you want cache tooling

- Small utilities miss cache a lot sometimes

  - `Vec3 SceneObject::GetPosition() { return m_ObjToWorld.v[3]; }`

- Can't optimize a single return statement! /tableflip

  - But you can optimize callers of that function to be less naïve

- Need to track call stacks as well, pass blame up the stack

- Found that 12k out of 14k misses in GetPosition() came from one gameplay system ☺

# Who's calling SceneObject::GetPosition() with cold data?

```cpp
uint32_t num_groups = 0;
// Get a ton of stuff to work on
ComponentHandle *groups = g_PlacedPedestrianSystem.GetGroups(num_groups);

for (uint32_t idx = 0; idx < num_groups; ++idx)
{
  PlacedPedestrianGroup *group = (PlacedPedestrianGroup*)groups[idx].Resolve();
  if(!group)
    continue;

  const Vec3 &group_pos = group->GetActorPosition();
  // ...
}
```

# Who's calling SceneObject::GetPosition() with cold data?

```cpp
uint32_t num_groups = 0;
// Get a ton of stuff to work on
ComponentHandle *groups = g_PlacedPedestrianSystem.GetGroups(num_groups);

for (uint32_t idx = 0; idx < num_groups; ++idx)
{
  PlacedPedestrianGroup *group = (PlacedPedestrianGroup*)groups[idx].Resolve();
  if(!group)
    continue;

  const Vec3 &group_pos = group->GetActorPosition();
  // ...
}
```

~12,000 L2 misses in a frame

# Don't ask about things you already know

- About 1 ms/frame on PC
  - Thousands of (random order) items being processed every frame

# Don't ask about things you already know

- About 1 ms/frame on PC
  - Thousands of (random order) items being processed every frame
- Digging in: They never move

# Don't ask about things you already know

- About 1 ms/frame on PC

  - Thousands of (random order) items being processed every frame

- Digging in: They never move

- Fix: Keep local cache of positions in parallel array

  - Avoiding the GetPosition() call entirely

# Don't ask about things you already know

- About 1 ms/frame on PC

  - Thousands of (random order) items being processed every frame

- Digging in: They never move

- Fix: Keep local cache of positions in parallel array

  - Avoiding the GetPosition() call entirely

- Result: 650 us/frame saved

# Don't ask about things you already know

- About 1 ms/frame on PC

  - Thousands of (random order) items being processed every frame

- Digging in: They never move

- Fix: Keep local cache of positions in parallel array

  - Avoiding the GetPosition() call entirely

- Result: 650 us/frame saved

- Time investment: 2 hours

# Instructions, instructions

- All the data we want is right there, in the instruction stream
  - We "just" need a way to look at each instruction as it executes

# Instructions, instructions

- All the data we want is right there, in the instruction stream

  - We "just" need a way to look at each instruction as it executes

- Plan:

  - Flip a switch upon reaching a point of interest

  - Trace every instruction (somehow)

  - Update a simulated cache for each memory access

  - Turn off trace (say at end of frame) & report!

# Instructions, instructions

- All the data we want is right there, in the instruction stream
  - We "just" need a way to look at each instruction as it executes
- Plan:
  - Flip a switch upon reaching a point of interest
  - Trace every instruction (somehow)
  - Update a simulated cache for each memory access
  - Turn off trace (say at end of frame) & report!
- But how?

# Where to start?

- First: It really helps to have an understanding boss

# Where to start?

- First: It really helps to have an understanding boss

# Where to start?

- First: It really helps to have an understanding boss

- Pitched 2 week project to dig in..

  - Subject: "I can see crazy town from here"

# Where to start?

- First: It really helps to have an understanding boss

- Pitched 2 week project to dig in..
  - Subject: "I can see crazy town from here"

- Approved!
  - Ok, let's start pitching a tent outside crazy town

UBM

# Reasoning about instructions

- Binary instrumentation frameworks exist (off the shelf)
  - DynamoRIO, Intel PIN, others

# Reasoning about instructions

- Binary instrumentation frameworks exist (off the shelf)
  - DynamoRIO, Intel PIN, others
- Quickly discarded this approach
  - Massive performance problems instrumenting a AAA game executable

# Reasoning about instructions

- Binary instrumentation frameworks exist (off the shelf)

  - DynamoRIO, Intel PIN, others

- Quickly discarded this approach

  - Massive performance problems instrumenting a AAA game executable

- Could have value for other things in our space

  - More guided dynamic instrumentation without code changes

  - "How often is this value zero at this spot?"

  - "What are the min/max input values to this function?"

# Idea #1

- Somehow write void TraceFunction(func_ptr)

# Idea #1

- Somehow write void TraceFunction(func_ptr)
- Somehow, for each instruction:
  - Disassemble the instruction
  - Find memory derefs, update a simulated cache
  - Copy instruction to temp buffer, run in isolation
- Sounds easy!

# Idea #1 – not going anywhere

- Branches need to be special handled

# Idea #1 – not going anywhere

- Branches need to be special handled
- Implicit uses of RIP (instruction pointer) are everywhere

# Idea #1 – not going anywhere

- Branches need to be special handled

- Implicit uses of RIP (instruction pointer) are everywhere

- Win64 exception handling has rules we're violating

  - OutputDebugString uses exceptions..

# Idea #1 – not going anywhere

- Branches need to be special handled

- Implicit uses of RIP (instruction pointer) are everywhere

- Win64 exception handling has rules we're violating

  - OutputDebugString uses exceptions..

- Super intrusive – need a top-level call to our trace function

  - On every thread

# Suddenly: EFLAGS

# Suddenly: EFLAGS

# Idea #2: Leveraging EFLAGS

- Single stepping is a CPU feature
  - It's how F11 in the debugger works
  - Set TRAP bit in EFLAGS

# Idea #2: Leveraging EFLAGS

- Single stepping is a CPU feature
  - It's how F11 in the debugger works
  - Set TRAP bit in EFLAGS
- Routes as an exception through the Windows SEH machinery
  - You can install a handler for it!

# Idea #2: Leveraging EFLAGS

- Single stepping is a CPU feature
  - It's how F11 in the debugger works
  - Set TRAP bit in EFLAGS
- Routes as an exception through the Windows SEH machinery
  - You can install a handler for it!
- But how do you install SEH handlers for all threads?
  - Vectored exception handlers

# Revised plan of attack

- To start tracing:
  - Install a VEH to filter TRAP exceptions
  - Set TF EFLAGS bit for all threads we want to capture

# Revised plan of attack

- To start tracing:
  - Install a VEH to filter TRAP exceptions
  - Set TF EFLAGS bit for all threads we want to capture
- In the handler:
  - Disassemble instruction, find memory operands
  - Update cache simulation
  - Re-set the TF bit before leaving to keep tracing

# Revised plan of attack

- To start tracing:
  - Install a VEH to filter TRAP exceptions
  - Set TF EFLAGS bit for all threads we want to capture

- In the handler:
  - Disassemble instruction, find memory operands
  - Update cache simulation
  - Re-set the TF bit before leaving to keep tracing

- To stop tracing:
  - Set some flag and (ultimately) remove the VEH

# Every plan has problems

- Good news: It basically works

# Every plan has problems

- Good news: It basically works

- Problem #1: The debugger is really unhappy

  - "What? You want to break in?"

  - Solution: Run detached

# Every plan has problems

- Good news: It basically works

- Problem #1: The debugger is really unhappy

  - "What? You want to break in?"

  - Solution: Run detached

- Problem #2: Massive amounts of deadlocks in ntdll.dll

  - Hanging on contended SRW lock protecting the VEH dispatch list

  - Threads waiting on wakeups for locks

  - But no one owns the lock, so it can never wake up

# Deadlock woes

- Using SEH way, way more than anyone had anticipated at MS

    - Every thread, every instruction will exercise the exception handling

- Nothing was obviously wrong in the code

- Suspect problem is reentrantly messing with critical sections

# Solution: Disable the locking code in ntdll

- VEH are an exotic feature, typically no handlers are installed
  - This is a debugging feature, not something we're shipping to players
  - Just smash ntdll!RtlpCallVectoredHandlers with a jump to our handler
- It's ugly, but it gets the job done
- Also: take no OS locks internally, just spinlocks

[cue evil laugh]

memes.com

# It works!

- We can inspect instructions one by one!

# It works!

- ## We can inspect instructions one by one!

- ## Need a disassembler that knows about memory operands

  - ### Radare2 fork of udis86 fit the bill

# It works!

- We can inspect instructions one by one!

- Need a disassembler that knows about memory operands

  - Radare2 fork of udis86 fit the bill

```
uintptr_t rip = ExcInfo->ContextRecord->Rip;

ud_set_input_buffer(ud, (const uint8_t*) rip, 16);
ud_set_pc(ud, rip);
int ilen = ud_disassemble(ud);
GenerateMemoryAccesses(core_index, ud, rip, ilen, ExcInfo->ContextRecord);
```

# Generating memory accesses

- Easy, right? Just look at memory operands
  - `mov dword ptr [rax], ebx` => write 4 bytes at `rax`

# Generating memory accesses

- Easy, right? Just look at memory operands

  - `mov dword ptr [rax], ebx` => write 4 bytes at `rax`

- Yeah, in theory..

  - Lots of things in x64 have memory operands

  - Some access memory and don't have memory operands!

# Lots of special cases to consider

- String instructions, e.g. LODSB, MOVSD (implicit RSI/RDI accesses)
- Stack push/pop
- CALL, RET (also write/read the stack respectively)
- LEA – super common, has mem operand but doesn't touch memory
- Crazy "long nop" instructions can have memory operands
- FXSTOR/FXRSTOR
- Prefetches and non-temporal loads/stores

## Poking the cache simulation

```cpp
// Generate I-cache traffic.
CacheSim::AccessResult r = g_Cache.Access(core_index, rip, ilen, CacheSim::kCodeRead);
stats->m_Stats[r] += 1;
```

# Poking the cache simulation

```cpp
// Generate I-cache traffic.
CacheSim::AccessResult r = g_Cache.Access(core_index, rip, ilen, CacheSim::kCodeRead);
stats->m_Stats[r] += 1;

// Generate D-cache traffic.
for (int i = 0; i < read_count; ++i)
{
  CacheSim::AccessResult r = g_Cache.Access(core_index, reads[i].ea, reads[i].sz, CacheSim::kRead);
  stats->m_Stats[r] += 1;
}
```

# Poking the cache simulation

```cpp
// Generate I-cache traffic.
CacheSim::AccessResult r = g_Cache.Access(core_index, rip, ilen, CacheSim::kCodeRead);
stats->m_Stats[r] += 1;

// Generate D-cache traffic.
for (int i = 0; i < read_count; ++i)
{
  CacheSim::AccessResult r = g_Cache.Access(core_index, reads[i].ea, reads[i].sz, CacheSim::kRead);
  stats->m_Stats[r] += 1;
}

for (int i = 0; i < write_count; ++i)
{
  CacheSim::AccessResult r = g_Cache.Access(core_index, writes[i].ea, writes[i].sz, CacheSim::kWrite);
  stats->m_Stats[r] += 1;
}
```

# Can model a set-assoc cache as a 2D array

Take apart input address:

0100011...1111101001...110011

Sets

Ways

# Can model a set-assoc cache as a 2D array

Take apart input address:

0100011...1111101001...110011

Locate the set

| ADDRESS | ADDRESS | ADDRESS | ... |

Sets

Ways

# Can model a set-assoc cache as a 2D array

Take apart input address:

0100011...1111101001...110011

Locate the set

| ADDRESS | ADDRESS | ADDRESS | ... |

Compare addr against each way to see if cached

Sets

Ways

# Simulating a Jaguar cache

- Console Jaguar has 2 modules
  - Each modules has a shared L2, 4 cores
  - Each core has its own D1, I1 caches

# Simulating a Jaguar cache

- Console Jaguar has 2 modules
  - Each modules has a shared L2, 4 cores
  - Each core has its own D1, I1 caches
- Jaguar cache is *inclusive* (lines in D1/I1 must also exist in L2)

# Simulating a Jaguar cache

- Console Jaguar has 2 modules
  - Each modules has a shared L2, 4 cores
  - Each core has its own D1, I1 caches
- Jaguar cache is *inclusive* (lines in D1/I1 must also exist in L2)
- Map our set associativity for our array structures
  - I1: 512 lines (32 KB), 2 ways, 256 sets
  - D1: 512 lines (32 KB), 8 ways, 64 sets
  - L2: 32,768 lines (2 MB), 16 ways, 2,048 sets

# Defining our caches

```cpp
// Simulating a Jaguar cache

//                          size in byte   |   assoc
using JaguarD1 = Cache<      32 * 1024,        8>;
using JaguarI1 = Cache<      32 * 1024,        2>;
using JaguarL2 = Cache<2 * 1024 * 1024,       16>;
```

# Defining our caches

```cpp
// Simulating a Jaguar cache

//                        size in byte  |  assoc
using JaguarD1 = Cache<      32 * 1024,       8>;
using JaguarI1 = Cache<      32 * 1024,       2>;
using JaguarL2 = Cache<2 * 1024 * 1024,      16>;


struct JaguarModule
{
  JaguarD1     m_CoreD1[4];   // 4 cores per module, with private D1 & I1
  JaguarI1     m_CoreI1[4];
  JaguarL2     m_Level2;       // A shared L2
  JaguarModule* m_OtherModule; // Pointer to other module for invalidations
};
```

# Updating the cache, in pseudocode

```
For each cache line accessed:
```

# Updating the cache, in pseudocode

```
For each cache line accessed:

    If we're writing:
        Kick line out of every other core
        Kick line out of other module's L2
```

# Updating the cache, in pseudocode

```
For each cache line accessed:

    If we're writing:
        Kick line out of every other core
        Kick line out of other module's L2

    Hit1 = Lookup+Record Line in D1/I1
    Hit2 = Lookup+Record Line in L2
```

# Updating the cache, in pseudocode

```
For each cache line accessed:

    If we're writing:
      Kick line out of every other core
      Kick line out of other module's L2

    Hit1 = Lookup+Record Line in D1/I1
    Hit2 = Lookup+Record Line in L2


    If Hit1 && Hit2:
      return kL1Hit
    Else If Hit2:
      return kL2Hit
    Else:
      return kL2Miss
```

UBM

# Running it

- Hook up trace machinery to keyboard shortcut in main loop
  - Automatically disable at end of frame

# Running it

- Hook up trace machinery to keyboard shortcut in main loop

  - Automatically disable at end of frame

- Data collection takes about 2-3 minutes

  - Depends on workload

# Running it

- Hook up trace machinery to keyboard shortcut in main loop

  - Automatically disable at end of frame

- Data collection takes about 2-3 minutes

  - Depends on workload

- Stash results in binary file

  - About 100-150 MB of data for our use case

# Running it

- Hook up trace machinery to keyboard shortcut in main loop
  - Automatically disable at end of frame
- Data collection takes about 2-3 minutes
  - Depends on workload
- Stash results in binary file
  - About 100-150 MB of data for our use case
- Game resumes running at full framerate after collection!
  - Analyze dump offline

# Analysis

- Stats collected are associated with an instruction
  - Call stack is also captured, used to disambiguate
- Current stats captured by our setup:
  - L1 hit (separate tracking of I1/D1)
  - L2 hit
  - L2 miss (tracks instruction/data separately – better than HW can!)
  - # of explicit prefetches that hit D1 or L2
  - #Instructions executed

# Tooling – Flat profile of stats

# Tooling – Top-down tree

# Tooling – Reverse trees

# Tooling – Source annotation

# Tooling considerations

- ## Want to make "obviously bad" data pop

  - ### Badness Factor = L2miss^2 / #instructions

- ## Worth investing in a bunch of different views

  - ### Maximizes value by allowing more analysis on dumps

# PushBuffer::SetTextureAssets revisited

# PushBuffer::SetTextureAssets revisited

# CacheSim Pros

- Gathers data for every memory access in the program

- Non-intrusive

- Non-encumbered

- Works on Windows

  - Deeply instruments even graphics drivers, OS calls down to syscall level

- Open Source

  - Can easily extend to more scenarios

# CacheSim Cons

- Capture speed could be better

- Only works on Windows

  - Can still simulate a Jaguar cache for console workflows (ignore OS stuff)

- Not 100% hardware accurate (and can't be)

  - Treats the CPU as an in-order CPU – no OOO scheduling

  - Must use virtual addresses to index cache (minor issue)

  - Array prefetchers not simulated (overly pessimistic about arrays)

  - MESI/Store forwarding buffers/…

# Future

- Hardware prefetch simulation

- Non-temporal store simulation

- Speeding up captures

- Extensions

# Thanks + Q & A

- Special thanks:
  - Mike Acton, Jonathan Adamczewski & Elan Ruskin
  - Mark Cerny
- http://github.com/insomniacgames/ig-cachesim
  - Play with it and submit your own improvements
- Get in touch
  - afredriksson@insomniacgames.com
  - @deplinenoise