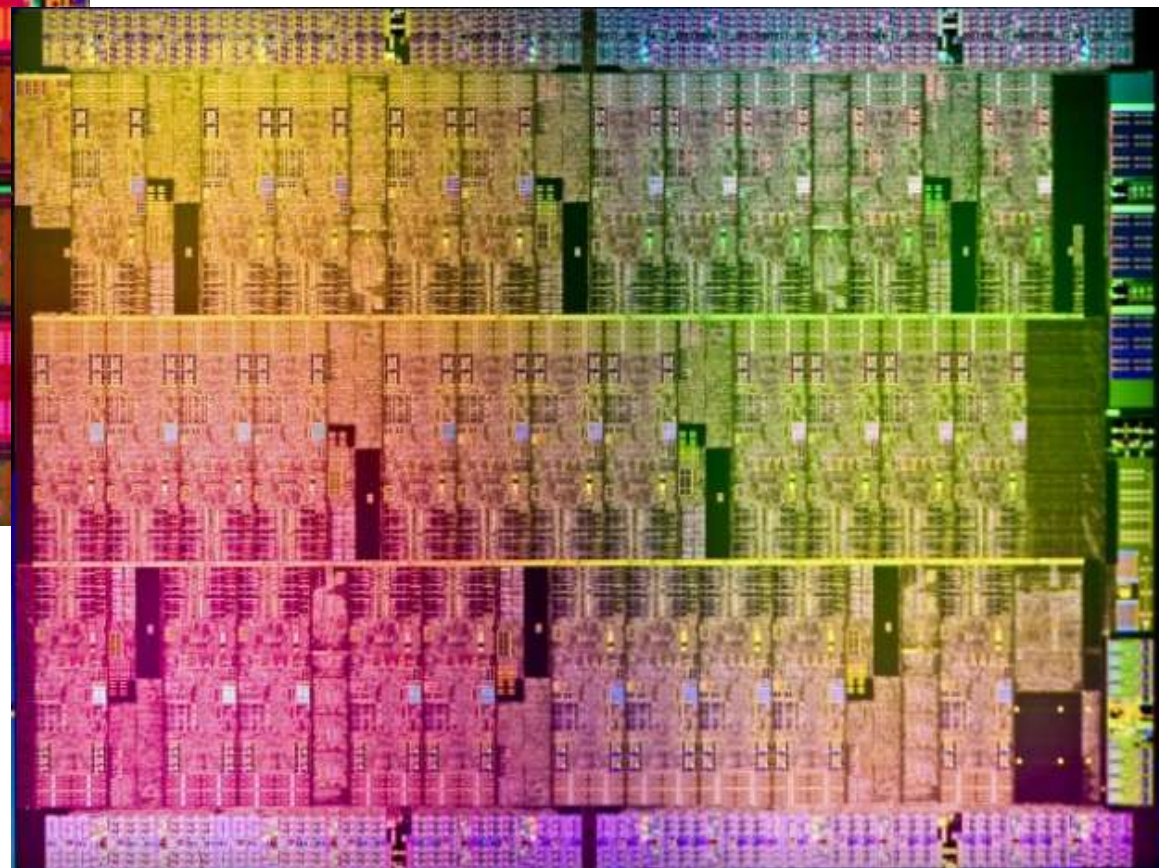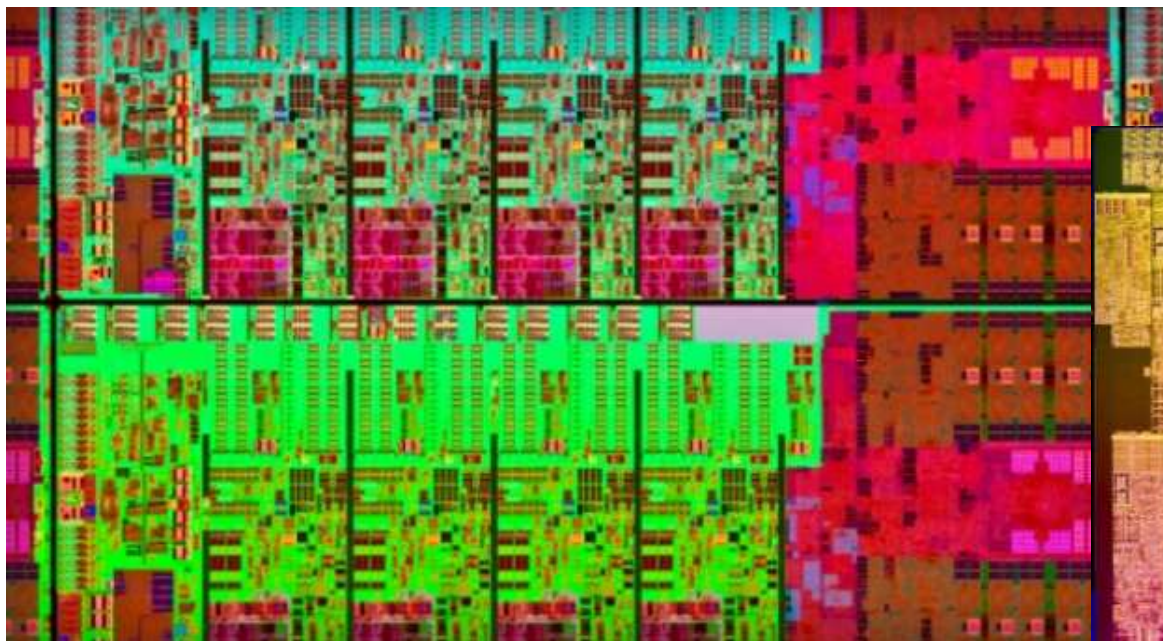# SIMD + NativeAOT + Dynamic PGO

Jiří Činčura (engineer, x̄ size) | Karel Zikmund (manager, tall)
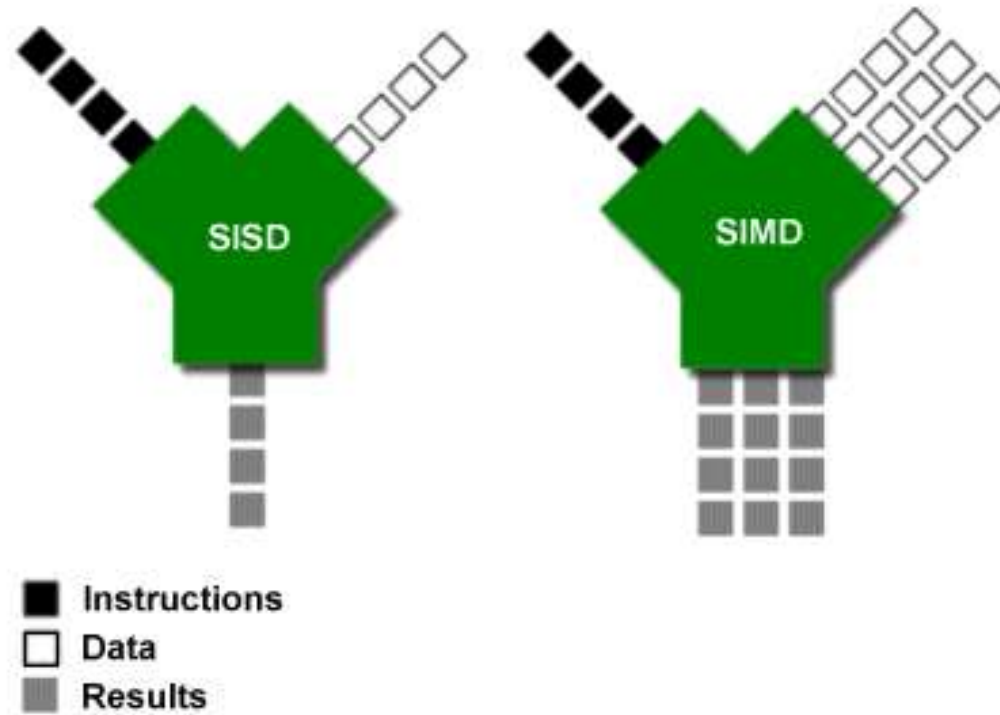
# A Modern Processor

# SISD

- Single Instruction Single Data

- Instruction Level Parallelism (ILP)
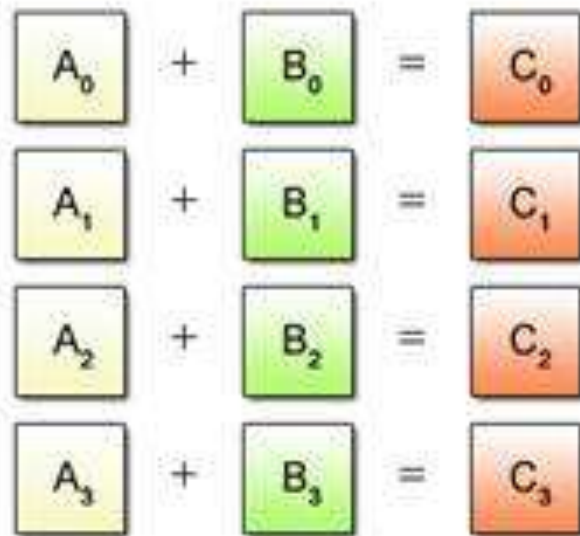
- Parallel processing

# SIMD

- Single Instruction Multiple Data

# SIMD

# SIMD



(a) Scalar Operation

$A_0 + B_0 = C_0$
$A_1 + B_1 = C_1$
$A_2 + B_2 = C_2$
$A_3 + B_3 = C_3$
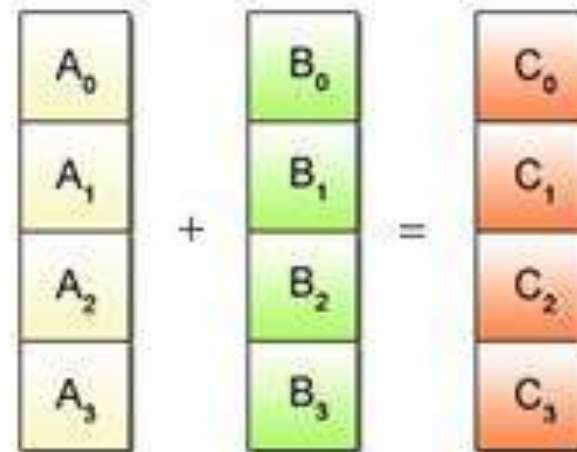
(b) SIMD Operation

# NativeAOT

- Ahead-of-Time compilation
- JIT, startup time, memory/working set
- Benefits
  - Faster startup time
  - Smaller memory footprint
  - Self-contained executable
  - Restricted environments (i.e. iOS)
- Limitations
  - No dynamic loading
  - No runtime code gen (System.Reflection.Emit, limited reflection)
  - Bigger application binary

# C# -> IL -> CPU

- C# code is compiled into IL (MSIL, CIL)
    - Stack based
    - Object oriented
- C# -> IL = Roslyn compiler (C# compiler)
- IL -> CPU = RyuJIT (JIT)
- Straightforward assembly code is not the fastest one
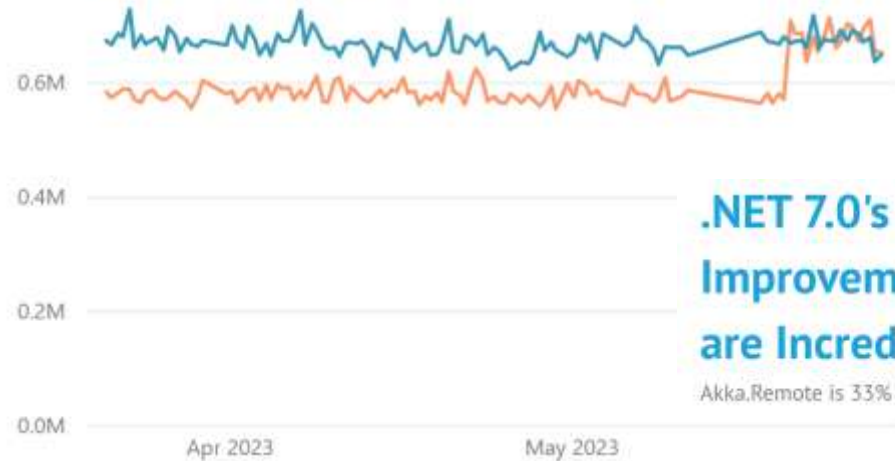- JIT must generate great code and do it fast
    - 🤷‍♂️

# PGO

- Profile Guided Optimization
- Static
  - Collect data from representative run and store along the executable
  - Is the data up-to-date?
- Dynamic
  - In-proc, no training or special builds
  - Uses Tiered Compilation by instrumenting code in initial tiers
  - Collected data used later for better or more optimization

# Benefits

- About 15% or up

# Key Optimizations

- Guarded Devirtualization (GDV)
  - At virtual and interface call sites, introduce tests for specific types
  - If the test succeeds, we know exactly which method will be called
    - Also try and inline the method
    - If the method is on a value class, inline the unboxing stub and the method
      - If source is a box attempt to optimize away the box too
  - If the test fails, just do the normal virtual / interface call
- .NET 8: extends GDV to handle some delegate invokes as well
- Opt-in: Multiple guesses GDV

# GDV

```
void RegisterUser(IUserService service, User user)
{
    service.Register(user); // virtual call
}
```

⬇

```
void RegisterUser(IUserService service, User user)
{
    CORINFO_HELP_CLASSPROFILE32(service.GetType());
    service.Register(user);
}
```

⬇

```
void RegisterUser(IUserService service, User user)
{
    if (service is UserServiceImpl impl)
        impl.Register(user); // direct call, can be inlined
    else
        service.Register();  // still virtual (fallback)
}
```

# GDV

```
void RegisterUser(IUserService service, User user)
{
    service.Register(user); // virtual call
}
```

⬇

```
void RegisterUser(IUserService service, User user)
{
    CORINFO_HELP_CLASSPROFILE32(user.GetType());
    service.Register(user);
}
```

⬇

```
void RegisterUser(IUserService service, User user)
{
    if (service is UserServiceImpl impl)
        impl.Register(user);
    else if (service is GenericUserService1<User> impl)
        impl.Register(user);
    else if (service is GenericUserService2<User> impl)
        impl.Register(user);
    else
        service.Register();
}
```

# Profile-Driven Inlining

- Use profile data to ensure key methods are inlined

- Relaxed thresholds for IL size and number of basic blocks

- Waste less energy on (semi-) cold call sites

# Profile-Driven Inlining

```csharp
bool IsPrimitiveType(Type type) =>
    type == typeof(bool)   ||
    type == typeof(char)   ||
    type == typeof(sbyte)  ||
    type == typeof(byte)   ||
    type == typeof(short)  ||
    type == typeof(ushort) ||
    type == typeof(int)    ||
    type == typeof(uint)   ||
    type == typeof(long)   ||
    type == typeof(ulong)  ||
    type == typeof(float)  ||
    type == typeof(double) ||
    type == typeof(nint)   ||
    type == typeof(nuint);
```

```
// Methods
.method public hidebysig static
    bool IsPrimitiveType (
        class [System.Runtime]System.Type 'type'
    ) cil managed
{

    .custom instance void System.Runtime.CompilerServices.NullableContextAttribu
        01 00 01 00 00
    )
    // Method begins at RVA 0x20a0
    // Code size 271 (0x10f)
    .maxstack 2

    IL_0000: ldarg.0
    IL_0001: ldtoken [System.Runtime]System.Boolean
    IL_0006: call class [System.Runtime]System.Type [System.Runtime]System.Type::
    IL_000b: call bool [System.Runtime]System.Type::op_Equality(class [System.Rur
    IL_0010: brtrue IL_010d

    IL_0015: ldarg.0
    IL_0016: ldtoken [System.Runtime]System.Char
    IL_001b: call class [System.Runtime]System.Type [System.Runtime]System.Type::
    IL_0020: call bool [System.Runtime]System.Type::op_Equality(class [System.Rur
    IL_0025: brtrue IL_010d

    IL_002a: ldarg.0
    IL_002b: ldtoken [System.Runtime]System.SByte
    IL_0030: call class [System.Runtime]System.Type [System.Runtime]System.Type::
    IL_0035: call bool [System.Runtime]System.Type::op_Equality(class [System.Rur
    IL_003a: brtrue IL_010d
```

# Profile-Driven Inlining

```csharp
for (int i = 0; i < 100; i++)
{
    Test<int, float>();
    Thread.Sleep(16);
}

[MethodImpl(MethodImplOptions.NoInlining)]
static bool Test<T1, T2>() =>
    IsPrimitiveType(typeof(T1)) &&
    IsPrimitiveType(typeof(T2));

static bool IsPrimitiveType(Type type) =>
    type == typeof(bool) ||
    type == typeof(char) ||
    type == typeof(sbyte) ||
    type == typeof(byte) ||
    type == typeof(short) ||
    type == typeof(ushort) ||
    type == typeof(int) ||
    type == typeof(uint) ||
    type == typeof(long) ||
    type == typeof(ulong) ||
    type == typeof(float) ||
    type == typeof(double) ||
    type == typeof(nint) ||
    type == typeof(nuint);
```

```asm
; Assembly listing for method Program:Test[int,float]():bool
; Tier-1 compilation
; No PGO data
        sub     rsp, 40
        mov     rcx, 0x11B802000B8 ; 'System.Int32'
        call    [Program:IsPrimitiveType(System.Type):bool]
        test    eax, eax
        je      SHORT G_M27198_IG05
        mov     rcx, 0x11B80205090 ; 'System.Single'
        call    [Program:IsPrimitiveType(System.Type):bool]
        nop
        add     rsp, 40
        ret
G_M27198_IG05:
        xor     eax, eax
        add     rsp, 40
        ret
; Total bytes of code 53

Inliner: too many IL bytes
```

# Profile-Driven Inlining

```csharp
for (int i = 0; i < 100; i++)
{
    Test<int, float>();
    Thread.Sleep(16);
}

[MethodImpl(MethodImplOptions.NoInlining)]
static bool Test<T1, T2>() =>
    IsPrimitiveType(typeof(T1)) &&
    IsPrimitiveType(typeof(T2));

static bool IsPrimitiveType(Type type) =>
    type == typeof(bool) ||
    type == typeof(char) ||
    type == typeof(sbyte) ||
    type == typeof(byte) ||
    type == typeof(short) ||
    type == typeof(ushort) ||
    type == typeof(int) ||
    type == typeof(uint) ||
    type == typeof(long) ||
    type == typeof(ulong) ||
    type == typeof(float) ||
    type == typeof(double) ||
    type == typeof(nint) ||
    type == typeof(nuint);
```

```asm
; Assembly listing for method Program:Test[int,float]():bool
; Tier-1 compilation
; Optimized with Dynamic PGO
        mov       eax, 1
        ret
; Total bytes of code 6

Inliner:
• Inline candidate has 13 foldable branches.
• Inline has 28 foldable intrinsics.
• Callsite has profile data: 1.0.
```

# Instrumentation Overhead

- [Dynamic PGO startup improvements in NET 8 · Issue #76969](#)
- Sparse , scalable edge profiles enabled for all methods
- GDV random state now in TLS
- Scalable profile mode
- More cases where we bypass instrumentation
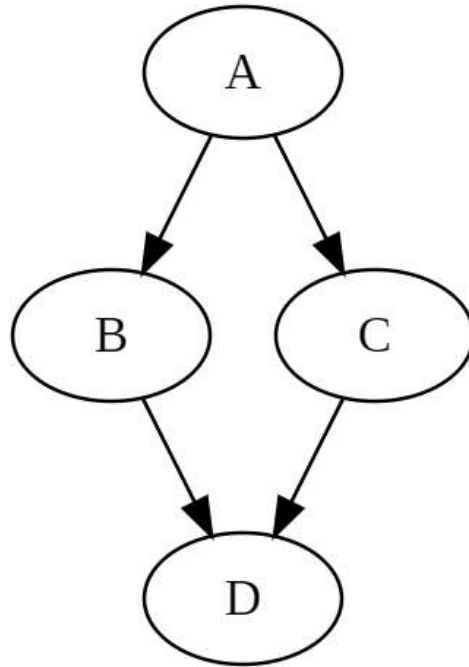- Enable intrinsic expansion in Tier0

# Class and Method Profiles

- *Reservoir sampling* used to create approximate histograms of target classes (for virtual/interface calls) and target methods (for indirect/delegate calls)
  - Fixed-sized table per site (currently 32 entries (was 8))
    - One global table per site
  - Each call adds entry to table, until table is full, then
    - Each call may randomly replace some table entry, with probability
    - This also keeps contention low
- When optimizing, this data is used to drive GDV, testing for the most likely outcome(s)

# Profiling Blocks

**Dense**

- Each block gets a counter

- Quite a bit of redundancy
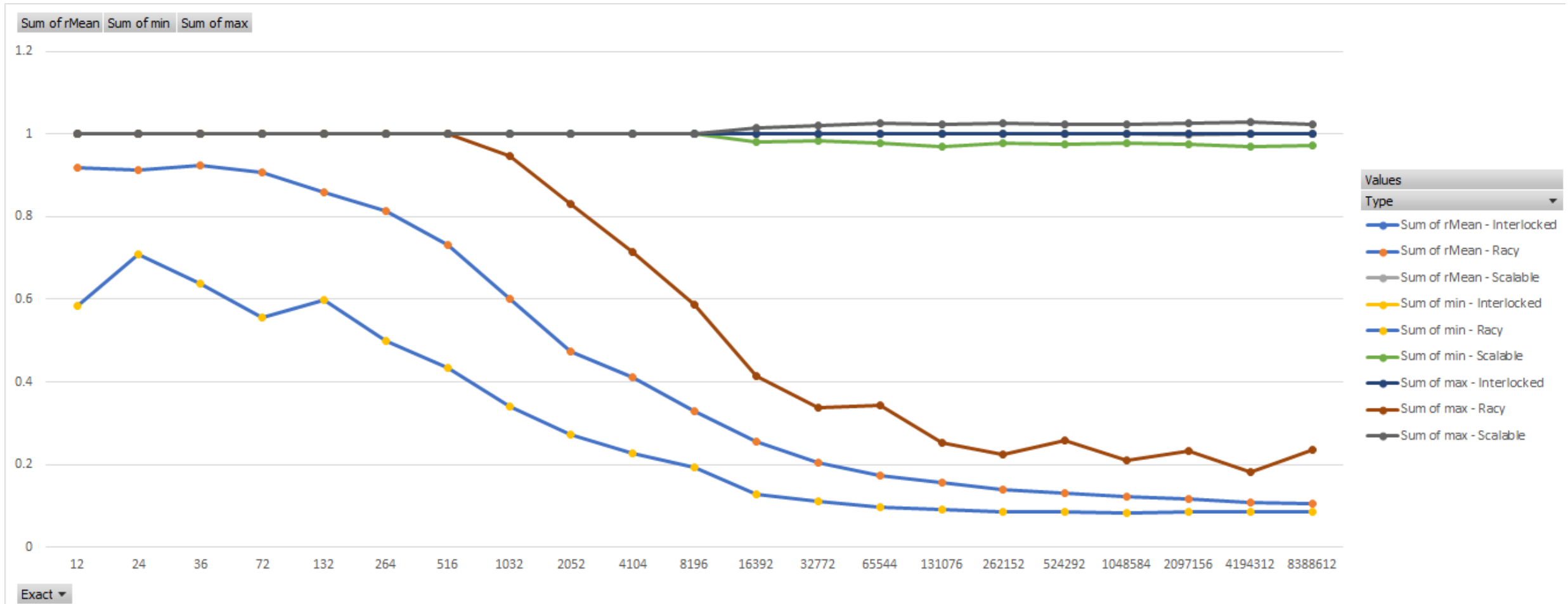  - Simple diamond: four blocks, two independent counts

**Sparse**

- Subset of edges get counters
  - Need to add in pseudo-edges

- Block counts reconstructed via "simple" math

# Scalable Counters

- Pre .NET 8, instrumentation was using a shared counter for the sparse edge counts, and not interlocking ("racy") updates.
  - When app is heavily multithreaded:
    - Heavy contention on some counters (very slow Tier0-instr code)
    - Poor accuracy as many updates are lost due to races
  - Interlocked adds fix the accuracy issue, but contention is even worse
  - Not feasible to shard the counters (i.e., TLS) both because of the space required and the need to aggregate across shards

# Scalable Counters



Blue and red lines show the ratio of a "racy" contended counter's value to the true value.
Note it can lose upwards of 90% of the counts (this was on a 12 core machine)

# Scalable Counters

- .NET 8 introduces scalable counters
  - Use interlocked add for first $2^N$ counts
  - Add randomly after that...
    - Add by 2 with probability ½ for count in $(2^N, 2^{N+1}]$
    - Then by 4 with probability ¼ for count in $(2^{N+1}, 2^{N+2}]$
- With suitable threshold (N=13) count value is very likely within 2% of true value
- Number of writes to "hot" (potentially contended) counters drops dramatically
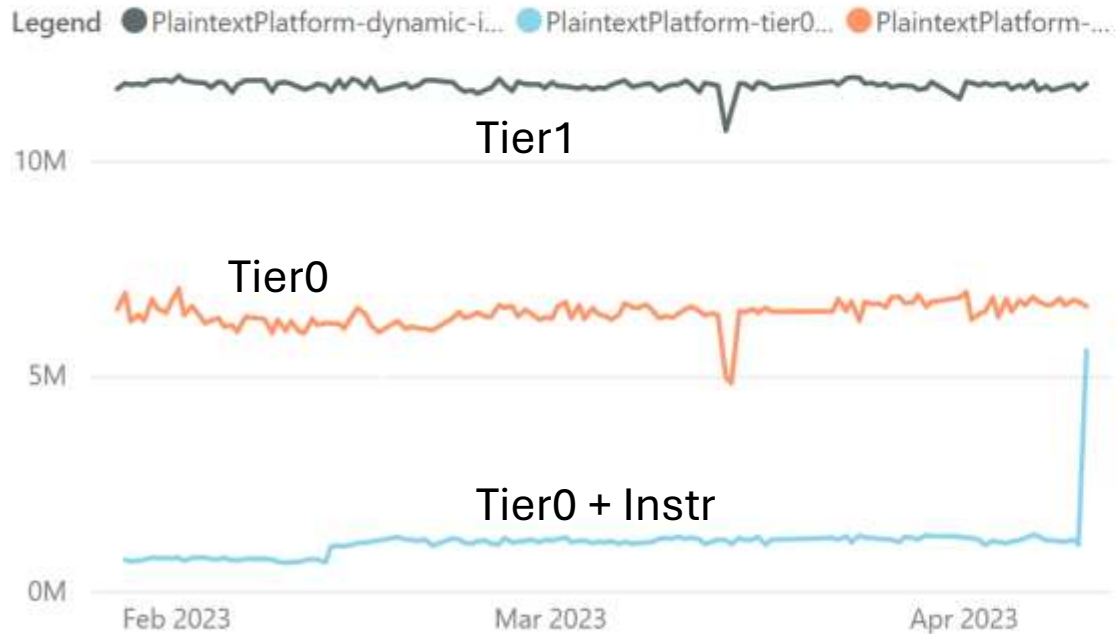
# Scalable Counters



Deviation of scalable counter from true value.

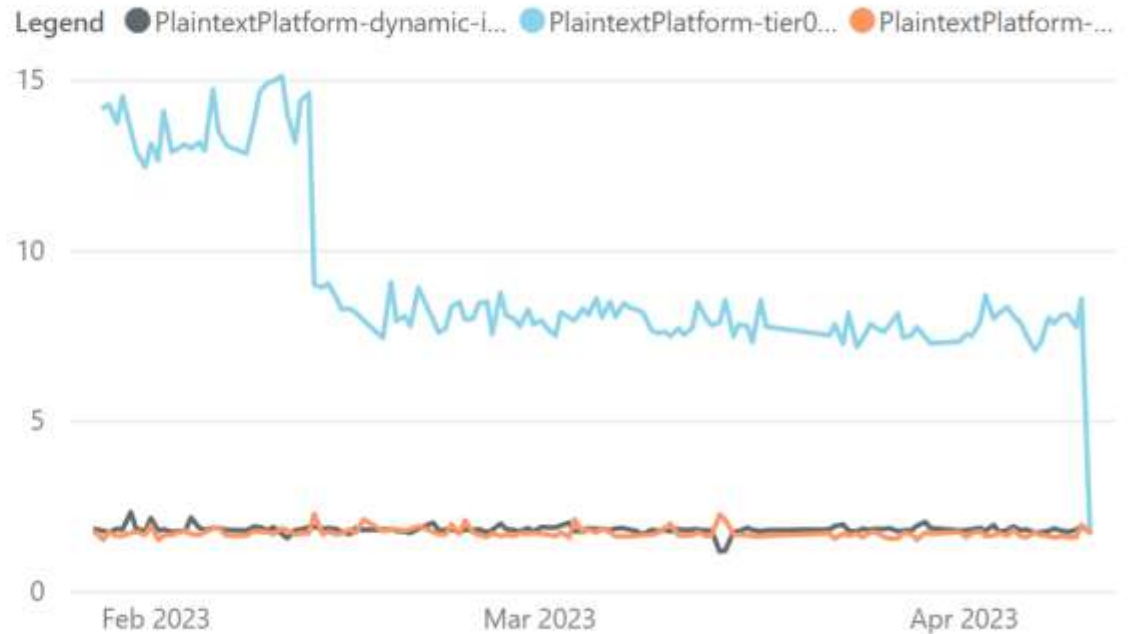Counts exactly up to 8192, then randomly for higher values

5-95 spread about +/- 2%

# Impact of improvements to instrumentation on Tech Empower RPS / Latency



Requests Per Second

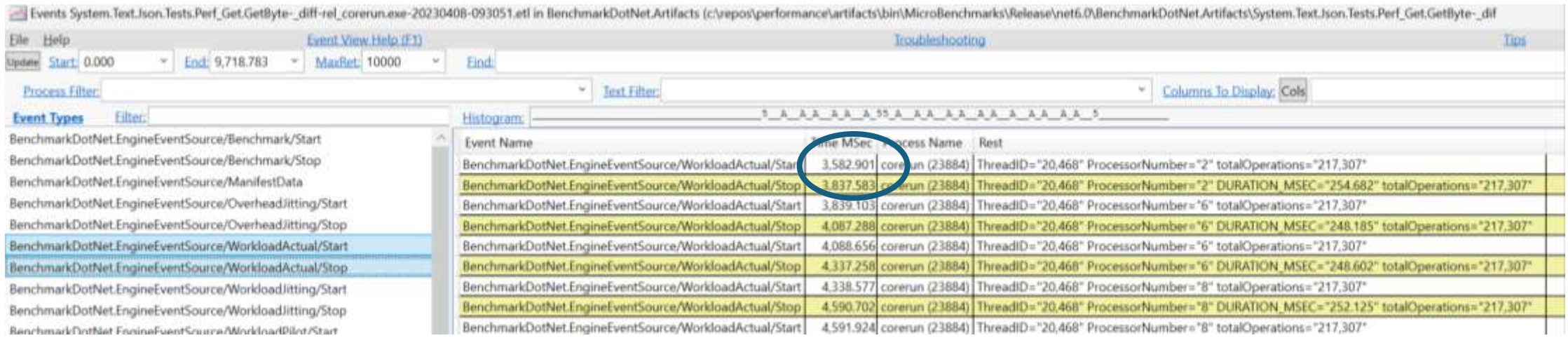Legend ● PlaintextPlatform-dynamic-i... ● PlaintextPlatform-tier0... ● PlaintextPlatform-...

Tier1

Tier0

Tier0 + Instr

Latency - Average (ms)

Legend ● PlaintextPlatform-dynamic-i... ● PlaintextPlatform-tier0... ● PlaintextPlatform-...

# Randomness

- Instrumentation relies quite a bit on *randomness*
  - GDV profiles use randomness for Reservoir Sampling to build approximate histograms
  - Count profiles use randomness to improve scalability
- PGO data will likely not be the same from one run to the next
  - But typically, there are enough observations that the overall behavior is still stable and repeatable
  - There is already a fair amount of non-determinism when running code, but now the jitted codegen depends on it in a fundamental way.
- If you suspect a bug, try running with `DOTNET_TieredPGO=0`

# PGO, BDN & PerfView



| Name | Exc % | Exc | Inc % | Inc | Fold | When | First |
|---|---|---|---|---|---|---|---|
| system.text.json![OptimizedTier1]System.Text.Json.Utf8JsonReader.TryGetByte(unsigned int8&) | 41.3 | 3,144 | 59.2 | 4,507.0 | 0 | _____266661666666666676666661_____ | 1,991 |
| microbenchmarks![OptimizedTier1]System.Text.Json.Tests.Perf_Get.GetByte() | 28.5 | 2,167 | 89.3 | 6,801.0 | 0 | _49999199999999999999999999_ | 1,991 |
| System.Private.CoreLib.il![OptimizedTier1]System.Buffers.Text.Utf8Parser.TryParseByteD(value class System.ReadOnlySpan`1< | 17.7 | 1,351 | 17.8 | 1,354.0 | 0 | _____02111012111112121112120_____ | 1,994 |
| coreclr!? | 2.1 | 159 | 99.4 | 7,564.0 | 0 | ____068A9AAA9AA9AAAAAA99999AA93_____ | 1,172 |
| ntoskrnl!? | 1.4 | 107 | 2.7 | 204.0 | 0 | _130oooooo_oo_oooo_o___oo_ooo_ | 1,092 |
| clrjit!? | 1.2 | 91 | 2.8 | 216.0 | 0 | _120oooo_0o_ooo00oo_____oooo_ | 1,330 |
| ci!? | 1.2 | 88 | 1.4 | 109.0 | 0 | _020_____ | 1,148 |
| system.text.json![QuickJitted]System.Text.Json.Utf8JsonReader.TryGetByte(unsigned int8&) | 0.9 | 70 | 2.1 | 161.0 | 0 | _____31_____ | 1,660 |
| system.text.json![OptimizedTier1]System.Text.Json.Utf8JsonReader..ctor(value class System.ReadOnlySpan`1<unsigned int8> | 0.8 | 59 | 0.9 | 72.0 | 0 | _oooooooooooo0oooo0oo0ooo_ | 2,136 |
| system.text.json![OptimizedTier1]System.Text.Json.Utf8JsonReader.TryGetNumber(value class System.ReadOnlySpan`1<unsig | 0.6 | 43 | 0.6 | 43.0 | 0 | _____ooooooooooooooo_o0o_oo0o_____ | 2,087 |
| ntdll!? | 0.5 | 41 | 99.8 | 7,600.0 | 0 | _168A9AAA9AA9AAAAAA99999AA93_ | 1,092 |
| System.Private.CoreLib.il![7]System.Buffers.Text.Utf8Parser.TryParseByteD(value class System.ReadOnlySpan`1<unsigned int8 | 0.5 | 40 | 0.5 | 40.0 | 0 | _____1_____ | 1,890 |

You can get ETL traces from BDN via -p ETW (on windows) and view these with PerfView.

Note the same method now appears as two entries. From the time chart we can see that the QuickJitted (Tier0) version ran early on, and then the Tier1 version took over.

Make sure to filter the profile to just the time that BDN is actually making measurements.

# PGO, BDN & PerfView



Open the events view, select BD's WorkloadActual events, verify the intervals show consistent times, pick one and set the time limits on your profile view. Here: 3582..3837