

# .NET Core: Performance Revolution

Adam Sitnik

# About myself

## Open Source:

- BenchmarkDotNet
- Awesome .NET Performance
- Core CLR (Span<T>)
- CoreFx Lab (Span<T>)
- & more

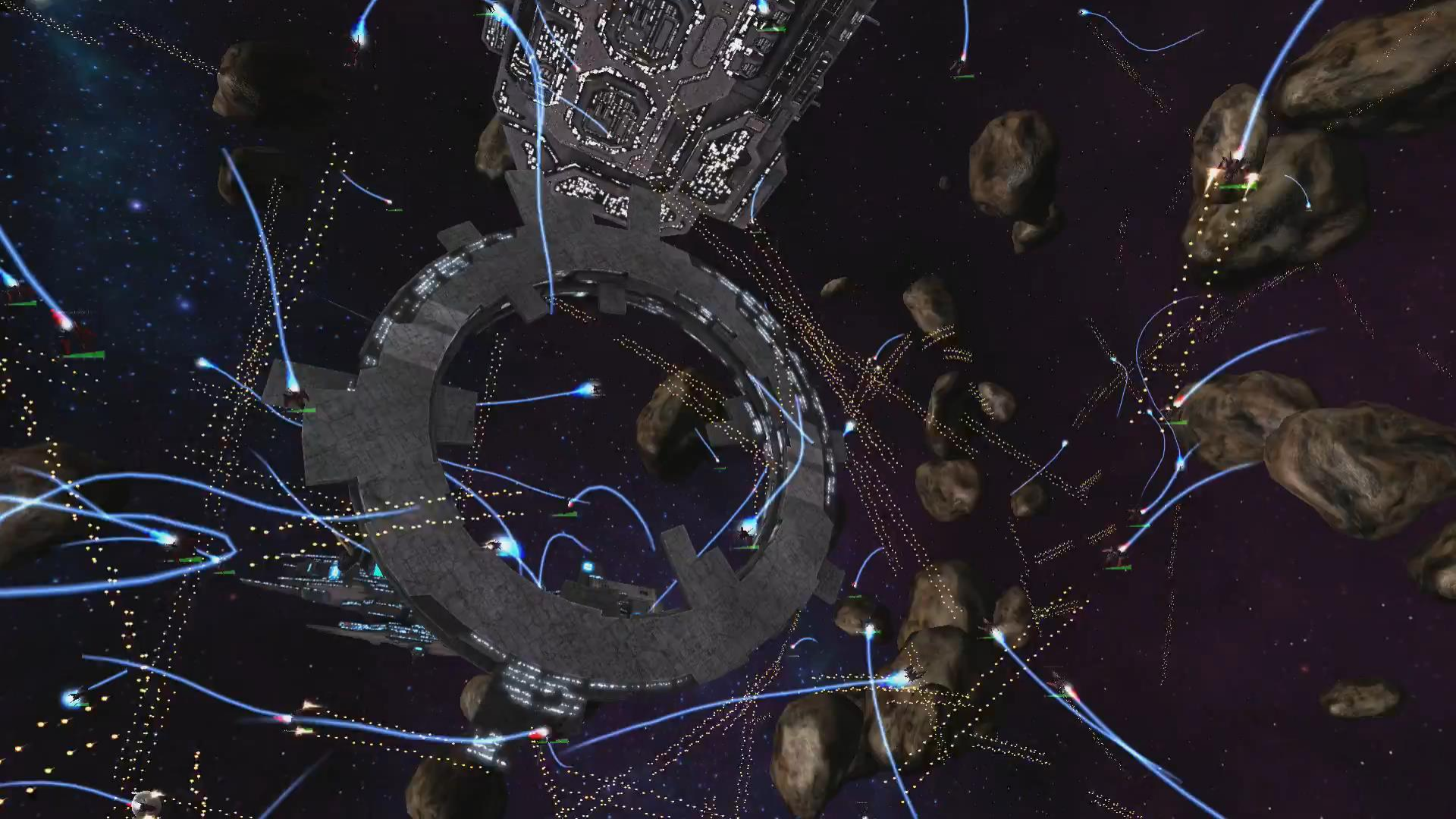
## Work:



- Energy Trading  
(.NET Core running in production since July 2016)

# ASP.NET Core road to high-performance

- [Performance Improvements in .NET Core](#)
- [Performance Improvements in RyuJIT](#)
- .NET Native & CoreRT (later today)
- New tools:
  - Span<T>
  - ArrayPool<T>
  - ValueTask<T>
  - Unsafe



# How to avoid GC?

- Reduce allocations
- Eliminate all managed allocations:
  - Pool the memory
  - Use unmanaged memory

# Async on hotpath

```
Task<T> SmallMethodExecutedVeryVeryOften()  
{  
    if(CanRunSynchronously()) // true most of the time  
    {  
        return Task.FromResult(ExecuteSynchronous());  
    }  
    return ExecuteAsync();  
}
```

# Sample ValueTask usage

```
[MethodImpl(MethodImplOptions.AggressiveInlining)]  
ValueTask<int> SampleUsage()  
    => IsFastSynchronousExecutionPossible()  
        ? new ValueTask<int>(  
            result: ExecuteSynchronous()) // INLINEABLE!!!  
        : new ValueTask<int>(  
            task: ExecuteAsync());
```

```
int ExecuteSynchronous() { }
```

```
Task<int> ExecuteAsync() { }
```

# How **not** to consume ValueTask

```
async ValueTask<int> ConsumeWrong(int repeats)
{
    int total = 0;
    while (repeats-- > 0)
        total += await SampleUsage();

    return total;
}
```



# Async Task Method Builder

```
[AsyncStateMachine(typeof(DemoInt.<ConsumeWrong>d__4))]
private Task ConsumeWrong(int repeats)
{
    DemoInt.<ConsumeWrong>d__4 <ConsumeWrong>d__;
    <ConsumeWrong>d__.<>4__this = this;
    <ConsumeWrong>d__.repeats = repeats;
    <ConsumeWrong>d__.<>t__builder = AsyncTaskMethodBuilder.Create();
    <ConsumeWrong>d__.<>1__state = -1;
    AsyncTaskMethodBuilder <>t__builder = <ConsumeWrong>d__.<>t__builder;
    <>t__builder.Start<DemoInt.<ConsumeWrong>d__4>(ref <ConsumeWrong>d__);
    return <ConsumeWrong>d__.<>t__builder.Task;
}
```

# How to consume ValueTask

```
async ValueTask<int> ConsumeProperly(int repeats)
{
    int total = 0;
    while (repeats-- > 0)
    {
        ValueTask<int> valueTask = SampleUsage(); // INLINEABLE

        total += valueTask.IsCompleted
            ? valueTask.Result // hot path
            : await valueTask.AsTask();
    }

    return total;
}
```

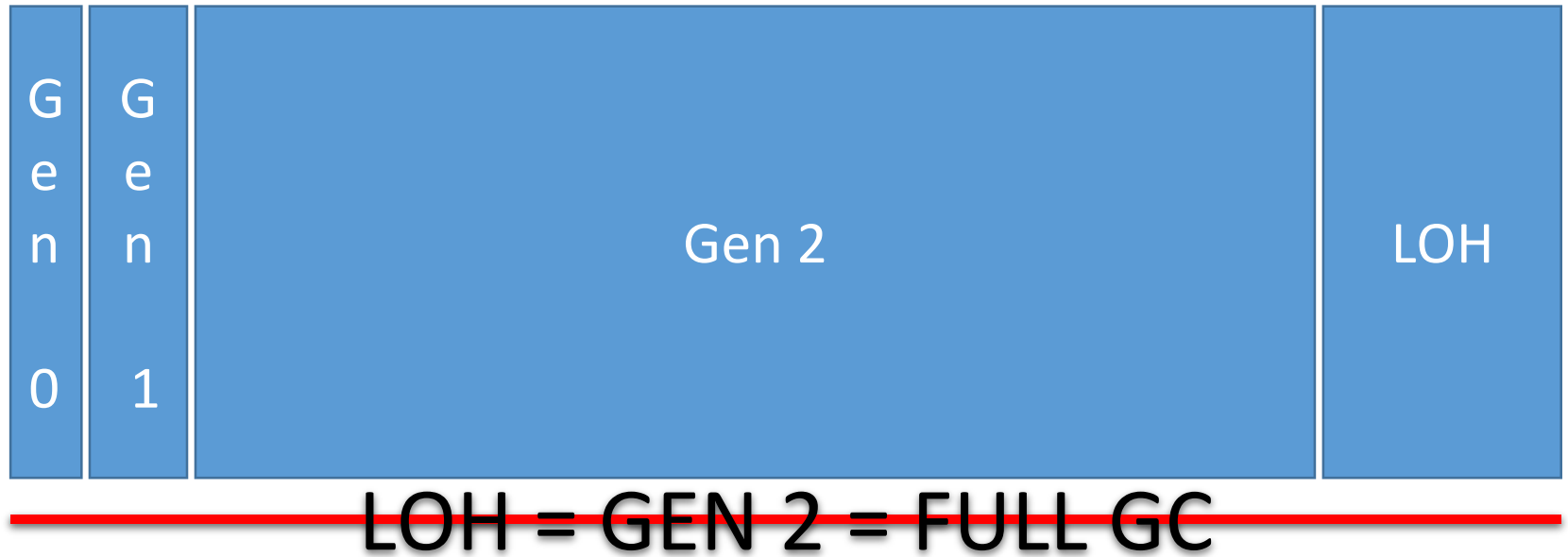
# ValueTask vs Task: Overhead Only

Method	Repeats	Mean	Scaled	Gen 0	Gen 1	Allocated
Task	100	720.9 ns	1.49	3.4674	0.0001	7272 B
ValueTask_Wrong	100	1,097.4 ns	2.27	-	-	0 B
ValueTask_Properly	100	482.9 ns	1.00	-	-	0 B

# Value Task: Summary

- It's not about replacing Task
- It has a **single purpose**: reduce heap allocations in async hot path where common synchronous execution is possible
- You can benefit from inlining, but not for free
- Use the `.IsCompleted` and `.Result` for getting best performance

# .NET Managed Heap\*



\* - simplified, Workstation mode or view per logical processor in Server mode

# ArrayPool

- **Pool of reusable managed arrays**
- The default maximum length of each array in the pool is  $2^{20}$  (1024\*1024 = 1 048 576)
- System.Buffers package

# ArrayPool: Sample

```
var samePool = ArrayPool<byte>.Shared;  
byte[] buffer = samePool.Rent(minLength);  
try  
{  
    Use(buffer);  
}  
finally  
{  
    samePool.Return(buffer);  
}
```

# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>



# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>
Allocate	1 000	41.122 ns	0.0812 ns	0.4880	0.0000	-	1024 B
RentAndReturn_Shared	1 000	<b>42.535 ns</b>	0.0621 ns	-	-	-	<b>0 B</b>

# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>
Allocate	1 000	41.122 ns	0.0812 ns	0.4880	0.0000	-	1024 B
RentAndReturn_Shared	1 000	<b>42.535 ns</b>	0.0621 ns	-	-	-	<b>0 B</b>
Allocate	10 000	<b>371.113 ns</b>	3.2994 ns	4.7847	0.0000	-	10024 B
RentAndReturn_Shared	10 000	<b>42.565 ns</b>	0.0450 ns	-	-	-	<b>0 B</b>

# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>
Allocate	1 000	41.122 ns	0.0812 ns	0.4880	0.0000	-	1024 B
RentAndReturn_Shared	1 000	<b>42.535 ns</b>	0.0621 ns	-	-	-	<b>0 B</b>
Allocate	10 000	371.113 ns	3.2994 ns	4.7847	0.0000	-	10024 B
RentAndReturn_Shared	10 000	<b>42.565 ns</b>	0.0450 ns	-	-	-	<b>0 B</b>
Allocate	100 000	3,625.029 ns	17.2533 ns	31.2497	31.2497	<b>31.2497</b>	100024 B
RentAndReturn_Shared	100 000	<b>42.426 ns</b>	0.0555 ns	-	-	-	<b>0 B</b>

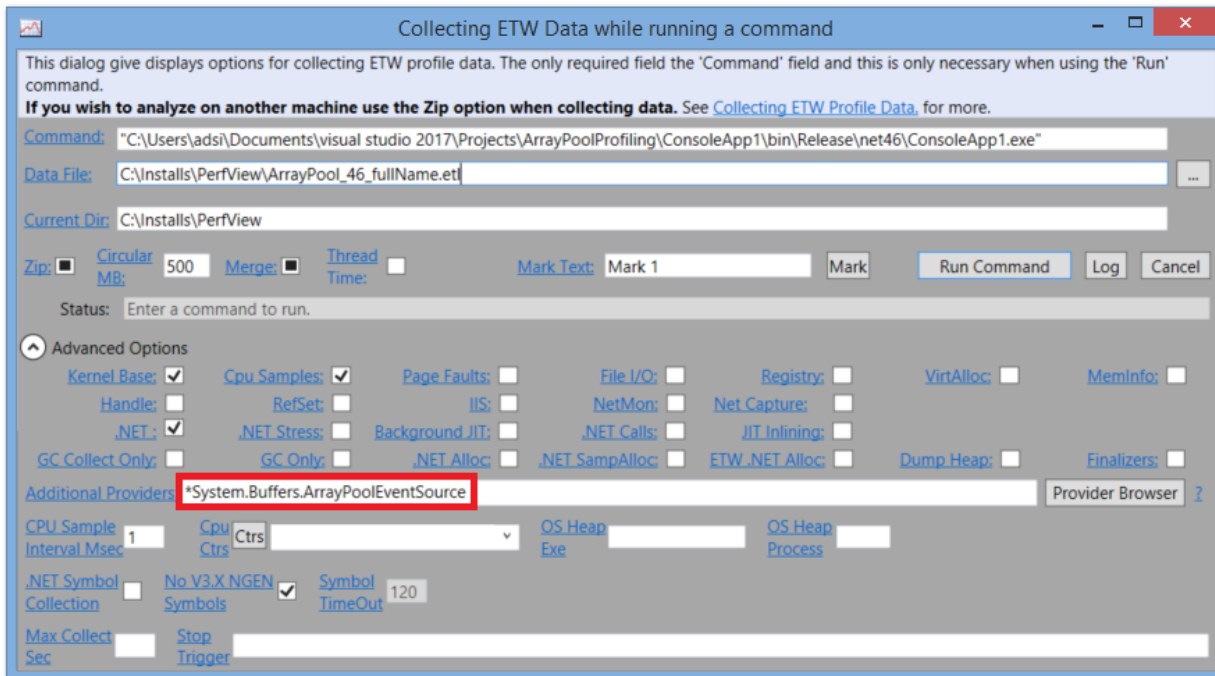
# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>
Allocate	1 000	41.122 ns	0.0812 ns	0.4880	0.0000	-	1024 B
RentAndReturn_Shared	1 000	<b>42.535 ns</b>	0.0621 ns	-	-	-	<b>0 B</b>
Allocate	10 000	371.113 ns	3.2994 ns	4.7847	0.0000	-	10024 B
RentAndReturn_Shared	10 000	<b>42.565 ns</b>	0.0450 ns	-	-	-	<b>0 B</b>
Allocate	100 000	3,625.029 ns	17.2533 ns	31.2497	31.2497	31.2497	100024 B
RentAndReturn_Shared	100 000	<b>42.426 ns</b>	0.0555 ns	-	-	-	<b>0 B</b>
Allocate	1 000 000	18,769.792 ns	60.4307 ns	249.9980	249.9980	249.9980	1000024 B
RentAndReturn_Shared	1 000 000	<b>41.979 ns</b>	0.0555 ns	-	-	-	<b>0 B</b>

# Allocate? Pool?

Method	Size	Mean	StdDev	Gen 0	Gen 1	Gen 2	Allocated
Allocate	100	8.149 ns	0.0169 ns	0.0610	-	-	128 B
RentAndReturn_Shared	100	<b>43.446 ns</b>	0.0908 ns	-	-	-	<b>0 B</b>
Allocate	1 000	41.122 ns	0.0812 ns	0.4880	0.0000	-	1024 B
RentAndReturn_Shared	1 000	<b>42.535 ns</b>	0.0621 ns	-	-	-	<b>0 B</b>
Allocate	10 000	371.113 ns	3.2994 ns	4.7847	0.0000	-	10024 B
RentAndReturn_Shared	10 000	<b>42.565 ns</b>	0.0450 ns	-	-	-	<b>0 B</b>
Allocate	100 000	3,625.029 ns	17.2533 ns	31.2497	31.2497	31.2497	100024 B
RentAndReturn_Shared	100 000	<b>42.426 ns</b>	0.0555 ns	-	-	-	<b>0 B</b>
Allocate	1 000 000	18,769.792 ns	60.4307 ns	249.9980	249.9980	249.9980	1000024 B
RentAndReturn_Shared	1 000 000	<b>41.979 ns</b>	0.0555 ns	-	-	-	<b>0 B</b>
Allocate	10 000 000	521,016.536 ns	55,326.9203 ns	211.2695	211.2695	211.2695	10000024 B
<b>RentAndReturn_Shared</b>	<b>10 000 000</b>	<b>639,916.968 ns</b>	<b>116,288.7309 ns</b>	<b>206.3623</b>	<b>206.3623</b>	<b>206.3623</b>	<b>10000024 B</b>
RentAndReturn_Aware	10 000 000	<b>47.200 ns</b>	0.0407 ns	-	-	-	<b>0 B</b>

# System.Buffers.ArrayPoolEventSource



# BufferAllocated event

Events ArrayPool\_46\_fullName.etl.zip in PerfView (C:\Installs\PerfView\ArrayPool\_46\_fullName.etl.zip)

Update: Start: 0,000 End: 4,337,767 MaxRet: 10000 Find:

Process Filter: Text Filter: Columns To Display: Cols

Event Types Filter: ArrayPool Histogram: A9

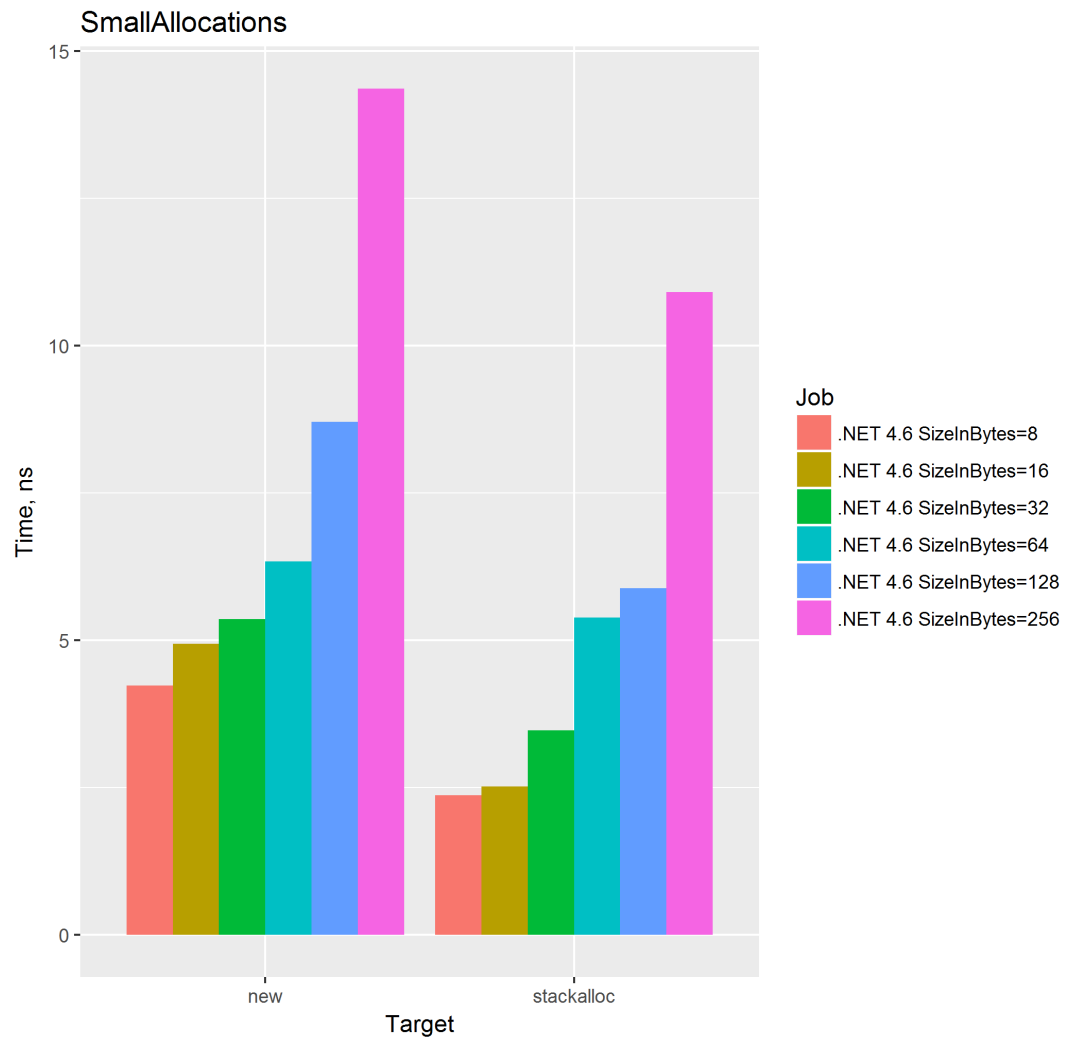
Event Name	Time MS	Process N	Rest
System.Buffers.ArrayPoolEventSource/BufferAllocated	333,566	ConsoleA	ThreadID="11.516" bufferId="15.368.010" bufferSize="524.288" poolId="21.083.178" bucketId="4.094.363" reason="Pooled"
System.Buffers.ArrayPoolEventSource/BufferRented	333,938	ConsoleA	ThreadID="11.516" bufferId="36.849.274" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferReturned	335,140	ConsoleA	ThreadID="11.516" bufferId="63.208.015" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/ManifestData	335,663	ConsoleA	ThreadID="11.516" bufferId="32.001.227" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	335,680	ConsoleA	ThreadID="11.516" bufferId="19.575.591" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	335,692	ConsoleA	ThreadID="11.516" bufferId="41.962.596" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	335,701	ConsoleA	ThreadID="11.516" bufferId="42.119.052" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	335,711	ConsoleA	ThreadID="11.516" bufferId="43.527.150" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	336,131	ConsoleA	ThreadID="11.516" bufferId="56.200.037" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	337,170	ConsoleA	ThreadID="11.516" bufferId="36.038.289" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,093	ConsoleA	ThreadID="11.516" bufferId="55.909.147" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,108	ConsoleA	ThreadID="11.516" bufferId="33.420.276" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,119	ConsoleA	ThreadID="11.516" bufferId="32.347.029" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,138	ConsoleA	ThreadID="11.516" bufferId="22.687.807" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,158	ConsoleA	ThreadID="11.516" bufferId="2.863.675" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,176	ConsoleA	ThreadID="11.516" bufferId="25.773.083" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,196	ConsoleA	ThreadID="11.516" bufferId="30.631.159" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"
System.Buffers.ArrayPoolEventSource/BufferAllocated	338,215	ConsoleA	ThreadID="11.516" bufferId="7.244.975" bufferSize="2.097.152" poolId="21.083.178" bucketId="-1" reason="OverMaximumSize"

# ArrayPool: Summary

- **LOH = Gen 2 = Full GC**
- ArrayPool was designed for best possible performance
- Pool the memory if you can control the lifetime
- Use **Pool.Shared** by default
- Pool allocates the memory for buffers > maxSize
- **The fewer pools, the smaller LOH, the better!**



Stackalloc is the fastest way to allocate small chunks of memory in .NET



	Allocation	Deallocation	Usage
Managed < 85 KB	Very fast	<ul style="list-style-type: none"> <li>• Non-deterministic</li> <li>• Blocking</li> <li>• Very slow</li> </ul>	<ul style="list-style-type: none"> <li>• Very easy</li> <li>• Common</li> <li>• Safe</li> </ul>
Managed: LOH	Fast		
Stackalloc	Super fast	<ul style="list-style-type: none"> <li>• Deterministic</li> <li>• Super fast</li> </ul>	<ul style="list-style-type: none"> <li>• Unsafe</li> <li>• Not common</li> <li>• Limited</li> </ul>
Marshal	Fast	<ul style="list-style-type: none"> <li>• Deterministic</li> <li>• Fast</li> </ul>	

# APIs before Span: parsing integer

```
int Parse(string input);
```

```
int Parse(string input, int startIndex, int length);
```

```
int Parse(string input, long startIndex, int length);
```

```
unsafe int Parse(char* input, int length);
```

```
unsafe int Parse(char* input, long startIndex, int length);
```

# Span<T>

It provides a uniform API for working with:

- Unmanaged memory buffers
- Arrays and subarrays
- Strings and substrings

It's fully **type-safe** and **memory-safe**.

Almost no overhead.

It's a stack only Value Type.

# Supports **any** memory

```
byte* pointerToStack = stackalloc byte[256];  
Span<byte> stackMemory = new Span<byte>(pointerToStack, 256);  
  
IntPtr unmanagedHandle = Marshal.AllocHGlobal(256);  
Span<byte> unmanaged = new Span<byte>(unmanagedHandle.ToPointer(), 256);  
  
char[] array = new char[] { 'i', 'm', 'p', 'l', 'i', 'c', 'i', 't' };  
Span<char> fromArray = array; // implicit cast  
  
ReadOnlySpan<char> fromString = "State of the .NET Performance".AsSpan();
```

# Simple API\*

```
public int Length { get; }  
public T this[int index] { get; set; }  
  
public Span<T> Slice(int start);  
public Span<T> Slice(int start, int length);  
  
public void Clear();  
public void Fill(T value);  
  
public void CopyTo(Span<T> destination);  
public bool TryCopyTo(Span<T> destination);  
  
public ref T DangerousGetPinnableReference();
```

\* It's not the full list

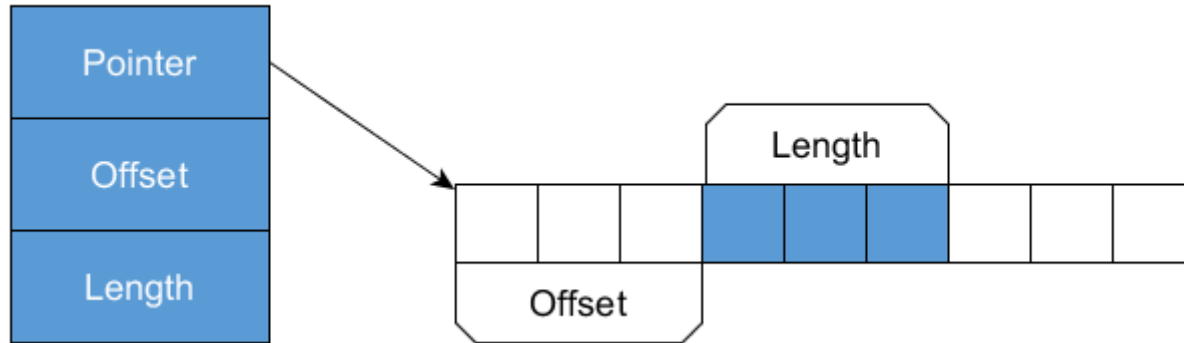
# API Simplicity!

```
int Parse(Span<char> input)
```

```
void Copy<T>(Span<T> source, Span<T> destination)
```

# Span for existing runtimes

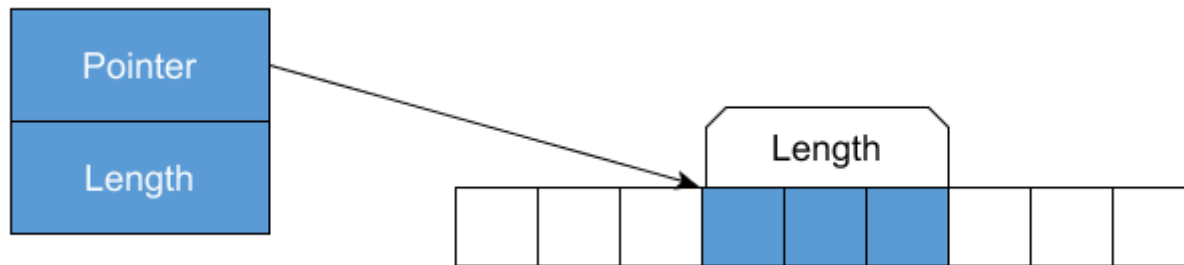
.NET Standard 1.0 (.NET 4.5+)





# Span for new runtimes

.NET Core 2.0 and any other runtime supporting by-ref fields



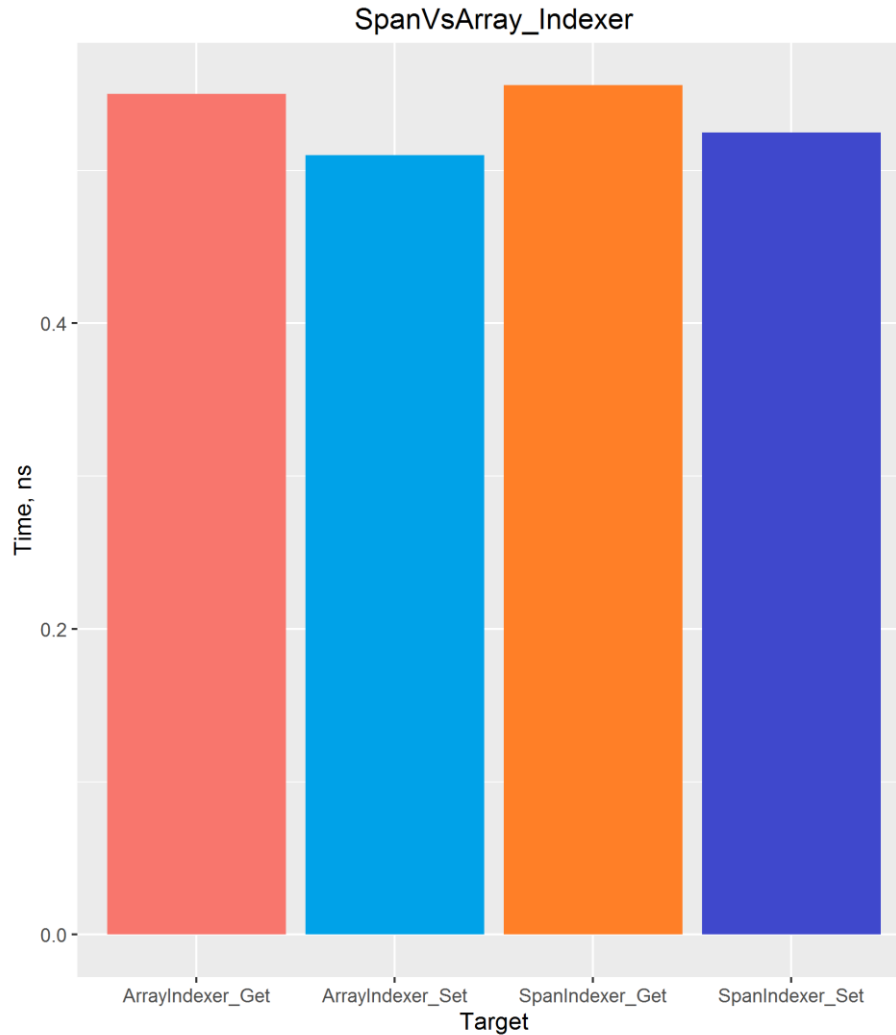
# “Fast” vs “Slow” Span

Method	Job	Mean	Scaled
SpanIndexer_Get	.NET 4.6	0.6119 ns	1.14
SpanIndexer_Get	.NET Core 1.1	0.6092 ns	1.13
<b>SpanIndexer_Get</b>	<b>.NET Core 2.0</b>	<b>0.5368 ns</b>	<b>1.00</b>
SpanIndexer_Set	.NET 4.6	0.6117 ns	1.13
SpanIndexer_Set	.NET Core 1.1	0.6082 ns	1.12
<b>SpanIndexer_Set</b>	<b>.NET Core 2.0</b>	<b>0.5417 ns</b>	<b>1.00</b>

There is some place for further improvement!

Span  
is on  
par  
with  
Array\*!

\*.NET Core 2.0



## Creating substrings **before Span** (pseudocode)

```
string Substring(string text, int startIndex, int length)
{
    string result = new string(length); // ALLOCATION!

    Memory.Copy(text, result, startIndex, length); // COPYING

    return result;
}
```

## Creating substrings **without allocation!** (pseudocode)

```
ReadOnlySpan<char> Slice(string text, int startIndex, int length)
    => new ReadOnlySpan<char>(
        ref text[0] + (startIndex * sizeof(char)),
        length);
```

# Substring vs Slice

Method	Chars	Mean	StdDev	Scaled	Gen 0	Allocated
Substring	10	<b>8.277 ns</b>	0.1938 ns	<b>4.54</b>	<b>0.0191</b>	<b>40 B</b>
Slice	10	<b>1.822 ns</b>	<b>0.0383 ns</b>	1.00	-	<b>0 B</b>
Substring	1000	<b>85.518 ns</b>	1.3474 ns	<b>47.22</b>	<b>0.4919</b>	<b>1032 B</b>
Slice	1000	<b>1.811 ns</b>	<b>0.0205 ns</b>	1.00	-	<b>0 B</b>

# Possible usages

- Parsing without allocations
- Formatting
- Base64/Unicode encoding
- HTTP Parsing/Writing
- Compression/Decompression
- XML/JSON parsing/writing
- Binary reading/writing
- & more!!

# Stack Only

- Instances can reside only on the stack
- Which is accessed by one thread at the same time

## **Advantages:**

- Few pointers for GC to track
- Safe Concurrency (no Struct Tearing)
- Safe lifetime. Method ends = memory can be returned to the pool or released

**The Limitations** <http://adamsitnik.com/Span/#the-limitations>



# Span: Summary

- Allows to work with **any** type of memory.
- It makes working with native memory much easier.
- Simple abstraction over Pointer Arithmetic.
- **Avoid allocation and copying of memory with Slicing.**
- Supports .NET Standard 1.0+
- It's performance is on par with Array for new runtimes.
- It's limited due to stack only requirements.

# System.Runtime.CompilerServices.Unsafe

## Overcoming C# limitations:

- Managed Pointer Arithmetic
- Casting w/o constraints
- Copy/Init Block
- Read/Write w/o constraints
- SizeOf(T)

```
ref T AddByteOffset<T>(ref T source, IntPtr byteOffset)
ref T Add<T>(ref T source, int elementOffset)
ref T Add<T>(ref T source, IntPtr elementOffset)
bool AreSame<T>(ref T left, ref T right)
void* AsPointer<T>(ref T value)
ref T AsRef<T>(void* source)
T As<T>(object o) where T : class
ref TTo As<TFrom, TTo>(ref TFrom source)
IntPtr ByteOffset<T>(ref T origin, ref T target)
void CopyBlock(ref byte destination, ref byte source, uint byteCount)
void CopyBlock(void* destination, void* source, uint byteCount)
void CopyBlockUnaligned(ref byte destination, ref byte source, uint byteCount)
void CopyBlockUnaligned(void* destination, void* source, uint byteCount)
void Copy<T>(void* destination, ref T source)
void Copy<T>(ref T destination, void* source)
void InitBlock(ref byte startAddress, byte value, uint byteCount)
void InitBlock(void* startAddress, byte value, uint byteCount)
void InitBlockUnaligned(ref byte startAddress, byte value, uint byteCount)
void InitBlockUnaligned(void* startAddress, byte value, uint byteCount)
T Read<T>(void* source)
T ReadUnaligned<T>(void* source)
T ReadUnaligned<T>(ref byte source)
int SizeOf<T>()
ref T SubtractByteOffset<T>(ref T source, IntPtr byteOffset)
ref T Subtract<T>(ref T source, int elementOffset)
ref T Subtract<T>(ref T source, IntPtr elementOffset)
void Write<T>(void* destination, T value)
void WriteUnaligned<T>(void* destination, T value)
void WriteUnaligned<T>(ref byte destination, T value)
```

# .NET Standard

Package name	.NET Standard	.NET Framework
System.Memory	1.0	4.5
System.Buffers	1.1	4.5.1
System.Threading.Tasks.Extensions	1.0	4.5
System.Runtime.CompilerServices.Unsafe	1.0	4.5

# Summary

- Use ValueTask only if it can help you!
- Pool the memory with ArrayPool
- Use Span and slicing to avoid allocations
- Use Span to take advantage of the native memory
- Use the “Unsafe” api to use C# only

# Sources

- [Span<T> design document](#)
- [Compile time enforcement of safety for ref-like types](#)
- [ValueTask doesn't inline well- GitHub issue](#)

# Děkuji!

Slides: <http://adamsitnik.com/files/Prague.pdf>

Code: <https://github.com/adamsitnik/StateOfTheDotNetPerformance>

@SitnikAdam

[Adam.Sitnik@gmail.com](mailto:Adam.Sitnik@gmail.com)