# Altimesh Hybridizer™
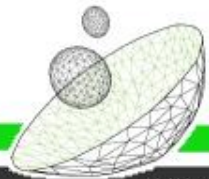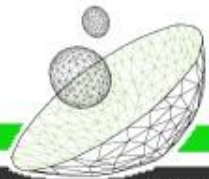
Enabling Accelerators in .Net and more
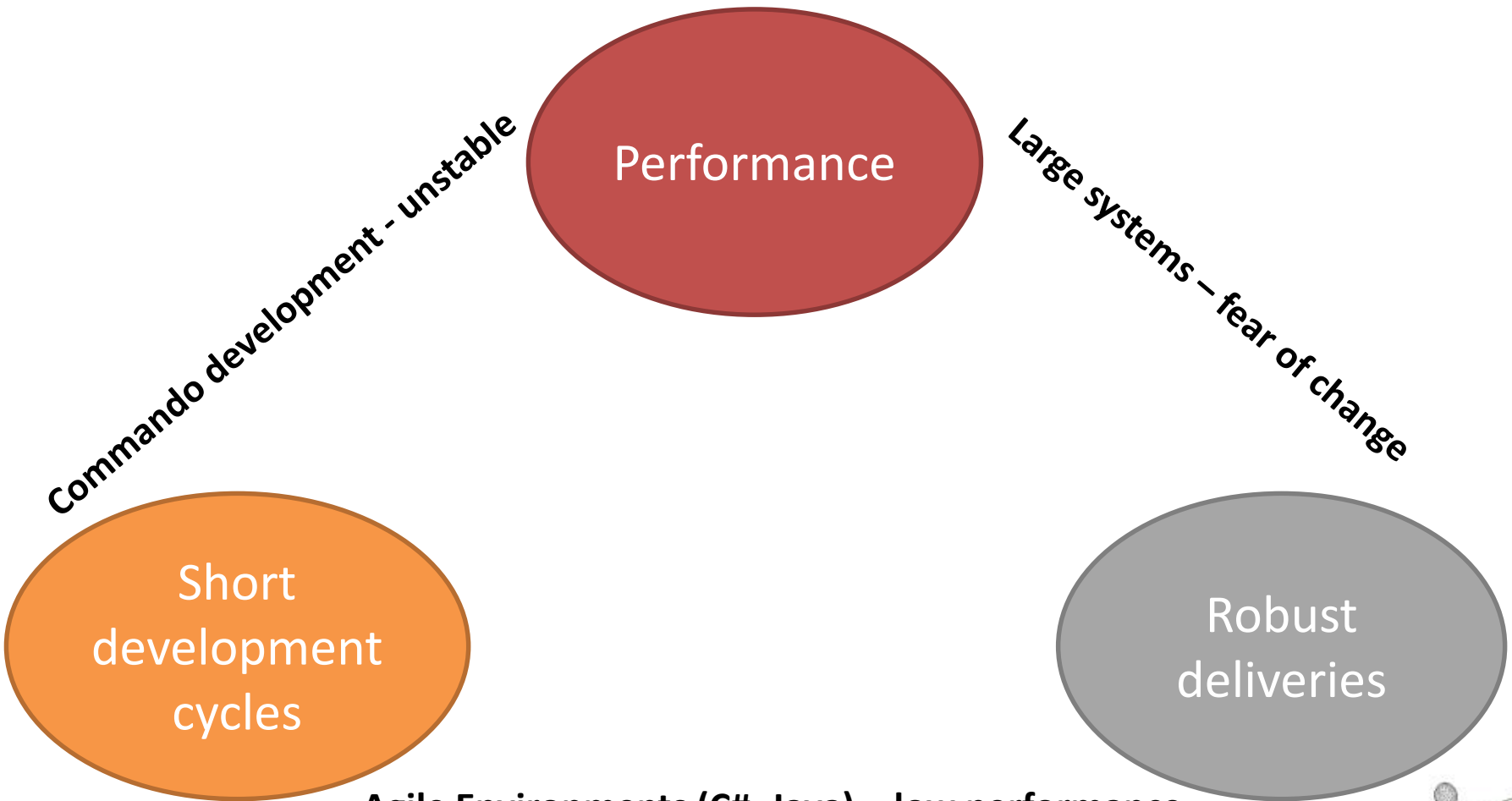
So many platforms, so few experts…

# WHY THE HYBRIDIZER ?

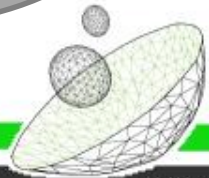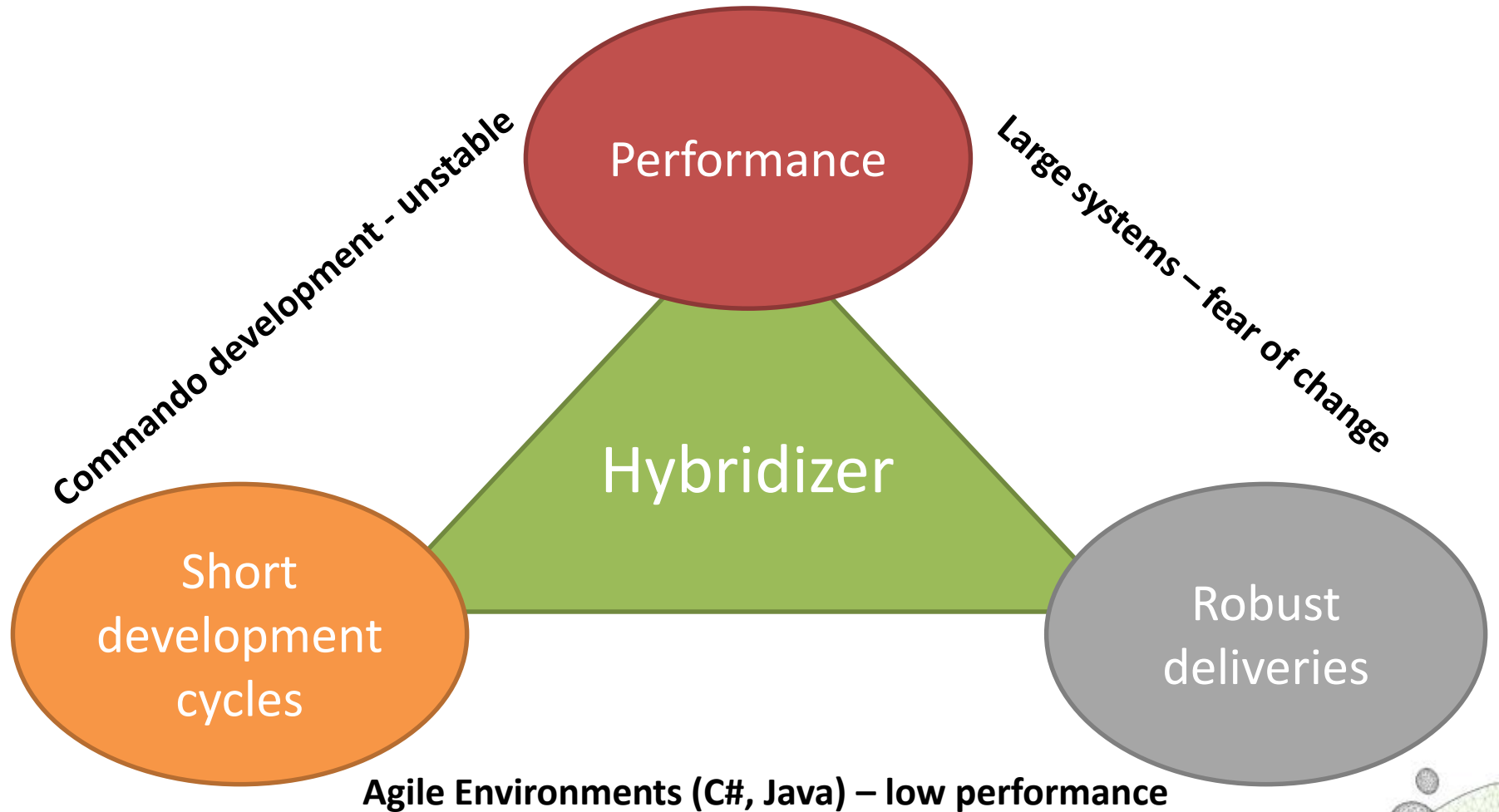# Software development teams accommodate external constraints

**Commando development - unstable**

**Large systems – fear of change**

Performance

Short development cycles

Robust deliveries

**Agile Environments (C#, Java) – low performance**

# Software development teams accommodate external constraints

**Performance**

Commando development - unstable

Large systems – fear of change

Hybridizer

**Short development cycles**

**Robust deliveries**

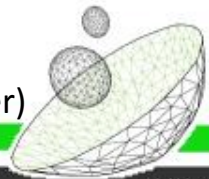**Agile Environments (C#, Java) – low performance**

# Why the Hybridizer?

- Develop in a managed environment (C#/Java)
  - Fast developments (fast compile time, edit and continue…)
  - Testing and refactoring ecosystem
  - Glitch-safe memory management
  - Embrace Change

  HIGHER PRODUCTIVITY    REDUCE TCO    OF APPLICATION DEVELOPMENT

- Benefit from manycore architectures
  - *With single version of the source code*
  - Obtain first grade performances (use >80% of peak)
  - Fine tune optimizations with debugger/profiler integration
  - Variety of execution platforms
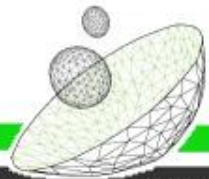  - Change execution target without rewriting code

  MORE EFFICIENT HARDWARE    REDUCE TCO    OF DATA CENTERS

IT spending : approx 30% in hardware and approx 20% in application development (Source : Gartner)
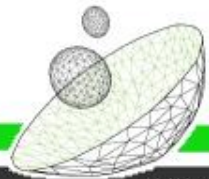
# What the Hybridizer is not

- Hybridizer is not a magic wand: some hints have to be given
  - Memory management is performed either in a naïve way, or needs to be done by hand
  - Memory level usages need to be defined
  - Some execution behaviors cannot be guessed

- Work distribution needs to be explicit
  - Loop parallelization is not automatic
  - Concurrency needs to be handled by hand
  - Code patterns need to be changed from sequential to parallel

# What the Hybridizer does

- Generates **source code** from **binaries**
  - Input is dot net binary (C#, VB.Net, Managed C++, other MSIL languages, Java)
  - Output is source code that can be used in various environments (plain C/C++ projects, CUDA projects, Windows/Linux, DotNet / Java runtimes)

- Supports the following language constructs
  - Virtual functions, generic types
  - Use of external libraries with seamless integration (e.g. CUBLAS, CURAND for CUDA environment) – user-extensible
  - Perform debugging within original source code – say C#. (this feature needs pdb)
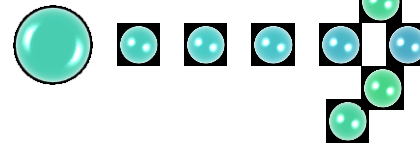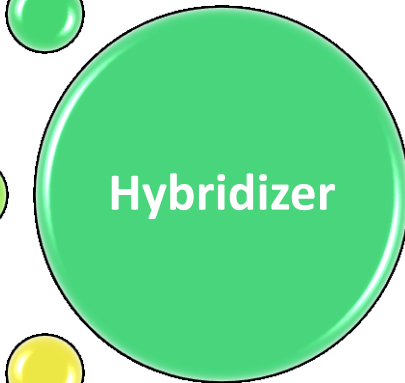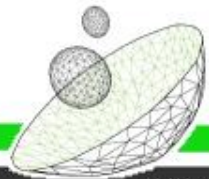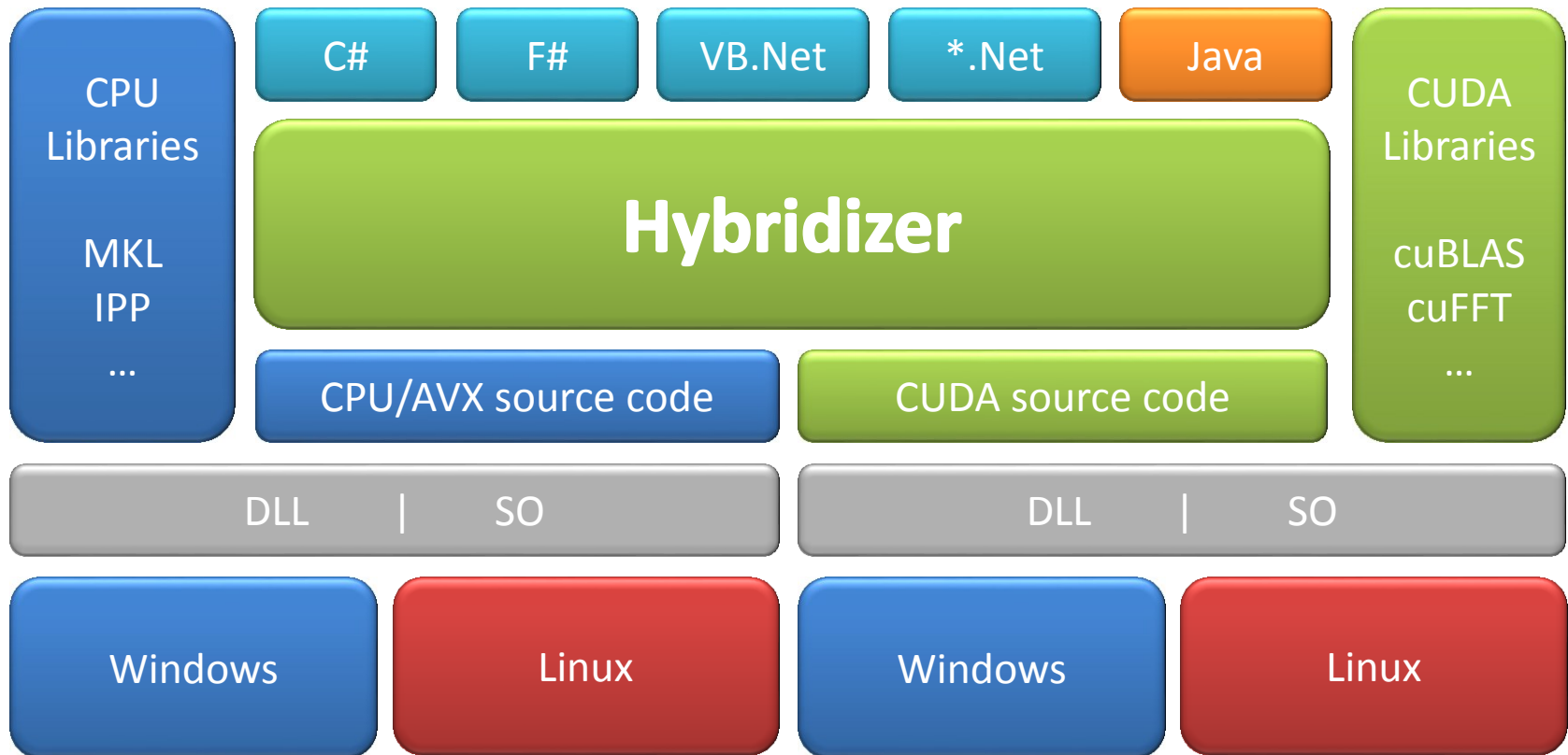
# What the Hybridizer does

.Net binaries (from C#,
F#, VB.Net, binary MSIL)

External libraries
(e.g. cuBLAS,
cuRAND, cuFFT)

**Hybridizer**

Customizations /
Optimizations

AVX

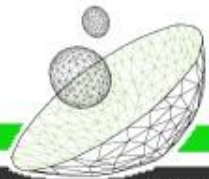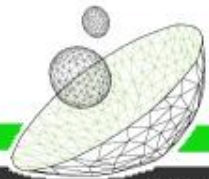Source code

Debug information

# Software Stack

Flexibility of managed environments,
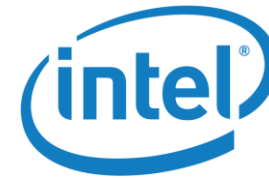80%+ usage of hardware

# HYBRIDIZER IN ACTION

# Basic features

- CUDA-style work distribution
- Seamless integration (attribute-based)
- Extensibility:
  - Usage of existing functions (erfc, hand-written, …)
  - Usage of external libraries (cuBLAS, cuRand, …)
  - Printf available using Console.Out / System.out
  - System.Math maps to <cmath> functions
- Customizable memory management
  - Zero copy arrays
  - Resident array (single copy for multiple kernel calls)

# Performances bandwidth & double precision

KEPLER – K20C

i7-3610 QM - AVX

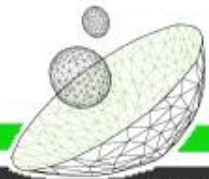| Compute | GCFLOPS | usage | GFLOPS | usage |
|---|---|---|---|---|
| whetstone | 541 | **92%** | 43.2 | **87%** |
| peak | 587 | - | 49.6 | - |

NOTE : Whetstone is our internal naive reproduction of the basic Whetstone test operating on doubles

| Memory | GB/s | usage | GB/s | usage |
|---|---|---|---|---|
| stream | 162 | 78% | 20.4 | 80% |
| peak | 208 | - | 25.6 | - |

NOTE : FMA IS COUNTED AS 1 FLOP HENCE REDUCING PEAK TO HALF
1 CFLOP = 1$^e$9 FMA DP – MEASURES ON K20C

1GB/s = 1$^e$9bytes /s here – MEASURES ON K20C – ECC OFF – CUDA 5.0

CORE i7-3610 QM (HT activated) @ 2.3 GHz
TurboBoost @ 3.1 GHz (observed using monitor)
AVX - OpenMP with 8 threads (4 cores)

## Virtual Functions

Support for Virtual functions

Function overriding : using inheritance

Use of Interfaces (single or multiple interfaces on classes or structs)

Native integration: no dedicated code needed.

```csharp
public interface ISimple
{
    int f();
}
public class Answer : ISimple
{

    [Kernel]
    public int f()
    {
        return 42 ;
    }
}
public class Other : ISimple
{

    [Kernel]
    public int f()
    {
        return 12;
    }
}
```
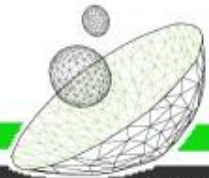
# Performances virtual functions

## KEPLER – K20C

| Expm1² | GFLOPS | GCFLOPS | usage |
|--------|--------|---------|-------|
| Local | 975 | 538 | **92%** |
| Dispatch | 478 | 263 | 45% |
| **peak** | 1174 | 587 | - |

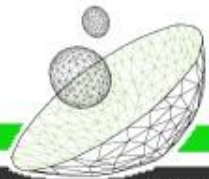Virtual functions suffer significant performance penalty

NOTE : FMA IS COUNTED AS 1 FLOP HENCE REDUCING PEAK TO HALF
1 GCFLOP = $1^e9$ FMA DP – MEASURES ON K20C

²: EXPM1 IS A TAYLOR EXPANSION OF EXP(X)-1: (1 ADDITION, 13 FUSED MULTIPLY ADD, 2 MULTIPLY)

# Improve performances with Generics

Generics can be converted to Templates

Generic constraints lead to usage of template functions (no virtual call)

Performances are very close to performance obtained with local functions (no inheritance/interface)

```csharp
[HybridTemplateConcept]
public interface IMyArray {
    double this[int index] { get; set; }
}


[HybridRegisterTemplate(Specialize=typeof(MyAlgorithm<MyArray>))]
public struct MyArray : IMyArray
{

    double[] _data;
    [Kernel] public double this[int index] {
        get { return _data[index]; }
        set { _data[index] = value; }
    }
}


public class MyAlgorithm<T> where T : struct, IMyArray
{
    T a, b;
    [Kernel] public void Add(int n) {
        for (int k = threadIdx.x + blockDim.x * blockIdx.x;
            k < n; k += blockDim.x * gridDim.x)
            a[k] += b[k];
    }
}
```

# Performances generics

## KEPLER – K20C

| Expm1² | GFLOPS | GCFLOPS | usage |
|--------|--------|---------|-------|
| **Local** | 975 | 538 | **92%** |
| **Dispatch** | 478 | 263 | 45% |
| **Generics** | 985 | 544 | **93%** |
| **peak** | 1174 | 587 | - |

Mapping generics to templates restores performances

NOTE : FMA IS COUNTED AS 1 FLOP HENCE REDUCING PEAK TO HALF
1 GCFLOP = 1$^{e}$9 FMA DP – MEASURES ON K20C

²: EXPM1 IS A TAYLOR EXPANSION OF EXP(X)-1: (1 ADDITION, 13 FUSED MULTIPLY ADD, 2 MULTIPLY)

# Performances single precision

## KEPLER – GTX 680
1536 cores @ 1006 GHz = 1545 GCFLOPS

## MAXWELL – GTX 750Ti
640 cores @ 1.085 GHz = 694.4 GCFLOPS

| Expm1² benchmark | GCFLOPS | Usage | GCFLOPS | usage |
|---|---|---|---|---|
| **Local** | 953.6 - 1234 | 61% - **80%** | 450.8 - 660.0 | 65% - **95%** |
| **Dispatch** | 392.3 - 632.7 | **25%** - 41% | 171.0 - 343.2 | **25%** - 49% |
| **Template** | 958.1 - 1069 | 62% - **69%** | 440.3 - 539.3 | 63% - **78%** |
| **peak** | 1545 | - | 694.4 | - |

without - *with* vectorization

without - *with* vectorization

NOTE : FMA IS COUNTED AS 1 FLOP HENCE REDUCING PEAK TO HALF : 1 GCFLOP = $1^e9$ FMA SP

²: EXPM1 IS A TAYLOR EXPANSION OF EXP(X)-1: (1 ADDITION, 13 FUSED MULTIPLY ADD, 2 MULTIPLY)

Developers perspective

# INTEGRATION WITH VISUAL STUDIO

# Debugging session using NSIGHT for Visual Studio [2010]



Breakpoint is set and hit in C# code

Values can be explored using Watch

Execution is on GPU

# Profiling session using
# NSIGHT for Visual Studio [2010]

Compilation with line-info allows dot net source-level profiling (also in release mode)

# Profiling session using
# VTune Amplifier for Visual Studio [2010]



See line association between original sequential C# code and vectorized x86/AVX assembly instructions

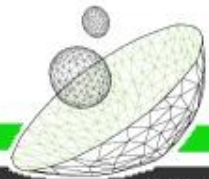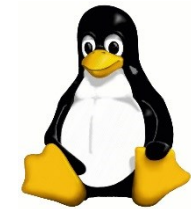# Usages – runtimes – execution environments

**Generated From**
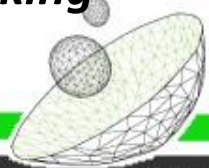
**Target**

**Use From**

**Run on**

"We have been using the Hybridizer for more than a year now with very satisfactory results. With **no prior knowledge of GPU programming,** we have been able to achieve **significant speedups in a large scale application** with unexcessive effort. Hybridizer enabled **rapid integration of GPU** within our development environment, with limited impact on a team of hundred programmers. It took **nine months to a handful of developers to go from early testing to production** on our first perimeter, and six more months to cover some of our most compute intensive calculations."

*Régis FRICKER - GPU project leader at Société Générale Investment Banking*

[Florent.Duguet@altimesh.com](mailto:Florent.Duguet@altimesh.com)

[Guillaume.de-Roujoux@altimesh.com](mailto:Guillaume.de-Roujoux@altimesh.com)

**THANK YOU**