# QuantAlea

Dynamic Cuda with F#

GTC 2013

March 21

San Jose, California

Dr. Daniel Egloff
Xiang Zhang
+41 44 520 01 17
+41 79 430 03 61

# Today

- Build awareness for F# and CUDA with F#

- Briefly present what we have achieved so far

- Show the potential of F# as a future core language for CUDA

- Do some live coding examples including CUDA scripting in Excel

- Point you to more resources and free licenses

# About Us

- Software development and consulting company
- Based in Zurich
- Core competence
    - Quantitative finance and risk management
    - Derivative pricing and modeling
    - Numerical computing
    - High performance computing (clusters, grid, GPUs)
    - Software engineering (C++, F#, Scala, ...)
- Early adopter of GPUs
    - First project with GPUs in finance back in 2007

# Why should you care about F#?

**QuantAlea**

- F# is now a first class CUDA language

- F# is a strongly typed functional first language
- Design goals of the F# language
  - High productivity
  - Development of robust and correct code
  - Efficiency
- F# is highly flexible and extensible
  - Monads, type providers
  - DSLs
- F# is ideal for numerical computing and therefore for CUDA
- F# has potential to play an important role in computational finance and big data
- Vast .NET ecosystem

# Alea.cuBase

QuantAlea

Alea.cuBase extends F# to a first class CUDA language

- Based on LLVM and CUDA 5 technology
- Noninvasive single language solution for host and GPU programming
- No additional language additions required, in particular no <<<...>>>
- Extensible
- Basis for creating higher level GPU aware DSLs

```
cuda {
    texture
    kernel
    ...
    launch logic
}
```

| Dynamic code generation | GPU algorithm scripting | Industry grade performance |
|---|---|---|
| Rapid development | Solid framework for reusability | Advanced CUDA programming |

## Dynamic code generation

- Generate GPU code programmatically at run-time
- Use .NET generics and F# code quotation splicing for flexible kernels
- Foundation to develop GPU aware domain specific languages
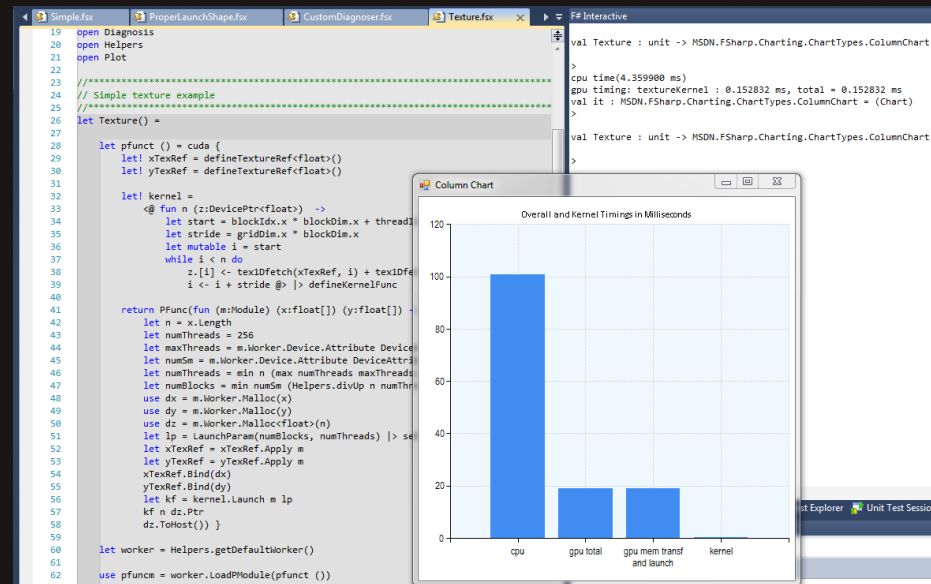
```fsharp
let init = <@ fun () -> 0.0 @>
```

```fsharp
let op = <@ (+) @>
```

```fsharp
let transf = <@ fun x -> x*x @>
```

```fsharp
let inline reduceKernel (initExpr:Expr<unit -> 'T>)
                        (opExpr:Expr<'T -> 'T -> 'T>)
                        (transfExpr:Expr<'T -> 'T>)
                        blockSize nIsPow2 =
<@ fun level n (inp:DevicePtr<'T>) (outp:DevicePt
    let init = %initExpr
    let op = %opExpr
    let transf = %transfExpr
```

QuantAlea

## Rapid development

- Easy and quick setup of development environment, no need to install NVIDIA nvcc compiler tools
- Rapid prototyping in F# interactive
- Iteratively improve CUDA kernel algorithms without time consuming build cycles
- Simple deployment

# Benefits

**QuantAlea**

## GPU algorithm scripting

- Execute F# scripts with GPU algorithms on command line or in F# interactive
- GPU scripting in Excel
- Integrate Alea.cuBase directly with Python

**Solid framework for reusability**

- Framework for type-safe definition of GPU resources
- CUDA monad to specify GPU resources together with launch logic in unified manner
- Reuse GPU kernel code and compose them to modular GPU kernel libraries

```
cuda {
  kernel_A
  launch logic
}
```

```
cuda {
  kernel_B
  launch logic
}
```

**Advanced CUDA programming**

- Support for texture, constant and shared memory
- Pointer operations to partition array data
- Special pointer types such as volatile pointers
- Runtime compilation control e.g. fast math
- Multiple streams
- Thread safe use of multiple GPUs
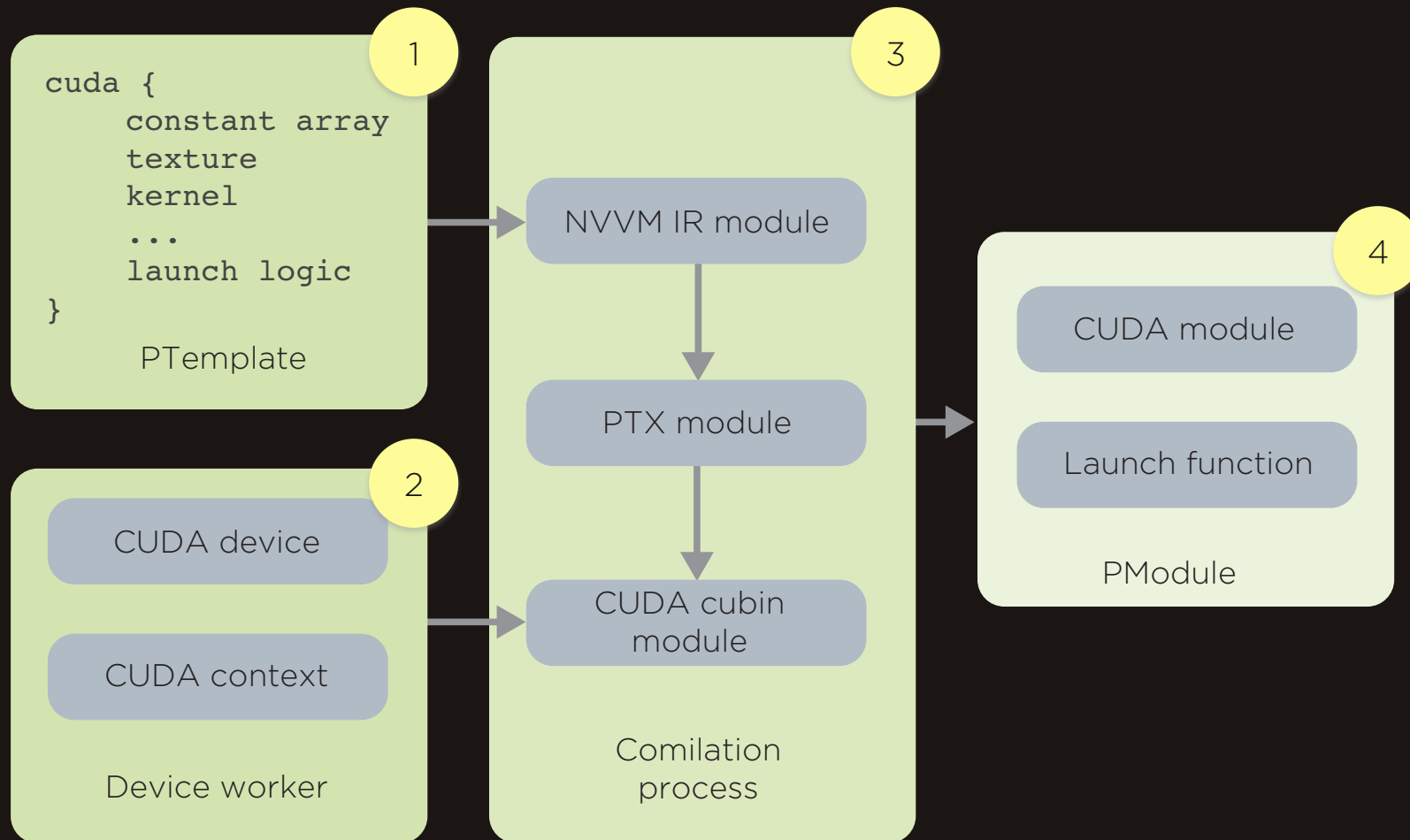- Inline PTX assembly instructions

# How does

# Alea.cuBase

# work ?

# Development Process

QuantAlea

## Four steps to a CUDA kernel with F# and Alea.cuBase

**1**

```
cuda {
    constant array
    texture
    kernel
    ...
    launch logic
}
```

PTemplate

**2**

CUDA device

CUDA context

Device worker

**3**

NVVM IR module

↓

PTX module

↓

CUDA cubin module

Comilation process

**4**

CUDA module

Launch function

PModule

How easily can I use

Alea.cuBase

?

# Live Coding

- Basic kernel programming

- Excel GPU scripting with Alea.cuBase and Alea.cuExtension, in Tsunami IDE and FCell
    - Excel based Monte Carlo simulation
    - PDE solver for 2d heat equation in GPU

# More Resources and Free Licenses     QuantAlea

- **More resources**
  - https://www.quantalea.net/products/resources/
  - https://github.com/quantalea

- **How to set up**
  - Fermi device or higher
  - Windows with .NET 4 and F# 2.0
  - CUDA 5 driver
  - Install Alea.cuBase
  - No need for CUDA toolkit or NVCC compiler

- **Apply for free licenses**
  - https://www.quantalea.net/news/22/

QuantAlea

Thank you