# QuantLib on Intel Xeon Phi
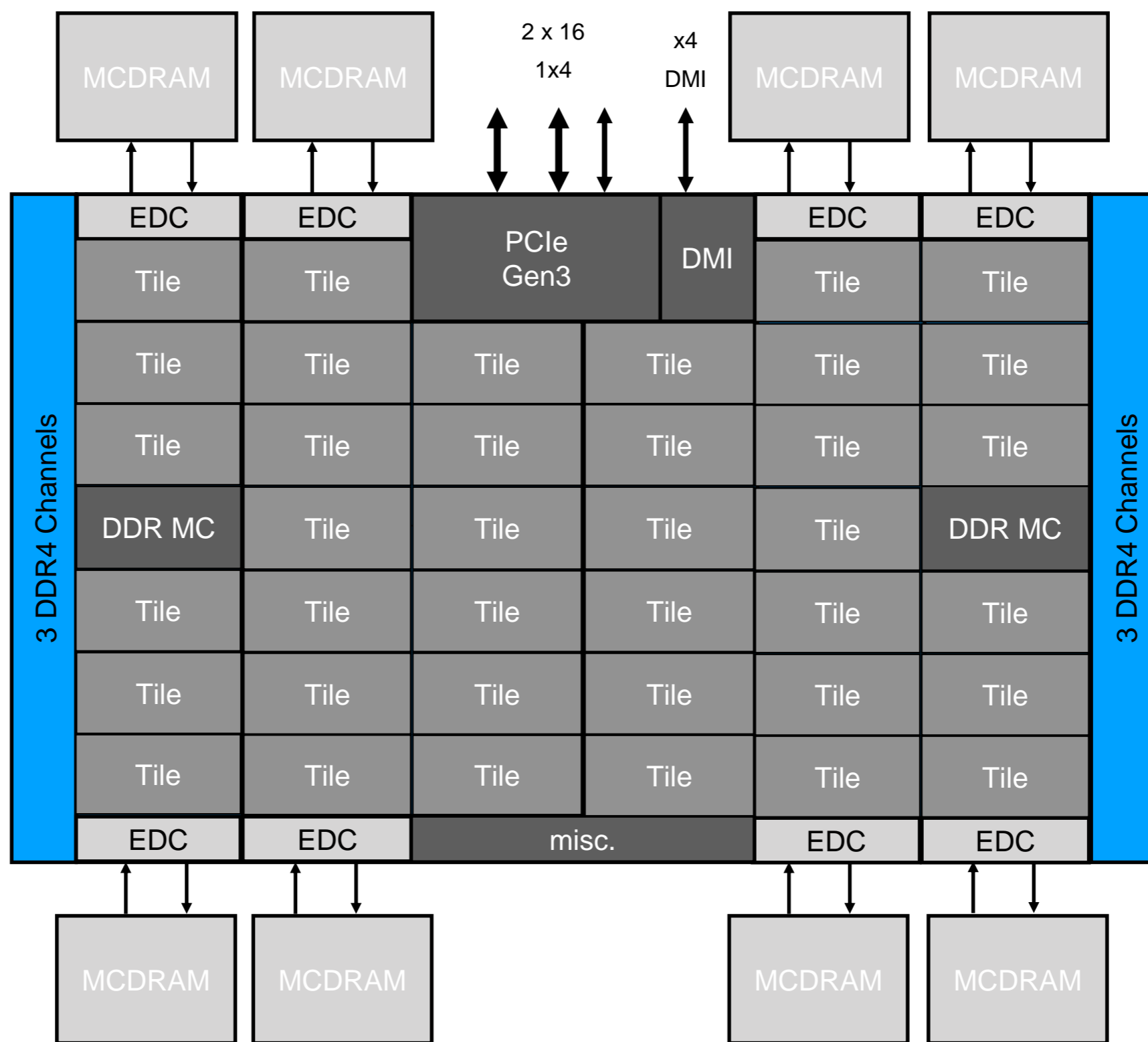
Thomas Rodgers
FCA Design
rodgert@twrodgers.com

# My background

- Disclaimer - I am not a quant

- C++ developer with > 20 years of experience in Finance

- Member of the C++ Standards Committee

  - SG1 - Concurrency and Parallelism

  - SG14 - Games, HFT, Low-Latency systems

- My interest in this project is how to make C++ a better tool for heterogenous multi-core development

# Project Background

- Work in support of research by Prof. Matthew Dixon, Stuart School of Business, Illinois Institute of Technology

  - Help Quants target highly parallel, or vector-parallel architectures, without needing a deep background in vector/parallel software engineering

  - Xeon Phi is one such target architecture

- Need a source of publicly releasable financial model codes

  - QuantLib is one possibility

Knights Landing

Die Layout

# Knights Landing
# Die Features

- Chip organized around the concept of a Tile, each consists of -

  - 2 cores, 32kb L1 Instruction and 32kb L1 Data cache per core

  - 2 vector processing units per core, total 4 per Tile

  - 1M of L2 cache shared between the two cores

- All Tiles are connected to a 2D mesh with > 700GB/s bandwidth

  - Organized into rows and columns of "half" rings that fold upon themselves at the endpoints

  - Enforces a "YX routing" rule, transactions travel vertically to target row, then travel horizontally to destination.

# Knights Landing
# Die Features
# (cont)

- 2 Memory controllers with 3 DDR4 channels each

  - Maximum of 384GB of DDR4

  - Aggregate DDR4 bandwidth is ~90 GB/s

- 8 MCDRAM (High Bandwidth Memory) devices, each 2GB capacity

  - Each MCDRAM is connected to it's own memory controller (EDC)

  - Aggregate MCDRAM bandwidth is > 450 GB/s

  - Can be used in cache mode, flat mode (shares address space with DDR), or a mix of the two

# Knights Landing Core details

- Each Knights Landing die contains 36 tiles of 2 cores each for a total of 72 cores per die

- Each KNL core is based on a significantly modified Atom processor (Silvermont)

- Features -

  - ISA compatible with Haswell (except TSX support)

  - Support for up to 4 simultaneous threads (hyperthreading)

  - Out of order execution

# Knights Landing
# Core details
# (cont)

- Each Knights Landing tile includes two Vector Processing Units per core for a total of 144 VPUs per KNL die

- Each VPU can perform 16 double precision, or 32 single precision floating point operations

- Each VPU supports an extended 512 bit instruction set, aka AVX-512

  - Includes "foundation" support for floating point and scalar vector operations, as well as dedicated support for exponentiation, prefetch, and conflict detection

  - AVX-512 ISA support is available via compiler intrinsics in ICC, GCC, and Clang.

  - Support for AVX-512 compiler intrinsics are also available in ICC, GCC, and Clang

    - see https://software.intel.com/sites/landingpage/IntrinsicsGuide for more information

# Knights Landing Configurations

- The Knights Landing processor, differs from the previous generation Knights Corner processor -

  - Standard Intel Architecture device

  - Fully capable of booting stock operating systems

    - Currently using CentOS 7.4

# Knights Landing Boot-time Configuration

- Unlike standard Intel Xeon CPUs the Knights Landing CPU supports several boot-time configuration parameters

  - Cluster Mode - Configures the on chip mesh in one of

    - all-to-all mode (used with mixed DRAM configurations)

    - quadrant mode (default)

    - Sub-NUMA clustering (SNC)

  - Memory Mode -

    - cache, flat, hybrid

Benchmark Index (MIPS)

# Benchmark Sample

Benchmark Sample

# Observations

- The benchmark suite is single threaded

- Aggregate single core performance for Phi vs. a typical Intel Core/Xeon is lower, this result is generally expected, but…

  - Some individual models >2x slower

- This appears at least in part due to how the QuantLib benchmark suite is structured

  - Significant amount of scalar set-up code inside deeply nested loops

# It's Full of Cores…

# The challenge is using them.

# Scaling up

- Parallelization

  - 72 physical scalar cores

  - High bandwidth on-chip mesh

  - High bandwidth CPU adjacent memory (MCDRAM)

- Vectorization

  - Up to 16 double or 32 single precision floating point operations per vector unit, per clock

  - Two vector units per physical core

# Intel Compiler Collection Auto Parallelization

- Using #pragma directives to guide code generation

- Can parallelize candidate loops with either OpenMP or TBB

  - OpenMP for homogenous workloads

  - TBB for variable workloads

# Auto parallelization is not always a win

# What's going on here?

- Initial guess is that there's not enough work per thread to offset the communication overhead

# Intel Compiler Collection

- Compiler vectorization reports

  - -qopt-report, -qopt-report-phase

- Vectorization Advisor

  - Extensive tools for analyzing and recommending vectorization opportunities

  - Not explicitly tied to code generated by the Intel Compiler

- VTune profiler

  - Sample based and "uncore" counter profiling

# What's going on here?

- VTune Says -

  - ~21 of 272 logical cores utilized, ~7.8% utilization

  - ~57% time spent in serial code

  - ~50% scalar/50% SIMD instruction mix

  - Significant fraction of parallel time is load-imbalanced

  - 83% of pipeline slots remain empty

# What's going on here?

- Top 3 serial hotspots -

  - QuantLib::TrinomialTree::TrinomialTree

  - QuantLib::TreeLattice<QuantLib::OneFactorModel::ShortRateTree>::computeStatePrices

  - malloc

# QuantLib is a non-trivial codebase

- Approximately 2280 source files, ~300k lines of code

- Makes heavy use of allocations, and shared pointers

  - Difficult to optimize without significant restructuring

- Detailed optimization of QuantLib is outside the scope of the work for this project

# Narrowing the scope

- Currently evaluating optimization of

  - Heston model

  - SVD model w/parallel RNG

- Stripped down implementations, not part of QuantLib

- Explicit use of Intel's Math Kernel Libraries where appropriate

# Early results
# Scaling individual models

- Heston Model

  - Peak core utilization ~14%,

  - 75% scalar to 25% SIMD instruction mix

- SVD & Parallel RNG

  - Peak core utilization ~19%

  - 100% Packed SIMD instruction mix

# Next steps

- Long term goal is not focused on individual model optimization

    - Enable Quants to describe compute intensive problems in terms of high level composition of numerical algorithms

- Focus evaluation on scaling up compute utilization via explicitly parallel work partitioning

# Next steps (cont)

- Custom C++ Allocators

  - Using C++ vector types backed by MCDRAM for frequently accessed data

  - Align vector types on cache-line boundary

- Non-Temporal Store

  - Avoid perturbing cache with data that is only ever written (e.g. result sets)