

# Towards ATLAS 4.0

R. Clint Whaley

rcwhaley@lsu.edu

[www.csc.lsu.edu/~whaley](http://www.csc.lsu.edu/~whaley)

Rakib Hasan (rhasan@cct.lsu.edu)

Majedul Sujon (msujon@cct.lsu.edu)

Louisiana State University

Computer Science and Engineering Division (CSC) &  
Center for Computation & Technology (CCT)

May 19, 2016

## Improving Kernels and Exploiting SIMD Vectorization

- ➊ Generic vectorization using `atlas_simd.h` & `atlas_cplxsimd.h`
- ➋ New tuning framework supporting SIMD-friendly storage
- ➌ Extending tuning framework for additional kernels
- ➍ Using FKO/iFKO for backend compilation

## Improving Thread-Level Parallelism

- ➎ Cache-based communication (CBC) for hardware-speed comms  
→ req careful R&D for weakly-ordered caches
- ➏ Thread-pool rather than launch & join to reduce overhead
- ➐ Investigate blocking & scheduling for extreme-scale shared mem  
? heterogenous (big/little), fast scheduling (gatmcctr, gbitvec)

## atlas\_simd.h: Generic VLEN-agnostic SIMD primitives

Must avoid kernel rewrite when VLEN or machine changes:

- Macro ATL\_VLEN tells you how long each vector is
  - ATL\_vvrsuml (l all  $pwr2 \leq VL$ ; eg., vvrsum4/2/1) provides parallel reduction for accumulators (needed for vectorization along dot product dim)
  - Basic ops supported for: ARM32(NEON)/ARM64(AdvSIMD), POWER (VSX), x86 (AVX[2]/SSE[1-3]), and gcc's builtins
  - Model may break down in some cases, so generic macros can lose performance in edge cases
    - Unaligned access handled differently to optimize, esp. POWER
- I need more experience on PWR & ARM

## Key Ideas

- ① All timing/tunings in parallel
  - ② Huge # of NBs: sml probs; target last private/LLC
  - ③ Presently supports M- or K-vec access-major storage
  - ④ Code gens (cpy&amm) using `atlas_simd.h`
- ⇒ Allow user to build custom  $\mu$ kernels and call directly (PCA)

## Challenges

- ① Search not yet optimized
  - ② Serial haswell: MKL:92.1%, ATL:89.2%
  - ③ 24-core haswell-e: MKL:86.6%; ATL:82.8%/85.6%
- ⇒ Multilvl blking, bad tuning, diff format?
- looking at supporting arbitrary formats

GEMM  $\mu$ kernel can probably help all L3 except TRSM

- With very large NB, diagonal blks remain important for perf
  - Add ability to generate register-blocked block-lower storage
  - Now SYRK/TRMM/SYRK supported by changing one loop param in amm  $\mu$ kern
  - Have support for SYRK using k-vec storage  $M_u = N_u$
- Need more info to see if this is worth complexity

## TRSM

- Idea is recursive/multlvl gemm-based
- Not clear we can keep overhead sufficiently low

# iterative Floating Point Kernel Optimizer

iFKO is an iterative and empirical compiler framework

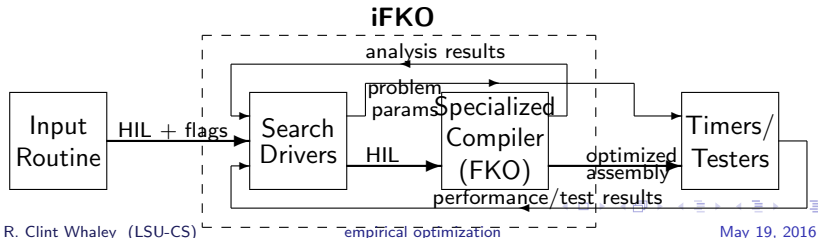
Attempt to extend impact past LA to all HPC computation.

## iFKO composed of:

- ① A collection of search drivers,
- ② a compiler specialized for empirical floating point kernel optimization (FKO)
  - Specialized in analysis, HIL, type & flexibility of transforms

## Key design decisions:

- Iterative & empirical for auto-adaptation
- Transforms done at backend, allowing arch exploitation (eg. SIMD vect, CISC inst formats)
- Kernel annotation markup
- Narrow & deep focus



## Recently published work

- Majedul Haque Sujon, R. Clint Whaley and Qing Yi, "Vectorization Past Dependent Branches Through Speculation", in *Proceedings of the 22nd International Conference on Parallel Architectures and Compilation Techniques (PACT)*, pages 353-362, Edinburgh, Scotland, September, 9-11, 2013.

## Ongoing & recent work in iFKO

- Auto-vectorize mixed-type inst (eg., IMAAX)
- Effectively handle nested loops
- Support for 2-D arrays
- Improved support for register exhaustion
- Outer loop vectorization & vector reduction parallelization
- Half incorporated into ATLAS framework