# Towards Heterogeneous Solvers for Large-Scale Linear Systems

Stylianos I. Venieris, Grigorios Mingas, Christos-Savvas Bouganis stylianos.venieris10@imperial.ac.uk

FPL 2015, London 2 Sept 2015

# Introduction – Solving Linear Systems

• Given  $A \in R^{(m \times n)}$ ,  $b \in R^m$  and  $m \ge n$ :

$$\min_{\boldsymbol{x}\in R^n} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2$$

• Find vector **x** 



# **Introduction – Solving Linear Systems**

• Given  $A \in R^{(m \times n)}$ ,  $b \in R^m$  and  $m \ge n$ :

$$\min_{\boldsymbol{x}\in R^n} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2$$

• Find vector **x** 



2

# **Introduction – Solving Linear Systems**



## Introduction – Data Systems with different structure



# Introduction – Linear Systems in Genetic Analysis

- Real-life Example:
  - Genetic Analysis [1]
  - Search for gene combinations by solving lots of linear systems

Genes Search Space

- Application Characteristics:
  - A lot of linear systems
  - Linear systems of varying size
  - Up to 100,000 rows and a few hundred columns

[1] L. Bottolo and S. Richardson, "Evolutionary Stochastic Search for Bayesian Model Exploration", Bayesian Analysis, vol. 5, no. 3, pp. 583–618, 09 2010.5

# **Our** focus

- Towards heterogeneity for high performance across matrix sizes
- Novel FPGA solver for tall-skinny linear systems
- Modelling framework
  - Performance and resource estimation (compile and runtime)
  - Optimal hardware configuration (compile-time)
- Up to 18x speed-up in GFLOPS across matrix sizes compared to existing works

## **Route Map**

- Background
- Towards Heterogeneity for Performance
- An Enhanced FPGA Solver
- Evaluation Results
- Conclusions

# **Background – Solving Linear Systems**

- The QR factorisation-based methods are dominant because of their properties
- A = QR Orthogonal Q, upper-triangular R

• 
$$Ax = b \Rightarrow QRx = b$$

• 
$$\boldsymbol{Q}^T \boldsymbol{b}$$
 – Matrix-vector product

• 
$$x = R^{-1}Q^Tb$$
 Solution using back-substitution

# The Challenge

- Existing solvers
  - Target CPUs, GPUs and FPGAs
  - Employ different algorithms
  - Tailored to specific matrix sizes
- Key Challenge:

Sustain high performance across matrix sizes

## **Heterogeneity for Performance – Different Solvers**



Anything in between

### **Heterogeneity for Performance – Different Solvers**



Anything in between















# **Compute Engines**

# **Compute Engines:** FPGA Solver

- Based on existing architecture for tall-skinny QR factorisations [2]
- Functionality Extension
  - From QR to Linear Systems workloads
  - Exploited an algorithmic property for acceleration
    - "Concurrent Solution and Factorisation"
- Any no. of rows
- Up to a max no. of columns
  - 5x the max no. of columns of existing FPGA work for the same device

[2] A. Rafique et al., FPL 2012.

# **Compute Engines:** FPGA Solver



- Configurable parameters
  - Size of arithmetic units
  - No. of blocks to be active in parallel in the architecture

## **Concurrent Solution and Factorisation**

$$1. A = QR$$

- Numerically stable algorithms do not return *Q* explicitly
- Additional computations for the reconstruction of Q and for  $Q^T b$

2. 
$$Ax = b \Rightarrow QRx = b$$

3. 
$$Q^T b$$
 – Matrix-vector product

$$4. \quad x = R^{-1}Q^Tb$$

# **Concurrent Solution and Factorisation**

• Compute  $Q^T b$  without forming Q explicitly



Q reconstructed from partial results

# **Compute Engines:** GPU and CPU Solvers

- GPU Solver
  - Based on the state-of-the-art work on tall-skinny QR factorisations [3]
  - Extended its functionality to Linear Systems
    by means of the Concurrent Solution and Factorisation
- CPU Solver
  - Optimised multithreaded linear algebra library (OpenBLAS)

[3] M. Anderson, C. Ballard, J. Demmel, and K. Keutzer, "Communication-Avoiding QR Decomposition for GPUs", in Parallel Distributed Processing Symposium (IPDPS), 2011 IEEE International, May 2011, pp.48-58.

# **Modelling Framework**

- Compile-time Framework
  - Performance in GFLOPS and resource estimation models for the FPGA solver
  - Optimal hardware configuration of the FPGA solver
- Runtime Framework
  - Workload allocation among the available solvers



























# **Evaluation**

# **Experimental Setup**

- FPGA
  - Xilinx Virtex-6 SX475 at 200 MHz with 2016 DSPs
  - Double-precision floating-point
  - Up to 275 columns (n = 275) 5x the max no. of columns of existing FPGA work
  - Any number of rows
  - 84.23% BRAM utilization post place-and-route
  - 99.45% DSP utilization post place-and-route
- GPU
  - NVIDIA Tesla K20
  - 2496 cores at 706 MHz
- *CPU* 
  - Intel i7-4770 at 3.40 GHz, 16 GB RAM, 8 MB cache
  - 4 cores, 8 threads

## **Internal Comparisons**



No. of Columns (n) = 51

### **Internal Comparisons**



No. of Columns (n) = 51

## **Internal Comparisons**



No. of Columns (n) = 51

### **Internal Comparisons**



*No. of Rows* (m) = 6400

## **Internal Comparisons**



*No. of Rows* (m) = 6400

### **External Comparisons**



*No. of Columns* (*n*) = 51 [4] *A. Rafique et al., FPL 2012.* 

### **External Comparisons**



## **Conclusions**

- Different solvers perform better on different matrix sizes
- Using heterogeneous solvers in a complementary way enable the high-performance solution of complex problems in fields such as genetic analysis

# Thank You & Questions ?