

Can Numerical Linear Algebra make it in Nature?

Advice on collaborations with computational biologists

Paolo Bientinesi

in collaboration with Diego Fabregat, Elmar Peise and Yurii Aulchenko

AICES, RWTH Aachen
pauldj@aices.rwth-aachen.de

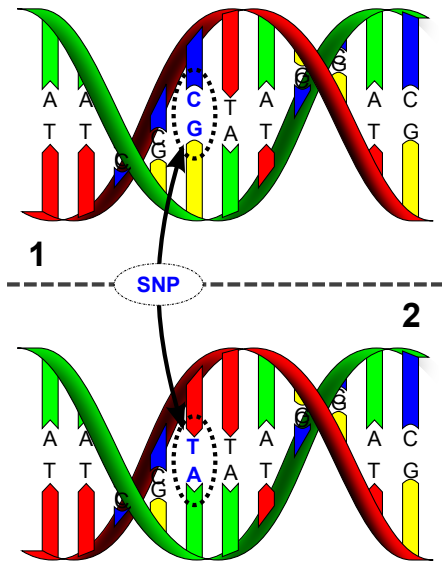
June 9, 2014
Householder Symposium XIX
Spa, Belgium



Deutsche
Forschungsgemeinschaft

DFG

The problem: Genome-Wide Association Studies



Source: David Hall

GWAS:

Correlation between a difference in the genome sequence (SNP) and a difference in the phenotype (observations)



Source: Teri Manolio

The method

- “Mixed models”

The method

- “Mixed models”
- Linear regression with non-independent outcomes

The method

- “Mixed models”
- Linear regression with non-independent outcomes
- Generalized least-square problems

The method

- “Mixed models”
- Linear regression with non-independent outcomes
- Generalized least-square problems

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$$

- y : **phenotype** (outcome; vector of observations)
E.g.: height, blood pressure for a set of people
- X : **genome measurements** and covariates
(design matrix; predictors)
E.g.: sex and age over height
- M : **dependencies** between observations
E.g.: tall parents have tall children
- b : **relation** between a variation in the outcome (y)
and a variation in the genome sequence (X)

$$b := \begin{pmatrix} \text{---} \\ X^T \\ \begin{matrix} \square \\ M \\ \square \end{matrix} \\ \text{---} \\ X \\ \text{---} \end{pmatrix}^{-1} \begin{pmatrix} \text{---} \\ X^T \\ \begin{matrix} \square \\ M \\ \square \end{matrix} \\ \text{---} \end{pmatrix}^{-1} y$$

- $X \in \mathcal{R}^{n \times p}$

- $y \in \mathcal{R}^n$

- $b \in \mathcal{R}^p$

- $M \in \mathcal{R}^{n \times n}$

“SNP”

“trait”

“genetic effect”

“covariance matrix”

- $n \approx 1,000 - 50,000$

- $p \in [1, \dots, 20]$

- M : SPD

(Wrong) Problem definition

Isolated problem instances

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$$

“to be repeated millions of times”

(Wrong) Problem definition

Isolated problem instances

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$$

“to be repeated millions of times”

↓

for $i = 1, \dots, m$ $m \approx 10^6 - 10^7$

$$b_i := (X_i^T M_i^{-1} X_i)^{-1} X_i^T M_i^{-1} y_i$$

(Wrong) Problem definition

Isolated problem instances

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$$

“to be repeated millions of times”

⇓

for $i = 1, \dots, m$ $m \approx 10^6 - 10^7$

$$b_i := (X_i^T M_i^{-1} X_i)^{-1} X_i^T M_i^{-1} y_i$$

⇓

for $i = 1, \dots, m$

$$L_i L_i^T = M_i \quad \text{CHOL}$$

$$X_i := L_i^{-1} X_i \quad \text{TRSM}$$

$$y_i := L_i^{-1} y_i \quad \text{TRSV}$$

$$b_i := \text{OLS}(X_i, y_i) \quad O(n^3 m)$$

(Wrong) Problem definition

Isolated problem instances

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$$

“to be repeated millions of times”

⇓

for $i = 1, \dots, m$

$m \approx 10^6 - 10^7$

$$b_i := (X_i^T M_i^{-1} X_i)^{-1} X_i^T M_i^{-1} y_i$$

⇓

for $i = 1, \dots, m$

$$L_i L_i^T = M_i \quad \text{CHOL}$$

$$X_i := L_i^{-1} X_i \quad \text{TRSM}$$

$$y_i := L_i^{-1} y_i \quad \text{TRSV}$$

$$b_i := \text{OLS}(X_i, y_i)$$

$O(n^3 m)$

Problem definition #1

One-dimensional sequence of GLS problems

for $i = 1, \dots, m$ $m \approx 10^6 - 10^7$

$$b_i := (X_i^T M^{-1} X_i)^{-1} X_i^T M^{-1} y$$

and

$$X_i = [X_L | X_{Ri}], \quad \text{where } X_{Ri} \in \mathcal{R}^{n \times 1}$$

Problem definition #1

One-dimensional sequence of GLS problems

for $i = 1, \dots, m$ $m \approx 10^6 - 10^7$

$$b_i := (X_i^T M^{-1} X_i)^{-1} X_i^T M^{-1} y$$

and

$$X_i = [X_L | X_{Ri}], \quad \text{where } X_{Ri} \in \mathcal{R}^{n \times 1}$$

↓

$$b_i := \left(\begin{bmatrix} X_L^T \\ X_{Ri}^T \end{bmatrix} M^{-1} [X_L | X_{Ri}] \right)^{-1} \begin{bmatrix} X_L^T \\ X_{Ri}^T \end{bmatrix} M^{-1} y$$

$$b_i := \left[\begin{array}{c|c} \overline{X}_L^T \overline{X}_L & \star \\ \hline \overline{X}_{Ri}^T \overline{X}_L & \overline{X}_{Ri}^T \overline{X}_{Ri} \end{array} \right]^{-1} \begin{bmatrix} \overline{X}_L^T \\ \overline{X}_{Ri}^T \end{bmatrix} L^{-1} y$$

Problem definition #1

One-dimensional sequence of GLS problems

for $i = 1, \dots, m$ $m \approx 10^6 - 10^7$

$$b_i := (X_i^T M^{-1} X_i)^{-1} X_i^T M^{-1} y$$

and

$$X_i = [X_L | X_{Ri}], \quad \text{where } X_{Ri} \in \mathcal{R}^{n \times 1}$$

↓

$$LL^T = M;$$

$$\bar{X}_L := L^{-1} X_L;$$

$$\bar{y} := L^{-1} y;$$

$$S_{TL} := \bar{X}_L^T \bar{X}_L;$$

$$\bar{b}_T := \bar{X}_L^T y;$$

for $i = 1, \dots, m$

$$\bar{X}_{Ri} := L^{-1} X_{Ri} \quad \text{TRSV}$$

$$S_{BLi} := X_{Ri}^T \bar{X}_L$$

$$S_{BRi} := \bar{X}_{Ri}^T \bar{X}_{Ri}$$

⋮

$$b_i := S_i^{-1} \bar{b}_i \quad O(n^2 m)$$

Problem definition #2

Two-dimensional sequence of GLS problems

$$b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j$$

↓

for $i = 1, \dots, m$

for $j = 1, \dots, t$

$$b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j$$

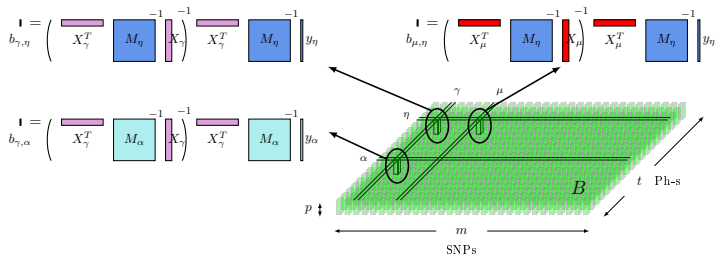
and

$X_i = [X_L | X_{Ri}]$, with $X_{Ri} \in \mathcal{R}^{n \times 1}$

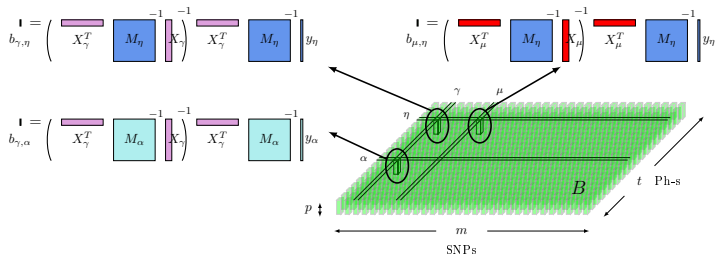
and $\text{SPD}(M_j)$

$m \approx 10^6 - 10^7$; $t = 1$ or $\approx 10^3 - 10^5$

Overview of the full problem



Overview of the full problem



Problem size

$M \in \mathbb{R}^{n \times n}$	$1000 \leq n \leq 100k$	7.5MBs – 74.5GBs
$X_{Ri}, y_j \in \mathbb{R}^n$		8 – 780KBs
$b_{ij} \in \mathbb{R}^p$	$3 \leq p \leq 20$	24 – 160 Bytes
Total	$m \leq 10^8, t \leq 10^5$	1.5 – 100s Terabytes

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = ? \end{cases}$$

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = \sigma_j^2 (h_j^2 \Phi + (1 - h_j^2) I) \end{cases}$$

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = \sigma_j^2 (h_j^2 \Phi + (1 - h_j^2) I) \end{cases}$$

$$\Phi = Q \Lambda Q^T$$

$$\Rightarrow M_j = Q (\alpha_j \Lambda + \beta_j I) Q^T$$

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = \sigma_j^2 (h_j^2 \Phi + (1 - h_j^2) I) \end{cases}$$

$$\Phi = Q \Lambda Q^T$$

$$\Rightarrow M_j = Q (\alpha_j \Lambda + \beta_j I) Q^T$$

$$\Rightarrow M_j^{-1} = Q (\alpha_j \Lambda + \beta_j I)^{-1} Q^T$$

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = \sigma_j^2 (h_j^2 \Phi + (1 - h_j^2) I) \end{cases}$$

$$\Phi = Q \Lambda Q^T$$

$$\Rightarrow M_j = Q (\alpha_j \Lambda + \beta_j I) Q^T$$

$$\Rightarrow M_j^{-1} = Q (\alpha_j \Lambda + \beta_j I)^{-1} Q^T$$

$$b_{ij} :=$$

$$(X_i^T Q D_j^{-1} Q^T X_i)^{-1} X_i^T Q D_j^{-1} Q^T y_j$$

$$\begin{cases} b_{ij} := (X_i^T M_j^{-1} X_i)^{-1} X_i^T M_j^{-1} y_j \\ M_j = \sigma_j^2 (h_j^2 \Phi + (1 - h_j^2) I) \end{cases}$$

$$\Phi = Q \Lambda Q^T$$

$$\Rightarrow M_j = Q (\alpha_j \Lambda + \beta_j I) Q^T$$

$$\Rightarrow M_j^{-1} = Q (\alpha_j \Lambda + \beta_j I)^{-1} Q^T$$

$$b_{ij} :=$$

$$(X_i^T Q D_j^{-1} Q^T X_i)^{-1} X_i^T Q D_j^{-1} Q^T y_j$$

```

1  QΛQT = Φ
2  for 1 ≤ i ≤ m
3      X'i := QT Xi           GEMM
4  for 1 ≤ j ≤ t
5      y'j := QT yj         GEMV
6  for 1 ≤ j ≤ t
7      Dj := σj2 (hj2Λ + (1 - hj2)I)
8      KjKjT = Dj-1/2
9      vj := KjT y'j
10     for 1 ≤ i ≤ m
11         Wij := KjT X'i
12         Sij := WijT Wij
13         bij := WijT vj
14         bij := Sij-1 bij

```

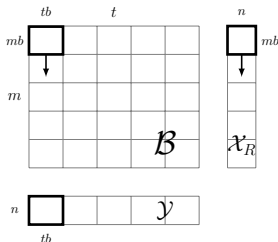
Cost: $O(n^2 mt)$ vs. $O(nmt)$

(Big) Data management

operands

X'_s	input	100s GBs – 2 TBs	streaming from disk
y'_s	input	1 – 10 GBs	streaming from disk
M	input	MBs – 80 GBs	read once
b'_s	output	100s MBs or 10s TBs	streaming to disk

Tiling: $tb, mb?$



Does M fit in memory?

Does M fit in memory?

- YES \Rightarrow single node + multithreading
streaming HD \leftrightarrow CPU, double buffering, in-core implementation

Does M fit in memory?

- YES \Rightarrow single node + multithreading
streaming HD \leftrightarrow CPU, double buffering, in-core implementation

Does M fit in GPU-memory?

- Yes \Rightarrow accelerator
streaming HD \leftrightarrow CPU \leftrightarrow GPU, triple+double buffering, CPU+GPU implementation

Does M fit in memory?

- YES \Rightarrow single node + multithreading
streaming HD \leftrightarrow CPU, double buffering, in-core implementation

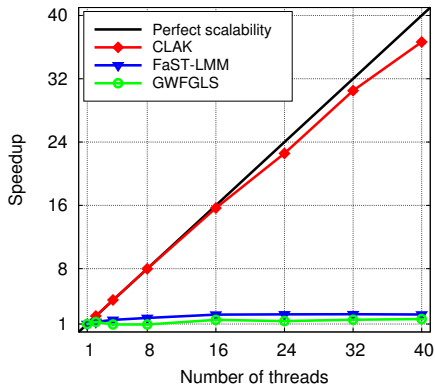
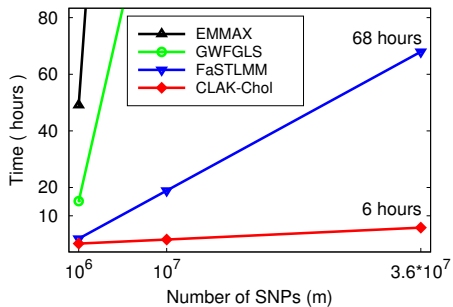
Does M fit in GPU-memory?

- Yes \Rightarrow accelerator
streaming HD \leftrightarrow CPU \leftrightarrow GPU, triple+double buffering, CPU+GPU implementation
- NO \Rightarrow distributed memory + hybrid parallelism
partitioning + streaming HD \leftrightarrow CPUs, double buffering, data distribution

Results

$t = 1$

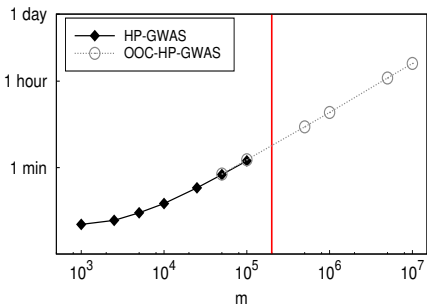
Single-trait analysis: one-dimensional sequence



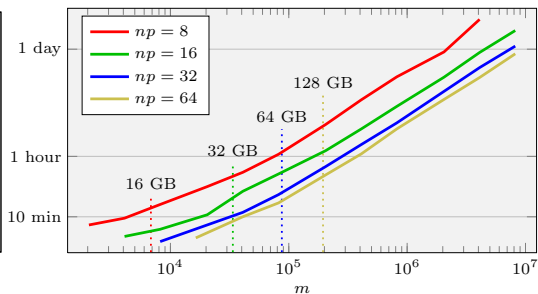
Results beyond memory capacity

Single-trait analysis: one-dimensional sequence

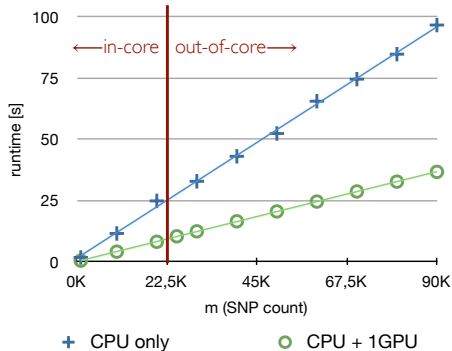
Single node



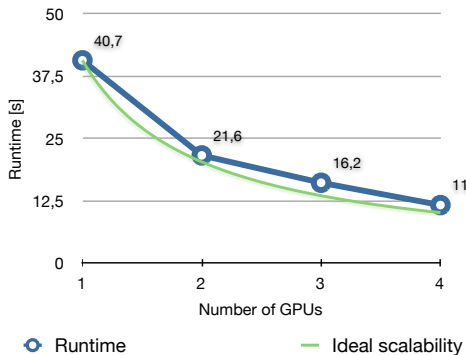
MPI



1 GPU



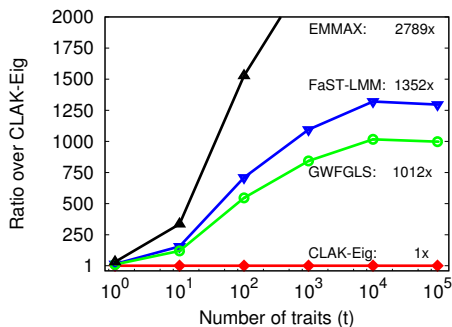
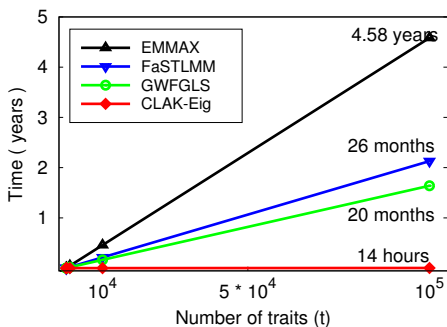
Scalability



Results

$t \gg 1$ – full grid

Multi-trait analysis: two-dimensional sequence



“Can Numerical Linear Algebra make it in Nature?”

...

“Can Numerical Linear Algebra make it in Nature?”

almost

“Can Numerical Linear Algebra make it in Nature?”

almost

- Multi-trait analysis:
20 months → 14hours
- 1000x speedups
- No need for supercomputers
- Elegant solution: knowledge

“Can Numerical Linear Algebra make it in Nature?”

almost

- Multi-trait analysis:
20 months → 14hours
- 1000x speedups
- No need for supercomputers
- Elegant solution: knowledge
- “Real” data vs. random data
- “Algorithms” vs. ”methods”
- (Too) new methodology
- Size of readership
- New findings

“Can Numerical Linear Algebra make it in Nature?”

almost

- Multi-trait analysis:
20 months → 14hours
- 1000x speedups
- No need for supercomputers
- Elegant solution: knowledge
- “Real” data vs. random data
- “Algorithms” vs. ”methods”
- (Too) new methodology
- Size of readership
- New findings

- D.F. and P.B., “Computing Petaflops over Terabytes of Data: The Case of GWAS”, TOMS 2014.
- D. F., Y.A. and P.B., “Sequences of Generalized. Least-Squares Problems on SMP Arch.”, AMC 2014.
- L.B. and P.B., “Streaming Data from HDD to GPUs for Sustained Peak Performance”, Euro-Par 2013.
- E.P., D.F., Y.A. and P.B., “Large-scale Whole Genome Association Analysis”, PBio 2013 (EuroMPI'13).