LAPACK-Style Codes for the QR Factorization of Banded Matrices

Alfredo Remón-Gómez, Enrique S. Quintana-Ortí, Gregorio Quintana-Ortí {remon,quintana,gquintan}@icc.uji.es

Universidad Jaume I de Castellón (Spain)

9th International Workshop on State-of-the-Art in Scientific and Parallel Computing





1 The QR factorization

• Householder transformations for the QR factorization

Routines for the QR factorization of general band matrices Routine GBBQRF



QR factorization



QR factorization for general matrices

Given $A \in \mathbb{R}^{m \times n}$, with $m \ge n$, its QR factorization is given by

$$A = Q \cdot R,$$

Where $Q \in \mathbb{R}^{m \times m}$ is orthogonal and $R \in \mathbb{R}^{m \times n}$ is singular

QR factorization for general band matrices

QR factorization for general band matrices

Given $A \in \mathbb{R}^{m \times n}$, with upper and lower bandwidth ku and kl respectively, then if QR factorization with Householder transformations is computed carefully,

- $\bullet \ Q$ is lower triangular with lower bandwidth kl
- R is upper triangular with upper bandwidth ku + kl

Householder transformations



The Householder transformations can be employed to annihilate elements from a vector or matrix

Householder transformation

• Consider vector
$$x = \begin{pmatrix} \chi_0 \\ x_1 \end{pmatrix}$$

• Then, $(I - \beta u u^T) x = \begin{pmatrix} \bar{\chi}_0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} -sign(\chi_0) \|x\|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$
With $\beta = \frac{2}{u^T u}$, $u = \begin{pmatrix} \frac{1}{x_1/v_0} \end{pmatrix}$ and $v_0 = \chi_0 + sign(x_0) \|x\|_2$

Where $\left(I - \beta u u^T\right) x$ is the Householder transformation

Householder transformations for the QR factorization

j i

Objective

Transform the general band matrix \boldsymbol{A} into an upper triangular band matrix



$$(I - \beta u_1 u_1^T) (I - \beta u_2 u_2^T) \dots (I - \beta u_j u_j^T) \cdot A$$

Householder transformations for the QR factorization

j,

Objective

Transform the general band matrix \boldsymbol{A} into an upper triangular band matrix



$$(I - \beta u_1 u_1^T) \cdots (I - \beta u_j u_j^T) \cdots (I - \beta u_k u_k^T) \cdot A$$

Householder transformations for the QR factorization

j i

Objective

Transform the general band matrix \boldsymbol{A} into an upper triangular band matrix



$$(I - \beta u_1 u_1^T) (I - \beta u_2 u_2^T) \dots (I - \beta u_n u_n^T) \cdot A$$

8

Routines for the QR factorization of general band matrices

New routines for the QR factorization of general band matrices

New routines implemented

- GBBQR2
 - Unblocked routine
 - Based on BLAS-2 kernels
- GBBQRF
 - Blocked routine
 - Based on BLAS-3 kernels

We will focus on the blocked routine ${\rm GBBQRF}$

ji

Routine GBBQRF

Functionality

Computes the QR factorization of a general band matrix

Main characteristics

- Exploits the matrix structure so the memory requirements and the computational cost are reduced
- Use of Householder transformations
- Based on BLAS-3 operations

GBBQRF (Householder transformations)

Householder transformations make the resultant matrix R becomes an upper triangular band with bandwidth $\leq ku+kl$



Routines for the QR factorization of general band matrices Routine GBBQRF

GBBQRF (Storage scheme)

- Matrix A is stored following a variant of the compact storage scheme employed by BLAS and LAPACK routines
- In this variant, kl extra upper diagonals are stored, initially padded with zeros



 $m = n = 6, \, kl = 2 \text{ and } ku = 1$

 This scheme is analogous to the one employed in LAPACK routines for the LU factorization

QR Factorization of Banded Matrices

GBBQRF (Storage scheme)

- $\bullet\,$ Following the LAPACK-Style, matrices the reflectors and $R\,$ replace matrix $A\,$
- The extra kl upper-diagonals padded with zeros are necessary to store all the elements of ${\cal R}$

*	*	*	0	0	0
*	*	0	0	0	0
*	$\alpha_{_{01}}$	α_{12}	$\alpha_{_{23}}$	$\alpha_{_{34}}$	$\alpha_{_{45}}$
α_{00}	α_{11}	$\alpha_{_{22}}$	$\alpha_{_{33}}$	$\alpha_{_{44}}$	α_{55}
α_{10}	$\alpha_{_{21}}$	$\alpha_{_{32}}$	$\alpha_{_{43}}$	$\alpha_{_{54}}$	*
α_{20}	$\alpha_{_{31}}$	$\alpha_{_{42}}$	α_{53}	*	*



Routines for the QR factorization of general band matrices Routine GBBQRF

GBBQRF (Algorithm)

$_{\rm GBBQRF}$ implements an iterative algorithm with a 5×5 partitioning



The operations performed in each iteration are

- The QR factorization of nb columns of A is computed
- The next ku + kl columns of A are updated

ji

Routines for the QR factorization of general band matrices Routine GBBQRF

GBBQRF (Algorithm)

$_{\rm GBBQRF}$ implements an iterative algorithm with a 5×5 partitioning



- Blocks in gray have been computed
- Blocks in blue form the active region (updated during the iteration)
- Blocks in green have not been accessed

GBBQRF (Algorithm)

During the current iteration, only blocks in the active region will be accessed



Note that the lower triangular part of the $(nb \times nb)$ block positioned at the bottom of A_{21} is out of the band

Routines for the QR factorization of general band matrices Routine GBBQRF

GBBQRF (Algorithm)



Step 1: QR factorization of blocks A_{11} and A_{21}

• Using the new non-blocked BLAS-2 routine GBBQR2, the QR factorization of blocks A₁₁ and A₂₁ is obtained



$$(\text{GBBQR2}) \quad \left[\begin{array}{c} A_{11} \\ A_{21} \end{array} \right] \quad = \text{QR} \left(\left[\begin{array}{c} A_{11} \\ A_{21} \end{array} \right], \ r \right),$$

GBBQRF (Algorithm)





• Copy A_{11} and A_{21} into a rectangular workspace where the whole block A_{21} is stored



$$U = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix},$$

STRIL(U) = 0,

Routines for the QR factorization of general band matrices Routine GBBQRF

GBBQRF (Algorithm)



Step 2.2: Apply the Householder transformations to blocks A_{12} and A_{22}

• Compute of the Householder transformation and update of ${\cal A}_{12}$ and ${\cal A}_{22}$



(LARFT)
$$T = \text{LARFT}(U, r),$$

(LARFB) $\begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} = (I - UTU^T) \cdot \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix}$

GBBQRF (Algorithm)

- The lower triangular part of the $nb \times nb$ block at the bottom of A_{21} contains only null elements
- No attemption to exploit this structure is done because:
 - As $b \ll kl$ no huge savings are expected Dealing with this structure would increment the number of invocations to BLAS routines

Step 3: Move boundaries

• The active region moves along the diagonal *nb* rows and columns

Experimental Results

١	L	5	
1	1	Ì.	
		L Í	ע

Platform	Architecture	Frequency (GHz)	L2 cache (KBytes)	L3 cache (MBytes)	RAM (GBytes)
RA	Intel Xeon	2.4	512	-	1
CATON	Intel Itanium2	1.5	256	4	4

Platform	Compiler	Optimization	BLAS	Operating
		flags		System
RA	gcc 3.3.5	-03	Goto BLAS 1.15	Linux 2.4.27
			MKL 8.1	
CATON	icc 9.0	-03	Goto BLAS 1.15	Linux 2.4.21
			MKL 8.0	

Results correspond to matrices of order n=5,000 with different bandwidths (from 10 to 1250) and block sizes (from 1 to 100)

Experimental results RA (1 proc.)



QR factorization of Band Matrices on RA (1 processor)



Experimental results RA (2 proc.)



QR factorization of Band Matrices on RA (2 processors)



Experimental results CATON (1 proc.)



Experimental results CATON (4 proc.)



QR factorization of Band Matrices on CATON (4 processors)



Conclusions



- Two new LAPACK-Style routines for the computation of the QR factorization of band matrices have been presented
- Both routines exploit the matrix structure to reduce the memory requirements and the number of computations
- The new codes have been evaluated in two current platforms, showing the higher performance of the blocked routine
- Results show that the operation presents a limited degree of parallelism

Thanks.