

## InterFLOP Meeting

6 October 2021

# Adaptive Precision Sparse Matrix–Vector Product and its Application to Krylov Solvers

**Theo Mary**

Sorbonne Université, CNRS, LIP6

Joint work with

**Stef Graillat**, **Fabienne Jézéquel**, and **Roméo Molina**

# Today's floating-point landscape

		Bits			
		Signif. ( $t$ )	Exp.	Range	$u = 2^{-t}$
bfloat16	B	8	8	$10^{\pm 38}$	$4 \times 10^{-3}$
fp16	H	11	5	$10^{\pm 5}$	$5 \times 10^{-4}$
fp32	S	24	8	$10^{\pm 38}$	$6 \times 10^{-8}$
fp64	D	53	11	$10^{\pm 308}$	$1 \times 10^{-16}$
fp128	Q	113	15	$10^{\pm 4932}$	$1 \times 10^{-34}$

- Low precision increasingly supported by hardware
- **Great benefits:**
  - Reduced **storage**, data movement, and communications
  - Increased **speed** on emerging hardware (**16x** on A100 from fp32 to fp16/bfloat16)
  - Reduced **energy** consumption (**5x** with fp16, **9x** with bfloat16)
- **Some risks too:**
  - Low precision (large  $u$ )
  - Narrow range

Mix several precisions in the same code with the goal of

- Getting the **performance benefits of low precisions**
- While preserving the **accuracy and stability of the high precision**

**Terminology varies:** Mixed precision, Multiprecision, Adaptive precision, Variable precision, Transprecision, Dynamic precision, ...

Mix several precisions in the same code with the goal of

- Getting the **performance benefits of low precisions**
- While preserving the **accuracy and stability of the high precision**

**Terminology varies:** Mixed precision, Multiprecision, Adaptive precision, Variable precision, Transprecision, Dynamic precision, ...

**How** to select the right precision for the right variable/operation

- **Precision tuning:** autotuning based on the source code
  - ▲ Does not need any understanding of what the code does
  - ▼ Does not have any understanding of what the code does
- **This work: exploit as much as possible the knowledge we have about the code**

## Algorithm

Example:  $Ax = b$   
in  $p$  precisions  $u_1, \dots, u_p$



Factorize  $A = LU$  **at low precision**

Solve  $Ax_1 = b$  via  $x_1 = U^{-1}(L^{-1}b)$

**repeat at high precision**

$$r_i = b - Ax_i$$

Solve  $Ad_i = r_i$  via  $d_i = U^{-1}(L^{-1}r_i)$

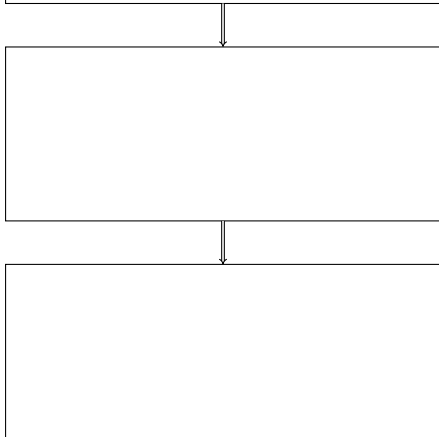
$$x_{i+1} = x_i + d_i$$

**until** converged

- $O(n^3)$  flops at low precision,  $O(n^2)$  flops per iteration at high precision
- What guarantees of convergence? What does "low" and "high" mean?

## Algorithm

Example:  $Ax = b$   
in  $p$  precisions  $u_1, \dots, u_p$



## Algorithm

Example:  $Ax = b$   
in  $p$  precisions  $u_1, \dots, u_p$

## Analysis

$$(A + \Delta A)\hat{x} = b$$
$$\|\Delta A\| \leq f(u_1, \dots, u_p)\|A\|$$
$$\|\hat{x} - x\| \leq g(u_1, \dots, u_p)\|x\|$$



Factorize  $A = LU$  **at precision  $u_f$**

Solve  $Ax_1 = b$  via  $x_1 = U^{-1}(L^{-1}b)$  **at precision  $u_f$**

**repeat**

$r_i = b - Ax_i$  **at precision  $u_r$**

Solve  $Ad_i = r_i$  via  $d_i = U^{-1}(L^{-1}r_i)$  **at precision  $u_f$**

$x_{i+1} = x_i + d_i$  **at precision  $u$**

**until** converged

- Theorem from [Carson and Higham \(2018\)](#) : provided that  $\kappa(A)u_f < 1$ , we reach  $\|\hat{x} - x\| \leq (u + u_r\kappa(A))\|x\|$   
 $\Rightarrow$  **Up to  $\kappa(A) = O(10^3)$  for fp16 factorization**

Factorize  $A = LU$  **at precision  $u_f$**

Solve  $Ax_1 = b$  via  $x_1 = U^{-1}(L^{-1}b)$  **at precision  $u_f$**

**repeat**

$r_i = b - Ax_i$  **at precision  $u_r$**

Solve  $Ad_i = r_i$  with preconditioned GMRES

**at precision  $u_g$**  except matvecs **at precision  $u_p$**

$x_{i+1} = x_i + d_i$  **at precision  $u$**

**until** converged

- Theorem from [Carson and Higham \(2018\)](#) : provided that  $\kappa(A)u_f < 1$ , we reach  $\|\hat{x} - x\| \leq (u + u_r\kappa(A))\|x\|$   
 $\Rightarrow$  **Up to  $\kappa(A) = O(10^3)$  for fp16 factorization**
- Theorem from [Amestoy et al. \(2021\)](#) : provided that  $(u_g + \kappa(A)u_p)\kappa(A)^2u_f^2 < 1$  we reach the same accuracy  
 $\Rightarrow$  **Up to  $\kappa(A) = O(10^{11})$  for fp16 factorization**

## Algorithm

Example:  $Ax = b$   
in  $p$  precisions  $u_1, \dots, u_p$

## Analysis

$$(A + \Delta A)\hat{x} = b$$
$$\|\Delta A\| \leq f(u_1, \dots, u_p)\|A\|$$
$$\|\hat{x} - x\| \leq g(u_1, \dots, u_p)\|x\|$$

## Algorithm

Example:  $Ax = b$   
in  $p$  precisions  $u_1, \dots, u_p$

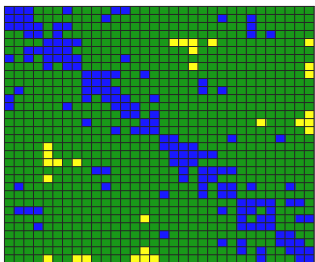
## Analysis

$$(A + \Delta A)\hat{x} = b$$
$$\|\Delta A\| \leq f(u_1, \dots, u_p)\|A\|$$
$$\|\hat{x} - x\| \leq g(u_1, \dots, u_p)\|x\|$$

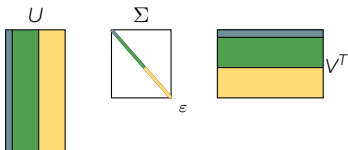
## Application

Take into account any  
**special property of  $A$**  in the  
specific application at hand

# Block low-rank (BLR) matrices



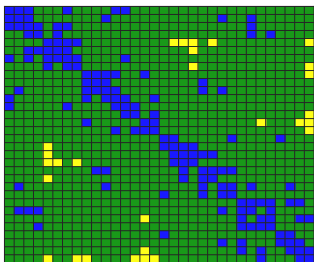
double single half



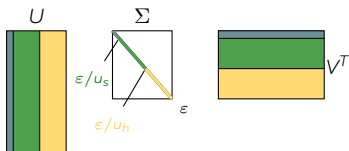
- BLR matrices arise in several applications (PDEs, covariance matrices, etc.)
- Blocks are compressed with low-rank approximations, e.g. via truncated SVD

**Why does mixed precision work here ?**

# Block low-rank (BLR) matrices



double single half



- BLR matrices arise in several applications (PDEs, covariance matrices, etc.)
- Blocks of norm less than  $\epsilon/u_i$  can be stored in precision  $u_i$
- Blocks are compressed with low-rank approximations, e.g. via truncated SVD
- Singular vectors associated with singular values less than  $\epsilon/u_i$  can be stored in precision  $u_i$

**Why does mixed precision work here ?**

 Amestoy et al. (2021b)

- Why can we store “less important” data in lower precision ?  
Because small elements produce small errors :

$$|fl(a \text{ op } b) - a \text{ op } b| \leq u |a \text{ op } b|, \quad \text{op} \in \{+, -, *, \div\}$$

## ⇒ Opportunity for mixed precision !

- BLR matrices [Amestoy et al. \(2021b\)](#) [Abdulah et al. \(2021\)](#)
- Low-rank approximations, SVD [Amestoy et al. \(2021b\)](#)
- Randomized SVD [Connolly, Higham, and Pranesh \(SIAM AN 2021\)](#)
- Quantized dot product [Diffenderfer, Osei-Kuffuor, & Menon \(2021\)](#)
- SpMV [Ahmad, Sundar and Hall \(2020\)](#)
- Multiword matrix multiplication [Fasi et al. \(SIAM LA 2021\)](#)
- Random matrix multiplication [Higham and Mary \(2020\)](#)
- Block Jacobi and SPAI preconditioners [Anzt et al. \(2019\)](#)
- Runge Kutta [Crocì and de Souza \(2021\)](#)

# Adaptive precision at the variable level ?

- Pushing adaptive precision to the extreme: can we benefit from storing **each variable** in a (possibly) different precision?
  - Same granularity as precision autotuning, but different method to select precisions
  - Example:  $Ax = b$  with adaptive precision for each  $A_{ij}$ 
    - **Is it worth it ?**  
Need to have elements of **widely different magnitudes**, and yet **not structured** in any obvious way (by blocks or columns, etc.)
    - **Is it practical ?**  
Probably not for compute-bound applications, but could it work for **memory-bound** ones?
- ⇒ Natural candidate: **sparse matrices**



# Sparse matrix–vector product (SpMV)

$$y = Ax, A \in \mathbb{R}^{m \times n}$$

```
for  $i = 1:m$  do
   $y_i = 0$ 
  for  $j \in \text{nnz}_i(A)$  do
     $y_i = y_i + a_{ij}x_j$ 
  end for
end for
```

- Standard backward error analysis: if  $y = Ax$  is performed in a uniform precision  $\epsilon$ , we obtain

$$|\hat{y}_i - y_i| \leq n_i \epsilon \sum_{j \in \text{nnz}_i(A)} |a_{ij}x_j|$$

- **Idea:** store elements of  $A$  in a precision inversely proportional to their magnitude (**smaller elements in lower precision**)

```
for  $i = 1:m$  do
   $y_i = 0$ 
  for  $k = 1:p$  do
     $y_i^{(k)} = 0$ 
    for  $j \in \text{nnz}_i(A)$  do
      if  $a_{ij}x_j \in B_{ik}$  then
         $y_i^{(k)} = y_i^{(k)} + a_{ij}x_j$  at precision  $u_k$ 
      end if
    end for
     $y_i = y_i + y_i^{(k)}$ 
  end for
end for
```

- Split row  $i$  of  $A$  into  $p$  buckets  $B_{ik}$  and sum elements of  $B_{ik}$  in precision  $u_k$
- Backward error analysis:  $|\hat{y}_i^{(k)} - y_i^{(k)}| \leq n_i^{(k)} u_k \sum_{a_{ij}x_j \in B_{ik}} |a_{ij}x_j|$

# Building the buckets

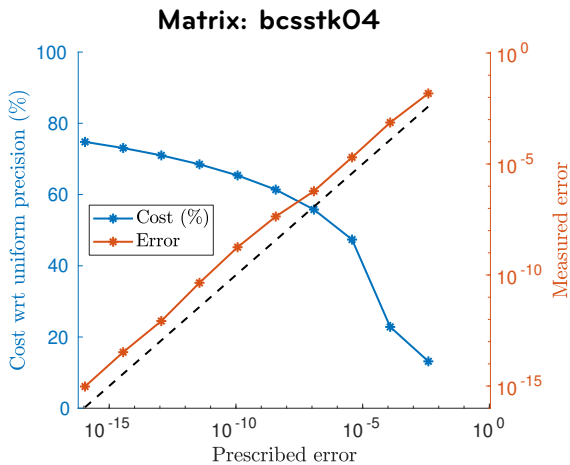
- $|\hat{y}_i^{(k)} - y_i^{(k)}| \leq n_i^{(k)} u_k \sum_{a_{ij}x_j \in B_{ik}} |a_{ij}x_j|$

⇒ Build the buckets such that  $u_k \sum_{a_{ij}x_j \in B_{ik}} |a_{ij}x_j| \approx \varepsilon \sum_j |a_{ij}x_j|$

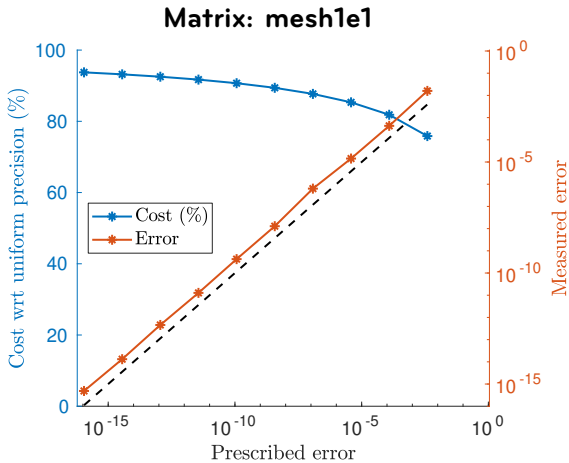
- By setting  $B_{ik}$  to the interval  $(\varepsilon\beta_i/u_{k+1}, \varepsilon\beta_i/u_k]$ , we obtain

$$|\hat{y}_i^{(k)} - y_i^{(k)}| \leq n_i^{(k)} \varepsilon\beta_i \text{ and so } |\hat{y}_i - y_i| \leq n_i \varepsilon\beta_i$$

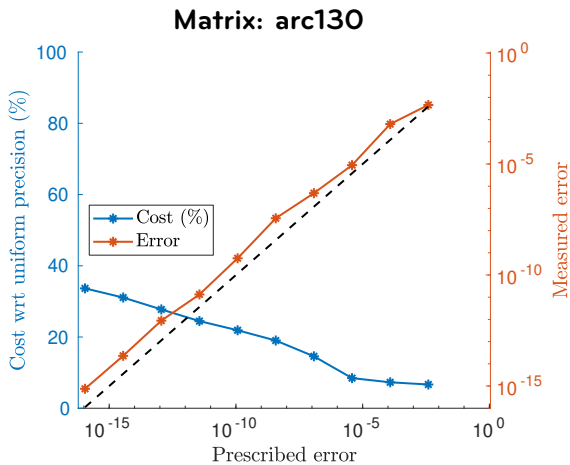
- Two possible choices for  $\beta_i$ :
  - $\beta_i = \sum_j |a_{ij}x_j| \Rightarrow$  guarantees  $O(\varepsilon)$  **componentwise** backward error
  - $\beta_i = \|A\| \|x\| \Rightarrow$  guarantees  $O(\varepsilon)$  **normwise** backward error



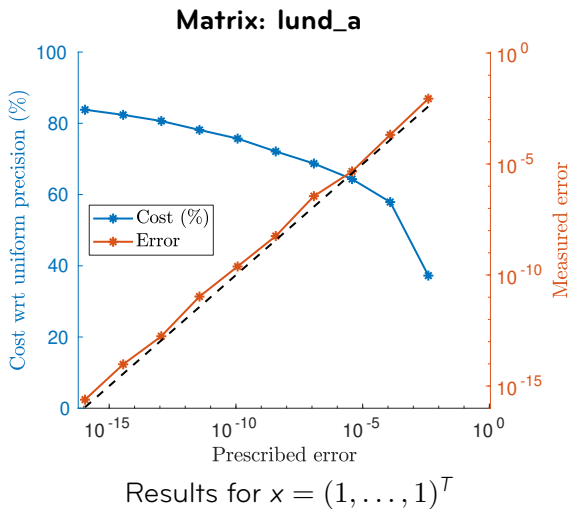
Results for  $x = (1, \dots, 1)^T$

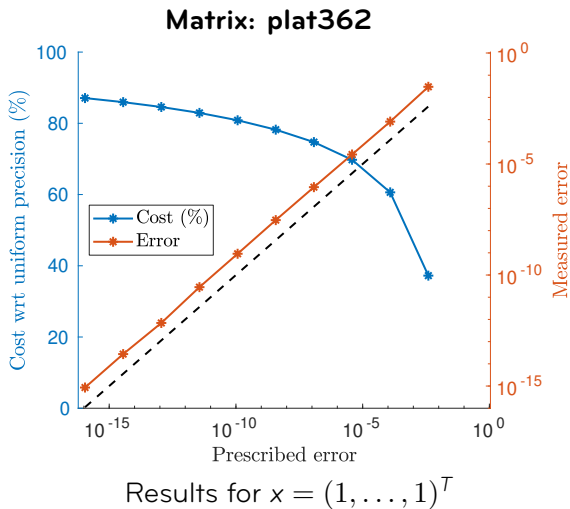


Results for  $x = (1, \dots, 1)^T$

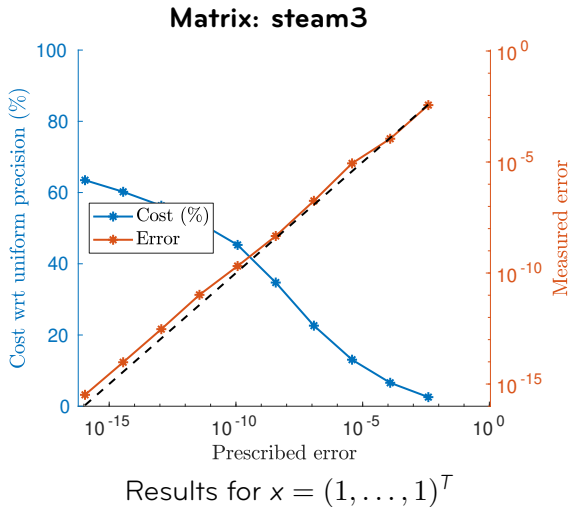


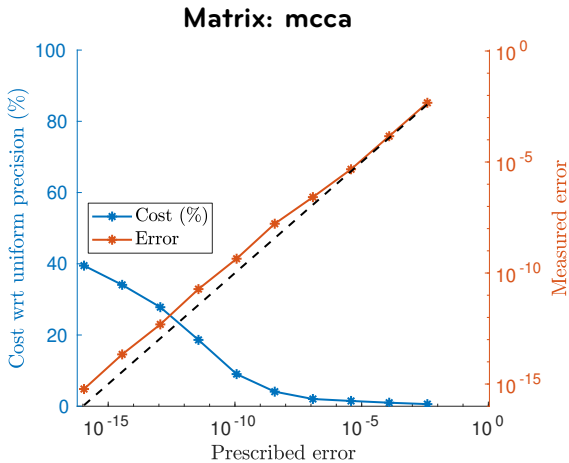
Results for  $x = (1, \dots, 1)^T$











Results for  $x = (1, \dots, 1)^T$

## Role of vector $x$

- Critical issue: accuracy of SpMV depends on  $x$ , but not practical to change precision of  $A$  based on  $x$
- Can still use it and cross fingers ...
- More promising avenue: use it in a setting where  $x$  is guaranteed to be “nice”

## Role of vector $x$

- Critical issue: accuracy of SpMV depends on  $x$ , but not practical to change precision of  $A$  based on  $x$
- Can still use it and cross fingers ...
- More promising avenue: use it in a setting where  $x$  is guaranteed to be "nice"

```
r = b - Ax0
β = ||r||2
q1 = r/β
for k = 1, 2, ... do
  y = Aqk
  for j = 1: k do
    hjk = qjTy
    y = y - hjkqj
  end for
  hk+1,k = ||y||2
  qk+1 = y/hk+1,k
  Solve the least squares problem minck ||Hck - βe1||2
  xk = x0 + Qkck
end for
```

⇒ Krylov solvers! In GMRES,  $x$  is orthonormal

Difficulties in assessing the potential of adaptive precision for GMRES:

- **Highly matrix dependent**, need to cover a wide range of applications
- For a given matrix, **hard to know what a good accuracy is**
  - What storage precision?
  - What tolerance threshold for GMRES convergence?
  - Normwise or componentwise stable SpMV?  
We have both  $\eta_{\text{fwd}} \leq \kappa(A)\eta_{\text{bwd}}^{\text{norm}}$  and  $\eta_{\text{fwd}} \leq \text{cond}(A)\eta_{\text{bwd}}^{\text{cmp}}$ , where  $\text{cond}(A) = \| \|A^{-1}\| \|A\| \leq \kappa(A) = \|A\| \|A^{-1}\|$   
If  $\text{cond}(A) \ll \kappa(A)$ , componentwise stability may be critical
  - How small should the forward error be?
- Comparison further muddled by possible use of
  - Preconditioners
  - Iterative refinement (i.e., restarted GMRES)

- Results with matrix arc130 ( $n = 130$ ,  $\kappa(A) = 1.2 \times 10^{12}$ ,  $\text{cond}(A) = 2.2 \times 10^6$ )
- Use unpreconditioned unrestarted GMRES with **A stored in precision  $\varepsilon$**  and with **convergence tolerance  $\tau$**
- Compare three variants
  - U: store  $A$  in uniform precision  $\varepsilon$
  - AC: store  $A$  in adaptive precision with  $\eta_{\text{bwd}}^{\text{cmp}} = \varepsilon$  (compwise stability)
  - AN: store  $A$  in adaptive precision with  $\eta_{\text{bwd}}^{\text{norm}} = \varepsilon$  (normwise stability)

# Application to GMRES: an example

$\varepsilon$	$\tau$	Iter.	Cost (% U)		$\eta_{\text{bwd}}^{\text{norm}}$		$\eta_{\text{bwd}}^{\text{cmp}}$		$\eta_{\text{fwd}}$	
			AC	AN	U/AC/AN	U/AC/AN	U/AC/AN	U/AC/AN		
$2^{-53}$	$10^{-14}$	15	57	37	$10^{-16}$		$10^{-10}$		$10^{-5}$	
	$10^{-12}$	13	57	37	$10^{-13}$		$10^{-7}$		$10^{-3}$	
	$10^{-10}$	10	57	37	$10^{-11}$		$10^{-5}$		$10^{-1}$	
					U/AC	AN	U/AC	AN	U/AC	AN
$2^{-37}$	$10^{-14}$	15	50	28	$10^{-16}$	$10^{-11}$	$10^{-9}$	$10^{-5}$	$10^{-5}$	$10^{-1}$
	$10^{-12}$	13	50	28	$10^{-13}$	$10^{-11}$	$10^{-7}$	$10^{-5}$	$10^{-3}$	$10^{-1}$
	$10^{-10}$	10	50	28	$10^{-11}$	$10^{-11}$	$10^{-5}$	$10^{-5}$	$10^{-1}$	$10^{-1}$
$2^{-24}$	$10^{-14}$	15	43	20	$10^{-12}$	$10^{-7}$	$10^{-7}$	$10^{-1}$	$10^{-2}$	$10^2$
	$10^{-12}$	13	43	20	$10^{-12}$	$10^{-7}$	$10^{-7}$	$10^{-1}$	$10^{-2}$	$10^2$
	$10^{-10}$	10	43	20	$10^{-11}$	$10^{-7}$	$10^{-5}$	$10^{-1}$	$10^{-1}$	$10^2$

- Use of AC/AN does not increase iterations
- Both AC/AN achieve significant gains, AN more so
- As  $\varepsilon$  increases, U/AC become more accurate than AN
- As  $\tau$  increases, gap between U/AC and AN closes

⇒ **No clear winner**

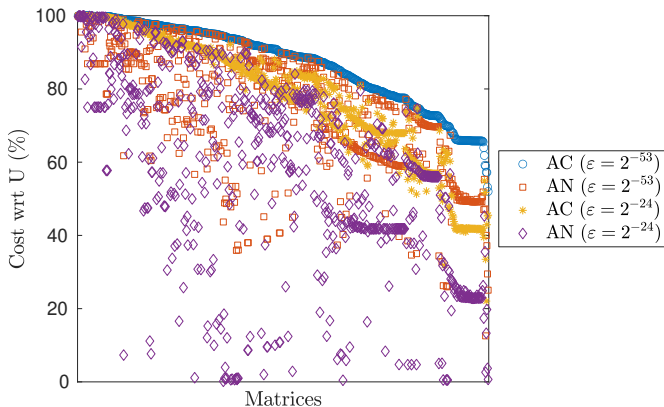
# Application to GMRES: an example

$\varepsilon$	$\tau$	Iter.	Cost (% U)		$\eta_{\text{bwd}}^{\text{norm}}$		$\eta_{\text{bwd}}^{\text{cmp}}$		$\eta_{\text{fwd}}$	
			AC	AN	U/AC/AN		U/AC/AN		U/AC/AN	
$2^{-53}$	$10^{-14}$	15	57	37	$10^{-16}$		$10^{-10}$		$10^{-5}$	
	$10^{-12}$	13	57	37	$10^{-13}$		$10^{-7}$		$10^{-3}$	
	$10^{-10}$	10	57	37	$10^{-11}$		$10^{-5}$		$10^{-1}$	
					U/AC	AN	U/AC	AN	U/AC	AN
$2^{-37}$	$10^{-14}$	15	50	28	$10^{-16}$	$10^{-11}$	$10^{-9}$	$10^{-5}$	$10^{-5}$	$10^{-1}$
	$10^{-12}$	13	50	28	$10^{-13}$	$10^{-11}$	$10^{-7}$	$10^{-5}$	$10^{-3}$	$10^{-1}$
	$10^{-10}$	10	50	28	$10^{-11}$	$10^{-11}$	$10^{-5}$	$10^{-5}$	$10^{-1}$	$10^{-1}$
$2^{-24}$	$10^{-14}$	15	43	20	$10^{-12}$	$10^{-7}$	$10^{-7}$	$10^{-1}$	$10^{-2}$	$10^2$
	$10^{-12}$	13	43	20	$10^{-12}$	$10^{-7}$	$10^{-7}$	$10^{-1}$	$10^{-2}$	$10^2$
	$10^{-10}$	10	43	20	$10^{-11}$	$10^{-7}$	$10^{-5}$	$10^{-1}$	$10^{-1}$	$10^2$

- If target  $\eta_{\text{fwd}} = 10^{-5}$ :
  - AN costs  $53 \times 15 \times 37$ , AC costs  $37 \times 15 \times 50$
  - ⇒ Similar cost, both about 50% better than U
- If target  $\eta_{\text{fwd}} = 10^{-1}$ :
  - AN costs  $37 \times 10 \times 28$ , AC costs  $24 \times 10 \times 43$
  - ⇒ Similar cost, both about 57% better than U



# Application to GMRES: results on real-life matrices



- Results on 659 SuiteSparse matrices of order  $n \in [100, 5000]$
- Keep only matrices for which AC ( $\epsilon = 2^{-53}$ ) achieves some gain (70% of the matrices)

# Conclusion: take-home messages

- The AAA approach:
  - Tailor the use of mixed precision to a specific algorithm class
  - **Use error analysis as a guide** to choose the precisions
  - **Adapt precisions to the data** at hand
- Illustration of this methodology for NLA
  - **Adaptive precision SpMV**
  - Application to Krylov solvers: **significant reductions of the data movement at equivalent accuracy**
  - Article in preparation

**Thank you! Any questions?**