

Responsibly Reckless Matrix Algorithms for HPC Scientific Applications

Hatem Ltaief, Marc G Genton, Damien Gratadour, David E Keyes, Matteo

Ravasi

► To cite this version:

Hatem Ltaief, Marc G Genton, Damien Gratadour, David E Keyes, Matteo Ravasi. Responsibly Reckless Matrix Algorithms for HPC Scientific Applications. Computing in Science and Engineering, 2022, 24 (4), pp.12 - 22. 10.1109/mcse.2022.3215477 . hal-04286094

HAL Id: hal-04286094 https://hal.science/hal-04286094

Submitted on 15 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Queries to the Author

Q1. Please provide the postal code in the affiliation of the author Damien Gratadour.

Q2. Please provide the year of the refs. [19] and [20].

Q3. Please update Refs. [21] and [24].

Q4. Please provide page range for Refs. [23] and [27]. Q5. Please provide DOIs for all the references.

THEME ARTICLE: CELEBRATING DONGARRA

Responsibly Reckless Matrix Algorithms for HPC Scientific Applications

Hatem Ltaief ¹⁰ and Marc G. Genton ¹⁰, King Abdullah University of Science and Technology, Thuwal, 23955,
 Saudi Arabia

Q1 8 Damien Gratadour, Observatoire de Paris, Paris, France

9 David E. Keyes ^(D) and Matteo Ravasi ^(D), King Abdullah University of Science and Technology, Thuwal, 23955,
 10 Saudi Arabia

11 High-performance computing (HPC) achieved an astonishing three orders of

12 magnitude performance improvement per decade for three decades, thanks to

hardware technology scaling resulting in an exponential improvement in the rate of

14 floating point executions, though slowing in the most recent. Captured in the

15 Top500 list, this hardware evolution cascaded through the software stack,

16 triggering changes at all levels, including the redesign of numerical linear algebra

17 libraries. HPC simulations on massively parallel systems are often driven by matrix

18 computations, whose rate of execution depends on their floating point precision.

19 Referred to by Jack Dongarra, the 2021 ACM A.M. Turing Award Laureate, as

²⁰ "responsibly reckless" matrix algorithms, we highlight the implications of mixed-

21 precision (MP) computations for HPC applications. Introduced 75 years ago, long

²² before the advent of HPC architectures, MP numerical methods turn out to be

23 paramount for increasing the throughput of traditional and artificial intelligence

24 (AI) workloads beyond riding the wave of the hardware alone. Reducing precision

comes at the price of trading away some accuracy for performance (reckless

26 behavior) but in noncritical segments of the workflow (responsible behavior) so that

the accuracy requirements of the application can still be satisfied. They offer a

valuable performance/accuracy knob and, just as they are in AI, they are now

²⁹ indispensable in the pursuit of knowledge and discovery in simulations. In

30 particular, we illustrate the MP impact on three representative HPC applications

31 related to seismic imaging, climate/environment geospatial predictions, and

32 computational astronomy.

iscoveries in computational science and engi-33 neering are often the result of multidisciplinary 34 research that synergistically combines efforts 35 from experts in hardware architecture, numerical librar-36 ies, system software, and domain science. The incen-37 tives for hardware and domain science experts are 38 often orthogonal: extracting the expected performance 39 for the former, and getting high-throughput accurate 40

1521-9615 © 2022 IEEE Digital Object Identifier 10.1109/MCSE.2022.3215477 scientific outcomes for the latter. The developers of 41 software libraries usually function as a bridge between 42 the two communities, which ultimately requires all 43 actors to move outside of their comfort zones. This may 44 translate into taking "shortcuts" to achieve the desired 45 outcomes, but these must still be rigorously verified. 46

Mixed-precision (MP) matrix algorithms are numerical methods that employ low/high precisions for storage 48 and/or computations in noncritical/critical sections of 49 the algorithm, respectively. Introduced 75 years ago,¹ 50 then revisited in the context of solving a system of linear 51 equations² and eigensolvers,³ these MP methods reduce 52 time-to-solution while recovering the lost numerical 53



accuracy via an iterative refinement procedure. These 54 early works on MP algorithms were substantial in provid-55 ing rigorous error analysis and numerical stability studies. 56 It is interesting that the development of MP algorithms 57 occurred well before the emergence of high-perfor-58 mance computing (HPC) systems featuring dedicated 59 hardware units for MP computations. Indeed, HPC his-60 tory shows that hardware technologies are typically 61 deployed first before HPC applications receive enough 62 attention to effectively run on them. Putting the cart 63 before the horse is not always smooth; hence recent 64 efforts to establish a roadmap for software/hardware co-65 design.4 66

Although MP algorithms have existed for decades, 67 their adoption into mainstream applications has gained 68 momentum only lately, driven by the ever-increasing 69 industry-led AI market. Figure 1 shows new floating-70 point (FP) representations designed by vendors specifi-71 cally for AI workloads in addition to the existing IEEE 72 754 formats. The number of FP representations indi-73 cates the opportunity for the HPC community to 74 embrace more flexibility in the software stack. Chip 75 manufacturers have deployed GPU accelerators with 76 special hardware support for fast MP arithmetic, attain-77 ing up to 30X performance speedup compared to FP64 78 arithmetic (see Table 1 in https://images.nvidia.com/ 79 aem-dam/en-zz/Solutions/data-center/nvidia-ampere-80 architecture-whitepaper.pdf⁵). This unprecedented per-81 formance improvement has initiated a forced march 82 toward integrating MP matrix algorithms into traditional 83 HPC scientific applications. This trend has been further 84 accelerated by considerations related to energy con-85 86 sumption due to expensive data movement. Unfortunately, this may sometimes lead to the reckless usage 87 of MP algorithms, without proper numerical validation, 88 especially in situations where multidisciplinary collabo-89 ration may not be fostered. Recent work⁶ provides 90 mathematical insight and error analyses that will 91

eventually engender a wider adoption from domain 92 scientists. 93

Reckless but responsible MP matrix algorithms are 94 what is needed to successfully bring together all actors 95 in the exciting upcoming HPC era in which approximate 96 computing will play a major role for scaling up scientific 97 computing. The Innovative Computing Laboratory at 98 the University of Tennessee, Knoxville, USA, led by Jack 99 Dongarra, is a pioneer research institution in the devel- 100 opment of numerical libraries for massively parallel sys- 101 tems, including MP matrix algorithms,⁷ as implemented 102 in LAPACK,⁸ ScaLAPACK,⁹ DPLASMA,¹⁰ PLASMA/ 103 MAGMA,¹¹ and more recently SLATE.¹² The development 104 of these libraries is more the result of a marathon than a 105 sprint, with contributions from the community, perpet-106 ual algorithmic innovations, and performance portabil- 107 ity across hardware vendors and generations. 108

JACK IN THE XBOX

The gaming industry represents one of the main 110 markets for GPU hardware accelerators. While the 111 main duty of these devices is to render graphics at 112 high resolutions, GPUs have been used from early 113 on as a commodity for performing traditional 114 computational simulations. They were initially provi- 115 sioned with low support for FP64 (in favor of FP32) 116 with error correcting code memory disabled. This 117 did not prevent a team of linear algebra experts, led 118 by Jack Dongarra, to exploit the performance of 119 FP32 arithmetic in obtaining FP64 accuracy. Early 120 investigations started with using Intel's Pentium 121 processors, AMD's Opteron architectures, and the 122 IBM's Cell Broad Engine processor [12], which pow- 123 ered the famous Xbox and PlayStation gaming con- 124 soles. Their MP approach used an iterative 125 refinement procedure to recover the precision loss 126 while solving systems of linear equations, but keep- 127 ing the bulk of the computation in FP32 to maximize 128 the GPU's throughput. This work was among the 129 first to democratize GPUs for traditional HPC work- 130 loads. But the road ahead was still bumpy for 131 leveraging MP algorithms into applications. 132

"ISN'T IT DONE YET?"

As often asked by Jack Dongarra—but it takes a village 134 to raise a child, and it takes the right ecosystem for sus-135 tainable research to develop. For MP computations, it 136 was important for the hardware/software ecosystem to 137 reach a certain level of maturity before MP numerical 138 methods could become mainstream. MP algorithms 139 realize the greatest advantage for their extra complexity 140 when the execution rates widen between successive 141

133



precisions (see Table 1 in https://images.nvidia.com/ 142 aem-dam/en-zz/Solutions/data-center/nvidia-ampere-143 architecture-whitepaper.pdf⁵). Al algorithms revolution-144 ized the way we do simulations, casting most of the 145 underlying operations into low-precision matrix-matrix 146 multiplications (GEMM), using FP8/FP16/BF16 FP repre-147 sentations. Chip manufacturers supported this trend by 148 provisioning dedicated hardware units, which translated 149 into unprecedented performance improvement. This is 150 where research on MP algorithms began to flourish and 151 impact the sustained performance of traditional HPC 152 applications. 153

At the same time, the linear algebra community 154 had to address the computational challenges brought 155 to the fore by the manycore era with massively parallel 156 hardware architectures. A profound redesign of the 157 dense block algorithms available in the legacy LAPACK 158 codes occurred to expose more fine-grained parallel-159 ism and mitigate the artifactual synchronization bar-160 riers. The resulting "tile algorithms" operate on the 161 162 dense matrix split into small tiles using a task-based programming model. These computational tasks and 163 their data dependencies constitute the vertices and 164 edges of a directed acyclic graph, which is used to 165 characterize the dataflow of the algorithm. A dynamic 166 runtime system is then employed to orchestrate the 167 task scheduling and execute them in an asynchronous 168 fashion while ensuring data dependencies are not vio-169 lated. This fine-grained task scheduling permits maxi-170 mizing hardware occupancy and strong scaling in the 171 presence of a large number of resources. As imple-172 mented in PLASMA¹¹/DPLASMA¹⁰/SLATE,¹² these tile 173 algorithms became standard for solving dense linear 174 algebra problems. 175

To integrate MP into tile algorithms, the key idea is to extend the mathematical philosophy of iterative refinement^{1,2,3} and to leverage the fine-grained task scheduling for performance, while reducing the overhead of data movement by transferring fewer bytes. This procedure consists in using low FP arithmetic for the bulk of the computation (i.e., matrix factorization) 182 and refining/correcting the residual/solution using a 183 higher FP arithmetic. A recent study¹³ demonstrates the 184 robustness of the procedure using three IEEE 754 FP16/ 185 FP32/FP64 precisions, even in presence of matrices 186 with high condition numbers. This may make the MP 187 with iterative refinement procedure agnostic to the 188 application⁷ in practice. These mathematical founda-189 tions are needed to design MP algorithms responsibly 190 reckless. The MP with iterative refinement procedure 191 does however require storing the whole matrix in the dif-192 ferent precisions, which incurs a significant cost in 193 terms of memory footprint. 194

NEW MP PERSPECTIVE

A new family of MP tile algorithms has emerged that 196 exploits the data sparsity structure of the matrix, e.g., 197 arising from Schur complements within discretizations 198 of elliptic and parabolic PDEs, radial basis functions 199 from unstructured meshing, and covariances in statis- 200 tics.^{14,15} In particular, the latter class of problems 201 appears when modeling major scientific applications 202 related to seismic imaging, climate/environmental geo- 203 spatial predictions, and computational astronomy. This 204 involves a symmetric matrix that represents correla- 205 tions between data values of the physical phenomenon 206 of interest. With proper ordering, strong correlations are 207 typically located around the matrix diagonal and they 208 start fading out as they are further away from the diago- 209 nal. Intuitively, when weak elements are combined with 210 strong ones, their less significant bits fall off the edge of 211 the result. They can therefore be approximated using 212 lower FP arithmetic. Figure 2 shows a standard tile svm- 213 metric matrix (only the lower part is represented) stored 214 in FP64 double precision. By leveraging the insights 215 from the application, one can take advantage of the gra- 216 dient pattern of the correlations and define band 217 regions for applying corresponding FP arithmetic. This 218 approach is referred to as cautious in Figure 2. This 219

approach is suboptimal in that the band tiles exhibiting 220 strong correlations may be far from the diagonal when 221 solving 3-D problems, hence requiring broad bands. A 222 more adequate approach is to prescribe the FP preci-223 sion on a tile-by-tile basis, referred to as responsible in 224 Figure 2. This tile-centric approach relies on comparing 225 226 the Frobenius norm of a tile against the Frobenius norm of the overall matrix, as explained in Higham and Mary's 227 work.⁶ Depending on the ratio of norms and the required 228 numerical accuracy by the application, a tile-centric 229 decision is made on the FP precision before the matrix 230 operations proceed. To further overcome challenges of 231 extreme scale, tile low-rank (TLR) matrix approximations 232 can be applied in addition to MP techniques. Each tile 233 can be independently compressed up to a prescribed 234 accuracy threshold to reduce the overall memory foot-235 print. In fact, TLR matrix approximations and MP can be 236 237 combined, resulting in a representation of the matrix that is rather challenging to manipulate. We refer to this 238 novel approach as reckless in Figure 2. The "reckless-239 ness" is not from the lack of mathematical foundations, 240 but from the "wild west" of marshaling matrix operations 241 242 on a tile data structure whose tiles can be stored in either dense or low-rank formats in any number of preci-243 sions. This is where tools based on runtime systems 244 like PaRSEC¹⁶ are key. They increase the user-pro-245 ductivity for code development and deal transpar-246 ently with data movement of this reckless MP+TLR 247 algorithms, beyond their traditional role of monitor-248 ing task scheduling. 249

Effectively and safely employing MP+TLR algo-250 rithms requires multidisciplinary expertise, engaging 251 hardware architects, linear algebra algorithm devel-252 253 opers, software engineers, and domain scientists. We present three scientific applications to highlight 254 the matrix operations required to leverage these 255 cautious, responsible, and reckless variants of MP 256 algorithms: 257

- imaging the Earth's subsurface using seismic
 redatuming, which requires fast matrix-vector
 multiplication;
- 261 **2)** modeling climate/environment with geostatis-262 tics; and
- a) understanding the origin of life and the Universe,
 which both employ the Cholesky factorization as
 a preliminary step toward solving a linear system
 of equations and evaluating a matrix determinant.

We thus highlight opportunities and implications for applications that model physical phenomena from the Earth's subsurface, the Earth's atmosphere, and beyond the Earth.

4

IMAGING THE EARTH'S SUBSURFACE

As early as the ancient Greeks, the curiosity of human-273 kind for what lies beneath our feet has led to the devel-274 opment of the field of geophysics, a discipline that 275 studies the Earth's structure and behavior. Like a cam-276 era, reflection seismology allows geophysicists to 277 image the Earth's subsurface at meter-scale resolution. 278

The process of acquiring, processing, and imaging 279 seismic data for large fields can take up to several 280 weeks to months. Moreover, while a seismic acquisi- 281 tion campaign back in the 90s may have created data 282 on the order of a few gigabytes, acquisition surveys 283 are now producing terabytes (or even petabytes) of 284 data. This creates computational challenges that 285 make seismology of great interest to the HPC commu-186 nity. It is no surprise that some of the fastest supercomputers in the Top500 list such as Total's Pangea 288 and ENI's HPC4 supercomputers have been specifi-289 cally designed to tackle such challenges.¹⁷ 290

As far as data storage and manipulation are con- 291 cerned, the oil industry has been very conservative 292 throughout the years. The main storage format, SEG- 293 Y,¹⁸ was developed in the mid 70s, a different geologi- 294 cal era when we look at it through the lens of scientific 295 computing: however, one specification of the SEG-Y 296 format still dictates how seismic data are stored and 297 processed to this day: each recorded seismic ampli- 298 tude is in fact delivered to us in a 32-bit (FP32) format. 299 When it comes to visualization and interpretation, 300 seismic volumes are usually optimized for perfor- 301 mance (e.g., DUG Insight User Manual-Optimizing Vol- 302 umes for Performance¹⁹), meaning that copies of the 303 same data are created in lower precision-usually 16- 304 bit or even 8-bit. While not space-friendly, today's low 305 cost of storing data motivates such an approach to 306 ultimately improve user experience. But for quantita- 307 tive analysis, geophysicists seem to be more conser- 308 vative. Quoting field expert Matt Hall, "...for seismic 309 analysis, 8-bit data is probably not precise enough. 310 Opinions vary, but I usually keep my 32-bit volume on 311 disk, but make all derivative volumes and attribute 312 volumes 16-bit. I think 65,536 values is enough, and 313 because of noise and other uncertainties in the data, 314 any precision beyond that is spurious."²⁰ 315

Recent research in the area of seismic processing 316 suggests that Matt may have been too conservative in 317 his statement. Various seismic processing algorithms 318 have been shown to be robust against storing the fre- 319 quency domain representation of the seismic datasets 320 in low-precision, alongside performing algebraic com- 321 pression in the form of TLR approximations.²¹ This is 322



FIGURE 3. Schematic diagram of a responsibly reckless seismic processing workflow, where the input dataset is preprocessed into a set of low precision U and V bases to enable TLR based matrix-vector multiplications.

exactly the domain where such algorithms operate,²² 323 where a multidimensional convolution modeling oper-324 ator that involves the evaluation of an extremely 325 expensive complex-valued, batched matrix-vector 326 multiplication operation is repeatedly applied. One 327 328 can, therefore, afford to spend a certain amount of time upfront to apply a series of preprocessing steps 329 such that the data is optimally arranged for subse-330 quent efficient computations; Hong et al.'s work^{21,23,24} 331 show that by using TLR compression, matrix rear-332 333 rangement based on Hilbert sorting, and low-precision storage as a way to reduce the sheer size of seismic 334 data could lead to an extraordinary reduction in both 335 computational resources and time. Figure 3 illustrates 336 the complex seismic workflow using TLR-MVM, which 337 eventually calls an MP procedure for further optimiz-338 ing execution and data movement. 339

Using the synthetic seismic dataset created in 340 Ravasi and Vasconcelos's work,²² Hong et al.²⁴ com-341 pared the time-to-solution of TLR-MVM with and with-342 out Hilbert sorting and using different levels of 343 precisions on a variety of architectures. In general, a 344 speedup is observed in all cases compared to the 345 standard dense MVM operation with FP32 precision. 346 In particular, our TLR-MVM outperforms dense MVM 347 by more than 6X on NVIDIA A100 GPU using a variety 348 of precisions in compute and storage (i.e., INT8/FP16/ 349 FP32). Moreover, considering the entire process of 350 redatuming (or virtually moving) seismic data from the 351 Earth's surface²⁵ to an area of about 1.3 km² at a 352 depth of 650 m below a complex overburden, the time 353 breakdown in Figure 4(b) reveals that compression is 354 just a tiny fraction of the overall cost of such a seismic 355 processing step. And by continuing to develop this MP 356 algorithm, the TLR-MVM time profile can be further 357 358 reduced, while minimally impacting the overall quality of the results. 359

Was anything lost in terms of accuracy? Yes, as 360 there is no free lunch in MP computations. Neverthe-361 less, Figure 5 shows that the solution obtained at one 362 subsurface point using a TLR compressed version of 363

the seismic data and then mixing INT8/FP16/FP32 pre- 364 cisions in compute and storage is very similar to that 365 obtained from its dense, FP32 counterpart. The signal- 366 to-noise ratio is only 0.4-dB smaller in the low-preci- 367 sion case, compared to the "ground truth," although 368 the data size has been compressed by a factor of over 369 100X! This reckless but responsible seismic processing 370 workflow indicates that the geophysicist community 371 may have been overcomputing and storing for years. 372 Similar MP opportunities may be lying ahead for the 373 seismic field at large. 374

MODELING CLIMATE/ENVIRONMENT 375 WITH GEOSTATISTICS 376

Geostatistics models and predicts quantities of inter- 377 est from data distributed in space-time based on sta- 378 tistical assumptions. It can be seen as complementary 379



FIGURE 4. Performance results in terms of time-to-solution (a) for a single MVM using dense and TLR compressed matrices with different precisions, and (b) for an entire seismic redatuming workflow using the best choice of TLR compressed matrix kernel.



FIGURE 5. (a) Seismic redatuming wavefield estimate using dense, high-precision (FP32) matrix kernels. (b) Error introduced when switching to TLR compressed bases stored in INT8 precision.

to modeling approaches based on first principles 380 rooted in conservation laws and physics-based mod-381 els commonly expressed by PDEs. Climate and envi-382 ronmental applications, e.g., soil moisture variables 383 384 recorded at the topsoil of the Mississippi River basin in Figure 6, are among the main workloads keeping 385 supercomputers busy worldwide and are intended for 386 exascale computers, thus even minor enhancements 387 for production applications may provide significant 388 rewards. 389



FIGURE 6. Soil moisture residuals at the topsoil of the Mississippi River basin.

A key ingredient of geostatistics is the covariance 390 matrix, often based on an underlying spatial covari- 391 ance function of the Matérn class, which appears in 392 the likelihood function of Gaussian random fields, in 393 the optimal spatial interpolation coined "Kriging" and 394 its uncertainty quantification, and in the simulation of 395 realizations from Gaussian random fields. This dense 396 covariance matrix has a symmetric and positive-defi- 397 nite form, and its algebraic dimension equals the num- 398 ber of spatially distributed data values times the 399 number of time steps. In the aforementioned tasks, 400 two fundamental operations on the covariance matrix 401 are the application of its inverse and the computation 402 of its determinant. These operations can all be 403 obtained through the celebrated Cholesky factoriza- 404 tion and triangular solution, but are characterized by 405 cubic/square complexity in the number of data values 406 in flops/memory, respectively, and hence become 407 unfeasible for large-scale problems. Indeed, a covari- 408 ance matrix for 1M spatial data values would require 409 4-TB memory in (symmetric) DP format and on the 410 order of 10¹⁸ flops to factor. 411

*ExaGeoStat*²⁶ is a software package built to provide 412 user-controlled approximations to extreme-scale geo-413 statistical problems by introducing innovative algorith-414 mic, architectural, and programming model features. A 415 MP tile Cholesky algorithm is introduced to speed up 416 the factorization in the key geostatistical tasks. It is 417 then deployed on large-scale heterogeneous systems 418 with the help of *PaRSEC* dynamic runtime system.¹⁶ 419



FIGURE 7. Gaussian maximum likelihood estimation performance breakdown across GPU generations and precision arithmetic.

With a suitable ordering, the algorithm works with dou-420 421 ble-precision arithmetic on tiles neighboring the main diagonal, while operating with single-precision or lower 422 arithmetic for tiles far enough, leading to a three-preci-423 sion FP16/FP32/FP64 approximation algorithm for the 424 Cholesky factorization.¹⁵ Referred to as cautious algo-425 426 rithms in Figure 2, ExaGeoStat leverages the inherent band data sparsity structure of the covariance matrix 427 to exploit all FP representations accordingly. Figure 7 428 demonstrates the impact of MP tile algorithms across 429 various NVIDIA GPU generations. The dashed curves 430 show the progress that can be obtained by riding the 431 hardware alone, through three generations of NVIDIA 432 GPUs, i.e., Pascal, Volta, and Ampere, respectively. The 433 red arrow shows a fivefold improvement from Pascal to 434 Ampere GPUs. The solid curves show the performance 435 obtained on the latter two architectures supporting a 436 437 mixed DP and HP covariance matrix, with twofold improvements as shown by the green arrows. The com-438 bined improvement coming from hardware and algo-439 rithm, shown by the blue arrow, is tenfold. Numerical 440 accuracy assessment¹⁵ shows that with a proper band 441 structure capturing the regions of strong/medium/ 442 weak correlations, ExaGeoStat is able to compute the 443 relevant statistical parameters as if all computations 444 were performed in FP64. 445

Finally, one can also leverage TLR approximations in 446 addition to MP tile algorithms to further reduce the 447 memory footprint and time-to-solution, leading to a 448 reckless but responsible algorithm that can still ensure 449 adequate accuracy.⁶ With the help of PaRSEC dynamic 450 runtime system,¹⁶ the Cholesky factorization relies on a 451 hybrid data distribution that mitigates the load imbal-452 453 ance between tasks next to and far from the diagonal with high/low algorithmic complexity, respectively. PaR-454 SEC then marshals the data movement of tiles stored 455 in dense and compressed formats. It also performs 456

precision conversion on-the-fly and uses advanced 457 look-ahead techniques to shorten the critical path, 458 beyond its original duty of orchestrating task schedul- 459 ing. This MP+TLR combination permits to take the best 460 of the two worlds: 1) MP tile-dense algorithms applied 461 on computational tasks around the critical path to 462 shorten it and 2) TLR algorithms applied on the remain- 463 ing computational tasks to reduce memory footprint 464 and address big data problems. Preliminary results of 465 the resulting responsibly reckless MP+TLR algorithms 466 show that the Cholesky factorization may achieve an 467 order of magnitude performance higher than if MP is 468 applied alone on the original FP64 dense covariance 469 problems. The numerical validation is even more robust 470 with this algorithmic variant than the cautious version, 471 especially when dealing with 3-D problems. Moving for- 472 ward, this creates new opportunities to study parallel 473 space-time likelihood optimization on large-scale sys- 474 tems,²⁷ which would be otherwise intractable. 475

CHASING THE ORIGIN OF LIFE AND 476 THE UNIVERSE 477

The superior angular resolution provides exciting 478 opportunities for astronomy, making possible major 479 scientific breakthroughs by enabling better photomet- 480 ric and astrometric precision and better contrast. The 481 2020 Nobel Prize in Physics was, for instance, awarded 482 to a group of astronomers "for the discovery of a 483 supermassive compact object at the center of our gal- 484 axy," believed to be a giant black hole.²⁸ Within the 485 observational arsenal, adaptive optics (AO) stood out 486 as a game changer to enable this significant leap in 487 our understanding of the Universe. High angular reso- 488 lution is key to making detailed studies of both our 489 Earth's neighborhood, distant reaches revealing the 490 early Universe, and everything in between. Moreover, 491 when coupled to high contrast techniques, getting 492 sharper images allows astronomers to study exopla- 493 nets in extrasolar planetary systems, over a range of 494 evolutionary stages in order to probe the initial condi- 495 tions for planetary formation, the evolution of plane- 496 tary systems over various time-frames and possibly 497 the emergence of life. 498

While the largest ground-based telescopes will 499 soon reach 40-m diameter and provide the angular resolution and collecting area required to detect the first 501 stars and first galaxies as well as faint rocky exoplanets 502 through direct imaging, they must be equipped with 503 appropriate apparatus to overcome optical distortions 504 induced by atmospheric turbulence. AO technologies, 505 dating back to the late 1980s, were developed for this 506 purpose and are now essential for the largest optical 507



FIGURE 8. AO loop is composed of the DM, the WFS, and the RTC. A typical WFS image is shown on the upper right panel and images of a star in open and close loop operation models are shown below that.

telescopes. In its simplest form, an AO system is com-508 posed of a wavefront sensor (WFS) used to measure 509 atmospheric distortions at a high frame rate, which are 510 compensated with a deformable mirror (DM). The sub-511 system linking those components, responsible for 512 interpreting wavefront measurements into actual com-513 mands to actuators of the DMs, is the real-time con-514 troller (RTC), as shown in Figure 8. It must operate at 515 high speed (kHz rate) to catch up with the rapidly 516 changing optical turbulence. 517

One of the limitations of classical AO is that the 518 correction is only valid in a very small patch of sky, the 519 size of which depends on the observing wavelength, 520 521 from a few arcseconds in the visible to a few tens of arcseconds in the near infrared. Multiconjugate adap-522 tive optics (MCAO) solves this problem by using a 523 series of DMs to compensate the turbulence in vol-524 ume, enabling AO correction over a wide field of view. 525 MCAO uses several guide stars and associated WFSs 526 to probe the light wave aberrations in several direc-527 tions, and the RTC, using tomographic reconstruction, 528 determines the best commands to apply to the DMs.²⁹ 529 The classical approaches to aberrations retrieval 530 are based on linear models of the relationship between 531 the sensor(s) data and the phase of the light wave, rely-532 ing on careful modeling of the AO system error budget 533 and solved through a linear approach, sometimes regu-534 535 larized with a given prior on the turbulence statistics. For instance, in present-day conventional AO systems, 536 537 the RTC follows a well-defined linear control scheme: input measurements from sensors are multiplied by 538 a control matrix to produce an output DM control 539 command. 540



FIGURE 9. Proportion of the different GEMM precision and associated ToR accuracy on an 8-m telescope with 17K measurements (left) and a 40-m telescope with 50K measurements (right).

The computation of the tomographic reconstruc- 541 tor (ToR), which consists of the Cholesky factorization 542 of the dense symmetric covariance matrix of WFS 543 measurements followed by a backward and forward 544 substitution, is at the core of operations for all tomo- 545 graphic AO instruments and must be updated regu- 546 larly to take into account the evolution of the 547 atmosphere's light-bending structure. The most time- 548 consuming kernel during the factorization and the 549 solve phase is the general matrix-matrix multiplication 550 (GEMM), which makes the compute-bound algorithm 551 run close to the system's sustained peak performance. 552 Performance results have been reported on shared- 553 memory systems equipped with hardware accelera- 554 tors³⁰ as well as distributed-memory systems³¹ using 555 single precision FP arithmetic. The measurements 556 used to generate the matrix operator come from 16- 557 bit unsigned integer signals from WFS cameras, which 558 has to be converted to 32-bit in order to perform FP 559 computations of the covariance matrix. In addition, 560 the dense covariance matrix, which may be of size as 561 large as 100K, has a data-sparse structure, due to 562 weak interactions between some of the measure- 563 ments taken by the WFS. This is expected, since meas- 564 urements taken by WFS subapertures physically 565 located next to each other exhibit higher correlations, 566 while as we move away from the matrix diagonal, 567 weak interactions between measurements taken by 568 remote WFS subapertures are expected. 569

The ToR computation has been tested for two different ranges of systems dimensions, including a state-ofthe-art 8-m telescope instrument and a future 40-m 572 telescope instrument up to a total of 17K and 50K measurements, respectively. Figure 9 shows the proportion 574 of tiles operated on as single precision or the different 575 variants of half-precision, as well as the performance of 576 the computed MP ToR. This corresponds to the *cau-* 577 *tious* MP variant of Figure 2 based on a band structure 578 for determining the FP32 or FP16 arithmetic. The graph 579 entries are sorted given the type and proportion of half-580 precision, and the ToR performance is expressed as 581



FIGURE 10. ToR performance breakdown across GPU generations and precision arithmetic.

Strehl ratio (SR): a measure of image quality as the ratio 582 of the maximum value in the image over its theoretical 583 maximum, 1 being the best achievable ratio. These SRs 584 585 are obtained with the end-to-end simulation tool COM-PASS³² generating long exposure images from a full 586 model of the system, including turbulence, telescope, 587 and AO instrument. SR values are missing when the 588 Cholesky factorization did not succeed (due to the loss 589 of positiveness engendered by the low precision), but in 590 a successful case, the ToR accuracy is almost equal to 591 the full single precision approach, regardless of the con-592 sidered instrument dimensioning. 593

Figure 10 tracks the same type of improvements 594 across hardware generations and algorithmic improve-595 ments as in Figure 7, this time between SP and HP. The 596 MP ToR computation on the newest Ampere GPU 597 scores a threefold speedup compared to the reference 598 SP ToR implementation on the previous Pascal GPU 599 generation. The sustained performance achieved by 600 601 using FP16/FP32 hardware support from the latest NVI-DIA A100 GPU enables the computational astronomy 602 603 community of ground-based telescopes to improve science insights during nightly on-sky observations. This 604 result has implications for all instruments requiring 605 606 real-time processing needed in AO.

Moreover, by using a tile-centric responsible MP 607 variant, the ToR computation can capture the needed 608 accuracy per tile⁶ and ensure numerical robustness 609 without requiring a priori knowledge on the data spar-610 sity structure of the covariance matrix. Combining MP 611 612 with TLR matrix approximations (i.e., reckless but responsible MP+TLR algorithms) is an interesting ave-613 nue to further satisfy real-time requirements and will 614 615 be investigated in future work.

CONCLUSION

MP matrix algorithms have evolved significantly since 617 their initial introduction 75 years ago. Recent work 618

provides new analyses to take into account multiple 619 lower precisions of FP arithmetic, which emerged with 620 the advent of hardware technologies that support Al. 621 Further challenges await with new non-IEEE FP repre- 622 sentations (e.g., BF16 and FP8) fostered by chip manu- 623 facturers. Rigorous numerical validations are critical 624 to ensure these reckless algorithms remain responsi- 625 ble. The new tile-centric MP approach for extreme- 626 scale applications offers further opportunities to 627 economize on storage and execution time while meet- 628 ing user-specified accuracy requirements, which were 629 often exceeded by default double precision. However, 630 multidisciplinary expertise is required to navigate the 631 software stack: the more reckless the MP algorithm, 632 the more responsible the users must be. 633

Dense linear algebra is just one of the "seven 634 dwarfs" identified by Phil Colella in a famous 2004 pre- 635 sentation³³ titled "Defining Software Requirements for 636 Scientific Computing." It is time to review the amena- 637 bility of the other six to exploiting the new precisions 638 available in hardware in order to further stretch the 639 capacity of HPC systems for extreme applications. 640

ACKNOWLEDGMENTS

The authors would like to thank Fujitsu/NVIDIA/NEC 642 for the remote access to their respective systems and 643 Intel/AMD for their hardware donations. 644

REFERENCES

1. L. Fox, H. D. Huskey, and J. H. Wilkinson, "Notes on the 646 **O5** solution of algebraic linear simultaneous equations," 647 Quart. J. Mechanics Appl. Math., vol. 1, no. 1, pp. 149–173, 648 1948, doi: 10.1093/qjmam/1.1.149. 649 2. C. B. Moler, "Iterative refinement in floating point," J. 650 ACM, vol. 14, no. 2, pp. 316-321, Apr. 1967, doi: 10.1145/ 651 321386.321394. 652 3. J. J. Dongarra, C. B. Moler, and J. H. Wilkinson, 653 "Improving the accuracy of computed eigenvalues and 654 eigenvectors," SIAM J. Numer. Anal., vol. 20, no. 1, 655 pp. 23-45, 1983, doi: 10.1137/0720002. 656 4. J. Dongarra et al., "The International Exascale Software 657 Project roadmap," Int. J. High Perform. Comput. Appl., 658 vol. 25, no. 1, pp. 3-60, Feb. 2011, doi: 10.1177/1094342 659 010391989. 660 5. "NVIDIA A100 Tensor Core GPU Architecture V1.0." 661 2020. [Online]. Available: https://images.nvidia.com/ 662 aem-dam/en-zz/Solutions/data-center/nvidia-ampere-663 architecture-whitepaper.pdf 664 O2 N. J. Higham and T. Mary, "Mixed precision algorithms 665 in numerical linear algebra," Acta Numerica, vol. 31,

616

pp. 347-414, 2022, doi: 10.1017/S0962492922000022.

666

667

641

668	7.	A. Abdelfattah et al., "A survey of numerical linear	21.	Y. Hong, H. Ltaief, M. Ravasi, D. Keyes, and D. Vargas,	721
669		algebra methods utilizing mixed-precision arithmetic,"		"Large-scale marchenko imaging with distance-aware	722
670		Int. J. High Perform. Comput. Appl., vol. 35, no. 4,		matrix reordering, tile low-rank compression, and	723
671		pp. 344–369, 2021, doi: 10.1177/10943420211003313.		mixed-precision computations," SEG Tech. Abstract,	724
672	8.	E. Anderson et al., LAPACK Users' Guide, vol. 9.		2022, accepted.	725
673		Philadelphia, PA, USA: SIAM, 1999.	22.	M. Ravasi and I. Vasconcelos, "An open-source	726
674	9.	L. S. Blackford et al., ScaLAPACK Users' Guide.		framework for the implementation of large-scale	727
675		Philadelphia, PA, USA: SIAM, 1997.		integral operators with flexible, modern hpc solutions-	728
676	10.	G. Bosilca et al., "Flexible development of dense linear		enabling 3d marchenko imaging by least-squares	729
677		algebra algorithms on massively parallel architectures		inversion "Geophysics vol 86 pp. WC177–WC194 2021	730
678		with DPI ASMA" in Proc. IEEE Int. Symp. Parallel		doi: 10.1190/geo2020-0796.1	731
679		Distrib Process Workshops Ph D Forum 2011	23	Y Hong H I taief M Rayasi L Gatineau and D Keyes	732
680		nn 1432–1441	20.	"Accelerating seismic redatuming using tile low-rank	732
681	11	F Aguillo et al "Numerical linear algebra on emerging		approximations on NEC SX-aurora TSUBASA "	734
682		architectures: The PLASMA and MAGMA projects " 1		Supercomputing Front Innovations vol 8 2021	734
662		Bhue Conf Ser vol 180 pp 1 2000 Art pp 012027		doi: 10.14520/iof210201	735
683	10	Phys. Conf. Ser., vol. 180, 110. 1, 2009, Alt. 110. 012037.	24	V Hang H Ltaisf M Daviesi and D Keyes "HDC	736
684	12.	M. Gales, J. Kurzak, A. Charara, A. Yarkhan, and J.	24.	Y. Hong, H. Llaiei, W. Ravasi, and D. Reyes, HPC	737
685		Dongarra, SLATE: Design of a modern distributed and		seismic redatuming by inversion with algebraic	738
686		accelerated linear algebra library," in Proc. Int. Conf.		compression and multiple precisions," ACM Trans.	739
687		High Perform. Comput., Netw., Storage Anal., 2019,		Math. Softw., 2022, in preparation.	740
688		pp. 1–18, doi: 10.1145/3295500.3356223.	25.	K. Wapenaar, J. Thorbecke, J. van der Neut, F. Broggini,	, 741
689	13.	E. Carson and N. J. Higham, "Accelerating the solution		E. Slob, and R. Snieder, "Marchenko imaging,"	742
690		of linear systems by iterative refinement in three		Geophysics, vol. 79, pp. WA39–WA57, 2014,	743
691		precisions," SIAM J. Sci. Comput., vol. 40, no. 2,		doi: 10.1190/geo2013-0302.1.	744
692		pp. A817–A847, 2018, doi: 10.1137/17M1140819.	26.	S. Abdulah, H. Ltaief, Y. Sun, M. G. Genton, and D. E. Keyes,	, 745
693	14.	S. Abdulah, H. Ltaief, Y. Sun, M. G. Genton, and D. E.		"ExaGeoStat: A high performance unified software for	746
694		Keyes, "Geostatistical modeling and prediction using		geostatistics on manycore systems," IEEE Trans. Parallel	747
695		mixed precision tile cholesky factorization," in Proc.		Distrib. Syst., vol. 29, no. 12, pp. 2771–2784, Dec. 2018.	748
696		IEEE 26th Int. Conf. High Perform. Comput., Data,	27.	M. L. Salvana, S. Abdulah, H. Ltaief, Y. Sun, M. Genton,	749
697		Analytics, 2019, pp. 152–162.		and D. Keyes, "Parallel space-time likelihood	750
698	15.	S. Abdulah et al., "Accelerating geostatistical modeling		optimization for air pollution prediction on large-scale	751
699		and prediction with mixed-precision computations: A		systems," in Proc. Platform Adv. Sci. Comput. Conf.,	752
700		high-productivity approach with PaRSEC," IEEE Trans.		2022.	753
701		Parallel Distrib. Syst., vol. 33, no. 4, pp. 964–976, Apr.	28.	"Press release: The nobel prize in physics 2020,"	754
702		2022.		nobelprize.org, 2020. [Online]. Available: https://www.	755
703	16.	G. Bosilca, A. Bouteiller, A. Danalis, M. Faverge, T.		nobelprize.org/prizes/physics/2020/press-release/	756
704		Hérault, and J. J. Dongarra, "PaRSEC: Exploiting	29.	F. Rigaut and B. Neichel, "Multiconjugate adaptive	757
705		heterogeneity to enhance scalability," Comput. Sci.		optics for astronomy," Annu. Rev. Astron. Astrophys.,	758
706		Eng., vol. 15, no. 6, pp. 36–45, 2013.		vol. 56, no. 1, pp. 277–314, 2018, doi: 10.1146/annurev-	759
707	17.	B. Felix, "Oil group total hopes new supercomputer will		astro-091916-055320.	760
708		help it find oil faster and more cheaply, 2019. [Online].	30.	H. Ltaief, D. Gratadour, A. Charara, and E. Gendron,	761
709		Available: https://www.reuters.com/article/us-total-		"Adaptive optics simulation for the world's largest	762
710		supercomputer-idUSKCN1TJ0FQ		telescope on multicore architectures with multiple	763
711	18.	K. Barry, D. Cavers, and C. Kneale, "Recommended		GPUs," in Proc. Platform Adv. Sci. Comput. Conf., 2016.	764
712		standards for digital tape formats." Geophysics, vol. 40.		pp. 1–12. doi: 10.1145/2929908.2929920.	765
713		no. 2. pp. 344–352. 1974.	31.	H. I taief et al. "Real-time massively distributed multi-	766
714	19	DLIG Insight Liser Manual—Ontimising Volumes for	0.1	object adaptive ontics simulations for the european	767
715	10.	Performance [Online] Available: https://help.dugeo		extremely large telescope " in Proc. IEEE Int. Parallel	768
716		com/m/Insight/I/438665-ontimising-volumes-for-		Distrib Process Symp May 2018 np 75–84	760
717		nerformance	32	F Ferreira D Gratadour A Sevin and N Doucet	770
719	20	M Hall "B is for Bit denth" [Online] Available: https://	52.	"Compase: An efficient GPI Leased simulation	770
710	20.	agilescientific com/blog/2011/2/4/b.is.for.bit.dorth		compass. An encient of orbased simulation	770
719				Conf High Darform Comput Simul 2019 pp 190 197	772
720		num		Conj. mgn Ferjonn. Comput. Sinui., 2010, pp. 180–187.	113

736 03

746 Q4

774	33.	P. Colella, "Defining software requirements for
775		scientific computing. Slide of 2004 presentation
776		included in David Patterson's 2005 talk," 2004. [Online].
777		Available: http://www.lanl.gov/orgs/hpc/salishan/
778		salishan2005/davidpatterson.pdf

HATEM LTAIEF is the principal research scientist of the
Extreme Computing Research Center, King Abdullah University
of Science and Technology, Thuwal, 23955, Saudi Arabia. His
research interests include parallel numerical algorithms, parallel
programming models, and performance optimizations for multicore architectures and hardware accelerators. He is a Member
of IEEE. Contact him at Hatem.Ltaief@kaust.edu.sa.

MARC G. GENTON is a distinguished professor of statistics 786 with King Abdullah University of Science and Technology, 787 Thuwal, 23955, Saudi Arabia. His research interests include 788 statistical analysis, flexible modeling, prediction, and uncer-789 tainty quantification of spatio-temporal data, with applica-790 tions in environmental and climate science, renewable 791 energies, geophysics, and marine science. He is a fellow of 792 the ASA, of the IMS, and the AAAS, and is an elected member 793 of the ISI. Contact him at Marc.Genton@kaust.edu.sa. 794

DAMIEN GRATADOUR is an associate professor with Labora toire d'Etudes Spatiales et d'Instrumentation en Astrophysique,

Observatoire de Paris, Paris, France. His research interests 797 include bridges astronomy with high-performance computing 798 and artificial intelligence, applied to modeling, signal processing 799 and instrumentation for large telescopes. Contact him at 800 damien.gratadour@obspm.fr. 801

DAVID E. KEYES is the director of the Extreme Computing 802 Research Center with King Abdullah University of Science and 803 Technology, Thuwal, 23955, Saudi Arabia. He works with the 804 interface between parallel computing and the numerical anal-905 ysis of PDEs with a focus on scalable implicit solvers, such as 806 the Newton–Krylov–Schwarz and the Additive Schwarz Pre-807 conditioned Inexact Newton methods, which he co-devel-808 oped. He is a fellow of the SIAM, AMS, and AAAS. He is also a 809 Member of IEEE. Contact him at David.Keyes@kaust.edu.sa. 810

MATTEO RAVASI is an assistant professor of geophysics with 811 King Abdullah University of Science and Technology, Thuwal, 812 23955, Saudi Arabia. His research interests include geophysical 813 inverse problems with applications to seismic acquisition and 814 processing, imaging, quantitative interpretation, and time-lapse 815 monitoring. He is also interested in the use of machine learning 816 and high-performance computing and heavily involved in the 817 development of open-source scientific software. Contact him 818 at Matteo.Ravasi@kaust.edu.sa. 819