

35 exponential integrators can be derived by optimizing the coefficients of an “ansatz” of the form

$$(0.2) \quad \varphi_0(A)\mathbf{v}_0 + \varphi_1(A)\mathbf{v}_1 + \varphi_2(A)\mathbf{v}_2 + \dots + \varphi_p(A)\mathbf{v}_p,$$

where the so-called φ -functions can be defined by the Taylor series

$$\varphi_k(A) = \sum_{i=0}^{\infty} \frac{A^i}{(i+k)!}.$$

36 In other words, the ansatz eq. (0.2) is just a linear combination of exponential-like functions
 37 evaluated at A that act on a set of vectors. Computationally, evaluating $\varphi_i(A)\mathbf{v}_i$ is the most
 38 expensive step in exponential time integrators, and it has therefore been the subject of a
 39 considerable amount research.

40 **0. The workshop.** In September 2023, members of the $f(A)\mathbf{b}$ community met at the
 41 Max Planck Institute for Dynamics of Complex Technical Systems in Magdeburg, Germany,
 42 to attend the first “ $f(A)$ bulous workshop on matrix functions and exponential integrators”,¹
 43 organized by Kathryn Lund, Stéphane Gaudreault, and Marcel Schweitzer. The event featured
 44 traditional-style scientific talks covering recent algorithmic advances as well as applications,
 45 but a significant portion of the two and a half days was reserved for moderated discussion
 46 sessions.

47 For these sessions, the workshop participants split into three focus groups (each led by
 48 one of the organizers of the workshop) in order to consider a broad challenge the community
 49 is facing, and they spent one afternoon looking at the issue and assembling potential solutions.
 50 For each group, a list of key questions was provided to foster and stimulate a lively—but also
 51 focused—discussion. For the guiding slides with questions, see [43].

52 The next day, each group reported the main discussion points to the other attendees,
 53 and the floor was then opened for further comments and questions. Several participants (in
 54 particular, Thomas Mach and Yannis Voet) took thorough notes on the discussions and shared
 55 them with other participants via the workshop website.

56 The three challenges that were selected for these focus groups were as follows:

- 57 1. knowledge transfer between the “ $f(A)\mathbf{b}$ community” and researchers from other
 58 areas (moderated by Marcel Schweitzer),
- 59 2. high-performance and energy-aware computing (moderated by Stéphane Gaudreault),
 60 and
- 61 3. benchmark problems and FAIR comparisons (moderated by Kathryn Lund).

62 The next three sections of this report summarize the main conclusions of the three focus
 63 groups and the main points that were raised in the discussion that ensued.

64 **1. Knowledge transfer.** The first group looked at ways to ensure that the algorithms
 65 developed within the community can reach the end users who need them.

66 As discussed before, items **T1** and **T2** are fundamentally different from a computational
 67 perspective, and they generally require different techniques. In particular, when only $f(A)\mathbf{b}$
 68 is sought, computing $f(A)$ is typically an unduly expensive and totally unnecessary step.

69 Anecdotal evidence suggests that the distinction between the two problems is not as
 70 clear outside the $f(A)\mathbf{b}$ community, and that researchers that wish to compute $f(A)\mathbf{b}$ in an
 71 application domain often rely on the simple (but extremely inefficient) approach of multiplying
 72 \mathbf{b} by $f(A)$ after having computed the latter explicitly.

73 There are several reasons behind this phenomenon. For one, there is no recent, authorita-
 74 tive source summarizing the existing methods for evaluating $f(A)\mathbf{b}$. The only comprehensive

¹<https://indico3.mpi-magdeburg.mpg.de/event/30/>

75 survey of the literature [40] is now over 16 years old and does not reflect the breadth of
 76 methods currently available; more recent surveys and a thesis [46, 47, 69] only deal with
 77 Krylov subspace methods and focus especially on limited-memory scenarios. The situation
 78 is similar for exponential integrators, where the last comprehensive survey [56] is 14 years
 79 old. This is in stark contrast to the literature on $f(A)$, which boasts a book [51],² a survey
 80 paper dedicated to computational aspects [52], and even a survey of existing software, which
 81 is periodically updated (and also includes software for $f(A)\mathbf{b}$) [28, 53, 54].

82 An additional obstacle in knowledge transfer is the absence of $f(A)\mathbf{b}$ in the standard
 83 academic curriculum. In effect, most practitioners learn about the topic either during their grad-
 84 uate studies or through self-study, and their learning is hindered by the lack of a comprehensive
 85 review or textbook.

86 For most functions f of interest, it is relatively easy to find robust and efficient software
 87 toolboxes to evaluate $f(A)$. These are available for most programming environments and
 88 work out-of-the-box for a large selection of test problems. The situation is very different in
 89 the $f(A)\mathbf{b}$ case, as the performance of the algorithm depends on a number of factors, and,
 90 with the software currently available, a certain level of experience is needed to find the right
 91 combination of parameters for a given computation.

92 In most cases, the missing cornerstone is a reliable stopping criterion. A stopping criterion
 93 requires a way of measuring or bounding the approximation error, but at present there is a
 94 knowledge gap in this area: available results typically only apply to normal matrices and are
 95 restricted to specific classes of matrices (e.g., Kronecker sums [10, 11]) or specific functions
 96 (e.g., Cauchy–Stieltjes functions [35, 39, 48, 49, 62], Laplace transforms [31, 37] or functions
 97 which can be related to an underlying ODE initial-value problem [13, 14, 15, 16]). The
 98 error bounds available for more general cases may be very pessimistic, and they are typically
 99 not fit for use as stopping criteria in practice. Furthermore, the derivation of error bounds
 100 quickly becomes interdisciplinary, as often function-specific, analytical results are necessary
 101 for deriving error expressions. This is in stark contrast to, say, linear systems of the form
 102 $A\mathbf{x} = \mathbf{b}$, where a notion of residual $\mathbf{r} := \mathbf{b} - A\tilde{\mathbf{x}}$ for an approximation $\tilde{\mathbf{x}}$ is readily available
 103 from the data, can be cheaply approximated in Krylov subspace methods like GMRES, and is
 104 widely used as a reliable stopping criterion [68].

105 The situation is further complicated by the lack of comprehensive comparisons of the
 106 performance of different algorithms, which is made arduous by the large number of imple-
 107 mentation details and parameters that have to be selected. We address this point further in
 108 Section 3.

109 This focus group also identified a few timely challenges for the community, which
 110 represent, at the same time, opportunities for research advances.

- 111 • Current methods may not exploit the full breadth of available techniques, a case in
 112 point being randomized methods, which have just started to be considered [18, 27,
 113 50, 66, 67].
- 114 • Existing implementations are not able to fully leverage the variety of hardware on
 115 modern computers and supercomputers, which makes them potentially less attractive.
 116 See also Section 2.
- 117 • Often, unlocking the full potential of methods might require “intermingling” the
 118 algorithm used for approximating $f(A)\mathbf{b}$ with the surrounding ecosystem, taking into
 119 account specifics of the application at hand, instead of just treating it as a black-box

²Incidentally, in the preface of [51], the author states “The problem of computing a function of a matrix times a vector, $f(A)\mathbf{b}$, is of growing importance, though as yet numerical methods are relatively undeveloped”. Due to this growing importance, a lot of developments have taken place since then, which are of course not covered by the last “ $f(A)\mathbf{b}$ survey” [40], which was published in the same year as the book [51].

120 that returns an approximate solution (see, e.g., [26] for an example of a “Krylov-
 121 aware” approach in trace estimation, or [44] for an algorithm exploiting the intimate
 122 connection to exponential integrators when approximating linear combinations of
 123 φ -functions).

124 *Next steps.* Although a clear path to solve all the challenges affecting knowledge transfer
 125 is hard to pin down, some steps the community can take to make some progress on these issues
 126 are relatively easy to trace.

127 First and foremost, there is a need for effective benchmarking standards. All methods to be
 128 compared should be implemented from the same building blocks, and the same implementation
 129 choices should be applied consistently. The performance of the methods should be measured
 130 in a uniform way, and this is not necessarily a simple task: estimating the execution time and
 131 memory usage of an algorithm is simple, but assessing its accuracy is not. Accuracy is usually
 132 measured in terms of the forward error of the computed result, but since the exact solution is
 133 often not available, a reference solution computed using a different algorithm—potentially run
 134 in higher-than-working precision—is typically employed. Attention should be paid to how the
 135 reference solution is computed, and how the error is estimated from it.

136 Ideally, one could identify classes of problems where certain methods perform well, so
 137 that precise and easy-to-follow recommendations—based on the structure of A , the behavior
 138 of f , or a combination of both—can be made.

139 In order for the comparison to be useful, it is also crucial that a representative set of
 140 benchmark problems be established. Difficult questions that should be addressed regard:

- 141 • what information should be collected, in addition to the obvious f , A , and \mathbf{b} , and
- 142 • what format should be used to store this information.

143 These points were discussed in more detail by another focus group—see Section 3.

144 These efforts would put the community in a better position to summarize the literature and
 145 produce easy-to-follow guidance for all those interested in computing $f(A)\mathbf{b}$ without delving
 146 into algorithmic and theoretical details: a general consensus is that drafting a modern survey of
 147 numerical methods for computing $f(A)\mathbf{b}$ should be a priority.³ The lack of a comprehensive
 148 literature review for more specific problems, such as exponential integrators, is also a shared
 149 concern which should be addressed in coming years.

150 Finally, a question that was raised is whether understanding better the sensitivity of
 151 the proposed algorithms, as well as offering ways of estimating the conditioning of a given
 152 problem, could help practitioners feel more confident about the use of a chosen algorithm for
 153 $f(A)\mathbf{b}$ —after all, this is one of the aspects that make the BLAS and LAPACK stand out.

154 There does indeed exist a lot of work on estimating the condition number of the com-
 155 putation of $f(A)$ and $f(A)\mathbf{b}$; see, e.g., [51, Chapter 3] for a general overview of the topic.
 156 However, these condition number estimates are intimately related to the *Fréchet derivative*
 157 $L_f(A, \cdot)$ of $f(A)$, an object that is typically several times more costly to compute than $f(A)$
 158 itself. For the $f(A)$ problem, there exist several algorithms (each for a specific function f) that
 159 allow for the computation of $f(A)$ and its Fréchet derivative simultaneously and reuse certain
 160 computations in the process [3, 4]. Building on this, [3, Algorithm 7.4] computes the matrix
 161 exponential e^A together with a (quite reliable) condition number estimate at a cost of roughly
 162 17 times that of computing e^A alone.⁴ Thus, even when using cleverly designed algorithms
 163 hand-tailored to a specific function, the overhead induced by the condition number estimator
 164 is quite substantial.

165 Additionally, it is currently unclear how to extend such approaches to, e.g., Krylov
 166 subspace algorithms for $f(A)\mathbf{b}$, as computing $f(A)\mathbf{b}$ and $L_f(A, \cdot)$ has much less in common

³Indeed, efforts in this direction are already underway.

⁴The factor 17 can be reduced to 9 if a slightly reduced reliability is acceptable.

167 than computing $f(A)$ and $L_f(A, \cdot)$. Recently, there has been some progress in Krylov subspace
 168 algorithms for low-rank approximations of the Fréchet derivative [57, 58, 61], which might
 169 facilitate making a first step in this direction.

170 **2. High-performance and energy-aware computing.** The second focus group looked at
 171 the challenges surrounding the applications of $f(A)\mathbf{b}$ in high-performance computing (HPC).
 172 In many domains within the natural and social sciences and engineering disciplines, there is
 173 a need to compute $f(A)\mathbf{b}$ where A is extremely large and sparse. Such large problems are
 174 typically solved using supercomputers, which are machines composed of many nodes with
 175 distributed memory, sometimes employing heterogeneous computing hardware. Each node is
 176 typically equipped with a number of CPUs and accelerators, such as GPUs, and to achieve
 177 peak performance, a routine must make the best use of all available resources.

178 Aside from those for $A^{-1}\mathbf{b}$, numerical methods to compute $f(A)\mathbf{b}$ have seldom been
 179 used in large-scale parallel applications. One of the most active fields of research in this area
 180 is the solution of differential equations using exponential time integrators with Krylov and
 181 Leja point methods. Various factors impede the application of certain algorithms developed by
 182 the $f(A)\mathbf{b}$ community, and we will provide a brief overview in this section.

183 When the A matrix is large and sparse, it is often impossible to store it explicitly in memory.
 184 Fortunately, in many cases one can use “matrix-free” algorithms, which converge without the
 185 cost of forming or storing the matrix. These are frequently used in HPC applications because
 186 they allow the solution of problems that would otherwise be intractable. In the context of
 187 $f(A)\mathbf{b}$, most matrix-free algorithms require only the action of the matrix (or an approximation
 188 to it) in the form of matrix–vector products. For example, instead of storing the sparse Jacobian
 189 J of a vector-valued function $F(\mathbf{u})$, its action on a vector can be approximated using the finite
 190 difference $J\mathbf{b} \approx [F(\mathbf{u} + \epsilon) - F(\mathbf{u})] / \epsilon$, where ϵ is a small perturbation [17]. Other matrix-
 191 free approaches, such as the complex-step approximation [70] or automatic differentiation
 192 [45], are often also used.

193 Requirements in terms of parallelism and memory storage considerably restrict the choice
 194 of possible algorithms. The finite difference approximation of the Jacobian action illustrated
 195 above, for example, does not allow operations such as transposition, slicing, or pivoting
 196 without a prohibitive computational cost. For the most part, it is not problematic to use a
 197 matrix–vector product routine instead of matrices for methods based on the Krylov subspace,
 198 the Taylor series, or Leja points. However, difficulties arise when information about the norm
 199 or the spectrum of A is needed to 1) compute the parameters that make these methods efficient,
 200 or 2) determine the stopping criterion. Without the matrix representation, it can be expensive
 201 to compute the operator norm or to estimate eigenvalues. While the spectral radius $\rho(A)$ can
 202 be cheaply approximated using the power method on a single CPU, the communication cost
 203 in parallel implementations renders this idea inefficient in many HPC applications. Further
 204 research will be necessary to develop numerical algorithms more suitable to this kind of
 205 problems.

206 Another important consideration is the capability of an algorithm to scale and optimize
 207 energy efficiency as the amount of computing resources increases. Clearly, existing imple-
 208 mentations are not yet ready for the future exascale machines, i.e., supercomputers capable
 209 of performing at least 10^{18} binary64⁵ floating-point operations per second. The problem of
 210 implementing Krylov subspace methods efficiently on a GPU has been considered from a
 211 theoretical point of view [34], but studies focusing on high-performance implementations
 212 suggest that the use of GPUs is most beneficial when A is dense [5], which is not often the
 213 case for $f(A)\mathbf{b}$ problems, or when A has a very specific sparsity pattern that can be mapped

⁵Previously known as “double precision”.

214 efficiently to GPU architectures [33]. Therefore, software for this problem primarily targets
 215 CPUs and can only rely on GPUs in a limited number of cases or for a subset of the relevant
 216 operations.

217 One obstacle is the fact that most implementations target binary64 accuracy, but binary64
 218 arithmetic is not very efficient on GPUs. When using the tensor cores on the latest NVIDIA
 219 H100 SXM5 GPUs [64, Table 1], for example, the theoretical peak performance of the 19-bit
 220 TensorFloat-32 arithmetic is 494.7 trillion floating-point operations per second (TFLOPS),
 221 which improves for BFLOAT16 and binary16 arithmetic (989.4 TFLOPS) and breaks the
 222 PFLOPS barrier for the fp8 formats (1978.9 TFLOPS). The peak performance of binary64
 223 arithmetic is almost 30 times slower, with just 66.0 TFLOPS when tensor cores are used for
 224 matrix–matrix multiplications. Using low-precision arithmetic is key to harnessing the full
 225 potential of GPUs, but only binary64 accuracy is typically sufficient for a range of applications.
 226 In other areas of numerical linear algebra, this challenge has been addressed effectively by
 227 developing mixed-precision algorithms [1, 55]. Efforts have been made recently in the context
 228 of exponential integrators to design such schemes [9], but further research will be necessary.

229 The issue is not solely with implementations: the algorithms themselves do not seem
 230 ready to address large-scale problems either. Many methods rely on matrix operations that
 231 do not scale well in a distributed environment. For example, full-basis orthogonalization,
 232 central to many Krylov subspace methods, necessitates numerous communication operations
 233 (i.e., message passing and synchronization) throughout the computation. This can result in
 234 unacceptable latencies, with most processes idly waiting for the slowest process to complete [8].
 235 There are a number of new developments on this front, including but not limited to low-
 236 synchronization orthogonalization [12, 21, 22, 23, 24, 41, 42, 65, 71, 76, 78] and s -step
 237 methods [19, 20, 75, 76], which can also be combined with one another. Furthermore, a number
 238 of well established techniques are being rediscovered as communication-reducing, such as the
 239 natural short-term recurrences of Lanczos or batching vectors into tall-skinny matrices (block
 240 vectors) to take better advantage of BLAS Level 3 [29, 30]. Sketching, randomization, and low
 241 precision can also be leveraged to reduce memory movement by shrinking the size of vectors
 242 to be stored and manipulated [6, 7, 27, 50, 66, 77]. The performance of these techniques
 243 has been and is being thoroughly explored for linear systems solvers, but their transfer to
 244 matrix functions requires a better understanding of their backward stability, how they can be
 245 integrated into extended and rational Krylov subspace methods, as well as the conditioning of
 246 $f(A)\mathbf{b}$ itself (cf. Section 1).

247 All these issues remain true for the computation of the full matrix $f(A)$ as well.

248 A significant push in this direction could arise from an increased interest in exponential
 249 integrators. However, there are many factors that hinder their use in practical applications.
 250 Firstly, they are seldom featured in textbooks and are often absent from university curricula,
 251 resulting in many practitioners being unfamiliar with them. In addition, despite their better
 252 stability properties, they are generally more complicated to implement than explicit methods.
 253 Even when stability is important, practitioners tend to be more attracted by implicit or implicit-
 254 explicit schemes, because of the availability of techniques that deal with the stiffness of their
 255 particular problems. Furthermore, highly optimized libraries that implement algorithms for
 256 solving linear or nonlinear problems are readily available. This is not the case for exponential
 257 integrators, and the necessity to implement a parallel solver for the φ -functions often discour-
 258 ages the use of these methods in HPC applications. In recent years, there has been a renewed
 259 interest in exponential integrators, and this can largely be attributed to advances in numerical
 260 algorithms for the computation of $f(A)\mathbf{b}$. While this is an encouraging trend, more work is
 261 needed to make these time integration schemes easier to use in applications.

262 *Next steps.* Readying current $f(A)b$ work for exascale presents a number of significant
 263 challenges, but the community is well equipped to make some progress towards this goal. It is
 264 clear that the focus should be on two distinct fronts, since not only the implementations, but
 265 also the algorithms, will require significant work in order to leverage the full computational
 266 power of next-generation supercomputers.

267 In terms of rethinking existing algorithms, there is a clear need for reducing the number
 268 and frequency of communication operations. In particular, parallel inner products are a known
 269 communication bottleneck on distributed systems. Numerical algorithms with high arithmetic
 270 intensity should be favored over those requiring a high degree of data movements. Some
 271 work has already been done in this area (see, for example, block methods [38], truncated
 272 orthogonalization [50, 60, 66] or restarts [2, 16, 32, 36]), but a completely different solution
 273 may be needed.

274 In terms of implementations, a significant challenge is to leverage the untapped potential
 275 of GPUs, which, as they become faster and more prevalent, represent an increasingly large
 276 share of the overall performance of a supercomputer.

277 Writing high-performance numerical linear algebra code that can target GPUs presents
 278 various difficulties. First and foremost, the variation in capabilities between different models
 279 of GPUs, especially those from different vendors, is dramatically larger than the variation
 280 between CPUs. As a consequence, there is no unified implementation of the BLAS and
 281 LAPACK for GPUs. Vendors provide highly-optimized libraries for their own hardware, but
 282 these require very different frameworks, which means that porting an implementation from one
 283 GPU to another requires a significant human effort. For example, NVIDIA provides cuBLAS,⁶
 284 which is part of the CUDA Toolkit,⁷ while AMD provides support through rocBLAS,⁸ which
 285 is part of the ROCm Platform.⁹

286 Potential solutions, which include the C++ runtime API HIP,¹⁰ also part of the ROCm
 287 Platform, the programming model SYCL¹¹ [59], and libraries such as MAGMA¹² [73, 74],
 288 are not yet mature enough to be used in production code.

289 At present, the community should attempt to rewrite existing algorithms to ensure optimal
 290 performance on HPC architectures. They should seek to minimize communications and use
 291 low precision (binary32, binary16, or lower) for the bulk of the computation, switching to
 292 higher precision (typically binary64) only when strictly necessary.

293 For research reproducibility, the $f(A)b$ community should adopt and promote open
 294 science best practices. This entails authors sharing the code and data that would allow to
 295 replicate the results presented in their publications.

296 **3. Benchmarking.** The last focus group discussed best practices for sharing code and
 297 data sets so that they can be easily reused, in accordance with FAIR guidelines.¹³ The main
 298 goal is to simplify two important steps of the algorithm development process:

- 299 • evaluating new implementations on established test problems, and
- 300 • comparing their performance with that of existing algorithms in the literature.

301 A welcome side effect, which comes at no additional cost, is the reproducibility of experimental
 302 results. This idea promises to solve a number of problems that commonly arise when new
 303 algorithms are proposed in the literature.

⁶<https://docs.nvidia.com/cuda/cublas/>

⁷<https://developer.nvidia.com/cuda-toolkit/>

⁸<https://rocm.docs.amd.com/projects/rocBLAS/>

⁹<https://rocm.docs.amd.com>

¹⁰<https://rocm.docs.amd.com/projects/HIP/>

¹¹<https://www.khronos.org/sycl/>

¹²<https://icl.utk.edu/magma/>

¹³<https://www.go-fair.org/>

304 Unless the authors decide to compare their proposed new method with all state-of-the-art
 305 algorithms for the same problem, it is impossible for the reader to understand how the new
 306 method compares with existing alternatives. A new implementation could easily perform
 307 worse than a much simpler and well established one, but the reader would have to spend a
 308 significant amount of effort to check whether this is the case, especially if the code used in the
 309 original publication is not available.

310 The peer review process can help with this, but there are limitations. Reviewers can
 311 recommend that new approaches be compared with the most relevant existing alternatives, but
 312 it is difficult to ensure that the comparison is fair, and most journals in numerical analysis and
 313 numerical linear algebra do not yet require submission of software or reproducibility of the
 314 experimental results. Moreover, as a test set of representative $f(A)\mathbf{b}$ problems is not currently
 315 available, it is difficult for a reviewer—and for the reader, later on—to make sure that the
 316 numerical experiments reported in a publication provide an impartial representation of the
 317 merits and drawbacks of new algorithms. Not all methods are suitable for a given choice of
 318 f , A , and \mathbf{b} , and having a battery of tests with clear classes of functions and matrices can
 319 help identify what types of problems a certain algorithm can deal with effectively. This can
 320 help corroborate theoretical results, in addition to providing a quick and standardized way of
 321 comparing all relevant algorithms for a specific choice of f , A , and \mathbf{b} .

322 The metrics against which these algorithms should be compared are also not uniquely
 323 determined, and authors are free, within reason, to choose the ones that suit them best. In
 324 some cases, the metric itself is poorly defined and can depend on a range of factors that are
 325 not within the control of who is performing the test. A case in point is runtime, which is
 326 commonly used to assess the performance of different implementations on a same test set.
 327 Runtime is very sensitive to the hardware configuration, as well as some low-level details of
 328 the software libraries being used, so that algorithm α_1 can easily be faster than algorithm α_2
 329 on a machine and slower on another for the same test problem.

330 When an algorithm cannot be implemented in the most efficient way possible, for example,
 331 because of limitations of existing hardware, an appealing alternative is to rely on the number of
 332 floating-point operations being performed. This metric is only meaningful for large matrices,
 333 and it can be very inaccurate on modern hardware and especially in distributed-computing
 334 settings, as it focuses on arithmetic intensity when, in practice, the performance of most
 335 algorithms is bounded by the memory bandwidth.

336 A final difficulty is represented by the lack of clear licensing for code and test problems
 337 alike. This prevents reuse and, in many cases, hinders reproducibility of existing results. For
 338 more information on research data management in mathematics, see a recent white paper by
 339 the Germany-based Mathematical Research Data Initiative (MaRDI) [72].

340 *Next steps.* It is a priority for the community to produce a set of representative test cases
 341 whereon new and old algorithms can be compared. Ideally, one would want access to a remote
 342 facility that is capable of testing submitted implementations against known and unknown
 343 benchmark problems, providing overall scores for a number of metrics including accuracy,
 344 stability, and runtime performance. Similar services exist for machine learning research [25],
 345 where the unknown problems are used to prevent authors from overfitting their models to the
 346 test set.

347 The main difficulty to address in order to deliver this golden standard is to ensure a
 348 fair comparison among implementations written using different languages, as Julia, MAT-
 349 LAB/GNU Octave, and Python are all well established in this community, and being able to
 350 compare code in these different languages is likely to pose significant challenges in terms of
 351 software engineering.

352 A more modest but attainable result would be the development of a curated reference
 353 collection of test problems. Authors testing their code could then simply choose which parts
 354 of the collection to include, and they could justify their choices by pointing out which classes
 355 of problems are not suitable in their context. Reviewers could equally rely on such a collection
 356 to ensure that authors are providing a fair picture of the merits of their algorithms.

357 Building and maintaining this infrastructure would come with some logistical challenges.
 358 A sufficient number of examples should be included, so that the collection represents the main
 359 applications in which $f(A)\mathbf{b}$ appears. As test examples can be quite large, the collection might
 360 require a hosting service with sufficient storage space and bandwidth, or standardized protocol
 361 to point to resources like Zenodo, from which data could be downloaded. A possible remedy
 362 for the latter issue is to promote the use of so-called procedural examples, whereby a problem
 363 is specified mathematically and the matrix A and vector \mathbf{b} can be generated with a desired size
 364 or other properties via a script.

365 It is necessary to ensure that the test cases remain relevant, and that the collection grows
 366 and remains representative despite hardware and algorithm improvements that may make
 367 problems that are difficult today trivial in the near future. The test cases will likely come from
 368 a number of researchers in various research domains, and will have to be collected and added
 369 to the collection by a number of volunteers. A standard license—or set of licenses—should be
 370 adopted that ensure reproducibility and that fair credit is given to test problem creators and
 371 curators.

372 Although it will not be possible for the community to enforce such a requirement, pub-
 373 lishing and advertising a reasonable set of recommendations should be one of the priorities of
 374 the group working on this.

375 **4. Conclusions.** It is easy to feel overwhelmed looking at the long to-do lists we have
 376 outlined. A change of perspective may lessen the anxiety: these are exciting opportunities,
 377 some of them even so-called “low-hanging fruit”, and the impact of addressing them is huge,
 378 even for such a small field. Matrix functions continue to surface in diverse applications, and
 379 many of the techniques developed for $f(A)\mathbf{b}$ can cross-pollinate work in linear systems, matrix
 380 equations, Fréchet derivatives, and other problems we are not yet aware of. Furthermore, the
 381 development of comprehensive surveys and language-agnostic benchmarking workflows for
 382 $f(A)\mathbf{b}$ can set an example for other mathematical fields that are struggling to modernize and
 383 keep up with an ever-increasing publication load. Our primary aim is that this manuscript
 384 builds on the momentum of a successful workshop and inspires new, meaningful projects in
 385 $f(A)\mathbf{b}$ and beyond.

386 **Acknowledgments.** The workshop itself was funded in part by DFG Project Number
 387 529315380. We thank all of the workshop participants (excluding the present authors) for
 388 their contributions to the substance of this manuscript: Francesca Arrigo, Michele Benzi, Kai
 389 Bergermann, Philipp Birken, Liam Burke, Marco Caliari, Benjamin Carrel, Fabio Cassini,
 390 Ranjan Kumar Das, Vladimir Druskin, Andreas Frommer, Oswald Knoth, Patrick Kürschner,
 391 Thomas Mach, David Persson, Helmut Podhaisky, Michele Rinelli, Jonas Schulze, Roger
 392 Sidje, Igor Simunec, Martin Stoll, Mayya Tokman, Manuel Tsolakis, Paul Van Dooren, and
 393 Yannis Voet.

REFERENCES

- 394
 395 [1] A. ABDELFAH, H. ANZT, E. G. BOMAN, E. CARSON, T. COJEAN, J. DONGARRA, A. FOX, M. GATES,
 396 N. J. HIGHAM, X. S. LI, J. LOE, P. LUSZCZEK, S. PRANESH, S. RAJAMANICKAM, T. RIBIZEL,
 397 B. F. SMITH, K. SWIRYDOWICZ, S. THOMAS, S. TOMOV, Y. M. TSAI, AND U. M. YANG, *A survey*

- 398 *of numerical linear algebra methods utilizing mixed-precision arithmetic*, Int. J. High Performance
 399 Computing Applications, 35 (2021), pp. 344–369.
- 400 [2] M. AFANASJEW, M. EIERMANN, O. G. ERNST, AND S. GÜTTEL, *Implementation of a restarted Krylov*
 401 *subspace method for the evaluation of matrix functions*, Linear Algebra Appl., 429 (2008), pp. 2293–2314.
- 402 [3] A. H. AL-MOHY AND N. J. HIGHAM, *Computing the Fréchet derivative of the matrix exponential, with an*
 403 *application to condition number estimation*, SIAM J. Matrix Anal. Appl., 30 (2009), pp. 1639–1657.
- 404 [4] A. H. AL-MOHY, N. J. HIGHAM, AND S. D. RELTON, *Computing the Fréchet derivative of the matrix*
 405 *logarithm and estimating the condition number*, SIAM J. Sci. Comput., 35 (2013), pp. C394–C410.
- 406 [5] N. AUER, L. EINKEMMER, P. KANDOLF, AND A. OSTERMANN, *Magnus integrators on multicore CPUs and*
 407 *GPUs*, Comput. Phys. Comm., 228 (2018), p. 115–122.
- 408 [6] O. BALABANOV AND L. GRIGORI, *Randomized Gram–Schmidt Process with Application to GMRES*, SIAM J.
 409 Sci. Comput., 44 (2022), pp. A1450–A1474.
- 410 [7] ———, *Randomized block Gram–Schmidt process for solution of linear systems and eigenvalue problems*,
 411 e-print 2111.14641, arXiv, 2023.
- 412 [8] G. BALLARD, E. C. CARSON, J. W. DEMMEL, M. HOEMMEN, N. KNIGHT, AND O. SCHWARTZ, *Com-*
 413 *munication lower bounds and optimal algorithms for numerical linear algebra*, Acta Numer., 23 (2014),
 414 pp. 1–155.
- 415 [9] C. J. BALOS, S. ROBERTS, AND D. J. GARDNER, *Leveraging mixed precision in exponential time integration*
 416 *methods*, in 2023 IEEE High Performance Extreme Computing Conference (HPEC), 2023, pp. 1–8.
- 417 [10] M. BENZI AND V. SIMONCINI, *Decay bounds for functions of Hermitian matrices with banded or Kronecker*
 418 *structure*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 1263–1282.
- 419 [11] ———, *Approximation of functions of large matrices with Kronecker structure*, Numer. Math., 135 (2016),
 420 pp. 1–26.
- 421 [12] D. BIELICH, J. LANGOU, S. THOMAS, K. ŚWIRYDOWICZ, I. YAMAZAKI, AND E. G. BOMAN, *Low-synch*
 422 *Gram–Schmidt with delayed reorthogonalization for Krylov solvers*, Parallel Computing, 112 (2022),
 423 p. 102940.
- 424 [13] M. A. BOTCHEV, V. GRIMM, AND M. HOCHBRUCK, *Residual, restarting, and Richardson iteration for the*
 425 *matrix exponential*, SIAM J. Sci. Comput., 35 (2013), pp. A1376–A1397.
- 426 [14] M. A. BOTCHEV, L. KNIZHNERMAN, AND M. SCHWEITZER, *Krylov subspace residual and restarting for*
 427 *certain second order differential equations*, SIAM J. Sci. Comput., (2023), pp. S223–S253.
- 428 [15] M. A. BOTCHEV, L. KNIZHNERMAN, AND E. E. TYRTYSHNIKOV, *Residual and restarting in Krylov*
 429 *subspace evaluation of the φ function*, SIAM J. Sci. Comput., 43 (2021), pp. A3733–A3759.
- 430 [16] M. A. BOTCHEV AND L. A. KNIZHNERMAN, *ART: Adaptive residual-time restarting for Krylov subspace*
 431 *matrix exponential evaluations*, J. Comput. Appl. Math., 364 (2020), p. 112311.
- 432 [17] P. N. BROWN, H. F. WALKER, R. WASYK, AND C. S. WOODWARD, *On using approximate finite differences*
 433 *in matrix-free Newton–Krylov methods*, SIAM Journal on Numerical Analysis, 46 (2008), pp. 1892–1911.
- 434 [18] L. BURKE AND S. GÜTTEL, *Krylov subspace recycling with randomized sketching for matrix functions*,
 435 arxiv:2308.02290 [math.NA], Aug. 2023.
- 436 [19] E. CARSON, T. GERGELITS, AND I. YAMAZAKI, *Mixed precision s-step Lanczos and conjugate gradient*
 437 *algorithms*, Numer. Linear Algebra Appl., 29 (2022), p. e2425.
- 438 [20] E. C. CARSON, *An adaptive s-step conjugate gradient algorithm with dynamic basis updating*, Appl. Math.,
 439 65 (2020), pp. 123–151.
- 440 [21] E. C. CARSON, K. LUND, Y. MA, AND E. OKTAY, *On the loss of orthogonality of low-synchronization*
 441 *variants of reorthogonalized block Gram–Schmidt*, tech. rep., In preparation, 2024.
- 442 [22] E. C. CARSON, K. LUND, AND E. OKTAY, *Reorthogonalized Pythagorean variants of block classical Gram*
 443 *Schmidt*, tech. rep., In preparation, 2024.
- 444 [23] E. C. CARSON, K. LUND, AND M. ROZLOŽNÍK, *The stability of block variants of classical Gram–Schmidt*,
 445 SIAM J. Matrix Anal. Appl., 42 (2021), pp. 1365–1380.
- 446 [24] E. C. CARSON, K. LUND, M. ROZLOŽNÍK, AND S. THOMAS, *Block Gram–Schmidt algorithms and their*
 447 *stability properties*, Linear Algebra Appl., 638 (2022), pp. 150–195.
- 448 [25] R. CHAN, K. LIS, S. UHLEMAYER, H. BLUM, S. HONARI, R. SIEGWART, P. FUA, M. SALZMANN, AND
 449 M. ROTTMANN, *SegmentMelfYouCan: A benchmark for anomaly segmentation*, in Proceedings of the
 450 Neural Information Processing Systems Track on Datasets and Benchmarks, J. Vanschoren and S. Yeung,
 451 eds., vol. 1, 2021, p. 13.
- 452 [26] T. CHEN AND E. HALLMAN, *Krylov-aware stochastic trace estimation*, SIAM J. Matrix Anal. Appl., 44
 453 (2023), pp. 1218–1244.
- 454 [27] A. CORTINOVI, D. KRESSNER, AND Y. NAKATSUKASA, *Speeding up Krylov subspace methods for*
 455 *computing $f(a)b$ via randomization*, arXiv:2212.12758 [math.NA], Dec. 2022. Revised June 2023.
- 456 [28] E. DEADMAN AND N. J. HIGHAM, *Testing matrix function algorithms using identities*, ACM Trans. Math.
 457 Software, 42 (2016), pp. 4:1–4:15.
- 458 [29] N.-A. DREIER AND C. ENGWER, *Strategies for the Vectorized Block Conjugate Gradients method*, in Numer.
 459 Math. Adv. Appl. ENUMATH 2019, F. J. Vermolen and C. Vuik, eds., vol. 139 of Lect. Notes Comput.

- 460 Sci. Eng., Springer, Cham, 2020, pp. 381–388.
- 461 [30] N.-A. DREIER AND C. ENGWER, *A Hardware-aware and Stable Orthogonalization Framework*, e-print
462 2204.13393, arXiv, 2022.
- 463 [31] V. DRUSKIN, *On monotonicity of the Lanczos approximation to the matrix exponential*, Linear Algebra and its
464 Applications, 429 (2008), pp. 1679–1683.
- 465 [32] M. EIERMANN AND O. G. ERNST, *A restarted Krylov subspace method for the evaluation of matrix functions*,
466 SIAM J. Numer. Anal., 44 (2006), pp. 2481–2504.
- 467 [33] L. EINKEMMER AND A. OSTERMANN, *Exponential integrators on graphic processing units*, in Proceedings
468 of the 2013 International Conference on High Performance Computing & Simulation (HPCS), Institute of
469 Electrical and Electronics Engineers, July 2013.
- 470 [34] M. E. FARQUHAR, T. J. MORONEY, Q. YANG, AND I. W. TURNER, *GPU accelerated algorithms for
471 computing matrix function vector products with applications to exponential integrators and fractional
472 diffusion*, SIAM J. Sci. Comput., 38 (2016), p. C127–C149.
- 473 [35] A. FROMMER, S. GÜTTEL, AND M. SCHWEITZER, *Convergence of restarted Krylov subspace methods for
474 Stieltjes functions of matrices*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 1602–1624.
- 475 [36] A. FROMMER, S. GÜTTEL, AND M. SCHWEITZER, *Efficient and stable Arnoldi restarts for matrix functions
476 based on quadrature*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 661–683.
- 477 [37] A. FROMMER, K. KAHL, M. SCHWEITZER, AND M. TSOLAKIS, *Krylov subspace restarting for matrix
478 Laplace transforms*, SIAM J. Matrix Anal. Appl., 44 (2023), pp. 693–717.
- 479 [38] A. FROMMER, K. LUND, AND D. B. SZYLD, *Block Krylov subspace methods for functions of matrices*,
480 Electron. Trans. Numer. Anal., 47 (2017), pp. 100–126.
- 481 [39] A. FROMMER AND M. SCHWEITZER, *Error bounds and estimates for Krylov subspace approximations of
482 Stieltjes matrix functions*, BIT Numerical Mathematics, 56 (2015), pp. 865–892.
- 483 [40] A. FROMMER AND V. SIMONCINI, *Matrix functions*, in Model Order Reduction: Theory, Research Aspects
484 and Applications, vol. 13 of Math. Ind., Springer-Verlag, Aug. 2008, pp. 275–303.
- 485 [41] T. FUKAYA, R. KANNAN, Y. NAKATSUKASA, Y. YAMAMOTO, AND Y. YANAGISAWA, *Shifted Cholesky
486 QR for computing the QR factorization of ill-conditioned matrices*, SIAM J. Sci. Comput., 42 (2020),
487 pp. A477–A503.
- 488 [42] T. FUKAYA, Y. NAKATSUKASA, Y. YANAGISAWA, AND Y. YAMAMOTO, *CholeskyQR2: A Simple and
489 Communication-Avoiding Algorithm for Computing a Tall-Skinny QR Factorization on a Large-Scale
490 Parallel System*, in 2014 5th Workshop Latest Adv. Scalable Algorithms Large-Scale Syst., 2014, pp. 31–
491 38.
- 492 [43] S. GAUDREULT, K. LUND, AND M. SCHWEITZER, *Introduction to focus groups topics*, 2023.
493 [https://indico3.mpi-magdeburg.mpg.de/event/30/attachments/182/265/
494 focus_group_slides.pdf](https://indico3.mpi-magdeburg.mpg.de/event/30/attachments/182/265/focus_group_slides.pdf).
- 495 [44] S. GAUDREULT, G. RAINWATER, AND M. TOKMAN, *KIOPS: A fast adaptive Krylov subspace solver for
496 exponential integrators*, J. Comput. Phys., 372 (2018), pp. 236–255.
- 497 [45] A. GRIEWANK AND A. WALTHER, *Evaluating derivatives: principles and techniques of algorithmic differen-
498 tiation*, SIAM, 2008.
- 499 [46] S. GÜTTEL, *Rational Krylov approximation of matrix functions: Numerical methods and optimal pole selection*,
500 GAMM Mitteilungen, 36 (2013), pp. 8–31.
- 501 [47] S. GÜTTEL, D. KRESSNER, AND K. LUND, *Limited-memory polynomial methods for large-scale matrix
502 functions*, GAMM Mitteilungen, 43 (2020), p. e202000019.
- 503 [48] S. GÜTTEL AND L. KNIZHNERMAN, *A black-box rational Arnoldi variant for Cauchy–Stieltjes matrix
504 functions*, BIT, 53 (2013), pp. 595–616.
- 505 [49] S. GÜTTEL AND M. SCHWEITZER, *A comparison of limited-memory Krylov methods for Stieltjes functions of
506 Hermitian matrices*, SIAM J. Matrix Anal. Appl., 42 (2021), p. 83–107.
- 507 [50] S. GÜTTEL AND M. SCHWEITZER, *Randomized sketching for Krylov approximations of large-scale matrix
508 functions*, SIAM J. Matrix Anal. Appl., 44 (2023), pp. 1073–1095.
- 509 [51] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, Society for Industrial and Applied Mathe-
510 matics, Philadelphia, PA, USA, 2008.
- 511 [52] N. J. HIGHAM AND A. H. AL-MOHY, *Computing matrix functions*, Acta Numerica, 19 (2010), pp. 159–208.
- 512 [53] N. J. HIGHAM AND E. DEADMAN, *A catalogue of software for matrix functions. Version 1.0*, MIMS EPrint
513 2014.8, Manchester Institute for Mathematical Sciences, The University of Manchester, UK, Feb. 2014.
- 514 [54] N. J. HIGHAM AND E. HOPKINS, *A catalogue of software for matrix functions. Version 3.0*, MIMS EPrint
515 2020.7, Manchester Institute for Mathematical Sciences, The University of Manchester, UK, Mar. 2020.
- 516 [55] N. J. HIGHAM AND T. MARY, *Mixed precision algorithms in numerical linear algebra*, Acta Numerica, 31
517 (2022), pp. 347–414.
- 518 [56] M. HOCHBRUCK AND A. OSTERMANN, *Exponential integrators*, Acta Numerica, 19 (2010), pp. 209–286.
- 519 [57] P. KANDOLF, A. KOSKELA, S. D. RELTON, AND M. SCHWEITZER, *Computing low-rank approximations of
520 the fréchet derivative of a matrix function using Krylov subspace methods*, Numer. Linear Algebra Appl.,
521 28 (2021), p. e2401.

- 522 [58] P. KANDOLF AND S. D. RELTON, *A block Krylov method to compute the action of the fréchet derivative of a*
 523 *matrix function on a vector with applications to condition number estimation*, SIAM J. Sci. Comput., 39
 524 (2017), pp. A1416–A1434.
- 525 [59] KHROSOS[©] SYCL[™] WORKING GROUP, *SYCL[™] Specification*, no. Git revision: tags/SYCL-1.2.1/final-rev7-
 526 0-g7145c7006, The Khronos[©] Group, Apr. 2020. Version 1.2.1.
- 527 [60] A. KOSKELA, *Approximating the matrix exponential of an advection-diffusion operator using the incomplete*
 528 *orthogonalization method*, in Numerical Mathematics and Advanced Applications-ENUMATH 2013: Pro-
 529 ceedings of ENUMATH 2013, the 10th European Conference on Numerical Mathematics and Advanced
 530 Applications, Lausanne, August 2013, Springer, 2014, pp. 345–353.
- 531 [61] D. KRESSNER, *A Krylov subspace method for the approximation of bivariate matrix functions*, Structured
 532 Matrices in Numerical Linear Algebra: Analysis, Algorithms and Applications, (2019), pp. 197–214.
- 533 [62] S. MASSEI AND L. ROBOL, *Rational Krylov for Stieltjes matrix functions: Convergence and pole selection*,
 534 BIT, 61 (2020), pp. 237–273.
- 535 [63] B. V. MINCHEV AND W. WRIGHT, *A review of exponential integrators for first order semi-linear problems*,
 536 (2005).
- 537 [64] NVIDIA CORPORATION, *NVIDIA H100 tensor core GPU architecture*, tech. rep., 2022.
- 538 [65] E. OKTAY AND E. CARSON, *Using Mixed Precision in Low-Synchronization Reorthogonalized Block Classical*
 539 *Gram-Schmidt*, PAMM, 23 (2023), p. e202200060.
- 540 [66] D. PALITTA, M. SCHWEITZER, AND V. SIMONCINI, *Sketched and truncated polynomial Krylov subspace*
 541 *methods: Evaluation of matrix functions*, arXiv:2306.06481 [math.NA], 2023.
- 542 [67] D. PERSSON AND D. KRESSNER, *Randomized low-rank approximation of monotone matrix functions*, SIAM
 543 J. Matrix Anal. Appl., 44 (2023), pp. 894–918.
- 544 [68] Y. SAAD, *Iterative methods for sparse linear systems*, SIAM, Philadelphia, 2nd ed ed., 2003.
- 545 [69] M. SCHWEITZER, *Restarting and error estimation in polynomial and extended Krylov subspace methods for*
 546 *the approximation of matrix functions*, PhD thesis, Fakultät für Mathematik und Naturwissenschaften,
 547 Bergische Universität Wuppertal, 2015.
- 548 [70] W. SQUIRE AND G. TRAPP, *Using complex variables to estimate derivatives of real functions*, SIAM review,
 549 40 (1998), pp. 110–112.
- 550 [71] K. ŚWIRYDOWICZ, J. LANGOU, S. ANANTHAN, U. YANG, AND S. THOMAS, *Low synchronization*
 551 *Gram–Schmidt and generalized minimal residual algorithms*, Numer Linear Algebra Appl, 28 (2021),
 552 p. e2343.
- 553 [72] THE MARDI CONSORTIUM, *Research data management planning in mathematics*, e-print, 2023.
- 554 [73] S. TOMOV, J. DONGARRA, AND M. BABOULIN, *Towards dense linear algebra for hybrid GPU accelerated*
 555 *manycore systems*, Parallel Comput., 36 (2010), p. 232–240.
- 556 [74] S. TOMOV, R. NATH, H. LTAIEF, AND J. DONGARRA, *Dense linear algebra solvers for multicore with*
 557 *GPU accelerators*, in Proceedings of the 2010 IEEE International Symposium on Parallel & Distributed
 558 Processing, Workshops and Phd Forum (IPDPSW), Institute of Electrical and Electronics Engineers, Apr.
 559 2010.
- 560 [75] I. YAMAZAKI, E. CARSON, AND B. KELLEY, *Mixed Precision s-step Conjugate Gradient with Residual*
 561 *Replacement on GPUs*, in 2022 IEEE Int. Parallel Distrib. Process. Symp. IPDPS, 2022, pp. 886–896.
- 562 [76] I. YAMAZAKI, S. THOMAS, M. HOEMMEN, E. G. BOMAN, K. ŚWIRYDOWICZ, AND J. J. EILLIOT, *Low-*
 563 *synchronization orthogonalization schemes for s-step and pipelined Krylov solvers in Trilinos*, in Proc.
 564 2020 SIAM Conf. Parallel Process. Sci. Comput. PP, 2020, pp. 118–128.
- 565 [77] I. YAMAZAKI, S. TOMOV, J. KURZAK, J. J. DONGARRA, AND J. L. BARLOW, *Mixed-precision block Gram*
 566 *Schmidt orthogonalization*, Proc. ScalA 2015 6th Workshop Latest Adv. Scalable Algorithms Large-Scale
 567 Syst. - Held Conjunction SC 2015 Int. Conf. High Perform. Comput. Netw. Storage Anal., (2015).
- 568 [78] Q. ZOU, *A flexible block classical Gram–Schmidt skeleton with reorthogonalization*, Numerical Linear Algebra
 569 with Applications, 30 (2023), p. e2491.