

The Numerical Evaluation of Slater Integrals on GPU

Duy-Khoi Dang, Leighton W. Wilson, and Paul M. Zimmerman*

University of Michigan

930 N. University Ave.

Ann Arbor, MI 48109

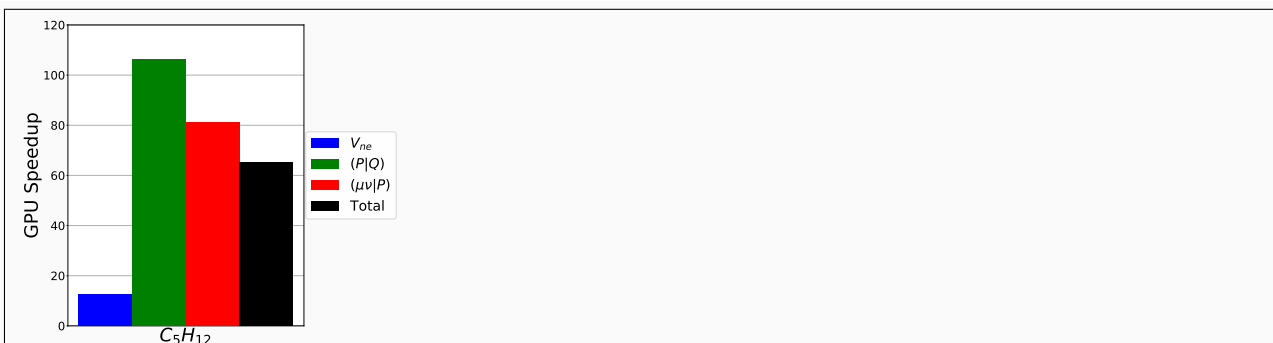
June 30, 2022

Abstract

This article presents SlaterGPU, a GPU accelerated library that uses OpenACC to numerically compute Slater-type orbital (STO) integrals. The electron repulsion integrals (ERI) are computed under the RI approximation using the Coulomb potential of the Slater basis function. To fully realize the performance capabilities of modern GPUs, the Slater integrals are evaluated in mixed-precision, resulting in speedups for the ERIs of over $80\times$. Parallelization on multiple GPUs allows for integral throughput of over 3 million integrals per second. This places STO integral throughput within reach of single-threaded, conventional Gaussian integration schemes. To test the quality of the integrals, the fluorine exchange reaction barrier in fluoromethane was computed using heat-bath configuration interaction (HBCI). In addition, the singlet-triplet gap of cyclobutadiene was examined using HBCI in a triple- ζ , polarized basis set. These benchmarks demonstrate the library's ability to generate the full set of integrals necessary for configuration interaction with up to $6h$ functions in the auxiliary basis.

Keywords: Slater orbitals, integrals, configuration interaction, GPU

*paulzim@umich.edu



SlaterGPU, a GPU accelerated library for numerically computing Slater-type orbital (STO) integrals, is presented in this study. The library achieves speedups over CPU for the electron repulsion integrals of over $80\times$. By utilizing mixed-precision arithmetic and multi-GPU parallelism, SlaterGPU achieves STO integral throughput of over 3 million integrals per second. SlaterGPU also generates the full complement of electron integrals needed for methods such as full configuration interaction.

INTRODUCTION

Advances in quantum chemistry and computer hardware have facilitated the routine use of electronic structure simulations for chemical applications. Some of the most widely used theories make use of one-electron, atom-centered basis functions¹ to represent the electron density. The simplest wave function that approximately solves the Schrödinger equation is Hartree-Fock (HF), which represents the wave function using a single Slater determinant. While HF is not a quantitatively accurate theory, it forms the basis for more sophisticated theories. In many canonical, post-HF methods, evaluation of Hamiltonian elements in the Schrödinger picture requires computing integrals of the form

$$O_{\mu\nu} = \langle \chi_\mu | \hat{O}_1 | \chi_\nu \rangle = \int \chi_\mu(\mathbf{r}) \hat{O}_1(\mathbf{r}) \chi_\nu(\mathbf{r}) d\mathbf{r}, \quad (1)$$

$$\begin{aligned} O_{\mu\nu\lambda\sigma} &= \langle \chi_\mu(1) \chi_\nu(1) | \hat{O}_2 | \chi_\lambda(2) \chi_\sigma(2) \rangle \\ &= \int \int \chi_\mu(\mathbf{r}_1) \chi_\nu(\mathbf{r}_1) \hat{O}_2 \chi_\lambda(\mathbf{r}_2) \chi_\sigma(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2. \end{aligned} \quad (2)$$

The first equation denotes 1-electron quantities such as the overlap $\hat{O}_1 = 1$, the kinetic energy $\hat{O}_1 = -\frac{1}{2}\nabla^2$, and the nuclear attraction $\hat{O}_1 = \frac{Z_A}{R_{1A}}$ operators. 2-electron operators include the Coulomb repulsion $\hat{O}_2 = \frac{1}{r_{12}}$, where r_{12} is the distance between electrons 1 and 2. Derivatives of these terms, for example with respect to nuclear position, are also quantities of interest.

Amongst these integrals, the electron repulsion integrals (ERIs) are the most difficult (and numerous) to evaluate, being 2-electron quantities that require six-dimensional integration. In addition, the $\frac{1}{r_{12}}$ operator contains a singularity at every point in three-dimensional space, further challenging their integration. Consequently, the choice of basis is important for not only accurate representation of the molecular wave function, but also for computational evaluation of integrals.

One physically-motivated choice are Slater-type orbitals (STOs), which are hydrogen-like orbitals of the form

$$S(\zeta, n, l, m, r, \theta, \phi) = N^{\text{STO}} r^{n-1} e^{-\zeta r} Y_{lm}(\theta, \phi), \quad (3)$$

47 where ζ is the exponent, n, l, m are the usual atomic quantum numbers, r, θ, ϕ are spherical co-
48 ordinates, N^{STO} is the normalization constant, and Z_{lm} are the spherical harmonics.¹⁻³ The STOs
49 satisfy the Kato cusp and exponential decay of atomic wave functions,^{1,4,5} making them a natu-
50 ral basis choice for quantum chemical calculations. However, the ERIs over STOs do not have a
51 known general analytic form.

52 The difficulties of STO integration led to the expansion of STOs in terms of Gaussian-type
53 orbitals⁶ (GTOs)

$$54 \quad G(\alpha, n, l, m, r, \theta, \phi) = N^{\text{GTO}} e^{-\alpha r^2} S_{lm}(r, \theta, \phi), \quad (4)$$

55 where the S_{lm} are the real solid harmonics.¹ The GTOs benefit from the Gaussian product rule,
56 i.e. the product of two GTOs is again a GTO, which simplifies Equation 2 from a 4-center, 2-
57 electron integral to 2-center, 2-electron. The 2-center, 2-electron integral can be evaluated over
58 the Coulomb potential of one GTO reducing a 6-dimensional integral to 3-dimensions. These nice
59 analytical properties of GTOs facilitated the development of fast analytical integral evaluation.⁷⁻¹¹

60 While GTOs can be quickly evaluated using modern integral libraries, they do not contain the
61 correct short- and long-range behaviors expected in molecular wave functions.¹² For example, the
62 cusp near the nucleus is important for computing properties such as nuclear magnetic resonance
63 shifts and polarizabilities,^{13,14} but the cusp is not present in the GTO basis, and only crudely
64 treated by using contracted sets of GTOs. Exponential decay of the wave function for an accurate
65 description is required for precise quantification of the HOMO energy, but this behavior is also
66 absent in GTOs.¹⁵

67 The imperfections of GTO basis sets have left room for the continued development and use
68 of STOs for quantum chemical applications. Several schemes have been developed to compute
69 general STO integrals. One approach is to expand each STO in a very large number of GTOs
70 and compute the GTO integrals analytically.¹⁶⁻¹⁸ Additionally, Monte Carlo has been used to
71 correct integrals over Gaussian expansions to evaluate the Slater quantity.¹⁹ These schemes are
72 prohibitively expensive for routine use. While the focus of this article is on the use of STOs in
73 integrals such as Equations 1 and 2, STOs have seen frequent use in quantum Monte Carlo (QMC)

74 wave functions, where the 1- and 2-electron integrals are not important.^{20–23}

75 An attractive alternative to explicit evaluation of STO ERIs involves density fitting—in partic-
76 ular the resolution-of-the-identity (RI) approximation (see Theory and Computational Details)—
77 which allows Equation 2 to be approximated as a tensor product of 2- and 3-index ERIs. Within
78 the RI approximation, one of the two electrons is described by a single basis function. This facili-
79 tates the use of a Coulomb potential to represent one electron without relying on an explicit basis
80 set product rule—which does not exist for STOs—to condense multiple centers. This simplifica-
81 tion, which is only necessary for systems with at least four distinct atomic centers, allows STO
82 integration of ERIs to be amenable to numerical quadrature schemes. The Amsterdam Density
83 Functional (ADF) package and other density functional theory (DFT) codes implement a density
84 fitting approach to use STOs in DFT.^{3,13,24,25} Other density fitting frameworks have allowed STOs
85 to be used used in approximate MP2,²⁶ double-hybrid DFT,²⁷ and Green’s function methods.^{28,29}
86 These previous STO studies, however, did not generate the full complement of ERIs required for
87 multiconfigurational methods such as those based on configuration interaction,^{30–33} multiconfigu-
88 rational self-consistent field^{34–38} and coupled cluster.^{39–42}

89 This study introduces and benchmarks a graphics processing unit (GPU) library for evaluating
90 STO integrals for wave function theories. The article will show that these can be accurately and
91 efficiently evaluated using numerical integration by combining the RI approximation with the STO
92 Coulomb potential. The large number of processing cores and high memory bandwidth make mod-
93 ern GPUs the architecture of choice for evaluating and summing numerical grids. For additional
94 performance, the integrals are also computed using mixed-precision evaluation. Timings suggest
95 that this library allows STOs to be useful alongside strongly correlated wave function theories.
96 Accuracy benchmarks indicate minimal loss in accuracy from using mixed-precision relative to
97 double-precision. The resulting code, called SlaterGPU,⁴³ is the first reported library to use GPUs
98 to accelerate STO integrals and evaluate the full set of 1- and 2-electron STO integrals up to the $6h$
99 subshell as well as $5g$ for first derivatives for the auxiliary basis.

100 THEORY AND COMPUTATIONAL DETAILS

101 The present STO integral scheme relies on numerical integration over atom-centered grids. Grid-
102 based integration can make use of single instruction, multiple data (SIMD) parallelism and there-
103 fore can leverage GPU hardware for acceleration. Even with this acceleration, the 6-dimensional
104 ERIs remain too costly for routine computations. The dimensionality of integration can be reduced,
105 however, by employing the resolution-of-the-identity (RI) approximation,^{44–46} where the Coulomb
106 potentials for the auxiliary basis functions are known analytically. The various components of the
107 STO integral algorithm are explained in the following sections: the Resolution of the Identity, Grid
108 Construction, Implementation on GPU, and Computational Details.

109 Resolution of the Identity

110 This section focuses on simplifying the challenging ERIs for numerical evaluation. The expres-
111 sions for numerically evaluating the 1-electron integrals are listed in Section S1 of the Supporting
112 Information. In the RI approximation, the 4-index ERIs $(\mu\nu|\lambda\kappa)$ are decomposed into tensor prod-
113 ucts of 2- and 3-index integrals by representing the density in terms of an auxiliary basis. Using
114 the Coulomb metric, the integral can be approximated with the expression^{44–46}

$$115 (\mu\nu|\lambda\kappa) \approx \sum_{PQ} (\mu\nu|P)(PQ)^{-1}(Q|\lambda\kappa) = \sum_Q B_{\mu\nu}^Q B_{\lambda\kappa}^Q, \quad (5)$$

116 where

$$117 B_{\mu\nu}^Q = \sum_P (\mu\nu|P)(PQ)^{-1/2}. \quad (6)$$

118 In a numerical integration scheme, Equation 5 not only reduces the count of numerical integrals
119 for a given basis set size (N) from $O(N^4)$ to $O(N^3)$, it also has a secondary consequence that is
120 useful in the context of STO basis functions. Specifically, the integral

$$121 (P|\mu\nu) = \int \int \chi_P(\mathbf{r}_1) \frac{1}{r_{12}} \chi_\mu(\mathbf{r}_2) \chi_\nu(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad (7)$$

122 can be simplified to

$$123 \quad (P|\mu\nu) = \int V_C^P(\mathbf{r})\chi_\mu(\mathbf{r})\chi_\nu(\mathbf{r})d\mathbf{r} \quad (8)$$

124 by using the known analytical form of the single-Slater Coulomb potential. In spherical coordi-
125 nates, this potential has the form³

$$126 \quad V_C(\zeta, n, l, m, r, \theta, \phi) = \frac{4\pi(2\zeta)^{n+(1/2)}}{\sqrt{(2n)!(2l+1)}} Z_{lm}(\theta, \phi) I_{nl}(r), \quad (9)$$

127 where

$$128 \quad I_{nl}(r) = r^{-l-1} \int_0^r (r')^{n+l+1} e^{-\zeta r'} dr' + r^l \int_r^\infty (r')^{n-l} e^{-\zeta r'} dr'. \quad (10)$$

129 I_{nl} has analytic expressions using finite Laurent polynomials for each n, l of interest.

130 For large angular momentum l , the Laurent expressions (see Section S2 of the Supporting
131 Information)—and especially their derivatives—exhibit numerical instability, especially when us-
132 ing mixed-precision arithmetic, which is essential for high performance integral evaluation. This
133 can result in non-smooth integrals as shown in Figure 1, which can in turn result in non-smooth or
134 discontinuous energies. Instead of applying the Laurent expressions, Equation 10 can be evaluated
135 using lower incomplete gamma functions, which have fast, numerically precise implementations.⁴⁷

136 The final form of Equation 10 used in the current implementation of SlaterGPU is

$$137 \quad I_{nl}(r) = r^{-l-1} \zeta^{-l-n-2} \left\{ (r\zeta)^{2l+1} [(-l+n)! - \gamma(-l+n+1, r\zeta)] + \gamma(l+n+2, r\zeta) \right\}, \quad (11)$$

138 where $\gamma(s, x)$ is the lower incomplete gamma function,

$$139 \quad \gamma(s, x) = \int_0^x t^{s-1} e^{-t} dt. \quad (12)$$

140 After evaluation of all 2- and 3-center Coulomb integrals, the full set of 4-index ERIs can
141 be reconstructed using Equation 5. SlaterGPU therefore uses the RI approximation for Slater
142 integrals, similar to prior implementations for DFT applications,^{24,25} but further provides all 4-

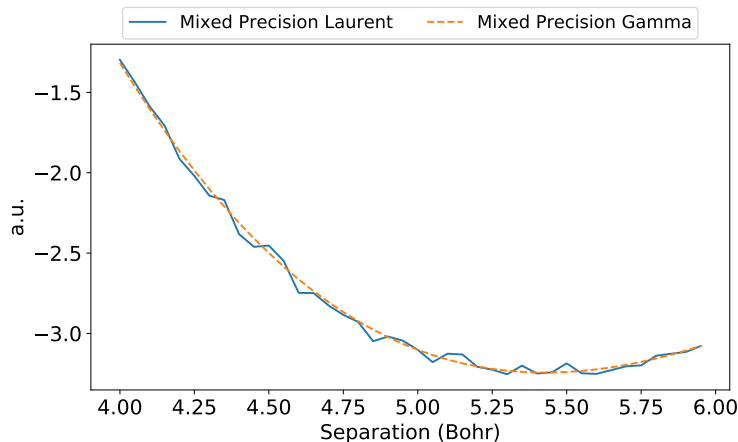


Figure 1: The $(6H|6H)$ integral is scanned in the $(0.370, 0.370, 0.853)$ direction with the left center at the origin. Evaluations are in mixed precision using either the Laurent polynomial expansion of Equation 10 or the lower incomplete gamma function, where mixed precision is defined similarly to Equation 14. Both basis functions have $m = 0$ and $\zeta = 1$.

143 index integrals, $(ij|kl)$, which are not generated or required for DFT. This allows the SlaterGPU
 144 library to be useful for wave function theories, which require a larger set of ERIs. In particular,
 145 while prior codes demonstrated applicability to $l \leq 3$,^{25,29} SlaterGPU is shown here to be useful
 146 for $l \leq 5$.

147 Grid Construction

148 When numerically evaluating integrals over atomic orbitals (AO), the choice of grid is important.
 149 The atom-centered grids used here borrow their core concept from prior studies, especially those
 150 involving integration of DFT functionals.^{24,48–52} The accepted route for integrating the exchange-
 151 correlation energy is to build atomic grids as products of radial and angular grids, then reweight
 152 these using Voronoi polyhedra centered about the nuclei. The atom-centered grids are neces-
 153 sary to capture the spherical harmonics and radial decay of atomic orbitals, while partitioning
 154 3-dimensional space into polyhedra divides the grid into volumes centered around each nucleus.
 155 The Voronoi boundaries are smoothed and reweighted to avoid double counting of volume ele-
 156 ments.^{49,50,52} This same framework is used in SlaterGPU, though only a maximum of three atom-
 157 centered grids are required for any given integral, since the ERIs only involve up to three centers at

158 a time in the RI approximation. In a polyatomic system, this greatly simplifies the form of the inte-
 159 gration grid, keeping each integral grid small enough to be efficiently evaluated. The grids chosen
 160 for this implementation are the "Log3" grid from Mura and Knowles⁵² for the radial component,
 161 and the Lebedev grid⁵³ for the angular component. Both grids are widely used in electronic struc-
 162 ture codes. Once each atom-centered grid is generated, and the Becke partitioning scheme⁴⁹ is
 163 applied, the 3-center integral $(\mu\nu|P)$ can be evaluated over the grid points x and grid weights $w(x)$
 164 as

$$165 \quad (\mu\nu|P) = N_{V_C}^{\text{STO}} N_{\chi_\mu}^{\text{STO}} N_{\chi_\nu}^{\text{STO}} \sum_x \bar{V}_C(x) \bar{\chi}_\mu(x) \bar{\chi}_\nu(x) w(x), \quad (13)$$

166 where \bar{V}_C , $\bar{\chi}_\mu$, $\bar{\chi}_\nu$ are the Coulomb potential and basis functions with their respective normalization
 167 constants, $N_{V_C}^{\text{STO}}$, $N_{\chi_\mu}^{\text{STO}}$, and $N_{\chi_\nu}^{\text{STO}}$, factored out. While Lebedev and Mura-Knowles grids are used
 168 with Becke weights, any quadrature grid and weighting scheme can be used in Equation 13. The
 169 2-index integrals $(P|Q)$ are evaluated in a similar manner.

170 Implementation on GPU

171 All integral code in this study is written in C++ using OpenACC for GPU acceleration, which has
 172 the advantage of being based on pragma directives allowing the same code base to be compiled
 173 to run on CPUs or GPUs. When evaluating Equation 10, a modified version from the Cephis
 174 library⁵⁴ was used for the lower incomplete gamma function, noting that OpenACC allows these
 175 implementations to be used directly. Most GPUs contain more single precision compute units
 176 than double precision, so mixed precision operations are an attractive choice in a practical im-
 177 plementation.⁵⁵⁻⁵⁸ For example, the 2080-Ti contains $\frac{1}{32}$ the double precision units compared to
 178 single precision units, while the GV100 contains $\frac{1}{2}$ the double precision units compared to single
 179 precision units. In SlaterGPU, mixed precision is available, where evaluations over the grid are
 180 performed using single-precision arithmetic, and the final summation occurs in double precision.

181 In mixed precision, Equation 13 becomes

$$182 \quad (\mu\nu|P)_{64} = N_{V_C}^{\text{STO}} N_{\chi_\mu}^{\text{STO}} N_{\chi_\nu}^{\text{STO}} \sum_x \bar{V}_P(x)_{32} \bar{\chi}_\mu(x)_{32} \bar{\chi}_\nu(x)_{32} w(x)_{32}, \quad (14)$$

183 where the subscript refers to the bits of precision of the quantity. Factorization of the normalization
 184 constants reduces the number of floating point operations required, which is essential for high
 185 performance. In addition to making use of the greater quantity of single-precision compute units
 186 in GPUs, single precision also reduces storage and memory bandwidth demands by a factor of 2. A
 187 generalization of the mixed precision procedure would be to adaptively determine which integrals
 188 to evaluate at each level of precision, as has been done in (analytic) GTO integration.⁵⁸ This is
 189 not done here; instead, the accuracy of the mixed precision approach is evaluated in comparison to
 190 double precision integration.

Algorithm 1 GPU compute structure for generating 3-center ERIs

```

1: #pragma omp parallel for schedule(dynamic)
2: for A, B, C in Atom List // A, B over all Atoms, C ≥ B
3:   Generate xA, xB, xC, w(xA), w(xB), w(xC)
4:   x ← xA ∪ xB ∪ xC
5:   w(x) ← w(xA) ∪ w(xB) ∪ w(xC)
6:   for Pi in Aux(A)
7:     Compute VPi(x)
8:     for χμj in Basis(B)
9:       Compute χμj(x)
10:    for χνk in Basis(C)
11:      Compute χνk(x)
12:    for Pi ∈ Aux(A), μj ∈ Basis(B), νk ∈ Basis(C)
13:      (Pi|μjνk) ← ∑x VPi(x) χμj(x) χνk(x) w(x)

```

191 Aux(·) and Basis(·) denote the set of auxiliary and main basis functions centered at ·, respectively.

192 In the GPU computing framework, data transfers between CPU and GPU incur large overhead
 193 penalties, and thus it is necessary to minimize these transactions for maximum performance. As
 194 such, all quantities in Equation 14 are generated and evaluated directly on GPU. The grid x and its
 195 weights $w(x)$ only depend on the set of atoms and not the basis functions, so these are generated
 196 once for each unique triad of atoms as described in Algorithm 1. Additional computations may be

197 avoided by evaluating each V_P , χ_μ and χ_ν on the grid only once per triad of atoms. In other words,
198 when evaluating $(P_i|\mu_j\nu_k)$, the quantities $V_P(x)$, $\chi_\mu(x)$ and $\chi_\nu(x)$ are all computed and stored as
199 GPU arrays for all P_i on atom A , μ_j on atom B and ν_k on atom C to avoid duplicating computations.
200 These arrays can then be contracted all at once in a single tensor operation as indicated in lines
201 11-12 of Algorithm 1. With OpenACC, the contraction on line 12 can be handled using a single
202 *pragma* directive containing the parallel and reduction clauses. Sample OpenACC code is pro-
203 vided in Section S8 of the Supporting Information. As the grid and weights are generated directly
204 on GPU, reuse of the grid benefits from the high memory bandwidth of the GPU (~ 600 GB/s on
205 the 2080-Ti). Once the integrals are computed on GPU, a single data transfer step returns the in-
206 tegrals to CPU memory. The code for numerically computing the STO integrals is freely available
207 on GitHub under an Apache 2.0 license with Commons Clause as noted in the Data Availability
208 Statement.

209 Multi-GPU parallelization is also implemented for a single node, using OpenMP to manage
210 the GPU processes. Each OpenMP thread is assigned a GPU, and a manager-worker scheme is
211 used for load balancing, where the work is partitioned using sets of atoms to take advantage of
212 grid/weight reuse. The parallelization occurs over the loop in Line 2 of Algorithm 1 and can be
213 accomplished with a single *pragma* directive, shown in Line 1.

214 Computational Details

215 An all-electron double-zeta STO basis set with polarization functions⁵⁹ was used (denoted DZP)
216 as the primary atomic orbital basis. The auxiliary basis sets were taken from the same source.
217 Full specification for the primary and auxiliary basis sets are provided in Section S3 of the Sup-
218 porting Information. Unless otherwise specified, the integration grid was a direct product of 60
219 radial points and 770 angular points (Lebedev order 18). Hartree-Fock and heat-bath configuration
220 interaction^{30,60-63} (HBCI) were used as representative electronic structure methods. The HBCI
221 parameters used are detailed in the following section. For GTOs, the 6-31G* basis with the RI-
222 cc-pVTZ auxiliary basis was used. All GTO integral evaluation was performed using the Libcint

223 library.⁶⁴ Molecules were placed in standard nuclear orientations.⁶⁵ The Nvidia HPC SDK 20.7
224 compiler suite with CUDA 11.0 was used to compile all code. CPU code was run on Intel Xeon
225 Gold 6242 processors clocked at 2.8 GHz and GPU code was executed using the Nvidia RTX
226 2080-Ti and GV-100 GPUs.

227 RESULTS AND DISCUSSION

228 Performance Analysis

229 High throughput integral evaluation is necessary for any electronic structure theory code, regard-
230 less of basis set type. Grid-based numerical integration, however, requires orders of magnitude
231 more floating point operations than analytical integration. To achieve the integral performance
232 required, GPUs are used in this study for numerical integration of STO integrals. These inte-
233 grals include all of the common 2-center integrals (overlap, electron-nuclear attraction, and kinetic
234 energy) as well as the 2- and 3-center Coulomb integrals needed for the RI approximation. The rel-
235 ative speedup for numerical GPU integration compared to CPU integration is visualized in Figure
236 2. In double precision, the 2080-Ti can achieve over $30\times$ speedup and the GV100 achieves $\sim 70\times$
237 speedup, allowing for tractable wall times for the integrals as listed in Table 1. Even further per-
238 formance can be gained by utilizing mixed precision, showing speedups of over $60\times$ and overall
239 integral throughput increasing by a factor of ~ 4 for the 2080-Ti (see Table 2). The speedup rela-
240 tive to CPU drops slightly to $\sim 55\times$ on the GV100. The performance behavior is a consequence of
241 the hardware configuration and the integral kernels being compute bound, with a detailed analysis
242 provided in the Supporting Information (Section S5).

243 Faster integral evaluation is also possible by distributing the workload across multiple GPUs.
244 To test multi-GPU scaling, 3-center ERIs were evaluated for C_9H_{20} , which has 76.9 million ERIs,
245 taking 125s to compute in mixed precision and 542s in double-precision. Figure 3 shows the
246 strong scaling performance when evaluating the 3-center ERIs for C_9H_{20} , which maintains paral-
247 lel efficiency greater than 75% on up to 5 GPUs for mixed-precision evaluation and greater than

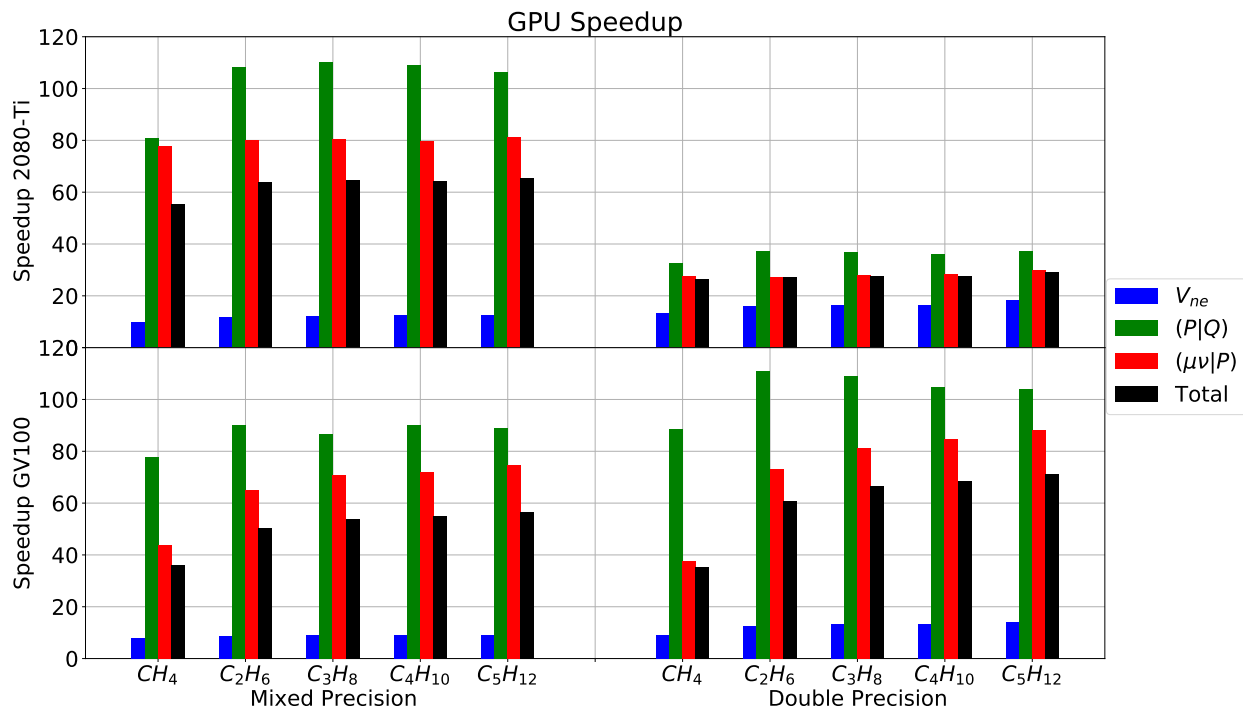


Figure 2: The GPU speedups for integral evaluation over CPU for the 2080-Ti (top) and GV100 (bottom) are shown for various alkanes using the DZP basis from ADF. The speedups are partitioned into the various integrals. Speedups for mixed (left) and double (right) precision evaluations are also shown.

248 90% for double-precision evaluations in this benchmark. Due to the reduced computational de-
 249 mand of mixed-precision integration, the serial components and communication overhead take
 250 up a proportionally larger amount of computational time. Consequently, the parallel efficiency
 251 for mixed-precision integral evaluation drops off more rapidly than for double precision in strong
 252 scaling tests. However, this parallelization scheme still allows STO integration to achieve greater
 253 than 75% parallel efficiency and overall integral throughput greater than 3 million integrals per
 254 second in mixed precision on 5 GPUs. For comparison, Sun reported a throughput of approxi-
 255 mately 6-8 million explicitly calculated integrals per second per thread with Libcint,⁶⁴ thus plac-
 256 ing STO integral throughput within reach of analytical integral evaluation for GTOs. While the
 257 Libcint performance was reported for 4-center ERIs, the comparison demonstrates the feasibility
 258 of Slater integration under the RI approximation. Additional developments in code optimization⁶⁶
 259 and screening protocols⁵⁶ may narrow this gap further. Other grid-based STO integral evalua-

260 tion implementations^{13,25–29} do not report timings nor do they use GPU acceleration. The closest
 261 available performance comparison is an example where a 9-Gaussian expansion was used to ap-
 262 proximate STO integrals,¹⁸ which would reduce throughput, relative to GTOs, by approximately a
 263 factor of 700 under the RI approximation.

Table 1: Double-precision timing data (in seconds) for various alkanes. Each atom contributes 46,200 grid points.

Molecule	Basis size		CPU Time			2080-Ti Time			V100 Time		
	Main	Aux	V_{ne}	$(P Q)$	$(\mu v Q)$	V_{ne}	$(P Q)$	$(\mu v Q)$	V_{ne}	$(P Q)$	$(\mu v Q)$
CH ₄	35	224	3.282	14.02	56.88	0.2483	0.4316	2.062	0.3583	0.1581	1.513
C ₃ H ₈	85	516	40.28	67.98	786.2	2.452	1.835	28.00	2.991	0.6237	9.693
C ₅ H ₁₂	135	808	166.1	163.5	3286	9.125	4.385	110.1	11.71	1.570	37.20

Table 2: Mixed-precision timing data (in seconds) for various alkanes. Each atom contributes 46,200 grid points.

Molecule	Basis size		CPU Time			2080-Ti Time			V100 Time		
	Main	Aux	V_{ne}	$(P Q)$	$(\mu v Q)$	V_{ne}	$(P Q)$	$(\mu v Q)$	V_{ne}	$(P Q)$	$(\mu v Q)$
CH ₄	35	224	2.095	8.686	34.34	0.2139	0.1076	0.4415	0.2708	0.1117	0.7867
C ₃ H ₈	85	516	24.95	42.34	484.8	2.037	0.3838	6.029	2.795	0.4892	6.832
C ₅ H ₁₂	135	808	95.62	98.58	1957	7.680	0.9275	24.12	10.58	1.109	26.21

264 The large performance gain ($\sim 4\times$ speedup) when using mixed precision on the 2080-Ti units
 265 necessarily comes with some loss in accuracy compared to double precision arithmetic. There-
 266 fore tests of the mixed-precision integral evaluation are needed, in order to gauge the quantitative
 267 tradeoff between accuracy and speed.

268 Mixed-Precision Evaluation

269 Numerical evaluation of integrals, whether done in single, double, or mixed precision, will nec-
 270 essarily contain some residual error with any finite grid. While this is expected with grid-based
 271 integration, estimates of the error and smoothness of the resulting integrals are necessary to test
 272 the accuracy of the procedure. First, a selection of 2-center ERIs were evaluated to determine the
 273 relative loss in precision when using the mixed-precision procedure. For each $(P|Q)$, the center Q
 274 was scanned radially away from center P in 16 directions corresponding to all-positive vectors of

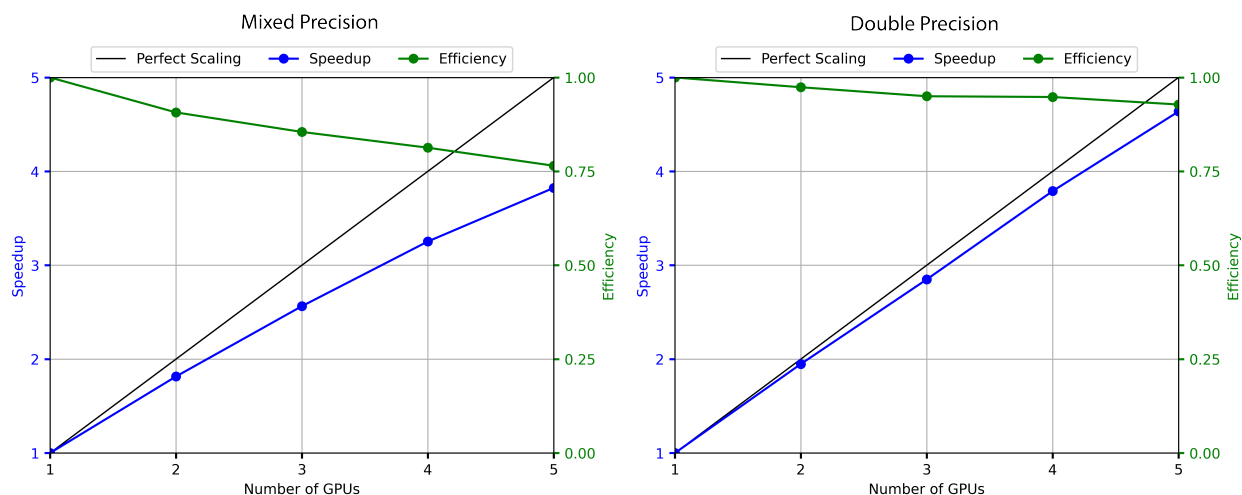


Figure 3: Multi-GPU speedups over single GPU and parallel efficiency for mixed (left) and double (right) precision evaluation of the 3-center ERIs for C_9H_{20} . There are a total of 76,873,200 3-center integrals. Perfect scaling is plotted as a solid black line. All GPUs are co-located on a single compute node. Single GPU run times were 125s and 542s for mixed- and double-precision implementations, respectively.

275 an 86-point Lebedev grid. Figure 4 plots the max and average absolute errors between mixed- and
 276 double-precision integrals at each distance. This indicates that the error of individual integrals are
 277 similarly sized across various distances, with the errors all being less than 8×10^{-5} .

278 The next measure of performance for mixed-precision integration is to evaluate the smoothness
 279 of the integrals with respect to changes in nuclear position. Therefore 2-center ERIs were evaluated
 280 in mixed precision as center Q is scanned radially for the same 16 directions as before. These yield
 281 qualitatively smooth plots, as seen in Figure 5. Additional plots for other basis set pairings are
 282 provided in Section S6 of the Supporting Information and show the same qualitative behavior as
 283 this figure.

284 Hartree-Fock and HBCI

285 Two levels of wave function-based electronic structure theory were selected to provide practical
 286 tests for the Slater GPU integrals. First, the Hartree-Fock energies for a set of benchmark molecules
 287 were computed and these are listed in Table 3 (see Table S2 of the Supporting Information for
 288 timing information). The DZP basis set, corresponding auxiliary basis sets, and grid described in

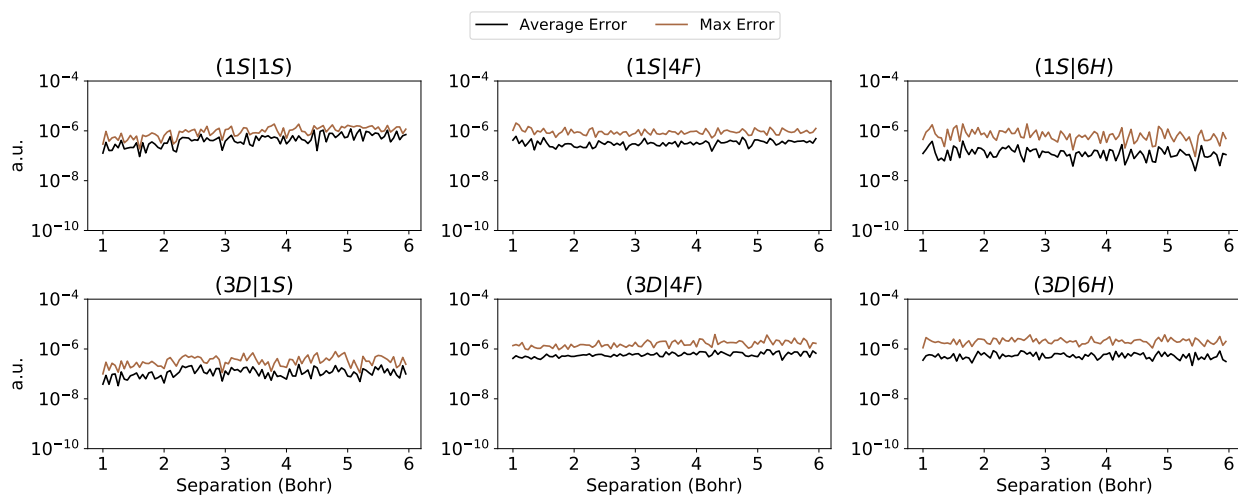


Figure 4: The max and average errors between mixed- and double-precision integral evaluation are plotted for various basis functions. All basis functions have $\zeta = 1$ and $m = 0$. The max and average errors are computed over internuclear distance scans based on the 16 all-positive directions of a Lebedev grid.

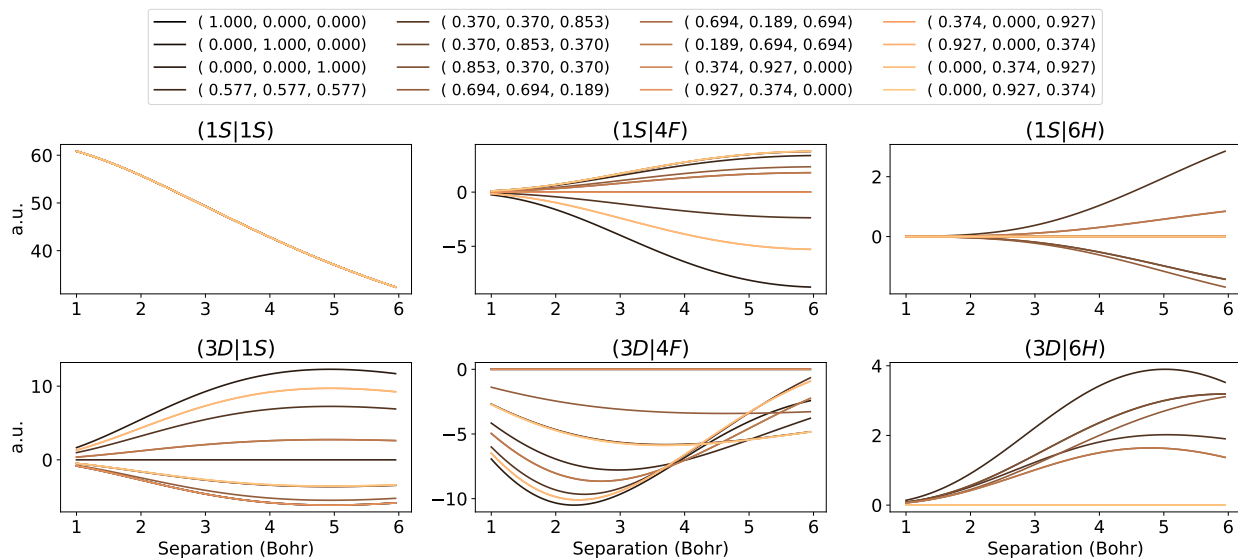


Figure 5: The value of 2-center ERIs are evaluated in mixed precision. All basis functions have $\zeta = 1$ and the right basis is scanned radially away from the origin in various directions. The directions selected are provided in the legend and were selected using the 16 all-positive directions of an 86-point Lebedev grid. For all integrals shown here, $m = 0$. The legend entries are direction unit vectors.

289 the Computational Details were used for these tests. Energies using the 6-31G* and RI-cc-pVTZ
 290 auxiliary GTO basis sets are also reported, to provide a baseline for comparison. The HF results
 291 for alkanes (C_nH_{2n+2}) show a slight increase in the mixed-precision error as the chain length
 292 increases. This is shown in Figure 6, which depicts the relative error of the HF energy when using
 293 mixed- and double-precision at various grid sizes. The roughly constant relative error as system
 294 size grows suggests that the mixed-precision error is size extensive. Combined with Figures 4 and
 295 5, Figure 6 indicates that errors due to using mixed-precision integrals may largely result in error
 296 cancellation.

Table 3: HF energies computed for several small molecules are listed. The STO basis sets and grid used are described in the Computational Details. The numbers in parenthesis in the header denotes the bits of precision used for integral evaluation.

Molecule	DZP (32)	DZP (64)	6-31G*
CH ₄	-40.199728	-40.199730	-40.194806
C ₂ H ₆	-79.232585	-79.232593	-79.227194
C ₃ H ₈	-118.269381	-118.269391	-118.261168
C ₄ H ₁₀	-157.305952	-157.305972	-157.294705
C ₅ H ₁₂	-196.342295	-196.342322	-196.328158
BH ₃	-26.395615	-26.395617	-26.390665
BF ₃	-323.166750	-323.166775	-323.142633
CF ₄	-435.667561	-435.667608	-435.642948
Cr(CO) ₆	-1714.832816	-1714.832901	-1714.469310

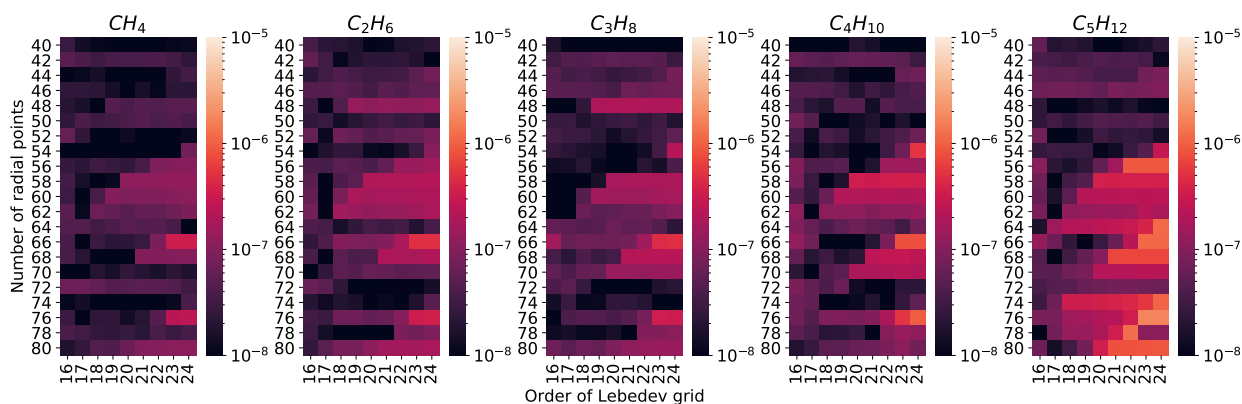


Figure 6: Heatmaps of the relative error of HF energies when using mixed- vs double-precision integral evaluation are shown for various alkanes using different angular and radial grid sizes.

297 The small error margins for STO integrals—as measured at the HF level of theory—suggest

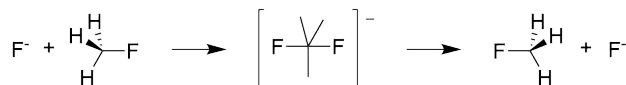


Figure 7: The S_N2 reaction for fluoride exchange of fluoromethane.

Table 4: HF and HBCI activation energies (kcal mol^{-1}) of CH_3F fluoride exchange at various grid sizes using single- and double-precision integral evaluations. The number of radial points and Lebedev order are provided for the radial and angular grids. The size of the angular grid is given in parenthesis next to the Lebedev order.

		Double Precision			Mixed Precision		
		HF					
Radial	Angular	17(590)	18(770)	19(974)	17(590)	18(770)	19(974)
	50	18.4	18.4	18.4	18.4	18.4	18.4
	60	18.4	18.4	18.4	18.4	18.4	18.4
		HBCI					
	50	13.7	13.6	13.8	13.6	13.8	13.7
	60	13.9	13.8	13.7	13.8	13.6	13.6

that thermochemical properties can be precisely evaluated. To test this hypothesis, an S_N2 reaction involving fluoride exchange in fluoromethane was evaluated (Figure 7). Since the HF level of theory is not expected to be quantitative, activation energies were computed not only with HF, but also with the heat-bath configuration interaction (HBCI) method, with ϵ_1 set to 1.0 mHa and ϵ_2 set to 1.0 μHa . HBCI provides a close approximation to full CI, and importantly, is tractable for the $20e^-$ in 60 orbital system of interest here. The activation energies of the exchange reaction using various grids are reported in Table 4. At the grids considered, the change in activation energy at the HF level is negligible between mixed and double precision as well as between grid sizes. At the HBCI level, more integrals contribute to the total energy. Consequently, the variation in the activation energy is larger for HBCI relative to HF. However, the range of activation barriers for HBCI is still less than half a kcal mol^{-1} .

Another test using the HBCI method was the calculation of the singlet-triplet gaps of cyclobutadiene at its D_{2h} and D_{4h} geometries using a triple- ζ polarized (denoted TZP) basis set. Cyclobutadiene has a multireference singlet ground state, due to its degenerate π orbitals in the D_{4h} geometry. The results, for HBCI parameters of $\epsilon_1 = 1.0$ mHa and $\epsilon_2 = 0.1$ μHa , are given in Table 5. These demonstrate mixed-precision errors of less than 0.1 kcal mol^{-1} . Furthermore, the

314 singlet-triplet gap at the D_{4h} is consistent with prior full-CI using GTOs.³¹

Table 5: Relative energies of cyclobutadiene at D_{2h} and D_{4h} geometries (kcal mol⁻¹).

	Double Precision		Mixed Precision	
	D_{2h}	D_{4h}	D_{2h}	D_{4h}
Singlet	0.0	9.4	0.0	9.3
Triplet	36.2	14.2	36.2	14.2
Gap	36.2	4.8	36.2	4.9

315 One final test will further show the utility of the Slater GPU integrals in quantum chemistry.
316 Specifically, the geometric gradients—which are essential in studying chemical reactions—were
317 evaluated using analytical nuclear derivatives of the quantities V_P , μ , and ν in Equation 13. As a
318 benchmark, the fully symmetric BF_3 , BH_3 , CF_4 , and CH_4 molecules were symmetrically stretched.
319 The HF energies as well as the projection of the mixed-precision HF geometric gradient onto each
320 A-X bond (A=B,C;X=H,F) are plotted in Figure 8. As before, the mixed- and double-precision
321 energies overlap with one another. As for the gradient, the magnitude along each A-X bond should
322 be identical for all distances. This is largely achieved in these test cases, however, there is some
323 variation when fluorine is present. For CF_4 , using the ADF fitting basis led to large gradient
324 errors, thus the auxiliary basis of fluorine was extended with additional functions (see Section S3
325 of the Supporting Information for additional details). Since this addition resulted in substantially
326 improved gradients, the remaining variations for BF_3 and CF_4 are attributed to an incomplete RI
327 auxiliary basis. While this work has not examined the choice of RI basis in detail, this subject will
328 need to be revisited in a future study.

329 CONCLUSIONS

330 The SlaterGPU integral code is herein shown capable of evaluating the full complement of ERIs
331 needed for HF and post-HF theories. Modern computer architectures combined with the RI approx-
332 imation have allowed STO integrals to be feasible even though analytic expressions are currently
333 unavailable. The use of mixed-precision integration allows further performance gains—achieving
334 speedups greater than 80× for the ERIs—with minimal loss to accuracy. In the future, comput-

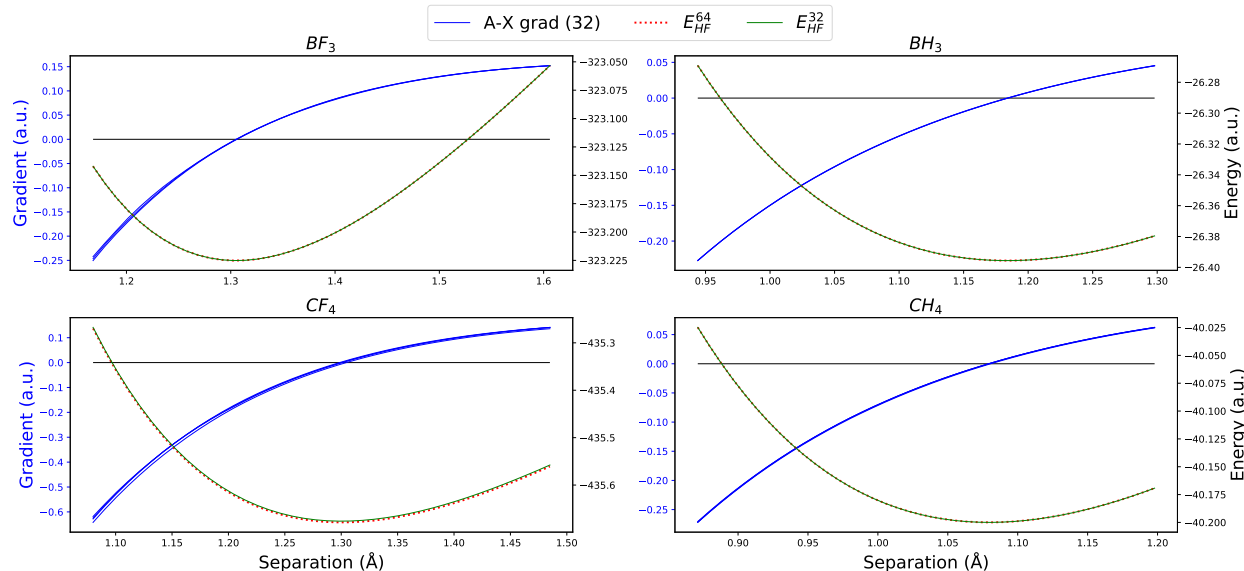


Figure 8: The Hartree-Fock geometric gradient projections (solid blue lines) of molecules with D_{3h} and T_d point groups are plotted as the A-X bond distance is scanned, where A=B,C and X=H,F. Gradients were computed in mixed precision. The mixed-precision (solid green) and double-precision (dotted red) Hartree-Fock energies at each point are also plotted. For CF_4 , the auxiliary basis for fluorine is extended with additional $2p$, $3d$, $4f$ and $5g$ functions.

335 ing select integrals in double precision may mitigate errors due to using mixed-precision integrals.
 336 The combination of GPU acceleration, multi-GPU parallelization, and mixed-precision integration
 337 make SlaterGPU competitive with single-threaded GTO integration with the possibility of tuning
 338 SlaterGPU for additional performance.

339 The current implementation and basis sets are adequate for performing correlated electronic
 340 structure computations at the full CI level, however room for improvement remains in the STO
 341 RI gradients, where the available auxiliary basis sets appear to be inadequate. Further develop-
 342 ment of auxiliary basis sets will be required before STO integrals are generally useful for gradient
 343 computations.

344 ACKNOWLEDGEMENTS

345 The authors acknowledge David Braun for loaning a GV100 for testing. Kris Keipert and Marc
 346 West of Nvidia were instrumental in assisting the authors with OpenACC. This work is supported

347 by a DOE CCS award, DE-SC0022241. D.-K. D thanks the National Science Foundation for
348 providing a Graduate Research Fellowship under grant number DGE 1841052.

349 DATA AVAILABILITY STATEMENT

350 Basis sets and geometries are available in the Supporting Information. The SlaterGPU integral
351 library is available on Github (<https://github.com/ZimmermanGroup/SlaterGPU>) under an Apache
352 2.0 License with Commons Clause.

353 References

- 354 (1) Helgaker, T; Jørgensen, P; Olsen, J, *Molecular Electronic Structure Theory*; John Wiley &
355 Sons, LTD: 2000.
- 356 (2) Slater, J. C. *Phys. Rev.* **1928**, *31*, 333–343.
- 357 (3) Cohen, A. J.; Handy, N. C. *The Journal of Chemical Physics* **2002**, *117*, 1470–1478.
- 358 (4) Kato, T. *Communications on Pure and Applied Mathematics* **1957**, *10*, 151–177.
- 359 (5) Reinhardt, P.; Hoggan, P. E. *International Journal of Quantum Chemistry* **2009**, *109*, 3191–
360 3198.
- 361 (6) Boys, S. F.; Egerton, A. C. *Proceedings of the Royal Society of London. Series A. Mathe-*
362 *matical and Physical Sciences* **1950**, *200*, 542–554.
- 363 (7) McMurchie, L. E.; Davidson, E. R. *Journal of Computational Physics* **1978**, *26*, 218–231.
- 364 (8) Obara, S.; Saika, A. *The Journal of Chemical Physics* **1986**, *84*, 3963–3974.
- 365 (9) Obara, S.; Saika, A. *The Journal of Chemical Physics* **1988**, *89*, 1540–1559.
- 366 (10) Head-Gordon, M.; Pople, J. A. *The Journal of Chemical Physics* **1988**, *89*, 5777–5786.
- 367 (11) Gill, P. M. W.; Head-Gordon, M.; Pople, J. A. *The Journal of Physical Chemistry* **1990**, *94*,
368 5564–5572.

- 369 (12) Kanungo, B.; Zimmerman, P. M.; Gavini, V. *Nature Communications* **2019**, *10*, 4497.
- 370 (13) Watson, M. A.; Handy, N. C.; Cohen, A. J.; Helgaker, T. *The Journal of Chemical Physics*
371 **2004**, *120*, 7252–7261.
- 372 (14) Hoggan, P. E. *International Journal of Quantum Chemistry* **2004**, *100*, 214–220.
- 373 (15) Ayers, P. W.; Morrison, R. C.; Parr, R. G. *Molecular Physics* **2005**, *103*, 2061–2072.
- 374 (16) Hehre, W. J.; Stewart, R. F.; Pople, J. A. *The Journal of Chemical Physics* **1969**, *51*, 2657–
375 2664.
- 376 (17) Hehre, W. J.; Ditchfield, R.; Stewart, R. F.; Pople, J. A. *The Journal of Chemical Physics*
377 **1970**, *52*, 2769–2773.
- 378 (18) Fernández Rico, J.; López, R.; Aguado, A.; Ema, I.; Ramírez, G. *International Journal of*
379 *Quantum Chemistry* **2001**, *81*, 148–153.
- 380 (19) Caffarel, M. *The Journal of Chemical Physics* **2019**, *151*, 064101.
- 381 (20) Nightingale, M.; Umrigar, C., *Quantum Monte Carlo methods in physics and chemistry*;
382 Kluwer Academic Publishers: 1999.
- 383 (21) Zhang, S.; Krakauer, H. *Phys. Rev. Lett.* **2003**, *90*, 136401.
- 384 (22) Zimmerman, P. M.; Toulouse, J.; Zhang, Z.; Musgrave, C. B.; Umrigar, C. J. *The Journal of*
385 *Chemical Physics* **2009**, *131*, 124103.
- 386 (23) Austin, B. M.; Zubarev, D. Y.; Lester, W. A. *Chemical Reviews* **2012**, *112*, PMID: 22196085,
387 263–288.
- 388 (24) Koch, W.; Holthausen, M. C., *A Chemist's Guide to Density Functional Theory*; Wiley -
389 VCH: Weinheim - New York, 2nd edition, 2001.
- 390 (25) Watson, M. A.; Handy, N. C.; Cohen, A. J. *The Journal of Chemical Physics* **2003**, *119*,
391 6475–6481.
- 392 (26) Förster, A.; Franchini, M.; van Lenthe, E.; Visscher, L. *Journal of Chemical Theory and*
393 *Computation* **2020**, *16*, PMID: 31930915, 875–891.

- 394 (27) Förster, A.; Visscher, L. *Journal of Computational Chemistry* **2020**, *41*, 1660–1684.
- 395 (28) Förster, A.; Visscher, L. *Journal of Chemical Theory and Computation* **2020**, *16*, PMID:
396 33174743, 7381–7399.
- 397 (29) Förster, A.; Visscher, L. *Journal of Chemical Theory and Computation* **0000**, *0*, PMID:
398 34236172, null.
- 399 (30) Holmes, A. A.; Tubman, N. M.; Umrigar, C. J. *Journal of Chemical Theory and Computa-*
400 *tion* **2016**, *12*, PMID: 27428771, 3674–3680.
- 401 (31) Zimmerman, P. M. *The Journal of Physical Chemistry A* **2017**, *121*, PMID: 28530830,
402 4712–4720.
- 403 (32) Zimmerman, P. M. *The Journal of Chemical Physics* **2017**, *146*, 104102.
- 404 (33) Zimmerman, P. M. *The Journal of Chemical Physics* **2017**, *146*, 224104.
- 405 (34) Roos, B. O. *International Journal of Quantum Chemistry* **1980**, *18*, 175–189.
- 406 (35) Roos, B. O.; Taylor, P. R.; Sigbahn, P. E. *Chemical Physics* **1980**, *48*, 157–173.
- 407 (36) Roos, B. O.; Taylor, P. R.; Sigbahn, P. E. *Chemical Physics* **1980**, *48*, 157–173.
- 408 (37) Zimmerman, P. M.; Rask, A. E. *The Journal of Chemical Physics* **2019**, *150*, 244117.
- 409 (38) Dang, D.-K.; Zimmerman, P. M. *The Journal of Chemical Physics* **2021**, *154*, 014105.
- 410 (39) Stanton, J. F.; Bartlett, R. J. *The Journal of Chemical Physics* **1993**, *98*, 7029–7039.
- 411 (40) Friedrich, J.; Hanrath, M.; Dolg, M. *The Journal of Chemical Physics* **2007**, *126*, 154110.
- 412 (41) Evangelista, F. A.; Prochnow, E.; Gauss, J.; Schaefer, H. F. *The Journal of Chemical Physics*
413 **2010**, *132*, 074107.
- 414 (42) Evangelista, F. A.; Gauss, J. *The Journal of Chemical Physics* **2011**, *134*, 114102.
- 415 (43) SlaterGPU <https://github.com/ZimmermanGroup/SlaterGPU>.
- 416 (44) Dunlap, B. I.; Connolly, J. W. D.; Sabin, J. R. *The Journal of Chemical Physics* **1979**, *71*,
417 3396–3402.

- 418 (45) Werner, H.-J.; Manby, F. R.; Knowles, P. J. *The Journal of Chemical Physics* **2003**, *118*,
419 8149–8160.
- 420 (46) Distasio JR., R. A.; Steele, R. P.; Rhee, Y. M.; Shao, Y.; Head-Gordon, M. *Journal of Com-*
421 *putational Chemistry* **2007**, *28*, 839–856.
- 422 (47) Press, W. H.; Teukolsky, S. A.; Vetterling, W. T.; Flannery, B. P., *Numerical Recipes 3rd*
423 *Edition: The Art of Scientific Computing*, 3rd ed.; Cambridge University Press: USA, 2007.
- 424 (48) Boerrigter, P. M.; Te Velde, G.; Baerends, J. E. *International Journal of Quantum Chemistry*
425 **1988**, *33*, 87–113.
- 426 (49) Becke, A. D. *The Journal of Chemical Physics* **1988**, *88*, 2547–2553.
- 427 (50) Murray, C. W.; Handy, N. C.; Laming, G. J. *Molecular Physics* **1993**, *78*, 997–1014.
- 428 (51) Treutler, O.; Ahlrichs, R. *The Journal of Chemical Physics* **1995**, *102*, 346–354.
- 429 (52) Mura, M. E.; Knowles, P. J. *The Journal of Chemical Physics* **1996**, *104*, 9848–9858.
- 430 (53) Lebedev, V. *USSR Computational Mathematics and Mathematical Physics* **1976**, *16*, 10–24.
- 431 (54) Implementations of the gamma functions <http://www.netlib.org/cephes/>.
- 432 (55) Ufimtsev, I. S.; Martinez, T. J. *Journal of Chemical Theory and Computation* **2009**, *5*, PMID:
433 26609609, 1004–1015.
- 434 (56) Luehr, N.; Ufimtsev, I. S.; Martínez, T. J. *Journal of Chemical Theory and Computation*
435 **2011**, *7*, PMID: 26606344, 949–954.
- 436 (57) Pokhilko, P.; Epifanovsky, E.; Krylov, A. I. *Journal of Chemical Theory and Computation*
437 **2018**, *14*, PMID: 29969560, 4088–4096.
- 438 (58) Seritan, S.; Bannwarth, C.; Fales, B. S.; Hohenstein, E. G.; Kokkila-Schumacher, S. I. L.;
439 Luehr, N.; Snyder, J. W.; Song, C.; Titov, A. V.; Ufimtsev, I. S.; Martínez, T. J. *The Journal*
440 *of Chemical Physics* **2020**, *152*, 224110.
- 441 (59) Van Lenthe, E.; Baerends, E. J. *Journal of Computational Chemistry* **2003**, *24*, 1142–1156.

- 442 (60) Smith, J. E. T.; Mussard, B.; Holmes, A. A.; Sharma, S. *Journal of Chemical Theory and*
443 *Computation* **2017**, *13*, PMID: 28968097, 5468–5478.
- 444 (61) Chien, A. D.; Holmes, A. A.; Otten, M.; Umrigar, C. J.; Sharma, S.; Zimmerman, P. M. *The*
445 *Journal of Physical Chemistry A* **2018**, *122*, PMID: 29473750, 2714–2722.
- 446 (62) Li, J.; Otten, M.; Holmes, A. A.; Sharma, S.; Umrigar, C. J. *The Journal of Chemical Physics*
447 **2018**, *149*, 214110.
- 448 (63) Brorsen, K. R. *Journal of Chemical Theory and Computation* **2020**, *16*, PMID: 32083870,
449 2379–2388.
- 450 (64) Sun, Q. *Journal of Computational Chemistry* **2015**, *36*, 1664–1671.
- 451 (65) Gill, P. M.; Johnson, B. G.; Pople, J. A. *Chemical Physics Letters* **1993**, *209*, 506–512.
- 452 (66) Khalilov, M.; Timoveev, A. *Journal of Physics: Conference Series* **2021**, *1740*, 012056.