

Software paper for submission to the Journal of Open Research Software

Please submit the completed paper to: editor.jors@ubiquitypress.com

1 (1) Overview

2 Title

spinsim: a GPU optimised simulator of spin half and spin one quantum systems

3 Paper Authors

1. Tritt, Alex;
2. Morris, Joshua;
3. Hochstetter, Joel;
4. Anderson, R. P.;
5. Saunderson, James;
6. Turner, L. D.;

4 Paper Author Roles and Affiliations

1. School of Physics & Astronomy, Monash University, Victoria 3800, Australia.
Primary author of the released packages.
2. School of Physics & Astronomy, Monash University, Victoria 3800, Australia.
Present address: Faculty of Physics, University of Vienna, 1010 Vienna, Austria.
Author of first version of code.
3. School of Physics & Astronomy, Monash University, Victoria 3800, Australia.
Present address: School of Physics, University of Sydney, NSW 2006, Australia.
Optimization and extension to spin one of first version of code.
4. School of Molecular Sciences, La Trobe University, PO box 199, Bendigo, Victoria 3552, Australia.
Original conception of first version of code.
5. Department of Electrical and Computer Systems Engineering, Monash University, Victoria 3800, Australia.
Advice on numerical analysis.
6. School of Physics & Astronomy, Monash University, Victoria 3800, Australia.
Original conception of released version of algorithm.

5 Abstract

spinsim is a *python* package that simulates spin half and spin one quantum mechanical systems following a time dependent Shroedinger equation. It makes use of **numba.cuda** [1], which is an *LLVM* (Low Level Virtual Machine) [2] compiler, and other optimisations, to allow for fast and accurate evaluation on *Nvidia Cuda* [3] compatible systems using GPU parallelisation. **spinsim** is available for installation on *PyPI*, and the source code is available on *github*. The initial use case for the package will be to simulate quantum sensing-based Bose Einstein condensate (BEC) experiments for the Monash University School of Physics and Astronomy spinor BEC lab, but we anticipate it will be useful in simulating any range of spin half or spin

one quantum systems with time dependent Hamiltonians that cannot be solved analytically. These appear in the fields of nuclear magnetic resonance (NMR), nuclear quadrupole resonance (NQR) and magnetic resonance imaging (MRI) experiments and quantum sensing, and with the spin one systems of nitrogen vacancy centres (NVCs) and BECs.

6 Keywords

Time dependent Schroedinger equation; Spin one; Spin half; Integrator; GPU; Solver; python; numba;

7 Introduction

7.1 Motivation

Ultracold rubidium atoms have proven their effectiveness in state of the art technologies in quantum sensing [4], the use of quantum mechanics to make precise measurements of small signals. The rotation of these atoms can be modelled as quantum spin systems, which is the quantum mechanical model for objects with angular momentum. The simplest spin system, spin half (ie spin quantum number of $\frac{1}{2}$, also referred to as a qubit), is quantised into just two quantum spin levels, and this describes the motion of some fundamental particles such as electrons. However, systems more practical for sensing, such as ultracold rubidium atoms, are more accurately described as a spin one quantum system (ie spin quantum number of 1, also referred to as a qutrit), which is quantised into three quantum spin levels.

The design of sensing protocols requires many steps of verification, including simulation. This is especially important, since running real experiments can be expensive and time consuming, and thus it is more practical to debug such protocols quickly and cheaply on a computer. In general, any design of experiments using spin systems could benefit from a fast, accurate method of simulation.

In the past, the spinor Bose Einstein condensate (spinor BEC) lab at Monash University used an in-house, *cython* based script, on which this package is based, and standard differential equation solvers (such as *Mathematica*'s function `NDSolve`) to solve the Schroedinger equation for quantum sensing spin systems. Spin one systems are sometimes approximated to spin half for a faster execution time, at the cost of modelling all effects of the system. However, these methods are not completely optimised for our use case, and therefore come with some issues.

First, while the execution time for these solvers is acceptable for running a small number of experiments, for certain experiments involving large arrays of independent atom clouds (which require many thousands of simulations to be run), this time accumulates to the order of many hours, or even multiple days.

Second, the Schroedinger equation has the geometric property of being norm preserving. In other words, the time evolution operator for a system between two points in time must be unitary. As such, numerical solutions to the Schroedinger equation should also preserve this property. For many numerical methods like those in the Runge Kutta family, the approximations used might not be norm preserving, and the evaluated quantum state may diverge towards an infinite norm, or converge to zero if run for many iterations.

Third, our system (and similar spin systems) can be very oscillatory. In standard conditions for our application, the expected spin projection of a system that we want to solve for can rotate in physical space (alternatively viewed as a point rotating around an abstract object known as a *Bloch sphere*) at a rate of 700kHz. Standard integration methods require very small time steps in order to accurately depict these oscillations.

8 Implementation and architecture

8.1 Quantum mechanics background

`spinsim` solves the time-dependent Schrödinger equation

$$i\hbar \frac{d\psi(t)}{dt} = H(t)\psi(t), \quad (1)$$

where the quantum state $\psi(t) \in \mathbb{C}^N$ is assumed normalised, and the Hamiltonian $H(t) \in \mathbb{C}^{N \times N}$ is Hermitian. Here N is the number of levels in the quantum system. Often we are considering systems with spin S . Here spin-half is the $N = 2$ case, and spin-one is $N = 3$. We set $\hbar = 1$, so that the Hamiltonian has units

$$E = \hbar\omega \quad (2)$$

$$= 2\pi\hbar f, \quad (3)$$

Rather than ... parametrizing the problem ... we consider time varying real coefficients in a linear combination of fixed operators,

$$H(t) = \sum_{j=1}^{N^2-1} a_j(t)A_j. \quad (4)$$

It is well known that ... spin in magnetic field has We exclude the identity ... energy zero point.

In the spin one case there is no standard basis Choices include Gell Mann, quadrupole basis, In general ... Lie algebra ... $\mathfrak{su}(N)$ admits a basis where $i^2 = -1$, $\psi(t) \in \mathbb{C}^2$, and the spin half spin projection operators are given by DONT WRITE OUT, sigma

$$J_x = \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad J_y = \frac{1}{2} \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad \text{and } J_z = \frac{1}{2} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}. \quad (5)$$

Call omega, no implementation detail

Similarly, when `spinsim` is set to spin one mode, it can solve the Schroedinger equation of the form

$$\frac{d}{dt}\psi(t) = -i2\pi(f_x(t)J_x + f_y(t)J_y + f_z(t)J_z + f_q(t)Q)\psi(t). \quad (6)$$

where now $\psi(t) \in \mathbb{C}^3$, and the spin one operators are given by

$$\begin{aligned} J_x &= \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}, & J_y &= \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & -i & 0 \\ i & 0 & -i \\ 0 & i & 0 \end{pmatrix}, \\ J_z &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -1 \end{pmatrix}, & \text{and } Q &= \frac{1}{3} \begin{pmatrix} 1 & 0 & 0 \\ 0 & -2 & 0 \\ 0 & 0 & 1 \end{pmatrix}. \end{aligned} \quad (7)$$

The matrices J_x, J_y, J_z are regular spin operators, and Q is a quadrupole operator. Note that Q is proportional to Q_{zz} as defined by Hamley et al [5], and Q_0 as defined by Di et al [6]. **INLINE Q DIAG. WHY WE'VE PICKED THESE FOUR.** For ground state spin one atoms in a magnetic field, ... Breit Rabi **COULD BE READILY EXTENDED TO INCLUDE FURTHER ELEMENTS OR TO HIGHER SPINS.**

Frequently when one calculates the times series of the quantum state $\psi(t)$, they are also interested in its expected spin projection $\langle J \rangle(t)$. It is given by the vector

$$\langle J \rangle(t) = \begin{pmatrix} \psi(t)^\dagger J_x \psi(t) \\ \psi(t)^\dagger J_y \psi(t) \\ \psi(t)^\dagger J_z \psi(t) \end{pmatrix}, \quad (8)$$

where \cdot^\dagger is the adjoint operator (also referred to as Hermitian conjugate, and is implemented as complex conjugate transpose). The expected spin projection can be interpreted as the average direction that the spin system is oriented in space when many systems of the equivalent state are measured (due to the Heisenberg uncertainty principle, only one component of the exact orientation of a quantum state can be known at any point in time, which is why we need to deal with expected values rather than absolute values). However, it can also be used to model the overall orientation of an ensemble of many quantum systems, like a BEC, for instance. `spinsim` has the functionality to calculate the expected spin projection of a system from its state.

8.2 Unitary time evolution and the Magnus expansion

Given that $\psi(t)$ is a unit vector, it is possible to write $\psi(t)$ in terms of a unitary transformation $U(t, t_0)$ of the state $\psi(t_0)$, for any time t_0 . In other words, for any $t_0, t \in \mathbb{R}$, there is a unitary transformation $U(t, t_0)$ such that

$$\psi(t) = U(t, t_0)\psi(t_0). \quad (9)$$

From Equation (1), $U(t, t_0)$ also follows the Schroedinger equation via

$$\frac{d}{dt}U(t, t_0)\psi(t_0) = -i2\pi H(t)U(t, t_0)\psi(t_0) \text{ for any } \psi(t_0) \in \mathbb{C}^N, \text{ so} \quad (10)$$

$$\frac{d}{dt}U(t, t_0) = -i2\pi H(t)U(t, t_0). \quad (11)$$

Thus, to solve Equation (1) for $\psi(t)$ given a particular $\psi(t_0)$, one only needs to solve Equation (11) for $U(t, t_0)$. Note that because $U(t, t_0)$ is unitary, the calculated value of $\psi(t)$ using this method will not change in magnitude, conserving probability. Also since $U(t, t_0)$ is unitary (and continuous in time t), it can be written as a matrix exponential of a skew Hermitian matrix. For example, if the Hamiltonian $H(t)$ is constant in time (and a real linear combination of Hermitian matrices defined by Equations (5) or (7)), then the solution would be

$$U(t, t_0) = \exp(-i2\pi H(t - t_0)). \quad (12)$$

The act of integrating a linear differential equation by exponentiation can be generalised for time dependent Hamiltonians using the Magnus Series [7], where

$$U(t, t_0) = \exp(-i2\pi G(t - t_0)) \quad (13)$$

$$= \exp(-i2\pi \sum_{m=1}^{\infty} G_m(t - t_0)), \quad (14)$$

where $G(t)$ is a series of $G_m(t)$, each of which are a function of H and its commutators. $G_1(t)$, for example, is just $G_1(t) = H_{\text{av}}t$, where H_{av} is the average value of H over the interval being integrated over, which is the solution for when H is constant in time. The terms $G_m(t)$ for $m > 1$ are corrections to this approximation. One can truncate the Magnus series before infinity to make a Magnus expansion, which can give an approximation for $U(t, t_0)$. Importantly, any truncation of the Magnus expansion is unitary, meaning that probability will be conserved when used. Something else to note about unitary time evolution operators, is that they can be broken into a product of multiple other time evolution operators via

$$U(t, t_0) = U(t, t_1)U(t_1, t_0). \quad (15)$$

Each of these parts can be approximated as a Magnus expansion. So, if the Magnus expansion of time evolution operator is inaccurate, it can be broken down into a product of many time evolution operators which each make a much smaller step between times. These components will each have a much more accurate Magnus expansion.

In fact, many unitary time stepping techniques have been developed based on the Magnus expansion [8]. Some of these techniques use Gauss-Legendre quadrature sampling to avoid the integration of commutators of the Hamiltonian, producing simple expressions that can be used for unitary time stepping. In particular, we used the commutator free 4 (CF4) from [8] as our integration method. Suppose we wish to evaluate a short time step of $U(t_f + dt, t_f)$ by taking a short step dt from a fine sample time t_f . Then as part of the CF4 method, the field functions are sampled at particular times based on the second order Gauss-Legendre quadrature, given by

$$t_1 = t_f + \frac{1}{2}dt \left(1 - \frac{1}{\sqrt{3}} \right), \text{ and} \quad (16)$$

$$t_2 = t_f + \frac{1}{2}dt \left(1 + \frac{1}{\sqrt{3}} \right). \quad (17)$$

Then given the matrices

$$\overline{G}_1 = 2\pi dt \left(\frac{3+2\sqrt{3}}{12}H(t_1) + \frac{3-2\sqrt{3}}{12}H(t_2) \right), \text{ and} \quad (18)$$

$$\overline{G}_2 = 2\pi dt \left(\frac{3-2\sqrt{3}}{12}H(t_1) + \frac{3+2\sqrt{3}}{12}H(t_2) \right), \quad (19)$$

the unitary time step can be calculated using

$$U(t_f + dt, t_f) = \exp(-i\overline{G}_2) \exp(-i\overline{G}_1). \quad (20)$$

8.3 Exponentiator

From Equation (20), evaluating matrix exponentials is a core part of the integration algorithm. However, the **spinsim** package uses the user provided field functions $f(t)$ as a vector representation of the Hamiltonian of the spin system being simulated, rather than generic matrix elements in an array. This allows for implementation of matrix exponentiation that takes advantage of this representation. For all implementations of exponentiation, the **spinsim** computes the matrix exponential

$$E(g) = E(g_x, g_y, g_z, g_q) \quad (21)$$

$$= \exp(-i(g_x J_x + g_y J_y + g_z J_z + g_q Q)). \quad (22)$$

In practice, g is the Lie algebra (spin and quadrature) basis representation of \overline{G}_1 or \overline{G}_2 defined in Equations (18) and (19), which need to be exponentiated in Equation (20). For spin half, the default exponentiator is in an analytic form. For spin one, an exponentiator based on the Lie Trotter product formula [9] is used instead. Explicitly, this formula is

$$\exp(A + B) = \lim_{n \rightarrow \infty} \left(\exp\left(\frac{A}{n}\right) \exp\left(\frac{B}{n}\right) \right)^n \quad (23)$$

Where $A, B \in \mathbb{C}^{N \times N}$ are matrix operators, and $n \in \mathbb{N}$. The equality holds for without the limit for all n if the operators commute, like it does for real and complex numbers. This tells us that we can approximate a matrix exponential of a linear combination of spin operators by reducing the size of the coefficients of

these operators, taking the analytically simple exponentials of each of the operators $-iJ_x, -iJ_y, -iJ_z$ and $-iQ$ separately, multiplying each individual result together, and then raising the combination T to the power of the original reduction.

For our spin one case, the exponential $E(g)$ can be approximated as, for large 2^τ ,

$$E(g) = \exp(-ig_x J_x - ig_y J_y - ig_z J_z - ig_q Q) \quad (24)$$

$$= \exp\left(2^{-\tau}(-ig_x J_x - ig_y J_y - ig_z J_z - ig_q Q)\right)^{2^\tau} \quad (25)$$

$$\begin{aligned} &\approx \left(\exp\left(-i\frac{1}{2}(2^{-\tau}g_z J_z + 2^{-\tau}g_q Q)\right) \right. \\ &\quad \cdot \exp\left(-i(2^{-\tau}g_\phi J_\phi)\right) \\ &\quad \left. \cdot \exp\left(-i\frac{1}{2}(2^{-\tau}g_z J_z + 2^{-\tau}g_q Q)\right) \right)^{2^\tau} \end{aligned} \quad (26)$$

$$= \begin{pmatrix} \left(\frac{\Gamma}{PZ}\right)^2 & \frac{SP}{Z\Phi} & -\left(\frac{\Sigma P}{\Phi}\right)^2 \\ \frac{Z}{SP\Phi} & CP^4 & \frac{SPZ}{\Phi} \\ -\left(\frac{\Sigma\Phi}{P}\right)^2 & \frac{SPZ\Phi}{1} & \left(\frac{\Gamma Z}{P}\right)^2 \end{pmatrix}^{2^\tau} \quad (27)$$

$$= T^{2^\tau}. \quad (28)$$

where

$$\begin{aligned} g_\phi &= \sqrt{g_x^2 + g_y^2} \\ J_\phi &= \frac{g_x}{g_\phi} J_x + \frac{g_y}{g_\phi} J_y \\ \Gamma &= \cos\left(2^{-\tau}\frac{g_\phi}{2}\right) & \Phi &= \exp\left(i2^{-\tau}g_\phi\right) \\ \Sigma &= \sin\left(2^{-\tau}\frac{g_\phi}{2}\right) & Z &= \exp\left(i2^{-\tau}\frac{g_z}{2}\right) \\ C &= \cos\left(2^{-\tau}g_\phi\right) & P &= \exp\left(i2^{-\tau}\frac{g_q}{6}\right) \\ S &= \frac{-i}{\sqrt{2}} \sin\left(2^{-\tau}g_\phi\right) \end{aligned} \quad (29)$$

Once T is calculated, it is then recursively squared τ times to obtain $E(g)$. The approach used for spin one exponentiation means that the package cannot solve arbitrary spin one quantum systems, as that would require the ability to exponentiate a point in the full, 8 dimensional Lie algebra of $\mathfrak{su}(3)$, rather than just the four dimensional subspace spanned by the subalgebra $\mathfrak{su}(2)$ spanned by $\{J_x, J_y, J_z\}$, and the single quadratic operator Q . Including the full algebra could be possible as a feature update if there is demand for it, though just including this subspace is sufficient for our application, and many others, and has the advantage of being able to use this faster, more specialised method of matrix exponentiation.

Note that, the methods for both spin half and spin one use analytic forms of exponentials to construct the result, meaning that all calculated time evolution operators are unitary. This guarantees that the results of `spinsim` maintain unitary and conserve probability as required.

8.4 Discretisation and parallelisation

Consider quantising time with

$$t_k = t_0 + Dt \cdot k, \quad (30)$$

where Dt is the time step of the time series (in contrast to dt , a smaller time step for integration). If we define

$$\psi_k = \psi(t_k) \text{ and} \quad (31)$$

$$U_k = U(t_k, t_{k-1}), \quad (32)$$

then the time series of states ψ_k and time evolution operators U_k satisfies

$$\psi_k = U_k \psi_{k-1} \quad (33)$$

This presents an opportunity for parallelism. While each of the ψ_k must be evaluated sequentially, the value of the U_k is independent of the value of any ψ_{k_0} , or any other U_{k_0} . This means that the time evolution operators U_k can all be calculated in parallel, and it allows **spinsim** to use GPU parallelisation on the level of time sample points, so a speed up is achieved even if just a single simulation is run.

In summary, **spinsim** splits the time evolution of the full simulation into time evolution U_k within small time intervals $[t_{k-1}, t_k]$, which are each calculated massively in parallel on a GPU. When all the U_k are calculated, the CPU then multiplies the U_k together (a comparatively less demanding job than calculating them) using Equation (33) to determine the ψ_k .

Beyond discretisation for parallelisation, we discretise our time evolution operator into individual integration time steps. So, each of the U_k are split into products of time evolution operators between times separated by the integration time step via

$$U(t_k, t_{k-1}) = U(t_k, t_k - dt) \cdots U(t_{k-1} + 2dt, t_{k-1} + dt) U(t_{k-1} + dt, t_{k-1}) \quad (34)$$

$$U_k = u_{L-1}^k \cdots u_0^k, \text{ where} \quad (35)$$

$$u_{L-1}^k = U(t_0 + (k-1)Dt + (l+1)dt, t_0 + (k-1)Dt + ldt) \quad (36)$$

with dt being the integration level time step. Note that the time steps are related to each other by $Dt = Ldt$, where $L \in \mathbb{N}$.

8.5 Rotating frame

If the rotating frame option is selected, the U_k are first calculated within a rotating frame of reference as U_k^r , which, in some situations, reduces the size of the field functions used in the calculation, increasing accuracy. The rotation speed of the rotating frame is calculated locally for each parallel time step U_k , and only for rotations around the z axis. The rotating from field functions f_x^r, f_y^r, f_z^r , and f_q^r are related to the field function from the user input via

$$f_x^r(t) + if_y^r(t) = e^{-i2\pi f_r t}(f_x(t) + if_y(t)), \quad (37)$$

$$f_z^r(t) = f_z(t) - f_r, \text{ and} \quad (38)$$

$$f_q^r(t) = f_q(t), \text{ for spin one.} \quad (39)$$

Where $f_r = f_z(t_k + \frac{1}{2}Dt)$ is sampled the midpoint value of the fields over the interval $[t_{k-1}, t_k]$. This, assuming that a midpoint sample is representative of an average value over the time interval, decreases the magnitude of $f_z(t)$, while leaving the other field components at an equivalent magnitude. The rotation is then applied to obtain the lab frame time evolution operator U_k via

$$U_k = \exp(-i2\pi f_r J_z Dt) U_k^r. \quad (40)$$

Specifically, this relationship is

$$U_k = \begin{pmatrix} e^{-i2\pi f_r Dt} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & e^{i2\pi f_r Dt} \end{pmatrix} U_k^r, \text{ for spin one, and} \quad (41)$$

$$U_k = \begin{pmatrix} e^{-i\pi f_r Dt} & 0 \\ 0 & e^{i\pi f_r Dt} \end{pmatrix} U_k^r, \text{ for spin half.} \quad (42)$$

It is a common technique in solving quantum mechanical problems to enter rotating frames, and their more abstract counterparts of interaction pictures [10]. Note that this is typically done in conjunction with the Rotating Wave Approximation (RWA), which is an assumption that the oscillatory components of f_x^r, f_y^r, f_z^r , and f_q^r on an average of many cycles do not have a large contribution to time evolution of the solution and can be ignored. In some cases, this allows for analytic solutions to the approximate quantum system to be obtained. Note that the RWA is *not* invoked in **spinsim**, as doing this would reduce the accuracy of simulation results, defeating our purpose of using a rotating frame.

8.6 Integrator architecture

The integrator in the **spinsim** package calls a `numba.cuda.jit()`ed kernel to be run on a *Cuda* capable *Nvidia* GPU in parallel, with a different thread being allocated to each of the U_k . This returns when each of the U_k have been evaluated.

The thread starts by calculating t_k and, if the rotating frame is being used, f_r . The latter is done by sampling a (`numba.cuda.jit()`ed version of a) user provided *python* function f describing how to sample the Hamiltonian. The code then loops over each integration time step dt to calculate the integration time evolution operators u_l^k .

Within the loop, the integrator enters a device function (ie a GPU subroutine, which is inline for speed) to sample $f(t)$, as well as calculate $e^{-i2\pi f_r t}$, at the sample times needed for the integration method. After this, it enters a second device function,

which makes a rotating wave transformation as needed in a third device function, before calculating g values, and finally taking the matrix exponentiation in a fourth device function. u_l^k is premultiplied to U_k^r (which is initialised to 1), and the loop continues.

When the loop has finished, if the rotating frame is being used, U_k^r is transformed to U_k as in Equation (33), and this is returned. Once all threads have executed, the state ψ_k is calculated in a (CPU) `numba.jit()`ed function from the U_k and an initial condition ψ_{init} .

8.7 Compilation of integrator

The `spinsim` integrator is constructed and compiled just in time, using `numba.cuda.jit()`. The particular device functions used are not predetermined, but are instead chosen based on user input to decide a closure. This structure has multiple advantages. First, the field function f is provided by the user as a plain python function (that must be `numba.cuda.jit()` compatible). This allows users to define f in a way that compiles and executes fast, does not put many restrictions on the form of the function, and returns the accurate results of analytic functions (compared to the errors seen in interpolation). Compiling the simulator also allows the user to set meta parameters, and choose the features they want to use, in a way that does not require experience with the `numba.cuda` library. This was especially useful for running benchmarks comparing old integration methods to the new ones, like CF4. The default settings should be optimal for most users, although tuning the values of *Cuda* meta parameters `max_registers` and `threads_per_block` could improve performance for GPUs with a differing number of registers and *Cuda* cores to the mobile *GeForce GTX 1070* mainly used in testing here. Finally, just in time compilation also allows the user to select a target device other than *Cuda* for compilation, so the simulator can run, using the same algorithm, on a multicore CPU in parallel instead of a GPU, if the user so chooses.

This functionality is interfaced through an object of class `spinsim.Simulator`. The *Cuda* kernel is defined as per the user's instructions on construction of the instance, and it is used by calling the method `spinsim.Simulator.evaluate()`, which returns a results object including the time, state, time evolution operator, and expected spin projection (that is, Bloch vector). Note that the expected spin projection is calculated as a lazy parameter if needed, rather than returned by the simulator object.

`spinsim` is designed so that a single simulator compilation can be used to execute many simulations, sweeping through a particular parameter, while not needing to be recompiled. This is done through the `sweep_parameter` argument to the user provided field function f . The user uses `sweep_parameter` to determine a variable parameter for the Hamiltonian. After compiling the `spinsim.Simulator` object, the user can set the value for `sweep_parameter` for a particular simulation as the first argument when calling the function `spinsim.Simulator.evaluate()`. For a use case example, one might want to simulate many spin systems at different locations in a magnetic field gradient in an MRI experiment. To do this they could choose to set $f_z = \text{sweep_parameter}$, which would vary the magnetic field bias in

the z direction for each experiment. This feature is demonstrated with examples in the documentation.

9 Quality control

9.1 Evaluation of accuracy

All accuracy benchmarks were run in *Cuda* mode, on the desktop computer with the *Ryzen 7 5800X* and *GeForce RTX 3080*, which from the tests in Figure 4 are the fastest CPU and GPU from the devices tested. Note that these benchmarks do not take into account the amount of time required to JIT compile simulation code for the first time (order of seconds), as we want to look at this in the limit of running many simulations while sweeping through a parameter that differs in each one.

Benchmarks were performed using `neural-sense.sim.benchmark` (where `neural-sense` [11] is the quantum sensing package that `spinsim` was written for). This simulation involves continuously driving transitions in the system for 100ms, while exposing it to a short 1ms pulsed signal that the system should be able to sense. We first wanted to test the accuracy of the different integration techniques for various integration time steps. Here we wanted to test the advantages, if any of using a Magnus based integration method. Accuracy was calculated by taking the quantum state simulation evaluations of a typical quantum sensing experiment and finding the Root Mean Squared (RMS) error from a baseline simulation run by `scipy.integrate.ivp_solve()` as part of the *SciPy* python package, via

$$\epsilon = \frac{1}{K} \sqrt{\sum_{k=0}^{K-1} \sum_{m_j=-j}^j |\psi_{k,(m_j)} - \psi_{k,(m_j)}^{\text{baseline}}|^2}, \quad (43)$$

where $j \in \{\frac{1}{2}, 1\}$ is the spin quantum number of the system.

This baseline was computed in 3.8 hours (in comparison to order of 10ths of a second these `spinsim` tests were executed in), and was also used for comparisons to other software packages. These are shown in Figures 1a, and 2a. For each of these simulations we measured the time of execution, as although it is more accurate compared to Euler integration, the CF4 method is slower for any fixed integration time step. These are shown in Figures 1b, and 2b. In all of these comparisons, errors above 10^{-3} were counted as a failed simulation, as the maximum possible error for a quantum state saturates given that it is a point on a unit complex sphere. Additionally, errors below 10^{-11} were excluded, as this was the order of magnitude of the errors in the reference simulation.

The integration techniques tested were the Magnus based CF4, as well as two Euler based sampling methods. A midpoint Euler method was chosen as the simplest (and fastest for a given integration time step) possible sampling method, whereas a Heun Euler sampling method was used as a comparison to previous versions of this code, which also used it. We benchmarked these methods both while using and not using the rotating frame option. This was done separately for spin one and spin half systems, to ensure they both yield accurate results.

From Figures 1 and 2, we find that overall, the results that `spinsim` gives are accurate to those of `SciPy`. Figures 1a, and 2a show that using the Magnus based integration method is up to 3 orders of magnitude more accurate when compared the Euler based methods. Also, using the rotating frame increased the accuracy here by 4 orders of magnitude for any individual integration method. From Figures 1b, and 2b, although the Magnus based method is slower than the midpoint Euler based method, it makes up for this in terms of its accuracy. Thus, by default `spinsim` sets the integrator to CF4, and uses the rotating frame. These can be modified using optional arguments when instantiating the `spinsim.Simulator` object.

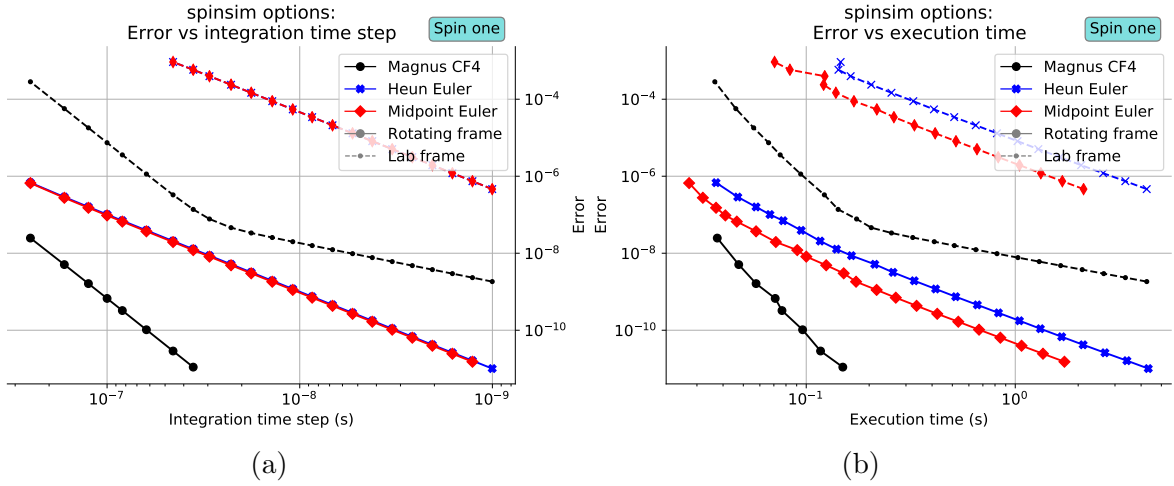


Figure 1: Speed and accuracy of the spin one options of `spinsim`. A simulation of a typical neural sensing experiment was run for every integration time step, for each of the possible integration techniques (Magnus based commutator free 4, and two Euler methods). In the simulation, transitions are continuously driven in the spin system for a duration of 100ms, and an additional small signal is injected for 1ms. Each technique was tested while both using and not using a transformation into a rotating frame. Both execution time and error were recorded for each of the simulations. Error is RMS error compared to a long running `SciPy` baseline.

9.2 Comparison to alternatives

We ran the same error (using Equation (43)) and execution time benchmarks on some alternative packages to compare `spinsim`'s performance to theirs. The packages compared were:

- `spinsim` running on the *Cuda* device, using the CF4 integrator and the rotating frame mode, which were the highest performing options when comparing `spinsim` integration methods.
- `NDSolve` from the *Mathematica* [12] software. This was chosen as it is popular with our lab group for simulating magnetometry experiments.
- `scipy.integrate.ivan_solve()` from the *python* library `SciPy` [13]. This was chosen as a generic solver from within the python ecosystem.

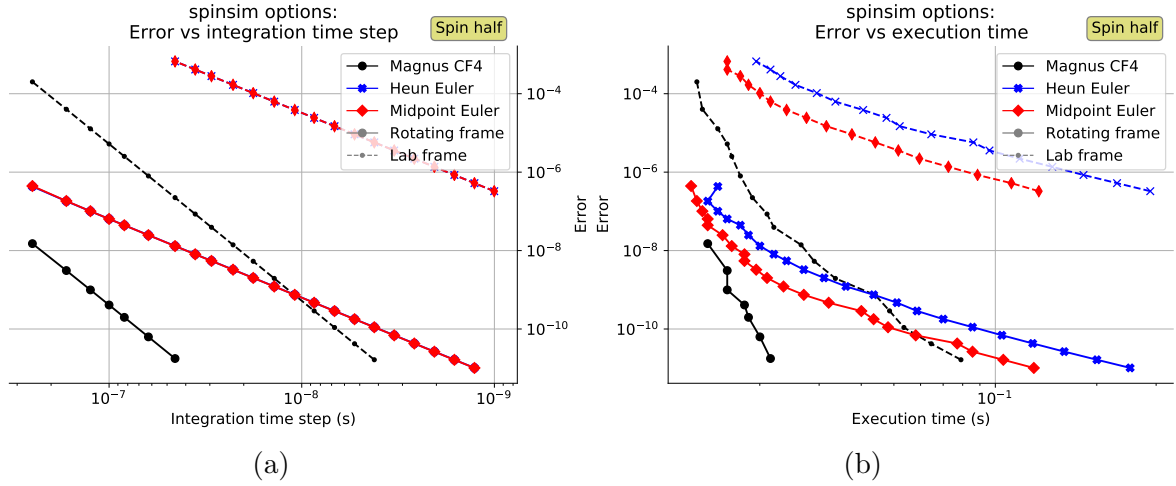


Figure 2: Speed and accuracy of the spin half options of `spinsim`. A simulation of a typical neural sensing experiment was run for every integration time step, for each of the possible integration techniques (Magnus based commutator free 4, and two Euler methods). In the simulation, transitions are continuously driven in the spin system for a duration of 100ms, and an additional small signal is injected for 1ms. Each technique was tested while both using and not using a transformation into a rotating frame. Both execution time and error were recorded for each of the simulations. Error is RMS error compared to a long running *SciPy* baseline.

- We had also planned to benchmark against `qutip.sesolve()`, a solver in the popular quantum mechanics *python* library, QuTip [14]. However, due to a known bug with the library's dependencies, this was not installable on Windows 10, the operating system being used for testing, and so benchmarks for it could not be run.

In each case, the step sizes of the alternative integrators were limited to a maximum value obtain simulation results of different accuracies. Apart from that, the integrator settings were left untouched from the default values, as a representation of what a user would experience using a generic solver for spin system problems. Similarly to with the internal `spinsim` benchmarks, the expected spin projection was evaluated in each case, but the states were compared to calculate a relative error. Also like with the internal benchmarks, we used the longest running *SciPy* simulation as a ground truth for comparison, as the accuracy of *Mathematica* plateaus at small time steps.

Each benchmark was run one simulation at a time. However, it might be possible to increase the average speed of many benchmarks from alternative packages using multithreading to run multiple benchmarks at a time. When this was attempted using *Mathematica*, the kernels crashed as the 32GiB of RAM was not enough to run them all at once. Multithreading was also not attempted using *SciPy*, due to the fact that running the full set of benchmarks of only a single simulation per integration time step consumes a day of computational time. But to be fair, both *Mathematica* and *SciPy* results are plotted with an artificial reduction in execution time by a factor of 4 and 8, which is an upper bound for the speed increase that

could be obtained by running them parallel on a 4 and 8 core processor, respectively. These are respectively shown in plots as both faded and dotted lines.

From Figure 3, for any given error tolerance, **spinsim** is over 3 orders of magnitude faster than *Mathematica*, and 4 orders of magnitude more accurate than *SciPy*. In practice, this means that an 8 minute *SciPy* simulation is reduced to 50ms, and a full week long *SciPy* batch simulation of 1000 separate systems (a realistic situation for testing quantum sensing protocols) would take less than one minute in **spinsim**.

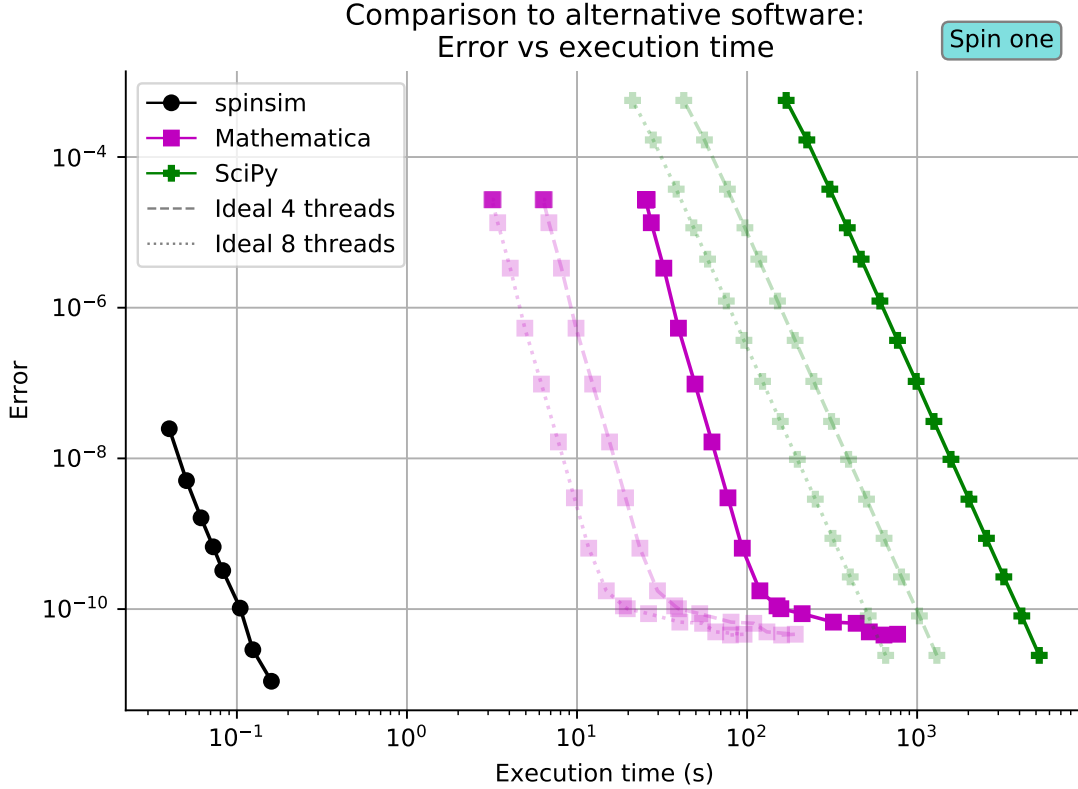


Figure 3: Speed vs accuracy of two alternative integration packages. A simulation of a typical neural sensing experiment was run for every integration time step, for each of alternative packages (`NDSolve` from *Mathematica*, and `scipy.integrate.ivan_solve()` from *SciPy*). In the simulation, transitions are continuously driven in the spin system for a duration of 100ms, and an additional small signal is injected for 1ms. Both execution time and error were recorded for each of the simulations. Error is RMS error compared to a long running *SciPy* baseline. The *Mathematica* and *SciPy* results are also shown with a speed up by a factor of 4 and 8 to represent the upper bound of hypothetical parallelisation across a 4 and 8 core CPU respectively.

9.3 Parallelisation performance

Once the algorithm behind **spinsim** was developed, we wanted to check its execution speed while running on various devices. Part of this test was also to quantify the

speed increase of parallelisation by comparing execution speeds on highly parallel devices (being GPUs), compared highly procedural devices (being CPUs). Speed benchmarks were performed using `sense.sim.benchmark`, by comparing evaluation speed of typical spin one sensing experiments on different devices. This is shown in Figure 4. The integration code was compiled by `numba` for multicore CPUs, CPUs running single threaded, and *Nvidia Cuda* compatible GPUs, and run on different models of each of them. These test devices are given in Table 1.

Table 1: Devices used in the parallelisation speed test. These devices are part of individual computers, which are separated here by horizontal lines.

Device	Type	RAM (GiB)	Cores	Cooling
Core i7-6700	Intel CPU	16	4	Air
Quadro K620	Nvidia GPU	2	384	Air
Core i7-8750H	Intel CPU	16	6	Air
GeForce GTX 1070	Nvidia GPU	8	2048	Air
Ryzen 9 5900X	AMD CPU	32	12	Air
GeForce RTX 3070	Nvidia GPU	8	5888	Air
Ryzen 7 5800X	AMD CPU	32	8	Liquid
GeForce RTX 3080	Nvidia GPU	10	8704	Air

The results shown in Figure 4 show the benefit to using parallelisation when solving a spin system problem. Moving from the 6 core *Core i7-8750H* CPU to the 12 core *Ryzen 9 5900X* CPU doubles the execution speed, as expected. Moving from a single core processor to a high end GPU increases performance by well over an order of magnitude on three of the four computers tested on. Even the low end *Quadro K620* was an improvement over the *Core i7-6700* used by the same computer. Execution speed vs number of cuda cores starts to plateau as the number of cores increases. This happens as the time it takes to transfer memory from RAM to VRAM (dedicated graphics memory), which is independent on the number of cores of the GPU, becomes a comparable to the execution time of the simulator logic. However, there is still a large improvement from using the high end *GeForce RTX 3070* to the *GeForce RTX 3080*, with the latter simulating the experiment in almost half the time of the simulated experiment duration.

Surprisingly, we found that the *Ryzen 7 5800X* 8 core CPU was able to execute the benchmark faster than the *Ryzen 9 5900X* 12 core CPU. This can be explained by the fact that the *Ryzen 7 5800X* was liquid cooled, rather than being air cooled, meaning it was likely able to boost to a higher core clock, and resist thermal throttling.

Figure 4 can be used by potential `spinsim` users for finding the relative performance for devices of varying abilities of parallelisation. We would recommend Another factor not shown in the plot is the fact that, in practice, running highly code parallel code on a CPU on a personal computer will severely limit the responsiveness of other applications, as it can utilise the entire CPU (as it should). In contrast, this does not happen when running a GPU based program, as it requires very little CPU utilisation to function. This can be convenient when running simulations on

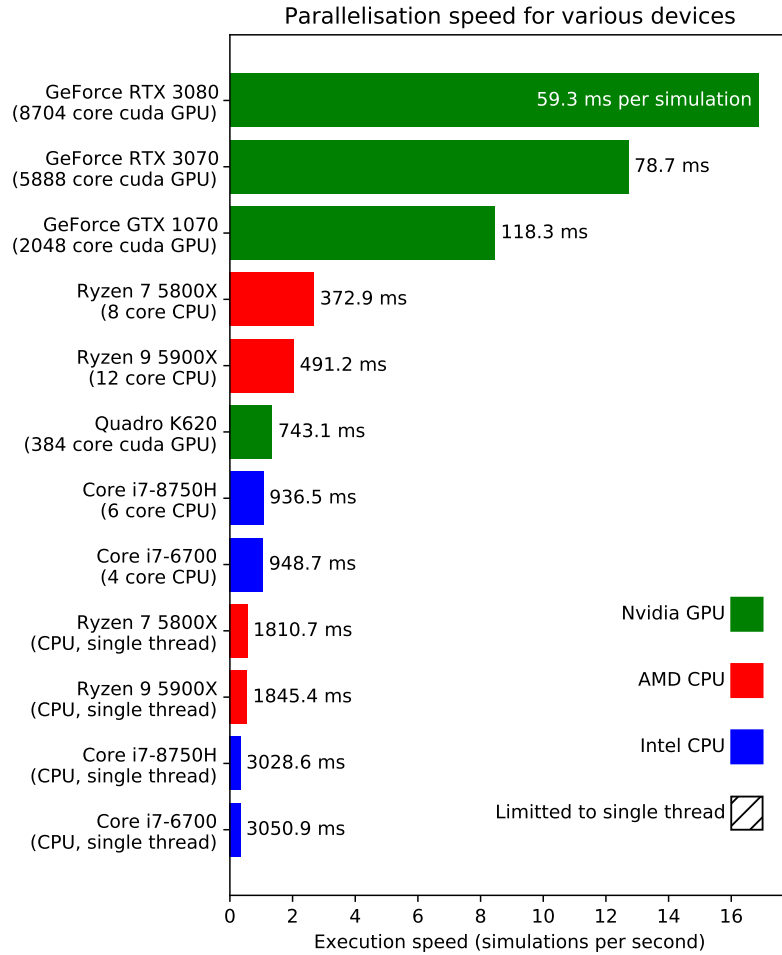


Figure 4: Evaluation speed of a simulation of a typical spin one sensing experiment on both CPUs and GPUs. Integration time step is set to 100ns. Transitions are continuously driven in the spin system for a duration of 100ms, and an additional small signal is injected for 1ms. Evaluation time is determined by an average of 100 similar simulations for each device, where each individual simulation varies in dressing amplitude (transition frequency).

a personal laptop or desktop, as other work on the computer does not have to halt while simulations are being run.

9.4 Testing

During the accuracy tests, it was confirmed that all possible modes of `spinsim` agree with a standard *SciPy* simulation up to an arbitrarily small error. The Lie Trotter matrix exponentiator was tested separately from the full system, as well as benchmarked separately. These tests and benchmarks were run as part of the `neural_sense` package. The simulator has also been used as part of the measurement protocol being developed there, and it has been tested as part of those algorithms as well.

The kernel execution was profiled thoroughly, and changes were made to optimise VRAM and register usage and transfer. This was done specifically for the development hardware of the *GeForce GTX 1070*, so one may get some performance increases by changing some GPU specific meta parameters when instantiating the `spinsim.Simulator` object.

A good way to confirm that `spinsim` is functioning properly after an installation is to run the tutorial code provided and compare the outputs. Otherwise, one can reproduce the benchmarks shown here using `neural_sense.sim.benchmark`.

10 (2) Availability

11 Operating system

Developed and tested on Windows 10. CPU functionality tested on MacOS Big Sur (note that modern Mac computers are not compatible with *Cuda* software). All packages referenced in `spinsim` are compatible with Linux, but functionality has not been tested.

12 Programming language

Python (3.7 or greater)

13 Additional system requirements

To use the (default) *Nvidia Cuda* GPU parallelisation, one needs to have a *Cuda* compatible *Nvidia* GPU [15]. For *Cuda* mode to function, one also needs to install the *Nvidia Cuda* toolkit [16]. If *Cuda* is not available on the system, the simulator will automatically parallelise over multicore CPUs instead.

14 Dependencies

numba (0.50.1 or greater)

numpy (1.19.3)

matplotlib (for example code, 3.2)

neuralsense (for benchmark code)

15 List of contributors

1. Alex Tritt

School of Physics & Astronomy, Monash University, Victoria 3800, Australia.

Primary author of the released packages.

2. Joshua Morris

School of Physics & Astronomy, Monash University, Victoria 3800, Australia.

Present address: Faculty of Physics, University of Vienna, 1010 Vienna, Austria.

Author of first version of code.

3. Joel Hockstetter

School of Physics & Astronomy, Monash University, Victoria 3800, Australia.

Present address: School of Physics, University of Sydney, NSW 2006, Australia.

Optimization and extension to spin one of first version of code.

4. Russell P. Anderson

School of Molecular Sciences, La Trobe University, PO box 199, Bendigo, Victoria 3552, Australia.

Original conception of first version of code.

5. James Saunderson

Department of Electrical and Computer Systems Engineering, Monash University, Victoria 3800, Australia.

Advice on numerical analysis.

6. Lincoln D. Turner

School of Physics & Astronomy, Monash University, Victoria 3800, Australia.

Original conception of released version of algorithm.

16 Software location:

Archive

Name: Monash Bridges

Persistent identifier: 10.26180/13285460

Licence: Apache 2.0

Publisher: Alex Tritt

Version published: 1.0.0

Date published: [dd/mm/yy](#)

Code repository

Name: GitHub

Persistent identifier: <https://github.com/alexander-tritt-monash/spinsim>

Licence: BSD 3 Clause

Date published: 18/11/20

17 Language

English.

18 (3) Reuse potential

18.1 Use potential and limitations

`spinsim` will be useful for any research group needing quick, accurate, and / or large numbers of simulations involving spin half or spin one systems. This is immediately

relevant to developing new quantum sensing protocols with spin half and spin one systems. This package is being used in the context of Bose Einstein Condensate (BEC) magnetic sensing protocol design by our lab.

This project is to be able to measure neural signals using BECs. The electrical pulses made by neurons are currently measured using electrical probes, which is intrusive and damages the cells. We instead propose to sense the small magnetic fields that these electrical currents produce. Rubidium BECs can potentially be made sensitive enough to these tiny magnetic fields that they can be measured by them. **spinsim** was written to simulate possible measurement protocols for this, showing the behaviour of the array of spin one atoms interacting with the magnetic fields of the neurons, control signals, and noise. The package is also now being used to simulate other BEC magnetometry experiments by the lab group.

Another example of spin based magnetic field sensing is the use of Nitrogen Vacancy Centres (NVCs). These are spin one structures found in diamond doped with Nitrogen atoms. This leaves a vacancy in a position adjacent to the Nitrogen atom, which pairs of electrons occupy to obtain the spin one properties. Similar to BECs, NVCs can be placed and addressed in 2D arrays in order to take many samples in one measurement. A paper was only recently released covering simulation experiments of magnetic neural pulse sensing using NVCs [17], which is something that **spinsim** could be useful for.

spinsim is designed to simulate small dimensional quantum systems, including large arrays of non-interacting spin systems. This means that it would not be able to integrate large arrays of entangled states or interacting particles. As a result, despite being fast at simulating qubits, it is inappropriate for the package to be used for quantum computing. In addition, **spinsim** is currently designed to integrate the time evolution of pure states only. This means that it may not be adequate for use in some Nuclear Magnetic Resonance (NMR) applications where relaxation [18] is important (or other kinds of simulations involving decoherence).

With these restrictions in mind, **spinsim** could be used for some simplified simulations in various areas of NMR. There are many atomic nuclei with spins of half (eg protons, Carbon 13) and, and fewer that have spins of one (eg Lithium 6, Nitrogen 14) [19], which, if relaxation and interactions between systems are not important for the application, **spinsim** could be used to simulate for spectroscopy experiments, for example. The inclusion of a quadrupole operator means that, with the same level of simplifications, **spinsim** should be able to simulate Nuclear Quadrupole Resonance (NQR) spectroscopy for spin one nuclei [20], such as Nitrogen 14, provided a suitable coordinate system is chosen. This technique measures energy level differences between levels split by electric field gradients, rather than static magnetic bias fields. Another possible use case could be for Magnetic Resonance Imaging (MRI) simulation and pulse sequence design. MRI uses measures the response of spins of an array of spin half protons to a spatially varying pulse sequence [21], which essentially just corresponds to many separate **spinsim** simulations of spins at different positions in space. While this package offers some advantages over state of the art simulators in the field [22], with its use of quantum mechanics over classical mechanics, and its absence of rotating wave approximations, its parametrised pulse sequence definitions and geometric integrator, again, the lack of interacting

particles and decoherence features are may limit its use in this area.

18.2 Support

Documentation for `spinsim` is available on *Read the Docs*. This documentation contains a thorough tutorial on how to use the package, and installation instructions. For direct support with the `spinsim` package, one can open an issue in the *github* repository. One can also use this contact to suggest extensions to the package. `spinsim` is planned to be maintained by the Monash University spinor BEC lab into the future.

19 Acknowledgements

Thank you to the Monash University School of Physics and Astronomy spinor BEC lab group, particularly Hamish Taylor and Travis Hartley, who have started using `spinsim` for their own projects and have given useful feedback of their user experience with the package.

20 Funding statement

If the software resulted from funded research please give the funder and grant number.

21 Competing interests

The authors declare that they have no competing interests.

References

- [1] Lam SK, Pitrou A, Seibert S. Numba: a LLVM-based Python JIT compiler. In: Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC - LLVM '15. Austin, Texas: ACM Press; 2015. p. 1–6. Available from: <http://dl.acm.org/citation.cfm?doid=2833157.2833162>.
- [2] Lattner C, Adve V. LLVM: a compilation framework for lifelong program analysis transformation. In: International Symposium on Code Generation and Optimization, 2004. CGO 2004.; 2004. p. 75–86.
- [3] Nickolls J, Buck I, Garland M, Skadron K. Scalable Parallel Programming with CUDA: Is CUDA the parallel programming model that application developers have been waiting for? Queue. 2008 Mar;6(2):40–53. Available from: <https://doi.org/10.1145/1365490.1365500>.
- [4] Degen CL, Reinhard F, Cappellaro P. Quantum sensing. Reviews of Modern Physics. 2017 Jul;89(3):035002. Available from: <https://link.aps.org/doi/10.1103/RevModPhys.89.035002>.
- [5] Hamley CD, Gerving CS, Hoang TM, Bookjans EM, Chapman MS. Spin-nematic squeezed vacuum in a quantum gas. Nature Physics. 2012 Apr;8(4):305–308. Number: 4 Publisher: Nature Publishing Group. Available from: <https://www.nature.com/articles/nphys2245>.

- [6] Di Y, Wang Y, Wei H. Dipole–quadrupole decomposition of two coupled spin 1 systems. *Journal of Physics A: Mathematical and Theoretical*. 2010 Jan;43(6):065303. Publisher: IOP Publishing. Available from: <https://doi.org/10.1088%2F1751-8113%2F43%2F6%2F065303>.
- [7] Blanes S, Casas F, Oteo JA, Ros J. The Magnus expansion and some of its applications. *Physics Reports*. 2009 Jan;470(5):151–238. Available from: <http://www.sciencedirect.com/science/article/pii/S0370157308004092>.
- [8] Auer N, Einkemmer L, Kandolf P, Ostermann A. Magnus integrators on multi-core CPUs and GPUs. *Computer Physics Communications*. 2018 Jul;228:115–122. Available from: <http://www.sciencedirect.com/science/article/pii/S0010465518300584>.
- [9] Moler C, Van Loan C. Nineteen Dubious Ways to Compute the Exponential of a Matrix, Twenty-Five Years Later. *SIAM Review*. 2003;45(1):3–49. Publisher: Society for Industrial and Applied Mathematics. Available from: <https://www.jstor.org/stable/25054364>.
- [10] J J Sakurai (Jun John). *Modern quantum mechanics*. Rev. ed. ed. Reading, Mass.: Addison-Wesley PubCo; 1994.
- [11] alexander-tritt monash. alexander-tritt-monash/neural-sense; 2020. Original-date: 2020-10-06T00:58:21Z. Available from: <https://github.com/alexander-tritt-monash/neural-sense>.
- [12] Wolfram Research I. *Mathematica*. Champaign, Illinois: Wolfram Research, Inc.; 2020. Available from: <https://www.wolfram.com/mathematica>.
- [13] Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*. 2020 Mar;17(3):261–272. Number: 3 Publisher: Nature Publishing Group. Available from: <https://www.nature.com/articles/s41592-019-0686-2>.
- [14] Johansson JR, Nation PD, Nori F. QuTiP 2: A Python framework for the dynamics of open quantum systems. *Computer Physics Communications*. 2013 Apr;184(4):1234–1240. Available from: <http://www.sciencedirect.com/science/article/pii/S0010465512003955>.
- [15] CUDA GPUs; 2012. Available from: <https://developer.nvidia.com/cuda-gpus>.
- [16] CUDA Toolkit; 2013. Available from: <https://developer.nvidia.com/cuda-toolkit>.
- [17] Parashar M, Saha K, Bandyopadhyay S. Axon hillock currents enable single-neuron-resolved 3D reconstruction using diamond nitrogen-vacancy magnetometry. *Communications physics*. 2020;3(1):174.

- [18] Veshtort M, Griffin RG. SPINEVOLUTION: A powerful tool for the simulation of solid and liquid state NMR experiments. *Journal of Magnetic Resonance*. 2006 Feb;178(2):248–282. Available from: <http://www.sciencedirect.com/science/article/pii/S1090780705002442>.
- [19] Fuller GH. Nuclear Spins and Moments. *Journal of Physical and Chemical Reference Data*. 1976 Oct;5(4):835–1092. Publisher: American Institute of Physics. Available from: <https://aip.scitation.org/doi/abs/10.1063/1.555544>.
- [20] Bain AD, Khasawneh M. From NQR to NMR: The complete range of quadrupole interactions. *Concepts in Magnetic Resonance Part A*. 2004;22A(2):69–78. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cmr.a.20013>. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cmr.a.20013>.
- [21] McKinnon G. The Physics of Ultrafast MRI. In: Debatin JF, McKinnon GC, editors. *Ultrafast MRI: Techniques and Applications*. Berlin, Heidelberg: Springer; 1998. p. 1–51. Available from: https://doi.org/10.1007/978-3-642-80384-0_1.
- [22] Kose R, Setoi A, Kose K. A Fast GPU-optimized 3D MRI Simulator for Arbitrary k -space Sampling. *Magnetic Resonance in Medical Sciences*. 2019;18(3):208–218.

Copyright Notice

Authors who publish with this journal agree to the following terms:

Authors retain copyright and grant the journal right of first publication with the work simultaneously licensed under a Creative Commons Attribution License that allows others to share the work with an acknowledgement of the work's authorship and initial publication in this journal.

Authors are able to enter into separate, additional contractual arrangements for the non-exclusive distribution of the journal's published version of the work (e.g., post it to an institutional repository or publish it in a book), with an acknowledgement of its initial publication in this journal.

By submitting this paper you agree to the terms of this Copyright Notice, which will apply to this submission if and when it is published by this journal.