

# Understanding the Rational Approximation of the Exponential Integrator (REXI)

Martin Schreiber <M.Schreiber@exeter.ac.uk>  
Pedro S. Peixoto <pedrosp@ime.usp.br>  
et al.

August 19, 2015

This document serves as the basis for implementing the Rational approximation of the EXponential Integrator (REXI). Here, we purely focus on the linear part of the shallow-water equations (SWE) and show the different steps to approximate solving this linear part with an exponential integrator. This paper mainly summarises previous work on REXI.

## 1 Problem formulation

We use linearised shallow water equations (SWE) with respect to a rest state with mean water depth of  $\eta_0$  and defined for perturbations of height  $\eta$  (see [1]). The linear operator ( $L$ ) may be written as

$$L(U) := \begin{pmatrix} 0 & \eta_0 \partial_x & \eta_0 \partial_y \\ g \partial_x & 0 & -f \\ g \partial_y & f & 0 \end{pmatrix} U$$

where  $U := (\eta, u, v)^T$ . Here, we neglect all non-linear terms and consider  $f$  constant (f-plane approximation). The time evolution of the PDE, with the subscript  $t$  denoting the derivative in time, is given by

$$U_t = L(U).$$

It is further worth noting, that this system describes an oscillatory system (2D wave equation), hence the operator  $L$  is hyperbolic and has imaginary eigenvalues. Linear initial value differential problems are well known to be solvable with exponential integrators for arbitrary time step sizes via

$$U(t) = e^{Lt} U(0).$$

see e.g. [4]. However, this is typically quite expensive to compute and analytic solutions only exist for some simplified system of equations, see e.g. [1]

for f-plane shallow-water equations. These exponential integrators can be approximated with rational functions and this paper is on giving insight into this approximation.

## 2 1D rational approximation

Terry et. al. [2] developed a rational approximation of the exponential integrator. First, we like to get more insight into it with a one-dimensional formulation before applying REXI to a rational approximation of a linear operator. Our main target is to find an approximation of an operator with a *complex exponential shape*, in our case  $e^{ix}$ , which (in one-dimension) is given as a function  $f(x)$ . We will end up in an approximation given by the following rational approximation:

$$e^{ix} \approx \sum_{n=-N}^N \frac{\beta_n}{ix - \alpha_n}$$

with complex coefficients  $\alpha_n$  and  $\beta_n$ . We point out that the coefficients  $\alpha_n$  will always have non zero real part, so no singularity occurs with the rational function.

### 2.1 Step A: Approximation of solution space

First, we assume that we can use Gaussian curves as basis functions for our approximation. So first we find an approximation of one of our underlying Gaussian basis function

$$\psi_h(x) := (4\pi)^{-\frac{1}{2}} e^{-x^2/(4h^2)}$$

In this formulation,  $h$  can be interpreted as the horizontal “stretching” of the basis function. Note the similarities to the Gaussian distribution, but by dropping certain parts of the vertical scaling as it is required for probability distributions. We can now approximate our function  $f(x)$  with a superposition of basis functions  $\psi_h(x)$  by

$$f(x) \approx \sum_{m=-M}^M b_m \psi_h(x + mh)$$

with  $M$  controlling the interval of approximation ( $\sim$ size of “domain of interest”) and  $h$  will be related to the accuracy of integration ( $\sim$ resolution in “domain of interest”). We choose  $h$  small enough so that the support of the Fourier transform of  $f$  is mainly localised within  $[-1/(2h), 1/(2h)]$ , i.e. almost zero outside this interval.  $M$  is chosen such that the approximation will be adequate in the interval  $|x| < Mh$ .

To compute the coefficients  $b_m$ , we rewrite the previous equation in Fourier space with

$$\frac{\hat{f}(\xi)}{\hat{\psi}_h(\xi)} = \sum_{m=-\infty}^{\infty} b_m e^{2\pi i m h \xi},$$

where the  $\hat{\cdot}$  symbols indicate the Fourier transforms of the respective functions. The  $b_m$  are now the Fourier coefficients of the series for the function  $\frac{\hat{f}(\xi)}{\hat{\psi}_h(\xi)}$  and can be calculated as (see [2], page 11),

$$b_m = h \int_{-\frac{1}{2h}}^{\frac{1}{2h}} e^{-2\pi i m h \xi} \frac{\hat{f}(\xi)}{\hat{\psi}_h(\xi)} d\xi,$$

for  $m$  integer and  $1/h$  defines the periodicity of the trigonometric basis function.

Since we are interested in approximating  $f(x) = e^{ix}$ , we can simplify the equation by using the response in frequency space  $\hat{f}(\xi) = \delta(\xi - \frac{1}{2\pi})$ , where here  $\delta$  is the Dirac distribution, and

$$b_m = h e^{-imh} \hat{\psi}_h\left(\frac{1}{2\pi}\right)^{-1}.$$

The Fourier transform of the Gaussian function is well known and given by

$$\begin{aligned} \hat{\psi}_h(\xi) &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{4\pi}} e^{-\left(\frac{x}{2h}\right)^2} e^{-2\pi i x \xi} dx \\ &= \frac{1}{\sqrt{4\pi}} \int_{-\infty}^{\infty} e^{-\left(\left(\frac{x}{2h}\right)^2 + 2\pi i x \xi + (2h\pi i \xi)^2 - (2h\pi i \xi)^2\right)} dx \\ &= \frac{1}{\sqrt{4\pi}} e^{-(2h\pi \xi)^2} \int_{-\infty}^{\infty} e^{-\left(\frac{x}{2h} + 2h\pi i \xi\right)^2} dx \\ &= \frac{1}{\sqrt{4\pi}} e^{-(2h\pi \xi)^2} \int_{-\infty}^{\infty} e^{-\left(\frac{x}{2h}\right)^2} dx \\ &= h e^{-(2h\pi \xi)^2} \end{aligned}$$

where we used that  $\int_{-\infty}^{\infty} e^{-\left(\frac{x}{2h}\right)^2} dx = h\sqrt{4\pi}$  and completed squares in the exponential term. For the case  $\xi = \frac{1}{2\pi}$ , we get

$$\hat{\psi}_h\left(\frac{1}{2\pi}\right) = h e^{-h^2}.$$

Finally, one can obtain the equation

$$b_m = h e^{-imh} \frac{1}{h e^{-h^2}} = e^{-imh} e^{h^2}$$

to compute the coefficients  $b_m$  for  $f(x) = e^{ix}$ .

## 2.2 Step B: Approximation of basis function

The second step is the approximation of the basis function  $\psi_h(x)$  itself with a rational approximation, see [5]. Our basis function is given by

$$\psi_h(x) := (4\pi)^{-\frac{1}{2}} e^{-x^2/(4h^2)}$$

and a close-to-optimal approximation of  $\psi_1(x)$  with a sum of rational functions is given by

$$\psi_1(x) \approx \operatorname{Re} \left( \sum_{l=-L}^L \frac{a_l}{ix + (\mu + il)} \right)$$

with the  $\mu$  and  $a_l$  given in [5], Table 1. We can generalise this approximation to arbitrary chosen  $h$  via

$$\psi_h(x) \approx \operatorname{Re} \left( \sum_{l=-L}^L \frac{a_l}{i \frac{x}{h} + (\mu + il)} \right).$$

The theory of how these coefficients are calculated are in [5] and will not be described here. Therefore, we assume that the coefficients  $a_l$  are given.

## 2.3 Step C: Approximation of the approximation

We then combine the approximation (B) into the approximation (A), yielding the approximation  $\tilde{f}$  for  $f$  given by

$$\begin{aligned} \tilde{f}(x) &= \sum_{m=-M}^M b_m \psi_h(x + mh) = \sum_{m=-M}^M b_m \operatorname{Re} \left( \sum_{l=-L}^L \frac{a_l}{i \frac{x+mh}{h} + (\mu + il)} \right) \\ &= \sum_{m=-M}^M b_m \sum_{l=-L}^L \operatorname{Re} \left( \frac{h a_l}{ix + h(\mu + i(m+l))} \right). \end{aligned}$$

We further like to simplify this equation and we observe, that for  $n := m + l$ , the denominator is equal. We can hence express parts of the denominator in terms of  $n := m + l$  by

$$\alpha_n := s(\mu + in). \tag{1}$$

Now, we merge the  $b_m$  and  $a_l$  coefficients and first have a look at the  $b_m$  which is complex valued. We observe the following property: Assuming that we want

to compute the real value of  $f(x)$ , only the real value of  $b_m$  has to be merged with the sum, since the imaginary component would be dropped afterwards. This allows us to move the  $Re(b_m)$  values inside the  $\sum_L$ :

$$Re(f(x)) := Re \left( \sum_{m=-M}^M \sum_{l=-L}^L \frac{Re(b_m) h a_l}{ix + h(\mu + i(m+l))} \right).$$

Now we can collect all nominators with equivalent denominator (if  $n = m + l$  and by using  $\delta$  as the Kronecker delta), yielding

$$\beta_n^{Re} := h \sum_{m=-M}^M \sum_{l=-L}^L Re(b_m) h a_l \delta(n, m+l) \quad (2)$$

for real values  $f(x)$  and

$$\beta_n^{Im} := h \sum_{m=-M}^M \sum_{l=-L}^L Im(b_m) a_l \delta(n, m+l) \quad (3)$$

for complex values of  $f(x)$ . Note, that the complexity of this operation is  $O(NL)$ , which is negligible for small  $N$  and  $L$ . This can be optimized by using  $L_1 = \max(-L, n - M)$  and  $L_2 = \min(L, n + M)$  (see [2]) and we can compute  $\beta_n^{Re}$  with

$$\beta_n^{Re} := h \sum_{k=L_1}^{L_2} Re(b_{n-k}) a_k, \quad (4)$$

and  $\beta_n^{Im}$  correspondingly.

Finally, we get the REXI approximation

$$e^{ix} \approx \sum_{n=-N}^N Re \left( \frac{\beta_n^{Re}}{ix + \alpha_n} \right) + i Re \left( \frac{\beta_n^{Im}}{ix + \alpha_n} \right).$$

### 3 REXI on linear operators

#### 3.1 Matrix exponential

We would like to apply REXI to a formulation such as

$$U(t) = e^{tL} U(0).$$

To see the relationship between the approximation of  $e^{ix}$  with  $e^{tL}$  we assume that  $L$  is skew hermitian and therefore has only purely imaginary eigenvalues, and maybe decomposed as  $\Sigma \Lambda \Sigma^H$ , yielding

$$e^{tL} = \sum_{k=0}^{\infty} \frac{tL}{k!} = \Sigma \left( \sum_{k=0}^{\infty} \frac{t\Lambda}{k!} \right) \Sigma^H = \Sigma e^{t\Lambda} \Sigma^H,$$

where we used the orthonormality of  $\Sigma$  to remove it from the summation, and

$$e^{t\Lambda} = \begin{pmatrix} \cdots & & \\ & e^{i\lambda_n t} & \\ & & \cdots \end{pmatrix},$$

where we have explicitly detached the imaginary unit from the eigenvalues, therefore  $\lambda_n$  are assumed real. Since  $e^{t\Lambda}$  is diagonal, it can be eigenvalue-wise approximated in the same way as in  $e^{ix}$ .

Some important points about the choice of  $M$  and  $h$  have to be made at this point. We know that  $e^{ix}$  is accurately approximated with REXI for the interval  $|x| < hM$ , where  $h$  is chosen small enough to obtain a good approximation in step (A), and  $M$  will define the interval size and number of approximation points. In the matrix case,  $M$  has to be chosen so that  $hM > t\bar{\lambda}$ , where  $\bar{\lambda} = \max_n |\lambda_n|$ , in order to capture all wavelengths of  $L$ . In other words,  $hM$  need to be set to capture the fastest wave. Note that if this is used as a time stepping method, with time step  $t = \tau$ , then, the larger the timestep, the larger  $M$  will be. Exact evaluations of the choices for  $h$  and  $M$  may be done based on equation (3.6) of [2].

We want to evaluate  $e^{\tau L}U(0)$  with REXI, where  $\tau$  will be a time step size and  $U(0)$  the initial condition for this time step. We will assume  $\tau$  a-priori fixed, which implies that the coefficients in REXI will not change and may be pre-computed.

Although  $L$  has imaginary eigenvalues, we wish to evaluate the  $e^{\tau L}U(0)$ , which is real valued, therefore, we will use the real approximation of  $e^{ix}$

$$\exp(\tau L) \approx \text{Re} \left( \sum_{n=-N}^N \beta_n (\tau L + \alpha_n)^{-1} \right), \quad (5)$$

where  $\beta_n$  is given by equation (2) and  $\alpha_n$  by equation (1). These coefficient may be pre-computed if  $L$  and  $\tau$  are fixed.

Note, that for debugging purpose, their *imaginary values have to cancel out*. **(TODO - PP: are you sure? MaS: No, but this is what Terry mentioned).**

Note an important property (see Sec. 3.3 in [2]). There's an anti-symmetry in the  $\alpha_i$  coefficients, which avoids computing half of the inverses,

$$\overline{(L + \alpha)^{-1}U(0)} = (L + \bar{\alpha})^{-1}U(0).$$

### 3.2 Handling $\tau$ in REXI

We reformulate the REXI approximation scheme given by

$$(\tau L + \alpha)^{-1}U(\tau) = U(0)$$

and by factoring  $\tau$  out, yielding

$$(L + \frac{\alpha}{\tau})^{-1}U(\tau)\tau^{-1} = U(0)$$

So instead of solving for  $U(\tau)$ , we are solving for  $U^\tau(\tau) := U(\tau)\tau^{-1}$  as well as  $\alpha^\tau := \frac{\alpha}{\tau}$ .

To summarize, we have to solve the system of equations given by

$$(L + \alpha^\tau)^{-1}U^\tau(\tau) = U(0) \quad (6)$$

with  $U(0)$  the initial conditions. For sake of simplicity, we stick to the formulation without the  $\tau$  notation.

### 3.3 Computing inverse of $(L + \alpha)^{-1}$

**@Pedro:** I've now also uploaded the jpeg ( IMG\_20150806\_002902.jpg ) which I've sent you 10 days ago.

This is how it is currently implemented and knowing the problems in this formulations would greatly simplify debugging. Maybe you can rearrange these equations in a form which match closely this formulations and check where the problem in this formulation (see jpeg) is?

**@Pedro:** The equations above were also derived with  $(L+\backslash\alpha)$ , not with  $(L-\backslash\alpha)$ . Sorry, but I forgot to update also the rest of this document with  $(L+\backslash\alpha)$ .

**@Martin:** Why not stick with Terry's paper standard?

**@Pedro:**  $\delta$  is already used above as the Kronecker delta and the derivative. Maybe we should avoid overutilising it :-).

**@Martin:** We should not be using  $\delta$  for derivative (where are we doing this?)

**@Pedro:** Consistency with formulations of SWE:  $H = \eta_0$ ,  $\eta_0 = \eta(0)$ ,  $\eta_1 = \eta(1)$ ,  $u_0 = u(0)$ ,  $v_0 = v(0)$ ...

**@ Martin:** I used  $\bar{\eta}$  instead of  $\eta_0$  for the mean fluid depth - needs to be changed in formulation doc

For computing the inverse, arbitrary solvers can be used. However we like to note, that  $\alpha$  is a complex number. Hence, requiring solvers with support for solving in complex space.

We wish to solve the differential problem for each time step

$$(L + \alpha)U = U(0)$$

so that  $U = (L + \alpha)^{-1}U_0$ . We will do this converting the problem into and elliptic equation.

First, lets expand the equations with the definition of  $L$ ,

$$-fv + g\eta_x + \alpha u = u_0, \quad (7)$$

$$fu + g\eta_y + \alpha v = v_0, \quad (8)$$

$$\bar{\eta}(u_x + v_y) + \alpha\eta = \eta_0. \quad (9)$$

Let  $f$  be constant (f-plane approximation),  $\delta := u_x + v_y$  be the wind divergence,  $\zeta := v_x - u_y$  be the wind (relative) vorticity and  $\Delta\eta := \eta_{xx} + \eta_{yy}$  the Laplacian of the fluid depth. We will re-write the problem in a divergence-vorticity formulation by taking 2 steps. First, sum the  $\partial_x$  of equation (7) and the  $\partial_y$  of equation (8), yielding

$$-f\zeta + g\Delta\eta + \alpha\delta = \delta_0. \quad (10)$$

Then subtract the  $\partial_y$  of equation (7) from the  $\partial_x$  of equation (8), yielding

$$f\delta + \alpha\zeta = \zeta_0. \quad (11)$$

Using equation (10) in equation (11) gives us

$$\delta = -\frac{1}{f^2 + \alpha^2} (\alpha g \Delta\eta - \alpha \delta_0 - f \zeta_0).$$

Finally, substituting  $\delta$  in equation (9), that reads  $\bar{\eta}\delta + \alpha\eta = \eta_0$ , results in

$$-\frac{\bar{\eta}}{f^2 + \alpha^2} (\alpha g \Delta\eta - \alpha \delta_0 - f \zeta_0) + \alpha\eta = \eta_0,$$

which may be simplified into the elliptic equation by multiplying by  $-\frac{f^2 + \alpha^2}{\bar{\eta}\alpha g}$

$$\Delta\eta - \kappa^2\eta = r_0 \quad (12)$$

where

$$\kappa^2 = \frac{f^2 + \alpha^2}{\bar{\eta}g}$$

and

$$r_0 = -\frac{\kappa^2}{\alpha}\eta_0 + \frac{1}{g}\delta_0 + \frac{f}{\alpha g}\zeta_0.$$

Multiplying the equation (12) by  $-g\bar{\eta}$  gives

$$((\alpha^2 + f^2) - g\bar{\eta}\Delta)\eta - \kappa^2\eta = \frac{f^2 + \alpha^2}{\alpha}\eta_0 - \bar{\eta}\delta_0 - \frac{f\bar{\eta}}{\alpha}\zeta_0$$

(This is exactly what is in the jpg)

Once  $\eta$  is calculated from any elliptic solver, we need to retrieve the velocities by solving the  $2 \times 2$  system formed by equations (7) and (8), which reads

$$A_\alpha U = U_0 - g\nabla\eta$$



with

$$A_\alpha = \begin{pmatrix} \alpha & -f \\ f & \alpha \end{pmatrix}.$$

The solution is

$$U = A_\alpha^{-1}(U_0 - g\nabla\eta)$$

where

$$A_\alpha^{-1} = \frac{1}{f^2 + \alpha^2} \begin{pmatrix} \alpha & f \\ -f & \alpha \end{pmatrix}$$

### 3.4 Interpretation of $\tau$

We like to close this section with a brief discussion of  $\tau$  by having a look on the REXI reformulation

$$(L - \frac{\alpha}{\tau})^{-1}U(\tau)\tau^{-1} = U(0)$$

We see, that for an increasing  $\tau$ , hence an integration in time over a larger time period, the poles given by  $\alpha$  are getting closer. This can possibly lead to a loss in accuracy for the data sampled by the outer poles  $\alpha_{-N}$  and  $\alpha_N$ . Therefore, the number  $N$  of poles is expected to scale linearly with the size of the coarse time step,

$$|N| \propto \tau.$$

Indeed, we saw in section 3.1 that for larger  $\tau$ ,  $M$  needs to be larger.

### 3.5 Complex solver

We are interested in solving a system of the form,

$$(L + \alpha)U = U_0,$$

where  $\alpha$  is a complex number, therefore we must allow complex solutions for  $U$ . The system can be transformed to have only real arithmetic in the following way.

First we decompose  $U$  into its real and imaginary parts,  $U = U^r + iU^i$ , and allow  $U_0$  to be decomposed in the same way. Although  $\alpha$  can be a general complex number, its real part is always constant in REXI, given by  $\mu$ . We will therefore write  $\alpha = \mu + i\alpha^i$ , and we can absorb  $\mu$  into  $L$  writing  $D = L + \mu I$ , where  $I$  is the identity matrix. Now

$$(D + i\alpha^i I)(U^r + iU^i) = U_0^r + iU_0^i,$$

$$DU^r - \alpha^i U^i + i(\alpha^i U^r + DU^i) = U_0^r + iU_0^i,$$

therefore

$$DU^r - \alpha^i U^i = U_0^r \tag{13}$$

$$\alpha^i U^r + DU^i = U_0^i, \tag{14}$$

which in matrix notation gives

$$\begin{pmatrix} D & -\alpha^i I \\ \alpha^i I & D \end{pmatrix} \begin{pmatrix} U^r \\ U^i \end{pmatrix} = \begin{pmatrix} U_0^r \\ U_0^i \end{pmatrix},$$

or,

$$\begin{pmatrix} L + \mu I & -\alpha^i I \\ \alpha^i I & L + \mu I \end{pmatrix} \begin{pmatrix} U^r \\ U^i \end{pmatrix} = \begin{pmatrix} U_0^r \\ U_0^i \end{pmatrix}.$$

A similar approach may be taken with the elliptic equation

$$\Delta\eta + \theta\eta = r_0, \quad (15)$$

where  $\eta$ ,  $\theta$  and  $r_0$  may be complex. The resulting system is given by

$$\Delta\eta^r + \theta^r\eta^r - \theta^i\eta^i = r_0^r, \quad (16)$$

$$\Delta\eta^i + \theta^i\eta^r + \theta^r\eta^i = r_0^i, \quad (17)$$

which can also be written in matrix notation and solved with arbitrary elliptic equation solvers.

## 4 Filtering

The method described in the previous section is well defined for skew hermitian  $L$ . If  $L$  is not skew hermitian, the real eigenvalues might cause the REXI to have absolute values larger than 1, which can lead to instabilities if used as time stepping method.

To ensure that the REXI is bounded by unit, a filtering process is proposed in [2]. REXI is prone to exceed unit in the neighbourhood of  $|t\lambda| \approx hM$ , therefore in the highest frequencies. The idea is to construct a rational function  $S(ix)$  that is approximately 1 in a smaller interval  $|t\lambda| < hM_0$ , with  $M_0 < M$ , and decays very fast to zero outside this interval. Then we multiply this filters function to the original REXI, which will lead to a unit bounded REXI.

Further details of how  $S(ix)$  is computed will be added later.

## 5 Bringing everything together

Using the spectral methods (e.g. in SWEET), we can directly solve the Helmholtz problem for the height in Eq. (??) and then solver for the velocity in Eqs. (??,??). Note that the Helmholtz problem is in complex space, as  $\alpha_n$  are complex. This is straightforward with spectral methods. For finite difference/element methods, the problems needs to be split into its real and imaginary parts.

Then, the problem is reduced to computing the REXI as given in Eq. (5). We like to note again, that the  $\alpha_n$  and  $\beta_n$  are independent of the system  $L$  to solve, and the number of coefficients only depends on the accuracy and the resolution.

## 6 Notes on HPC

- The terms in REXI to solve are all independent. Hence, for latency avoiding, the communication can be interleaved with computations.
- The iterative solvers are memory bound. Instead of computing  $c := a * b$  for the stencil operations, we could compute  $\vec{c} := a\vec{b}$  with  $a$  one coefficient in the stencil. This allows vectorization over  $c$  and  $b$  on accelerator cards with strided memory access.
- It is unknown which method is more efficient to solve the system of equations:
  - iterative solvers have low memory access,
  - inverting the system and storing it as a sparse matrix allows fast direct solving but can yield more memory access operations.
- Splitting the solver into real and complex number would store them consecutively in memory. This has a potential to avoid non-strided memory access and using the same SIMD operations (Just a rough idea, TODO: check if this is really the case).

## 7 Acknowledgements

Thanks to Terry for the feedback & discussions!

## References

- [1] Formulations of the shallow-water equations, M. Schreiber, P. Peixoto et al.
- [2] High-order time-parallel approximation of evolution operators, T. Haut et al.
- [3] An asymptotic parallel-in-time method for highly oscillatory PDEs, T. Haut et al.
- [4] Nineteen Dubious Ways to Compute the Exponential of a Matrix, Twenty-Five Years Later, Cleve Moler and Charles Van Loan, SIAM review
- [5] Near optimal rational approximations of large data sets, Damle, A., Beylkin, G., Haut, T. S. & Monzon