

# Finite-dimensional approximations to QG dynamics

Grooms, Julien, Watwood

Last Update: July 17, 2018

The traditional inviscid, unforced continuously-stratified QG equations are

$$\partial_t \vartheta^+ + J[\psi^+, \vartheta^+] = 0 \quad (1)$$

$$\partial_t q + J[\psi, q] + \beta v = 0 \quad (2)$$

$$\partial_t \vartheta^- + J[\psi^-, \vartheta^-] = 0 \quad (3)$$

$$\nabla_h^2 \psi + \partial_z (S(z) \partial_z \psi) = q, \quad S(z) = \frac{f_0^2}{H^2 N^2(z)} \quad (4)$$

where

$$\vartheta^\pm = S(z) \partial_z \psi \text{ at } z = 0, 1 \quad (5)$$

Notation follows Rocha et al. (JPO 2016; RYG16). I have assumed that the dimensional  $z$  goes from 0 to  $H$ , and the above equations use a nondimensional  $z$  that goes from 0 to 1.

## Energy conservation and Galerkin approximation

We want an approximate solution that has the following form

$$q_{\mathcal{N}}^G = \sum_{n=1}^{\mathcal{N}} \check{q}_n(x, y, t) p_n^q(z), \quad \psi_{\mathcal{N}}^G = \sum_{i=1}^{\mathcal{N}} \check{\psi}_i(x, y, t) p_n^\psi(z) \quad (6)$$

where  $p_n^q(z)$  and  $p_n^\psi(z)$  are a basis functions. E.g. we could use the modal expansion of RYG16, or any basis for the space of polynomials of degree  $\leq n$ , or finite elements, etc. The basis for  $q$  does not need to be the same as the basis for  $\psi$ . If we insert these above we get

$$\partial_t q_{\mathcal{N}}^G + J[\psi_{\mathcal{N}}^G, q_{\mathcal{N}}^G] + \beta \partial_x \psi_{\mathcal{N}}^G = r_{qt} \neq 0 \quad (7)$$

$$q_{\mathcal{N}}^G - (\nabla_h^2 \psi_{\mathcal{N}}^G + \partial_z (S(z) \partial_z \psi_{\mathcal{N}}^G)) = r_q \neq 0 \quad (8)$$

$$\vartheta^+ + r_\vartheta^+ = S(z) \partial_z \psi_{\mathcal{N}}^G \text{ at } z = 1, \text{ and } \vartheta^- + r_\vartheta^- = S(z) \partial_z \psi_{\mathcal{N}}^G \text{ at } z = 0. \quad (9)$$

I allow some error in the boundary conditions. We don't need to allow for errors in the evolution of  $\vartheta^\pm$  because that equation does not directly depend on the vertical structure. To put it another way, it is straightforward to enforce the following exact evolution equations:

$$\partial_t \vartheta^+ + J[\psi_{\mathcal{N}}^G, \vartheta^+] = 0 \quad (10)$$

$$\partial_t \vartheta^- + J[\psi_{\mathcal{N}}^G, \vartheta^-] = 0. \quad (11)$$

Note that because the Galerkin approximation  $\psi_{\mathcal{N}}^G$  is not equal to the true solution  $\psi$  there will be a different kind of errors in the evolution of  $\vartheta^\pm$ , but we can still enforce the above equations to hold exactly.

We want our equations to conserve energy in the form  $1/2 \int |\nabla \psi_N^G|^2 + S(z)(\partial_z \psi_N^G)^2$ . The standard approach to proving energy conservation is to multiply equation (7) for  $\partial_t q_N^G$  by  $-\psi_N^G$ , then use equation (8) to replace  $q_N^G$ , then perform multiple integrations by parts to arrive at an expression for the evolution of energy:

$$\int -\psi_N^G \partial_t (\nabla_h^2 \psi_N^G + \partial_z (S(z) \partial_z \psi_N^G) + r_q) = \int \psi_N^G (J[\psi_N^G, q_N^G] - r_{qt}) = - \int \psi_N^G r_{qt}.$$

The hope is that we can impose appropriate conditions on  $r_q$  and  $r_{qt}$  such that the energy is conserved. E.g. a standard Galerkin condition would be to require both  $r_q$  and  $r_{qt}$  to be orthogonal to the span of the basis functions.

We simplify the above expression piece by piece. First, assuming appropriate lateral boundary conditions, we have

$$\int -\psi_N^G \partial_t \nabla_h^2 \psi_N^G = \frac{1}{2} \frac{d}{dt} \int |\nabla_h \psi_N^G|^2.$$

Next

$$\int -\psi_N^G \partial_z (S(z) \partial_{zt} \psi_N^G) = - \int_x [S(z) \psi_N^G \partial_{zt} \psi_N^G]_-^+ + \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2 \quad (12)$$

$$= - \int_x [\psi_N^G (\partial_t \vartheta + \partial_t r_\vartheta)]_-^+ + \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2 \quad (13)$$

$$= - \int_x [\psi_N^G (-J[\psi_N^G, \vartheta] + \partial_t r_\vartheta)]_-^+ + \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2 \quad (14)$$

$$= - \int_x [\psi_N^G \partial_t r_\vartheta]_-^+ + \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2. \quad (15)$$

Putting it all together,

$$\frac{1}{2} \frac{d}{dt} \int |\nabla_h \psi_N^G|^2 + S(z) (\partial_z \psi_N^G)^2 = \int_x [\psi_N^G \partial_t r_\vartheta]_-^+ + \int \psi_N^G \partial_t r_q - \int \psi_N^G r_{qt}.$$

The idea is that the coefficients  $\check{q}_n$  and the boundary values  $\vartheta^\pm$  are known, but the coefficients  $\check{\psi}_n$  and  $\partial_t \check{q}_n$  are unknown. We can try to specify these coefficients to achieve the dual purposes of accuracy and energy conservation. Suppose we wish to impose the usual Galerkin conditions: First that  $r_q$  is orthogonal to the  $\mathcal{N}$  basis functions  $p_n^\psi$ , which gives us  $\mathcal{N}$  constraints on the PV inversion, then that  $r_{qt}$  is also orthogonal to the  $\mathcal{N}$  basis functions  $p_n^\psi$ , which gives us another  $\mathcal{N}$  constraints on the PV evolution. That will leave zero degrees of freedom to impose the boundary conditions, i.e. it won't be possible to set the error on the boundary  $r_\vartheta^\pm$  to zero, and we won't be able to conserve energy. This seems to be a severe obstacle to energy conservation, since we can set the terms corresponding to  $r_q$  and  $r_{qt}$  to zero in the energy conservation equation, but we appear to be unable to control the term corresponding to  $r_\vartheta^\pm$ . 'Approximation B' from RYG16 is the Tulloch & Smith 2009 model; it enforces  $r_\vartheta^\pm = 0$  and  $r_q = 0$ , but then doesn't have enough remaining degrees of freedom to set the terms corresponding to  $r_{qt}$  to zero.

On the other hand, RYG16 was able to achieve an energy-conserving Galerkin formulation; how was this possible? First, the 'modal' basis functions in that paper have  $\partial_z p_n = 0$  at the surfaces, which means that  $r_\vartheta^\pm = -\vartheta^\pm$ . As a result,  $[\int_x \psi_N^G \partial_t r_\vartheta]_-^+ = -[\int_x \psi_N^G \partial_t \vartheta]_-^+ = 0$ . In 'approximation C' from RYG16 the boundary information is included in the PV inversion via delta-function sheets of PV at the boundary. This relies on the following fact, first noted by Bretherton (1966):

The solution  $\psi$  to (4)-(5) (with inhomogeneous Neumann boundary conditions) is the same as the solution to the following PV inversion

$$q - \vartheta^+ \delta(z-1) + \vartheta^- \delta(z) = \nabla_h^2 \psi + \partial_z (S(z) \partial_z \psi) \quad (16)$$

with homogeneous Neumann boundary conditions  $\partial_z \psi = 0$ .

We will now generalize ‘Approximation C’ from RYG16 to a general basis. Suppose only that our basis  $p_n^\psi$  has  $dp_n^\psi(z)/dz = 0$  at the boundaries. Inserting our ansatz into the Bretherton PV inversion we find a new residual  $r_q$

$$q_N^G - \vartheta^+ \delta(z-1) + \vartheta^- \delta(z) = \nabla_h^2 \psi_N^G + \partial_z (S(z) \partial_z \psi_N^G) + r_q^B. \quad (17)$$

The superscript  $B$  stands for ‘Bretherton’ and emphasizes that this residual is different from  $r_q$ . Repeating the above analysis with the new definition of the residual  $r_q^B$ :

$$\int -\psi_N^G \partial_t q_N^G = \int -\psi_N^G \partial_t (\nabla_h^2 \psi_N^G + \partial_z (S(z) \partial_z \psi_N^G) + \vartheta^+ \delta(z-1) - \vartheta^- \delta(z) + r_q^B).$$

Again,

$$\int -\psi_N^G \partial_t \nabla_h^2 \psi_N^G = \frac{1}{2} \frac{d}{dt} \int |\nabla_h \psi_N^G|^2.$$

As before

$$\int -\psi_N^G \partial_z (S(z) \partial_{zt} \psi_N^G) = - \int_x [S(z) \psi_N^G \partial_{zt} \psi_N^G]_0^1 + \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2 = \frac{1}{2} \frac{d}{dt} \int S(z) (\partial_z \psi_N^G)^2. \quad (18)$$

This time we used the fact that  $\partial_z \psi_N^G = 0$  on the boundaries. We now have a new term

$$\int -\psi_N^G (\partial_t \vartheta^+ \delta(z-1) - \partial_t \vartheta^- \delta(z)) = - \int_x [\psi_N^G \partial_t \vartheta]_0^1 = 0.$$

The zero is because  $\partial_t \vartheta^\pm = -J[\psi_N^G, \vartheta^\pm]$ ; multiplying by  $\psi_N^G$  and integrating yields 0. This leaves

$$\int -\psi_N^G \partial_t q_N^G = \frac{1}{2} \frac{d}{dt} \int [|\nabla_h \psi_N^G|^2 + S(z) (\partial_z \psi_N^G)^2] - \int \psi_N^G \partial_t r_q^B.$$

The energy equation is now

$$\frac{1}{2} \frac{d}{dt} \int |\nabla_h \psi_N^G|^2 + S(z) (\partial_z \psi_N^G)^2 = \int \psi_N^G \partial_t r_q^B - \int \psi_N^G r_{qt}.$$

To conserve energy we simply impose the usual Galerkin conditions that the residuals  $r_q^B$  and  $r_{qt}$  are orthogonal to the span of the basis functions  $p_n^\psi$ .

To be precise, if the basis for  $q$  is the same as the basis for  $\psi$  then we are imposing a traditional Galerkin condition. But there is no reason why we should enforce  $\partial_z q_N^G = 0$  at the boundary, so it is probably advantageous to use a more general basis for  $q$  than for  $\psi$ . In this case we are enforcing one Galerkin condition and one Petrov-Galerkin condition. A Galerkin condition requires a residual to be orthogonal to the space in which an approximate solution is sought. The residual in the PV equation  $r_q^B$  is required to be orthogonal to the span of  $p_n^\psi$ , and when solving the PV inversion we are seeking a solution  $\psi_N^G$  in the span of the  $p_n^\psi$ , so this is a Galerkin condition. The solution we obtain for  $\psi_N^G$  by imposing the Galerkin condition is optimal in the sense that it minimizes

$$- \int (\psi - \psi_N^G) [\nabla^2 (\psi - \psi_N^G) + \partial_z (S(z) \partial_z (\psi - \psi_N^G))]$$

over all functions  $\psi_N^G$  in the subspace. RYG16 makes it look like we’re minimizing the  $L^2$  norm of the error in  $\psi$ , which is not the case.

A Petrov-Galerkin condition requires a residual to be orthogonal to a *different* subspace than the subspace in which a solution is sought. In the PV evolution equation we are seeking a solution for  $\partial_t q_N^G$  that is in the span of  $p_n^q$  but requiring the residual  $r_{qt}$  to be orthogonal to the span of the  $p_n^\psi$ , so this is a Petrov-Galerkin condition. In RYG16 the same basis is used for  $\psi$  and  $q$ , so it’s a Galerkin condition, and in that case the approximation is optimal in the  $L^2$  norm. I.e. the approximation  $\partial_t q_N^G$  is chosen to be as close as possible to  $-J[\psi_N^G, q_N^G]$  in the  $L^2$  norm. We could do that here if we wanted to by setting  $p_n^\psi = p_n^q$ .

---

### Implementation with a generic Galerkin basis

The following discussion covers how to implement this method for any Galerkin basis, i.e. it applied to the modal basis, to polynomials, to finite elements including DG, etc. For simplicity of exposition assume that the horizontal directions will be periodic and make use of the Fourier tranform so that

$$\hat{q}_n(\mathbf{k}, t), \quad \hat{\psi}_n(\mathbf{k}, t)$$

are the Fourier transforms of  $\check{q}_n$  and  $\check{\psi}_n$ , respectively, and define vectors  $\hat{\mathbf{q}}$  and  $\hat{\boldsymbol{\psi}}$  whose  $n^{\text{th}}$  elements are  $\hat{q}_n(\mathbf{k}, t)$  and  $\hat{\psi}_n(\mathbf{k}, t)$ , respectively, for  $n = 1, \dots, \mathcal{N}$ . Using the Galerkin conditions, the PV inversion takes the form

$$\mathbf{B}\hat{\mathbf{q}} - \hat{\vartheta}^+ \mathbf{p}^+ + \hat{\vartheta}^- \mathbf{p}^- = -k^2 \mathbf{M}\hat{\boldsymbol{\psi}} - \mathbf{L}\hat{\boldsymbol{\psi}}$$

where

$$(\mathbf{p}^+)_j = p_j^\psi(1), \quad (\mathbf{p}^-)_j = p_j^\psi(0),$$

$$\mathbf{B}_{ij} = \int p_i^\psi(z) p_j^q(z) dz, \quad \mathbf{M}_{ij} = \int p_i^\psi(z) p_j^\psi(z) dz, \quad \mathbf{L}_{ij} = \int S(z) (\partial_z p_i^\psi(z)) (\partial_z p_j^\psi(z)) dz$$

(in the right expression I assumed  $\partial_z p_i^\psi(z) = 0$  at the boundaries and integrated by parts.) This is very similar to the standard finite-difference approach (see Grooms and Nadeau, Fluids 2017) where  $\mathbf{M}$  and  $\mathbf{B}$  would be the identity and  $\mathbf{L}$  would be tridiagonal. In the modal basis (not polynomial) the basis functions are orthogonal so  $\mathbf{M} = \mathbf{B}$  is diagonal, and they are eigenfunctions so  $\mathbf{L}$  is just  $\kappa_n^2 \mathbf{I}$ . Regardless of which basis you choose, you should evaluate the elements of  $\mathbf{B}$ ,  $\mathbf{M}$  and  $\mathbf{L}$  to machine precision, either analytically if possible, or with adaptive quadrature or Gaussian quadrature with sufficient nodes. Notice that the matrix  $\mathbf{L}$  is a Gram matrix based on the functions  $\partial_z p_n^\psi$ . If  $p_1^\psi = 1$  is a basis function then the set of functions  $\partial_z p_n^\psi$  is linearly dependent and  $\mathbf{L}$  will be positive semi-definite. Similarly,  $\mathbf{M}$  is a Gram matrix for the functions  $p_n^\psi$ , so it is symmetric positive definite.

We need a fast way to repeatedly solve the system for  $\psi$ . We need to solve for lots of different values of  $k^2$ , as well as repeatedly in time. To rapidly solve the PV inversion we can first compute the Cholesky factorization of  $\mathbf{M}$  and store it:  $\mathbf{M} = \mathbf{G}\mathbf{G}^T$ . Then note

$$-k^2 \mathbf{M} - \mathbf{L} = -\mathbf{G}(k^2 \mathbf{I} + \mathbf{G}^{-1} \mathbf{L} \mathbf{G}^{-T}) \mathbf{G}^T.$$

The matrix  $\mathbf{G}^{-1} \mathbf{L} \mathbf{G}^{-T}$  is symmetric (and positive semi-definite) so it has an orthogonal eigenvector decomposition

$$\mathbf{G}^{-1} \mathbf{L} \mathbf{G}^{-T} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$$

where  $\mathbf{D}$  is diagonal with non-negative elements and  $\mathbf{Q}$  is an orthogonal matrix. This allows us to write

$$\mathbf{B}\hat{\mathbf{q}} + (\text{boundary terms}) = -\mathbf{G}\mathbf{Q}(k^2 \mathbf{I} + \mathbf{D})\mathbf{Q}^T \mathbf{G}^T \hat{\boldsymbol{\psi}}.$$

To obtain the solution:

- Compute  $\mathbf{Q}^T \mathbf{G}^{-1} (\mathbf{B}\hat{\mathbf{q}} + (\text{boundary terms}))$ .
- Next multiply the previous result by  $-(k^2 \mathbf{I} + \mathbf{D})^{-1}$  from the left.
- Finally go back to the original basis: multiply the previous result by  $\mathbf{G}^{-T} \mathbf{Q}$  from the left.

This is analogous to the approach taken in finite-difference approximations of the vertical direction.

The above analysis shows that there is a diagonalizing basis. This diagonalizing basis approximates the modal basis. The elements of a column of the matrix  $\mathbf{G}^{-T} \mathbf{Q}$  are the coordinates of a diagonalizing basis function with respect to the basis  $p_n^\psi$ . As  $\mathcal{N}$  increases we expect these diagonalizing basis functions to converge to the baroclinic modes of the full problem. This is essentially the same as what happens in the finite-difference approximation, where the eigenvectors of the matrix  $\mathbf{L}$  are the ‘discrete’ baroclinic modes.

The (Petrov)-Galerkin conditions define how the PV coefficients  $\check{q}_n$  should evolve. Define the vector  $\check{\mathbf{q}}$  to have elements  $\check{q}_n$ , and define the vector  $\mathbf{NL}$  to have elements

$$\mathbf{NL}_n = \int p_n^\psi(z) J[\psi_{\mathcal{N}}^G, q_{\mathcal{N}}^G] dz.$$

The PV evolution then takes the form

$$\mathbf{B} \frac{d}{dt} \check{\mathbf{q}} + \mathbf{NL} (+\beta i k_x \hat{\psi}) = 0.$$

In a fully-nonlinear implementation, one would need to repeatedly evaluate the integrals defining the elements of  $\mathbf{NL}$ . This could be done via quadrature. The LU factorization of the matrix  $\mathbf{B}$  could be computed once, then stored.

---

### Polynomial bases

The above considerations do not rely on any particular basis. We now specialize to polynomials. The set of polynomials of degree  $\leq \mathcal{N} + 1$  and with  $\partial_z p = 0$  at the boundaries is a vector space of dimension  $\mathcal{N}$ , with an infinite number of bases, any of which could be used in the above analysis. Shen (SIAM J Sci Comput 1994; section 4) gives a basis for the space of polynomials of degree  $\leq \mathcal{N} + 1$  and with  $\partial_z p_n^\psi = 0$  at the boundaries using a re-combination of Legendre polynomials. For this basis the matrix  $\mathbf{M}$  is pentadiagonal, with the further property that an even/odd permutation will bring it into a block-tridiagonal form. This matrix has condition number about  $6 \times 10^5$  for  $\mathcal{N} = 1000$ , which is quite good. The matrix  $\mathbf{L}$  depends on the stratification  $S(z) = f_0^2/N^2(z)$  and in general will be dense. (If we used a finite-element basis then  $\mathbf{L}$  would be sparse.) If we use the Shen basis for  $p_n^\psi$  and the standard Legendre basis for  $p_n^q$  then the matrix  $\mathbf{B}$  is upper-triangular with upper bandwidth 2, so we wouldn't even need to compute the LU factorization and solving for the time-tendency of PV would only take  $\mathcal{O}(\mathcal{N})$  flops. There's no benefit to using Chebyshev since energy conservation requires us to use the standard  $L^2$  inner product, and Chebyshev polynomials are not orthogonal in the standard  $L^2$  inner product.

The overall cost of the PV inversion using the Shen and Legendre bases is as follows. Pre-computing the Cholesky factor of  $\mathbf{M}$  is  $\mathcal{O}(\mathcal{N})$  because the matrix is banded. Computing the eigenvalue decomposition of  $\mathbf{G}^{-1} \mathbf{L} \mathbf{G}^{-T}$  with the basic QR algorithm should converge quickly and not cost much per iteration because the matrix is symmetric and presumably has separated eigenvalues.

- Move into the diagonalizing basis: Compute  $\mathbf{Q}^T \mathbf{G}^{-1} (\mathbf{B} \check{\mathbf{q}} + (\text{boundary terms}))$ . Multiplication by  $\mathbf{B}$  is  $\mathcal{O}(\mathcal{N})$ . The cost to invert the Cholesky is  $\mathcal{O}(\mathcal{N})$ . The cost to multiply by  $\mathbf{Q}^T$  is  $\mathcal{O}(\mathcal{N}^2)$ .
- Invert in the diagonalizing basis: Multiply the previous result by  $-(k^2 \mathbf{I} + \mathbf{D})^{-1}$  from the left. This converts from  $q$  to  $\psi$  in the diagonalizing basis. Cost is  $\mathcal{O}(\mathcal{N})$ .
- Finally go back to the original basis: multiply the previous result by  $\mathbf{G}^{-T} \mathbf{Q}$  from the left. Cost is  $\mathcal{O}(\mathcal{N}^2)$  to multiply by  $\mathbf{Q}$  and  $\mathcal{O}(\mathcal{N})$  to invert the Cholesky.

Once the requisite decompositions have been pre-computed the inversion cost is  $\mathcal{O}(\mathcal{N}^2)$ . Of course, there's no particular reason why we need to go back and forth from the Shen basis to the diagonalizing basis. We could just start with the Shen basis, then compute the diagonalizing basis, then stick with the diagonalizing basis from then on. If so, the cost to invert is just  $\mathcal{O}(\mathcal{N})$ .

Note that the integral for  $\mathbf{NL}$  is a product of three polynomials. If the basis  $p_n^q$  goes to degree  $\mathcal{N} - 1$  and the basis  $p_n^\psi$  goes to degree  $\mathcal{N} + 1$  then the product can have degree up to  $3\mathcal{N} + 1$ . These integrals can be evaluated exactly using Gauss-Legendre quadrature with  $1.5\mathcal{N} + 1$  quadrature nodes. But we need to evaluate  $\psi$  at the boundaries so that we can evolve  $\vartheta^\pm$  on the boundaries. We could use Gauss-Legendre-Lobatto quadrature instead; overall it would require  $1.5\mathcal{N} + 2$  quadrature nodes. If we used Gauss-Legendre then we would need to evaluate  $\psi$  at  $1.5\mathcal{N} + 1$  interior points *and* at the boundaries, and we would need to evaluate  $q$  at  $1.5\mathcal{N} + 1$  interior points for a total of  $3\mathcal{N} + 4$  polynomial evaluations. If we used Gauss-Legendre-Lobatto then we would need to evaluate  $\psi$  and  $q$  at  $1.5\mathcal{N} + 2$  points for a total of  $3\mathcal{N} + 4$  polynomial evaluations. So ultimately it's the same number of polynomial evaluations. Gauss-Legendre is easier than Lobatto.

The one drawback to using something other than Chebyshev polynomials is that it costs  $\mathcal{O}(\mathcal{N}^2)$  flops to evaluate the polynomial at  $\mathcal{N}$  points (vs  $\mathcal{N} \log(\mathcal{N})$  for Chebyshev). We should, as mentioned above, just use the diagonalizing basis instead of the Shen basis. If we do that then we need to compute the matrix corresponding to the map from coordinates of a polynomial in the diagonalizing basis to values of the polynomial at the quadrature nodes. Then, to move from the coordinates in the diagonalizing basis to the values at the quadrature nodes will cost  $\mathcal{O}(\mathcal{N}^2)$  and can be achieved via matrix/vector multiplication.

Overall the algorithm would be

- Start from the Shen (1994) and Legendre bases, then compute the diagonalizing basis. (The columns of the matrix  $\mathbf{G}^{-T}\mathbf{Q}$  are the coefficients [in the Shen (1994) basis] of the diagonalizing basis that approximates the baroclinic modes.)
- Construct the matrix that maps from (coordinates in the diagonalizing basis)  $\rightarrow$  (values on the Gauss-Legendre quadrature nodes)
- Evolve the system in time by updating  $q$  and  $\vartheta^\pm$ , then solving for  $\psi$ , etc.

In addition to the linear instability problem we should look at how rapidly the modes and deformation radii converge in the Galerkin vs FD methods and at the convergence of the interaction coefficients.

### Linear stability problem

The following is an exact solution of the fully-nonlinear QG equations:

$$\bar{\psi} = -\bar{u}(z)y, \quad \bar{q} = -y \frac{d}{dz} \left( S(z) \frac{d\bar{u}}{dz} \right), \quad \bar{\vartheta}^\pm = -y S(z) \frac{d\bar{u}^\pm}{dz}.$$

We can linearize the PDE about this equilibrium solution to see whether it is stable/unstable to small perturbations. Instability of this kind of equilibrium in the QG equations is one example of something called ‘baroclinic’ instability. The linearized equations are

$$\partial_t \vartheta^+ + \bar{u}^+ \partial_x \vartheta^+ + (\partial_x \psi^+) \partial_y \bar{\vartheta}^+ = 0 \quad (19)$$

$$\partial_t q + \bar{u}(z) \partial_x q + (\partial_x \psi) \partial_y \bar{q} + \beta (\partial_x \psi) = 0 \quad (20)$$

$$\partial_t \vartheta^- + \bar{u}^- \partial_x \vartheta^- + (\partial_x \psi^-) \partial_y \bar{\vartheta}^- = 0 \quad (21)$$

$$\nabla^2 \psi + \frac{d}{dz} \left( S(z) \frac{d\psi}{dz} \right) = q - \vartheta^+ \delta(z-1) + \vartheta^- \delta(z) \quad (22)$$

Coefficients don’t vary in the horizontal, so we take the Fourier transform

$$\partial_t \hat{\vartheta}^+ + ik_x \bar{u}^+ \hat{\vartheta}^+ + ik_x (\partial_y \bar{\vartheta}^+) \hat{\psi}^+ = 0 \quad (23)$$

$$\partial_t \hat{q} + ik_x \bar{u}(z) \hat{q} + ik_x (\partial_y \bar{q}) \hat{\psi} + ik_x \beta \hat{\psi} = 0 \quad (24)$$

$$\partial_t \hat{\vartheta}^- + ik_x \bar{u}^- \hat{\vartheta}^- + ik_x (\partial_y \bar{\vartheta}^-) \hat{\psi}^- = 0 \quad (25)$$

$$-(k_x^2 + k_y^2) \hat{\psi} + \frac{d}{dz} \left( S(z) \frac{d\hat{\psi}}{dz} \right) = \hat{q} - \hat{\vartheta}^+ \delta(z-1) + \hat{\vartheta}^- \delta(z). \quad (26)$$

We look for exponential growth so we try to find solutions with  $\partial_t \hat{q} = -ik_x c \hat{q}$

$$-c \hat{\vartheta}^+ + \bar{u}^+ \hat{\vartheta}^+ + (\partial_y \bar{\vartheta}^+) \hat{\psi}^+ = 0 \quad (27)$$

$$-c \hat{q} + \bar{u}(z) \hat{q} + (\partial_y \bar{q}) \hat{\psi} + \beta \hat{\psi} = 0 \quad (28)$$

$$-c \hat{\vartheta}^- + \bar{u}^- \hat{\vartheta}^- + (\partial_y \bar{\vartheta}^-) \hat{\psi}^- = 0 \quad (29)$$

$$-(k_x^2 + k_y^2) \hat{\psi} + \frac{d}{dz} \left( S(z) \frac{d\hat{\psi}}{dz} \right) = \hat{q} - \hat{\vartheta}^+ \delta(z-1) + \hat{\vartheta}^- \delta(z). \quad (30)$$

For some careful choices of  $\bar{u}(z)$  and  $N^2(z)$  the equations can be solved analytically. For example if  $\bar{u} = z$ ,  $\beta = 0$ , and  $N^2(z) = N^2$  it is the ‘Eady’ problem. More generally you have to discretize and then solve an eigenvalue problem to find  $c$ .

The previous derivation of the linear stability problem did not adhere correctly to the philosophy of ‘Approximation C’ from RYG16. The previous derivation considered the background velocity  $\bar{u}$  to be the fundamental quantity from which others are derived. The natural viewpoint, taken by Approximation C, is that  $q$  and  $\vartheta^\pm$  are fundamental and that  $\psi$  is uniquely derived from them. The following should be a correct formulation.

To use our Galerkin formulation we need to represent the equilibrium solution using our Galerkin bases. The first step would be to compute the Galerkin coefficients of  $\partial_y \bar{q}$ , which are

$$(\partial_y \bar{q})_n = \frac{\int_0^1 p_n^q(z) (\partial_y \bar{q}) dz}{\int_0^1 (p_n^q(z))^2 dz}. \quad (31)$$

(Note that the above expression assumes that the basis functions  $p_n^q$  are  $L^2$ -orthogonal.) The Eady equilibrium has  $\bar{q} = 0$ , which can be represented exactly in our Galerkin basis (and any basis). Our standard PV inversion says that we can obtain the Galerkin coefficients of  $\bar{u}$  (in the basis  $p_n^\psi$ ) by solving the system

$$\mathbf{L}\bar{u} = \mathbf{B}(\partial_y \bar{q}) - (\partial_y \bar{\vartheta}^+) \mathbf{p}^+ + (\partial_y \bar{\vartheta}^-) \mathbf{p}^-. \quad (32)$$

As noted previously, the  $\mathbf{L}$  matrix is singular. This is natural because  $\mathbf{L}$  is the discrete version of the operator  $\partial_z(S(z)\partial_z \cdot)$  with homogeneous Neumann boundary conditions, which has constant functions in its null space. The null space of  $\mathbf{L}$  is also one-dimensional and corresponds to constant functions (for a polynomial basis).

The Fredholm alternative tells us that a solution will exist whenever the right hand side is orthogonal to the cokernel, and since  $\mathbf{L}$  is symmetric the cokernel is the same as the kernel, which includes constant functions. More simply, the first row (and column) of  $\mathbf{L}$  is zero, so we need to ensure that the first entry of the RHS is also zero. The first row of the right hand side is the integral of  $\phi_0(z) = 1$  times  $\bar{q}_N^G - \bar{\vartheta}^+ \delta(z-1) + \bar{\vartheta}^- \delta(z)$ . It turns out that this is zero by construction (I will not make the argument here), though the quadrature might have something nonzero due to quadrature and/or roundoff errors. This shows that the RHS is always in the range of  $\mathbf{L}$ , so a solution always exists.

At this point we know that a solution to (32) always exists, but there are in fact an infinite number of solutions. To arrive at a unique solution we need to constrain the null space, i.e. we need to constrain the depth-independent (aka barotropic) part of  $\bar{u}$ . One natural choice would be to simply compute the first Galerkin coefficient of  $\bar{u}$  (i.e. the coefficient corresponding to the basis function  $p_1^\psi(z) = \phi_0(z) = 1$ ) using the integral definition. So the first element of the vector  $\bar{u}$  is

$$\int_0^1 p_1^\psi(z) \bar{u}(z) dz.$$

The remaining elements can be found by solving (32) ignoring the first row and column.

We next discretize (27)–(29). The streamfunction  $\hat{\psi}$  appears in these equations, but we only want evolution equations involving our basic/fundamental variables  $\hat{q}$  and  $\hat{\vartheta}^\pm$ , so we formally eliminate  $\hat{\psi}$  using

$$\hat{\psi} = -((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1} \mathbf{B}\hat{q} + \hat{\vartheta}^+ \mathbf{p}^+ - \hat{\vartheta}^- \mathbf{p}^-.$$

The vectors  $\mathbf{p}^\pm$  are

$$\mathbf{p}^+ = ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1} \mathbf{p}^+, \quad \mathbf{p}^- = ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1} \mathbf{p}^-.$$

We can now discretize the surface equations (27) and (29). These don’t require Galerkin projection, and have the form

$$\begin{aligned} (-c + \bar{u}_N^G(z=1))\hat{\vartheta}^+ + (\partial_y \bar{\vartheta}^+) \mathbf{p}^+ \cdot \hat{\psi} &= 0 \\ (-c + \bar{u}_N^G(z=1))\hat{\vartheta}^+ + (\partial_y \bar{\vartheta}^+) (\mathbf{p}^+ \cdot (-(k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1} \mathbf{B}\hat{q} + \hat{\vartheta}^+ \mathbf{p}^+ - \hat{\vartheta}^- \mathbf{p}^-)) &= 0 \end{aligned} \quad (33)$$

$$\begin{aligned} (-c + \bar{u}_N^G(z=0))\hat{\vartheta}^- + (\partial_y \bar{\vartheta}^-) \mathbf{p}^- \cdot \hat{\psi} &= 0 \\ (-c + \bar{u}_N^G(z=0))\hat{\vartheta}^- + (\partial_y \bar{\vartheta}^-) (\mathbf{p}^- \cdot (-(k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1} \mathbf{B}\hat{q} + \hat{\vartheta}^+ \mathbf{p}^+ - \hat{\vartheta}^- \mathbf{p}^-)) &= 0 \end{aligned} \quad (34)$$

Next we discretize the PV evolution equation (28) by inserting our Galerkin approximation and integrating against the basis functions  $p_n^\psi$

$$-c\mathbf{B}\hat{\mathbf{q}} + \left[ \bar{\mathbf{U}} - (\bar{\mathbf{Q}}_y + \beta\mathbf{M})((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B} \right] \hat{\mathbf{q}} + (\bar{\mathbf{Q}}_y + \beta\mathbf{M})(\hat{\vartheta}^+\boldsymbol{\psi}^+ - \hat{\vartheta}^-\boldsymbol{\psi}^-) = 0 \quad (35)$$

$$(\bar{\mathbf{U}})_{jk} = \int_0^1 p_j^\psi(z) p_k^q(z) \bar{u}_N^G(z) dz, \quad (\bar{\mathbf{Q}}_y)_{jk} = \int_0^1 (\partial_y \bar{q}_N^G) p_j^\psi(z) p_k^\psi(z) dz$$

We can now write a generalized eigenvalue problem for  $c$  as follows

$$\left[ \begin{array}{c|c|c} \bar{u}_N^G(z=1) + (\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^+) & -(\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B}) & -(\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^-) \\ \hline (\bar{\mathbf{Q}}_y + \beta\mathbf{M})\boldsymbol{\psi}^+ & \bar{\mathbf{U}} - (\bar{\mathbf{Q}}_y + \beta\mathbf{M})((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B} & -(\bar{\mathbf{Q}}_y + \beta\mathbf{M})\boldsymbol{\psi}^- \\ \hline (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^+) & -(\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B}) & \bar{u}_N^G(z=0) - (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^-) \end{array} \right] \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{q} \\ \hat{\vartheta}^- \end{pmatrix} = c \begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{B} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{q} \\ \hat{\vartheta}^- \end{pmatrix} \quad (36)$$

The linear instability analysis will in general proceed as follows

- Choose  $\bar{u}$  and  $N^2(z)$ . (and  $f_0$  and  $H$  and  $\beta$  and  $\mathcal{N}$ ) Use this to construct  $\bar{q}$  and  $\bar{\vartheta}^\pm$ . Do this step analytically.
- Use (31) and (32) to construct the Galerkin approximations  $\partial_y \bar{q}_N^G$  and  $\bar{u}_N^G$ , respectively. The code takes as input a function that evaluates  $\partial_y \bar{q}$ , and scalar values for  $\partial_y \bar{\vartheta}^\pm$ .
- Construct all the matrices  $\mathbf{B}$ ,  $\mathbf{M}$ ,  $\mathbf{L}$ ,  $\bar{\mathbf{U}}$ ,  $\bar{\mathbf{Q}}_y$
- Construct the matrices in (36) and pass the whole thing to a generalized eigenvalue solver, searching for the eigenvalue with largest imaginary part. If the imaginary part of the eigenvalue is positive then there is exponential growth with rate  $k_x$  times the imaginary part of  $c$ .

**Eady** The Eady and ‘Green’ problems simplify considerably, which is why they were considered in RYG16. In both the Eady and Green problems  $\bar{q} = 0$ , which implies  $\bar{\mathbf{Q}} = 0$ . In the Eady problem we also have  $\beta = 0$ . In the Eady problem the generalized eigenvalue problem reduces to

$$\left[ \begin{array}{c|c|c} \bar{u}_N^G(z=1) + (\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^+) & -(\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B}) & -(\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^-) \\ \hline 0 & \bar{\mathbf{U}} & 0 \\ \hline (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^+) & -(\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot ((k_x^2 + k_y^2)\mathbf{M} + \mathbf{L})^{-1}\mathbf{B}) & \bar{u}_N^G(z=0) - (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^-) \end{array} \right] \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{q} \\ \hat{\vartheta}^- \end{pmatrix} = c \begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{B} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{q} \\ \hat{\vartheta}^- \end{pmatrix} \quad (37)$$

The spectrum can be decomposed into two parts. The first part satisfies  $\bar{\mathbf{U}}\hat{\mathbf{q}} = c\mathbf{B}\bar{\mathbf{q}}$ . Once these eigenvalues and vectors are obtained, you can plug them back in to solve for  $\hat{\vartheta}^\pm$ . We expect this part of the problem to be stable, i.e. all eigenvalues  $c$  should be real, or have negative imaginary part. This can be tested just by computing the eigenvalues of the generalized eigenvalue problem  $\bar{\mathbf{U}}\hat{\mathbf{q}} = c\mathbf{B}\bar{\mathbf{q}}$ .

The second part has  $\hat{\mathbf{q}} = 0$ , leaving only the  $2 \times 2$  subsystem

$$\left[ \begin{array}{cc} \bar{u}_N^G(z=1) + (\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^+) & -(\partial_y \bar{\vartheta}^+)(\mathbf{p}^+ \cdot \boldsymbol{\psi}^-) \\ (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^+) & \bar{u}_N^G(z=0) - (\partial_y \bar{\vartheta}^-)(\mathbf{p}^- \cdot \boldsymbol{\psi}^-) \end{array} \right] \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{\vartheta}^- \end{pmatrix} = c \begin{pmatrix} \hat{\vartheta}^+ \\ \hat{\vartheta}^- \end{pmatrix}. \quad (38)$$

This is where we expect instability, and eigenvalues are much easier to compute for this  $2 \times 2$  system.



**Ocean-Charney** There are lots of other ‘canonical’ problems that we could do, like the Phillips problem which has  $\beta \neq 0$ ,  $N^2$  constant, and  $\bar{u} \propto \cos(\pi z)$ . We want a problem that will stress our method, i.e. where the solution is hard to resolve as a function of  $z$ . One way to accomplish this is with an ocean-Charney problem where the unstable modes are surface-intensified. We also want to demonstrate that we can handle non-constant stratification. One common non-constant stratification is to use an exponential profile for  $N^2(z)$ . Our ocean-Charney configuration with exponential stratification is

$$\bar{u} = \frac{1}{54} (3e^{6z}(6z-1) - 2e^6 - 1), \quad \bar{q} = -2y, \quad \beta = 1, \quad \frac{f_0^2}{H^2 N^2} = e^{-6z}, \quad \bar{\vartheta}^+ = -2y, \quad \bar{\vartheta}^- = 0.$$

There are two main questions: how rapidly does the eigenvalue converge as a function of  $\mathcal{N}$  at fixed  $k$  (e.g. for the fastest-growing mode), and how does the range of accurate growth rates grow with  $\mathcal{N}$ .

---

### Linear stability problem: Finite Difference

The standard vertical discretization method for both linear and nonlinear QG is a finite difference approximation. Details can be found almost anywhere, but there’s a treatment in Grooms & Nadeau (Fluids, 2016) that explicitly considers the treatment of surface buoyancy. The usual approach uses unequal spacing in the vertical direction with a goal of maximizing the accuracy for a fixed number of vertical levels by putting the resolution preferentially in places where it’s needed. The method has been shown rigorously to converge in the absence of surface buoyancy. The drawback is that the approximation reduces to first order, so we will only consider an equispaced vertical grid here. The main reasons why the finite difference method is ubiquitous in QG are that it is easy to implement and that it exactly conserves a discrete energy.

Let  $\Delta_z = 1/\mathcal{N}$  be the grid spacing where  $\mathcal{N}$  is the number of vertical levels. Both  $\psi$  and  $q$  are tracked at  $\mathcal{N}$  points starting at  $z_1 = \Delta_z/2$  and ending at  $z_{\mathcal{N}} = 1 - \Delta_z/2$ . The finite difference approximation to  $\nabla^2 \psi + \partial_z(S(z)\partial_z \psi)$  at an interior point  $z_k$  ( $k \neq 1, \mathcal{N}$ ) is

$$(\nabla^2 \psi + \partial_z(S(z)\partial_z \psi))|_{z=z_k} \approx \nabla^2 \psi_k + \frac{1}{\Delta_z} \left[ S_k \frac{\psi_{k+1} - \psi_k}{\Delta_z} - S_{k-1} \frac{\psi_k - \psi_{k-1}}{\Delta_z} \right] = q_k. \quad (39)$$

I’ve introduced the notation  $S_k = S(k\Delta_z)$ . At the boundaries we have the following approximations

$$(\nabla^2 \psi + \partial_z(S(z)\partial_z \psi))|_{z=z_1} \approx \nabla^2 \psi_1 + \frac{1}{\Delta_z} \left[ S_1 \frac{\psi_2 - \psi_1}{\Delta_z} - \vartheta^- \right] = q_1 \quad (40)$$

$$(\nabla^2 \psi + \partial_z(S(z)\partial_z \psi))|_{z=z_{\mathcal{N}}} \approx \nabla^2 \psi_{\mathcal{N}} + \frac{1}{\Delta_z} \left[ \vartheta^+ - S_{\mathcal{N}-1} \frac{\psi_{\mathcal{N}} - \psi_{\mathcal{N}-1}}{\Delta_z} \right] = q_{\mathcal{N}}. \quad (41)$$

As discussed in Grooms & Nadeau (2016), if we define

$$Q_1 = q_1 + \frac{\vartheta^-}{\Delta_z} = \nabla^2 \psi_1 + \frac{1}{\Delta_z} \left[ S_1 \frac{\psi_2 - \psi_1}{\Delta_z} \right], \quad (42)$$

$$Q_{\mathcal{N}} = q_{\mathcal{N}} - \frac{\vartheta^+}{\Delta_z} = \nabla^2 \psi_{\mathcal{N}} - \frac{1}{\Delta_z} \left[ S_{\mathcal{N}-1} \frac{\psi_{\mathcal{N}} - \psi_{\mathcal{N}-1}}{\Delta_z} \right] \quad (43)$$

Then the fully nonlinear system dynamics are controlled entirely by the following system

$$\partial_t Q_1 + \mathbf{J}[\psi_{\mathcal{N}}, Q_1] + \beta \partial_x \psi_1 = 0 \quad (44)$$

$$\partial_t q_k + \mathbf{J}[\psi_k, Q_k] + \beta \partial_x \psi_k = 0, \quad k = 2, \dots, \mathcal{N} - 1 \quad (45)$$

$$\partial_t Q_{\mathcal{N}} + \mathbf{J}[\psi_{\mathcal{N}}, Q_{\mathcal{N}}] + \beta \partial_x \psi_{\mathcal{N}} = 0. \quad (46)$$

The only caveat is that by evolving this system you can’t distinguish  $\vartheta^{\pm}$  or  $q_1, q_{\mathcal{N}}$ , but the dynamics of  $\psi_k$  are completely controlled by the above system: (42)–(44) for the dynamics and (37), (40), (41) for the PV inversion.

The discrete version of the linear stability problem is straightforward in the finite difference approximation. We can start with (27)–(30) and then discretize as described above. The discrete version is

$$[\bar{\mathbf{U}}_{FD} ((k_x^2 + k_y^2)\mathbf{I} + \mathbf{L}_{FD}) - (\bar{\mathbf{Q}}_{y,FD} + \beta\mathbf{I})] \vec{\psi} = c [(k_x^2 + k_y^2)\mathbf{I} + \mathbf{L}_{FD}] \vec{\psi}.$$

The matrix  $\bar{\mathbf{U}}_{FD}$  is diagonal with diagonal elements  $\bar{u}(z_k)$ . The matrix  $\mathbf{L}_{FD}$  is tridiagonal with the form

$$\mathbf{L}_{FD} = \frac{1}{\Delta_z^2} \begin{bmatrix} S_1 & -S_1 & 0 & \cdots & 0 \\ -S_1 & S_1 + S_2 & -S_2 & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ & -S_{k-1} & S_{k-1} + S_k & -S_k & \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & -S_{\mathcal{N}-1} & S_{\mathcal{N}-1} \end{bmatrix}$$

The matrix  $\bar{\mathbf{Q}}_{y,FD}$  is also diagonal. If we define a vector  $\bar{\mathbf{u}}$  whose elements are  $\bar{u}(z_k)$ , the diagonal elements of  $\bar{\mathbf{Q}}_{y,FD}$  are the elements of the vector  $\mathbf{L}_{FD}\bar{\mathbf{u}}$ .