

# Notes on computational fluid dynamics

Cesar B Rocha

June 10, 2015

This are notes I've been writing for my self study as part of the class "Computational Fluid Dynamics" (MAE 290C), taught by Prof. Juan C. del Alamo in Winter 2015. I'm also including material beyond the class syllabus.

I claim no originality to the content of these notes. In particular, I'm loosely following del Alamo's class notes, and the following books

- *Spectral methods: Fundamentals in single domains* by Canuto et al.;
- *Numerical Renaissance* by Bewley;
- *Numerical Methods for Fluid Dynamics: With applications to Geophysics* by Durran.
- *Numerical Methods for Conservation Laws* by LeVeque.

## 1 Introduction

We will study methods for solving the Navier-Stokes (NS) equations

$$\partial_t \vec{u} + \vec{u} \cdot \nabla \vec{u} = -\frac{\nabla p}{\rho} + \nu \Delta \vec{u}, \quad (1)$$

where the laplacian is

$$\Delta \stackrel{\text{def}}{=} \nabla \cdot \nabla \quad (2)$$

The momentum equation (1) is complemented by the conservation of mass

$$\nabla \cdot \vec{u} = 0. \quad (3)$$

To close the system, we will also need a thermodynamic equation, which we will introduce later.

The advective term  $\vec{u} \cdot \nabla \vec{u}$  gives a hyperbolic flavor to the NS equations. Its quadratic non-linearity typically prohibits analytic solutions. It gives rise to fascinating phenomena like turbulence. In contrast the linear viscous term  $\nu \Delta \vec{u}$  gives a parabolic flavour the NS equations.

## 2 Temporal discretizations of the NS equations

We consider a one-dimensional linear version of (1), and drop the pressure gradient term for simplicity,

$$\partial_t u + c \partial_x u = \nu \partial_{xx}^2 u, \quad (4)$$

where  $c$  is a constant speed. Fourier transforming in the  $x$ -direction we obtain

$$\partial_t \hat{u} = -\underbrace{(\mathrm{i} c k + \nu k^2)}_{\stackrel{\text{def}}{=} \lambda} \hat{u}, \quad (5)$$

where  $k$  is the wavenumber. (4) is an eigenproblem with eigenvalue  $\lambda$ . Decaying solutions have  $\text{Re}(\lambda) < 0$ , and numerical schemes must represent this behavior. Figure 1 shows the eigenvalue as a function of  $\nu$ . As  $\nu$  increases the real part of the eigenvalue becomes more negative, whereas for very small  $\nu$  the eigenvalue are almost purely imaginary. This example illustrates one of the challenges for time marching numerical schemes. These schemes should have stability regions that comprise large negative real number and the imaginary axis.

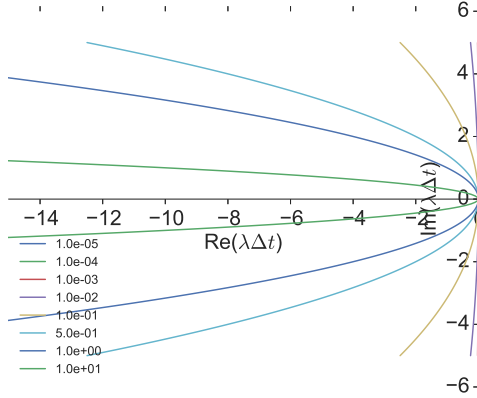


Figure 1: The eigenvalue  $\lambda$  as function of the parameter  $\nu$  with fixed speed  $c = .5$ .

One may be tempted to use implicit methods since they are unconditionally stable. However, from an implementation point of view, the implicit methods applied to the non-linear terms in (1) would lead to non-linear algebraic equations to be solved every time step, which can be very computationally costly.

For the advective terms we typically use two main schemes:

- Multi-step (e.g. Adams-Bashfort methods);
- Single-step, multiple sub-steps (e.g. Runge-Kutta methods).

Adams-Bashforth (AB) methods are more memory demanding, whereas Runge-Kutta (RK) schemes require more flops. The most efficient method depends on computer architecture. In recent computers memory, not CPU, is the main limitations, so we may use RK methods.

## 2.1 Runge-Kutta methods

Runge-Kutta (RK) are explicit, single-step, multi-sub-steps, schemes. The general form of an RK scheme for a non-linear equation

$$\partial_t u = F(u, t), \quad (6)$$

is

$$u^{n+1} = u^n + \Delta t (b_1 k_1 + \dots + b_M k_M) \quad (7)$$

$$\begin{aligned} k_1 &= F(u^n, t^n), \\ k_2 &= F(u^n + a_{21} k_1 \Delta t, t^n + c_2 \Delta t), \\ &\vdots \\ k_M &= F(u^n + a_{M1} k_1 \Delta t + \dots + a_{M,M-1} k_{M-1} \Delta t, t^n + c_M \Delta t). \end{aligned} \quad (8)$$

To systematically keep track of the coefficients, we introduce the Butcher table

Table 1: The Butcher table.

$c_2$	$a_{2,1}$	0	...	0
$c_3$	$a_{3,1}$	$a_{3,2}$	0	...
$\vdots$	$\vdots$	$\ddots$	$\ddots$	$\vdots$
$c_M$	$a_{M,1}$	$a_{M,2}$	...	$a_{M,M-1}$

### 2.1.1 RK2, the simplest example

With  $M=2$  we have

$$u^{n+1} = u^n + \Delta t (b_1 k_1 + b_2 k_2), \quad (9)$$

with

$$\begin{aligned} k_1 &= F(u^n, t^n), \\ k_2 &= F(u^n + a_{2,1} k_1 \Delta t, t^n + c_2 \Delta t). \end{aligned} \quad (10)$$

There are four parameters  $a_{2,1}$ ,  $b_1$ ,  $b_2$ , and  $c_2$ . We choose these parameters to maximize the order of accuracy of the scheme. We expand  $k_2$  about  $k_1$

$$k_2 = k_1 + a_{2,1} k_1 \Delta t \partial_u F(u^n, t^n) + c_2 \Delta t \partial_t F + \mathcal{O}(\Delta t^2). \quad (11)$$

Thus

$$u^{n+1} = u^n + \Delta t [(b_1 + b_2) F(u^n, t^n) + b_2 a_{2,1} \Delta t \partial_u F(u^n, t^n) + b_2 c_2 \Delta t \partial_t F(u^n, t^n)] \quad (12)$$

The Taylor expansion of  $u^{n+1}$  about  $u^n$  is

$$\begin{aligned} u^{n+1} &= u^n + \Delta t \partial_t u|_n + \frac{\Delta t^2}{2} \partial_{tt}^2 u|_n + \mathcal{O}(\Delta t^3) = \\ &u^n + \Delta t \partial_t u|_n + \partial_t F(u^n, t^n) + \partial_u F(u^n, t^n) F(u^n, t^n) + \mathcal{O}(\Delta t^3), \end{aligned} \quad (13)$$

where the equality follows from the governing equation (6). Matching terms in (13) and (12) gives

$$\begin{cases} b_1 + b_2 = 1; \\ b_2 c_2 = \frac{1}{2}; \\ b_2 a_{2,1} = \frac{1}{2}. \end{cases} \quad (14)$$

The system (14) is underdetermined: there are 4 unknowns but only 3 equations. Choosing  $c_2 \equiv c$  as a free parameter, we obtain

$$a_{2,1} = c, \quad b_2 = \frac{1}{2c}, \quad \text{and} \quad b_1 = 1 - \frac{1}{2c}. \quad (15)$$

The parameter  $c$  can be chosen to give the best accuracy. Typically  $\frac{1}{2} \leq c \leq 1$ . The linear stability is independent of  $c$ . For the linear problem  $F = \lambda u$ . Thus

$$u^{n+1} = u^n \left( 1 + \lambda \Delta t + \frac{(\lambda \Delta t)^2}{2} \right). \quad (16)$$

Stability requires

$$\sigma = \left| 1 + \lambda \Delta t + \frac{(\lambda \Delta t)^2}{2} \right| \leq 1. \quad (17)$$

Along the real axis, the stability the stability constrain (17) is

$$1 + \lambda_r \Delta t + \frac{(\lambda_r \Delta t)^2}{2} \leq 1, \quad (18)$$

which gives

$$-2 \leq \lambda_r \Delta t \leq 0. \quad (19)$$

Along the imaginary axis, the stability region only touches  $\lambda_i = 0$ . Figure 2 shows the linear stability region for the RK2 scheme.

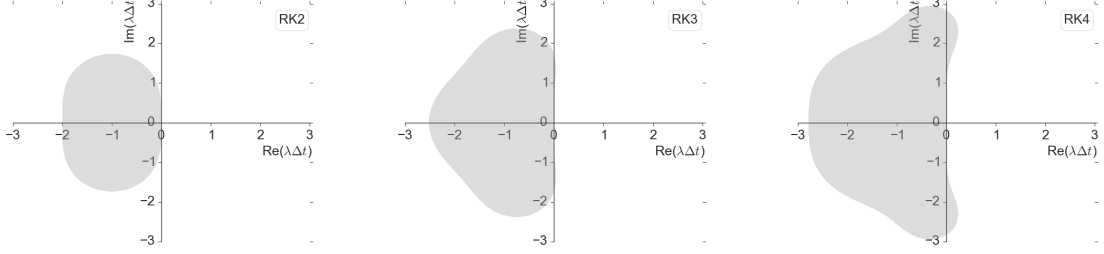


Figure 2: Stability region for Runge-Kutta schemes.

## 2.2 RK4

Derivation of higher-order RK schemes is straightforward but clumsy. A popular scheme is RK4

$$u^{n+1} = u^n + \frac{\Delta t}{6}(k_1 + 2k_2 + 2k_3 + k_4), \quad (20)$$

where

$$\begin{aligned} k_1 &= F(u^n, t^n), \\ k_2 &= F(u^n + \frac{1}{2}k_1\Delta t, t^n + \frac{1}{2}\Delta t), \\ k_3 &= F(u^n + \frac{1}{2}k_2\Delta t, t^n + \frac{1}{2}\Delta t), \\ k_4 &= F(u^n + k_3\Delta t, t^n + \Delta t). \end{aligned} \quad (21)$$

This RK4 is one of my favorite time marching schemes. It is easy to implement, and has good stability properties (see figure 2).

## 2.3 A low-storage RK3 scheme

To simplify notation we write

$$\partial_t u = \mathbf{L}(u) + \mathbf{H}(u), \quad (22)$$

where  $\mathbf{L}$  comprises the linear terms and  $\mathbf{H}$  contains the nonlinear terms. The low-storage RKW3 scheme is in the form

$$u^{n+1} = u^{**} + \Delta t[\mathbf{L}(\alpha_3 u^* + \beta_3 u^{n+1})] + \gamma_3 \mathbf{H}(u^{**}) + \xi_2 \mathbf{H}(u^*), \quad (23)$$

where

$$u^* = u^n + \Delta t[\mathbf{L}(\alpha_1 u^n + \beta_1 u^*) + \gamma_1 \mathbf{H}(u^n)], \quad (24)$$

and

$$u^{**} = u^* + \Delta t[\mathbf{L}(\alpha_2 u^* + \beta_2 u^{**}) + \gamma_2 \mathbf{H}(u^*) + \xi_1 \mathbf{H}(u^n)], \quad (25)$$

with

$$\gamma_1 = \frac{8}{15}, \quad \gamma_2 = \frac{5}{12}, \quad \text{and} \quad \gamma_3 = \frac{3}{4}, \quad (26)$$

$$\beta_1 = \frac{37}{160}, \quad \beta_2 = \frac{5}{24}, \quad \text{and} \quad \beta_3 = \frac{1}{6}, \quad (27)$$

$$\alpha_1 = \frac{29}{96}, \quad \alpha_2 = -\frac{3}{40}, \quad \text{and} \quad \alpha_3 = \frac{1}{6}, \quad (28)$$

and

$$\xi_1 = -\frac{17}{60}, \quad \text{and} \quad \xi_2 = -\frac{5}{12}. \quad (29)$$

### 3 $\theta$ -schemes

For the linear viscous terms, it is customary to use implicit schemes. The  $\theta$ -schemes avoid rapid oscillations for large  $\lambda\Delta t$ , while Crank-Nicolson's second-order. These schemes have the form

$$u^{n+1} = u^n + \Delta t[(1 - \theta)F^n + \theta F^{n+1}]. \quad (30)$$

Notice that CN is recovered with  $\theta = \frac{1}{2}$ , and implicit Euler is obtain with  $\theta = 1$ . Typically  $\frac{1}{2} \leq \theta \leq 1$ .

#### Linear stability

The growth factor is

$$\sigma = \frac{1 + (1 - \theta)\lambda\Delta t}{1 - \theta\lambda\Delta t}. \quad (31)$$

Thus

$$\lim_{\lambda\Delta t \rightarrow -\infty} \sigma = -\frac{1 - \theta}{\theta} \leq 1. \quad (32)$$

With  $\theta = \frac{3}{4}$ , the limit above is  $-\frac{1}{3}$ , a third of CN's limit.

#### Order of accuracy

From the  $\theta$ -scheme we have

$$\frac{u^{n+1}}{u^n} = 1 + \lambda\Delta t + \theta\lambda^2\Delta t^2 + \mathcal{O}(\lambda^3\Delta t^3). \quad (33)$$

Comparing with the expansion for the exponential  $e^{\lambda t}$ , we conclude that

$$\text{Err}_\theta \sim \left(\frac{1}{2} - \theta\right) (\lambda\Delta t)^2. \quad (34)$$

With  $\theta = \frac{3}{4}$  the error is

$$\text{Err}_{3/4} \sim \frac{1}{4}(\lambda\Delta t)^2, \quad (35)$$

which is half of the error of the implicit Euler scheme.

### 4 Spectral methods

Spectral methods are the golden choice for numerical methods. In problems with simple geometry we express the solution as a linear combination of basis functions

$$u(x_j, t) = \sum_{n=1}^N \hat{u}_n(t) \mathbf{b}_n(x_j), \quad (36)$$

where  $b_n(x)$  is the  $n$ 'th basis function. A few desired properties for the basis are

Table 2: A comparison of  $N^2$  versus  $5N \log_2 N$

N	$5N \log_2 N$	N
8	120	64
16	320	256
32	800	1024
64	1920	4096
128	4480	16384
256	10240	65536

### 1. Orthogonality

This property ensures accuracy. The basis formed by the set  $\{b_n\}$  usually derives from a Sturm-Liouville problem on the domain of interest. The basis functions are typically orthogonal with respect to an inner product  $\langle \rangle$ :

$$\langle b_n b_m \rangle = \alpha_{mn} \delta_{mn} , \quad (37)$$

where  $\delta$  is the Kronecker delta, and  $\alpha$  is a normalization constant, usually taken to be 1. The coefficient  $\hat{u}_n$  is then found by projection of the true solution on the  $n$ 'th basis function:

$$\hat{u}_p = \langle b_p u(x, t) \rangle . \quad (38)$$

### 2. Derivatives are known analytically

This property allows easy computation of derivatives, e.g.

$$\partial_x u(x_j, t) = \sum_{n=1}^N \hat{u}_n \partial_x b_n . \quad (39)$$

We assume that the series for the derivative converges within the domain. This is always the case when  $u(x, t)$  is periodic and sufficiently smooth within the domain. If  $u(x, t)$  has corners, discontinuities, etc differentiating term-by-term may result in a nonconvergent series. We will return to this issue later.

### 3. A computational efficient way to transform from the physical domain to the spectral domain.

We should be able to easily and efficiently compute  $u$  from  $\hat{u}_n$ , and vice-versa. An example of transform is

$$\hat{u}_n(t) = \frac{1}{N} \sum_{j=1}^{N-1} u(x_j, t) b_n(x_j) , \quad (40)$$

Evaluation of the sum in above requires  $\mathcal{O}(N)$  operations for every  $n$ . Thus the cost of the transform is  $\mathcal{O}(N^2)$ . For a spectral method to be efficient we should be able to come up with algorithms that reduce this cost, and that can be implemented in parallel. For instance, the celebrated Fast Fourier Transform (FFT) reduces this cost to  $\mathcal{O}(N \log_2 N)$ . This results in a significant reduction of number of operations for large  $N$  (see table 3). For instance, for  $N = 256$  this results in a saving of  $\sim 80\%$ ! Turbulence simulations would not be possible without the FFT.

## Fourier spectral methods

We will focus on the most used class of methods, namely the Fourier spectral methods. In Fourier methods, the set  $\{\mathbf{b}_n\}$  is formed by complex exponentials  $\{\exp(i\kappa_n x_j)\}$  where

$$\kappa_n \stackrel{\text{def}}{=} \frac{2\pi n}{L}, \quad n = -\frac{N}{2}, -\frac{N}{2} + 1, \dots, \frac{N}{2} - 2, -\frac{N}{2} - 1. \quad (41)$$

The inverse transform is defined by a truncated Fourier series

$$u(x_j) = \sum_{n=-N/2}^{N/2-1} \hat{u}_n e^{i\kappa_n x_j}, \quad (42)$$

where  $\hat{u}_n$  is the  $n$ 'th Fourier coefficient, which is generally complex-valued. Notice that the series (43) implies that  $u$  is periodic with period  $L$

$$u\left(-\frac{L}{2}\right) = u\left(\frac{L}{2}\right). \quad (43)$$

The discrete Fourier transform (DFT) is

$$\hat{u}_n = \frac{1}{N} \sum_{j=-N/2}^{N/2-1} u_j e^{-i\kappa_n x_j}. \quad (44)$$

We emphasize that one should sum only up to  $N/2 - 1$  since the  $u$  is periodic, so that the data at  $x_{N/2}$  add no information to the problem. This definition is consistent with the inverse transform (43), following from the orthogonality of the discrete complex exponentials. This property is fundamental in many calculations, so we shall formally prove it. Consider the summation

$$\sum_{j=0}^{N-1} e^{i\left(\frac{2\pi}{N}(n-s)\right)j}, \quad (45)$$

where we used  $x_j = \frac{L}{N}j$ . The summation  $S$  is a geometric series

$$s \stackrel{\text{def}}{=} 1 + r + r^2 + \dots + r^{N-1}, \quad \text{with} \quad r = e^{i\left(\frac{2\pi}{N}(n-s)\right)}. \quad (46)$$

Hence

$$rs = r + r^2 + \dots + r^N \implies s - rs = 1 - r^N, \quad (47)$$

and therefore

$$s = \frac{1 - r^N}{1 - r}, \quad r \neq 1. \quad (48)$$

Now

$$r^N = e^{i2\pi(n-s)} = 1, \quad \text{if} \quad n \neq s + mN, \quad (49)$$

where  $m$  is an integer. Thus

$$\sum_{j=0}^{N-1} e^{i\left(\frac{2\pi}{N}(n-s)\right)j} = \begin{cases} N & : \quad n = s + mN, \quad m = 0 \pm 1 \pm 2 \pm 3 \dots; \\ 0 & : \quad \text{otherwise.} \end{cases} \quad (50)$$

## Properties of the DFT

The discrete Fourier pair (43)-(44) satisfy the desired properties discussed above. This is the reason for the popularity of Fourier methods. The set  $\{\exp(i\kappa_n x_j)\}$  forms a complete orthogonal basis on any interval of size  $L$  (e.g.  $[0, L]$ ,  $[-L/2, L/2]$ , etc). Moreover, there exists the FFT algorithm performs the DFT very efficiently. In particular, the FFTW implementation in Fortran and C is very popular. A summary of important properties is

1. As mentioned above the complex exponential are orthogonal with respect to the simplest inner product

$$\langle e^{i\kappa_n x} e^{i\kappa_m x} \rangle \stackrel{\text{def}}{=} \frac{1}{L} \int_0^L e^{i\kappa_n x} e^{i\kappa_m x} dx = \delta_{mn}. \quad (51)$$

2. The function  $u$  is periodic with period  $L$

$$u(0) = \sum_{n=-N/2}^{N/2-1} \hat{u}_n e^{i\kappa_n \times 0} = \sum_{n=-N/2}^{N/2-1} \hat{u}_n e^{i\kappa_n \times L} = u(L) \quad (52)$$

3. The Fourier coefficients are periodic with period  $N$

$$\hat{u}_{n+N} = \frac{1}{N} \sum_{j=-N/2}^{N/2-1} u_j e^{-i \frac{2\pi}{L} (n+N) \frac{jL}{N}} = \frac{1}{N} \sum_{j=-N/2}^{N/2-1} u_j e^{-i 2\pi \frac{n}{N} j} = \hat{u}_n. \quad (53)$$

4. The zeroth mode represents the average of the function. With  $\kappa_0 = 0$ , we have

$$\hat{u}_0 = \frac{1}{N} \sum_{j=-N/2}^N u_j. \quad (54)$$

Note that some implementations of the DFT do not divide by  $N$  in the definition (43). This is simply a matter of definition; one must always check DocStrings and algorithm manuals to make sure which convention is being used.

## Truncation error, aliasing, and other potential problems

The definition of the inverse DFT assumes that  $u$  has nice properties such as continuity of the function and its derivatives. If the function is discontinuous, then Gibbs oscillation might be introduced. Gibbs oscillation, which amounts for a 10% error near the discontinuities. The rapid Gibbs oscillations can be significantly reduced by using filters.

A potential problem in spectral methods is aliasing of unresolved high-frequency modes. These modes alias back into low-frequency modes. For linear problems this is not an issue provided that the initial conditions do not contain unresolved high-frequency content. For non-linear problems, such as the NS equations, high-frequency modes are introduced every time step.

## Parseval's theorem

The Parseval's relation essentially says that the total energy in physical space is equal to the total energy in Fourier space

$$\sum_{j=0}^{N-1} u_j u_j^* = N \sum_{n=-N/2}^{N/2-1} |\hat{u}_n|^2. \quad (55)$$



The proof of this relationship is simple. We introduce the inverse DFT definition (43) into (55)

$$\begin{aligned} \sum_{j=0}^{N-1} u_j u_j^* &= \sum_{j=0}^{N-1} \left( \sum_{n=-N/2}^{N/2-1} \hat{u}_n e^{i\kappa_n x_j} \right) \left( \sum_{m=-N/2}^{N/2-1} \hat{u}_m^* e^{-i\kappa_m x_j} \right) \\ &= \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{u}_m^* \underbrace{\left( \sum_{j=0}^{N-1} e^{i(\kappa_n - \kappa_m) x_j} \right)}_{=N\delta_{nm}}. \end{aligned} \quad (56)$$

### Truncation error

To study the truncation error associated with a finite  $N$  in (43) we consider the true solution

$$u_e = \sum_{n=-\infty}^{\infty} \tilde{u}_n e^{i\kappa_n x_j}, \quad (57)$$

where the Fourier coefficient is

$$\tilde{u}_n = \frac{1}{L} \int_0^L u_e(x) e^{i\kappa_n x} dx. \quad (58)$$

Notice that, by aliasing, we have

$$\hat{u}_n = \tilde{u}_n + \sum_{m=-\infty}^{\infty} \tilde{u}_{n+mN}. \quad (59)$$

Thus

$$u_j = \sum_{n=-N/2}^{N/2-1} \left( \tilde{u}_n + \sum_{m=-\infty}^{\infty} \tilde{u}_{n+mN} \right) e^{i\kappa_n x_j}, \quad (60)$$

and the  $L_2$ -norm of the error is

$$\frac{1}{L} \int_0^L \left| u_e(x) - \sum_{n=-N/2}^{N/2-1} \hat{u}_n e^{i\kappa_n x_j} \right|^2 dx = \underbrace{\sum_{|n|>N/2} |\tilde{u}_n|^2}_{\text{truncation error}} + \underbrace{\sum_{-n=N/2}^{N/2-1} \left| \sum_{m=-\infty}^{\infty} \tilde{u}_{n+mN} \right|^2}_{\text{aliasing error}}. \quad (61)$$

### Convergence of the inverse DFT

Integrating (58) by parts we obtain

$$\begin{aligned} \tilde{u}_n &= \frac{1}{L} \int_0^L u_e(x) e^{i\kappa_n x} dx = -\frac{1}{L} \int_0^L \frac{u_e'(x) e^{i\kappa_n x}}{i\kappa_n} dx \\ &= \frac{1}{L} \int_0^L \frac{u_e''(x) e^{i\kappa_n x}}{(i\kappa_n)^2} dx. \end{aligned} \quad (62)$$

We can continue the process above as long as the derivate are smooth and periodic. Hence

$$|\tilde{u}_n|^2 \lesssim (\Delta x)^{-2p}, \quad (63)$$

for all  $p \in \mathbb{N}$ . Thus for smooth functions in physical space, spectral methods beats the most accurate finite differences scheme. Spectral accuracy is also know as exponential accuracy. A rule of thumb is that for a given data spacing  $\Delta x$ , a spectral scheme with  $3\Delta x$  does a better job than second-order finite-differences scheme.

### Gibbs phenomenon

Consider the step function

$$S_e(x) = \begin{cases} 1, & 0 \leq x < \pi; \\ 0, & \pi \leq x < 2\pi. \end{cases} \quad (64)$$

A truncated Fourier series representation is

$$S_T(x) = \sum_{n=-N/2}^{N/2} \tilde{S}_n e^{-inx}, \quad (65)$$

where the Fourier coefficients are

$$\tilde{S}_n = \frac{1}{2\pi} \int_0^\pi S_e(x) e^{inx} dx. \quad (66)$$

Thus

$$S_T(x) = \sum_{n=-N/2}^{N/2} \left( \frac{1}{2\pi} \int_0^\pi S_e(\xi) e^{-in\xi} d\xi \right) e^{inx} = \frac{1}{2\pi} \int_0^\pi \underbrace{S_e(\xi)}_{=1} \underbrace{\left( \sum_{n=-N/2}^{N/2} e^{in(x-\xi)} \right)}_{\stackrel{\text{def}}{=} D_N} d\xi, \quad (67)$$

where  $D_N$  is the Dirichlet Kernel

$$D_N(z) = \begin{cases} \sin\left(\frac{N+1}{2}z\right) : & \text{if } z \neq 2\pi j; \\ N+1 : & \text{if } z = 2\pi j. \end{cases} \quad (68)$$

Figure 3 shows the Dirichlet kernel for truncation indices  $N$  from 1 through 11. As  $N$  increases the kernel distributionally approaches a delta function at  $z = 0$ .

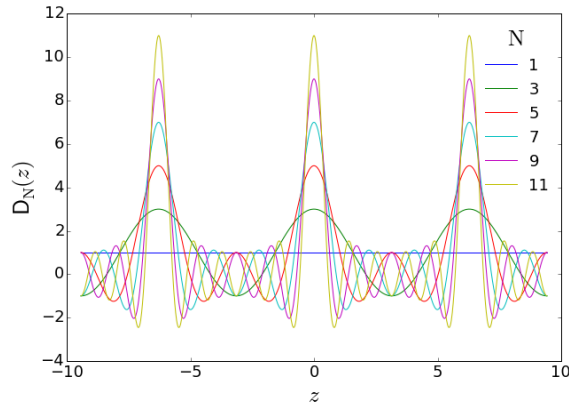


Figure 3: The Dirichlet kernel (68) for various truncation indices  $N$ .

We can rewrite the truncated series (67) by changing variables with  $y \stackrel{\text{def}}{=} x - \xi$

$$S_T(x) = \frac{1}{2\pi} \int_{x-\pi}^x D_N(y) dy = \frac{1}{2\pi} \left[ \int_0^x D_N(y) dy + \underbrace{\int_{-\pi}^0 D_N(y) dy}_{\approx \pi} + \underbrace{\int_{x-\pi}^\pi D_N(y) dy}_{<<1} \right]. \quad (69)$$

We have

$$\int_0^x D_N(y) dy \approx 1.08949\pi, \quad \text{for } x \text{ small.} \quad (70)$$

So that in the vicinity of the discontinuity at  $x = 0$  we have

$$S_T(x) \approx 1.09. \quad (71)$$

That is, the truncated series  $S_T$  overshoots the exact value by about 9%. This phenomenon is the Gibbs oscillation.

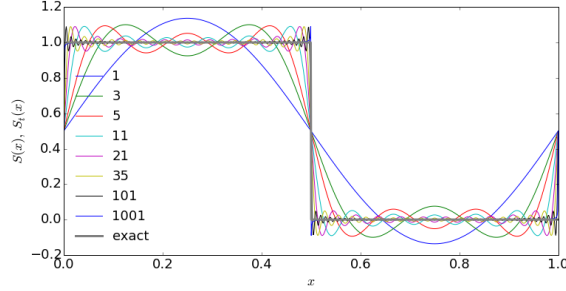


Figure 4: The representation of a step function by Fourier series under various truncation. Notice the Gibbs phenomenon near the discontinuity.

## Filters

To remove undesired Gibbs oscillations that may arise from discontinuities present in the initial conditions or developed by the equation (e.g., Burger's equation) we can use filters. The filtering process can operate in two levels

### 1. Level of equation

We introduce an additional term, e.g., artificial viscosity or hyperviscosity, that will filter out highwavenumber variability. For instance, in the one-dimensional Burger equation, we can add a fourth-order hyperviscosity

$$\partial_t u + u \partial_x u = D \partial_{x^4}^4 u. \quad (72)$$

In Fourier space the solution to the linear hyperviscous term is  $\exp(-k^4 Dt)$ , which makes clear its low-pass character.

### 2. Level of solution

Once the solution is obtained, one can filter out highwavenumber oscillations using selective filters. A celebrated example is the use of the Lanczos  $\sigma$ -factors. In this approach, one substitute the original Fourier series by

$$u_j = \sum_{n=-N/2}^{N/2-1} \sigma_n \hat{u}_n e^{ik_n x_j}, \quad (73)$$

where  $\sigma_n$  is a filter coefficient in the form

$$\sigma_n = \sigma(\theta), \quad \theta \stackrel{\text{def}}{=} \frac{2\pi n}{N}. \quad (74)$$

The factor  $\sigma(\theta)$  must satisfy the properties

- (a) Reality;
- (b) Symmetry about  $\theta = 0$ ;
- (c)  $p - 1$  continuous differentiable for  $p \geq 1$ ;
- (d)  $\sigma(\theta) = 0$ , if  $|\theta| \geq \pi$ ;  
 $\sigma(0) = 1$  (do not change the average);;  
 $\sigma^{(j)} = 0$ , for  $1 \leq j \leq p - 1$ .

Some examples of filter in the level of equation include

- Lanczos filter  
 $\sigma(0) = 1$ ;  
 $\sigma(\theta) = \frac{\sin \theta}{\theta}$ .
- Raised cosine  
 $\sigma(\theta) = (1 + \cos \theta)/2$ .
- Sharpened raised cosine  
 $\sigma(\theta) = \sigma_0^4(\theta) [35 - 84\sigma_0(\theta) + 70\sigma_0(\theta)^2 - 20\sigma_0(\theta)^3]$ ,  
where  $\sigma_0$  is the raised cosine filter.
- Exponential filter of order  $p$  (generalized viscosity)  
 $\sigma(\theta) = e^{-\alpha\theta^p}$ ,  
where  $\alpha$  is an arbitrary parameter.

Figure 5 shows the step function truncated series with  $N = 9$ , and various filtered solutions.

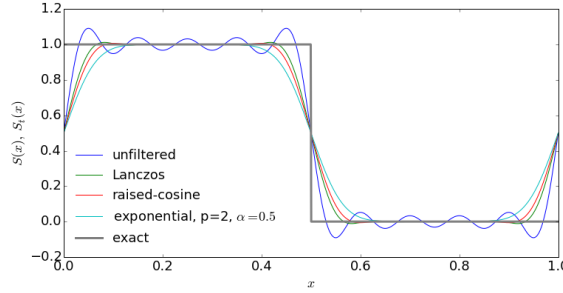


Figure 5: The representation of a step function by Fourier series with  $N = 9$  and various filtered solutions.

## Aliasing

The calculation of Fourier transform of nonlinear terms of terms with nonconstant coefficients is very expensive. Consider the two variables

$$u_e(x) = \sum_{-N/2}^{N/2} \tilde{u}_n e^{inx}, \quad \text{and} \quad v_e(x) = \sum_{-N/2}^{N/2} \tilde{v}_n e^{inx}, \quad (75)$$

and their product

$$w_e(x) = u_e(x)v_e(x) = \sum_{-N/2}^{N/2} \sum_{-N/2}^{N/2} \tilde{u}_n \tilde{v}_m e^{i(n+m)x}, \quad \text{and} \quad (76)$$

Hence

$$\begin{aligned} \tilde{w}_k &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \sum_{n=-N/2}^{N/2} \sum_{m=-N/2}^{N/2} \tilde{u}_n \tilde{v}_m e^{i(n+m-k)x} dx = \sum_{n=-N/2}^{N/2} \sum_{m=-N/2}^{N/2} \tilde{u}_n \tilde{v}_m \underbrace{\frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i(n+m-k)x} dx}_{= \begin{cases} 1 : & \text{if } m+n=k; \\ 0 : & \text{otherwise.} \end{cases}} \\ & \end{aligned} \quad (77)$$

Thus for series with  $N$  truncation index, the multiplication in physical space is equivalent to convolution in Fourier space. Computation of convolution sums requires  $\mathcal{O}(N^2)$  operations. The alternative is to use the so-called pseudo-spectral methods that requires  $\mathcal{O}(N \log_2 N)$ . In pseudo-spectral methods we compute the products in physical space, and then transform the products to Fourier space to march the system forward, then transform the variables back to physical space to compute the products, and so on. The problem with this approach is that it introduces aliasing. To see that, note that the discrete Fourier transform of (76) is

$$\begin{aligned} \hat{w}_k &= \frac{1}{N} \sum_{j=-N/2}^{N/2-1} \left( \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m e^{i(m+n)x_j} \right) e^{-ikx_j} \\ &= \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m \frac{1}{N} \sum_{j=-N/2}^{N/2-1} e^{i(m+n-k)x_j}. \end{aligned} \quad (78)$$

Now recall that

$$\frac{1}{N} \sum_{j=-N/2}^{N/2-1} e^{i(m+n-k)x_j} = \begin{cases} 1 : & \text{if } m+n=k+pN, p=0,1,2,\dots; \\ 0 : & \text{otherwise.} \end{cases} \quad (79)$$

Thus

$$\hat{w}_k = \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m + \sum \sum_{m+n=k \pm N} \hat{u}_n \hat{v}_m + \sum \sum_{m+n=k \pm 2N} \hat{u}_n \hat{v}_m \dots \quad (80)$$

The first term on the right of (80) is the discrete Fourier transform of the product  $uv$ . The double sums represent aliasing terms. These terms are due to unresolved high frequency modes that project back into the lower frequency modes.

For quadratic nonlinearity we only need to consider the first term on the right of (80). The extrema alias frequencies are

$$(m+n)_{\max} = N-2, \quad \text{and} \quad (m+n)_{\min} = -N. \quad (81)$$

Thus we only need to consider frequencies from  $-3N/2$  through  $3N/2$ .

## Dealiasing

There are two main techniques used to remove the aliased frequencies due to nonlinearities.

### Zero padding

We use discrete Fourier series for  $u$ ,  $v$ , and  $w$  with  $M \geq 3N/2$ , and set  $\hat{u}_n = \hat{v}_m = 0$  for  $m \geq N/2$  before going to physical domain. We also set  $\hat{w}_n$  for  $n \geq N/2$  after going back to Fourier space.

Notice that the worst case scenario is  $n = m = -N/2$ :

$$-\frac{N}{2} - \frac{N}{2} = \frac{N}{2} \pm M \Rightarrow M = \frac{3}{2}N. \quad (82)$$

The computational cost of zero padding is  $\mathcal{O}(3/2)^n$ , where  $n$  is the number of dimensions. Thus zero padding can be quite expensive, particularly in two- or three-dimensions.

### Phase shift

Consider the product of  $u$  and  $v$  computed at the original grid

$$w(x) = u(x)v(x), \quad (83)$$

and the products computed in a grid shifted by  $\xi$

$$w(x + \Delta x) = u(x + \xi)v(x + \xi). \quad (84)$$

The DFT of (83) is

$$\hat{w}_{o,k} = \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m + \sum_{m+n=k \pm N} \hat{u}_n \hat{v}_m, \quad (85)$$

whereas the DFT of (84) is

$$\hat{w}_{\xi,k} = \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m + e^{\pm iN\xi} \sum_{m+n=k \pm N} \hat{u}_n \hat{v}_m. \quad (86)$$

The idea is to combine (85) and (86) in such a way as to eliminate the aliasing terms

$$\hat{w}_k = \frac{\hat{w}_{o,k} + \hat{w}_{\xi,k}}{2} = \hat{w}_{\xi,k} = \sum_{n=-N/2}^{N/2-1} \sum_{m=-N/2}^{N/2-1} \hat{u}_n \hat{v}_m + \frac{1 + e^{\pm iN\xi}}{2} \sum_{m+n=k \pm N} \hat{u}_n \hat{v}_m. \quad (87)$$

Hence

$$e^{\pm iN\xi} = -1 \Rightarrow \xi N = \pm \pi \Rightarrow \xi = \frac{\pi}{N} = \frac{\Delta x}{2}. \quad (88)$$

Thus the algorithm for phase shifting is

1. Perform the pseudo-spectral without shift;
2. Shift the Fourier coefficients  $\hat{u}_k$  and  $\hat{v}_k$  in Fourier space by multiplying by  $\exp(iN\Delta x/2)$ . Then perform the pseudo-spectral calculation for this shifted coefficients;
3. Take the average of (1) and (2), and march the system forward.

The cost of the phase shifting is  $N \log_2 N$ , and therefore keeps the same order of operations as the aliased calculation. It may be advantageous if  $N$  is too large, since it does not require more memory than the original problem.

### Zero padding vs. phase shifting in one dimension

It is informative to compare the relative cost of the two dealiasing methods described above. Suppose the original grid has  $N$  elements. Hence zero padding requires computation of DFT on  $M = \frac{3}{2}N$  grids. Thus one cycle of the pseudo-spectral calculation requires  $15M \log_2 M = \frac{45}{2}N \log_2 N$  operations. The phase-shifting dealiasing method requires DFT computations on  $N$  grid but the pseudo-spectral cycle must be repeated twice (in one dimension). Thus the total number operations required for dealiasing using phase shifting is  $30N \log_2 N$ . Thus in one dimension phase shifting requires  $\frac{12}{5}$  more operations than zero padding. The latter, however, requires more memory. Thus, the question of the relative efficiency of these two methods may depend on the computer one wants to perform the calculations.

Table 3: Number of aliasing terms for quadratic nonlinearities in pseudo-spectral methods from 1 through three dimensions.

	1D	2D	3D
single	1	2	3
double	0	1	3
triple	0	0	1

### Dealiasing in higher dimensions

Real world problems are typically two- of three-dimensional. There is a qualitative difference between aliasing in one and two dimensions because for dimensions larger than one various combinations between wavenumbers produce aliasing. For instance if the original signal is

$$s = \sin x + \cos y, \quad (89)$$

then a quadratic product produces the terms

$$\sin^2 x = \frac{1 - \cos 2x}{2}, \quad (90)$$

$$\cos^2 y = \frac{1 + \cos 2y}{2}, \quad (91)$$

and also the cross-terms

$$2 \sin x \cos y = \sin 2x. \quad (92)$$

Thus aliasing in higher dimensions is richer.

Consider the representation of the discrete Fourier transform of  $\hat{u}_{\vec{n}}$

$$\hat{u}_{\vec{n}} = \sum_{j_1=0}^{N_1-1} \sum_{j_2=0}^{N_2-1} u_{j_1, j_2} \exp -i[k_{n_1} x_{j_1} + k_{n_2} x_{j_2}] \equiv \sum_{\vec{j}=0}^{\vec{N}-1} u_{\vec{j}} e^{i k_{\vec{n}} \cdot \vec{x}_{\vec{j}}}, \quad (93)$$

where

$$\vec{n} = (n_1, n_2), \quad \vec{j} = (j_1, j_2), \quad \text{and} \quad \vec{x} = x_{j_1} \vec{e}_1 + x_{j_2} \vec{e}_2, \quad (94)$$

with  $\vec{e}_1$  and  $\vec{e}_2$  the unit vector of the two-dimensional Cartesian grid.

Analogously to the one-dimensional case, the aliasing error is obtained by analyzing the DFT of the product of two variables  $u$  and  $v$ . We obtain

$$\sum_{\vec{m}+\vec{n}=\vec{p}} \tilde{u}_{\vec{m}} \tilde{v}_{\vec{n}} = \underbrace{\sum_{\vec{m}+\vec{n}=\vec{p}} \hat{u}_{\vec{m}} \hat{v}_{\vec{n}}}_{\text{Truncated series}} + \underbrace{\sum_{\vec{m}+\vec{n}=\vec{p} \pm N_1 \vec{e}_1} \hat{u}_{\vec{m}} \hat{v}_{\vec{n}} + \sum_{\vec{m}+\vec{n}=\vec{p} \pm N_2 \vec{e}_2} \hat{u}_{\vec{m}} \hat{v}_{\vec{n}}}_{\text{Single alias errors}} + \underbrace{\sum_{\vec{m}+\vec{n}=\vec{p} \pm N_1 \vec{e}_1 \pm N_2 \vec{e}_2} \hat{u}_{\vec{m}} \hat{v}_{\vec{n}}}_{\text{Double alias error}}. \quad (95)$$

Table ?? shows the number of aliasing term for pseudo-spectral problems in 1 through 3 dimensions.

### Zero padding in two dimensions

The zero padding approach discussed above in one dimension easily generalizes to higher dimensions. That is, we simply apply the 3/2 rule to each dimension. The memory cost increase is  $(\frac{3}{2})^p$  where  $p$  is the number of dimensions. Thus in higher dimensions zero padding may significantly slow down the code.

Table 4: Phase shift operations to zero out aliasing terms.

$\xi_1$	0	$\Delta x/2$	0	$\Delta x/2$
$\xi_2$	0	0	$\Delta y/2$	$\Delta y/2$
AE <sub>10</sub>	1	-1	1	-1
AE <sub>01</sub>	1	1	-1	-1
AE <sub>11</sub>	1	-1	-1	1

### Phase shifting in two dimensions

Phase shifting becomes trickier in higher dimensions. While phase shifting in one direction, it introduces error in other directions. Thus a combination of phase shifts are necessary to cancel all errors. Table 4 shows the combination of shifts that zeros out the aliasing terms. In two dimensions, phase shifts requires four evaluations of the pseudo-spectral cycle.

### Hybrid dealiasing schemes

The idea of the hybrid dealiasing method is to combine a polyhedral truncation with phase shift. The polyhedral truncation is an incomplete zero padding scheme. In particular, one realizes that the corner of the wavenumber domain are insignificant when it comes to studying flows with isotropic statistics. Thus only  $|k_{N_1} \pm k_{N_2}| < \frac{3}{2}$  is kept; the other wavenumbers are zeroed out. Similarly, in three dimensions we use the spherical truncation  $|k_{N_1}^2 + k_{N_2}^2 + k_{N_3}^2| < \frac{4}{9}N^2$ . This truncation removes double and triple aliasing terms (and somewhat more mores). After this truncation, one only needs to three evaluations of the pseudo-spectral cycle: one in the original grid and two with shifts in each direction to cancel out the single aliasing errors. Similarly, only 4 evaluations are needed in three dimensions.

### Some practical aspects of FFT routines in Fortran

We use FFT in the West (FFTW; for details see [www.fftw.org/documentation](http://www.fftw.org/documentation)). The application of the one dimensional transform is straight forward and follows a common intuition. In higher dimensions, however, the output is given in nonstandard form. In particular, we note that, for a real signal we have

$$\begin{aligned} 2D &\rightarrow \text{CFT} [\text{CRFT}(u)] ; \\ 3D &\rightarrow \text{CFT} \{ \text{CFT} [\text{CRFT}(u)] \} , \end{aligned} \quad (96)$$

where CFT represents the complex Fourier transform and CRFT represents the complex real Fourier transform. The latter takes advantage of the Hermitian symmetric of the Fourier coefficients to perform the calculation faster and requires less storage.

There are two ways of calling the FFTW routine:

- **Out-of-place**

The input array is not destroyed. That is, the routine returns an output array different than the input array. For a one dimensional real input with  $N$  points, the output has  $\frac{N+1}{2}$  non-redundant elements.

- **In-place**

The routine overwrites the input with the output, using the same memory allocation. Because the input is real, care has to be taken to interpret the output. In one dimension, if the input goes is distributed as  $0, 1, \dots, N-2, N-1$ . Then the output is  $R_0, R_1, \dots, R_{N-1}, I_{N/2-1}, \dots, I_2, I_1$ , where  $R$  stands for the real part and  $I$  for the imaginary part of the Fourier coefficients.



Things get more complicated in higher dimensions. So, please consult the manual for an accurate description.

## 5 Working with incompressible equations

The incompressible Navier-Stokes equations are

$$\partial_t \vec{u} + \vec{u} \cdot \nabla \vec{u} = -\frac{\nabla p}{\rho} + \nu \Delta \vec{u}, \quad (97)$$

and

$$\nabla \cdot \vec{u} = 0. \quad (98)$$

### 5.1 Fractional step method

In this approach, we split every time integration step in two substeps. In the first substep, we perform a preliminary estimation of the velocity  $\mathbf{u}^*$  neglecting the pressure term altogether. Using an implicit Euler scheme for the linear terms and explicit Euler for the nonlinear terms, we obtain

$$\frac{\mathbf{u}^* - \mathbf{u}^n}{\Delta t} - \nu \Delta \mathbf{u}^* = \mathbf{N}(\mathbf{u}^n), \quad \text{on the volume } \Omega, \quad (99)$$

with non slip boundary conditions

$$\mathbf{u}^* = 0, \quad \text{on the surface } \partial\Omega. \quad (100)$$

Notice that these boundary conditions are necessary in this substeps because it contains the higher derivative viscous terms. The second substep, the projection step, deals with the pressure term

$$\frac{\mathbf{u}^{n+1} - \mathbf{u}^*}{\Delta t} + \nabla p^{n+1} = 0, \quad (101)$$

where  $p$  is the dynamic pressure, and imposes the nondivergence constraint

$$\nabla \cdot \mathbf{u}^{n+1} = 0, \quad (102)$$

and no-normal flow boundary conditions

$$\mathbf{u}^{n+1} \cdot \hat{\mathbf{n}}, \quad \text{on the surface } \partial\Omega. \quad (103)$$

Taking the divergence of the second equation (101) we obtain an elliptic problem for the pressure

$$\Delta p^{n+1} = \frac{\nabla \cdot \mathbf{u}^*}{\Delta t}, \quad (104)$$

Using the divergence theorem, we obtain

$$\frac{\mathbf{u}^{n+1} \cdot \hat{\mathbf{n}}}{\Delta t} + \frac{\partial p}{\partial \mathbf{n}} = 0, \quad \text{on the surface } \partial\Omega. \quad (105)$$

That is, we impose no-slip on  $\mathbf{u}^*$  and no-normal flow on  $\mathbf{u}^*$  and  $\mathbf{u}^{n+1}$ . There is an  $\mathcal{O}(\Delta t)$  error in the no-slip boundary condition on  $\mathbf{u}^{n+1}$ . The final single time-step marching scheme has the form

$$\mathbf{u}^{n+1} = \mathbf{u}^* - \nabla p^{n+1} \Delta t. \quad (106)$$

## 5.2 Pressure correction methods

This scheme is also based on two substeps. In the first substep we solve the problem implicitly for pressure

$$\frac{\mathbf{u}^* - \mathbf{u}^n}{\Delta t} - \nu \Delta \mathbf{u}^* = \mathbf{N}(\mathbf{u}^n) - \nabla p^n, \quad (107)$$

with no-slip boundary conditions

$$\mathbf{u}^* = 0, \quad \text{on the surface} \quad \partial\Omega. \quad (108)$$

The second substep then makes a pressure correction

$$\frac{\mathbf{u}^{n+1} - \mathbf{u}^*}{\Delta t} + \nabla \cdot \phi = 0, \quad (109)$$

by imposing the nondivergence constraint on  $\mathbf{u}^{n+1}$ , which implies

$$\Delta \phi = \frac{\nabla \cdot \mathbf{u}^*}{\Delta t}, \quad \text{on the volume} \quad \Omega, \quad (110)$$

subject to no-slip boundary conditions i.e.,

$$\frac{\partial \phi}{\partial n} = 0, \quad \text{on the surface} \quad \partial\Omega. \quad (111)$$

Thus the corrected pressure is

$$p^{n+1} = p^n + \phi - \nu \nabla \cdot \mathbf{u}^*. \quad (112)$$

For this scheme, the no-slip boundary conditions are affected by a  $\Delta t^2$  error term.

## 6 Pressure correction with RKW3

Following the notation used above for the RKW3 scheme, the first substep of the pressure correction method has the form

$$\frac{\mathbf{u}^* - \mathbf{u}^n}{\Delta t} = -\gamma_1 \mathbf{H}^n + \frac{\alpha_1}{\text{Re}} \mathbf{L} \mathbf{u}^* - (\alpha_1 + \beta_1) \mathbf{G} p^n, \quad (113)$$

$$\frac{\mathbf{u}^2 - \mathbf{u}^*}{\Delta t} = -(\alpha_1 \beta_1) \mathbf{G} \phi, \quad (114)$$

$$\mathbf{D} \cdot \mathbf{u}^* = 0 \implies \mathbf{L} \phi^1 = \frac{\mathbf{D} \mathbf{u}^*}{(\alpha_1 + \beta_1) \Delta t}, \quad \text{on the volume} \quad \Omega, \quad (115)$$

and therefore

$$\partial_n \phi^1 = 0, \quad \text{on the surface} \quad \partial\Omega. \quad (116)$$

The pressure correction has the form

$$p^2 = p^n + \phi^1 - \frac{1}{\text{Re}} \frac{\beta_1}{\alpha_1 + \beta_1} \mathbf{D} \cdot \mathbf{u}^*. \quad (117)$$

The second substep is

$$\frac{\mathbf{u}^{**} - \mathbf{u}^2}{\Delta t} - \gamma_2 \mathbf{H}^2 - \xi_2 \mathbf{H}^n + \frac{\alpha_2}{\text{Re}} \mathbf{L} \mathbf{u}^2 + \frac{\beta_2}{\text{Re}} \mathbf{L} \mathbf{u}^* - (\alpha_2 + \beta_2) \mathbf{G} p^2, \quad (118)$$

$$\frac{\mathbf{u}^3 - \mathbf{u}^{**}}{\Delta t} = -(\alpha_2 + \beta_2) \mathbf{G} p^2, \quad (119)$$

$$\mathbf{D} \cdot \mathbf{u}^3 = 0 \implies \mathbf{L}\phi^2 = \frac{\mathbf{D}\mathbf{u}^{**}}{(\alpha_2 + \beta_2)\Delta t}. \quad (120)$$

The pressure correction in this substep is

$$p^3 = p^2 + \phi^2 - \frac{\beta_2}{(\alpha_2 + \beta_2)\text{Re}} \mathbf{D} \cdot \mathbf{u}^{**}. \quad (121)$$

The third substep is

$$\frac{\mathbf{u}^{***} - \mathbf{u}^3}{\Delta t} - \gamma_3 \mathbf{H}^3 - \xi_3 \mathbf{H}^2 + \frac{\alpha_3}{\text{Re}} \mathbf{L}\mathbf{u}^3 + \frac{\beta_3}{\text{Re}} \mathbf{L}\mathbf{u}^{***} - (\alpha_3 + \beta_3) \mathbf{G}p^3, \quad (122)$$

$$\frac{\mathbf{u}^{n+1} - \mathbf{u}^{***}}{\Delta t} = -(\alpha_3 + \beta_3) \mathbf{G}\phi^3, \quad (123)$$

$$\mathbf{D} \cdot \mathbf{u}^3 = 0 \implies \mathbf{L}\phi^3 = \frac{\mathbf{D}\mathbf{u}^{***}}{(\alpha_3 + \beta_3)\Delta t}. \quad (124)$$

Finally, the pressure at step  $n + 1$  is

$$p^{n+1} = p^3 + \phi^3 - \frac{1}{\text{Re}} \frac{\beta_3}{\alpha_3 + \beta_3} \mathbf{D}\mathbf{u}^{***}. \quad (125)$$

## 7 Solving large algebraic systems

As we have seen, working with incompressible Navier-Stokes equations give rise to a Poisson problem to be solve every time step. The discrete version of this problem is a large system of algebraic equations, whose solution is computationally costly. In this section, we explore some of the main iterative approaches to solve such systems.

### 7.1 Relaxation

The problem to solve is

$$\Delta\phi = f, \quad (126)$$

where  $f$  is a given function of space. In the relaxation method, we add time dependence to the problem

$$\partial_t\phi = \Delta\phi - f, \quad (127)$$

and look for the steady solution. Because we are not interested in the transients, we simply step this system forward using an implicit scheme, given an initial guess. Because the true solution is steady, this is equivalent of solving for the residual due to the initial guess

$$\partial_t r = \Delta r. \quad (128)$$

In Fourier space we have

$$\hat{r}_k = \hat{r} e^{-k^2 t}. \quad (129)$$

Clearly, high wavenumbers  $k$  converge faster than low wavenumbers, since the residual of the former decays faster.

## 7.2 Alternate direction implicit

For two-dimensional or three-dimensional problems, we generally simplify the problem by making an operator split

$$[I - \Delta t(\mathbf{L}_x + \mathbf{L}_y)] = [I - \Delta t\mathbf{L}_x][I - \Delta t\mathbf{L}_y] + 2\mathbf{L}_x\mathbf{L}_y\Delta t^2, \quad (130)$$

where  $\mathbf{L}_x$  and  $\mathbf{L}_y$  are the discrete Laplacian operator in each direction. The split above is accurate to second-order in time, which is the accuracy of a Crank-Nicolson scheme. We then solve the problem in two steps

$$[I - \Delta t\mathbf{L}_x]r^* = r^n, \quad (131)$$

and

$$[I - \Delta t\mathbf{L}_y]r^{n+1} = r^*. \quad (132)$$

These two step are much easier to solve than the original problem because they're series of tri-diagonal systems. This ADI scheme can be embedded in the Crank-Nicolson scheme

$$[I - \frac{\Delta t}{2}\mathbf{L}_x]r^* = [I + \frac{\Delta t}{2}\mathbf{L}_y]r^n \quad (133)$$

$$[I - \frac{\Delta t}{2}\mathbf{L}_y]r^{n+1} = [I + \frac{\Delta t}{2}\mathbf{L}_x]r^*. \quad (134)$$

Notice that the growth factor at each sub-steps is

$$\sigma_k = \frac{1 - k^2 \frac{\Delta t}{2}}{1 + k^2 \frac{\Delta t}{2}}. \quad (135)$$

Figure (6) shows the absolute values of the growth factor  $|\sigma_k|$  as a function of time step for various values of  $k\Delta x$ . As can be seen, different wavenumbers have distinct optimum time steps. The optimum time step for low wavenumbers is typically very large; the contrary is true for very small wavenumbers:

$$\delta_{max} = \Delta t \left( \frac{\pi}{\Delta x} \right)^2 = \Delta t k_{max}^2, \quad (136)$$

$$\delta_{min} = \Delta t \left( \frac{\pi}{\Delta x} \right)^2 = \Delta t k_{min}^2 = \frac{\delta_{max}}{N^2}, \quad (137)$$

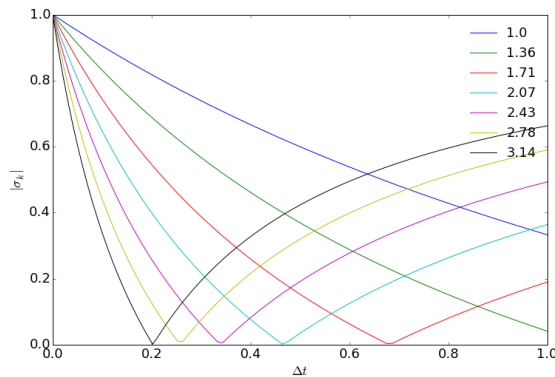


Figure 6: The growth factor of the ADI-CN method for various values of  $k\Delta x$ .

(SOME SUTFF AND DON'T UNDERSTAND WHERE IT COMES FROM. FOR IN-STANCE HOW TO OBTAIN THIS DELTA -8, ETC)

The optimum  $\Delta t$  is then

$$\Delta t_{\text{opt}} = 2 \frac{\Delta x^2}{\pi^2} N = 2 \frac{L_x \Delta x}{\pi^2}, \quad (138)$$

which is the geometric mean between the extrema. An estimative of the number of iterations required for convergence within machine precision is

$$\left| \frac{r^{n+1} - r^n}{r^0} \right| = \text{Rel. err.} . \quad (139)$$

Thus

$$\log \text{Rel. err.} = \log \text{tol.} = n \log(\epsilon - 1), \quad (140)$$

where

$$\epsilon = 1 - 2 \frac{\delta_{\max}}{N^2} = 1 - \frac{4}{N}. \quad (141)$$

We can always sweep randomly through the  $\Delta t$  space. In this case, it can be shown that we obtain a faster convergence ( $N^{1/2}$ ).