# Dealiased Convolutions without the Padding

John C. Bowman and Malcolm Roberts (University of Alberta)

June 15, 2010

`www.math.ualberta.ca/∼bowman/talks`

# Outline

- Discrete Convolutions

  – Cyclic vs. Linear

  – Standard vs. Centered

  – Complex vs. Hermitian

- Dealiasing

  – Zero Padding

  – Phase-shift dealiasing

- Implicit Padding in 1D, 2D, and 3D:

  – Standard Complex

  – Centered Hermitian

  – Ternary Convolutions

- Conclusions

# Discrete Convolutions

- Discrete linear convolution sums based on the fast Fourier transform (FFT) algorithm [Gauss 1866], [Cooley & Tukey 1965] have become important tools for:

    – image filtering;

    – digital signal processing;

    – correlation analysis;

    – pseudospectral simulations.

# Discrete Cyclic Convolution

- The FFT provides an efficient tool for computing the *discrete cyclic convolution*

$$\sum_{p=0}^{N-1} F_p G_{k-p},$$

  where the vectors $F$ and $G$ have period $N$.

- Define the $N$th primitive root of unity:

$$\zeta_N = \exp\left(\frac{2\pi i}{N}\right).$$

- The fast Fourier transform method exploits the properties that $\zeta_N^r = \zeta_{N/r}$ and $\zeta_N^N = 1$.

- The unnormalized backwards discrete Fourier transform of $\{F_k : k = 0, \ldots, N\}$ is

$$f_j \doteq \sum_{k=0}^{N-1} \zeta_N^{jk} F_k \qquad j = 0, \ldots, N-1,$$

- The corresponding forward transform is

$$F_k \doteq \frac{1}{N} \sum_{j=0}^{N-1} \zeta_N^{-kj} f_j \qquad k = 0, \ldots, N-1.$$

- The orthogonality of this transform pair follows from

$$\sum_{j=0}^{N-1} \zeta_N^{\ell j} = \begin{cases} N & \text{if } \ell = sN \text{ for } s \in \mathbb{Z}, \\ \dfrac{1 - \zeta_N^{\ell N}}{1 - \zeta_N^{\ell}} = 0 & \text{otherwise.} \end{cases}$$
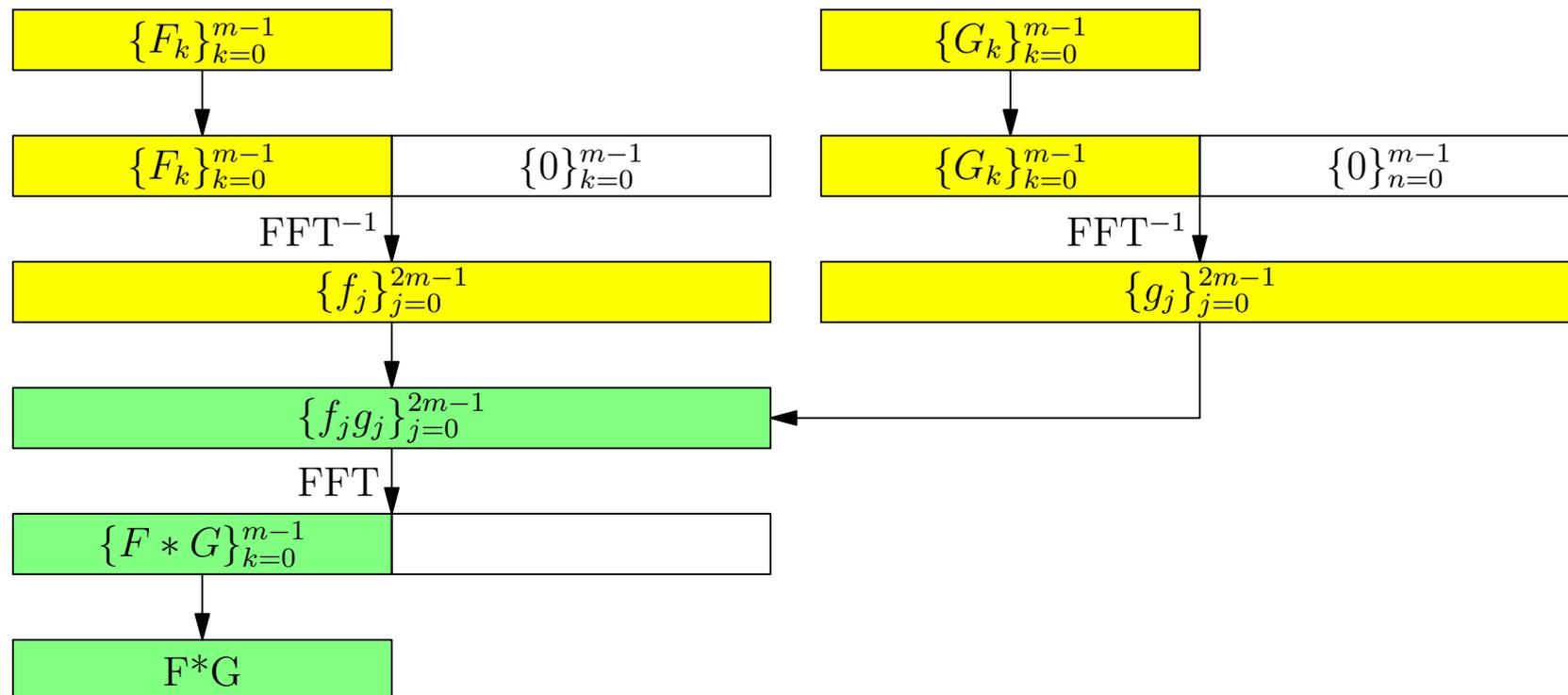
# Discrete Linear Convolution

- The pseudospectral method requires a *linear convolution* since wavenumber space is not periodic.

- The convolution theorem states:

$$\sum_{j=0}^{N-1} f_j g_j \zeta_N^{-jk} = \sum_{j=0}^{N-1} \zeta_N^{-jk} \left( \sum_{p=0}^{N-1} \zeta_N^{jp} F_p \right) \left( \sum_{q=0}^{N-1} \zeta_N^{jq} G_q \right)$$

$$= \sum_{p=0}^{N-1} \sum_{q=0}^{N-1} F_p G_q \sum_{j=0}^{N-1} \zeta_N^{(-k+p+q)j}$$

$$= N \sum_{s} \sum_{p=0}^{N-1} F_p G_{k-p+sN}.$$

- The terms indexed by $s \neq 0$ are called *aliases.*

- We need to remove the aliases by ensuring that $G_{k-p+sN} = 0$ whenever $s \neq 0$.

- If $F_p$ and $G_{k-p+sN}$ are nonzero only for $0 \leq p \leq m-1$ and $0 \leq k - p + sN \leq m - 1$, then we want $k + sN \leq 2m - 2$ to have no solutions for positive $s$.

- This can be achieved by choosing $N \geq 2m - 1$.

- That is, one must *zero pad* input data vectors of length $m$ to length $N \geq 2m - 1$:

| $\{F_k\}_{k=0}^{m-1}$ | | | $\{G_k\}_{k=0}^{m-1}$ | |
|---|---|---|---|---|

$$\text{FFT}^{-1} \qquad\qquad \text{FFT}^{-1}$$

| $\{F_k\}_{k=0}^{m-1}$ | $\{0\}_{k=0}^{m-1}$ | | $\{G_k\}_{k=0}^{m-1}$ | $\{0\}_{n=0}^{m-1}$ |
|---|---|---|---|---|

| $\{f_j\}_{j=0}^{2m-1}$ | | | $\{g_j\}_{j=0}^{2m-1}$ | |
|---|---|---|---|---|

| $\{f_j g_j\}_{j=0}^{2m-1}$ |
|---|

$$\text{FFT}$$

| $\{F * G\}_{k=0}^{m-1}$ | |
|---|---|

| F*G |
|---|

- Physically, *explicit zero padding* prevents mode $m - 1$ from beating with itself, wrapping around to contaminate mode $N = 0 \bmod N$.

- Since FFT sizes with small prime factors in practice yield the most efficient implementations, the padding is normally extended to $N = 2m$.

# Implicit Padding

- If $f_k = 0$ for $k \geq m$, one can easily avoid looping over the unwanted zero Fourier modes by decimating in wavenumber

$$f_{2\ell} = \sum_{k=0}^{m-1} \zeta_{2m}^{2\ell k} F_k = \sum_{k=0}^{m-1} \zeta_m^{\ell k} F_k,$$

$$f_{2\ell+1} = \sum_{k=0}^{m-1} \zeta_{2m}^{(2\ell+1)k} F_k = \sum_{k=0}^{m-1} \zeta_m^{\ell k} \zeta_N^k F_k \qquad \ell = 0, 1, \dots m-1.$$

- This requires computing two subtransforms, each of size $m$, for an overall computational scaling of order $2m \log_2 m = N \log_2 m$.

- Odd and even terms of the convolution can then be computed separately, multiplied term-by-term, and transformed again to Fourier space:

$$2mF_k = \sum_{j=0}^{2m-1} \zeta_{2m}^{-kj} f_j = \sum_{\ell=0}^{m-1} \zeta_{2m}^{-k2\ell} f_{2\ell} + \sum_{\ell=0}^{m-1} \zeta_{2m}^{-k(2\ell+1)} f_{2\ell+1}$$

$$= \sum_{\ell=0}^{m-1} \zeta_m^{-k\ell} f_{2\ell} + \zeta_{2m}^{-k} \sum_{\ell=0}^{m-1} \zeta_m^{-k\ell} f_{2\ell+1} \qquad k = 0, \ldots, m-1.$$

- No bit reversal is required at the highest level.

- An implicitly padded convolution is implemented as in our FFTW++ library (version 1.07) as cconv(f,g,u,v) computes an in-place implicitly dealiased convolution of two complex vectors f and g using two temporary vectors u and v, each of length $m$.

- This in-place convolution requires six out-of-place transforms, thereby avoiding bit reversal at all levels.

**Input**: vector f, vector g
**Output**: vector f
$u \leftarrow \texttt{fft}^{-1}(f);$
$v \leftarrow \texttt{fft}^{-1}(g);$
$u \leftarrow u * v;$
**for** $k = 0$ **to** $m - 1$ **do**
$\quad | \quad f[k] \leftarrow \zeta_{2m}^k f[k];$
$\quad | \quad g[k] \leftarrow \zeta_{2m}^k g[k];$
**end**
$v \leftarrow \texttt{fft}^{-1}(f);$
$f \leftarrow \texttt{fft}^{-1}(g);$
$v \leftarrow v * f;$
$f \leftarrow \texttt{fft}(u);$
$u \leftarrow \texttt{fft}(v);$
**for** $k = 0$ **to** $m - 1$ **do**
$\quad | \quad f[k] \leftarrow f[k] + \zeta_{2m}^{-k} u[k];$
**end**
**return** f/(2m);

11

# Implicit Padding in 1D

# Implicit Padding in 2D

# Implicit Padding in 3D

# Hermitian Convolutions

- *Hermitian convolutions* arise when the input vectors are Fourier transforms of real data:

$$f_{N-k} = \overline{f_k}.$$
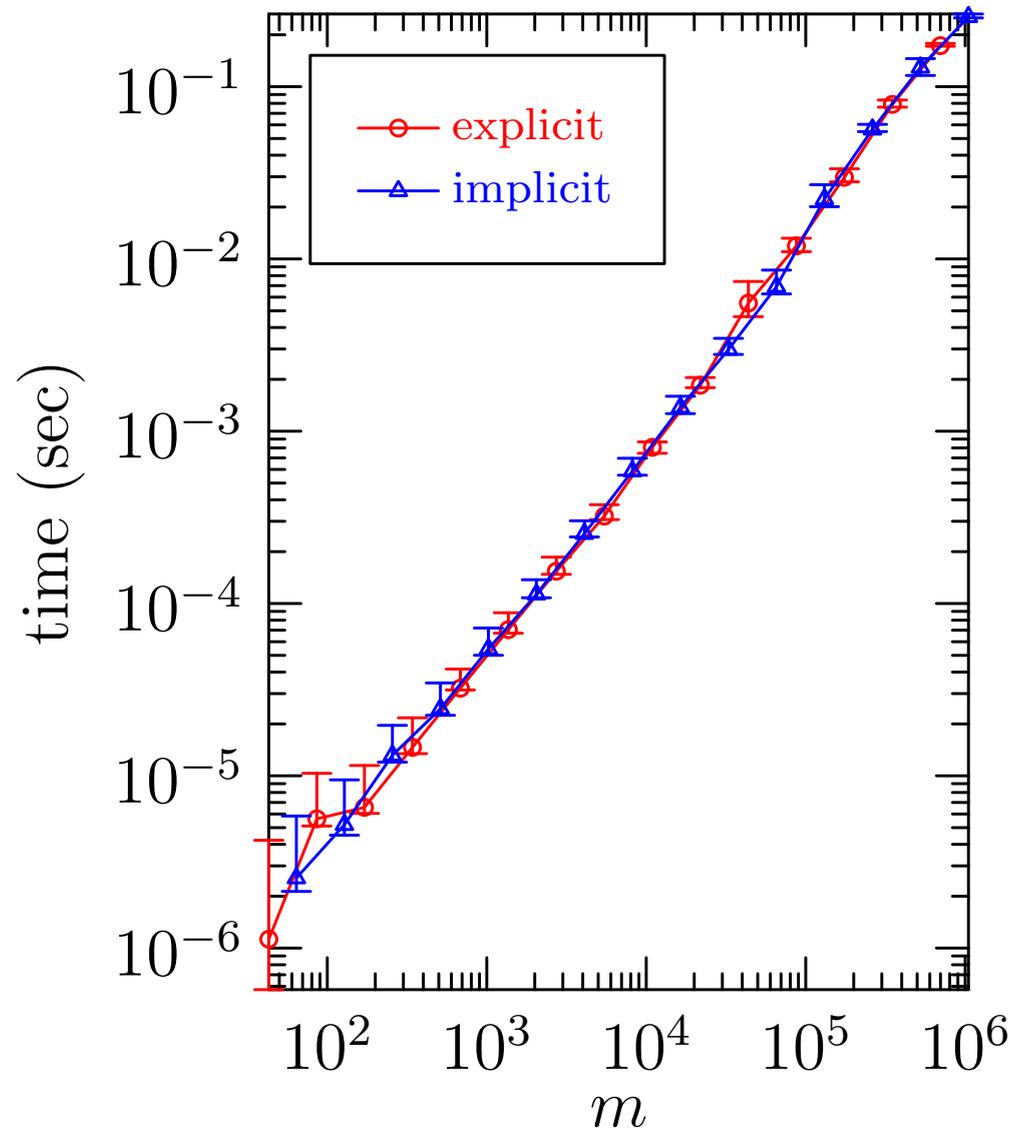
# Centered Convolutions

- For a *centered convolution*, the Fourier origin $(k = 0)$ is centered in the domain:
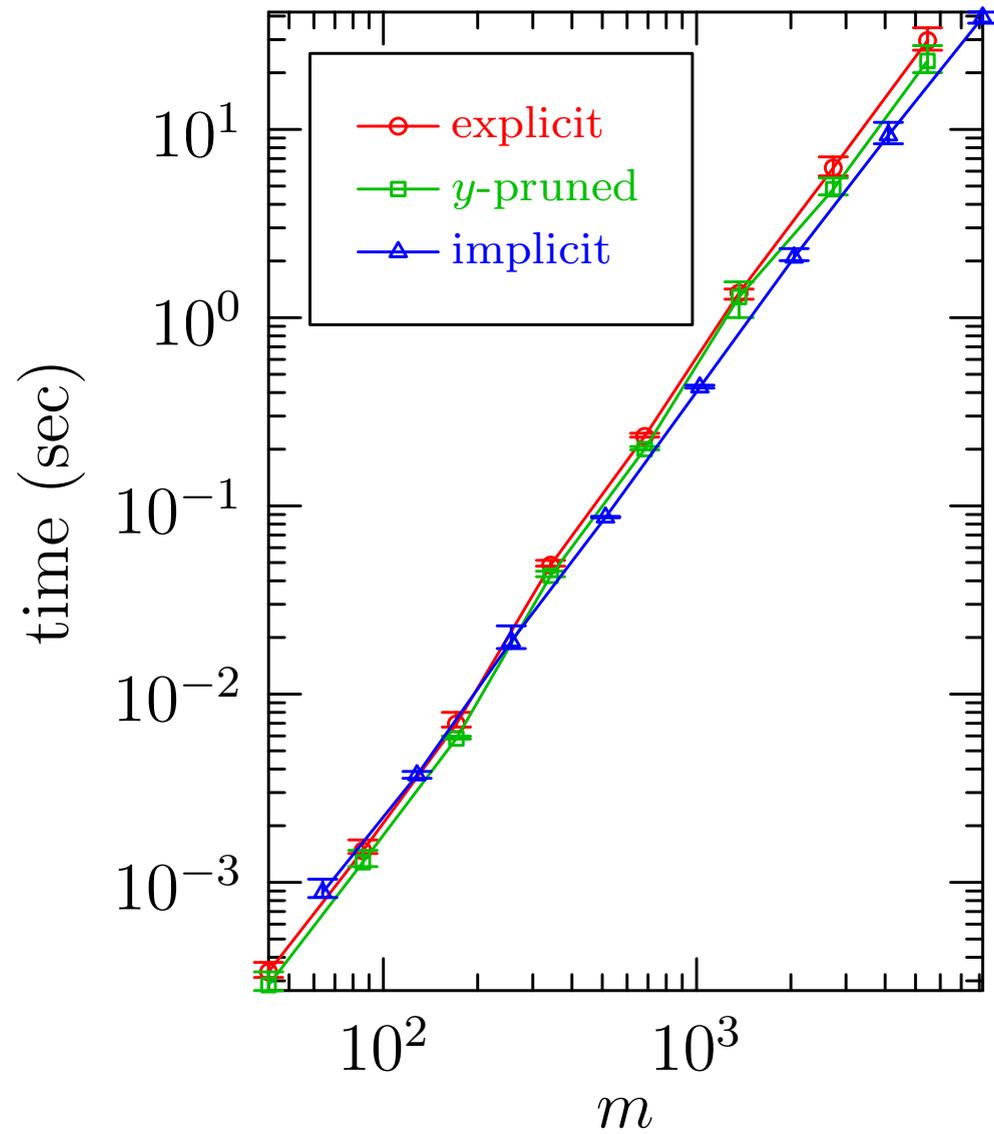
$$\sum_{p=k-m+1}^{m-1} f_p g_{k-p}$$

- Here, one needs to pad to $N \geq 3m - 2$ to prevent mode $m - 1$ from beating with itself to contaminate the most negative (first) mode, corresponding to wavenumber $-m + 1$. Since the ratio of the number of physical to total modes, $(2m - 1)/(3m - 2)$ is asymptotic to $2/3$ for large $m$, this padding scheme is often referred to as the *2/3 padding rule.*

- The Hermiticity condition then appears as

$$f_{-k} = \overline{f_k}.$$

# Implicit Hermitician Centered Padding in 1D

# Implicit Hermitician Centered Padding in 2D

# Ternary convolution

- The *ternary convolution* of three vectors $F$, $G$, and $H$ is

$$\sum_{p=0}^{N-1} \sum_{q=0}^{N-1} F_p G_q H_{k-p-q}.$$

- Computing the transfer function for $Z_4 = N^3 \sum_{\boldsymbol{j}} \omega^4(x_{\boldsymbol{j}})$ requires computing the Fourier transform of the cubic quantity $\omega^3$.
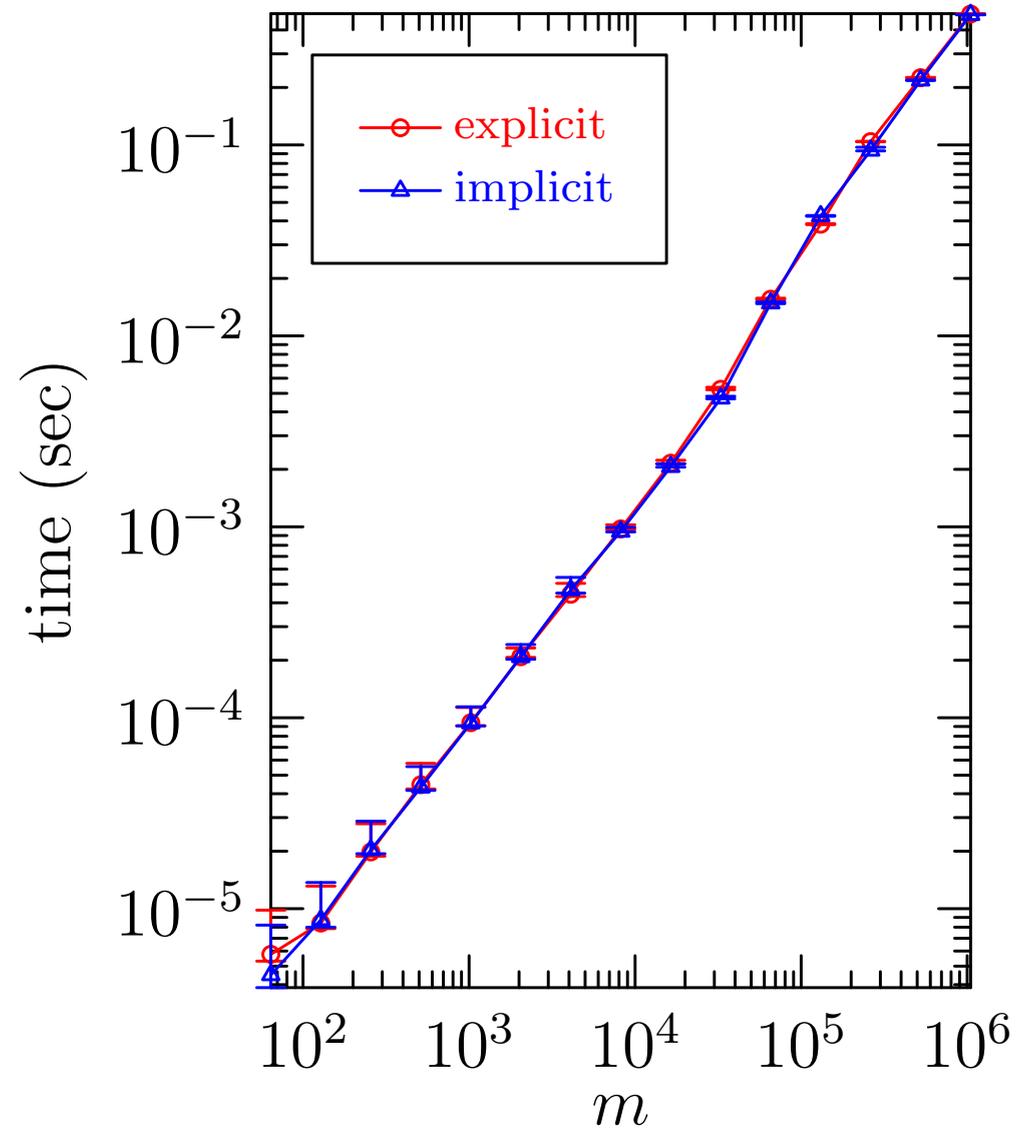
- This requires a centered Hermitian ternary convolution:

$$\sum_{p=-m+1}^{m-1} \sum_{q=-m+1}^{m-1} \sum_{r=-m+1}^{m-1} F_p G_q H_r \delta_{p+q+r,k}.$$

- Correctly dealiasing requires a 2/4 zero padding rule (instead of the usual 2/3 rule for a single convolution).
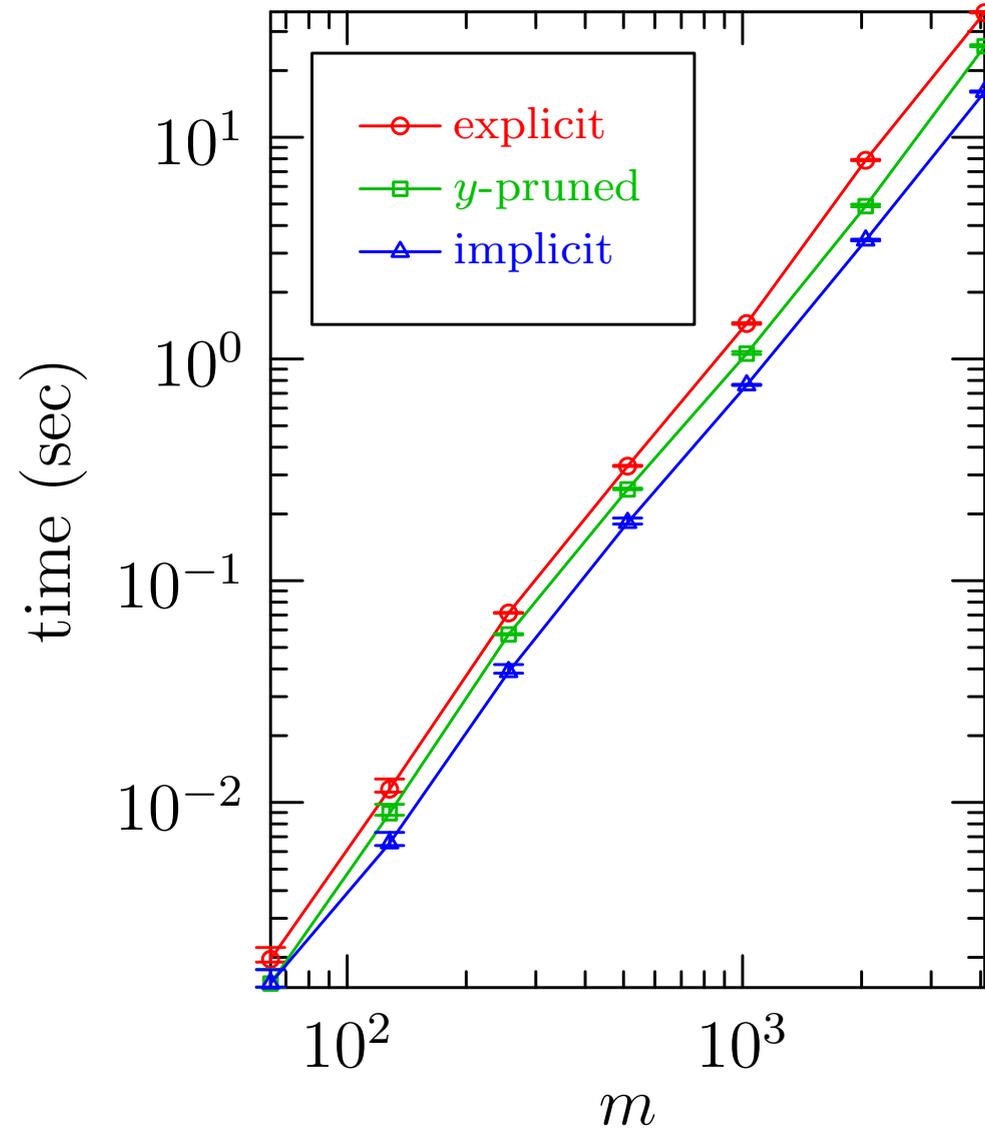
# 2/4 Padding Rule

- Computing the transfer function for $Z_4$ with a 2/4 padding rule means that in a $2048 \times 2048$ pseudospectral simulation, the maximum physical wavenumber retained in each direction is only 512.

- For a centered Hermitian ternary convolution, implicit padding is twice as fast and uses half of the memory required by conventional explicit padding.

# Implicit Ternary Convolution in 1D

# Implicit Ternary Convolution in 2D

# Conclusions

- Memory savings: in $d$ dimensions implicit padding asymptotically uses $1/2^{d-1}$ of the memory require by conventional explicit padding.

- Computational savings due to increased data locality: about a factor of two.

- Highly optimized versions of these routines have been implemented as a software layer `FFTW++` on top of the `FFTW` library and released under the Lesser GNU Public License.

- With the advent of this `FFTW++` library, writing a high-performance dealiased pseudospectral code is now a relatively straightforward exercise.
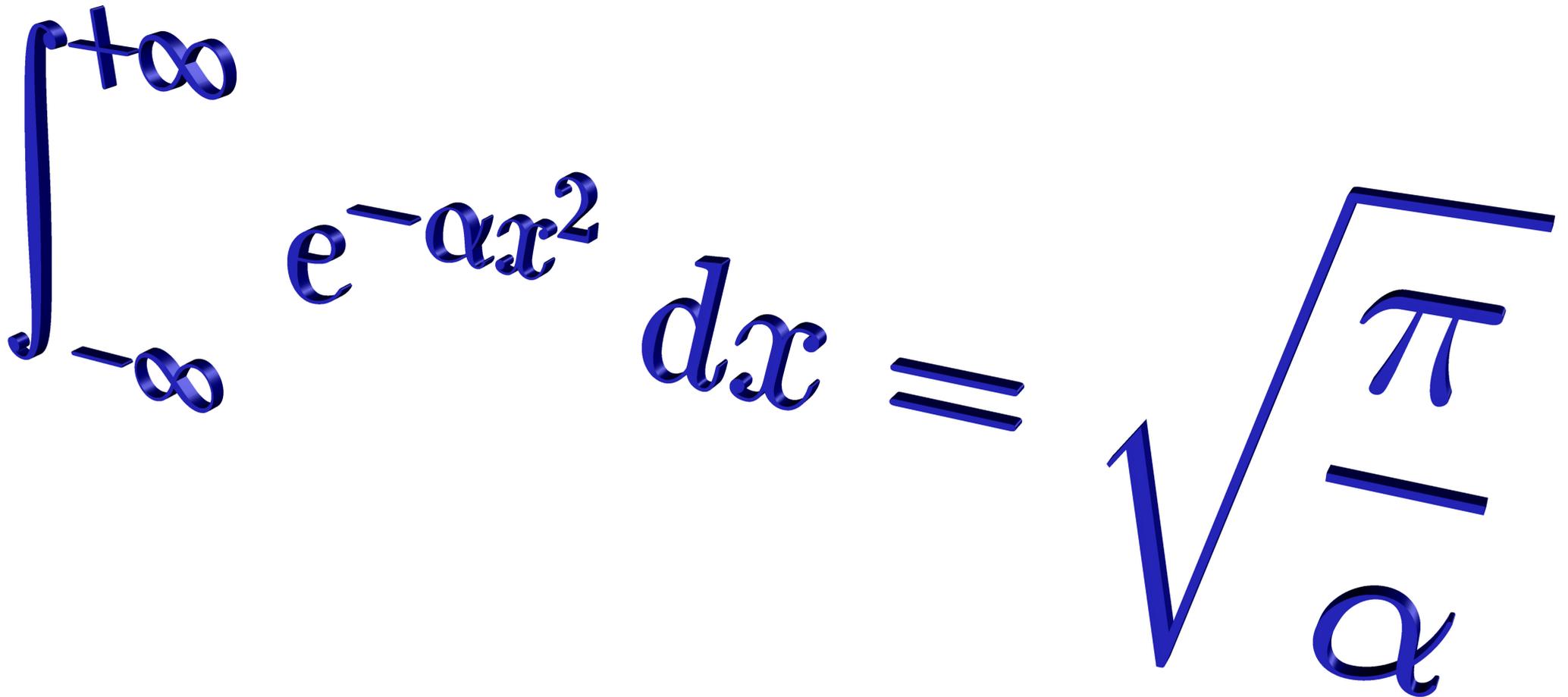
# Asymptote: 2D & 3D Vector Graphics Language



Andy Hammerlindl, John C. Bowman, Tom Prince

`http://asymptote.sf.net`

(freely available under the Lesser GNU Public License)

# Asymptote Lifts TeX to 3D

$$\int_{-\infty}^{+\infty} e^{-\alpha x^2}\, dx = \sqrt{\frac{\pi}{\alpha}}$$

http://asymptote.sf.net

# References

[Cooley & Tukey 1965]  J. W. Cooley & J. W. Tukey, Mathematics of Computation, **19**:297, 1965.

[Frigo & Johnson ]  "M. Frigo & S. G. Johnson, `http://www.fftw.org/pruned.html`.

[Gauss 1866]  C. F. Gauss, "Nachlass: Theoria interpolationis methodo nova tractata," in *Carl Friedrich Gauss Werke*, volume 3, pp. 265–330, Königliche Gesellschaft der Wissenschaften, Göttingen, 1866.