

Simulation and High-Performance Computing

Steffen Börm

October 11, 2020

Contents

1	Simulations based on ordinary differential equations	3
1.1	Timestepping for ordinary differential equations	3
1.2	Timestepping methods for special applications	8
2	Higher-order methods	14
2.1	Higher-order methods	14
2.2	Runge-Kutta methods and multi-step methods	24
3	Partial differential equations	29
3.1	Finite difference methods for partial differential equations	29
3.2	LR decomposition and BLAS	33
3.3	Method of lines for the heat equation and the wave equation	39
4	Iterative and semi-iterative solvers	43
4.1	Iterations	43
4.2	Krylov methods for linear systems of equations	49
5	Algorithms for large systems	54
5.1	Multigrid methods	54
5.2	Non-local force fields	58
6	Shared-memory parallelization	61
6.1	Multithreading and OpenMP	61
6.2	Tasks	66
7	Vectorization	70
7.1	Vectorization	70
7.2	Explicit vectorization	74
8	Computing with graphics cards	78
8.1	Computing with graphics cards using CUDA	78
8.2	Streams, blocks, and shared memory	82
9	Distributed computing	86
9.1	Message passing and MPI	86
9.2	Collective communication	91

1 Simulations based on ordinary differential equations

1.1 Timestepping for ordinary differential equations

Our goal: Predict the future, based on given natural laws.

Examples

Predator-prey system: (Alfred J. Lotka, Vito Volterra)

$$b'(t) = b(t)(\alpha - \beta r(t)), \quad r'(t) = r(t)(\gamma b(t) - \delta).$$

$b(t)$ is the number of prey, $r(t)$ the number of predators at time t .

Logistic differential equation: (Pierre-François Verhulst)

$$x'(t) = x(t)(1 - x(t)).$$

In population models, e.g., $x(t)$ could be the number of bacteria at time t .

Classical mechanics: (Isaac Newton)

$$x'(t) = v(t), \quad v'(t) = \frac{1}{m}f(t).$$

m is the mass, $x(t)$ is the position, $v(t)$ is the velocity, $f(t)$ is the force of a body.

Simple mass-spring system: (Isaac Newton, Robert Hooke)

$$x'(t) = v(t), \quad v'(t) = -\frac{c}{m}x(t).$$

c is the spring constant, m is the mass, $x(t)$ is the extension of the spring, $v(t)$ the velocity.

Many-body mass-spring system:

$$x'_i(t) = v_i(t), \quad v'_i(t) = \sum_{j \neq i} \frac{c_{ij}}{m_i} (L_{ij} - \|x_i(t) - x_j(t)\|) \frac{x_i(t) - x_j(t)}{\|x_i(t) - x_j(t)\|}.$$

c_{ij} is the spring constant for the spring between bodies i and j . L_{ij} is the resting length of this spring. m_i is the mass of the i -th body.

Many-body gravitational system: (Isaac Newton)

$$x'_i(t) = v_i(t), \quad v'_i(t) = \gamma \sum_{j \neq i} m_j \frac{x_j(t) - x_i(t)}{\|x_j(t) - x_i(t)\|^3}.$$

γ is the gravitational constant, m_j is the mass of the j -th body, $x_i(t)$ is the position and $v_i(t)$ is the velocity of the i -th body at time t .

Common form

All equations share the same form

$$y'(t) = f(t, y(t))$$

with a vector-valued function y and a right-hand side function f .

Difference quotient: Taylor expansion centered at t yields

$$\begin{aligned} y(t + \delta) &= y(t) + \delta y'(t) + \frac{\delta^2}{2} y''(\eta), \\ y(t + \delta) - y(t) &= \delta y'(t) + \frac{\delta^2}{2} y''(\eta), \\ \frac{y(t + \delta) - y(t)}{\delta} &= y'(t) + \frac{\delta}{2} y''(\eta) \end{aligned}$$

with $\eta \in [t, t + \delta]$, i.e., the *difference quotient* on the left side of the equation converges to $y'(t)$ with an error on the order of δ .

Explicit Euler method: (Leonhard Euler)

$$\begin{aligned} \frac{y(t + \delta) - y(t)}{\delta} &\approx y'(t) = f(t, y(t)), \\ y(t + \delta) - y(t) &\approx \delta f(t, y(t)), \\ y(t + \delta) &\approx \underbrace{y(t) + \delta f(t, y(t))}_{=: \tilde{y}(t + \delta)}. \end{aligned}$$

Accuracy determined by the accuracy of the difference quotient:

$$y(t + \delta) - \tilde{y}(t + \delta) = \frac{\delta^2}{2} y''(\eta)$$

with $\eta \in [t, t + \delta]$. Taking the maximum on the right-hand side yields

$$\|y(t + \delta) - \tilde{y}(t + \delta)\| \leq \frac{\delta^2}{2} \|y''\|_{\infty, [t, t + \delta]},$$

so we need small timesteps to reach a high accuracy.

Timestepping Consider the solution at times $t_k = t_0 + \delta k$, proceed iteratively

$$\begin{aligned}\tilde{y}(t_0) &\leftarrow y(t_0), \\ \tilde{y}(t_{k+1}) &\leftarrow \tilde{y}(t_k) + \delta f(t_k, \tilde{y}(t_k)).\end{aligned}$$

Example: logistic differential equation

$$y'(t) = y(t)(1 - y(t)),$$

applying Euler yields

$$\tilde{y}(t_{k+1}) \leftarrow \tilde{y}(t_k) + \delta \tilde{y}(t_k)(1 - \tilde{y}(t_k)).$$

Example: simple mass-spring system

$$x'(t) = v(t), \quad v'(t) = -\frac{c}{m}x(t),$$

equivalent with

$$y(t) = \begin{pmatrix} x(t) \\ v(t) \end{pmatrix}, \quad y'(t) = \begin{pmatrix} y_2(t) \\ -\frac{c}{m}y_1(t) \end{pmatrix},$$

and applying Euler yields

$$\tilde{y}(t_{k+1}) \leftarrow \tilde{y}(t_k) + \delta \begin{pmatrix} \tilde{y}_2(t) \\ -\frac{c}{m}\tilde{y}_1(t) \end{pmatrix},$$

which is equivalent to

$$\begin{aligned}\tilde{x}(t_{k+1}) &= \tilde{x}(t_k) + \delta \tilde{v}(t_k), \\ \tilde{v}(t_{k+1}) &= \tilde{v}(t_k) - \delta \frac{c}{m}\tilde{x}(t_k).\end{aligned}$$

Experiment: Convergence Simulate the mass-spring system starting at $t = 0$ and stopping at $t = 20$.

δ	error	ratio
1	1.02 ₊₃	
1/2	8.24 ₊₁	12.4
1/4	7.87 ₊₀	10.5
1/8	1.32 ₊₀	6.0
1/16	3.97 ₋₁	3.3
1/32	1.58 ₋₁	2.5
1/64	7.07 ₋₂	2.2
1/128	3.36 ₋₂	2.1
1/256	1.64 ₋₂	2.1

Under typical conditions, the Euler method exhibits first-order convergence if the solution is twice continuously differentiable, i.e., the error is proportional to the stepsize δ , therefore a high accuracy requires a large number of steps.

Central difference quotient: Taylor expansion centered at $t + \frac{\delta}{2}$ yields

$$\begin{aligned} y(t + \delta) &= y(t + \frac{\delta}{2}) + \frac{\delta}{2}y'(t + \frac{\delta}{2}) + \frac{\delta^2}{8}y''(t + \frac{\delta}{2}) + \frac{\delta^3}{48}y'''(\eta_+), \\ y(t) &= y(t + \frac{\delta}{2}) - \frac{\delta}{2}y'(t + \frac{\delta}{2}) + \frac{\delta^2}{8}y''(t + \frac{\delta}{2}) - \frac{\delta^3}{48}y'''(\eta_-), \\ y(t + \delta) - y(t) &= \delta y'(t + \frac{\delta}{2}) + \frac{\delta^3}{24}y'''(\eta), \\ \frac{y(t + \delta) - y(t)}{\delta} &= y'(t + \frac{\delta}{2}) + \frac{\delta^2}{24}y'''(\eta) \end{aligned}$$

with $\eta_+ \in [t, t + \frac{\delta}{2}]$, $\eta_- \in [t - \frac{\delta}{2}, t]$, and $\eta \in [\eta_-, \eta_+] \subseteq [t - \frac{\delta}{2}, t + \frac{\delta}{2}]$.

While the error of the standard difference quotient behaves like δ , the central difference quotient yields δ^2 , i.e., significantly faster convergence.

Runge's method: If we use the *central difference quotient*, we have

$$\begin{aligned} \frac{y(t + \delta) - y(t)}{\delta} &\approx y'(t + \frac{\delta}{2}), \\ y(t + \delta) &\approx y(t) + \delta y'(t + \frac{\delta}{2}), \\ y(t + \delta) &\approx y(t) + \delta f(t + \frac{\delta}{2}, y(t + \frac{\delta}{2})). \end{aligned}$$

Unfortunately, we do not know $y(t + \frac{\delta}{2})$, but we can use an explicit Euler step to approximate it:

$$y(t + \frac{\delta}{2}) \approx y(t) + \frac{\delta}{2}f(t, y(t)).$$

Combining both approximations yields

$$y(t + \delta) \approx y(t) + \delta f\left(t + \frac{\delta}{2}, y(t) + \frac{\delta}{2}f(t, y(t))\right) =: \tilde{y}(t + \delta)$$

Accuracy similar to that of *central* difference quotient, i.e., on the order of δ^2 , if the solution y is thrice continuously differentiable and additionally if the right-hand side f is Lipschitz continuous in the second argument, i.e., if

$$|f(t, x_1) - f(t, x_2)| \leq L_f |x_1 - x_2| \quad \text{for all } t, x_1, x_2.$$

Runge timestepping: Introduce intermediate steps $t_{k+1/2} = t_0 + \delta(k + 1/2)$ and use

$$\begin{aligned} \tilde{y}(t_0) &\leftarrow y(t_0), \\ \tilde{y}(t_{k+1/2}) &\leftarrow \tilde{y}(t_k) + \frac{\delta}{2}f(t_k, \tilde{y}(t_k)), \\ \tilde{y}(t_{k+1}) &\leftarrow \tilde{y}(t_k) + \delta f(t_{k+1/2}, \tilde{y}(t_{k+1/2})). \end{aligned}$$

Computational work doubled compared to Euler's method, but second-order convergence instead of first-order, i.e., an accuracy of ϵ requires $\sim 1/\epsilon$ timesteps with Euler, but only $\sim 2/\sqrt{\epsilon}$ timesteps with Runge.

Example: simple mass-spring system

$$\begin{aligned}\tilde{x}(t_{k+1/2}) &\leftarrow \tilde{x}(t_k) + \frac{\delta}{2} \tilde{v}(t_k), & \tilde{v}(t_{k+1/2}) &\leftarrow \tilde{v}(t_k) - \frac{\delta}{2} \frac{c}{m} \tilde{x}(t_k), \\ \tilde{x}(t_{k+1}) &\leftarrow \tilde{x}(t_k) + \delta \tilde{v}(t_{k+1/2}), & \tilde{v}(t_{k+1}) &\leftarrow \tilde{v}(t_k) - \delta \frac{c}{m} \tilde{x}(t_{k+1/2}).\end{aligned}$$

Experiment: Convergence Simulate the mass-spring system starting at $t = 0$ and stopping at $t = 20$.

δ	Euler		Runge	
	error	ratio	error	ratio
1	1.02 ₊₃		9.61 ₊₀	
1/2	8.24 ₊₁	12.4	8.70 ₋₁	11.1
1/4	7.87 ₊₀	10.5	1.85 ₋₁	4.7
1/8	1.32 ₊₀	6.0	4.61 ₋₂	4.0
1/16	3.97 ₋₁	3.3	1.17 ₋₂	4.0
1/32	1.58 ₋₁	2.5	2.94 ₋₃	4.0
1/64	7.07 ₋₂	2.2	7.39 ₋₄	4.0
1/128	3.36 ₋₂	2.1	1.85 ₋₄	4.0
1/256	1.64 ₋₂	2.1	4.64 ₋₅	4.0

Second-order convergence is clearly visible, Runge's method reaches higher accuracies than Euler's method, and the errors converge significantly faster.

1.2 Timestepping methods for special applications

Goal of the leapfrog method: Efficiency of Runge at the cost of Euler

Idea: Central difference quotient, as for Runge's method.

$$\begin{aligned}y'(t + \frac{\delta}{2}) &\approx \frac{y(t + \delta) - y(t)}{\delta}, \\y(t + \delta) &\approx y(t) + \delta y'(t + \frac{\delta}{2}), \\y(t + \delta) &\approx y(t) + 2\delta f(t + \frac{\delta}{2}, y(t + \frac{\delta}{2})).\end{aligned}$$

Leapfrog: We need $y(t)$ and $y(t + \frac{\delta}{2})$ to compute $y(t + \delta)$, so we start with a half-step using Euler's method to approximate $y(t + \frac{\delta}{2})$ and then proceed with central differences.

$$\begin{aligned}\tilde{y}(t_0) &= y(t_0), & \tilde{y}(t_{1/2}) &= \tilde{y}(t_0) + \frac{\delta}{2}f(t_0, \tilde{y}(t_0)), \\ \tilde{y}(t_{k+1}) &= \tilde{y}(t_k) + \delta f(t_{k+1/2}, \tilde{y}(t_{k+1/2})), \\ \tilde{y}(t_{k+3/2}) &= \tilde{y}(t_{k+1/2}) + \delta f(t_{k+1}, \tilde{y}(t_{k+1})) && \text{for all } k \in \mathbb{N}_0.\end{aligned}$$

Special case: In mass-spring systems and the gravitational system, the force depends only on the position of the bodies, not on the velocities, i.e., we have

$$x'(t) = v(t), \quad v'(t) = f(t, x(t)).$$

If we are only interested in $\tilde{x}(t_k)$, we can save work:

$$\begin{aligned}\tilde{x}(t_0) &= x(t_0), & \tilde{v}(t_{1/2}) &= v(t_0) + \frac{\delta}{2}f(t_0, \tilde{x}(t_0)), \\ \tilde{x}(t_{k+1}) &= \tilde{x}(t_k) + \delta \tilde{v}(t_{k+1/2}), \\ \tilde{v}(t_{k+3/2}) &= \tilde{v}(t_{k+1/2}) + \delta f(t_{k+1}, \tilde{x}(t_{k+1})) && \text{for all } k \in \mathbb{N}_0.\end{aligned}$$

Experiment: Convergence Simulate the mass-spring system starting at $t = 0$ and stopping at $t = 20$.

δ	Euler		Runge		Leapfrog	
	error	ratio	error	ratio	error	ratio
1	1.02 ₊₃		9.61 ₊₀		9.08 ₋₁	
1/2	8.24 ₊₁	12.4	8.70 ₋₁	11.1	2.04 ₋₁	4.5
1/4	7.87 ₊₀	10.5	1.85 ₋₁	4.7	4.84 ₋₂	4.2
1/8	1.32 ₊₀	6.0	4.61 ₋₂	4.0	1.19 ₋₂	4.1
1/16	3.97 ₋₁	3.3	1.17 ₋₂	4.0	2.98 ₋₃	4.0
1/32	1.58 ₋₁	2.5	2.94 ₋₃	4.0	7.43 ₋₄	4.0
1/64	7.07 ₋₂	2.2	7.39 ₋₄	4.0	1.86 ₋₄	4.0
1/128	3.36 ₋₂	2.1	1.85 ₋₄	4.0	4.64 ₋₅	4.0
1/256	1.64 ₋₂	2.1	4.64 ₋₅	4.0	1.16 ₋₅	4.0

We can see that the leapfrog method also converges like δ^2 , but is always more accurate than Runge's method, although it requires only half the computational work.

Goal of the implicit Euler method: Avoid unrealistic behaviour of the numerical approximation in the pre-asymptotic regime.

Example: Exponential function $y(t) = e^{-\lambda t}$ with $\lambda > 0$ satisfies $y'(t) = -\lambda y(t)$. Applying the explicit Euler method to this equation yields

$$\tilde{y}(t_{k+1}) = \tilde{y}(t_k) - \delta\lambda\tilde{y}(t_k) = (1 - \delta\lambda)\tilde{y}(t_k),$$

i.e., numerical solution changes its sign in each step if $\delta\lambda > 1$.

Implicit Euler method:

$$\begin{aligned} \frac{y(t + \delta) - y(t)}{\delta} &\approx y'(t + \delta) = f(t + \delta, y(t + \delta)), \\ y(t + \delta) &\approx y(t) + \delta f(t + \delta, y(t + \delta)). \end{aligned}$$

Approximation of $y(t + \delta)$ is given *implicitly* as the solution z of a fixed-point equation.

Timestepping requires us to solve fixed-point equations in each step.

$$\begin{aligned} \tilde{y}(t_0) &\leftarrow y(t_0), \\ \tilde{y}(t_{k+1}) &\text{ is solution of } \tilde{y}(t_{k+1}) = \tilde{y}(t_k) + \delta f(t_{k+1}, \tilde{y}(t_{k+1})). \end{aligned}$$

Example: exponential function We have to solve

$$\begin{aligned} \tilde{y}(t_{k+1}) &= \tilde{y}(t_k) - \delta\lambda\tilde{y}(t_{k+1}), \\ (1 + \delta\lambda)\tilde{y}(t_{k+1}) &= \tilde{y}(t_k), \\ \tilde{y}(t_{k+1}) &\leftarrow \frac{1}{1 + \delta\lambda}\tilde{y}(t_k). \end{aligned}$$

No oscillations, exponential decay, similar to the exact solution (but at a slower rate).

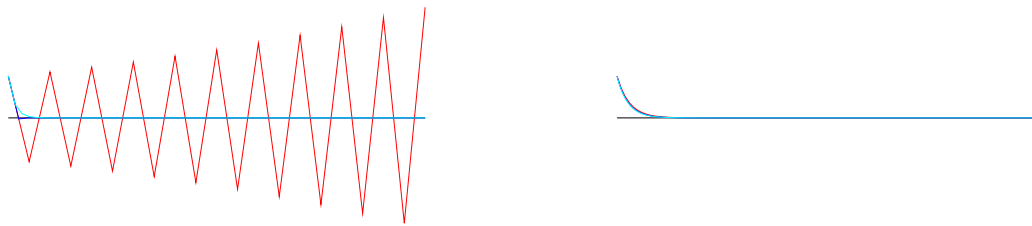


Figure 1.1: Numerical solutions of $y'(t) = -\lambda y(t)$ with explicit and implicit Euler.

Energy conservation: In the simple mass-spring system

$$x'(t) = v(t), \quad v'(t) = -\frac{c}{m}x(t),$$

the energy

$$E(t) := \frac{c}{2}x(t)^2 + \frac{m}{2}v(t)^2$$

is conserved over time, since the product rule of differentiation yields

$$E'(t) = cx(t)x'(t) + mv(t)v'(t) = cx(t)v(t) - mv(t)\frac{c}{m}x(t) = 0.$$

Since the energy is an important quantity, we are interested in finding timestepping methods that also conserve it.

Energy and Euler's methods: For the explicit Euler method, we have

$$\begin{aligned} \tilde{E}(t + \delta) &:= \frac{c}{2}(x(t) + \delta v(t))^2 + \frac{m}{2}\left(v(t) - \delta\frac{c}{m}x(t)\right)^2 \\ &= E(t) + \delta cx(t)v(t) + \delta^2\frac{c}{2}v(t)^2 - \delta cv(t)x(t) + \delta^2\frac{c^2}{2m}x(t)^2 \\ &= E(t) + \delta^2\frac{c}{m}\left(\frac{m}{2}v(t)^2 + \frac{c}{2}x(t)^2\right) = E(t) + \delta^2\frac{c}{m}E(t) \geq E(t), \end{aligned}$$

i.e., the energy increases with each timestep.

For the implicit Euler method, on the other hand, we have

$$\begin{aligned} \tilde{E}(t + \delta) &:= \frac{c}{2}(x(t) + \delta\tilde{v}(t + \delta))^2 + \frac{m}{2}\left(v(t) - \delta\frac{c}{m}\tilde{x}(t + \delta)\right)^2 \\ &= E(t) + \delta cx(t)\tilde{v}(t + \delta) + \delta^2\frac{c}{2}\tilde{v}(t + \delta)^2 - \delta cv(t)\tilde{x}(t + \delta) + \delta^2\frac{c^2}{2m}\tilde{x}(t + \delta)^2 \\ &= E(t) + \delta c(x(t) + \delta\tilde{v}(t + \delta))\tilde{v}(t + \delta) - \delta^2\frac{c}{2}\tilde{v}(t + \delta)^2 \\ &\quad - \delta c\left(v(t) - \delta\frac{c}{m}\tilde{x}(t + \delta)\right)\tilde{x}(t + \delta) - \delta^2\frac{c^2}{2m}\tilde{x}(t + \delta)^2 \\ &= E(t) + \delta\tilde{x}(t + \delta)\tilde{v}(t + \delta) - \delta\tilde{v}(t + \delta)\tilde{x}(t + \delta) - \delta^2\frac{c}{m}\tilde{E}(t + \delta) \\ &= E(t) - \delta^2\frac{c}{m}\tilde{E}(t + \delta) \leq E(t), \end{aligned}$$

i.e., the energy decreases with each timestep.

Maybe a combination of the explicit and the implicit method could conserve the energy and therefore preserve the behaviour of the original differential equation?

Crank-Nicolson method: Take the arithmetic mean of the forward and backward difference quotients.

$$\begin{aligned}
y'(t) &\approx \frac{y(t+\delta) - y(t)}{\delta}, \\
y'(t+\delta) &\approx \frac{y(t+\delta) - y(t)}{\delta}, \\
\frac{y'(t) + y'(t+\delta)}{2} &\approx \frac{y(t+\delta) - y(t)}{\delta}, \\
y(t+\delta) &\approx y(t) + \delta \frac{y'(t) + y'(t+\delta)}{2}, \\
y(t+\delta) &\approx y(t) + \frac{\delta}{2} (f(t, y(t)) + f(t+\delta, y(t+\delta))).
\end{aligned}$$

This is again an implicit method, since $y(t+\delta)$ appears on both sides of the fixed-point equation.

Timestepping

$$\begin{aligned}
\tilde{y}(t_0) &\leftarrow y(t_0), \\
\tilde{y}(t_{k+1}) &\text{ is solution of } \tilde{y}(t_{k+1}) = \tilde{y}(t_k) + \frac{\delta}{2} (f(t_k, \tilde{y}(t_k)) + f(t_{k+1}, \tilde{y}(t_{k+1}))).
\end{aligned}$$

Example: mass-spring system

$$x'(t) = v(t), \quad v'(t) = -\frac{c}{m}x(t)$$

leads to

$$\begin{aligned}
\tilde{x}(t_{k+1}) &= \tilde{x}(t_k) + \frac{\delta}{2} (\tilde{v}(t_k) + \tilde{v}(t_{k+1})), \\
\tilde{v}(t_{k+1}) &= \tilde{v}(t_k) - \frac{\delta}{2} \frac{c}{m} (\tilde{x}(t_k) + \tilde{x}(t_{k+1})).
\end{aligned}$$

Substituting $\tilde{v}(t_{k+1})$ in the first equation by the second equation yields

$$\begin{aligned}
\tilde{x}(t_{k+1}) &= \tilde{x}(t_k) + \frac{\delta}{2} \left(\tilde{v}(t_k) + \tilde{v}(t_k) - \frac{\delta}{2} \frac{c}{m} (\tilde{x}(t_k) + \tilde{x}(t_{k+1})) \right) \\
&= \tilde{x}(t_k) + \delta \tilde{v}(t_k) - \frac{\delta^2}{4} \frac{c}{m} \tilde{x}(t_k) - \frac{\delta^2}{4} \frac{c}{m} \tilde{x}(t_{k+1}), \\
\left(1 + \frac{\delta^2 c}{4m}\right) \tilde{x}(t_{k+1}) &= \left(1 - \frac{\delta^2 c}{4m}\right) \tilde{x}(t_k) + \delta \tilde{v}(t_k),
\end{aligned}$$

while substituting $\tilde{x}(t_{k+1})$ in the second equation by the first leads to

$$\tilde{v}(t_{k+1}) = \tilde{v}(t_k) - \frac{\delta}{2} \frac{c}{m} \left(\tilde{x}(t_k) + \tilde{x}(t_k) + \frac{\delta}{2} (\tilde{v}(t_k) + \tilde{v}(t_{k+1})) \right)$$

$$\begin{aligned}
&= \tilde{v}(t_k) - \delta \frac{c}{m} \tilde{x}(t_k) - \frac{\delta^2}{4} \frac{c}{m} \tilde{v}(t_k) - \frac{\delta^2}{4} \frac{c}{m} \tilde{v}(t_{k+1}), \\
\left(1 + \frac{\delta^2 c}{4m}\right) \tilde{v}(t_{k+1}) &= \left(1 - \frac{\delta^2 c}{4m}\right) \tilde{v}(t_k) - \delta \frac{c}{m} \tilde{x}(t_k).
\end{aligned}$$

Experiment: Convergence Simulate the mass-spring system starting at $t = 0$ and stopping at $t = 20$.

δ	Explicit		Implicit		Crank-Nicol	
	error	ratio	error	ratio	error	ratio
1	1.02 ₊₃		4.09 ₋₁		5.46 ₋₁	
1/2	8.24 ₊₁	12.4	3.97 ₋₁	1.0	3.24 ₋₁	1.7
1/4	7.87 ₊₀	10.5	3.43 ₋₁	1.2	9.19 ₋₂	3.5
1/8	1.32 ₊₀	6.0	2.63 ₋₁	1.3	2.36 ₋₂	3.9
1/16	3.97 ₋₁	3.3	1.77 ₋₁	1.5	5.93 ₋₃	4.0
1/32	1.58 ₋₁	2.5	1.05 ₋₁	1.7	1.49 ₋₃	4.0
1/64	7.07 ₋₂	2.2	5.78 ₋₂	1.8	3.71 ₋₄	4.0
1/128	3.36 ₋₂	2.1	3.03 ₋₂	1.9	9.29 ₋₅	4.0
1/256	1.64 ₋₂	2.1	1.55 ₋₂	2.0	2.32 ₋₅	4.0

We observe first-order convergence for the Euler methods and second-order convergence for the Crank-Nicolson method.

Energy conservation: the Crank-Nicolson method conserves the energy:

$$\begin{aligned}
\tilde{E}(t_{k+1}) - \tilde{E}(t_k) &= \frac{c}{2}(\tilde{x}(t_{k+1})^2 - \tilde{x}(t_k)^2) + \frac{m}{2}(\tilde{v}(t_{k+1})^2 - \tilde{v}(t_k)^2) \\
&= \frac{c}{2}(\tilde{x}(t_{k+1}) - \tilde{x}(t_k))(\tilde{x}(t_{k+1}) + \tilde{x}(t_k)) \\
&\quad + \frac{m}{2}(\tilde{v}(t_{k+1}) - \tilde{v}(t_k))(\tilde{v}(t_{k+1}) + \tilde{v}(t_k)) \\
&= \frac{c}{2} \frac{\delta}{2} (\tilde{v}(t_{k+1}) + \tilde{v}(t_k))(\tilde{x}(t_{k+1}) + \tilde{x}(t_k)) \\
&\quad - \frac{m}{2} \frac{\delta}{2} \frac{c}{m} (\tilde{x}(t_{k+1}) + \tilde{x}(t_k))(\tilde{v}(t_{k+1}) + \tilde{v}(t_k)) = 0.
\end{aligned}$$

Euler methods do not conserve the energy. For the explicit Euler method, we have

$$\begin{aligned}
\tilde{E}(t_{k+1}) &= \frac{c}{2} \tilde{x}(t_{k+1})^2 + \frac{m}{2} \tilde{v}(t_{k+1})^2 \\
&= \frac{c}{2} (\tilde{x}(t_k) + \delta \tilde{v}(t_k))^2 + \frac{m}{2} (\tilde{v}(t_k) - \delta \frac{c}{m} \tilde{x}(t_k))^2 \\
&= \frac{c}{2} (\tilde{x}(t_k)^2 + 2\delta \tilde{x}(t_k) \tilde{v}(t_k) + \delta^2 \tilde{v}(t_k)^2) \\
&\quad + \frac{m}{2} \left(\tilde{v}(t_k)^2 - 2\delta \frac{c}{m} \tilde{x}(t_k) \tilde{v}(t_k) + \delta^2 \frac{c^2}{m^2} \tilde{x}(t_k)^2 \right)
\end{aligned}$$

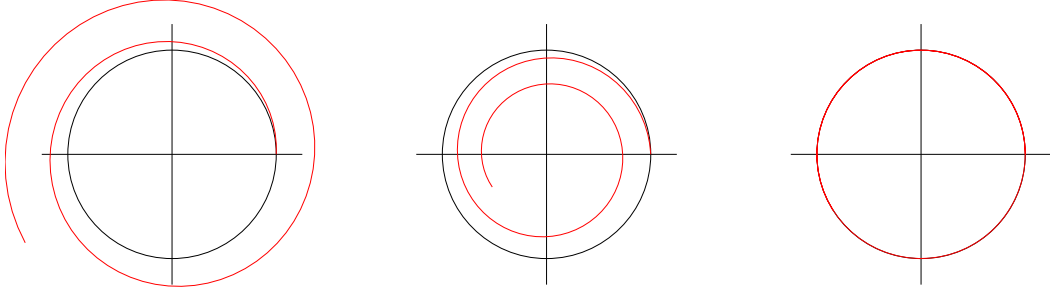


Figure 1.2: Phase diagrams for the spring-mass system plotting the curve $(\tilde{x}(t), \tilde{v}(t))$ in two-dimensional space with explicit and implicit Euler and Crank-Nicolson. For $c = m = 1$, the energy $x(t)^2 + v(t)^2$ should remain constant, i.e., the curve should be contained in the unit circle.

$$= \tilde{E}(t_k) + \frac{\delta^2}{2} \left(c\tilde{v}(t_k)^2 + \frac{c^2}{m}\tilde{x}(t_k)^2 \right) = \left(1 + \delta^2 \frac{c}{m} \right) \tilde{E}(t_k),$$

so each timestep will lead to an increase in energy (unless we start with zero energy).

For the implicit Euler method, we obtain

$$\begin{aligned} \tilde{E}(t_{k+1}) &= \frac{c}{2}\tilde{x}(t_{k+1})^2 + \frac{m}{2}\tilde{v}(t_{k+1})^2 \\ &= \frac{c}{2}(\tilde{x}(t_k) + \delta\tilde{v}(t_{k+1}))^2 + \frac{m}{2}\left(\tilde{v}(t_k) - \delta\frac{c}{m}\tilde{x}(t_{k+1})\right)^2 \\ &= \frac{c}{2}(\tilde{x}(t_k)^2 + 2\delta\tilde{x}(t_k)\tilde{v}(t_{k+1}) + \delta^2\tilde{v}(t_{k+1})^2) \\ &\quad + \frac{m}{2}\left(\tilde{v}(t_k)^2 - 2\delta\frac{c}{m}\tilde{v}(t_k)\tilde{x}(t_{k+1}) + \delta^2\frac{c^2}{m^2}\tilde{x}(t_{k+1})^2\right) \\ &= \frac{c}{2}(\tilde{x}(t_k)^2 + 2\delta\tilde{x}(t_{k+1})\tilde{v}(t_{k+1}) - \delta^2\tilde{v}(t_{k+1})^2) \\ &\quad + \frac{m}{2}\left(\tilde{x}(t_k)^2 - 2\delta\frac{c}{m}\tilde{v}(t_{k+1})\tilde{x}(t_{k+1}) - \delta^2\frac{c^2}{m^2}\tilde{x}(t_{k+1})^2\right) \\ &= \tilde{E}(t_k) - \frac{\delta^2}{2}\left(c\tilde{v}(t_{k+1})^2 + \frac{c^2}{m}\tilde{x}(t_{k+1})^2\right) = \tilde{E}(t_k) - \delta^2\frac{c}{m}\tilde{E}(t_{k+1}). \end{aligned}$$

2 Higher-order methods

2.1 Higher-order methods

Observation: In order to reduce the error by a factor of 100, a first-order method like Euler requires 100 times as many steps, while a second-order method like Runge or Crank-Nicolson requires only 10 times as many.

Goal: Construct higher-order methods.

Polynomials: Most higher-order methods are based on polynomials.

$$p(x) = c_0 + c_1x + \dots + c_mx^m = \sum_{k=0}^m c_k x^k$$

Example: Taylor expansion Given a function $f \in C^{m+1}[a, b]$ and $x, x_0 \in [a, b]$, we find $\eta \in (a, b)$ with

$$f(x) = \sum_{k=0}^m \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k + \frac{f^{(m+1)}(\eta)}{(m+1)!} (x - x_0)^{m+1},$$

i.e., the *Taylor polynomial* of degree m

$$p(x) = \sum_{k=0}^m \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k$$

is a $(m+1)$ -th order approximation of $f(x)$.

Interpolation Since computing derivatives all derivatives $f, f', \dots, f^{(m)}$ by hand is often tedious and prone to errors, we frequently prefer to use the values of f in different points $x_0, \dots, x_m \in [a, b]$ instead: We look for a polynomial p of degree m such that

$$p(x_k) = f(x_k) \quad \text{for all } k \in [0 : m].$$

Unfortunately, now p is given *implicitly* as the solution of these equations.

Lagrange polynomials In theory, we can represent p by *Lagrange polynomials*

$$p(x) = \sum_{j=0}^m f(x_j) \ell_j(x), \quad \ell_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^m \frac{x - x_i}{x_j - x_i},$$

since they satisfy

$$\ell_j(x_k) = \begin{cases} 1 & \text{if } k = j, \\ 0 & \text{otherwise.} \end{cases}$$

Actually evaluating p in this way is time-consuming and not particularly stable.

Application: Numerical differentiation We have already seen that we can construct approximations of the derivative of a function f from appropriately chosen point values.

With interpolation, we can approximate the derivative in a point y by

$$f'(y) \approx p'(y) = \sum_{j=0}^m f(x_j) \underbrace{\ell_j'(y)}_{=:w_j} = \sum_{j=0}^m w_j f(x_j).$$

If f is a polynomial of degree m , this equation is exact, and this property can be used to determine the weights w_0, \dots, w_m without the need for the derivatives of the Lagrange polynomials: we choose linear independent polynomials for which the derivative can be computed easily, e.g., monomials $f_i(x) = x^i$ for $i \in [0 : m]$ with

$$f_i'(y) = \begin{cases} 0 & \text{if } i = 0, \\ iy^{i-1} & \text{otherwise.} \end{cases}$$

For all of these monomials, our algorithm should deliver the exact derivative, i.e.,

$$\sum_{j=0}^m x_j^i w_j = \sum_{j=0}^m w_j f_i(x_j) = f_i'(y) = \begin{cases} 0 & \text{if } i = 0, \\ iy^{i-1} & \text{otherwise} \end{cases}$$

This is a system of $m + 1$ linear equations that can be solved to obtain the $m + 1$ weights w_0, \dots, w_m *without* the need for computing the derivatives of the Lagrange polynomials explicitly.

Example: Cubic approximation of the first derivative Let $f \in C^1[-2, 2]$. We want to approximate $f'(0)$ by $p'(0)$, where p is a cubic polynomial interpolating f in $x_0 = -3$, $x_1 = -1$, $x_2 = 1$, $x_3 = 3$. We obtain the linear system

$$\begin{aligned} 0 &= x_0^0 w_0 + x_1^0 w_1 + x_2^0 w_2 + x_3^0 w_3 = w_0 + w_1 + w_2 + w_3, \\ 1 &= x_0^1 w_0 + x_1^1 w_1 + x_2^1 w_2 + x_3^1 w_3 = -3w_0 - w_1 + w_2 + 3w_3, \\ 0 &= x_0^2 w_0 + x_1^2 w_1 + x_2^2 w_2 + x_3^2 w_3 = 9w_0 + w_1 + w_2 + 9w_3, \end{aligned}$$

$$0 = x_0^3 w_0 + x_1^3 w_1 + x_2^3 w_2 + x_3^3 w_3 = -27w_0 - w_1 + w_2 + 27w_3.$$

Combining the first and third equation yields $w_0 = -w_3$ and $w_1 = -w_2$, and the fourth equation gives us $w_1 = -27w_0$.

Substituting these identities in the second equation leads us to $w_0 = \frac{1}{48}$ and the approximation

$$f'(0) \approx p'(0) = \frac{f(-3) - 27f(-1) + 27f(1) - f(3)}{48}.$$

Scaling In order to control the accuracy, we fix a *scaling parameter* $h > 0$ and apply the approximation to

$$\hat{f}(x) := f(hx) \tag{2.1}$$

instead of f itself. The chain rule yields $\hat{f}'(x) = hf'(hx)$ and we obtain

$$f'(0) = \frac{1}{h} \hat{f}'(0) \approx \frac{f(-3h) - 27f(-h) + 27f(h) - f(3h)}{48h}.$$

Applying this approximation to $f(x) = e^x$ gives us fourth-order accuracy.

h	error	ratio
1	9.53 ₋₂	
1/2	4.98 ₋₃	19.2
1/4	2.97 ₋₄	16.7
1/8	1.84 ₋₅	16.2
1/16	1.15 ₋₆	16.0
1/32	7.15 ₋₈	16.0
1/64	4.47 ₋₉	16.0

Usually we would only expect third-order accuracy (fourth order for p , but we lose one order by differentiating), but our approximation happens to be accurate for quartic polynomials as well, giving us fourth order.

Example: Quadratic approximation of the second derivative Let $f \in C^2[-1, 1]$. We want to approximate $f''(0)$ by $p''(0)$, where p is a quadratic polynomial interpolating f in $x_0 = -1$, $x_1 = 0$, $x_2 = 1$.

We obtain the linear system

$$\begin{aligned} 0 &= x_0^0 w_0 + x_1^0 w_1 + x_2^0 w_2 = w_0 + w_1 + w_2, \\ 0 &= x_0^1 w_0 + x_1^1 w_1 + x_2^1 w_2 = -w_0 + w_2, \\ 2 &= x_0^2 w_0 + x_1^2 w_1 + x_2^2 w_2 = w_0 + w_2. \end{aligned}$$

The second row yields $w_0 = w_2$, the third $w_0 + w_2 = 2$, i.e., $w_0 = w_2 = 1$, so that the first row leads to $w_1 = -2$. We conclude

$$f''(0) \approx p''(0) = f(-1) - 2f(0) + f(1).$$

Using scaling again, we observe $\hat{f}''(x) = h^2 f''(hx)$ and arrive at

$$f''(0) = \frac{1}{h^2} \hat{f}''(0) \approx \frac{f(-h) - 2f(0) + f(h)}{h^2}.$$

Quartic approximation of the second derivative Let $f \in C^2[-2, 2]$. We want to approximate $f''(0)$ by $p''(0)$, where p is a quartic polynomial interpolating f in $x_0 = -2$, $x_1 = -1$, $x_2 = 0$, $x_3 = 1$, and $x_4 = 2$. Since the approximation has to be exact for quartic polynomials, we obtain the linear system

$$\begin{aligned} 0 &= x_0^0 w_0 + x_1^0 w_1 + x_2^0 w_2 + x_3^0 w_3 + x_4^0 w_4 = w_0 + w_1 + w_2 + w_3 + w_4, \\ 0 &= x_0^1 w_0 + x_1^1 w_1 + x_2^1 w_2 + x_3^1 w_3 + x_4^1 w_4 = -2w_0 - w_1 + w_3 + 2w_4, \\ 2 &= x_0^2 w_0 + x_1^2 w_1 + x_2^2 w_2 + x_3^2 w_3 + x_4^2 w_4 = 4w_0 + w_1 + w_3 + 4w_4, \\ 0 &= x_0^3 w_0 + x_1^3 w_1 + x_2^3 w_2 + x_3^3 w_3 + x_4^3 w_4 = -8w_0 - w_1 + w_3 + 8w_4, \\ 0 &= x_0^4 w_0 + x_1^4 w_1 + x_2^4 w_2 + x_3^4 w_3 + x_4^4 w_4 = 16w_0 + w_1 + w_3 + 16w_4. \end{aligned}$$

The second and fourth equations yield $w_0 = w_4$ and $w_1 = w_3$. The fifth equation gives us $w_1 = -16w_0$, and the third $2 = -24w_0$, i.e., $w_0 = -\frac{1}{12}$. With the first equation we find $30w_0 = w_2$, and therefore the approximation

$$f''(0) \approx p''(0) = \frac{-f(-2) + 16f(-1) - 30f(0) + 16f(1) - f(2)}{12}.$$

We can again use scaling to control the accuracy and arrive at

$$f''(0) = \frac{1}{h^2} \hat{f}''(0) \approx \frac{-f(-2h) + 16f(-h) - 30f(0) + 16f(h) - f(2h)}{12h^2}.$$

To verify the accuracy, we apply the approximations to $f(x) = e^x$.

h	2nd order		4th order	
	error	ratio	error	ratio
1	8.06 ₋₂		1.02 ₋₂	
1/2	2.07 ₋₂	3.9	6.79 ₋₄	15.0
1/4	5.20 ₋₃	4.0	4.32 ₋₅	15.7
1/8	1.30 ₋₃	4.0	2.71 ₋₆	15.9
1/16	3.25 ₋₄	4.0	1.69 ₋₇	16.0
1/32	8.14 ₋₅	4.0	1.06 ₋₈	16.0
1/64	2.03 ₋₅	4.0	6.62 ₋₁₀	16.0

Usually, we would only expect third-order accuracy (fifth order for p , but we lose two orders by differentiating), but our approximation happens to be accurate for quintic polynomials as well due to symmetry, giving us fourth order.

Application: Extrapolation Let $f \in C[0, 1]$. We want to approximate the limit $f(0) = \lim_{x \rightarrow 0} f(x)$, but we cannot simply evaluate f in zero, e.g., because evaluating $f(x)$ involves a division by x . With interpolation, we can approximate the limit by

$$f(0) = \lim_{x \rightarrow 0} f(x) \approx \lim_{x \rightarrow 0} p(x) = p(0) = \sum_{j=0}^m f(x_j) \underbrace{\ell_j(0)}_{=:w_j} = \sum_{j=0}^m w_j f(x_j).$$

As in the previous application, we can use the fact that the interpolation is exact for monomials $f_i(x) = x^i$ to obtain the linear system

$$\sum_{j=0}^m x_j^i w_j = \begin{cases} 1 & \text{if } i = 0, \\ 0 & \text{otherwise} \end{cases} \quad \text{for all } i \in [0 : m].$$

Example: Quadratic extrapolation Let $f \in C[0, 1]$. We want to approximate $f(0)$ by $p(0)$, where p is a quadratic polynomial interpolating f in $x_0 = 1/4$, $x_1 = 1/2$, $x_2 = 1$. We obtain the linear system

$$\begin{pmatrix} 1 & 1 & 1 \\ 1/4 & 1/2 & 1 \\ 1/16 & 1/4 & 1 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

with the solution $w_0 = 8/3$, $w_1 = -2$, $w_2 = 1/3$, i.e.,

$$f(0) \approx \frac{8f(1/4) - 6f(1/2) + f(1)}{3}$$

We can again use scaling (2.1) to control the accuracy:

$$f(0) = \hat{f}(0) \approx \frac{8f(h/4) - 6f(h/2) + f(h)}{3}$$

Example: Cubic extrapolation Let $f \in C[0, 1]$. We want to approximate $f(0)$ by $p(0)$, where p is a cubic polynomial interpolating f in $x_0 = 1/8$, $x_1 = 1/4$, $x_2 = 1/2$, and $x_3 = 1$. We obtain the linear system

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1/8 & 1/4 & 1/2 & 1 \\ 1/64 & 1/16 & 1/4 & 1 \\ 1/512 & 1/64 & 1/8 & 1 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

with the solution $w_0 = 64/21$, $w_1 = -8/3$, $w_2 = 2/3$, $w_3 = -1/21$, this leads to

$$f(0) \approx \frac{64f(1/8) - 56f(1/4) + 14f(1/2) - f(1)}{21}.$$

Scaling (2.1) yields

$$f(0) = \hat{f}(0) \approx \frac{64f(h/8) - 56f(h/4) + 14f(h/2) - f(h)}{21}.$$

Applying this approximation to $f(x) = \frac{e^x - 1}{x}$ yields

h	Quadratic		Cubic	
	error	ratio	error	ratio
1	7.48_{-3}		1.80_{-4}	
1/2	7.78_{-4}	9.6	9.54_{-6}	18.8
1/4	8.89_{-5}	8.8	5.50_{-7}	17.3
1/8	1.06_{-5}	8.4	3.31_{-8}	16.6
1/16	1.30_{-6}	8.2	2.03_{-9}	16.3
1/32	1.61_{-7}	8.1	1.26_{-10}	16.1
1/64	2.00_{-8}	8.0	7.74_{-12}	16.2

We can observe cubic convergence for the quadratic approximation and quartic convergence for the cubic approximation.

Application: Numerical integration We can also approximate the integral of f by

$$\int_a^b f(x) dx \approx \int_a^b p(x) dx = \sum_{j=0}^m f(x_j) \underbrace{\int_a^b \ell_j(x) dx}_{=:w_j} = \sum_{j=0}^m w_j f(x_j).$$

Once again, if f is a polynomial of degree m , the equation is exact, and we can compute the weights w_0, \dots, w_m by solving the system

$$\sum_{j=0}^m x_j^i w_j = \int_a^b x^i dx = \frac{b^{i+1} - a^{i+1}}{i+1} \quad \text{for all } i \in [0 : m].$$

Example: Trapezoidal quadrature We approximate the integral $\int_{-1}^1 f(x) dx$ using a linear polynomial p interpolating in $x_0 = -1$ and $x_1 = 1$. We obtain the linear system

$$\begin{aligned} 2 &= \int_{-1}^1 x^0 dx = x_0^0 w_0 + x_1^0 w_1 = w_0 + w_1, \\ 0 &= \int_{-1}^1 x^1 dx = x_0^1 w_0 + x_1^1 w_1 = -w_0 + w_1. \end{aligned}$$

The second row yields $w_0 = w_1$, and substituting this equation in the first row gives us $w_0 = w_1 = 1$. Our approximation of the integral takes the form

$$\int_{-1}^1 f(x) dx \approx \int_{-1}^1 p(x) dx = f(-1) + f(1).$$

Example: Simpson's quadrature We want to approximate the integral $\int_{-1}^1 f(x) dx$ using a quadratic polynomial p interpolating in $x_0 = -1$, $x_1 = 0$ and $x_2 = 1$. We obtain the linear system

$$2 = \int_{-1}^1 x^0 dx = x_0^0 w_0 + x_1^0 w_1 + x_2^0 w_2 = w_0 + w_1 + w_2,$$

$$0 = \int_{-1}^1 x^1 dx = x_0^1 w_0 + x_1^1 w_1 + x_2^1 w_2 = -w_0 + w_2,$$

$$\frac{2}{3} = \int_{-1}^1 x^2 dx = x_0^2 w_0 + x_1^2 w_1 + x_2^2 w_2 = w_0 + w_2,$$

that can be written in matrix-vector form as

$$\begin{pmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 0 \\ 2/3 \end{pmatrix}.$$

The second row yields $w_0 = w_2$, the third row yields $w_0 + w_2 = 2/3$, i.e., $w_0 = w_2 = 1/3$, and the first row leads to $w_1 = 2 - 2/3 = 4/3$. We conclude

$$\int_{-1}^1 f(x) dx \approx \int_{-1}^1 p(x) dx = \frac{1}{3}f(-1) + \frac{4}{3}f(0) + \frac{1}{3}f(1).$$

Composite quadrature The simple scaling approach used before cannot be applied directly to the task of numerical integration: if we use (2.1) and the substitution $y = hx$, we find

$$h \int_{-1}^1 \hat{f}(x) dx = \int_{-1}^1 f(hx)h dx = \int_{-h}^h f(y) dy,$$

i.e., the interval depends on h .

Since we usually want the interval to be fixed, we have to modify the scaling approach: assume that $a < b$ are given and we have to compute

$$\int_a^b f(x) dx.$$

In order to control the accuracy, we choose $k \in \mathbb{N}$, let

$$c_i := a + h(2i - 1), \quad h := \frac{b - a}{2k} \quad \text{for all } i \in [1 : k],$$

and apply the quadrature rule to the functions

$$\hat{f}_i(x) := f(hx + c_i) \quad \text{for all } i \in [1 : k].$$

Substitution yields

$$\begin{aligned} h \sum_{i=1}^k \int_{-1}^1 \hat{f}_i(x) dx &= \sum_{i=1}^k \int_{-1}^1 f_i(hx + c_i)h dx = \sum_{i=1}^k \int_{c_i-h}^{c_i+h} f(y) dy \\ &= \sum_{i=1}^k \int_{a+2h(i-1)}^{a+2hi} f(y) dy = \int_a^{a+2hk} f(y) dy = \int_a^b f(y) dy. \end{aligned}$$

Example: Composite trapezoidal and Simpson's quadrature For the trapezoidal rule, we obtain

$$\begin{aligned}
 \int_a^b f(y) dy &= h \sum_{i=1}^k \int_{-1}^1 \hat{f}_i(x) dx \approx h \sum_{i=1}^k \hat{f}_i(-1) + \hat{f}_i(1) \\
 &= h \sum_{i=1}^k f(c_i - h) + f(c_i + h) = h \sum_{i=1}^k f(a + 2h(i - 1)) + f(a + 2hi) \\
 &= h \sum_{i=0}^{k-1} f(a + 2hi) + h \sum_{i=1}^k f(a + 2hi) \\
 &= h \left(f(a) + 2 \sum_{i=1}^{k-1} f(a + 2hi) + f(b) \right).
 \end{aligned}$$

for Simpson's rule, we have

$$\begin{aligned}
 \int_a^b f(y) dy &= h \sum_{i=1}^k \int_{-1}^1 \hat{f}_i(x) dx \approx \frac{h}{3} \sum_{i=1}^k \hat{f}_i(-1) + 4\hat{f}_i(0) + \hat{f}_i(1) \\
 &= \frac{1}{3} \left(h \sum_{i=1}^k \hat{f}_i(-1) + \hat{f}_i(1) \right) + \frac{2}{3} \sum_{i=1}^k 2hf(c_i) \\
 &= \frac{h}{3} \left(f(a) + 2 \sum_{i=1}^{k-1} f(a + 2hi) + f(b) \right) + \frac{2h}{3} \sum_{i=1}^k 2f(a + h(2i - 1)),
 \end{aligned}$$

i.e., it can be computed by adding a correction term to the trapezoidal rule.

Applied to $\int_0^2 e^x dx$, we obtain the following approximations:

k	Trapezoidal		Simpson's	
	error	ratio	error	ratio
1	2.00 ₊₀		3.17 ₋₂	
2	5.24 ₋₁	3.8	2.15 ₋₃	14.7
4	1.33 ₋₁	4.0	1.38 ₋₄	15.7
8	3.32 ₋₂	4.0	8.65 ₋₆	15.9
16	8.32 ₋₃	4.0	5.41 ₋₇	16.0
32	2.08 ₋₃	4.0	3.38 ₋₈	16.0
64	5.20 ₋₄	4.0	2.12 ₋₉	16.0

We can see that the trapezoidal rule only shows second-order convergence, while Simpson's rule shows fourth-order convergence.

Example: Romberg integration Compared to other methods, the composite trapezoidal rule is very easy to apply, but unfortunately offers only second-order convergence.

We can improve the accuracy significantly by combining the trapezoidal rule with *extrapolation*, i.e., by approximating the trapezoidal rule with *infinitely many* subintervals.

We slightly change the notation: h denoted the radius of the subintervals before, now it denotes their diameter. The composite trapezoidal rule takes the form

$$T(h) := \frac{b-a}{2k} \left(f(a) + 2 \sum_{i=1}^{k-1} f\left(a + \frac{b-a}{k}i\right) + f(b) \right), \quad k := \frac{b-a}{h} \in \mathbb{N},$$

where we have to ensure that $b-a$ is always a multiple of h .

The error converges like h^2 , and the exact solution is the limit

$$\int_a^b f(x) dx = \lim_{h \rightarrow 0} T(h).$$

A closer look at the error reveals that we can find constants c_1, c_2, \dots satisfying

$$T(h) = \int_a^b f(x) dx + c_1 h^2 + c_2 h^4 + c_3 h^6 + \dots,$$

so approximating $T(h)$ directly would “waste” every second term in the polynomial.

To avoid this issue, we approximate the function $g(h) := T(\sqrt{h})$ that satisfies

$$g(h) = \int_a^b f(x) dx + c_1 h + c_2 h^2 + c_3 h^3 + \dots$$

and has no “gaps” in the polynomial expansion.

In order to construct an interpolating polynomial p , we have to evaluate g in interpolation points h_0, h_1, \dots , i.e., we need the values of $g(h_j) = T(\sqrt{h_j})$. $T(h)$ is only well-defined if $(b-a)/h$ is an integer, so we have to choose the interpolation points in a way that guarantees this property. A simple approach is to start with $h_0 := (b-a)^2/k^2$ for $k \in \mathbb{N}$ and continue with $h_1 := h_0/4, \dots, h_m := h_0/4^m$. The resulting algorithm is known as *Romberg integration*.

Neville-Aitken algorithm Romberg integration requires us to evaluate the interpolating polynomial p in just one point, $h = 0$. The idea of the *Neville-Aitken algorithm* is to define polynomials of increasing degree that interpolate the given function in increasing numbers of points: we define polynomials p_{ij} of degree $j-i$ such that they interpolate g in h_i, \dots, h_j . We obtain the recurrence relation

$$p_{ij}(x) = \begin{cases} g(h_i) & \text{if } i = j, \\ \frac{x-h_i}{h_j-h_i} p_{i+1,j}(x) + \frac{h_j-x}{h_j-h_i} p_{i,j-1}(x) & \text{otherwise} \end{cases} \quad \text{for all } i \leq j.$$

Since we are only interested in $x = 0$, we can make use of $h_k = h_0 4^{-k}$ and

$$\begin{aligned} \frac{-h_i}{h_j-h_i} &= \frac{-4^{-i}}{4^{-j}-4^{-i}} = \frac{-4^{j-i}}{1-4^{j-i}} = \frac{4^{j-i}}{4^{j-i}-1}, \\ \frac{h_j}{h_j-h_i} &= \frac{4^{-j}}{4^{-j}-4^{-i}} = \frac{1}{1-4^{j-i}} = \frac{-1}{4^{j-i}-1} \end{aligned}$$

to obtain the simpler relation

$$p_{ij}(0) = \begin{cases} g(h_i) & \text{if } i = j, \\ \frac{4^{j-i} p_{i+1,j}(0) - p_{i,j-1}(0)}{4^{j-i} - 1} & \text{otherwise} \end{cases} \quad \text{for all } i \leq j. \quad (2.2)$$

This leads to the triangular scheme

$$\begin{array}{cccccccc} g(h_0) = p_{00}(0) & & & & & & & & \\ & \searrow & & & & & & & \\ g(h_1) = p_{11}(0) & \rightarrow & p_{01}(0) & & & & & & \\ & \searrow & & \searrow & & & & & \\ g(h_2) = p_{22}(0) & \rightarrow & p_{12}(0) & \rightarrow & p_{02}(0) & & & & \\ & \searrow & & \searrow & & \searrow & & & \\ g(h_3) = p_{33}(0) & \rightarrow & p_{23}(0) & \rightarrow & p_{13}(0) & \rightarrow & p_{03}(0) & & \\ & \searrow & & \searrow & & \searrow & & \searrow & \\ g(h_4) = p_{44}(0) & \rightarrow & p_{34}(0) & \rightarrow & p_{24}(0) & \rightarrow & p_{14}(0) & \rightarrow & p_{04}(0) \end{array}$$

that can be evaluated efficiently by storing only the last row: if $p_{mm}(0), \dots, p_{0m}(0)$ are available, we can compute $g(h_{m+1}) = p_{m+1,m+1}(0)$ and then use the current values to find $p_{m,m+1}(0), \dots, p_{0,m+1}(0)$ successively by applying (2.2).

Applying the trapezoidal rule and Romberg integration to the exponential function in $[0, 2]$ yields the following errors:

k	Trapezoidal		Romberg	
	error	ratio	error	ratio
2	5.24_{-1}		5.24_{-1}	
4	1.33_{-1}	4.0	2.15_{-3}	243.1
8	3.32_{-2}	4.0	3.20_{-6}	647.1
16	8.32_{-3}	4.0	1.25_{-9}	2661.9
32	2.08_{-3}	4.0	1.22_{-13}	10251.7
64	5.20_{-4}	4.0	"0"	

We can see that extrapolation significantly improves the accuracy.

2.2 Runge-Kutta methods and multi-step methods

Ordinary differential equation in general explicit form

$$y'(t) = f(t, y(t)).$$

Goal: Develop higher-order methods that can be easily implemented.

Structure of method discussed so far:

$$\begin{array}{ll} y(t + \delta) \approx y(t) + \delta f(t, y(t)), & \text{explicit Euler,} \\ y(t + \delta) \approx y(t) + \delta f(t + \delta, y(t + \delta)), & \text{implicit Euler,} \\ y(t + \delta) \approx y(t) + \delta f\left(t + \frac{\delta}{2}, y(t) + \frac{\delta}{2} f(t, y(t))\right), & \text{Runge,} \\ y(t + \delta) \approx y(t) + \frac{\delta}{2} (f(t, y(t)) + f(t + \delta, y(t + \delta))), & \text{Crank-Nicolson.} \end{array}$$

Multiple evaluations of f , some of which depend on previous evaluations.

Implicit methods can be rewritten:

$$\tilde{y}(t + \delta) = y(t) + \delta k_1, \quad k_1 = f(t + \delta, y(t) + \delta k_1),$$

for the implicit Euler method and

$$\begin{array}{ll} \tilde{y}(t + \delta) = y(t) + \frac{\delta}{2}(k_1 + k_2), & k_1 = f(t, y(t)), \\ & k_2 = f\left(t + \delta, y(t) + \frac{\delta}{2}(k_1 + k_2)\right) \end{array}$$

for the Crank-Nicolson method.

General Runge-Kutta method with s stages:

$$\begin{aligned} k_i &= f\left(t + c_i \delta, y(t) + \delta \sum_{j=1}^s a_{ij} k_j\right), \\ y(t + \delta) &\approx y(t) + \delta \sum_{i=1}^s b_i k_i. \end{aligned}$$

Butcher tableaus provide compact notation for general Runge-Kutta methods:

$$\begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & \ddots & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline & b_1 & \cdots & b_s \end{array}$$

Tableaus for the methods we have seen so far:

0	0				
1	1				
1/2	1/2	0			
0	0				
1	1/2	1/2			
1/2	1/2	0			
1	1				

For explicit methods, we have $a_{ij} = 0$ for all $j \geq i$, i.e., only k_1, \dots, k_{i-1} are required to compute k_i . For implicit methods, k_i may depend on all k_1, \dots, k_s .

A special case are *semi-implicit* methods, characterized by $a_{ij} = 0$ for $j > i$: k_i may depend on k_1, \dots, k_i , so we still can compute k_1, k_2, \dots, k_s in sequence, but each step requires us to solve an equation.

Classical Runge-Kutta method of fourth order given by

0	0				
1/2	1/2	0			
1/2	0	1/2	0		
1	0	0	1	0	
1/6	1/3	1/3	1/6		

Experiment: Runge-Kutta for the mass-spring system If we use the Crank-Nicolson method and the classic Runge-Kutta method for the mass-spring system starting at time $t = 0$ and stopping at time $t = 20$, we obtain the following results:

δ	Crank-Nicol error	ratio	Runge-Kutta error	ratio
1	5.46 ₋₁		4.07 ₋₂	
1/2	3.24 ₋₁	1.7	6.91 ₋₃	5.9
1/4	9.19 ₋₂	3.5	5.26 ₋₄	13.1
1/8	2.36 ₋₂	3.9	3.52 ₋₅	14.9
1/16	5.93 ₋₃	4.0	2.26 ₋₆	15.6
1/32	1.49 ₋₃	4.0	1.43 ₋₇	15.8
1/64	3.71 ₋₄	4.0	9.02 ₋₉	15.9
1/128	9.29 ₋₅	4.0	5.65 ₋₁₀	16.0
1/256	2.32 ₋₅	4.0	3.54 ₋₁₁	16.0

We can see that the classic Runge-Kutta method indeed shows fourth-order convergence.

Multi-step methods: Runge-Kutta need s evaluations of f in each time step. *Multi-step methods* re-use results from previous steps to reduce the work.

Idea of the *Adams-Bashforth method*: use fundamental theorem of calculus to find

$$y(t_{i+1}) = y(t_i) + \int_{t_i}^{t_{i+1}} y'(s) ds,$$

replace y' by interpolating polynomial in $t_i, t_{i-1}, \dots, t_{i-m+1}$.

$$\begin{aligned} y(t_{i+1}) &\approx y(t_i) + \int_{t_i}^{t_{i+1}} p(s) ds = y(t_i) + \sum_{j=0}^{m-1} y'(t_{i-j}) \underbrace{\int_{t_i}^{t_{i+1}} \ell_{i,j}(s) ds}_{=: a_{ij}} \\ &= y(t_i) + \sum_{j=0}^{m-1} a_{ij} f(t_{i-j}, y(t_{i-j})). \end{aligned}$$

Multi-step method stores $f_i := f(t_i, \tilde{y}(t_i))$ to reduce computational work:

$$\begin{aligned} f_i &\leftarrow f(t_i, y(t_i)) && \text{for all } i \in [0 : m-1], \\ \tilde{y}(t_{i+1}) &\leftarrow \tilde{y}(t_i) + \sum_{j=0}^{m-1} a_{ij} f_{i-j}, \\ f_{i+1} &\leftarrow f(t_{i+1}, \tilde{y}(t_{i+1})) && \text{for all } i \in [m-1 : \infty). \end{aligned}$$

Only f_i, \dots, f_{i-m+1} have to be stored, so only m auxiliary vectors are required.

Cyclic storage: Since we only need f_i, \dots, f_{i-m+1} to compute $\tilde{y}(t_{i+1})$, we can subsequently overwrite the “oldest” vector f_{i-m+1} by f_{i+1} . This approach allows us to avoid copying vectors, but we have to be careful to keep track where which vector is stored.

When implementing multi-step methods in programming languages like C, we can use the *modulo* operator (returning the remainder of an integer division) for this purpose: If we have arrays $\mathbf{f}[0], \dots, \mathbf{f}[s-1]$, we keep f_i in the array $\mathbf{f}[i\%s]$.

Equidistant time steps allow us to reduce the number of coefficients. If we choose $t_i = i\delta$, we have

$$\begin{aligned} a_{ij} &= \int_{t_i}^{t_{i+1}} \ell_{i,j}(s) ds = \int_{t_i}^{t_{i+1}} \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{s - t_{i-k}}{t_{i-j} - t_{i-k}} ds \\ &= \int_{i\delta}^{(i+1)\delta} \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{s - i\delta + k\delta}{i\delta - j\delta - i\delta + k\delta} ds = \int_0^\delta \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{k\delta + s}{(k-j)\delta} ds \end{aligned}$$

We can even substitute $s = \hat{s}\delta$ to obtain

$$a_{ij} = \int_0^\delta \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{(k + \hat{s})\delta}{(k - j)\delta} d\hat{s} = \delta \underbrace{\int_0^1 \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{k + \hat{s}}{k - j} d\hat{s}}_{=: \hat{a}_j}$$

with coefficients \hat{a}_j that are independent of δ and i .

These coefficients are the integrals of Lagrange polynomials for the interpolation points $x_0 = 0, x_1 = -1, \dots, x_m = -(m - 1)$.

Multi-step method with equidistant steps:

$$\begin{aligned} f_i &\leftarrow f(t_i, y(t_i)) && \text{for all } i \in [0 : m - 1], \\ \tilde{y}(t_{i+1}) &\leftarrow \tilde{y}(t_i) + \delta \sum_{j=0}^{m-1} \hat{a}_j f_{i-j}, \\ f_{i+1} &\leftarrow f(t_{i+1}, \tilde{y}(t_{i+1})) && \text{for all } i \in [m - 1 : \infty). \end{aligned}$$

Coefficients can be again found by using the fact that interpolation reproduces polynomials exactly, i.e., we have

$$\sum_{j=0}^{m-1} \hat{a}_j p(-j) = \int_0^1 \sum_{j=0}^{m-1} p(-j) \prod_{\substack{k=0 \\ k \neq j}}^{m-1} \frac{k + \hat{s}}{k - j} d\hat{s} = \int_0^1 p(\hat{s}) d\hat{s}$$

for all polynomials p of degree $m - 1$, so we can compute the coefficients by applying this equation to monomials $p(x) = x^i$ and solving the resulting system.

Example: Quadratic Adams-Bashforth method For $m = 3$, we find

$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & -2 \\ 0 & 1 & 4 \end{pmatrix} \begin{pmatrix} \hat{a}_0 \\ \hat{a}_1 \\ \hat{a}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1/2 \\ 1/3 \end{pmatrix}$$

with the solution $\hat{a}_0 = 23/12, \hat{a}_1 = -16/12, \hat{a}_2 = 5/12$, and one step takes the form

$$\begin{aligned} \tilde{y}(t_{i+1}) &\leftarrow \tilde{y}(t_i) + \frac{\delta}{12} (23f_i - 16f_{i-1} + 5f_{i-2}), \\ f_{i+1} &\leftarrow f(t_{i+1}, \tilde{y}(t_{i+1})). \end{aligned}$$

Example: Cubic Adams-Bashforth method For $m = 4$, we find

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & -1 & -2 & -3 \\ 0 & 1 & 4 & 9 \\ 0 & -1 & -8 & -27 \end{pmatrix} \begin{pmatrix} \hat{a}_0 \\ \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1/2 \\ 1/3 \\ 1/4 \end{pmatrix}$$

with the solution $\hat{a}_0 = 55/24, \hat{a}_1 = -59/24, \hat{a}_2 = 37/24, \hat{a}_3 = -9/24$.

Experiment: Convergence for the spring-mass system in the time interval $[0, 20]$.

δ	Runge-Kutta		Adams-Bash 3		Adams-Bash 4	
	error	ratio	error	ratio	error	ratio
1	4.07_{-2}		4.23_{+3}		2.74_{+4}	
1/2	6.91_{-3}	5.9	3.96_{-1}	10685.8	2.20_{-1}	124462.5
1/4	5.26_{-4}	13.1	6.63_{-2}	6.0	1.76_{-2}	12.5
1/8	3.52_{-5}	14.9	7.52_{-3}	8.8	1.36_{-3}	13.0
1/16	2.26_{-6}	15.6	8.50_{-4}	8.8	9.15_{-5}	14.9
1/32	1.43_{-6}	15.8	1.00_{-4}	8.5	5.90_{-6}	15.5
1/64	9.02_{-9}	15.9	1.21_{-5}	8.3	3.74_{-7}	15.8
1/128	5.65_{-10}	16.0	1.49_{-6}	8.1	2.36_{-8}	15.9
1/256	3.54_{-11}	16.0	1.84_{-7}	8.1	1.48_{-9}	15.9

We can see that the errors of the Adams-Bashforth methods converge like δ^3 and δ^4 , respectively, although they only require one evaluation of the right-hand side f per step.

3 Partial differential equations

3.1 Finite difference methods for partial differential equations

Poisson's equation in the unit square $\Omega = (0, 1) \times (0, 1)$:

$$-\Delta u(x) = f(x), \quad \text{for all } x \in \Omega$$

with the *Laplace operator*

$$\Delta := \frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2}.$$

Challenge: Solution depends on multiple variables, and we have derivatives with respect to all of these variables.

→ simple timestepping not possible.

Idea: Essentially, the PDE is a linear system of equations, just with an *infinite* number of degrees of freedom, i.e., the values $u(x)$ for all $x \in \Omega$.

→ turn it into a system with finitely many variables.

Grid: Fix $N \in \mathbb{N}$ and define the *stepsize* $h := \frac{1}{N+1}$.

$$\begin{aligned} \bar{\Omega}_h &:= \{(ih, jh) : i, j \in [0 : N + 1]\} \subseteq \bar{\Omega} && \text{domain including boundary,} \\ \Omega_h &:= \{(ih, jh) : i, j \in [1 : N]\} \subseteq \Omega && \text{interior of the domain,} \\ \partial\Omega_h &:= \bar{\Omega}_h \setminus \Omega_h && \text{boundary of the domain.} \end{aligned}$$

Grid function: Replace a function $u: \bar{\Omega} \rightarrow \mathbb{R}$ by a grid function $u_h: \bar{\Omega}_h \rightarrow \mathbb{R}$.

Since u_h is characterized by only $(N + 2)^2$ values, we have replaced infinitely many variables by finitely many.

Representation: Since Ω_h is two-dimensional, a first idea could be to use a two-dimensional array `double u[N+2][N+2]`. Since the C programming language does not support multidimensional variable-length arrays, this would mean that we have to fix N at compilation time.

In order to be able to choose the grid resolution a runtime, we map the two-dimensional indices $(i, j) \in [0 : N + 1] \times [0 : N + 1]$ to one-dimensional indices $k := i + j(N + 2)$.

Our gridfunction is now represented by an array `u` with $(n + 2)^2$ coefficients, and the value $u_h(ih, jh)$ can be found at `u[i+j*(N+2)]`.

Finite differences: In order to replace u by a grid function u_h , we have to be able to replace the Laplace operator

$$\Delta = \frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2}$$

by an operator acting on grid functions.

We have already seen that we can approximate the second derivative by a difference quotient

$$f''(0) \approx \frac{f(h) - 2f(0) + f(-h)}{h^2},$$

and applying this approximation to the partial derivatives yields

$$\begin{aligned} \frac{\partial^2}{\partial x_1^2} u(x) &\approx \frac{u(x_1 + h, x_2) - 2u(x) + u(x_1 - h, x_2)}{h^2}, \\ \frac{\partial^2}{\partial x_2^2} u(x) &\approx \frac{u(x_1, x_2 + h) - 2u(x) + u(x_1, x_2 - h)}{h^2}, \\ \Delta_h u(x) &:= \frac{u(x_1 + h, x_2) + u(x_1 - h, x_2) + u(x_1, x_2 + h) + u(x_1, x_2 - h) - 4u(x)}{h^2}. \end{aligned}$$

We call this the *discrete Laplace operator*.

Discrete Poisson's equation:

$$-\Delta_h u_h(x) = f(x) \quad \text{for all } x \in \Omega_h.$$

Well-defined, since u_h provides values for all points in $\bar{\Omega}_h$, so all neighbours of a point in Ω_h are available.

Boundary conditions: Discrete equation gives us $n := \#\Omega_h = N^2$ equations for $\#\bar{\Omega}_h = (N+2)^2 = n + 4N + 4$ unknown values, so we cannot expect the solution to be unique.

Easy fix: Add *boundary conditions*, e.g., of *Dirichlet* type.

$$u_h(x) = g(x) \quad \text{for all } x \in \partial\Omega_h.$$

Since the $4N+4$ boundary values are no longer unknown now, we would like to eliminate them from the equation.

General index sets: Given general finite sets \mathcal{I}, \mathcal{J} , we can define general vectors $v = (v_i)_{i \in \mathcal{I}}$, $w = (w_j)_{j \in \mathcal{J}}$ and general matrices $A = (a_{ij})_{i \in \mathcal{I}, j \in \mathcal{J}}$.

The sets of these vectors and matrices are denoted by $\mathbb{R}^{\mathcal{I}}$, $\mathbb{R}^{\mathcal{J}}$, and $\mathbb{R}^{\mathcal{I} \times \mathcal{J}}$, as with standard vectors and matrices.

The matrix-vector multiplication $v = Aw$ is given by

$$v_i = \sum_{j \in \mathcal{J}} a_{ij} w_j \quad \text{for all } i \in \mathcal{I}.$$

Grid functions as vectors: We can see u_h as a vector in $\mathbb{R}^{\bar{\Omega}_h}$ and the operator $-\Delta_h$ as a matrix $L \in \mathbb{R}^{\Omega_h \times \bar{\Omega}_h}$ with the coefficients

$$\ell_{ij} = \begin{cases} 4h^{-2} & \text{if } i = j, \\ -h^{-2} & \text{if } |j_1 - i_1| = h, i_2 = j_2, \\ -h^{-2} & \text{if } |j_2 - i_2| = h, i_1 = j_1, \\ 0 & \text{otherwise} \end{cases} \quad \text{for all } i \in \Omega_h, j \in \bar{\Omega}_h.$$

Elimination of the boundary: Let $A := L|_{\Omega_h \times \Omega_h}$ and $B := L|_{\Omega_h \times \partial\Omega_h}$. We have

$$\begin{aligned} -\Delta u_h = f_h \text{ and } u_h|_{\partial\Omega} = g_h &\iff Lu_h = f_h \text{ and } u_h|_{\partial\Omega} = g_h \\ &\iff Au_h|_{\Omega_h} + Bu_h|_{\partial\Omega_h} = f_h \text{ and } u_h|_{\partial\Omega} = g_h \\ &\implies Au_h|_{\Omega_h} + Bg_h = f_h \\ &\iff Au_h|_{\Omega_h} = f_h - Bg_h \end{aligned}$$

with $f_h = f|_{\Omega_h}$ and $g_h = g|_{\partial\Omega_h}$.

Now A is a square matrix of dimension $n = N^2$, and we can apply standard linear algebra to solve the system.

Solvers:

- Gaussian elimination, LR factorization: Requires explicit elimination of boundary values, adds non-zero entries to the matrix, very slow for high resolutions
- Iterative solvers: Can work with grid functions directly, can handle boundary values elegantly, matrix is not changed, very fast for some applications.
- \mathcal{H} -matrix solvers: Can work with grid functions directly, require elimination of boundary values, very robust for a large class of applications.

Lexicographic order: If we want to represent A as a matrix with one-dimensional index sets, i.e., as $A \in \mathbb{R}^{n \times n}$, we have to describe how to map the indices $k \in [1 : n]$ to the grid points $(ih, jh) \in \Omega_h$.

A standard approach is the *lexicographic order*, i.e., numbering the gridpoints row by row: we let $k = (i - 1) + (j - 1)N + 1$.

This means that the left and right neighbours of a gridpoint $x = (ih, jh)$ with index k are found at $k - 1$ and $k + 1$, while the lower and upper neighbours are found at $k - N$ and $k + N$.

k	x
1	(h, h)
2	$(2h, h)$
\vdots	\vdots
N	(Nh, h)
$N + 1$	$(h, 2h)$
$N + 2$	$(2h, 2h)$
\vdots	\vdots
$2N$	$(Nh, 2h)$
\vdots	\vdots
$(N - 1)N + 1$	(h, Nh)
$(N - 1)N + 2$	$(2h, Nh)$
\vdots	\vdots
N^2	(Nh, Nh)

The entire matrix A can be represented as a *block tridiagonal matrix*

$$A = \frac{1}{h^2} \begin{pmatrix} T & -I & & & \\ -I & \ddots & \ddots & & \\ & \ddots & \ddots & -I & \\ & & -I & T & \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad T = \begin{pmatrix} 4 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & -1 & 4 & \end{pmatrix} \in \mathbb{R}^{N \times N},$$

where the tridiagonal matrix T corresponds to the interactions between grid points in the same row and the identity matrices I correspond to the interactions between adjacent rows.

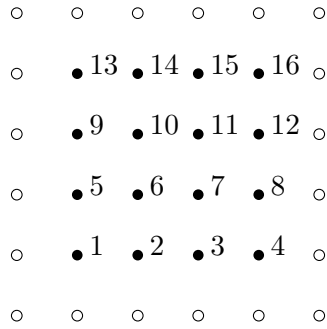


Figure 3.1: Lexicographic enumeration of grid points for $N = 4$

3.2 LR decomposition and BLAS

Goal: Since the discretization of partial differential equations usually leads to large systems of linear equations

$$Ax = b, \quad A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}, \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix},$$

we are interested in finding efficient algorithms for solving these systems (“solvers”).

LR decomposition: Classical direct solver, corresponds to Gaussian elimination. We use the factorization

$$A = LR$$

with triangular matrices

$$L = \begin{pmatrix} 1 & & & \\ \ell_{21} & 1 & & \\ \vdots & \ddots & \ddots & \\ \ell_{n1} & \cdots & \ell_{n,n-1} & 1 \end{pmatrix}, \quad R = \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & \ddots & \vdots \\ & & \ddots & r_{n-1,n} \\ & & & r_{nn} \end{pmatrix}.$$

The diagonal entries of L are always equal to one, so we can store the entire factorization in an $n \times n$ matrix

$$B = \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ \ell_{21} & r_{22} & \ddots & \vdots \\ \vdots & \ddots & \ddots & r_{n-1,n} \\ \ell_{n1} & \cdots & \ell_{n,n-1} & r_{nn} \end{pmatrix}.$$

Not every invertible matrix can be factorized in this way, but if an LR factorization exists, we can find an algorithm that overwrites A by its representation B .

LR construction: Split A , L , and R into the first rows and columns and a remainder.

$$A = \begin{pmatrix} a_{11} & A_{1*} \\ A_{*1} & A_{**} \end{pmatrix}, \quad L = \begin{pmatrix} \ell_{11} & \\ L_{*1} & L_{**} \end{pmatrix}, \quad R = \begin{pmatrix} r_{11} & R_{1*} \\ & R_{**} \end{pmatrix}.$$

The equation $A = LR$ is equivalent to

$$\begin{aligned} a_{11} &= \ell_{11}r_{11}, \\ A_{*1} &= L_{*1}r_{11}, \quad A_{1*} = \ell_{11}R_{1*}, \\ A_{**} &= L_{**}R_{**} + L_{*1}R_{1*}. \end{aligned}$$

The first equation yields $\ell_{11} = 1$, $r_{11} = a_{11}$. The second and third leads to $L_{*1} = A_{*1}/r_{11}$ and $R_{1*} = A_{1*}$. The fourth is equivalent to $L_{**}R_{**} = A_{**} - L_{*1}R_{1*}$, so we have to find an LR factorization of a matrix of dimension $n - 1$. This last task can be handled by recursion.

BLAS: The construction of an LR factorization requires approximately $\frac{2}{3}n^3$ operations, so it can be very time-consuming for larger systems.

In order to make the computation as fast as possible, the manufacturers of processors and compilers provide us with optimized functions for basic operations. The most important standard in this field is *BLAS* (Basic Linear Algebra Subroutines).

BLAS vectors are represented by a pointer \mathbf{x} to the first coefficient, an increment \mathbf{incx} , and the dimension k . We have

$$x_i = \mathbf{x}[\mathbf{i} * \mathbf{incx}] \quad \text{for all } i \in [0 : k - 1].$$

Note that k -dimensional vectors use coefficients from 0 to $k - 1$ instead of 1 to k .

BLAS matrices are represented by a pointer \mathbf{A} to the upper left coefficient, the *leading dimension* \mathbf{ldA} , the number n of rows, and the number m of columns. We have

$$a_{ij} = \mathbf{a}[\mathbf{i} + \mathbf{j} * \mathbf{ldA}] \quad \text{for all } i \in [0 : n - 1], j \in [0 : m - 1].$$

Setting the parameters correctly, we can access rows, columns, diagonals, and submatrices of A :

- i -th row: $\mathbf{x} = \mathbf{A} + \mathbf{i}$, $\mathbf{incx} = \mathbf{ldA}$, $\mathbf{k} = \mathbf{m}$
- j -th column: $\mathbf{x} = \mathbf{A} + \mathbf{j} * \mathbf{ldA}$, $\mathbf{incx} = 1$, $\mathbf{k} = \mathbf{n}$
- Diagonal: $\mathbf{x} = \mathbf{A}$, $\mathbf{incx} = \mathbf{ldA} + 1$, $\mathbf{k} = \min(\mathbf{n}, \mathbf{m})$
- Subdiagonal: $\mathbf{x} = \mathbf{A} + 1$, $\mathbf{incx} = \mathbf{ldA} + 1$, $\mathbf{k} = \min(\mathbf{n}, \mathbf{m}) - 1$
- Submatrix $C = A_{**}$: $\mathbf{C} = \mathbf{A} + 1 + \mathbf{ldA}$, $\mathbf{ldC} = \mathbf{ldA}$

BLAS functions are separated into three levels.

Level 1 functions are for actions on vectors like

- scaling $x \leftarrow \alpha x$,
`void scal(int n, real alpha, real *x, int incx)`
- adding $y \leftarrow y + \alpha x$.
`void axpy(int n, real alpha, const real *x, int incx,
real *y, int incy)`
- or computing the inner product $\langle x, y \rangle_2 = \sum_{i=1}^n x_i y_i$.
`real dot(int n, const real *x int incx,
const real *y, int incy)`

Level 2 functions are for interactions between vectors and matrices like

- the matrix-vector product $y \leftarrow y + \alpha Ax$,

```
void gemv(bool trans, int rows, int cols,
          real alpha, const real *A, int ldA,
          const real *x, int incx,
          real *y, int incy)
```
- or the rank-one update $A \leftarrow A + \alpha xy^T$.

```
void ger(int rows, int cols, real alpha,
          const real *x, int incx,
          const real *y, int incy,
          real *A, int ldA)
```

Level 3 functions are for interactions between matrices. We only need the most general and most important one,

- the matrix-matrix product $C \leftarrow \alpha AB + \beta C$.

```
void gemm(bool transA, bool transB,
          int rows, int cols, int k,
          real alpha, const real *A, int ldA,
          const real *B, int ldB,
          real beta, real *C, int ldC)
```

Examples: Programming with BLAS We can realize level-2 functions using only level-1 functions:

```
void
gemv(bool trans, int rows, int cols,
     real alpha, const real *A, int ldA,
     const real *x, int incx,
     real *y, int incy)
{
    int i;

    if(trans) {
        for(i=0; i<cols; i++)
            y[i*incy] += alpha * dot(rows, A+i*ldA, 1, x, incx);
    }
    else
        for(i=0; i<cols; i++)
            axpy(rows, alpha * x[i*incx], A+i*ldA, 1, y, incy);
}
```

In order to obtain the best possible performance, it is usually a good idea to ensure that we run through coefficients in the order in which they are stored in main memory, i.e., to try to work with the *columns* of matrices instead of rows.

```

void
ger(int rows, int cols, real alpha,
    const real *x, int incx,
    const real *y, int incy,
    real *A, int ldA)
{
    int i;

    for(i=0; i<cols; i++)
        axpy(rows, alpha * y[i*incy], x, incx, A+i*ldA, 1);
}

```

Example: LR factorization with BLAS

```

void
lrdecomp(int n, real *A, int ldA)
{
    int k;

    for(k=0; k+1<n; k++) {
        scal(n-k-1, 1.0/A[k+k*ldA], A+(k+1)+k*ldA, 1);
        ger(n-k-1, n-k-1, -1.0,
            A+(k+1)+k*ldA, 1,
            A+k+(k+1)*ldA, ldA,
            A+(k+1)+(k+1)*ldA, ldA);
    }
}

```

The call to `scal` scales A_{*1} by $1/r_{11} = 1/a_{11}$ to obtain L_{*1} , the call to `ger` subtracts $L_{*1}R_{1*}$ from A_{**} . Instead of performing a “true” recursion, we simply use the variable k to keep track of the recursion level, i.e., of how many rows and columns have already been computed.

Solving with the LR factorization: Once we have the factorization $A = LR$ at our disposal, $Ax = b$ is equivalent to $LRx = b$, and with $y := Rx$ we find

$$Ly = b, \quad Rx = y.$$

Forward substitution: To solve $Ly = b$, we consider

$$L = \begin{pmatrix} \ell_{11} & \\ L_{*1} & L_{**} \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_* \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ b_* \end{pmatrix},$$

and due to $\ell_{11} = 1$, $Ly = b$ is equivalent to

$$y_1 = \ell_{11}y_1 = b_1, \quad L_{**}y_* + L_{*1}y_1 = b_*.$$

It is customary to overwrite the coefficients of b successively with those of y .

```
void
lsolve(int n, const real *L, int ldL, real *b, int incb)
{
    int i;

    for(i=0; i+1<n; i++)
        axpy(n-i-1, -b[i*incb], L+(i+1)+i*ldL, 1, b+i+1, incb);
}
```

We can easily generalize this algorithm to solve $LY = B$ with $n \times m$ matrices X and B :

```
void
lsolve(int n, int m, const real *L, int ldL, real *B, int ldB)
{
    int i;

    for(i=0; i+1<n; i++)
        ger(n-i-1, m, -1.0,
            L+(i+1)+i*ldL, 1, B+i, ldB,
            B+(i+1), ldB);
}
```

Backward substitution: To solve $Rx = b$, we consider

$$R = \begin{pmatrix} R_{**} & R_{*n} \\ & r_{nn} \end{pmatrix}, \quad x = \begin{pmatrix} x_* \\ x_n \end{pmatrix}, \quad y = \begin{pmatrix} y_* \\ y_n \end{pmatrix},$$

and $Rx = y$ is equivalent to

$$r_{nn}y_n = x_n, \quad R_{**}x_* + R_{*n}x_n = y_*.$$

Again overwriting y with x leads to

```
void
rsolve(int n, const real *R, int ldR, real *y, int incy)
{
    int i;

    for(i=n; i-->0; ) {
        y[i*incy] /= R[i+i*ldR];
        axpy(i, -y[i*incy], R+i*ldR, 1, y, incy);
    }
}
```

For $RX = Y$ with $n \times m$ matrices X and Y , we obtain

```

void
rsolve(int n, int m, const real *R, int ldR, real *Y, int ldY)
{
    int i;

    for(i=n; i-->0; ) {
        scal(m, 1.0/R[i+i*ldR], Y+i, ldY);
        ger(i, m, -1.0,
            R+i*ldR, 1, Y+i, ldY,
            Y, ldY);
    }
}

```

Caching: Modern processors do not access main memory directly, but via *caches*, small and fast memory units that are typically contained directly on the chip.

The caches contain copies of addresses in main memory and allow the processor to avoid accessing main memory if the required data is already present in the fast cache.

It is generally a good idea to arrange programs in a way that works as frequently as possible with cached data. Concerning linear algebra, using level 3 BLAS is a good strategy: operations like matrix multiplication or LR factorization perform $\mathcal{O}(n^3)$ operations, but require only $\mathcal{O}(n^2)$ units of storage, so there is a chance to ensure that multiple operations are carried out on data moved into the cache.

Caches are typically organized in *cache lines*, frequently of 64 bytes, that are treated as indivisible. This means that even if just one byte is required, 64 bytes will be transferred, so it is a good idea to ensure that consecutive memory addresses are used consecutively in the program.

3.3 Method of lines for the heat equation and the wave equation

Heat equation in the domain $\Omega = (0, 1) \times (0, 1)$:

$$\begin{aligned} \frac{\partial u}{\partial t}(t, x) &= c\Delta u(t, x) && \text{for all } t \in \mathbb{R}, x \in \Omega, \\ u(t, x) &= 0 && \text{for all } t \in \mathbb{R}, x \in \partial\Omega. \end{aligned}$$

Looks like an ODE for $t \mapsto u(t, \cdot)$ at first glance, but function-valued instead of vector-valued.

Method of lines: Discretize in space, i.e., replace Ω by grid Ω_h .

$$\begin{aligned} \frac{\partial u_h}{\partial t}(t, x) &= c\Delta_h u_h(t, x) && \text{for all } t \in \mathbb{R}, x \in \Omega_h, \\ u_h(t, x) &= 0 && \text{for all } t \in \mathbb{R}, x \in \partial\Omega_h. \end{aligned}$$

Change of perspective: In order to be able to apply ODE algorithms, consider u_h a mapping from time into the space of grid functions with homogeneous boundary conditions, $u_h(t, x) \equiv u_h(t)(x)$.

$$u_h'(t) = c\Delta_h u_h(t) \quad \text{for all } t \in \mathbb{R}.$$

CFL condition: If we use the explicit Euler method, we have

$$\tilde{u}_h(t_{k+1}) \leftarrow \tilde{u}_h(t_k) + \delta c \Delta_h \tilde{u}_h(t_k) = (I + \delta c \Delta_h) \tilde{u}_h(t_k) \quad \text{for all } k \in \mathbb{N}_0.$$

Given $\nu, \mu \in [1 : N]$, consider the grid function

$$e_{\nu\mu}(x) := \sin(\pi\nu x_1) \sin(\pi\mu x_2).$$

Straightforward computation with the identities $\sin(\alpha + \beta) = \sin(\alpha) \cos(\beta) + \cos(\alpha) \sin(\beta)$ and $\cos(\alpha) = 1 - 2 \sin^2(\alpha/2)$ yield

$$-\Delta_h e_{\nu\mu} = \lambda_{\nu\mu} e_{\nu\mu}, \quad \lambda_{\nu\mu} = 4h^{-2}(\sin^2(\pi\nu h/2) + \sin^2(\pi\mu h/2)).$$

If we start the simulation with $u_h(t_0) = e_{\nu\mu}$, we find

$$u_h(t_k) = \exp(-c\lambda_{\nu\mu} t_k) e_{\nu\mu}, \quad \tilde{u}_h(t_k) = (1 - \delta c \lambda_{\nu\mu})^k e_{\nu\mu},$$

and in particular the Euler approximation will oscillate rapidly if $1 - \delta c \lambda_{\nu\mu} < 0$, and even grows exponentially if $1 - \delta c \lambda_{\nu\mu} < -1$, while the continuous solution converges smoothly and exponentially.

For the explicit Euler method, we can only avoid this problem by choosing the stepsize δ small enough to fulfill the *Courant-Friedrichs-Levy condition*

$$\delta c 8h^{-2} \leq 1, \quad \delta \leq \frac{h^2}{8c},$$

since $\lambda_{NN} \approx 8h^{-2}$, i.e., the number of timesteps has to grow rapidly if the grid is refined.

Experiment: Convergence for the heat equation discretized with a mesh parameter of $h = 1/65$ in the time interval $[0, 1]$.

δ	$8c\delta/h^2$	Expl. Euler		Runge	
		error	ratio	error	ratio
1/16	2.1 ₊₃	1.97 ₊₃₉		3.11 ₊₈₇	
1/32	1.1 ₊₃	4.58 ₊₈₂		4.12 ₊₁₆₉	
1/64	5.3 ₊₂	8.03 ₊₁₅₉		∞	
1/16384	2.1 ₊₀	∞		∞	
1/32768	1.0 ₊₀	2.73 ₋₈		5.49 ₋₁₂	
1/65536	5.2 ₋₁	1.36 ₋₈	2.0	1.37 ₋₁₂	4.0
1/131072	2.6 ₋₁	6.83 ₋₉	2.0	3.43 ₋₁₃	4.0
1/262144	1.3 ₋₁	3.41 ₋₉	2.0	8.57 ₋₁₄	4.0

As long as the time steps are too large, both methods diverge. The expected first- and second-order convergence sets in once the time steps are small enough.

Implicit methods: All explicit methods need to fulfill a CFL condition in order to avoid unrealistic oscillations. Implicit methods are better, e.g., the implicit Euler scheme

$$\begin{aligned}\tilde{u}_h(t_{k+1}) &= \tilde{u}_h(t_k) + \delta c \Delta_h \tilde{u}_h(t_{k+1}), \\ (I - \delta c \Delta_h) \tilde{u}_h(t_{k+1}) &= \tilde{u}_h(t_k).\end{aligned}$$

Now we have

$$u_h(t_k) = \exp(-c\lambda_{\nu\mu} t_k) e_{\nu\mu}, \quad \tilde{u}_h(t_k) = \left(\frac{1}{1 + \delta c \lambda_{\nu\mu}} \right)^k e_{\nu\mu},$$

and since $1 + \delta c \lambda$ is always positive, the Euler solution converges exponentially, just like the continuous solution.

For the Crank-Nicolson method, we obtain a factor of $\frac{2 - \delta c \lambda_{\nu\mu}}{2 + \delta c \lambda_{\nu\mu}}$, i.e., there will be oscillations if $\delta c \lambda_{\nu\mu} > 2$, but the approximation will not “blow up”.

Experiment: Convergence for the heat equation in the time interval $[0, 1]$.

δ	$8c\delta/h^2$	Impl. Euler		Crank-Nicolson	
		error	ratio	error	ratio
1/16	2.1 ₊₃	4.46 ₋₃		4.43 ₋₆	
1/32	1.1 ₊₃	3.56 ₋₄	12.5	2.23 ₋₆	2.0
1/64	5.3 ₊₂	5.33 ₋₅	6.7	6.74 ₋₇	3.3
1/128	2.6 ₊₂	1.37 ₋₅	3.9	1.77 ₋₇	3.8
1/256	1.3 ₊₂	4.88 ₋₆	2.8	4.48 ₋₈	4.0
1/512	6.6 ₊₁	2.06 ₋₆	2.4	1.12 ₋₈	4.0
1/1024	3.3 ₊₁	9.50 ₋₇	2.2	2.81 ₋₉	4.0
1/2048	1.7 ₊₁	4.56 ₋₇	2.1	7.02 ₋₁₀	4.0

We can see that both implicit methods converge at the expected rate. Although implicit methods require the often computationally expensive task of solving linear systems, the significantly better convergence and stability still makes them very attractive.

Wave equation in the domain $\Omega = (0, 1) \times (0, 1)$:

$$\begin{aligned}\frac{\partial^2 u}{\partial t^2}(t, x) &= c\Delta u(t, x) && \text{for all } t \in \mathbb{R}, x \in \Omega, \\ u(t, x) &= 0 && \text{for all } t \in \mathbb{R}, x \in \partial\Omega.\end{aligned}$$

Order reduction: So far, we have only considered *first-order* ordinary differential equations, now we have a *second* derivative in time.

Introducing the first time derivative as an auxiliary variable, in this case the *velocity*, yields a first-order system

$$\begin{aligned}\frac{\partial u}{\partial t}(t, x) &= v(t, x), \\ \frac{\partial v}{\partial t}(t, x) &= c\Delta u(t, x) && \text{for all } t \in \mathbb{R}, x \in \Omega, \\ u(t, x) &= v(t, x) = 0 && \text{for all } t \in \mathbb{R}, x \in \partial\Omega.\end{aligned}$$

Method of lines: Replace functions by time-dependent grid functions

$$\begin{aligned}u'_h(t) &= v_h(t), \\ v'_h(t) &= c\Delta_h u_h(t) && \text{for all } t \in \mathbb{R}.\end{aligned}$$

Leapfrog timestepping:

$$\begin{aligned}\tilde{u}_h(t_0) &= u_h(t_0), \\ \tilde{v}_h(t_1) &= v_h(t_0) + \delta c\Delta_h \tilde{u}_h(t_0), \\ \tilde{u}_h(t_{2\ell+2}) &= \tilde{u}_h(t_{2\ell}) + 2\delta \tilde{v}_h(t_{2\ell+1}), \\ \tilde{v}_h(t_{2\ell+3}) &= \tilde{v}_h(t_{2\ell+1}) + 2\delta c\Delta_h \tilde{u}_h(t_{2\ell+2}) && \text{for all } \ell \in \mathbb{N}_0.\end{aligned}$$

Crank-Nicolson timestepping: The trapezoidal rule leads to

$$\begin{aligned}\tilde{u}_h(t_{i+1}) &= \tilde{u}_h(t_i) + \frac{\delta}{2}(\tilde{v}_h(t_i) + \tilde{v}_h(t_{i+1})), \\ \tilde{v}_h(t_{i+1}) &= \tilde{v}_h(t_i) + \frac{\delta}{2}c(\Delta_h \tilde{u}_h(t_i) + \Delta_h \tilde{u}_h(t_{i+1})).\end{aligned}$$

Substituting the second equation in the first leads to

$$\begin{aligned}\tilde{u}_h(t_{i+1}) &= \tilde{u}_h(t_i) + \frac{\delta}{2}(\tilde{v}_h(t_i) + \tilde{v}_h(t_i) + \frac{\delta}{2}c(\Delta_h \tilde{u}_h(t_i) + \Delta_h \tilde{u}_h(t_{i+1}))) \\ &= \tilde{u}_h(t_i) + \delta \tilde{v}_h(t_i) + \frac{\delta^2}{4}c\Delta_h \tilde{u}_h(t_i) + \frac{\delta^2}{4}c\Delta_h \tilde{u}_h(t_{i+1}), \\ \left(I - \frac{\delta^2}{4}c\Delta_h\right) \tilde{u}_h(t_{i+1}) &= \left(I + \frac{\delta^2}{4}c\Delta_h\right) \tilde{u}_h(t_i) + \delta \tilde{v}_h(t_i),\end{aligned}$$

while substituting the first in the second leads to

$$\begin{aligned}\tilde{v}_h(t_{i+1}) &= \tilde{v}_h(t_i) + \frac{\delta}{2}c(\Delta_h \tilde{u}_h(t_i) + \Delta_h(\tilde{u}_h(t_i) + \frac{\delta}{2}(\tilde{v}_h(t_i) + \tilde{v}_h(t_{i+1})))) \\ &= \tilde{v}_h(t_i) + \delta c\Delta_h \tilde{u}_h(t_i) + \frac{\delta^2}{4}c\Delta_h \tilde{v}_h(t_i) + \frac{\delta^2}{4}c\Delta_h \tilde{v}_h(t_{i+1}), \\ \left(I - \frac{\delta^2}{4}c\Delta_h\right) \tilde{v}_h(t_{i+1}) &= \left(I + \frac{\delta^2}{4}c\Delta_h\right) \tilde{v}_h(t_i) + \delta c\Delta_h \tilde{u}_h(t_i).\end{aligned}$$

Experiment: Convergence for the wave equation with a mesh parameter $h = 1/65$ in the time interval $[0, 1]$ using implicit methods.

δ	Leapfrog		Impl. Euler		Crank-Nicol.	
	error	ratio	error	ratio	error	ratio
1/16	3.10 ₊₁₉		1.07 ₊₂		4.64 ₊₁	
1/32	1.63 ₊₃₃		8.60 ₊₁	1.2	1.17 ₊₁	4.0
1/64	4.19 ₊₃₅		5.50 ₊₁	1.6	2.94 ₊₀	4.0
1/128	3.68 ₋₁		3.11 ₊₁	1.8	7.36 ₋₁	4.0
1/256	9.20 ₋₂	4.0	1.66 ₊₁	1.9	1.84 ₋₁	4.0
1/512	2.30 ₋₂	4.0	8.53 ₊₀	1.9	4.60 ₋₂	4.0
1/1024	5.75 ₋₃	4.0	4.33 ₊₀	2.0	1.15 ₋₂	4.0
1/2048	1.44 ₋₃	4.0	2.18 ₊₀	2.0	2.87 ₋₃	4.0

The explicit leapfrog method diverges as long as the time steps are too large, while both implicit methods show the expected convergence for all time steps.

Energy: Using the inner product

$$\langle u_h, v_h \rangle_{\Omega_h} := h^2 \sum_{x \in \Omega_h} u_h(x) v_h(x),$$

the *energy* of our system is given by

$$E_h(u_h, v_h) := \frac{1}{2} \langle v_h, v_h \rangle_{\Omega_h} - \frac{c}{2} \langle u_h, \Delta_h u_h \rangle_{\Omega_h}.$$

Due to

$$\begin{aligned} \frac{\partial}{\partial t} E_h(u_h(t), v_h(t)) &= \langle v'_h(t), v_h(t) \rangle_{\Omega_h} - c \langle u'_h(t), \Delta_h u_h(t) \rangle_{\Omega_h} \\ &= c \langle v_h(t), \Delta_h u_h(t) \rangle_{\Omega_h} - \langle c \Delta_h u_h(t), v_h(t) \rangle_{\Omega_h} = 0, \end{aligned}$$

the energy of the continuous system is constant.

For the Euler methods, the energy increases with each step of the explicit method and decreases with each of the implicit one.

For the Crank-Nicolson method, the “arithmetic mean” of both Euler methods, the energy is constant:

$$\begin{aligned} &c \langle u_h(t_{k+1}), \Delta_h u_h(t_{k+1}) \rangle_{\Omega_h} - c \langle u_h(t_k), \Delta_h u_h(t_k) \rangle_{\Omega_h} \\ &= c \langle u_h(t_{k+1}) + u_h(t_k), \Delta_h (u_h(t_{k+1}) - u_h(t_k)) \rangle_{\Omega_h} \\ &= c \frac{\delta}{2} \langle u_h(t_{k+1}) + u_h(t_k), \Delta_h (v_h(t_{k+1}) + v_h(t_k)) \rangle_{\Omega_h} \\ &= \frac{\delta}{2} \langle c \Delta_h (u_h(t_{k+1}) + u_h(t_k)), v_h(t_{k+1}) + v_h(t_k) \rangle_{\Omega_h} \\ &= \langle v_h(t_{k+1}) - v_h(t_k), v_h(t_{k+1}) + v_h(t_k) \rangle_{\Omega_h} \\ &= \langle v_h(t_{k+1}), v_h(t_{k+1}) \rangle_{\Omega_h} - \langle v_h(t_k), v_h(t_k) \rangle_{\Omega_h}. \end{aligned}$$

In the third step, we have used the fact that Δ_h is *self-adjoint*, i.e., that $\langle x, \Delta_h y \rangle_{\Omega_h} = \langle \Delta_h x, y \rangle_{\Omega_h}$ holds for all grid functions x and y .

4 Iterative and semi-iterative solvers

4.1 Iterations

In order to be able to use implicit timestepping methods for nonlinear ODEs, we need algorithms for solving nonlinear systems of equations.

Frequently it is not possible to compute the exact solution of an equation in a finite number of steps. In this situation, *iterative methods* are employed that compute a sequence of increasingly accurate approximate solutions.

Newton method: Approximates zeros of a function $f \in C^2(\mathbb{R})$. Let x_* satisfy $f(x_*) = 0$. Taylor expansion yields

$$\begin{aligned} 0 &= f(x_*) = f(x) + (x_* - x)f'(x) + \frac{(x_* - x)^2}{2}f''(\eta), \\ 0 &= \frac{f(x)}{f'(x)} + x_* - x + \frac{f''(\eta)}{2f'(x)}(x_* - x)^2, \\ x_* &= x - \frac{f(x)}{f'(x)} - \frac{f''(\eta)}{2f'(x)}(x_* - x)^2 \end{aligned}$$

for an $\eta \in \mathbb{R}$, and neglecting the rightmost term yields

$$x_* \approx x - \frac{f(x)}{f'(x)}.$$

Iteration: We start with an *initial guess* x_0 and compute a sequence

$$x_{m+1} \leftarrow x_m - \frac{f(x_m)}{f'(x_m)} \quad \text{for all } m \in \mathbb{N}_0.$$

Quadratic convergence: Assume $c := \max \left\{ \frac{|f''(\eta)|}{2|f'(x)|} : x, \eta \in \mathbb{R} \right\}$ is finite. We have

$$|x_* - x_{m+1}| = \frac{|f''(\eta)|}{2|f'(x_m)|} |x_* - x_m|^2 \leq c|x_* - x_m|^2 \quad \text{for all } m \in \mathbb{N}_0.$$

If we assume $c|x_* - x_0| \leq 1/2$, we obtain

$$\begin{aligned} c|x_* - x_1| &\leq c^2|x_* - x_0|^2 \leq \frac{1}{4}, \\ c|x_* - x_2| &\leq c^2|x_* - x_1|^2 \leq \frac{1}{16}, \end{aligned}$$

$$c|x_* - x_3| \leq c^2|x_* - x_2|^2 \leq \frac{1}{256},$$

$$c|x_* - x_m| \leq 2^{-2^m} \quad \text{for all } m \in \mathbb{N}_0,$$

i.e., each step of the iteration will *double* the number of correct binary digits in the approximation.

Example: Square root The square root of $a \in \mathbb{R}_{>0}$ is usually computed by Newton's iteration. If we define it as the positive zero of $f(x) = x^2 - a$, we obtain $f'(x) = 2x$ and

$$x_{m+1} \leftarrow x_m - \frac{f(x_m)}{f'(x_m)} = x_m - \frac{x_m^2 - a}{2x_m} = \frac{x_m^2 + a}{2x_m} = \frac{1}{2} \left(x_m + \frac{a}{x_m} \right).$$

Applied to $a = 2$ and different initial guesses x_0 , we obtain the following errors:

m	$x_0 = 2$	$x_0 = 10$	$x_0 = -5$
1	8.58 ₋₂	3.69 ₊₀	1.29 ₊₀
2	2.45 ₋₃	1.33 ₊₀	3.06 ₋₁
3	2.12 ₋₆	3.23 ₋₁	2.72 ₋₂
4	1.59 ₋₁₂	3.00 ₋₂	2.57 ₋₄
5		3.12 ₋₄	2.34 ₋₈
6		3.44 ₋₈	2.22 ₋₁₆
7		2.22 ₋₁₆	

Example: Reciprocal square root The reciprocal square root $1/\sqrt{a}$ of $a \in \mathbb{R}_{>0}$ is the positive zero of $f(x) = x^{-2} - a$. With $f'(x) = -2x^{-3}$, we find

$$x_{m+1} \leftarrow x_m - \frac{f(x_m)}{f'(x_m)} = x_m - \frac{x_m^{-2} - a}{-2x_m^{-3}} = x_m + \frac{1}{2}(x_m - ax_m^3) = \frac{x_m}{2}(3 - ax_m^2).$$

Subtraction and multiplication are sufficient to compute the approximation, and these operations are *significantly* faster than divisions on modern processors.

Applied to $a = 2$ and different initial values x_0 , we obtain the following errors:

m	$x_0 = 1$	$x_0 = 1.5$	$x_0 = 2$
1	2.07 ₋₁	4.18 ₋₁	4.29 ₊₀
2	8.21 ₋₂	4.43 ₋₁	1.18 ₊₂
3	1.37 ₋₂	3.30 ₋₁	1.62 ₊₆
4	3.98 ₋₄	1.95 ₋₁	4.27 ₊₁₈
5	3.36 ₋₇	7.32 ₋₂	7.77 ₊₅₅
6	2.40 ₋₁₃	1.10 ₋₂	4.70 ₊₁₆₇
7	1.11 ₋₁₆	2.55 ₋₄	∞
8	0	1.38 ₋₇	

We can see that an unsuitable initial guess may lead to divergence. In practice, precomputed lookup tables or the ‘‘Carmack trick’’ can be used to obtain initial guesses that ensure convergence.

Multidimensional Newton: For multidimensional functions $f \in C^2(\mathbb{R}^n, \mathbb{R}^n)$, we can use the *Jacobi matrix*

$$Df(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(x) & \cdots & \frac{\partial f_1}{\partial x_n}(x) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(x) & \cdots & \frac{\partial f_n}{\partial x_n}(x) \end{pmatrix}$$

to define the *multidimensional Newton iteration*

$$\begin{aligned} x_{m+1} &\leftarrow x_m - Df(x_m)^{-1}f(x_m), \\ x_{m+1} &\leftarrow x_m + p_m, \quad Df(x_m)p_m = -f(x_m). \end{aligned}$$

Every step of the algorithm requires us to solve a linear system with the Jacobi matrix, and the Jacobi matrix can be expected to change between steps.

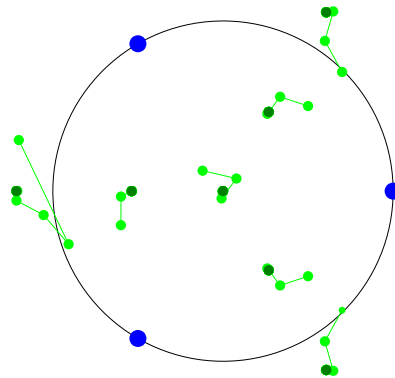
Example: Lagrange points Consider n planets with masses $m_1, \dots, m_n \in \mathbb{R}_{>0}$ at coordinates $y_1, \dots, y_n \in \mathbb{R}^2$ in a rotating coordinate system with angular velocity $\alpha \in \mathbb{R}$. The sum of gravitational forces and centrifugal force is given by

$$f(x) = \alpha^2 x + \gamma \sum_{j=1}^n m_j \frac{y_j - x}{\|y_j - x\|^3}.$$

Lagrange points are characterized by $f(x) = 0$ and can be found with the Newton iteration using the Jacobi matrix

$$Df(x) = \alpha^2 I + \gamma \sum_{j=1}^n 3m_j \frac{(y_j - x)(y_j - x)^T}{\|y_j - x\|^5} - \frac{m_j}{\|y_j - x\|^3} I.$$

Based on the chosen initial guess, the Newton iteration finds different Lagrange points.



Eigenvalue problems: Given $A \in \mathbb{R}^{n \times n}$, we are looking for $x \in \mathbb{R}^n \setminus \{0\}$ and $\lambda \in \mathbb{R}$ with $Ax = \lambda x$.

Searching for eigenvalues is equivalent to finding zeros of polynomials, so in principle the Newton iteration could be applied. In practice, there are better algorithms.

Power iteration: Assume that A is diagonalizable, i.e., that there is an invertible matrix T with

$$T^{-1}AT = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|.$$

Given $x_0 \in \mathbb{R}^n \setminus \{0\}$, consider the sequence

$$x_{m+1} \leftarrow Ax_m \quad \text{for all } m \in \mathbb{N}_0.$$

If we define

$$\hat{x}_m := T^{-1}x_m \quad \text{for all } m \in \mathbb{N}_0,$$

we find

$$\hat{x}_{m+1} = T^{-1}Ax_m = T^{-1}AT\hat{x}_m = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \hat{x}_m \quad \text{for all } m \in \mathbb{N}_0.$$

This leads us to

$$\hat{x}_m = \begin{pmatrix} \lambda_1^m \hat{x}_{0,1} \\ \lambda_2^m \hat{x}_{0,2} \\ \vdots \\ \lambda_n^m \hat{x}_{0,n} \end{pmatrix} \quad \text{for all } m \in \mathbb{N}_0,$$

i.e., if $|\lambda_1| > |\lambda_2|$ and $\hat{x}_{0,1} \neq 0$, the first component of \hat{x}_m will dominate for $m \rightarrow \infty$.

For sufficiently large m , we conclude

$$\begin{aligned} Ax_m &= TT^{-1}AT\hat{x}_m = T \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix} \begin{pmatrix} \hat{x}_{m,1} \\ \hat{x}_{m,2} \\ \vdots \\ \hat{x}_{m,n} \end{pmatrix} \\ &\approx T \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix} \begin{pmatrix} \hat{x}_{m,1} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \lambda_1 T \begin{pmatrix} \hat{x}_{m,1} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \approx \lambda_1 T\hat{x}_m = \lambda_1 x_m, \end{aligned}$$

so x_m is close to an eigenvector for the eigenvalue λ_1 .

Practical power iteration: In order to avoid an overflow, we work with unit vectors:

$$y_{m+1} \leftarrow Ax_m, \quad x_{m+1} \leftarrow \frac{y_{m+1}}{\|y_{m+1}\|} \quad \text{for all } m \in \mathbb{N}_0.$$

Rayleigh quotient: The power iteration approximates an eigenvector, but we usually also require the corresponding eigenvalue. We can recover it using the *Rayleigh quotient*

$$\Lambda_A(x) := \frac{\langle x, Ax \rangle}{\langle x, x \rangle}.$$

If x is an eigenvector for the eigenvalue λ , we have

$$\Lambda_A(x) = \frac{\langle x, Ax \rangle}{\langle x, x \rangle} = \frac{\langle x, \lambda x \rangle}{\langle x, x \rangle} = \lambda \frac{\langle x, x \rangle}{\langle x, x \rangle} = \lambda.$$

We can check the convergence of the power iteration by comparing Ax and $\Lambda_A(x)x$. A good stopping criterion is

$$\|Ax - \Lambda_A(x)x\| \leq \epsilon |\Lambda_A(x)| \|x\|$$

for a given relative accuracy $\epsilon \in \mathbb{R}_{>0}$.

Inverse iteration: Frequently we are not interested in the largest eigenvalue, but in the smallest one, or the eigenvalue closest to a given $\mu \in \mathbb{R}$. If x is an eigenvector for the eigenvalue λ and if μ is not an eigenvalue, we have

$$\begin{aligned} Ax = \lambda x &\iff Ax - \mu x = (\lambda - \mu)x \iff (A - \mu I)x = (\lambda - \mu)x \\ &\iff x = (\lambda - \mu)(A - \mu I)^{-1}x \iff \frac{1}{\lambda - \mu}x = (A - \mu I)^{-1}x. \end{aligned}$$

The eigenvalue that is closest to μ is therefore the eigenvalue of maximal absolute value of $(A - \mu I)^{-1}$, and we can apply the power iteration to the latter matrix to obtain

$$y_{m+1} \leftarrow (A - \mu I)^{-1}x_m, \quad x_{m+1} \leftarrow \frac{y_{m+1}}{\|y_{m+1}\|} \quad \text{for all } m \in \mathbb{N}_0.$$

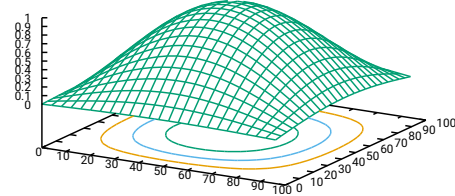
This is called the *inverse iteration* with the *shift* μ , and we can expect it to converge to an eigenvector for the eigenvalue closest to μ .

Rayleigh iteration: The inverse iteration converges quickly if the shift parameter μ is close to an eigenvalue. Given a suitable initial guess, we can use the Rayleigh quotient to provide us with approximations of an eigenvalue that can be used as shifts:

$$y_{m+1} \leftarrow (A - \mu_m I)^{-1}x_m, \quad x_{m+1} \leftarrow \frac{y_{m+1}}{\|y_{m+1}\|}, \quad \mu_{m+1} \leftarrow \Lambda_A(x_{m+1}) \quad \text{for all } m \in \mathbb{N}_0.$$

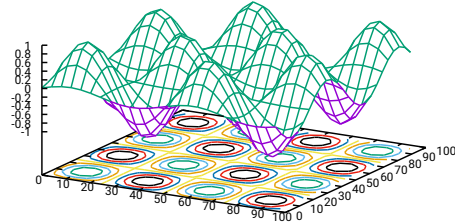
Example: Discrete Laplace operator In order to compute eigenvalues of the discrete Laplace operator $-\Delta_h$ on a grid with the mesh parameter $h = 1/101$, we can use the inverse iteration with a random initial guess.

m	$\sin \alpha_m$	$\tan \alpha_m$	ratio
0	1.00 ₊₀	6.97 ₊₁	
1	3.87 ₋₁	4.19 ₋₁	166.3
2	1.20 ₋₁	1.20 ₋₁	3.5
3	4.46 ₋₂	4.46 ₋₂	2.7
4	1.73 ₋₂	1.73 ₋₂	2.6
5	6.86 ₋₃	6.86 ₋₃	2.5
6	2.73 ₋₃	2.73 ₋₃	2.5
7	1.09 ₋₃	1.09 ₋₃	2.5
8	4.36 ₋₄	4.36 ₋₄	2.5



Example: Rayleigh iteration If we use the inverse iteration starting with a shift $\mu_0 = 310$, we obtain significantly faster convergence:

m	μ	$\sin \alpha_m$	$\tan \alpha_m$	ratio
0	3.10 ₊₂	1.00 ₊₀	5.43 ₊₁	
1	3.20 ₊₂	2.03 ₋₁	2.07 ₋₁	261.9
2	3.15 ₊₂	4.10 ₋₂	4.10 ₋₂	5.1
3	3.15 ₊₂	3.85 ₋₅	3.84 ₋₅	1067.6
4	3.15 ₊₂	3.06 ₋₁₃	3.06 ₋₁₃	1.25 ₊₈



4.2 Krylov methods for linear systems of equations

Goal: Efficient solver for linear systems $Ax = b$ where the matrix A is *sparse*, i.e., most of its coefficients are equal to zero.

For these matrices, the matrix-vector multiplication $y \mapsto Ay$ can be performed very efficiently, so it makes sense to look for algorithms based on this operation.

Minimization: Assume that $A \in \mathbb{R}^{\mathcal{I} \times \mathcal{I}}$ is symmetric and positive definite, i.e.,

$$\langle x, Ay \rangle_2 = \langle Ax, y \rangle_2 \quad \text{for all } x, y \in \mathbb{R}^{\mathcal{I}}, \quad (4.1a)$$

$$\langle x, Ax \rangle_2 > 0 \quad \text{for all } x \in \mathbb{R}^{\mathcal{I}} \setminus \{0\}. \quad (4.1b)$$

We define the function

$$J: \mathbb{R}^{\mathcal{I}} \rightarrow \mathbb{R}, \quad x \mapsto \frac{1}{2} \langle x, Ax \rangle_2 - \langle b, x \rangle_2,$$

and use (4.1a) to obtain

$$\begin{aligned} J(x + \theta p) &= \frac{1}{2} \langle x + \theta p, A(x + \theta p) \rangle_2 - \langle b, x + \theta p \rangle_2 \\ &= \frac{1}{2} \langle x, Ax \rangle_2 + \frac{\theta}{2} \langle p, Ax \rangle_2 + \frac{\theta}{2} \langle x, Ap \rangle_2 + \frac{\theta^2}{2} \langle p, Ap \rangle_2 - \langle b, x \rangle_2 - \theta \langle b, p \rangle_2 \\ &= J(x) + \frac{\theta}{2} \langle p, Ax \rangle_2 + \frac{\theta}{2} \langle Ax, p \rangle_2 + \frac{\theta^2}{2} \langle p, Ap \rangle_2 - \theta \langle b, p \rangle_2 \\ &= J(x) + \theta \langle p, Ax - b \rangle_2 + \frac{\theta^2}{2} \langle p, Ap \rangle_2 \end{aligned} \quad (4.2)$$

for all $p \in \mathbb{R}^{\mathcal{I}}$ and $\theta \in \mathbb{R}$.

We assume $p \neq 0$. If we cannot reduce the value of J by “moving” from x along the direction p , i.e., if

$$J(x) \leq J(x + \theta p) \quad \text{for all } \theta \in \mathbb{R}, \quad (4.3a)$$

holds, we can use

$$\theta := - \frac{\langle p, Ax - b \rangle_2}{\langle p, Ap \rangle_2}$$

to find

$$J(x) \leq J(x + \theta p) = J(x) - \frac{\langle p, Ax - b \rangle_2^2}{2 \langle p, Ap \rangle_2} \leq J(x)$$

and use (4.1b) to conclude

$$\langle p, Ax - b \rangle_2 = 0. \quad (4.3b)$$

If (4.3b) holds, on the other hand, we can use (4.1b) to find that (4.3a) holds.

This means that the minimality condition (4.3a) and the orthogonality condition (4.3b) are equivalent. This is also true for $p = 0$, and therefore for all directions $p \in \mathbb{R}^{\mathcal{I}}$.

If J takes its *global* minimum in x , we can choose the direction $p = Ax - b$ and obtain $0 = \langle p, Ax - b \rangle_2 = \|Ax - b\|_2^2$, i.e., $Ax = b$. This means that looking for a minimum of J is equivalent with looking for the solution of the linear system.

Iteration: Given an approximation $x_m \in \mathbb{R}^{\mathcal{I}}$ of the solution, choose a direction $p_m \in \mathbb{R}^{\mathcal{I}}$ and a stepsize $\theta_m \in \mathbb{R}$ and define $x_{m+1} \leftarrow x_m + \theta_m p_m$.

Residual: For small values of θ , we have

$$J(x + \theta p) = J(x) + \theta \langle p, Ax - b \rangle_2 + \frac{\theta^2}{2} \langle p, Ap \rangle_2 \approx J(x) + \theta \langle p, Ax - b \rangle_2,$$

and the *Cauchy-Schwarz inequality* indicates that the *residual* $p = b - Ax$ is the best choice for the direction, at least for small stepsizes θ .

Optimal stepsize: Given an approximation $x_m \in \mathbb{R}^{\mathcal{I}}$ and a direction $p_m \in \mathbb{R}^{\mathcal{I}}$, we look for the stepsize $\theta_m \in \mathbb{R}$ minimizing

$$\theta \mapsto J(x_{m+1}) = J(x_m + \theta p_m).$$

Due to (4.3), x_{m+1} is optimal with respect to the direction p_m if

$$0 = \langle p_m, Ax_{m+1} - b \rangle_2 = \langle p_m, A(x_m + \theta_m p_m) - b \rangle_2 = \langle p_m, Ax_m - b \rangle_2 + \theta_m \langle p_m, Ap_m \rangle_2$$

holds, i.e.,

$$\theta_m = -\frac{\langle p_m, Ax_m - b \rangle_2}{\langle p_m, Ap_m \rangle_2} = \frac{\langle p_m, b - Ax_m \rangle_2}{\langle p_m, Ap_m \rangle_2}.$$

Gradient method: We denote the residuals by $r_m := b - Ax_m$. If we choose $p_m = r_m$ and the optimal stepsizes, we arrive at the *gradient method*.

In order to avoid recomputing the residual in each step, we use the auxiliary vector $a_m := Ap_m$ and

$$r_{m+1} = b - Ax_{m+1} = b - A(x_m + \theta_m p_m) = r_m - \theta_m Ap_m = r_m - \theta_m a_m$$

to arrive at

$$\begin{aligned} p_0 &\leftarrow b - Ax_0, \\ a_m &\leftarrow Ap_m, \\ \theta_m &\leftarrow \frac{\langle p_m, p_m \rangle_2}{\langle p_m, a_m \rangle_2}, \\ x_{m+1} &\leftarrow x_m + \theta_m p_m, \\ p_{m+1} &\leftarrow p_m - \theta_m a_m \end{aligned} \quad \text{for all } m \in \mathbb{N}_0.$$

Observations:

- We only have to be able to multiply vectors by the matrix A , we do not need the matrix itself. This means that we can represent A , e.g., by a callback function in C or a virtual member function in C++ or Java.
- We only need one matrix-vector multiplication per step.
- We only need to store the vectors b , x_m , p_m , and a_m in each step.

Experiment: Convergence for a finite difference discretization of Poisson's equation with $n = 32$ mesh points per direction and stepwidth $h = \frac{1}{33}$:

m	$\ b - Ax_m\ $	$f(x_m)$
0	2.33 ₊₃	0.00 ₊₀
1	1.21 ₊₃	-1.27 ₊₃
2	8.92 ₊₂	-1.61 ₊₃
3	7.31 ₊₂	-1.80 ₊₃
4	6.34 ₊₂	-1.92 ₊₃
10	4.20 ₊₂	-2.25 ₊₃
100	1.89 ₊₂	-3.02 ₊₃
1000	3.04 ₊₀	-3.45 ₊₃
2000	3.14 ₋₂	-3.45 ₊₃

We can see that the gradient method converges *very* slowly.

Preservation of optimality: After the first step, x_1 is optimal with respect to the direction p_0 , i.e., it cannot be improved by adding a multiple of p_0 . Unfortunately this property is usually lost: x_2 is optimal with respect to p_1 , but not longer with respect to the direction p_0 .

We can preserve optimality by choosing the directions a little differently. Assume that x_m is optimal with respect to a direction p_ℓ , $\ell < m$. If we want x_{m+1} to still be optimal, we have to ensure

$$0 = \langle p_\ell, Ax_{m+1} - b \rangle_2 = \langle p_\ell, A(x_m + \theta_m p_m) - b \rangle_2 = \langle p_\ell, Ax_m - b \rangle_2 + \theta_m \langle p_\ell, Ap_m \rangle_2,$$

i.e., we can either choose $\theta_m = 0$, which is counterproductive, or ensure $\langle p_\ell, Ap_m \rangle_2 = 0$, i.e., that the vectors p_ℓ and p_m are *conjugates*.

Gram-Schmidt orthogonalization: Starting with r_m , we construct a direction p_m that is conjugate with respect to all p_ℓ with $\ell < m$:

$$p_m \leftarrow r_m - \sum_{\ell=0}^{m-1} \frac{\langle p_\ell, Ar_m \rangle_2}{\langle p_\ell, Ap_\ell \rangle_2} p_\ell. \quad (4.4)$$

Assuming that the directions p_0, \dots, p_{m-1} are already conjugate with respect to each other, we have

$$\begin{aligned} \langle p_k, Ap_m \rangle_2 &= \langle p_k, Ar_m \rangle_2 - \sum_{\ell=0}^{m-1} \frac{\langle p_\ell, Ar_m \rangle_2}{\langle p_\ell, Ap_\ell \rangle_2} \langle p_k, Ap_\ell \rangle_2 \\ &= \langle p_k, Ar_m \rangle_2 - \frac{\langle p_k, Ar_m \rangle_2}{\langle p_k, Ap_k \rangle_2} \langle p_k, Ap_k \rangle_2 = 0 \end{aligned}$$

for all $k \in [0 : m - 1]$, i.e., the new direction p_m will also be conjugate with respect to all previous directions.

Krylov spaces In the gradient method, new directions p_{m+1} are obtained by linear combinations between p_m and $a_m = Ap_m$, and a simple induction yields that p_m must be a linear combination of $p_0, Ap_0, \dots, A^m p_0$.

This leads us to define the *Krylov spaces*

$$\mathcal{K}_m := \text{span}\{p_0, Ap_0, \dots, A^m p_0\} \quad \text{for all } m \in \mathbb{N}_0.$$

Even if the directions are constructed by Gram-Schmidt orthogonalization (4.4), we have

$$\text{span}\{p_0, \dots, p_m\} = \text{span}\{r_0, \dots, r_m\} = \mathcal{K}_m$$

as long as $p_m \neq 0$.

The first inclusion $\text{span}\{p_0, \dots, p_m\} \subseteq \text{span}\{r_0, \dots, r_m\}$ follows directly from (4.4) by induction. The second inclusion $\text{span}\{r_0, \dots, r_m\} \subseteq \mathcal{K}_m$ follows from the first with $r_{m+1} = r_m - \theta_m Ap_m$, again by induction. As long as all directions are non-zero, they have to be linear independent, since they are pairwise conjugate, so all spaces are $(m+1)$ -dimensional and therefore identical.

If we assume that x_m is optimal with respect to all previous directions p_0, \dots, p_{m-1} , it is also optimal with respect to all vectors in \mathcal{K}_{m-1} .

For $\ell \in [0 : m-2]$, we have $p_\ell \in \mathcal{K}_{m-2}$ and therefore $Ap_\ell \in \mathcal{K}_{m-1}$, and the optimality together with the symmetry (4.1a) of A yields

$$\langle p_\ell, Ar_m \rangle_2 = \langle Ap_\ell, r_m \rangle_2 = 0.$$

This equation allows us to eliminate almost all terms in (4.4) and obtain

$$p_m \leftarrow r_m - \frac{\langle p_{m-1}, Ar_m \rangle_2}{\langle p_{m-1}, Ap_{m-1} \rangle_2} p_{m-1} = r_m - \frac{\langle Ap_{m-1}, r_m \rangle_2}{\langle p_{m-1}, Ap_{m-1} \rangle_2} p_{m-1}.$$

This makes the Gram-Schmidt orthogonalization far more efficient.

Conjugate gradient method: Adapting the gradient method for conjugate directions leads to

$$\begin{aligned} r_0 &\leftarrow b - Ax_0, & p_0 &\leftarrow r_0, \\ a_m &\leftarrow Ap_m, \\ \theta_m &\leftarrow \frac{\langle p_m, r_m \rangle_2}{\langle p_m, a_m \rangle_2}, \\ x_{m+1} &\leftarrow x_m + \theta_m p_m, \\ r_{m+1} &\leftarrow r_m - \theta_m a_m, \\ \mu_m &\leftarrow \frac{\langle a_m, r_{m+1} \rangle_2}{\langle p_m, a_m \rangle_2}, \\ p_{m+1} &\leftarrow r_{m+1} - \mu_m p_m \end{aligned} \quad \text{for all } m \in \mathbb{N}_0.$$

Observations:

- We still need only one matrix-vector multiplication per step.
- We need an additional vector to store r_m , one additional inner product, and one additional vector addition.
- In practice, the conjugate gradient method (“cg method”) often performs significantly better than the simple gradient method.

Experiment: Convergence for a finite difference discretization of Poisson’s equation with $n = 32$ mesh points per direction and stepwidth $h = \frac{1}{33}$:

m	$\ b - Ax_m\ $	$f(x_m)$
0	2.33 ₊₃	0.00 ₊₀
1	1.21 ₊₃	-1.27 ₊₃
2	8.64 ₊₂	-1.74 ₊₃
3	7.08 ₊₂	-2.02 ₊₃
10	5.44 ₊₂	-2.85 ₊₃
20	4.46 ₊₂	-3.35 ₊₃
30	4.95 ₊₀	-3.45 ₊₃
40	2.05 ₋₁	-3.45 ₊₃
50	4.95 ₋₄	-3.45 ₊₃
100	6.11 ₋₁₅	-3.45 ₊₃

We can see that the conjugate gradient method converges significantly faster than the “simple” gradient method.

GMRES: If A is just invertible, but neither symmetric nor positive definite, we can still look for approximations $x_m \in x_0 + \mathcal{K}_m$ that minimize the *residual norm* $\|b - Ax_m\|_2 = \|r_m\|_2$, i.e.,

$$\|b - Ax_m\|_2 \leq \|b - Ay\|_2 \quad \text{for all } y = x_0 + z, \quad z \in \mathcal{K}_m.$$

If we write $x_m = x_0 + z_m$ with $z \in \mathcal{K}_m$, we arrive at

$$\|r_0 - Az_m\|_2 \leq \|r_0 - Az\|_2 \quad \text{for all } z \in \mathcal{K}_m.$$

This is a least-square optimization problem that can be handled by constructing an orthonormal basis of the Krylov space called the *Arnoldi basis*.

The resulting algorithm is called GMRES, *general minimized residual*.

5 Algorithms for large systems

5.1 Multigrid methods

Relaxation methods are based on the idea of choosing a subset of equations and a subset of unknown variables and solving the corresponding restricted system.

The simplest version consists of choosing one index $i \in \mathcal{I}$ and determining the i -th component of the solution such that the i -th equation is satisfied, i.e., we look for \tilde{x}_i such that

$$a_{ii}\tilde{x}_i + \sum_{\substack{j \in \mathcal{I} \\ j \neq i}} a_{ij}x_j = b_i$$

holds. Assuming $a_{ii} \neq 0$, this is equivalent with

$$\tilde{x}_i = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j \in \mathcal{I} \\ j \neq i}} a_{ij}x_j \right).$$

Jacobi iteration If we apply the relaxation step to all indices $i \in \mathcal{I}$, we arrive at the *Jacobi iteration* for linear systems of equations.

It is frequently applied with a *dampening parameter* $\theta \in [0, 1]$ that can be used to ensure convergence even if the diagonal entries are fairly small:

$$\tilde{x}_i = (1 - \theta)x_i + \frac{\theta}{a_{ii}} \left(b_i - \sum_{\substack{j \in \mathcal{I} \\ j \neq i}} a_{ij}x_j \right).$$

If the matrix A is symmetric and positive definite (cf. (4.1)), the Jacobi method converges as long as θ is small enough.

Gauß-Seidel method If we apply the relaxation steps one after the other, we can replace the original values x_i on the right-hand side of the equation by the already updated values \tilde{x}_i of previous steps.

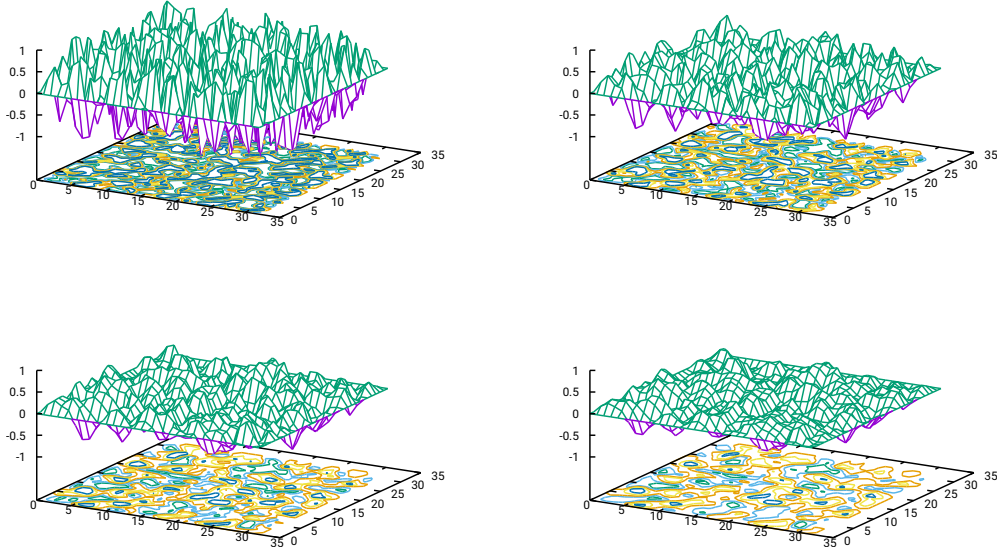
We write $j < i$ if the index j is treated before i and $j > i$ if it is treated later. The *Gauß-Seidel iteration* is given by

$$\tilde{x}_i = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j \in \mathcal{I} \\ j < i}} a_{ij}\tilde{x}_j - \sum_{\substack{j \in \mathcal{I} \\ j > i}} a_{ij}x_j \right).$$

If the matrix A is symmetric and positive definite (cf. again (4.1)), the Gauß-Seidel method converges.

Smoothing Unfortunately, both the Jacobi and the Gauß-Seidel iteration converge quite slowly when applied to linear systems resulting from partial differential equations.

If we look at the *errors* left by these iterations after a few steps, we see that they are getting more and more “smooth”.



In model problems, it is possible to prove this behaviour: high-frequency errors are reduced rapidly, low-frequency errors far more slowly.

If the error is smooth, there is a good chance that it can be approximated on a coarser grid, and computations on a coarser grid involve a lower number of grid points and therefore a lower amount of work. If we have a good approximation of the error $x - x_m$, adding this approximation to x_m will significantly improve the accuracy.

Grid transfer In order to make use of a coarser grid Ω_H to approximate the error on a finer grid Ω_h , we have to be able to transfer data between both grids. We introduce the *prolongation* $p \in \mathbb{R}^{\Omega_h \times \Omega_H}$ mapping coarse grid functions to fine grid functions and the *restriction* $r \in \mathbb{R}^{\Omega_H \times \Omega_h}$ for the opposite direction.

In the simple case $H = 2h$, we can define a suitable prolongation by simple linear interpolation: fine-grid points between coarse-grid points are assigned interpolated values.

$$p_{xy} = \begin{cases} 1 & \text{if } x = y, \\ 1/2 & \text{if } x - y = (h, 0) \text{ or } y - x = (h, 0), \\ 1/2 & \text{if } x - y = (0, h) \text{ or } y - x = (0, h), \\ 1/2 & \text{if } x - y = (h, h) \text{ or } y - x = (h, h), \\ 0 & \text{otherwise} \end{cases} \quad \text{for all } x \in \Omega_h, y \in \Omega_H.$$

For the restriction, a good choice is the scaled transposed matrix of the prolongation, i.e., $r = \frac{1}{4}p^T$.

Coarse-grid correction Assuming that the error $x - x_m$ is sufficiently smooth, we can approximate it by a coarse-grid function c defined via

$$\begin{aligned} x - x_m &\approx pc, \\ -\Delta_h(x - x_m) &\approx -\Delta_h(pc), \\ f_h + \Delta_h x_m &\approx -\Delta_h pc, \\ r(f_h + \Delta_h x_m) &= -r\Delta_h pc = -\Delta_H c, \end{aligned}$$

where we use the *Galerkin product identity* $\Delta_H = r\Delta_h p$ in the last step.

Two-grid method Combining a smoothing iteration and a coarse grid correction leads to the two-grid iteration:

```

procedure twogrid( $x_h, f_h, d_h$ );
for i=1 to  $\nu$  do
  Smoother( $x_h, f_h, d_h$ );
 $d_h \leftarrow f_h + \Delta_h x_h$ ;
 $f_H \leftarrow r d_h$ ;
Solve  $-\Delta_H x_H = f_H$ ;
 $x_h \leftarrow x_h + p x_H$ 

```

Multigrid method In practice, the coarse grid Ω_H may still contain too many grid points to be solved efficiently.

In this case, we can replace the solution step on the coarse grid recursively by smoothing and a coarse-grid correction on an even coarser grid. This leads to the *multigrid method*: we assume that a grid hierarchy $\Omega_0 \subseteq \Omega_1 \subseteq \dots \subseteq \Omega_L = \Omega_h$ is given with finite difference operators $\Delta_0, \dots, \Delta_L$ and prolongation and restriction matrices $p_\ell \in \mathbb{R}^{\Omega_\ell \times \Omega_{\ell-1}}$ and $r_\ell \in \mathbb{R}^{\Omega_{\ell-1} \times \Omega_\ell}$. A simple multigrid method takes the following form:

```

procedure multigrid( $\ell, x_\ell, f_\ell, d_\ell$ );
if  $\ell = 0$  then
  Solve  $-\Delta_\ell x_\ell = f_\ell$ 
else begin
  for  $i = 1$  to  $\nu$  do
    Smoother( $x_\ell, f_\ell, d_\ell$ );
 $d_\ell \leftarrow f_\ell + \Delta_\ell x_\ell$ ;
 $f_{\ell-1} \leftarrow r_\ell d_\ell$ ;  $x_{\ell-1} \leftarrow 0$ ;
  for  $i = 1$  to  $\gamma$  do
    multigrid( $\ell - 1, x_{\ell-1}, f_{\ell-1}, d_{\ell-1}$ );
 $x_\ell \leftarrow x_\ell + p_\ell x_{\ell-1}$ ;
  for  $i = 1$  to  $\nu$  do
    Smoother( $x_\ell, f_\ell, d_\ell$ )
end

```


Here we use $\nu \in \mathbb{N}$ steps of the smoothing iteration before and after the coarse-grid correction and $\gamma \in \mathbb{N}$ steps of the multigrid methods to approximate the coarse-grid correction.

Multigrid convergence An important property of multigrid methods is that their convergence rate is fairly independent of the grid. For our model problem, it is possible to prove that the rate can be bounded by a constant that does not depend on the grid parameter h .

m	$h = 1/16$		$h = 1/32$		$h = 1/64$		$h = 1/8192$	
	error	ratio	error	ratio	error	ratio	error	ratio
0	1.50 ₊₁		3.10 ₊₁		6.30 ₊₁		8.19 ₊₃	
1	2.52 ₊₀	5.94	5.39 ₊₀	5.75	1.11 ₊₁	5.68	1.46 ₊₃	5.62
2	4.49 ₋₁	5.62	9.68 ₋₁	5.57	1.99 ₊₀	5.56	2.62 ₊₂	5.57
3	8.09 ₋₂	5.55	1.76 ₋₁	5.51	3.62 ₋₁	5.51	4.64 ₊₁	5.52
4	1.47 ₋₂	5.52	3.21 ₋₂	5.47	6.62 ₋₂	5.47	8.65 ₊₀	5.48
5	2.67 ₋₃	5.49	5.89 ₋₃	5.45	1.22 ₋₂	5.44	1.59 ₊₀	5.45
10	5.46 ₋₇	5.46	1.27 ₋₆	5.40	2.64 ₋₆	5.39	3.44 ₋₄	5.38
20	2.35 ₋₁₄	5.45	6.17 ₋₁₄	5.38	1.32 ₋₁₃	5.37	1.73 ₋₁₁	5.36

Multigrid runtimes In theory, the number of operations of the multigrid method is proportional to the number of gridpoints n . In practice, we observe a slightly faster growth of the computing time as the grid is refined, possibly due to the data transfer between the four eight-core processors used in the parallel implementation used for this experiment.

h	dofs	conv	time	ratio
2^{-12}	16 769 025	5.36	1.9	
2^{-13}	67 092 481	5.36	5.4	2.84
2^{-14}	268 402 689	5.36	15.6	2.89
2^{-15}	1 073 676 289	5.36	77.9	4.99
2^{-16}	4 294 836 225	5.36	361.5	4.64
2^{-17}	17 179 607 041	5.36	1 727.1	4.78

5.2 Non-local force fields

Example: Gravitation in many-body systems Given n planets at positions $y_1, \dots, y_n \in \mathbb{R}^3$ with masses $m_1, \dots, m_n \in \mathbb{R}$, we want to evaluate the gravitational potential

$$f(x) = \gamma \sum_{j=1}^n m_j \frac{1}{\|y_j - x\|_2},$$

possibly for many different positions $x \in \mathbb{R}^3$.

Challenge: Evaluating $f(x)$ requires at least n operations, since every mass influences the entire gravitational field.

Degenerate approximation: We introduce the *kernel function*

$$g(x, y) := \frac{\gamma}{\|y - x\|_2},$$

and write

$$f(x) = \sum_{j=1}^n g(x, y_j) m_j.$$

If we can find a *degenerate approximation*

$$g(x, y) \approx \sum_{\nu=0}^k a_\nu(x) b_\nu(y),$$

we obtain

$$f(x) = \sum_{j=1}^n g(x, y_j) m_j \approx \sum_{j=1}^n \sum_{\nu=0}^k a_\nu(x) b_\nu(y_j) m_j = \sum_{\nu=0}^k a_\nu(x) \underbrace{\sum_{j=1}^n b_\nu(y_j) m_j}_{=: z_\nu},$$

so once we have prepared

$$z_\nu = \sum_{j=1}^n b_\nu(y_j) m_j \quad \text{for all } \nu \in [0 : k],$$

we can approximate $f(x)$ by

$$f(x) \approx \sum_{\nu=0}^k a_\nu(x) z_\nu.$$

If $k \ll n$, this latter step can be *far* more efficient than the original computation.

Three-dimensional interpolation Assume that the positions of all planets are contained in an axis-parallel *bounding box*

$$B = [a_1, b_1] \times [a_2, b_2] \times [a_3, b_3].$$

If we have interpolations points

$$\begin{aligned} \xi_{1,0}, \dots, \xi_{1,m} &\in [a_1, b_1], \\ \xi_{2,0}, \dots, \xi_{2,m} &\in [a_2, b_2], \\ \xi_{3,0}, \dots, \xi_{3,m} &\in [a_3, b_3] \end{aligned}$$

for the three intervals at our disposal, we can approximate the kernel function by interpolation to obtain

$$g(x, y) \approx \sum_{\nu_1=0}^m \sum_{\nu_2=0}^m \sum_{\nu_3=0}^m g \left(x, \begin{pmatrix} \xi_{1,\nu_1} \\ \xi_{2,\nu_2} \\ \xi_{3,\nu_3} \end{pmatrix} \right) \ell_{1,\nu_1}(y_1) \ell_{2,\nu_2}(y_2) \ell_{3,\nu_3}(y_3).$$

Introducing the abbreviations

$$\xi_{B,\nu} := \begin{pmatrix} \xi_{1,\nu_1} \\ \xi_{2,\nu_2} \\ \xi_{3,\nu_3} \end{pmatrix}, \quad \ell_{B,\nu}(y) := \ell_{1,\nu_1}(y_1) \ell_{2,\nu_2}(y_2) \ell_{3,\nu_3}(y_3) \quad \text{for all } \nu \in M := [0 : m]^3,$$

we arrive at

$$g(x, y) \approx \sum_{\nu \in M} g(x, \xi_{B,\nu}) \ell_{B,\nu}(y),$$

and this is a degenerate approximation for $k = \#M = (m + 1)^3$.

Admissibility The kernel function g has a singularity at $x = y$ that cannot be handled by a degenerate approximation.

We can only apply our technique to boxes satisfying the *admissibility condition*

$$\text{diam}(B) \leq \text{dist}(x, B).$$

This condition is sufficient to prove *exponential convergence* of the interpolation, i.e., the error behaves like $e^{-\alpha m}$ with a parameter $\alpha > 0$.

Recursion If the box B does not satisfy the admissibility condition, we split it into smaller boxes and check these boxes again. We can repeat this procedure recursively until the boxes are small enough to contain only a small number of planets that can be handled directly.

Cluster tree We construct a tree \mathcal{T} consisting of axis-parallel boxes such that

- the root box contains all points y_1, \dots, y_n ,
- the sons s' of a box s describe a partition of their father, and
- the leaf boxes contain only a small number of points.

We call this a *cluster tree*, its nodes are called *clusters*.

Setup We compute the cluster tree \mathcal{T} , the index sets

$$\hat{s} := \{j \in [1 : n] : y_j \in s\} \quad \text{for all } s \in \mathcal{T},$$

and the coefficients

$$z_{s,\nu} := \sum_{j \in \hat{s}} \ell_{s,\nu}(y_j) m_j \quad \text{for all } s \in \mathcal{T}, \nu \in M$$

of our degenerate approximation.

The sons of $s = [a_1, b_1] \times [a_2, b_2] \times [a_3, b_3]$ can be constructed using the midpoints $c_1 := (b_1 + a_1)/2$, $c_2 := (b_2 + a_2)/2$, $c_3 := (b_3 + a_3)/2$ as

$$\begin{aligned} s_1 &\leftarrow [a_1, c_1] \times [a_2, c_2] \times [a_3, c_3], & \hat{s}_1 &\leftarrow \{j \in \hat{s} : y_{j,1} < c_1, y_{j,2} < c_2, y_{j,3} < c_3\}, \\ s_2 &\leftarrow [a_1, c_1] \times [a_2, c_2] \times [c_3, b_3], & \hat{s}_2 &\leftarrow \{j \in \hat{s} : y_{j,1} < c_1, y_{j,2} < c_2, y_{j,3} \geq c_3\}, \\ s_3 &\leftarrow [a_1, c_1] \times [c_2, b_2] \times [a_3, c_3], & \hat{s}_3 &\leftarrow \{j \in \hat{s} : y_{j,1} < c_1, y_{j,2} \geq c_2, y_{j,3} < c_3\}, \\ s_4 &\leftarrow [a_1, c_1] \times [c_2, b_2] \times [c_3, b_3], & \hat{s}_4 &\leftarrow \{j \in \hat{s} : y_{j,1} < c_1, y_{j,2} \geq c_2, y_{j,3} \geq c_3\}, \\ s_5 &\leftarrow [c_1, b_1] \times [a_2, c_2] \times [a_3, c_3], & \hat{s}_5 &\leftarrow \{j \in \hat{s} : y_{j,1} \geq c_1, y_{j,2} < c_2, y_{j,3} < c_3\}, \\ s_6 &\leftarrow [c_1, b_1] \times [a_2, c_2] \times [c_3, b_3], & \hat{s}_6 &\leftarrow \{j \in \hat{s} : y_{j,1} \geq c_1, y_{j,2} < c_2, y_{j,3} \geq c_3\}, \\ s_7 &\leftarrow [c_1, b_1] \times [c_2, b_2] \times [a_3, c_3], & \hat{s}_7 &\leftarrow \{j \in \hat{s} : y_{j,1} \geq c_1, y_{j,2} \geq c_2, y_{j,3} < c_3\}, \\ s_8 &\leftarrow [c_1, b_1] \times [c_2, b_2] \times [c_3, b_3], & \hat{s}_8 &\leftarrow \{j \in \hat{s} : y_{j,1} \geq c_1, y_{j,2} \geq c_2, y_{j,3} \geq c_3\}. \end{aligned}$$

Evaluation We start with the root of the cluster tree and $v = 0$.

- If s satisfies the admissibility condition, update the approximation

$$v \leftarrow v + \sum_{\nu \in M} g(x, \xi_{s,\nu}) z_{s,\nu} \approx v + \sum_{j \in \hat{s}} g(x, y_j) m_j.$$

- Otherwise, if s has sons, consider the sons' contributions.
- Otherwise, s has to be a leaf, i.e., \hat{s} contains only a few indices, so we can afford to update the approximation directly via

$$v \leftarrow v + \sum_{j \in \hat{s}} g(x, y_j) m_j.$$

It can be proven that only a small number of clusters contribute on every level of the tree, and the total computational work is $\sim k \log(n)$.

6 Shared-memory parallelization

6.1 Multithreading and OpenMP

Multicore processors Modern processors contain multiple *cores*, each equipped with its own arithmetic and logical units, flow control logic, and caches. Each of these cores is able to execute programs independently of the others, only the interface to off-chip components like memory, graphics cards, and other peripherals is shared.

If we can share the work of a computation between multiple cores, we can significantly reduce the necessary runtime.

OpenMP A very popular standard for writing programs that can run on multiple cores is OpenMP, the *open multi-processing standard*. OpenMP is an extension of the C/C++ programming language.

- Programs are executed in *threads*, each thread has its own instruction pointer, registers, and stack (containing local variables and data for subroutine calls).
- The program starts with the *master thread* that executes the `main` function of a C program.
- Threads can new threads, organized in *teams*. The programmer can state how many threads a team should contain, but the decision is ultimately up to the implementation.
- All threads share the same address space, e.g., global variables and objects allocated on the heap.
- Most of OpenMP's functionality is accessed by `#pragma` directives, some by library functions defined in the `omp.h` header file.

Parallel sections If a thread encounters a parallel section, a team of threads is created:

```
int
main()
{
    int i=0, j=0;

    #pragma omp parallel firstprivate(j)
    {
        int k=0;
```

```

    i++; j++; k++;
}

return 0;
}

```

In this example, the variable `i` is *shared*, i.e., all threads can modify it. The variable `j` is *private* to each thread, since the `firstprivate` clause creates copies and initializes them with the current value of the original. The variable `k` is also private, since it is created in each thread.

Since all threads share the same address space, a thread may still be able to access private variables of other threads by misusing pointers.

Work-sharing via functions In order to make each thread work on a part of a problem, we can use the functions `omp_get_num_threads` to get the number of threads in the current team, and `omp_get_thread_num` to get the number of the current thread within the team, starting with zero:

```

#pragma omp parallel
{
    int start = n * omp_get_thread_num() / omp_get_num_threads();
    int end = n * (omp_get_thread_num()+1) / omp_get_num_threads();
    int i;

    for(i=start; i<end; i++)
        work(i);
}

```

While this approach offers us great freedom to assign work to threads, it works only if a compiler supports OpenMP and if the header file `omp.h` is included.

Work-sharing via directives For simple loops, we can use `#pragma` directives to share iterations among the threads of a team:

```

#pragma omp parallel
{
    int i;

    #pragma omp for
    for(i=0; i<n; i++)
        work(i);
}

```

If we want to parallelize only one loop, both pragmas can be combined:

```
#pragma omp parallel for
for(i=0; i<n; i++)
    work(i);
```

Creating a team of threads is a relatively lengthy procedure, so this latter approach should only be used if only one loop has to be parallelized or if this loop takes very long to complete.

Example: matrix-vector multiplication

```
#pragma omp parallel for collapse(2)
for(j=1; j<=n; j++)
    for(i=1; i<=n; i++)
        y[i+j*ld] += alpha * (4.0 * x[i+j*ld]
                               - x[(i-1)+j*ld] - x[(i+1)+j*ld]
                               - x[i+(j-1)*ld] - x[i+(j+1)*ld]);
```

The `collapse` clause causes the two nested loops to be treated as one larger loop.

Example: trapezoidal rule

```
double sum = 0.0;
double h = (b-a) / n;

#pragma omp parallel for reduction(+:sum)
for(i=1; i<n; i++)
    sum += f(a+h*i);

sum = h * (0.5 * f(a) + sum + 0.5 * f(b));
```

The `reduction` clause makes each thread create its own private copy of the `sum` variable, initialized with zero. After the loop has been completed, all private copies are accumulated in the original variable.

Interference If multiple threads work with the same shared variable, the results may be non-deterministic:

```
int i = 0;

#pragma omp parallel
i++;
```

If we execute this fragment with 16 threads, the shared variable `i` can take any value between 1 and 16 after the parallel section is complete: `i++` is split into reading, adding, and writing, and if multiple threads read the value of `i` before others have updated it, these updates will be lost.

One solution is to make the accesses to `i` *atomic* by the `#pragma omp atomic` directive. This works only for certain operations.

Critical sections Another solution is to define a *critical section*. Only one thread can execute the critical section at any time, so if a variable is only accessed in a critical section, it is safe from interference.

```
#pragma critical safe_print
{
    printf("Hello\n");
    printf("World\n");
}
```

Other threads trying to execute the critical section wait until the previous thread has left the section.

Critical sections can be given a name (`safe_print` in the example), and all critical sections with the same name are treated as if they were one section, i.e., once a thread enters one of them, no further thread can enter any of the sections with the same name.

Locks Critical sections have to be defined at compile time, so they cannot grow and shrink with dynamic data structures.

Locks, on the other hand, can be created at run time.

```
int *counter;
omp_lock_t *lock;

counter = (int *) malloc(sizeof(int) * n);
lock = (omp_lock_t *) malloc(sizeof(omp_lock_t) * n);

for(i=0; i<n; i++) {
    counter[i] = 0;
    omp_init_lock(lock+i);
}

#pragma omp parallel for
for(j=0; j<m; j++) {
    i = data[j];
    omp_set_lock(lock+i);
    counter[i]++;
    omp_unset_lock(lock+i);
}
```

This fragment counts how often numbers between 0 and $n-1$ appear in the array `data`. In order to avoid multiple threads accessing the same entry of the `counter` array, each entry is protected by a lock: if the lock is *unset*, `omp_set_lock` will set it and proceed. If the lock is *set*, `omp_set_lock` will wait until it becomes *unset*, then set it and proceed.

Forgetting a call to `omp_unset_lock` can lead to a *deadlock* where no thread can proceed and the program never finishes.

False sharing Even if two threads access different variables, they may feel the influence of each other.

```
int a[2];

#pragma omp parallel sections private(i)
{
    #pragma omp section
    for(i=0; i<n; i++)
        a[0]++;
    #pragma omp section
    for(i=0; i<n; i++)
        a[1]++;
}
```

If the variables `a[0]` and `a[1]` share the same *cache line*, the cores executing both loops have to exchange the cache lines after every update and cannot properly take advantage of their local caches.

NUMA architectures Multi-processor systems frequently use *non-uniform memory access*, i.e., each processor has its own memory and accessing the memory belonging to another processor requires time-consuming communication.

In order to obtain good performance, we should try to ensure that processors work in their own memory most of the time.

Memory is organized in *pages*, e.g., of 4096 bytes. If we allocate memory by `malloc`, only *logical* memory is allocated, no *physical* memory is assigned.

Physical memory is assigned when an address within the page is first accessed. In a NUMA system with the common *first touch* allocation strategy, the page is placed in the memory of the first processor accessing it.

This means, e.g., that it is not a good idea to have one processor initialize a large array, because it would lead to operation system to try to store the entire array in this processor's memory.

6.2 Tasks

Simple work-sharing constructs are too limited to handle computations with non-trivial data dependencies.

Task: A thread can define a *task*.

- A task is a segment of code that operates with its own instruction pointer, local variables, and stack.
- If a thread in the current team is idle, the runtime system can assign it a task to execute.
- Tasks can depend on the results of other tasks, and they can produce results that other tasks depend on.

Creating a task takes time, but not as much time as creating a team of threads.

Example: Parallel tree traversal

```
static void
work(tree *t)
{
    int i;

    if(t->sons > 0) {
        for(i=0; i<t->sons; i++)
            #pragma omp task firstprivate(t,i)
            work(t->son[i]);

        #pragma omp taskwait
    }

    /* ... do something ... */
}

void
work_tree(tree *t)
{
    #pragma omp parallel
    #pragma omp single
    work(t);
}
```

The `task` directive creates a task for each node in the tree, and the `firstprivate` clause ensures that the tasks have their own copies of *t* and *i*.

The `taskwait` directive waits for all task created by the current task to finish, so the father node could work with data provided by its sons.

The `single` directive ensures that only one thread executes for the root of the tree.

Data dependencies Tasks can depend on variables.

```
int a, b;

#pragma omp task depend(out: a)
a = 5;

#pragma omp task depend(out: b)
b = 7;

#pragma omp task depend(inout: a), depend(in: b)
a += b;
```

The first task has `a` as its output, the second `b`, while the third requires `a` and `b` as input and has `a` as its output.

The first two tasks can be executed in parallel, but the third task has to wait for both to complete.

Example: LR factorization Our goal is to parallelize the LR factorization. In order to ensure that tasks have sufficient work, we use a *block* algorithm: the matrices A , L , and R are split into $m \times m$ submatrices

$$A = \begin{pmatrix} A_{11} & \cdots & A_{1m} \\ \vdots & \ddots & \vdots \\ A_{m1} & \cdots & A_{mm} \end{pmatrix}, \quad L = \begin{pmatrix} L_{11} & & \\ \vdots & \ddots & \\ L_{m1} & \cdots & L_{mm} \end{pmatrix}, \quad R = \begin{pmatrix} R_{11} & \cdots & R_{1m} \\ & \ddots & \vdots \\ & & R_{mm} \end{pmatrix},$$

and we obtain the equations

$$A_{11} = L_{11}R_{11}, \quad A_{1j} = L_{11}R_{1j}, \quad A_{i1} = L_{i1}R_{11} \quad \text{for all } i, j \in [2 : m].$$

Due to the first equation, we can obtain L_{11} and R_{11} by a standard LR factorization.

Once L_{11} is available, we can compute R_{1j} from A_{1j} by forward substitution. Solving $A_{i1} = L_{i1}R_{11}$ is equivalent to $A_{i1}^T = R_{11}^T L_{i1}^T$, so we can again use forward substitution to obtain L_{i1} .

```
void
rtsolve(int n, int m, const real *R, int ldR, real *B, int ldB)
{
    int i;

    for(i=0; i+1<n; i++) {
```

```

    scal(m, 1.0/R[i+i*ldR], B+i*ldB, 1);
    ger(m, n-i-1, -1.0,
        B+i*ldB, 1, R+i+(i+1)*ldR, ldR,
        B+(i+1)*ldB, ldB);
}
}

```

Once L_{i1} and R_{1j} are available for all $i, j \in [2 : m]$, we can use

$$\begin{pmatrix} L_{22} & & \\ \vdots & \ddots & \\ L_{m2} & \cdots & L_{mm} \end{pmatrix} \begin{pmatrix} R_{22} & \cdots & R_{2m} \\ & \ddots & \vdots \\ & & R_{mm} \end{pmatrix} = \begin{pmatrix} A_{22} - L_{21}R_{12} & \cdots & A_{2m} - L_{21}R_{1m} \\ \vdots & \ddots & \vdots \\ A_{m2} - L_{m1}R_{12} & \cdots & A_{mm} - L_{m1}R_{1m} \end{pmatrix}$$

to reduce the problem size: now we only have to compute the LR factorization with $(m-1) \times (m-1)$ submatrices.

```

for(k=0; k<m; k++) {
    ok = n * k / m;          /* Offset block k */
    dk = n * (k+1) / m - ok; /* Dimension block k */

    #pragma omp task firstprivate(dk, ok),\
        depend(inout: A[ok+ok*ldC])
    lrdecomp(dk, A+ok+ok*ldA, ldA);

    for(j=k+1; j<m; j++) {
        oj = n * j / m;          /* Offset block j */
        dj = n * (j+1) / m - oj; /* Dimension block j */

        #pragma omp task firstprivate(dj, oj, dk, ok),\
            depend(in: A[ok+ok*ldA]), depend(inout: A[ok+oj*ldA])
        lsolve(dk, dj, A+ok+ok*ldA, ldA, A+ok+oj*ldA, ldA);

        #pragma omp task firstprivate(dj, oj, dk, ok),\
            depend(in: A[ok+ok*ldA]), depend(inout: A[oj+ok*ldA])
        rtsolve(dk, dj, A+ok+ok*ldA, ldA, A+oj+ok*ldA, ldA);
    }

    for(j=k+1; j<m; j++) {
        oj = n * j / m;          /* Offset block j */
        dj = n * (j+1) / m - oj; /* Dimension block j */

        for(i=k+1; i<m; i++) {
            oi = n * i / m;          /* Offset block i */
            di = n * (i+1) / m - oi; /* Dimension block i */

```

```
#pragma omp task firstprivate(di, oi, dj, oj, dk, ok), \  
    depend(in: A[oi+ok*ldA], A[ok+oj*ldA]), \  
    depend(inout: A[oi+oj*ldA])  
gemm(false, false, di, dj, dk,  
    -1.0, A+oi+ok*ldA, ldA, A+ok+oj*ldA, ldA,  
    1.0, A+oi+oj*ldA, ldA);  
    }  
    }  
}  
#pragma omp taskwait
```

7 Vectorization

7.1 Vectorization

Instead of adding more full processor cores to a chip, it is also possible to just add more arithmetic-logic units.

SIMD: Single instruction, multiple data.

The processor's instruction set is extended to allow one instruction to act on multiple sets of input data simultaneously.

The registers used by the processor are extended to *vector registers* that can hold multiple values simultaneously, and arithmetic and logical instructions are extended to work with these vector registers.

A simple example is the *vector add* operation

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \leftarrow \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ x_3 + y_3 \\ x_4 + y_4 \end{pmatrix}$$

that performs four additions and — on most processors — takes as long as just one addition.

Example: axpy operation The function `axpy` of BLAS level 1 can profit from SIMD units, since it has to perform one multiplication and one addition for each element of the array and these operations are all independent.

Using Intel[®] AVX intrinsics, a possible realization (for the special case `incx=1` and `incy=1`) would be

```
__m256 v_alpha = _mm256_set1_ps(alpha);
int i;

for(i=0; i+7<n; i+=8)
    _mm256_storeu_ps(y+i, _mm256_add_ps(
        _mm256_mul_ps(v_alpha, _mm256_loadu(x+i)),
        _mm256_loadu(y+i)));

for(; i<n; i++)
    y[i] += alpha * x[i];
```

Since AVX allows us to keep 8 single-precision floating point numbers in a 256-bit variable of type `__m256`, the first loop progresses in steps of 8, loads 8 coefficients from the arrays `x` and `y`, multiplies them by the scaling factor `alpha` (duplicated 8 times in `v_alpha`), and stores them in `y`.

In order to be able to handle dimensions that are not multiples of 8, we have to add a short loop that deals with the remainder.

Bandwidth In theory, we would expect an SIMD code with vectors of length 8 to run 8 times as fast as a standard code.

In practice, frequently the memory bandwidth is a limiting factor: with AVX, an `axpy` implementation take two clock cycles to treat 32 bytes of the arrays `x` and `y`, i.e., we could process 32 bytes per cycle. At a clock speed of 4 GHz, this means that we could theoretically process 128 GBytes per second. Unfortunately, even modern quad-channel memory only reaches 50 GBytes per second, so we will not see the expected speedup.

It is therefore very important to perform a large number of operations on each byte that has been read from memory, e.g., by rewriting algorithms using BLAS level 3.

Automatic vectorization Modern compilers are able to turn certain parts of programs written in standard C or FORTRAN into SIMD code.

Even for fairly simple situations, this involves a closer analysis of the program and is therefore usually only enabled for higher optimization levels.

Example: axpy If we use the GNU C compiler on an x86 processor, we can enable vectorization with the command-line parameter `-ftree-loop-vectorize` and use `-mavx` to choose the AVX command set.

```
for(i=0; i<n; i++)
    x[i] += y[i];
```

For this very simple vector addition, the compiler will detect that all components are independent and use SIMD instructions.

Aliasing If the addition loop seen above takes place in a function that receives the pointers `x` and `y` as parameters, *aliasing* may occur, i.e., `x` and `y` could point to the same array, and straightforward vectorization could change the result, e.g., `x=y+1` would introduce dependencies between different iterations. The compiler is not allowed to change the outcome of an algorithm by vectorization, so it may not vectorize the loop under these circumstances. In our simple example, the compiler can find out whether aliasing is possible by comparing the addresses of `x` and `y`: if they are further apart than 32, it is safe to use AVX instructions.

In general situations, we can help the compiler by using the `restrict` keyword: using this keyword in the declaration of pointer parameters, we assert that they will never point to the same arrays:

```

void
simple_axpy(int n, const float * restrict y, float * restrict x)
{
    int i;

    for(i=0; i<n; i++)
        x[i] += y[i];
}

```

Example: Rounding errors In some typical situations, rounding errors depend on the order in which floating-point operations are carried out:

```

float sum;
int i, n;

sum = 0.0f;
for(i=1; i<=n; i++)
    sum += i;

```

In this example, the computer has to convert `i` to a `float` value and perform a floating-point addition. Due to rounding, these additions are not associative, i.e., the sequence of their execution matters. Unfortunately, this means that SIMD instructions cannot be used without changing the result of the computation.

Some compilers offer us the possibility to state that reordering floating-point instructions is acceptable even if it changes the result (within the interval of rounding errors), e.g., by using the `-ffast-math` command-line parameter for the GNU C compiler. If this parameter is used in combination with `ftree-loop-vectorize`, the summation will be performed with SIMD instructions, and consequently run significantly faster.

SIMD with OpenMP OpenMP 4.0 introduces directives to help the compiler identify portions of the code that can be parallelized using SIMD instructions:

```

#pragma omp simd
for(i=0; i<n; i++)
    y[i] += alpha * x[i]

```

By using the `simd` directive, we state that the operations in the `for` loop are independent and can therefore be treated by SIMD instructions.

For our summation example, using OpenMP directives implies that we allow the compiler to reorder floating-point operations, even if it leads to different rounding errors:

```

#pragma omp simd reduction(+:sum)
for(i=1; i<=n; i++)
    sum += i;

```


If we want to call functions from within the vectorized loop, we should declare them using the `declare simd` directive, so the compiler can create vectorized versions of the functions:

```
#pragma omp declare simd
float
my_sinf(float x)
{
    float x2 = x*x;

    return x*(c0 + x2 * (c1 + x2 * (c2 + x2 * (c3 + x2 *
        (c4 + x2 * c5))))));
}
```

Using the `simd` directive does *not* guarantee that the compiler will vectorize the loop, e.g., because it considers it too short, because it is unable to determine its iteration count in advance, or because the loop body calls functions that have not been vectorized (unfortunately, this appears to currently include many functions in the standard math library).

Structure of arrays For sequential codes, it makes sense to arrange data that belongs to the same object in close proximity in memory, since this improves the utilization of caches and makes for a program that is easier to read.

An example could be the simulation of a gravitational field: each planet is represented by a `struct`, and all the planets as an array of these `structs`:

```
typedef struct {
    float x, y, z;
    float mass;
} planet;

planet pl[N];
```

In a vectorized code, this representation is unattractive, since we cannot load the x coordinates for several masses at once into a vector register.

If we plan to use a vectorized code, it is frequently advisable to use a *structure of arrays*, i.e., to define a `struct` that holds arrays for each variable:

```
typedef struct {
    float x[N], y[N], z[N];
    float mass[N]
} planets;
```

With this approach, it is straightforward to obtain the parameters for multiple planets in a form that can be directly used in an SIMD code.

7.2 Explicit vectorization

Example: Reciprocal square root We have already seen that the Newton iteration takes the particularly friendly form $\Phi(x) = x(\frac{3}{2} - \frac{1}{2}ax^2)$ if we approximate the reciprocal square root $1/\sqrt{a}$ by looking for a zero of $f(x) = x^{-2} - a$.

With AVX, we can obtain the following SIMD implementation:

```
__m256
avx_rsqrt(__m256 x)
{
    __m256 c15 = _mm256_set1_ps(1.5f);
    __m256 c05 = _mm256_set1_ps(0.5f);
    __m256 y;

    y = _mm256_rsqrt_ps(x);
    y = _mm256_mul_ps(y, _mm256_sub_ps(c15,
        _mm256_mul_ps(_mm256_mul_ps(c05, x),
            _mm256_mul_ps(y, y))));

    return y;
}
```

Compared to the simple C version `1.0f / sqrtf(x)`, we observe a speedup of more than 30, not only due to the SIMD approach, but also due to the fact that we can avoid computing the square root and the reciprocal separately.

Example: Sine function We consider approximating the sine function by its Taylor series

$$\sin(x) \approx x \left(1 - \frac{x^2}{3!} + \frac{x^4}{5!} - \frac{x^6}{7!} + \frac{x^8}{9!} - \frac{x^{10}}{11!} \right).$$

If we can ensure $x \in [-\pi/2, \pi/2]$, we can expect the error to be bounded by

$$\frac{(\pi/2)^{13}}{13!} \approx 5.7 \times 10^{-8},$$

which should be sufficient for single-precision computations. We can evaluate the truncated Taylor expansion by Horner's method

$$\sin(x) \approx x \left(1 - \frac{x^2}{2 \cdot 3} \left(1 - \frac{x^2}{4 \cdot 5} \left(1 - \frac{x^2}{6 \cdot 7} \left(1 - \frac{x^2}{8 \cdot 9} \left(1 - \frac{x^2}{10 \cdot 11} \right) \right) \right) \right) \right)$$

using the following AVX code:

```
const float c[] = { 1.0/6.0, 1.0/20.0, 1.0/42.0,
                    1.0/72.0, 1.0/110.0 };
```

```

__m256 x2 = _mm256_mul_ps(x, x);
__m256 one = _mm256_set1_ps(1.0f);
__m256 px;
int m = sizeof(c) / sizeof(float);
int i;

px = one;
while(m > 0)
    px = _mm256_sub_ps(one, _mm256_mul_ps(x2,
        _mm256_mul_ps(_mm256_set1_ps(c[--m]), px)));

px = _mm256_mul_ps(px, x);

```

Branching A major difficulty of SIMD programming is dealing with *data-dependent branches*, i.e., situations in which different instructions have to be executed depending on the data in the vectors.

Since SIMD units can only perform the same operation for all components of a vector, techniques like *masking* have to be used: either the operation can be disabled for certain components of the vector based on suitable binary flags, or the operation is always performed for all components and bitwise logical operations are applied to merge the appropriate results depending on the data.

Example: Avoiding branches The challenging part of the computation of the sine — with respect to vectorization — is not the Taylor series, but moving x to the interval $[-\pi/2, \pi/2]$ where it converges rapidly. In a first step we compute $z = \frac{x}{2\pi}$, round to the closest $k \in \mathbb{Z}$, and subtract k to get $\hat{z} = z - k \in [-1/2, 1/2]$. Now we have $\hat{x} := 2\pi\hat{z} = x - 2\pi k \in [-\pi, \pi]$, and since the sine is 2π -periodic, we have $\sin(x) = \sin(\hat{x})$.

If $\hat{x} > \pi/2$, i.e., $\hat{z} > 1/4$, we can replace \hat{x} with $\pi - \hat{x}$ (or \hat{z} with $1/2 - \hat{z}$), since

$$\sin(\pi - \hat{x}) = \sin(-\hat{x}) \cos(\pi) + \cos(-\hat{x}) \sin(\pi) = -\sin(-\hat{x}) = \sin(\hat{x}).$$

Similarly, if $\hat{x} < -\pi/2$, i.e., $\hat{z} < -1/4$, we can replace \hat{x} with $-\pi - \hat{x}$ (or \hat{z} with $-1/2 - \hat{z}$).

With AVX instructions, these operations can be realized as follows:

```

z = _mm256_mul_ps(x, r2pi);

z = _mm256_sub_ps(z, _mm256_round_ps(z, _MM_FROUND_TO_NEAREST_INT));

sgn = _mm256_and_ps(z, _mm256_castsi256_ps(
    _mm256_set1_epi32(0x80000000)));

mask = _mm256_cmp_ps(_mm256_xor_ps(z, sgn), quar, _CMP_GT_OS);
y = _mm256_sub_ps(_mm256_xor_ps(half, sgn), z);

```

```
z = _mm256_or_ps(_mm256_and_ps(mask, y),
                _mm256_andnot_ps(mask, z));
```

Here `r2pi`, `quad`, and `half` are the constants $\frac{1}{2\pi}$, $1/4$, and $1/2$, respectively.

The variable `sgn` stores the sign of z , and since this is just the highest bit in a 32-bit floating point number, we can use bitwise *xor* operations to flip the signs of z and $1/2$ to find out if $|z| > 1/2$ and to compute $y = \text{sgn}(z)/2 - z$.

The variable `mask` is of particular importance: the function `_mm256_cmp_ps` compares the components of a vector and sets all bits of the corresponding components of the result to one if the comparison yields true and to zero otherwise. In our case all bits will be one if and only if $|z| > 1/2$.

Once y has been computed, we can use bitwise *and* operations to clear the components with $|z| \leq 1/2$ and use a bitwise *or* operation to merge with the original values of z .

Latency and throughput Modern processors typically use an *instruction pipeline*, i.e., the execution of an instruction is split into several steps (e.g., decoding the instruction, computing addresses, obtaining data, computing, or storing the result), and each of these steps is performed by a separate *stage* of the processor.

This means that the processor can start decoding the next instruction before the previous instruction has finished, as long as the processor knows what the next instruction will be.

Every instruction has a *latency* and a *throughput*: the latency is the number of clock cycles required to complete processing the instruction, while the throughput is the number of clock cycles that need to pass before the next instruction can be moved into the pipeline.

On an Intel[®] Skylake processor, the floating-point addition has a throughput of 0.5, i.e., every cycle the processor can start working on two additions, but a latency of 4, i.e., we have to wait four cycles before we can use the result of the addition.

Some operations take significantly longer, e.g., Intel reports an upper bound of 18 for the latency of the double-precision square root and a throughput of less than 12, while older architectures could have a latency of up to 35 cycles and a throughput of up to 27.

Modern processors deal with latencies by executing instructions *out of order*, i.e., rearranging the order of the instructions without changing the result. An example could be the computation of Euclidean norms:

```
for(i=0; i+7<n; i+=8) {
    vx = _mm256_loadu_ps(x+i);
    vy = _mm256_loadu_ps(y+i);
    vz = _mm256_loadu_ps(z+i);

    sq = _mm256_add_ps(_mm256_mul_ps(vx, vx),
                      _mm256_add_ps(_mm256_mul_ps(vy, vy),
                                    _mm256_mul_ps(vz, vz)));
```

```

    _mm256_storeu_ps(nr+i, _mm256_sqrt_ps(sq));
}

```

While the square root is computed in the last step, a processor with *out-of-order* capabilities could already perform additions and multiplications for the next iteration of the loop and even already load the data for even later iterations.

Prefetching Since memory accesses have a *very* large latency, up to more than 100 cycles, relying on the processor to start fetching required data from memory in time can be unattractive.

This is particularly critical if predicting the addresses of future load instructions is not straightforward.

In these situations, *prefetching* might be of interest: using the function `_mm_prefetch`, we can indicate that we expect to use a certain cacheline soon, so that the processor can start fetching it into its cache while it is still busy executing other instructions.

```

for(i=0; i+7<n; i+=8) {
    _mm_prefetch(x+i+64);

    vx = _mm256_loadu_ps(x+i);

    _mm256_stream_ps(y+i, _mm256_mul_ps(vx, vx));
}

```

Using the `_mm256_stream_ps` function to write the result to memory indicates to the processor that we do not expect to use the data any time soon, so it does not have to be stored in the cache, leaving more cache space for important data.

8 Computing with graphics cards

8.1 Computing with graphics cards using CUDA

Since modern graphics algorithms require relatively complicated computations, i.e., to compute realistic shadows and lighting effects, modern graphics cards are far more flexible than their predecessors. This property makes them interesting for more general applications.

CUDA NVIDIA[®] was the first major manufacturer of graphics cards that offered a user-friendly way of programming their devices to handle general tasks, called CUDA, for *compute unified device architecture*.

CUDA programs are usually written in CUDA C, an extension of C++ that includes keywords and constructs for creating code for graphics cards, running this code, and transferring data between the computer (called the “host”) and the graphics card (called the “device”).

Cross-compilation CUDA C programs contain two kinds of functions: most functions are probably standard C/C++ functions that are compiled for the current host system, e.g., a Linux PC. Some functions are marked with the `__global__` or `__device__` keywords, these are compiled for NVIDIA’s graphics cards and cannot be called directly.

SIMT: Single instruction, multiple threads.

Since modern graphics processors have hundreds or thousands of “processing elements” that can execute code, parallelization takes the form of often hundreds of thousands of threads that are scheduled to run on these elements.

CUDA threads are significantly simpler than OpenMP threads:

- each thread has an instruction counter,
- the variables used by a thread are assigned statically to registers,
- there is no stack, so recursive function calls are impossible,
- only threads belonging to the same block can synchronize.

Execution is managed by *streaming multiprocessors*. Typically a graphics card contains several of these multiprocessors.

Every multiprocessor contains a number of processing elements that can carry out instructions, a number of registers that can be used for thread-local variables, and a small amount of local memory that can be used for inter-thread communication.

- Threads are organized in *blocks*. Only threads within the same block can communicate.
- Blocks are split into *warps*, consisting each of a fixed number of threads.

Each block is assigned to a multiprocessor, and the multiprocessor is responsible for executing the threads in this block. In every clock cycle, the multiprocessor chooses one of the block's warps, checks which instructions the threads in the warp have to execute next, and executes one of these instructions. If multiple threads need to execute the chosen instruction, all of them do.

In this regard, SIMT programming is much like SIMD: under optimal conditions, the same instruction is performed for all threads in the current warp.

Since every thread has its own instruction counter, SIMT programming can be significantly more user-friendly than SIMD programming.

Limits Since one multiprocessor is responsible for executing the threads in a block, and since every thread requires a certain number of registers, the blocks should not be too large.

Another limiting factor is that since one block is executed by one multiprocessor, we should at least have as many blocks as multiprocessors to ensure that all processors can work.

Since the multiprocessor can execute one instruction in one warp per clock cycle, we should make sure that at least one instruction is ready for execution, e.g., that the necessary data has been fetched from memory, therefore the blocks should also not be too small.

Another limiting factor is that since a multiprocessor has as many processing elements as the warp size (frequently 32), the blocks should not be smaller than the warp size in order to avoid idle processing elements.

Kernels Code that will be executed on the graphics card is frequently called a *kernel*. For the simple `axpy` operations, it could look like this:

```
__global__ void
axpy(int n, float alpha, const float *x, float *y)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;

    if(i < n)
        y[i] += alpha * x[i];
}
```

The keyword `__global__` tells the compiler that we intend to run the following function on the graphics card. Kernels cannot return results, so their return type always has to be `void`. Since graphics cards can handle a large number of threads with very low overhead, we can use kernels for each component of the vector.

The current thread is identified by the built-in read-only variables `blockIdx` and `threadIdx` that give us the index of the block and the index of the thread within the block. `blockDim` gives us the number of threads per block, so the first line of our function computes a global number for the current thread.

Since CUDA is able to handle up to three-dimensional index sets, we have to use `blockIdx.x`, `blockDim.x`, and `threadIdx.x` to denote the first component of the indices.

If the number of coefficients in the vectors `x` and `y` is not divisible by the block size, the thread number `i` could be larger than `n`. In this case, the thread simply returns. Otherwise, it updates the coefficient `y[i]`.

Kernels cannot be called like standard C functions, since they will be executed on the graphics card. In our case, we can use

```
axpy<<<(n+255)/256, 256>>>(n, alpha, dx, dy);
```

to create $256 \cdot \lceil n/256 \rceil \geq n$ threads with a block size of 256. Given that usually 32 threads are collected in a warp, this means that the multiprocessor can choose between eight warps in each cycle, which should be sufficient to hide latencies.

These threads are executed *asynchronously*, i.e., the host's main processor will execute the instructions following the call to `axpy` even if this call has not yet been completed. This is not a bug, but a feature: it allows the processor to do useful work while the graphics card is busy.

Device memory Kernels typically do not have access to the main memory managed by the host's processor, only to memory on the graphics card. This means that we have to allocate memory and transfer the relevant data:

```
cudaMalloc(&dx, sizeof(float) * n);
cudaMalloc(&dy, sizeof(float) * n);

cudaMemcpy(dx, x, sizeof(float) * n, cudaMemcpyHostToDevice);
cudaMemcpy(dy, y, sizeof(float) * n, cudaMemcpyHostToDevice);
```

In order to copy the kernel's results back to main memory, we can use

```
cudaMemcpy(y, dy, sizeof(float) * n, cudaMemcpyDeviceToHost);
```

Calls to `cudaMemcpy` are *blocking*, i.e., they only return after they have completed their work. They also wait for the completion of the kernel before they start. In particular, if we use

```
cudaMemcpy(dx, x, sizeof(float) * n, cudaMemcpyHostToDevice);
cudaMemcpy(dy, y, sizeof(float) * n, cudaMemcpyHostToDevice);
```



```
axpy<<<(n+255)/256, 256>>>(n, alpha, dx, dy);
```

```
cudaMemcpy(y, dy, sizeof(float) * n, cudaMemcpyDeviceToHost);
```

the kernel will not get started before the vectors `x` and `y` have been copied, and the result will not be copied back before the kernel has completed.

Managed memory Recent versions of CUDA offer an alternative to explicitly copying data between main and device memory: if we use `cudaMallocManaged` instead of `cudaMalloc`, the system automatically moves it where it is needed:

```
cudaMallocManaged(&x, sizeof(float) * n);  
cudaMallocManaged(&y, sizeof(float) * n);
```

```
for(i=0; i<n; i++) {  
    x[i] = i; y[i] = n-i;  
}
```

```
axpy<<<(n+255)/256, 256>>>(n, 1.0f, x, y);
```

```
cudaDeviceSynchronize();
```

```
for(sum=0.0f, i=0; i<n; i++)  
    sum += y[i];
```

In this setting, we have to call `cudaDeviceSynchronize` to ensure that the kernel completes its work before we start the final `for` loop on the main processor.

8.2 Streams, blocks, and shared memory

Transfer latency Although managed memory allows us to hide transfers from and to graphics memory in our program, these transfers still have to take place, and they may take a very long time: a NVIDIA GeForce[®] GTX 1080 Ti graphics card offers a memory bandwidth of 484 GB per second, but is connected to the host via a PCIe 3.0 bus with a maximum of 16 GB per second.

In order to avoid the relatively low speed of the PCIe bus hampering the performance of our code, we have to manage transfers between main and graphics memory carefully.

Hiding latencies If our code carries out a sufficient number of operations with each item of data transferred from main memory, most modern graphics cards allow us to hide memory transfers behind computations: the memory interface and the multiprocessors can work simultaneously, we only have to arrange our program in a way that allows it to take advantage of this property.

Streams CUDA allows us to create *streams*, sequences of operations that are scheduled by the system in order to try to minimize the overall runtime. The system guarantees that all operations within one stream are carried out in sequence, but operations in different streams can be interleaved.

If we use two streams containing read-compute-write sequences, one stream can perform transfers while the other performs computations.

Streams can be created using the function `cudaStreamCreate`, and some operations, particularly kernel executions and memory transfers, allow us to assign them to a stream for execution.

Pinned memory In order to utilize the hardware's ability to transfer data between main and graphics memory, the main memory area has to be *pinned*: typically programs work with *logical* addresses, and the operating system is free to move the corresponding data around.

In order for the graphics card to access the memory, it has to correspond to a fixed area of physical memory.

We can use the function `cudaMallocHost` to allocate pinned memory.

Example: Trapezoidal rule We assume that we have to compute integrals

$$v_i = \int_{a_i}^{a_i+1} f(x) dx \quad \text{for all } i \in [1 : n]$$

by the composite trapezoidal rule with m subintervals, but we do not have enough graphics memory to store all n values of v_i and a_i simultaneously.

We only have enough storage for `nb` of these integrals on the graphics card, so we use streaming to move the required data to and from graphics memory:

```

__global__ void
trapezoidal(int n, int m, const float *a, float *v)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j;
    float sum, h = 1.0f / m;

    sum = 0.5f * (func(a[i]) + func(a[i]+1.0f));
    for(j=1; j<m; j++)
        sum += func(a[i] + h * j);
    v[i] = h * sum;
}

cudaMallocHost(&a, sizeof(float) * n);
cudaMallocHost(&v, sizeof(float) * n);

cudaMalloc(da, sizeof(float) * nb);
cudaMalloc(da+1, sizeof(float) * nb);
cudaMalloc(dv, sizeof(float) * nb);
cudaMalloc(dv+1, sizeof(float) * nb);

cudaMemcpyAsync(da[0], a, sizeof(float) * nb,
                cudaMemcpyHostToDevice, stream[0]);
trapezoidal<<<nb, 256, 0, stream[0]>>>(n, m, da[0], dv[0]);
cudaMemcpyAsync(v, dv[0], sizeof(float) * nb,
                cudaMemcpyDeviceToHost, stream[0]);

current = 0;
for(i=0; i+nb<n; i+=nb) {
    next = 1-current;
    cudaMemcpyAsync(da[next], a+i+nb, sizeof(float) * nb,
                    cudaMemcpyHostToDevice, stream[next]);
    trapezoidal<<<nb, 256, 0, stream[next]>>>(nb, m,
        da[next], dv[next]);
    cudaMemcpyAsync(v+i+nb, dv[next], sizeof(float) * nb,
                    cudaMemcpyDeviceToHost, stream[next]);

    cudaStreamSynchronize(stream[current]);

    current = next;
}
cudaStreamSynchronize(stream[current]);

```

Blockwise collaboration If we are only computing a moderate number of outputs, it might be a good idea to spread the computation of one output among multiple threads, e.g., to let each thread sum up part of a composite quadrature rule.

Once the partial results are available, we have to combine them to obtain the final result. This is called a *reduction*, and instead of just summing up all results sequentially, we can employ multiple threads cooperating within a block.

In order to exchange partial results between the blocks, we use *shared memory* that is kept in the multiprocessor handling the block and is significantly faster (and smaller) than standard graphics memory.

```

__global__ void
midpoint(int n, int m, const float *a, float *v)
{
    __shared__ float p[256]; /* blockDim.x = 256 */
    int i = blockIdx.x;
    int it = threadIdx.x;
    int start = m * it / blockDim.x;
    int end = m * (it+1) / blockDim.x;
    float sum, h = 1.0f / m;
    int j, k;

    sum = 0.0f;

    for(j=start; j<end; j++)
        sum += func(a[i] + h * (j + 0.5f));

    p[it] = h * sum;

    k = blockDim.x;
    while(k > 1) {
        __syncthreads();
        k /= 2;
        if(it < k)
            p[it] += p[it+k];
    }

    if(it == 0)
        v[i] = p[0];
}

```

Each thread in a block computes part of the sum with the entries from `start` to `end-1` and stores the result in `p[it]`, i.e., in shared memory that is visible for all threads in the block.

Then follows the reduction: in the first step, the first 128 threads in the block read the partial results of the other 128 threads and add them together. Now the total sum

is the sum of the contributions of the first 128 threads.

In the second step, the first 64 threads read the results of the other 64 and add them. Now the total sum is the sum of the contributions of the first 64 threads.

We repeat this procedure until the first thread has the entire sum and stores it in `v[i]`.

In order to make sure that all threads are aware of the updated array `s` before each step, the function `__syncthreads()` is used to synchronize all threads in the block.

Since now each block corresponds to an entry of the result, we have to call the new kernel as

```
midpoint<<<n,256>>>(n, m, a, v);
```

so that n blocks are used.

9 Distributed computing

9.1 Message passing and MPI

Distributed computing Use several independent computers that can communicate only via a suitable network.

👍 Every computer can use its full memory bandwidth.

👍 If a computer breaks down, it can be easily replaced.

👎 Communication via network is slow.

👎 Exchange of data has to be performed explicitly.

MPI The *message-passing interface* is a well-established standard for distributing work among distributed computers.

Programs have to be compiled using a special compiler, usually named `mpicc`, and they have to be run using `mpirun` on simple systems or via a batch processing system on more sophisticated installations.

Example: Hello world We want two computers to produce a greeting. The first computer creates the zero-terminated string “Hello” and sends it to the second, the second receives it and prints it followed by the string “World”.

```
#include <stdio.h>
#include <mpi.h>

int
main(int argc, char **argv)
{
    char buf[6];
    int rank, size;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    if(rank == 0)
```

```

    MPI_Send("Hello", 6, MPI_CHAR, 1, 0, MPI_COMM_WORLD);

    if(rank == 1) {
        MPI_Recv(buf, 6, MPI_CHAR, 0, 0, MPI_COMM_WORLD,
                MPI_STATUS_IGNORE);
        printf("%s World\n", buf);
    }

    MPI_Finalize();
    return 0;
}

```

This simple program already uses a number of important MPI functions:

- `MPI_Init` initializes the MPI system.
- `MPI_Finalize` exits the MPI system.
- `MPI_Comm_rank` returns this process's number.
- `MPI_Comm_size` returns the number of processes.
- `MPI_Send` sends a message to a process.
- `MPI_Recv` receives a message from a process.

Communicators MPI processes are organized in *communicators* providing context for communication operations: they ensure that messages are received in the order in which they were sent, they provide unique “addresses” for all parties, and they have several other useful properties.

After the initialization, two communicators `MPI_COMM_WORLD` and `MPI_COMM_SELF` are available, containing all processes managed by the system and only the current process, respectively. More communicators containing arbitrary subsets of processes can be defined at runtime.

The number of processes, the *size* of the communicator, can be requested using `MPI_Comm_size`. Each process in the communicator has a unique number ranging from 0 to `size-1`.

Since a process can belong to multiple communicators, we have to state the communicator that we want to use for each MPI operation.

Messages MPI messages contain an array of data of a certain type, a communicator, the rank of the sender and the receiver within the context of the communicator, and a tag that may be used to distinguish different kinds of messages passed between the same processes.

Within a communicator, messages are received in the order they were sent in.

Blocking and non-blocking functions The function `MPI_Recv` used for receiving messages is *blocking*, i.e., it will return control to the caller only after the message has been received and the corresponding *receive buffer* is filled.

The function `MPI_Send`, on the other hand, may or may not be blocking, depending on the implementation. We only have to guarantee that it returns control after it has made sure that the contents of the *send buffer* can be overwritten without changing the message, e.g., because its contents have been copied to an internal buffer or because its contents have already reached their destination.

```
int x, y;

if(rank == 0) {
    x = func1();
    MPI_Send(&x, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
    MPI_Recv(&y, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
else if(rank == 1) {
    y = func2();
    MPI_Send(&y, 1, MPI_INT, 0, 0, MPI_COMM_WORLD);
    MPI_Recv(&x, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

This code fragment will work, i.e., exchange the values of `x` and `y` between the processes 0 and 1, if `MPI_Send` does not block in the current implementation, and it will wait forever if `MPI_Send` blocks, since neither of the two processes ever gets to call `MPI_Recv`, allowing the other to complete the send operation.

Asynchronous functions Most problems with colliding function calls can be resolved by using non-blocking versions of send and received operations:

```
int x, y;
MPI_Request in, out;

if(rank == 0) {
    x = func1();
    MPI_Isend(&x, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &out);
    MPI_Irecv(&y, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &in);
}
else {
    y = func2();
    MPI_Isend(&y, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, &out);
    MPI_Irecv(&x, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, &in);
}
MPI_Wait(&out, MPI_STATUS_IGNORE);
MPI_Wait(&in, MPI_STATUS_IGNORE);
```


These *immediate* versions of the send and receive operations return `MPI_Request` objects that can be used to wait for completion. Except for the separation between requesting the operation and waiting for its completion, the semantics are similar: for a send operation, we wait until the send buffer is no longer required, for a receive operation, we wait until the receive buffer has been filled.

Hiding communication Apart from avoiding deadlock situations, non-blocking operations also allow us to hide communication latencies: before processing a set of data, we can issue non-blocking requests for the next set. If the computation takes sufficiently long, the data may have arrived once the current data has been handled.

Example: Gravitation For the sake of simplicity, we assume that all planets have unit mass. Our task is to compute

$$p_i = \sum_{\substack{j=1 \\ j \neq i}}^n \frac{1}{\|x_j - x_i\|_2} \quad \text{for all } i \in [1 : n].$$

The computation itself is straightforward:

```
void
compute(int n, const double *x, const double *y,
        const double *m, double *p)
{
    double norm;
    int i, j;

    for(i=0; i<n; i++)
        for(j=0; j<n; j++) {
            norm = sqrt((x[3*i ]-y[3*i ])*(x[3*i ]-y[3*i ]) +
                        (x[3*i+1]-y[3*i+1])*(x[3*i+1]-y[3*i+1]) +
                        (x[3*i+2]-y[3*i+2])*(x[3*i+2]-y[3*i+2]));
            if(norm > 0.0)
                p[i] += 1.0 / norm;
        }
}
```

In order to parallelize the computation with m processes, we assume $n = n_p m$ and store only n_p masses in each process. The computation is split into n_p phases, and in the k -th phase, process α computes interactions between its masses and the masses of process $(\alpha + k) \bmod m$. We can concurrently pass these masses on the process $(\alpha + 1) \bmod m$, so they arrive in time for the next phase:

```
x = (double *) malloc(sizeof(double) * 3 * np);
y[0] = (double *) malloc(sizeof(double) * 3 * np);
```

```

y[1] = (double *) malloc(sizeof(double) * 3 * np);
p = (double *) malloc(sizeof(double) * np);

MPI_Isend(x, 3*np, MPI_DOUBLE, (rank+1)%size,
          0, MPI_COMM_WORLD, &out);
MPI_Irecv(y[0], 3*np, MPI_DOUBLE, (rank+size-1)%size,
          0, MPI_COMM_WORLD, &in);
compute(n, x, x, p);

current = 0;
for(i=1; i<size; i++) {
    MPI_Wait(&out, MPI_STATUS_IGNORE);
    MPI_Wait(&in, MPI_STATUS_IGNORE);

    MPI_Isend(y[current], 3*np, MPI_DOUBLE, (rank+1)%size,
              i, MPI_COMM_WORLD, &out);
    MPI_Irecv(y[1-current], 3*np, MPI_DOUBLE, (rank+size-1)%size,
              i, MPI_COMM_WORLD, &in);
    compute(n, x, y[current], p);

    current = 1-current;
}
MPI_Wait(&out, MPI_STATUS_IGNORE);
MPI_Wait(&in, MPI_STATUS_IGNORE);

```

Since the computation has a complexity of $\sim n_p^2$, while only data for n_p planets has to be transferred per phase, there is a good chance that the communication time can be hidden behind the computation.

9.2 Collective communication

The communication operations we have seen so far transfer one message from one source to one target process, these are called *point-to-point operations* in the MPI standard.

In some situations, it may be beneficial for all processes to participate in an operation. In this case, it makes sense to use what is called a *collective operation* in the standard.

Broadcast A typical example is the *broadcast* operation, i.e., the transfer of data from one process, called the *root*, to *all* other processes belonging to a communicator.

We could perform this operation by having the root process send the data via `MPI_Send` to all other processes, but there is a better approach: once the root has sent its data to one other process, a total of two processes has it, so *both* can send it to one other process *simultaneously*, bringing the total to four. If each process performs p send operations, we can reach 2^p processes, so for m processes, only $\lceil \log_2 m \rceil$ operations are required.

In some systems, communication between some processes is faster than between others, and taking the properties of the infrastructure into account may significantly improve performance. Since the user should not have to worry too much about the properties of the hardware, it is best to have an optimized implementation of MPI handle tasks like this. That's why the standard provides functions like `MPI_Bcast`:

```
int n, rank;

MPI_Comm_rank(MPI_COMM_WORLD, &rank);

if(rank == 0) {
    printf("Number of intervals?\n");
    scanf("%d", &n);
}

MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
```

Process zero requests data from the user and then shares it with all other processes.

Reduction Another important collective operation is the *reduction*: if we want to add numbers stored in all processes, or take the maximum of them, or perform any other associative operation, we can follow a similar pattern to the broadcast, only in reverse: in a first step, one half of the processes obtain data from the other half and applies the associative operations. Since the contributions of the second half has now been taken into account, we can repeat the procedure with the first half, leaving us with only a quarter of the processes. We again find that $\lceil \log_2 m \rceil$ operations are sufficient to combine data from m processes.

Typical applications are the computation of inner products

```
for(i=start; i<end; i++)
    sum_local += x[i] * y[i];
```

```

MPI_Reduce(&sum_local, &sum, 1, MPI_DOUBLE,
          MPI_SUM, 0, MPI_COMM_WORLD);

```

or of maximal runtimes

```

time_local = MPI_Wtime();

/* ... do something ... */

time_local = MPI_Wtime() - time_local;

MPI_Reduce(&time_local, &time, 1, MPI_DOUBLE,
          MPI_MAX, 0, MPI_COMM_WORLD);

```

If we consider Krylov methods like the conjugate gradient iteration, it would be useful to have the result of the reduction operation available in all processes instead of just the root:

```

for(i=start; i<end; i++)
    sum_local += x[i] * y[i];

MPI_Allreduce(&sum_local, &sum, 1, MPI_DOUBLE,
             MPI_SUM, MPI_COMM_WORLD);

```

Example: Clustering The tree algorithm for evaluating non-local potentials and forces relies on a decomposition of space into boxes. If we want to parallelize this algorithm, we can assign one box to each process, since this will reduce the interactions between processes and therefore the necessary amount of communication.

We can approach this task by first locally sorting the points into boxes, exchanging information about the size of these boxes, and finally gathering the points for each box in the responsible process. For the sake of simplicity, we consider only the one-dimensional case in this example, i.e., point coordinates are just one `double` variable.

We assume that each process has a certain number of points, not necessarily corresponding to “its” box. In a first step, we sort these points into the available boxes and obtain three arrays: `ny` contains the number of points in each box, `y` contains the points, ordered by box, and `yoff` contains the offset for each box in `y`.

We use a collective *all-to-all* operation to let each process know how many points in “its” box the other processes have:

```

nx = (int *) malloc(sizeof(int) * size);

MPI_Alltoall(ny, 1, MPI_INT,
             nx, 1, MPI_INT, MPI_COMM_WORLD);

```

This operation ensures that in process i , the value $ny[i]$ of process j can be found in $nx[j]$. In our case, $nx[j]$ is the number of points presently in process j that belong in the current process's box.

Now we can compute the total number of points in “our” box and offsets to the final array:

```
xoff = (int *) malloc(sizeof(int) * size);
xoff[0] = 0;
for(j=1; j<size; j++)
    xoff[j] = xoff[j-1] + nx[j-1];
m = xoff[size-1] + nx[size-1];
```

We allocate the required storage and use the *variable-length all-to-all* operation to gather the points from the other processes:

```
x = (double *) malloc(sizeof(double) * m);
MPI_Alltoallv(y, ny, yoff, MPI_DOUBLE,
             x, nx, xoff, MPI_DOUBLE, MPI_COMM_WORLD);
```

Now x contains all points in “our” box, first the ones from process 0, followed by those from process 1, and so on.

One-sided communication So far, communication via MPI required the sending and the receiving process to cooperate. In some situations, this poses a major difficulty, since the sender has to predict which data the receiver will need.

As an example, consider the evaluation of the potential in the tree algorithm: if we want to evaluate at a point x for a cluster s that is stored in another process, we have to obtain information on this cluster, check the admissibility condition, and decide whether we want to evaluate directly or switch to the sons. It can be hard for the process owning s to predict what “our” process will do.

This problem is solve by the *one-sided communication operations* introduce in version 2 of the MPI standard: all processes in one communicator collectively define *windows* into their own address spaces: a double array of size n could be made available using

```
MPI_Win_create(a, n*sizeof(double), sizeof(double), MPI_INFO_NULL,
              MPI_COMM_WORLD, &window);
```

The third argument is the stride of the array, and using `sizeof(double)` allows us to use array indices as in C.

Since the creation of an `MPI_Win` object is a collective operation, the variable `window` is now available in all processes.

In order to read the entries $a[5]$, $a[6]$ and $a[7]$ of process i into an array b , we use

```
MPI_Get(b, 3, MPI_DOUBLE, i, 5, 3, MPI_DOUBLE, window);
```

We can also write four entries of c into $a[10]$ to $a[13]$ by using

```
MPI_Put(c, 4, MPI_DOUBLE, i, 10, 4, MPI_DOUBLE, window);
```

Before these operations, we have to synchronize the windows in the different processes by calling `MPI_Win_fence`:

```
MPI_Win_create(a, n*sizeof(double), sizeof(double), MPI_INFO_NULL,  
              MPI_COMM_WORLD, &window);
```

```
MPI_Win_fence(0, window);
```

```
MPI_Get(b, 3, MPI_DOUBLE, i, 5, 3, MPI_DOUBLE, window);  
MPI_Put(c, 4, MPI_DOUBLE, i, 10, 4, MPI_DOUBLE, window);
```

```
MPI_Win_fence(0, window);
```