

**О р д е н а   Л е н и н а**  
**ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ**  
**имени М.В.Келдыша**  
**Р о с с и й с к о й   а к а д е м и и   н а у к**

**А.А.Белов, А.С.Вергазов, Н.Н.Калиткин**

**Контроль точности**  
**при численном интегрировании**  
**жестких систем**

**Москва — 2020**

**Белов А. А., Вергазов А. С., Калиткин Н. Н.**

### **Контроль точности при численном интегрировании жестких систем**

Ранее для численного решения жестких систем обыкновенных дифференциальных уравнений было предложено а) использовать в качестве аргумента длину дуги интегральной кривой и б) выбирать оптимальный шаг интегрирования пропорционально  $\kappa^{-2/5}$ , где  $\kappa$  – кривизна интегральной кривой. В данной работе построена тестовая задача, в которой точное решение представлялось через элементарные функции как аргумента времени  $t$ , так и аргумента дуги  $l$ . Это позволило провести количественное сравнение различных разностных схем. Показано, что при расчетах с оптимальным шагом удастся использовать даже явные схемы Рунге-Кутты. При этом схема первого порядка давала невысокую точность, но очень высокую надежность даже при огромной жесткости. С повышением порядка точности надежность схем ухудшалась.

Предложена смешанная стратегия. На первом этапе по надежной схеме первого порядка строится оптимальная сетка, адаптированная к решению. На втором этапе эта сетка сгущается по правилу дробления квазиравномерных сеток, а расчет выполняется по схеме четвертого порядка точности. Смешанная стратегия дает одновременно хорошую надежность и высокую точность расчета.

**Ключевые слова:** дифференциальные уравнения, задача Коши, жесткие задачи, оптимальный шаг, смешанная стратегия расчета

***Aleksandr Aleksandrovich Belov, Artem Sergeevich Vergazov, Nikolay Nikolaevich Kalitkin***

### **Accuracy control in stiff system integration**

Previously, for numerical solution of stiff systems of ordinary differential equations, it was proposed to a) use length of the integral curve as the argument and b) choose optimal integration step proportional to  $\kappa^{-2/5}$ , where  $\kappa$  is the curvature of the integral curve. In this work, we construct a test problem in which the exact solution is expressed via elementary functions for both time and arc length arguments. This permitted quantitative comparison of various differential schemes. We show that even explicit Runge-Kutta methods are applicable in calculations with optimal step. The first order scheme provides low accuracy but very high reliability even at enormous stiffness. As order of accuracy increases, reliability of the schemes decreases.

We propose mixed computation strategy. At the first stage, an optimal mesh adapted to solution is built via the first order scheme. At the second stage, this mesh is thickened according to the rule of quasi-uniform meshes splitting and the calculation is done via the scheme with the fourth order of accuracy. The mixed strategy allowed to achieve both good reliability and high accuracy of calculation.

**Key words:** differential equations, Cauchy problem, stiff problems, optimal step, mixed strategy of calculation

Работа выполнена при поддержке Российского фонда фундаментальных исследований, проект 18-01-00175.

# 1. Методы численного интегрирования жестких систем

**1.1. Длина дуги.** Численное интегрирование задачи Коши для жестких систем обыкновенных дифференциальных уравнений (ОДУ) является одной из очень трудных задач вычислительной математики. Формально эта задача имеет следующий вид:

$$\frac{d\mathbf{u}}{dt} = \mathbf{f}(\mathbf{u}, t), \quad 0 \leq t \leq T, \quad \mathbf{u}(0) = \mathbf{u}^0. \quad (1)$$

Здесь  $t$  – скаляр,  $\mathbf{u}(t) = \{u_m(t), 1 \leq m \leq M\}$  и  $\mathbf{f}(\mathbf{u}, t) = \{f_m(\mathbf{u}, t), 1 \leq m \leq M\}$  – векторные функции, а  $M$  – размерность системы.

Задачу традиционно считают жесткой, если  $T \|\mathbf{f}(\mathbf{u}, t)\| \gg 1$  хотя бы в части отрезка  $t \in [0, T]$ . Решению жестких задач посвящена обширная литература, наиболее подробный обзор которой дан в монографии [1-2].

Различают два типа прикладных задач. Первый тип – это так называемые большие задачи, или задачи со многими процессами. Они описываются системами уравнений в частных производных, к которым подключена одна или несколько систем ОДУ. Примером могут служить задачи горения и взрыва. В них процесс горения, то есть реакции химических веществ, описывается системой уравнений химической кинетики; это система ОДУ. Выделяющееся тепло приводит к движению вещества. Это движение описывается уравнениями газодинамики, то есть системой уравнений в частных производных.

Расчет уравнений в частных производных гораздо более трудоемок, чем решение ОДУ. Поэтому в больших задачах шаг по времени  $\tau$  определяется требованиями методов решения уравнений в частных производных. Этот же шаг  $\tau$  вычислитель вынужден использовать для решения сопутствующей системы ОДУ. В этом случае естественным аргументом для решения задачи (1) является время  $t$ .

Второй класс задач содержит только систему ОДУ. Для этого класса задач оказывается выгоднее выбрать другой аргумент – длину дуги интегральной кривой в многомерном пространстве. Выполним переход к этому аргументу в два этапа.

Во-первых, добавим новую неизвестную функцию  $u_0 \equiv t$ , для нее имеем  $du_0/dt = 1$ . Тогда система (1) преобразуется к виду

$$\frac{du_m}{dl} = f_m(u_0, u_1, \dots, u_M), \quad 0 \leq m \leq M; \quad f_0(u_0, u_1, \dots, u_M) \equiv 1. \quad (2)$$

Общее число неизвестных теперь равно  $M + 1$ . Однако теперь правые части формально не содержат аргумента  $t$  в правых, то есть система становится автономной. Такую автономизацию называют тривиальной.

Во-вторых, определим элемент длины дуги интегральной кривой соотношением

$$dl^2 = \sum_{m=0}^M du_m^2 = (1 + \sum_{m=1}^M f_m^2) dt^2. \quad (3)$$

Тогда система (2) преобразуется к виду

$$\frac{du_m}{dl} = \frac{f_m(u_0, u_1, \dots, u_M)}{\sqrt{\sum_{m=0}^M f_m^2(u_0, u_1, \dots, u_M)}} \equiv F_m(u_0, u_1, \dots, u_M), \quad 0 \leq m \leq M. \quad (4)$$

При аргументе  $l$  сумма квадратов правых частей (4) равна 1, то есть норма правой части никогда не бывает большой. Это облегчает численное интегрирование системы. Такой прием полезен даже для нежестких ОДУ, а для жестких он кардинально облегчает решение задачи. Переход к длине дуги и различные преимущества этого метода подробно описаны в монографии [3]. Заметим, что правые части (4) не содержат аргумента  $l$ , то есть эта система является автономной.

Переход к длине дуги практически не увеличивает трудоемкость одного шага численных расчетов. Обычно основное время расчетов уходит на вычисление правых частей  $f_m$ , которые могут быть достаточно сложными функциями своих аргументов. Переход от (1) к (4) включает лишь несколько дополнительных арифметических операций, трудоемкость которых невелика. Размерность системы увеличивается на 1, что так же мало существенно для прикладных задач, где  $M$  обычно довольно велико.

**1.2. Выбор шага.** Расчет с постоянным шагом по времени  $\tau$  или по длине дуги  $h$  обычно невыгоден. Шаг целесообразно уменьшать там, где решение быстро меняется, то есть правые части ОДУ велики. На участках слабого изменения решения шаг можно увеличивать. В [2] подробно описаны алгоритмы автоматического выбора шага, принятые в мировой литературе. Традиционно используют два основных метода выбора шага. В первом методе каждый шаг выполняют по некоторой схеме  $(p+1)$ -го порядка точности, в которую вложена схема  $p$ -го порядка точности. Результат вложенной схемы берут в качестве ответа. По разности результатов двух схем выбирают величину следующего шага. Во втором методе шаг  $\tau$  рассчитывают повторно, разбив его на два шага величиной  $\tau/2$ . По разности этих расчетов вычисляют новый шаг.

На основе этих методов написано много пакетов программ. Большинство из них хорошо работают на нежестких задачах. Однако проверка этих пакетов на тестах с известными точными решениями показывает, что реальная точность расчетов лишь по порядку величины близка к запросу пользователя. Она обычно оказывается в несколько раз хуже или лучше. Поэтому в расчетах

прикладных задач, где ответ неизвестен, пользователь не может быть вполне уверен в достижении заданной им точности.

Намного хуже ситуация для случая жестких задач. Для них расчеты на тестах показали [4-6], что реальная точность может быть до  $10^8$  раз хуже заявленной. Это относится даже к таким тщательно выверенным программам, как пакеты Гира или программы Дормана-Принса `dopri5`. Кроме того, на жестких задачах возможны “срывы” шага [2]: иногда на участках слабо меняющегося решения программа без видимых причин уменьшает шаг в 100-1000 раз. Затем шаг постепенно увеличивается, но снова срывается; это может повторяться неоднократно.

В [4] был предложен принципиально другой алгоритм автоматического выбора шага. Он основан на использовании длины дуги и кривизны  $\kappa$  интегральной кривой в многомерном пространстве. Интуитивно понятно, что чем больше кривизна  $\kappa$ , тем меньше должен быть шаг  $h$ . Но каким должен быть алгоритм, связывающий эти две величины?

**1.3. Оптимальный шаг.** Напомним определение кривизны кривой  $\mathbf{u}(t)$  в  $(M + 1)$ -мерном пространстве с координатами  $\{t, u_1, u_2, \dots, u_M\}$ . Касательная к этой кривой определяется через производную  $d\mathbf{u}/dt$ . Деля на длину этого вектора, получим единичный вектор направления касательной:

$$\mathbf{n} = \frac{d\mathbf{u}}{dt} / \left\| \frac{d\mathbf{u}}{dt} \right\|_2. \quad (5)$$

Вектор кривизны и кривизна определяются через производную вектора направления касательной по длине дуги:

$$\mathbf{\kappa} = \frac{d\mathbf{n}}{dl}, \quad \kappa = \|\mathbf{\kappa}\|_2. \quad (6)$$

Таким образом, кривизна является вектором даже в случае плоской кривой, то есть одного ОДУ; это радиус-вектор окружности, имеющей касание второго порядка с кривой  $\mathbf{u}(t)$ . Наряду с этим говорят о скалярной кривизне  $\kappa$ , которая равна величине радиуса этой окружности.

В ранних работах [4,5] с помощью многих численных экспериментов была подобрана неплохая эвристическая закономерность:  $h \sim \kappa^{-0.5}$ . Однако в пакетах прикладных задач [7] был использован постоянный шаг по длине дуги и несколько семейств, в которых конкретные схемы имели порядок точности от 1 до 4. Эти пакеты позволили провести расчеты ряда прикладных задач, например – образование окислов фосфора и серы при горении различных топлив.

Задача теоретического нахождения зависимости  $h(\kappa)$  крайне сложна. Из общих соображений понятно, что результат должен зависеть от того, какая именно разностная схема используется для интегрирования системы ОДУ.

Однако удалось найти такой случай [8], для которого теоретически обосновывается формула выбора оптимального шага. Пусть схема интегрирования имеет точность  $O(h)$ ; это может быть явная или неявная схема Эйлера или явно-неявная L1-устойчивая схема Розенброка. Тогда

$$h_{opt}(\kappa) = \text{const} \cdot \kappa^{-2/5}. \quad (7)$$

Качественный вид этой формулы совпадает с ранее найденным эвристическим видом, а оптимальный показатель степени лишь слабо отличается от эвристического.

Как часто оказывается, чисто теоретическая формула (7) для своего практического применения требует “кухонных” поправок.

Во-первых, ясно, что при  $\kappa = 0$  она дает  $h = \infty$ . Формально это правильно: если на некотором участке кривой  $\kappa = 0$ , этот участок есть прямая, то есть  $u(t)$  является линейной функцией. А для линейной функции любая численная схема дает точный ответ при сколь угодно большом шаге.

На практике  $\kappa = 0$  лишь в отдельных точках кривой (например, в точках перегиба), а попадание узла расчетной сетки именно в эту точку имеет нулевую вероятность. Однако возможно попадание счетного узла в малую окрестность точки перегиба и получение неприемлемо большого шага  $h$ . Надо ввести такую поправку, чтобы в любом случае число интервалов было не меньше некоторого разумного  $N_{min}$ .

Во-вторых, надо дать разумное определение константы в (7). В нее должна быть включена некоторая интегральная нормировка по длине дуги, обеспечивающая желательное количество интервалов сетки  $N_{max}$ .

Исходя из этого, в [8] был предложен следующий алгоритм. Пусть требуется решить задачу (4) на отрезке  $[0, L]$  с априорно заданными  $N_{max}$ ,  $N_{min}$ . Несколько видоизменим формулу (7) с учетом сделанных замечаний:

$$h(l) = \left[ \frac{N_{min}}{L} + \frac{N_{max} \kappa^{2/5}}{\int_0^L \kappa^{2/5}(l') dl'} \right]^{-1}. \quad (8)$$

В точках с очень малой кривизной  $\kappa(l) \approx 0$  она дает  $h(l) \approx L/N_{min}$ . В точках с очень большой кривизной формула (8) переходит в (7); при этом константа такова, что  $\int h^{-1}(l) dl \approx N_{max}$ . При любых  $\kappa(l)$  не может получиться  $h = 0$  или  $h = \infty$ . Это свидетельствует о разумности формулы (8).

**1.4. Двухэтапная стратегия.** Численный расчет может называться хорошим, если он гарантирует пользователю требуемую точность. Для нежестких задач, когда в качестве аргумента выбирают время  $t$ , этого нетрудно

добиться. Расчет до заданного момента  $T$  проводят на равномерной сетке сначала с некоторым умеренным числом шагов  $N$ . Затем увеличивают  $N$  в 2 раза, одновременно уменьшая  $\tau$  вдвое, и проводят новый расчет. При этом узлы  $t_n$  первой сетки точно совпадают с четными узлами удвоенной сетки. В совпадающих узлах находят разность двух сеточных решений на соседних сетках  $u(t)$  и вычисляют норму разности (обычно это нормы  $L_2$  или  $C$ ). В этом случае для оценки погрешности применим метод Ричардсона [1, 9-12]. Если заданная точность не будет достигнута, то процедуру удвоения сетки повторяют. Этот метод дает асимптотически точную оценку погрешности.

Однако для жестких задач равномерная сетка непригодна. Чтобы шаг равномерной сетки стал достаточно малым на участках быстрого изменения решения, нужны огромные числа узлов  $N$ . Даже переход к аргументу  $l$  не позволяет использовать равномерные сетки.

Существуют так называемые квазиравномерные сетки [11,12]. Для них разработана процедура удвоения  $N$ , позволяющая пользоваться методом Ричардсона. Но для построения такой сетки нужно заранее знать поведение всей функции  $u(l)$  и строить густую сетку там, где функция быстро меняется. Однако до начала решения задачи мы не знаем, где расположены такие участки. Поэтому в [4,5] была предложена двухэтапная процедура сгущения сеток. С учетом новейших работ она формулируется следующим образом.

**Первый этап.** Для начала расчета выбираем некоторые разумные значения  $N_{min}$  и  $N_{max}$ , входящие в формулу шага (8). В эту формулу входят также  $L$  и интеграл. До начала расчетов они неизвестны, поэтому задаются их правдоподобные оценки. После этого проводят расчет с выбором шага согласно (8) и продолжают его до тех пор, пока  $u_0(l_n) \equiv t_n$  не превысит значение  $T$ . Полученный номер узла есть  $N$  для построенной сетки  $l_n$ . При этом автоматически вычисляется значение полной длины дуги для данной сетки  $L = l_N$ . Попутно в ходе расчета вычисляют интеграл в (8) по какой-нибудь квадратурной формуле; допустима даже формула левых прямоугольников.

Найденное в этом расчете значение  $N$  может не совпадать с исходным значением  $N_{max}$ . Разумеется, вычисленные значения  $L$  и интеграла в (8) не будут совпадать с теми величинами, которые использовались для начала расчета.

Затем удвоим значения  $N_{min}$ ,  $N_{max}$  и построим сгущенную примерно вдвое сетку. В этом расчете будем использовать значения интеграла и  $L$ , найденные на предыдущей сетке. Практика показывает, что при этом  $N$  окажется примерно вдвое большим, чем на предыдущей сетке. Однако новые значения интеграла и  $L$  могут сильно отличаться от предыдущих значений, поскольку на первой сетке они были выбраны произвольно. Очевидно, что новая сетка будет примерно вдвое подробнее предыдущей. Однако ее четные узлы не будут совпадать с узлами предыдущей сетки. Поэтому оценка точности по правилу Ричардсона невозможна.

Снова удвоим  $N_{min}$ ,  $N_{max}$  и построим третью сетку. На третьей сетке значение  $N$  уже почти удваивается по сравнению со второй сеткой, а значения интеграла и  $L$  будут близки к соответствующим величинам второй сетки. Значения четных узлов третьей сетки так же будут существенно ближе к узлам второй сетки. Однако правило Ричардсона применять по-прежнему нельзя.

Эту процедуру удвоения сетки повторяем до тех пор, пока четные узлы новой сетки не станут достаточно близкими к узлам предыдущей сетки. Различные критерии близости сеток сравнивались в [13]; наилучшим был признан следующий критерий. Обозначим шаги исходной сетки через  $h_n = l_n - l_{n-1}$ ,  $1 \leq n \leq N$ . Шаги удвоенной сетки обозначим через  $\hat{h}_n$ ,  $1 \leq n \leq \hat{N}$ ; на практике  $\hat{N} \approx 2N$ , хотя точного равенства может не быть. Шагу  $h_n$  в удвоенной сетке соответствует сумма шагов  $\hat{h}_{2n-1} + \hat{h}_{2n}$ . Будем считать сетки близкими, если

$$\sqrt{\frac{1}{N} \sum_{n=1}^N \left( \sqrt{\xi_n} - \frac{1}{\sqrt{\xi_n}} \right)^2} \leq \eta, \quad \xi_n = \frac{\hat{h}_{2n-1} + \hat{h}_{2n}}{h_n}, \quad \eta = \text{const}; \quad (9)$$

разумную величину  $\eta$  надо подбирать на основе практических расчетов. При этом значения интеграла и  $L$  на этих сетках будут совпадать с приемлемой точностью. При выполнении критерия (9) первый этап заканчивается.

**Второй этап.** Будем считать, что последняя сетка первого этапа уже достаточно хорошо адаптирована к участкам быстрого изменения решения. Поэтому дальнейшие сгущения сеток будем производить так, чтобы можно было применять метод Ричардсона для оценки погрешности. Для этого надо проводить сгущения так, чтобы последовательность удвоенных сеток была квазиравномерной. Простейший способ такого сгущения имеет следующий вид.

Пусть исходная сетка имеет узлы  $l_n$ ,  $0 \leq n \leq N$  и шаги  $h_n = l_n - l_{n-1}$ . Пусть интервал сетки  $[l_{n-1}, l_n]$  является внутренним:  $2 \leq n \leq N-1$ . Тогда узлы этого интервала являются четными узлами новой сетки  $\hat{l}_{2n-2}$  и  $\hat{l}_{2n}$ , а шаг  $h_n$  делится на два шага новой сетки  $\hat{h}_{2n-1}$  и  $\hat{h}_{2n}$  по следующим формулам:

$$\hat{h}_{2n-1} = h_n \frac{\sqrt[4]{h_{n-1}}}{\sqrt[4]{h_{n-1}} + \sqrt[4]{h_{n+1}}}, \quad \hat{h}_{2n} = h_n \frac{\sqrt[4]{h_{n+1}}}{\sqrt[4]{h_{n-1}} + \sqrt[4]{h_{n+1}}}, \quad 2 \leq n \leq N-1. \quad (10)$$

Для левого граничного интервала шаг  $h_1$  делится по формулам

$$\hat{h}_1 = h_1 \frac{\sqrt{h_1}}{\sqrt{h_1} + \sqrt{h_2}}, \quad \hat{h}_2 = h_1 \frac{\sqrt{h_2}}{\sqrt{h_1} + \sqrt{h_2}}. \quad (11)$$



Правый граничный шаг  $h_N$  делится следующим образом:

$$\hat{h}_{2N-1} = h_N \frac{\sqrt{h_{N-1}}}{\sqrt{h_{N-1}} + \sqrt{h_N}}, \quad \hat{h}_{2N} = h_N \frac{\sqrt{h_N}}{\sqrt{h_{N-1}} + \sqrt{h_N}}. \quad (12)$$

При таких формулах сгущения сетки число интервалов точно удваивается, а полная длина дуги  $L$  не меняется.

Можно доказать, что описанный способ удвоения сетки порождает последовательность квазиравномерных сеток. На такой последовательности правомерно однократное применение метода Рундсона: по каждой паре соседних сеток можно находить разности функций в совпадающих узлах  $l_n = \hat{l}_{2n}$ , вычислять нормы этих разностей и определять асимптотически точное значение погрешности по известному порядку точности формулы интегрирования.

Однако этот способ имеет одно ограничение: нельзя пользоваться рекуррентным методом Рундсона, то есть повышать порядок точности, используя 3 и более соседние сетки.

## 2. Апробация на тестах

**2.1. Требования к тестам.** Общепринятым способом апробации различных численных методов является расчет тестовых задач. Хороший тест должен удовлетворять нескольким требованиям. Во-первых, он должен содержать типичные трудности, на преодоление которых ориентирован исследуемый численный метод. Во-вторых, у него должно существовать легко реализуемое точное решение. Обычно под этим подразумевают, что точное решение при любых значениях аргументов легко вычисляется с любой требуемой точностью. Реально для этого нужно, чтобы решение достаточно просто выражалось через элементарные функции; запись через специальные функции нежелательна, так как их реализация с любым требуемым числом значащих цифр далеко не всегда доступна. В-третьих, желательно, чтобы тест содержал один или несколько параметров, которыми можно регулировать его жесткость и другие качественные свойства.

При наличии такого теста апробация метода несложна. Проводят расчет исследуемым методом, причем этот метод выбирает сетку по своим естественным правилам. В узлах этой сетки вычисляют точное решение и непосредственно находят точное значение погрешности, то есть разности между сеточным и точным решением. При этом можно вычислить любую норму погрешности. Это будет достоверная оценка.

Задачи, в которых решение известно только в отдельных реперных точках, непригодны в качестве теста. Во-первых, они не позволяют найти нормы ошибки. Во-вторых, для сравнения с ними приходится проводить решение на

таких сетках, некоторые узлы которых попадают в реперные точки. Для большинства алгоритмов такое построение сеток неестественно.

Для обыкновенных дифференциальных уравнений с аргументом  $t$  сравнительно легко конструируются различные тесты, в том числе для случаев высокой и сверхвысокой жесткости. Однако при переходе к аргументу  $l$  проблема существенно усложняется. Недостаточно, чтобы в элементарных функциях выражалось  $\mathbf{u}(t)$ . Нужно наличие решения в элементарных функциях так же для  $\mathbf{u}(l)$ . При этом все формулы не только прямого вычисления решения, но и вычисления обратных функций так же должны выражаться в элементарных функциях, причем достаточно просто. В противном случае, вычисление погрешностей при переходе от одного аргумента к другому будет практически невозможным.

До сих пор такие тесты не были построены. Нам удалось сконструировать два требуемых теста.

**2.2. Гиперболический тест.** Рассмотрим тестовое уравнение, содержащее один свободный параметр

$$\frac{du}{dt} = f(u) \equiv \text{sh}(\lambda u), \quad \lambda > 0, \quad u(0) = u^0 > 0. \quad (13)$$

Точное решение этой задачи имеет следующий вид:

$$u(t) = \frac{1}{\lambda} \ln \frac{1+B(t)}{1-B(t)}, \quad B(t) = e^{\lambda t} \text{th}(\lambda u^0/2). \quad (14)$$

Подставляя выражение (14) в формулу (2) и интегрируя, получим выражение для длины дуги:

$$l(t) = \frac{1}{\lambda} \ln \frac{\text{sh}[\ln(1+B(t)) - \ln(1-B(t))]}{\text{sh}(\lambda u^0)}. \quad (15)$$

Далее увидим, что жесткость задачи (13) быстро увеличивается с ростом  $\lambda$ . В задачу (13) нетрудно ввести еще один параметр — множитель  $\alpha$  в правой части. Однако это не представляет интереса, так как заменой  $\tilde{u} = u/\alpha$  и  $\tilde{\lambda} = \alpha\lambda$  задача сводится к форме (13).

Исследуем качественный вид решения (14). Функция  $u(t)$  положительна и монотонно возрастает. Она определена на конечном отрезке, поскольку

$$u(t) \rightarrow +\infty \text{ при } t \rightarrow t_* = \frac{1}{\lambda} \ln(\text{cth}(\lambda u^0/2)) > 0. \quad (16)$$

Длина дуги при этом так же неограниченно возрастает:  $l(t) \rightarrow +\infty$  при  $t \rightarrow t_*$ . Качественный вид решения  $u(t)$  изображен на рис. 1. Решение имеет продолжение влево на полуось  $-\infty < t \leq 0$ .

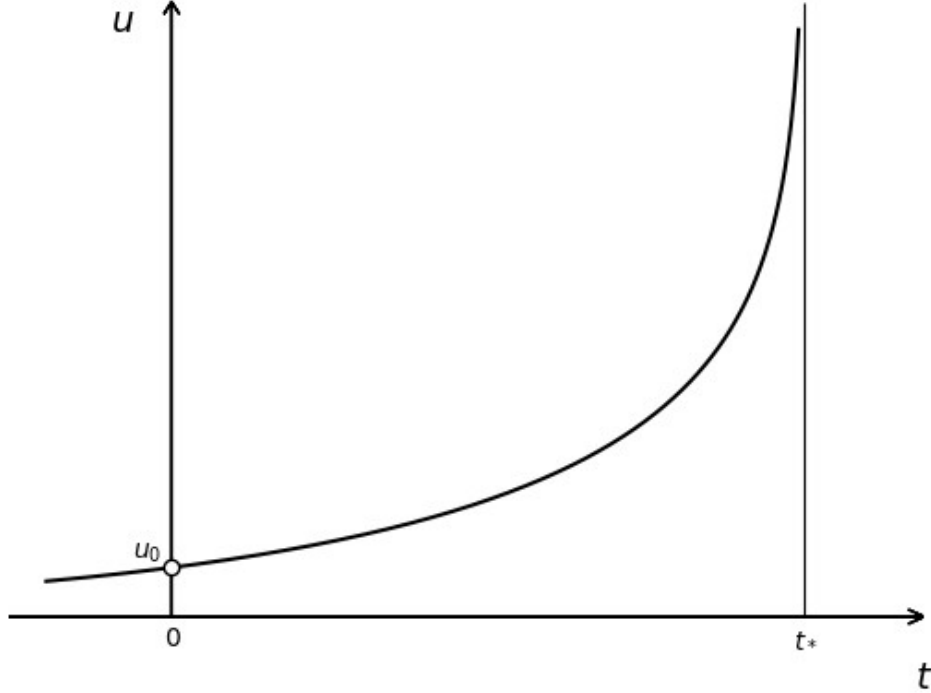


Рис. 1. Точное решение (14);  $\circ$  – начальная точка.

Перейдем в уравнении (13) к переменной  $l$ . Учтем, что  $1 + f^2(u) = \text{ch}^2(\lambda u)$ . Тогда форма записи (4) примет следующий вид:

$$\frac{du}{dl} = \text{th}(\lambda u), \quad \frac{dt}{dl} = \frac{1}{\text{ch}(\lambda u)}, \quad u(0) = u^0, \quad t(0) = 0, \quad 0 \leq l < +\infty. \quad (17)$$

Для первого уравнения системы (17) нетрудно написать точное решение:

$$u(l) = \frac{1}{\lambda} \ln(A(l) + \sqrt{A^2(l) + 1}), \quad A(l) = e^{\lambda l} \text{sh}(\lambda u^0), \quad 0 \leq l < +\infty. \quad (18)$$

Функция  $t(l)$  так же явно выражается через длину дуги:

$$t(l) = \frac{1}{\lambda} \ln \frac{\text{th} \left[ (1/2) \ln \left( A(l) + \sqrt{A^2(l) + 1} \right) \right]}{\text{th}(\lambda u^0/2)}, \quad 0 \leq l < +\infty. \quad (19)$$

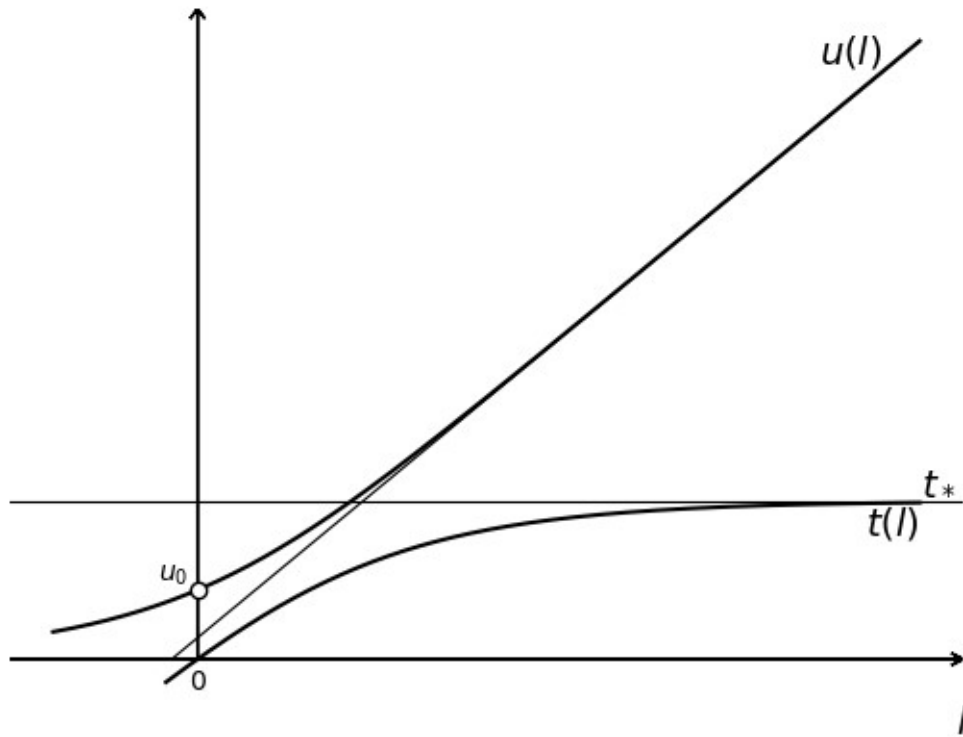


Рис. 2. Точное решение (18-19);  $\circ$  – начальная точка для  $u$ .

Качественный вид функции  $u(l)$  изображен на рис. 2. Интересно сопоставить эту зависимость с рис. 1. На рис. 1 решение резко возрастает при  $t \rightarrow t_*$ . Тем самым, численный расчет на этом участке возможен только очень малым шагом  $\tau$ . Однако на рис. 2 этому участку соответствует правая часть графика, где наклон  $du/dl \approx 1$ . Тем самым, при аргументе  $l$  эту область решения можно проходить большим шагом  $h$ .

Приведем еще некоторые явные формулы, из которых особенно важно выражение кривизны:

$$t(u) = \frac{1}{\lambda} \ln \frac{\text{th}(\lambda u/2)}{\text{th}(\lambda u^0/2)}, \quad (20)$$

$$l(u) = \frac{1}{\lambda} \ln \frac{\text{sh}(\lambda u)}{\text{sh}(\lambda u^0)}, \quad (21)$$

$$\kappa(u) = \frac{u_{tt}}{(1+u_t^2)^{3/2}} = \frac{\lambda \text{sh}(\lambda u)}{\text{ch}^2(\lambda u)}. \quad (22)$$

Кривизна немонотонно меняется вдоль интегральной кривой на рис. 1 или рис. 2. В левой части обоих графиков она возрастает, затем достигает максимума при условиях

$$\max \kappa = \frac{\lambda}{2}, \quad u_{extr} = \frac{1}{2\lambda} \ln(3 + 2\sqrt{2}), \quad t_{extr} = \frac{1}{2\lambda} \ln\left(\frac{1}{3 + 2\sqrt{2}} \frac{\text{ch}(\lambda u^0) + 1}{\text{ch}(\lambda u^0) - 1}\right) \quad (23)$$

и затем монотонно убывает до 0 в правой части графиков. Хорошо видна связь  $\lambda$  с максимальной кривизной интегральной кривой и жесткостью задачи.

Данный тест является уникальным. Во-первых, в нем все значимые величины выражаются друг через друга с помощью несложных комбинаций элементарных функций. Во-вторых, выбирая большие значения  $\lambda$ , можно сделать задачу сколь угодно жесткой.

**Трудность задачи.** Трудность определяется тем числом шагов  $N$ , которое необходимо сделать при оптимальном выборе шага. Из формулы оптимального шага (8) можно предположить, что это число шагов пропорционально входящему в нее интегралу

$$I = \int_0^L \kappa^{2/5}(l') dl'. \quad (24)$$

Величина этого интеграла зависит не только от кривизны  $\kappa(l)$ , но и от отрезка интегрирования  $[0, L]$ . Такой критерий определения трудности можно применять для любых тестовых или прикладных задач.

Для гиперболического теста зависимость  $\kappa(l)$  можно получить, подставляя формулу (18) в формулу (22):

$$\kappa(l) = \frac{\lambda A(l)}{1 + A^2(l)} = \frac{\lambda e^{\lambda l} \text{sh}(\lambda u^0)}{1 + e^{2\lambda l} \text{sh}^2(\lambda u_0)}. \quad (25)$$

Подстановка величины (25) в интеграл (24) и замена переменных  $\zeta = (e^{\lambda l} \text{sh}(\lambda u^0))^{2/5}$  дают

$$I = \int_0^L \frac{\lambda^{2/5} e^{(2/5)\lambda l} \text{sh}^{2/5}(\lambda u^0)}{(1 + e^{2\lambda l} \text{sh}^2(\lambda u_0))^{2/5}} dl = \frac{5}{2} \lambda^{-3/5} \int_{\zeta_0}^{\zeta_L} \frac{d\zeta}{(1 + \zeta^5)^{2/5}}. \quad (26)$$

Пределы интегрирования  $0 < \zeta_0 < \zeta_L < +\infty$ , поэтому последний интеграл в (26) не превышает 2. С учетом множителя  $\lambda^{-3/5}$  перед интегралом неясно, можно ли назвать данный тест по-настоящему трудным. Однако в приведенных далее численных расчетах многие схемы давали срывы при больших  $\lambda$ . Кроме того, более трудного теста с такими уникальными выражениями через элементарные функции найти не удалось.

**2.3. Тригонометрический тест.** Упомянем еще один тест, в котором явно вычисляется длина дуги, и все имеющиеся величины выражаются через элементарные функции как аргумента  $t$ , так и аргумента  $l$ .

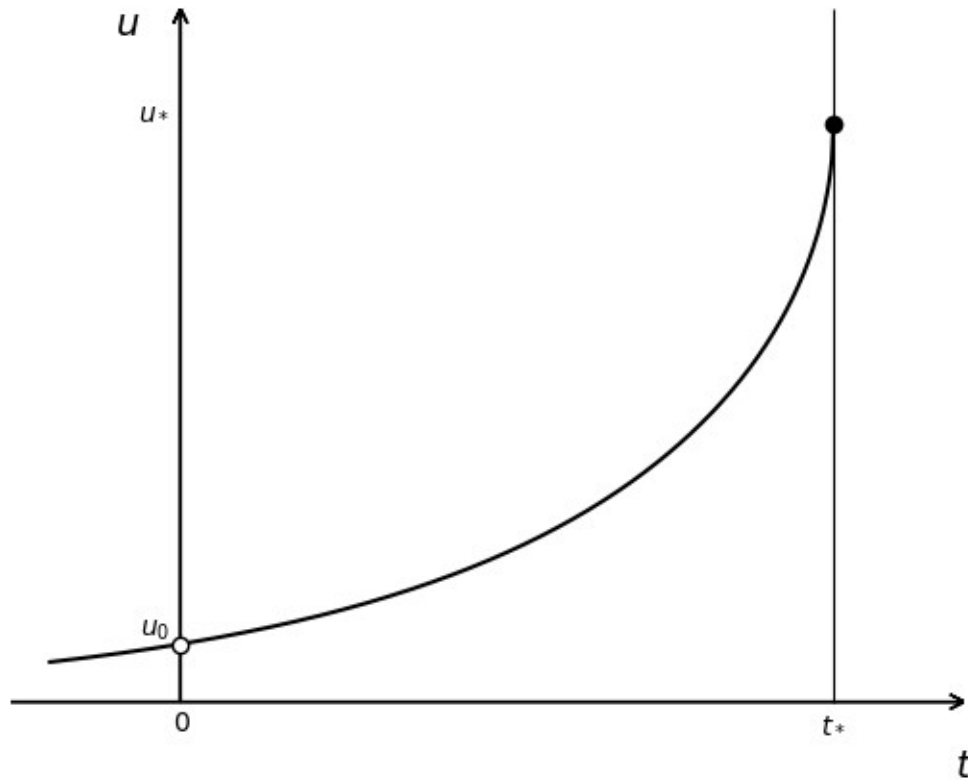


Рис. 3. Точное решение (28); ○ – начальное значение  $u$ , ● – конечное значение  $u$ .

Приведем набор формул этого теста (текстовые пояснения опускаем, поскольку они аналогичны гиперболическому тесту):

$$\frac{du}{dt} = \operatorname{tg}(\lambda u), \quad \lambda > 0, \quad u(0) = u^0 > 0; \quad (27)$$

$$u(t) = \frac{1}{\lambda} \operatorname{arcsin} \left[ e^{\lambda t} \sin(\lambda u^0) \right]; \quad (28)$$

$$l(t) = \frac{1}{\lambda} \ln \frac{\operatorname{tg} \left[ (1/2) \operatorname{arcsin} (e^{\lambda t}) \sin(\lambda u^0) \right]}{\operatorname{tg}(\lambda u^0/2)}; \quad (29)$$

$$\frac{du}{dl} = \sin \lambda u, \quad \frac{dt}{dl} = \cos \lambda u; \quad (30)$$

$$u(l) = \frac{1}{\lambda} 2 \operatorname{arctg} \left[ e^{\lambda l} \operatorname{tg}(\lambda u^0/2) \right]; \quad (31)$$

$$t(l) = \frac{1}{\lambda} \ln \frac{\sin \left[ 2 \operatorname{arctg} \left( e^{\lambda l} \operatorname{tg}(\lambda u^0/2) \right) \right]}{\sin(\lambda u^0)}; \quad (32)$$

$$t(u) = \frac{1}{\lambda} \ln \frac{\sin \lambda u}{\sin \lambda u^0}; \quad (33)$$

$$l(u) = \frac{1}{\lambda} \ln \frac{\operatorname{tg}(\lambda u/2)}{\operatorname{tg}(\lambda u^0/2)}. \quad (34)$$

Качественный вид решения при аргументе  $t$  изображен на рис. 3, а при аргументе  $l$  – на рис. 4. Поскольку мы положили  $\lambda > 0$ , то решение  $u(t)$  монотонно возрастает на конечном промежутке времени от 0 до  $t_* = -\frac{1}{\lambda} \ln [\sin \lambda u^0]$ . При этом функция достигает значения  $u_* \equiv u(t_*) = \frac{\pi}{2\lambda}$ , а производная в этой точке  $u_t(t_*) = +\infty$ . Решение не имеет продолжения за  $t_*$ .

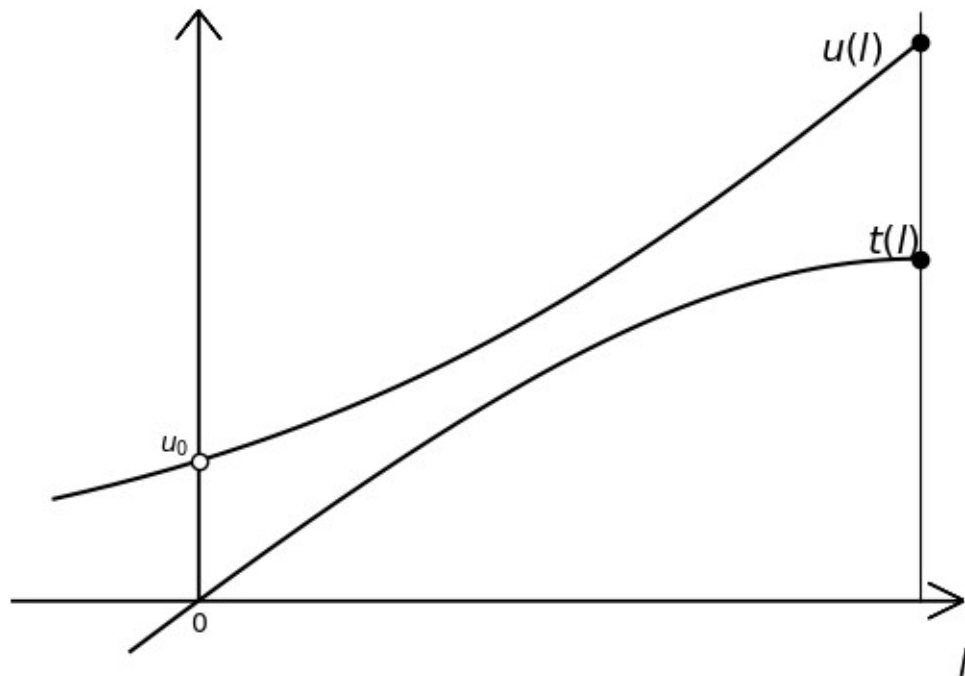


Рис. 4. Точное решение (29-30);  $\circ$  – начальное значение  $u$ ,  $\bullet$  – конечные значения  $u$  и  $t$ .

Если  $u^0 \rightarrow 0$ , то  $t_* \rightarrow +\infty$  и полная длина дуги  $l_* \rightarrow +\infty$ , хотя  $u_*$  остается конечным. Последнее объясняется тем, что при  $u^0 \rightarrow 0$  правая часть  $f(u^0) \rightarrow 0$ , и движение вдоль интегральной кривой начинается очень медленно, так что полное время движения стремится к бесконечности; соответственно стремится к бесконечности длина дуги.

Отметим, что можно рассмотреть случай  $\lambda < 0$ . В этом случае решение  $u(t)$  будет монотонно убывать, асимптотически стремясь к 0 при  $t \rightarrow +\infty$ . Это решение несколько напоминает известный тест Далквиста.

Тригонометрический тест не позволяет получить задачи произвольно высокой жесткости. Поэтому он менее интересен для тестирования методов, чем гиперболический тест.

### 3. Численные расчеты

**3.1. Условия тестирования.** Для численных иллюстраций был выбран гиперболический тест со следующим набором параметров:  $\lambda = 10^\nu$ ,  $\nu = 1, 2, \dots, 10$ ; для каждой используемой схемы  $\nu$  увеличивали до тех пор, пока расчет не срывался. Для каждого  $\lambda$  начальный и конечный моменты расчета определялись так, чтобы в них выполнялось условие  $\kappa = 1$ . Для  $t = 0$ , согласно формуле (22), полагалось

$$u^0 = \frac{1}{\lambda} \operatorname{arcsch} \left( \frac{\lambda - \sqrt{\lambda^2 - 4}}{2} \right). \quad (35)$$

Далее выполнялся расчет по выбранной схеме интегрирования; в нем вычислялись  $u$ ,  $L$  и  $\kappa$ . Кривизна  $\kappa$  сначала возрастала от 1 до  $\sim \lambda/2$ , а затем начинала убывать. Когда выполнялось условие  $\kappa < 1$ , то расчет прекращался.

Опишем расчет  $\kappa_n$ . Для всех использованных схем он выполнялся одинаково по формулам точности  $O(h)$ . В начальный момент времени  $t_0 = 0$ , согласно выбору начального условия,  $\kappa_0 = 1$ ,  $l_0 = 0$ , а вектор правых частей системы (4) для задачи (17) вычисляется по начальным данным. Мы можем вычислить величину первого шага  $h_1$  по кривизне  $\kappa_0$  в формуле (8), провести вычисление по выбранной схеме и найти узел  $l_1$  и величины  $u_1$ ,  $t_1$  в этом узле.

По найденным  $u_1$ ,  $t_1$  вычисляем правые части системы (17) в узле  $l_1$ . Находим длину вектора изменения правых частей и делим ее на величину шага по длине дуги; это дает нам кривизну в первом узле

$$\kappa_1 = \frac{\|\mathbf{F}_1 - \mathbf{F}_0\|_2}{h_1} \equiv \frac{\|\mathbf{F}(t_1, u_1) - \mathbf{F}(t_0, u_0)\|_2}{h_1} \quad (36)$$



с точностью  $O(h)$ . Аналогичная процедура повторяется на всех последующих шагах.

Для расчета необходимо выбрать настроечные параметры формулы (8) для первой сетки. Мы не стали проводить специальные настройки и взяли следующие значения:  $N_{min} = 6$ ,  $N_{max} = 20$ ,  $L = 1$ ,  $\int \dots dl = 1$ . Для перехода с первого этапа сгущения сеток на второй выбиралось значение критерия  $\eta = 0,1$ .

**Надежность программы.** В жестких задачах сравнительно легко происходит срыв расчета, то есть программа оказывается недостаточно надежной. Опишем здесь две типичные ситуации.

1° В жестких задачах при аргументе  $t$  правые части  $f_m$  нередко оказываются сопоставимыми с максимально допустимыми на машине числами. При аргументе  $l$  правые части невелики: сумма их квадратов есть 1. Однако при вычислении знаменателей этих правых частей суммируются квадраты величин  $f_m$ . При этом может произойти переполнение в промежуточных вычислениях. Это приводит к срыву расчетов.

Поэтому надо тщательно продумывать написание правых частей ОДУ, и при появлении слишком больших чисел автоматически производить явное сокращение максимальных множителей в числителях и знаменателях формул. По существу, это означает введение автоматического масштабирования. Этот прием полезен и для не особенно жестких задач.

2° Задачи малой жесткости все схемы считают успешно. Однако при достаточном увеличении жесткости каждая схема может сорваться. В наших расчетах схема ERK4 успешно работала при  $\lambda = 10^4$ , но срывалась при  $\lambda = 10^5$ . Схемы ERK2 и ERK1 работали при  $\lambda = 10^5$ , но срывались при  $\lambda = 10^6$ .

Анализ показал, что на первом этапе шаг  $h_1$  оказывается огромным, и дает сетку только с одним интервалом:  $N=1$ . При этом неправильно работают формулы дробления сетки на втором этапе. Удалось ликвидировать срывы, следующим образом исправив формулы сгущения сетки второго этапа (10)-(12).

Если  $N=1$ , то шаг  $h_1$  делится точно пополам. Если  $N=2$ , то производится расчет только для граничных интервалов (11)-(12). Если  $N \geq 3$ , то включается вычисление и для внутренних интервалов (10).

**Выбор схемы.** Формула выбора шага (7) геометрически согласовывалась лишь со схемами первого порядка точности. Однако мы провели расчет с несколькими схемами разного порядка точности. Наиболее просты и мало трудоемки явные схемы. Поэтому мы взяли три явные схемы Рунге-Кутты ERK1, ERK2, ERK4 точностей  $O(h)$ ,  $O(h^2)$  и  $O(h^4)$  соответственно. Обычно надежность схем уменьшается с повышением их порядка точности, причем на жестких задачах это сказывается особенно сильно. Поэтому такое тестирование различных схем интересно.

**3.2. Результаты расчетов.** Все расчеты проводились с аргументом  $l$ . На каждой сетке локальная погрешность в каждой точке сетки  $l_n$  находилась путем вычитания численного решения  $u_n$  и  $t_n$  из точного  $u(l)$  и  $t(l)$ . Затем вводилась интегральная погрешность

$$\Delta = \sqrt{\sum_{n=1}^N \frac{(u_n - u(l_n))^2 + (t_n - t(l_n))^2}{u^2(l_n) + t^2(l_n)} h_n} / \sum_{n=1}^N h_n. \quad (37)$$

Такое определение погрешности удобно следующим: а) оно учитывает не абсолютные погрешности, а относительные, то есть адаптируется к масштабам конкретного решения; б) одновременно учитывает ошибки обеих функций системы (17); в) автоматически настраивается на длину дуги; г) является интегральной, что облегчает визуализацию погрешности и проведение сравнений. На дальнейших рисунках будет показана зависимость нормы погрешности от числа шагов сетки  $N$ . Рисунки выполнены в двойном логарифмическом масштабе. Тогда прямая линия означает степенной характер зависимости, а тангенс угла ее наклона есть эффективный порядок точности. Это легко позволяет проверить, соответствует ли скорость сходимости теоретическому порядку точности схемы.

**Точность.** На рис. 5 показаны погрешности для схем разного порядка точности ERK1, ERK2, ERK4 при одинаковой жесткости  $\lambda = 10^4$ . Эта жесткость не слишком велика, и все схемы ведут расчет в обычном режиме. Видно, что начала всех линий близки и соответствуют очень большой погрешности  $\sim 10^0$ . Начала всех линий искривлены, но уже в пределах первого этапа линии переходят в прямые, а их наклоны соответствуют теоретическим порядкам точности.

Поэтому при увеличении числа узлов погрешность убывает тем быстрее, чем выше порядок точности схемы. Данные расчеты проводились до  $N \sim 10^4$  узлов. При таком числе узлов погрешность для схемы ERK1 убывает до  $\sim 10^{-3}$ , а для схемы ERK2 – до  $\sim 10^{-6}$ . Для схемы ERK4 погрешность успевает выйти на ошибки округления, которые для данного  $\lambda$  составляют  $\sim 10^{-10}$ . Такие большие ошибки округления (при вычислениях теряется 6 знаков) связаны с достаточно большой жесткостью задачи.

Линии вторых этапов являются прямыми и практически точно продолжают прямые первых этапов; не наблюдается никакого скачка или излома между первым и вторым этапами. Последнее обстоятельство является неожиданным. Во всех ранее проводимых работах между первым и вторым этапом наблюдались изломы или скачки. По-видимому, это связано с тем, что данная работа является первым расчетом, в котором использован тест с явными выражениями всех величин через длину дуги.

Заметим, что для всех схем на первом этапе число узлов при каждом сгущении увеличивается не в 2 раза, особенно на первом сгущении; но на последующих сгущениях первого этапа отношение чисел интервалов  $N$  приближается к 2. Это свидетельствует о том, что заключительная сетка первого этапа близка к квазиравномерной.

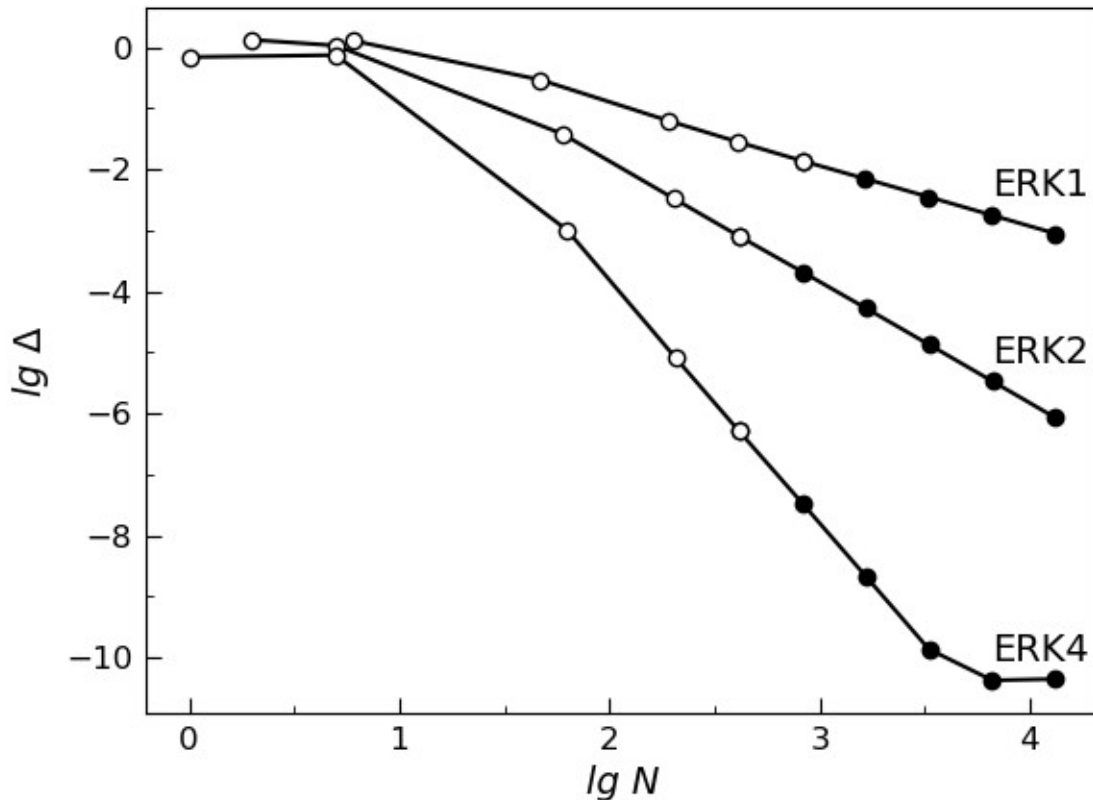


Рис. 5. Погрешности схем ERK для  $\lambda = 10^4$ . Около каждой кривой указан порядок точности схемы. ○ - первый этап, ● - второй этап.

Разумеется, на втором этапе числа шагов  $N$  увеличиваются вдвое при каждом сгущении. Точное удвоение начинается с последней сетки первого этапа.

Отметим важное обстоятельство. Формула (7) выбора оптимального шага согласована лишь со схемами первого порядка точности, и само  $k$  мы вычисляем по формуле точности  $O(h)$ . Однако это не препятствует тому, чтобы схемы ERK2 и ERK4 реализовали свой, более высокий порядок точности. Это свидетельствует об удачном выборе формулы оптимального шага.

**Надежность.** На каждом из рис. 6-8 даны графики погрешности только одной из схем ERK, зато для разных жесткостей. Значения жесткости указаны около каждой кривой. Сравним эти рисунки.

1° На рис. 6 приведены расчеты для схемы ERK1. Схема оказалась исключительно надежной, и даже при  $\lambda = 10^8$  не было обнаружено никаких срывов расчета; лишь при  $\lambda = 10^{10}$  расчет сорвался. При небольших  $\lambda < \sim 100$

линии погрешности с самого начала являются прямыми с наклоном -1, соответствующим теоретическому порядку точности. При большей жесткости начала кривых искривлены, и могут стать даже немонотонными. Однако уже в пределах первого этапа линии становятся прямыми с наклоном -1, а затем гладко переходят в прямые второго этапа.

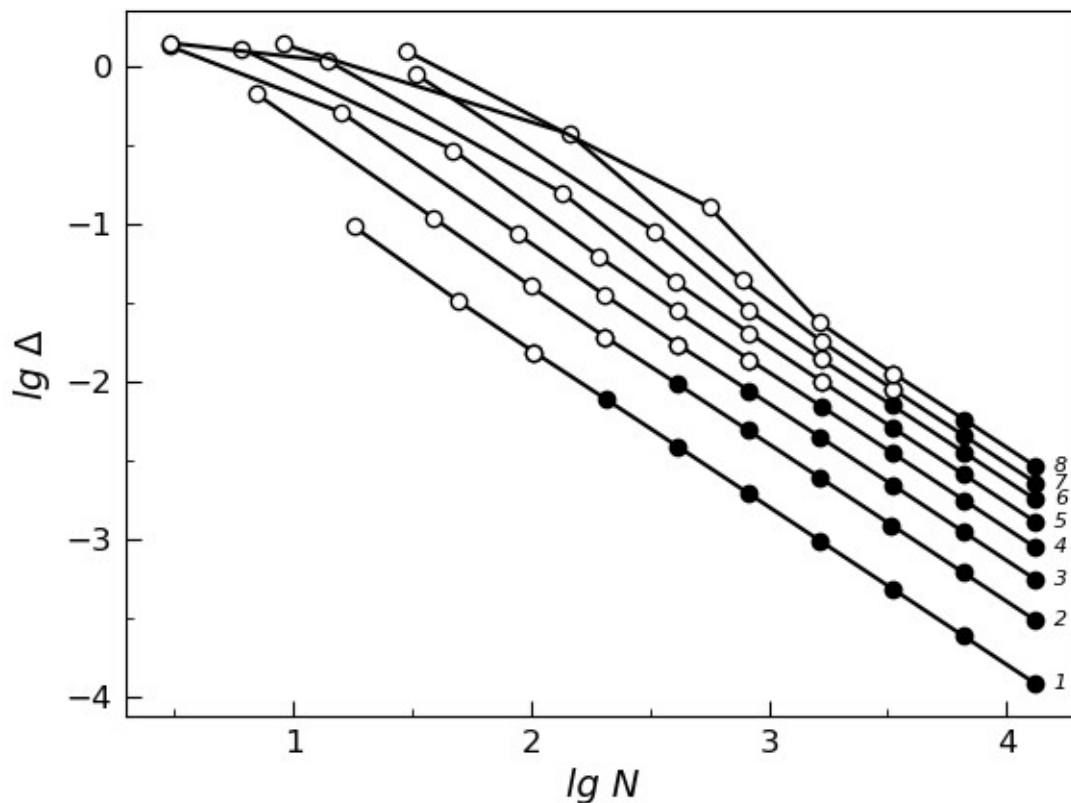


Рис. 6. Погрешности схемы ERK1 для различных жесткостей  $\lambda = 10^v$ ; около линий указаны значения  $v$ .  $\circ$  - первый этап,  $\bullet$  - второй этап.

При повышении жесткости от  $\lambda = 10$  до  $\lambda = 10^8$  происходит ухудшение точности расчета при одинаковом числе узлов. Однако наблюдаемую закономерность такого ухудшения мы не будем анализировать, так как она может относиться только к выбранному тесту.

Заметим, что при всех жесткостях переход с первого на второй этап происходит при почти одинаковом уровне погрешности  $\sim 0.02$  (но разумеется, при разных числах узлов). Остается открытым вопрос - является это обстоятельство случайным или закономерным? Ведь переход на второй этап производится не по уровню погрешности, а по критерию установления квазиравномерности сетки.

2° Расчеты по схеме ERK2 дали срыв при  $\lambda = 10^8$ , то есть несколько ранее, чем для схемы ERK1. Поэтому по надежности схема ERK2 уступает схеме

ERK1, как и следовало ожидать. Кривые погрешностей для жесткостей  $\lambda \leq 10^7$  приведены на рис. 7.

Поведение кривых качественно соответствует тому, что наблюдалось для схемы ERK1. Линии вторых этапов являются прямыми с наклоном -2, что соответствует теоретическому порядку точности. Конечные участки первых этапов так же являются прямыми с наклоном -2, и они гладко сопрягаются со вторыми этапами. Разумеется, поведение начальных участков первых этапов нерегулярно.

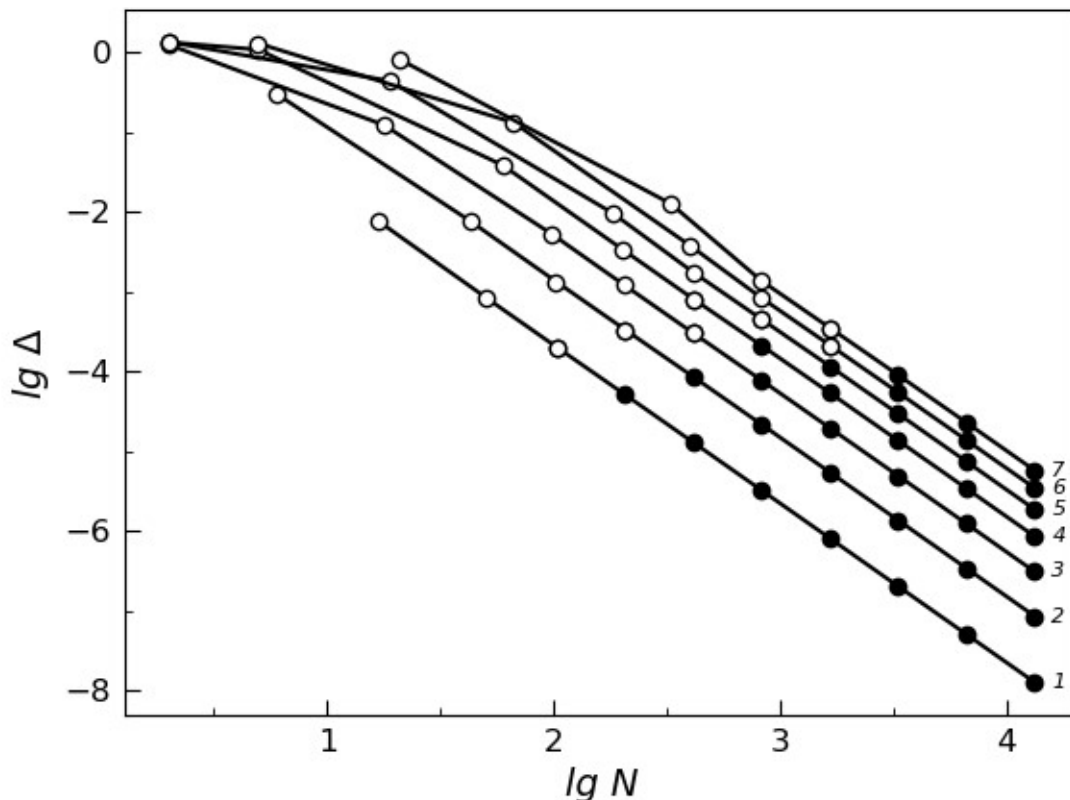


Рис. 7. Погрешности схемы ERK2; обозначения идентичны рис. 6.

Здесь так же переход с первого на второй этапы происходит при почти одинаковом уровне погрешности  $\sim 0.0003$ . Этот уровень погрешности существенно меньше, чем для схемы ERK1.

3° Расчеты по схеме ERK4 срывались при  $\lambda = 10^6$ . Таким образом, схема оказалась менее надежной, чем ERK2, не говоря уже о ERK1. Кривые погрешности для  $\lambda \leq 10^5$  приведены на рис. 8.

На рис. 8 для  $\lambda \leq 10^4$  все линии вторых этапов являются прямыми с наклоном -4, что соответствует теории. Концы линий первых этапов так же являются прямыми с наклоном -4, и гладко сопрягаются с линиями вторых этапов. Видно, что при этих жесткостях схема ERK4 работает очень надежно.

Однако при  $\lambda = 10^5$  начало линии является сильно искривленным, и лишь при  $N \sim 100$  линия становится прямой с теоретическим наклоном 4. При этом первый этап расчета фактически отсутствует ( $N=1$ ), поэтому сетка второго этапа оказывается не адаптированной, а просто равномерной. В результате при  $\lambda = 10^5$  не происходит построения оптимальной сетки. Можно сказать, что при такой жесткости схема ERK4 работает на грани срыва.

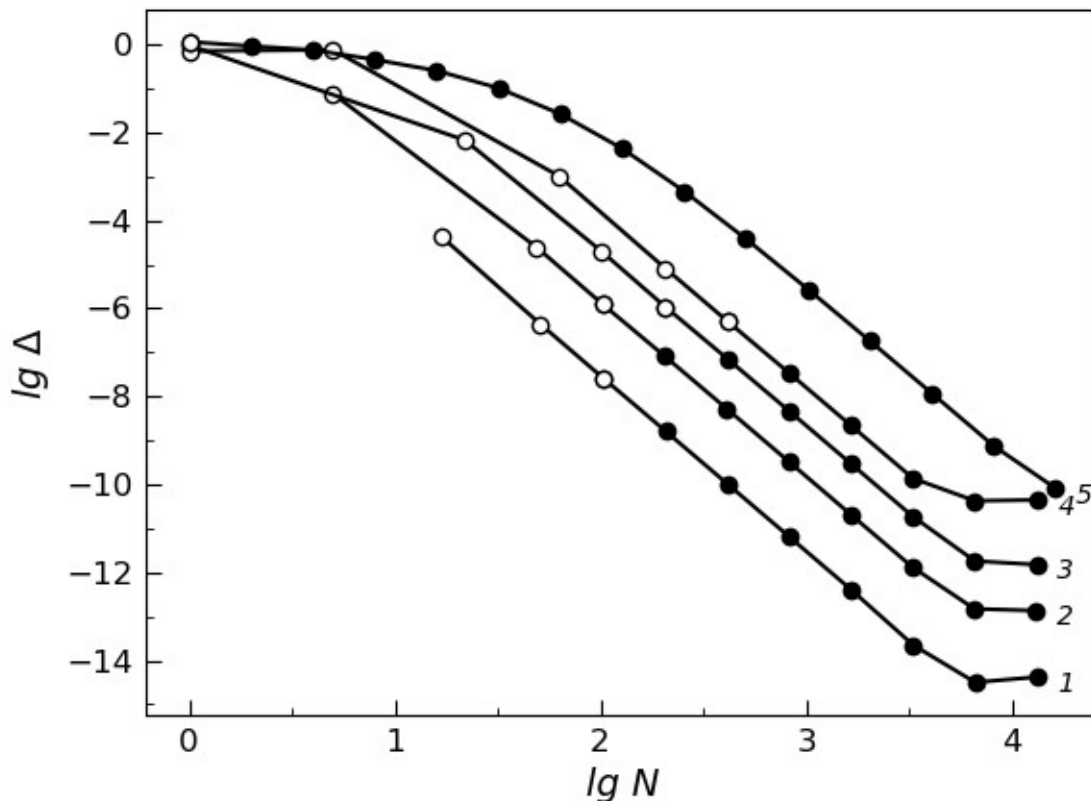


Рис. 8. Погрешности схемы ERK4; обозначения идентичны рис. 6.

Пока схема ведет себя надежно, расстояния между последовательными кривыми 1-2, 2-3, 3-4 уменьшаются. Однако расстояние между прямолинейными участками 4-5 сильно увеличивается. Это показывает, что если предел надежности схемы нарушен, то точность расчета сразу ухудшается.

В следствие высокого порядка точности схемы ERK4, погрешности быстро убывают при увеличении числа узлов. Поэтому каждая линия успевает выйти на ошибки округления. При всех жесткостях это происходит при примерно одинаковом числе узлов  $N \sim 10^4$ . Но уровень ошибок округления существенно возрастает при увеличении  $\lambda$ : это  $\sim 10^{-14}$  при  $\lambda=10$ ,  $\sim 10^{-13}$  при  $\lambda=100$ ,  $\sim 10^{-12}$  при  $\lambda=10^3$  и  $\sim 10^{-10}$  при  $\lambda=10^4$  и  $10^5$ .

**Замечание.** Срывы всех схем зависят не только от величины  $\lambda$ , но и от параметров для выбора начальной сетки (8):  $N_{min}$ ,  $N_{max}$ ,  $L$  и интеграла. Мы намеренно выбрали небольшие значения этих параметров, чтобы тестировать

схемы в трудных условиях. Если увеличивать  $N_{min}$  и  $N_{max}$ , а также найти более удачные формулы для выбора  $L$  и интеграла, можно построить алгоритмы более высокой надежности с использованием тех же схем.

**3.3. Смешанная стратегия.** Рассмотрим несложный способ повышения надежности алгоритма. Схема ERK1 наиболее надежна среди явных схем. Проведем расчеты первого этапа по этой схеме. Результатом расчета даже в случае очень высокой жесткости будет построение хорошей начальной сетки для второго этапа.

Поскольку на первом этапе построена хорошая сетка, второй этап можно выполнять по менее надежной, но гораздо более точной схеме ERK4. При этом следует ожидать гораздо более высокой точности, чем при использовании схемы ERK1 на обоих этапах. Такая стратегия может дать хорошие результаты в тех случаях, когда расчет первого этапа по схеме ERK4 срывается из-за большой жесткости.

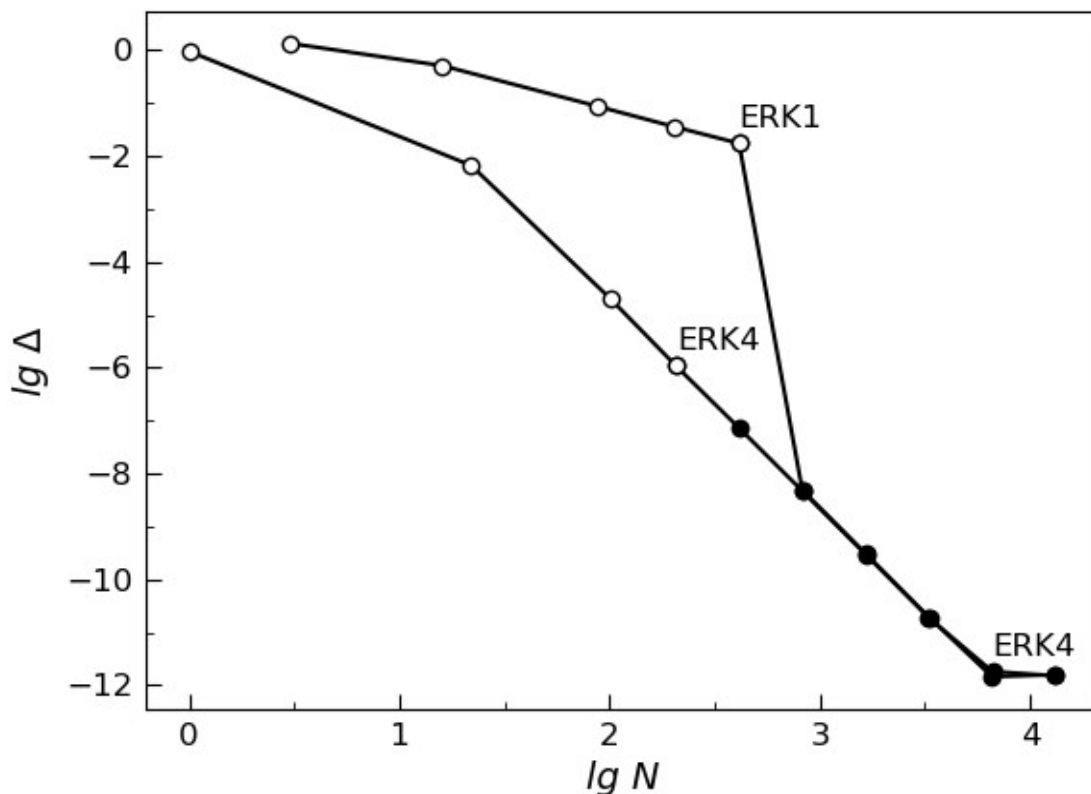


Рис. 9. Смешанная стратегия. Погрешности для  $\lambda = 10^3$ . Около линий подписаны схемы.  $\circ$  - первый этап,  $\bullet$  - второй этап.

Были проведены расчеты по сравнению смешанной и «чистой» стратегий. На рис. 9 показаны результаты при умеренной жесткости  $\lambda = 10^3$ . На нем видны две прямые линии. Верхняя прямая – расчет обоих этапов по схеме ERK1; ее наклон есть -1. Нижняя линия – расчет обоих этапов по схеме ERK4; ее наклон

равен -4. Если выполнить первый этап по схеме ERK1, и затем перейти на схему ERK4, то второй этап расчета точно ложится на нижнюю линию.

Поэтому при умеренной жесткости смешанная стратегия дает ту же самую точность, что «чистая» схема высокого порядка точности. Получается только два преимущества. Первое – небольшое уменьшение объема расчетов, так как первый этап выполняется по менее трудоемкой схеме ERK1. Второе – увеличение надежности, так как начало расчетов выполняется по более надежной схеме ERK1.

На рис. 10 показаны аналогичные расчеты для значительной жесткости  $\lambda = 10^5$ . Здесь картина иная. «Чистый» расчет по схеме ERK4 близок к срыву и прямолинейный участок его кривой лежит довольно высоко. Расчет по смешанной стратегии надежно проходит первый этап, и затем резким скачком переходит ко второму этапу. Его линия второго этапа лежит существенно ниже, уже почти сразу приближаясь к ошибкам округления. Поэтому в данном случае смешанная стратегия дает существенное превосходство в точности. Она дает также заметный выигрыш в трудоемкости, так как расчет можно закончить на одну-две сетки раньше.

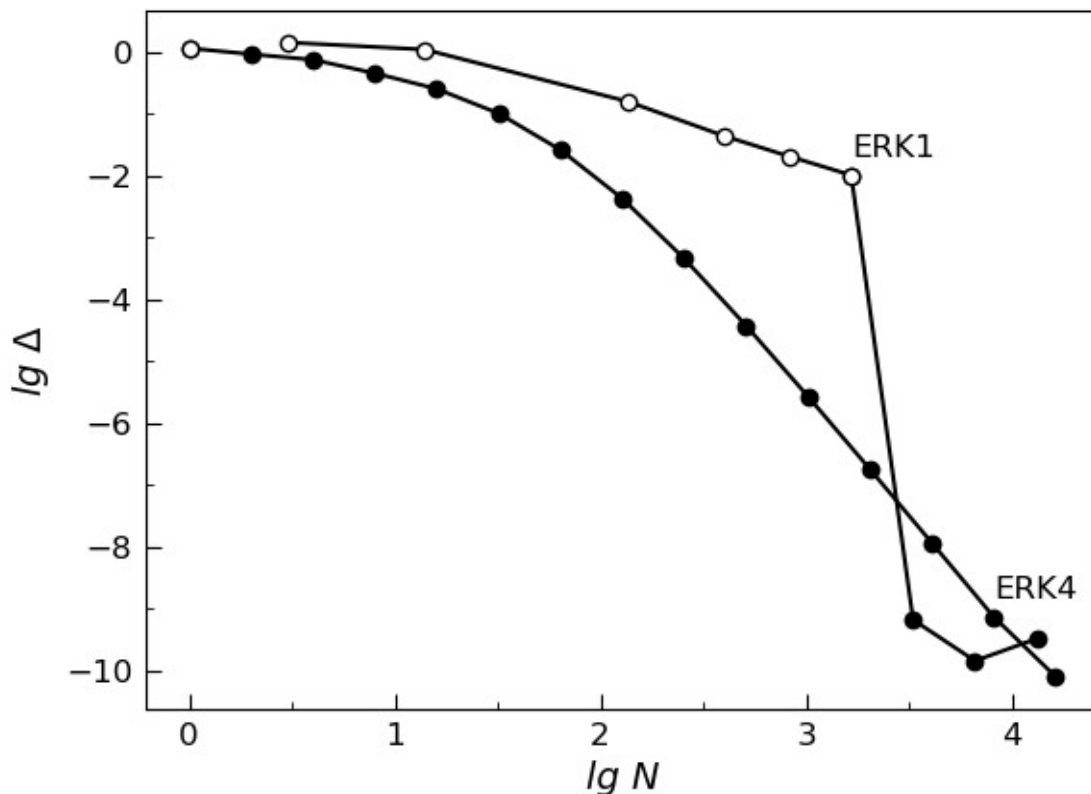


Рис. 10. Смешанная стратегия. Погрешности для  $\lambda = 10^5$ ; обозначения идентичны рис. 9.

На рис. 11 показан расчет для большей жесткости  $\lambda = 10^6$ . Расчет по «чистой» схеме ERK4 здесь срывается на первой же сетке. Однако расчет по



смешанной стратегии благополучно проходит. При этом после первого этапа происходит переход на схему ERK4, и сразу же достигаются ошибки округления. Для такого уровня при «чистом» использовании схем ERK1 или даже ERK2 пришлось бы провести еще много сгущений сеток. Эти примеры наглядно иллюстрируют преимущество смешанной стратегии.

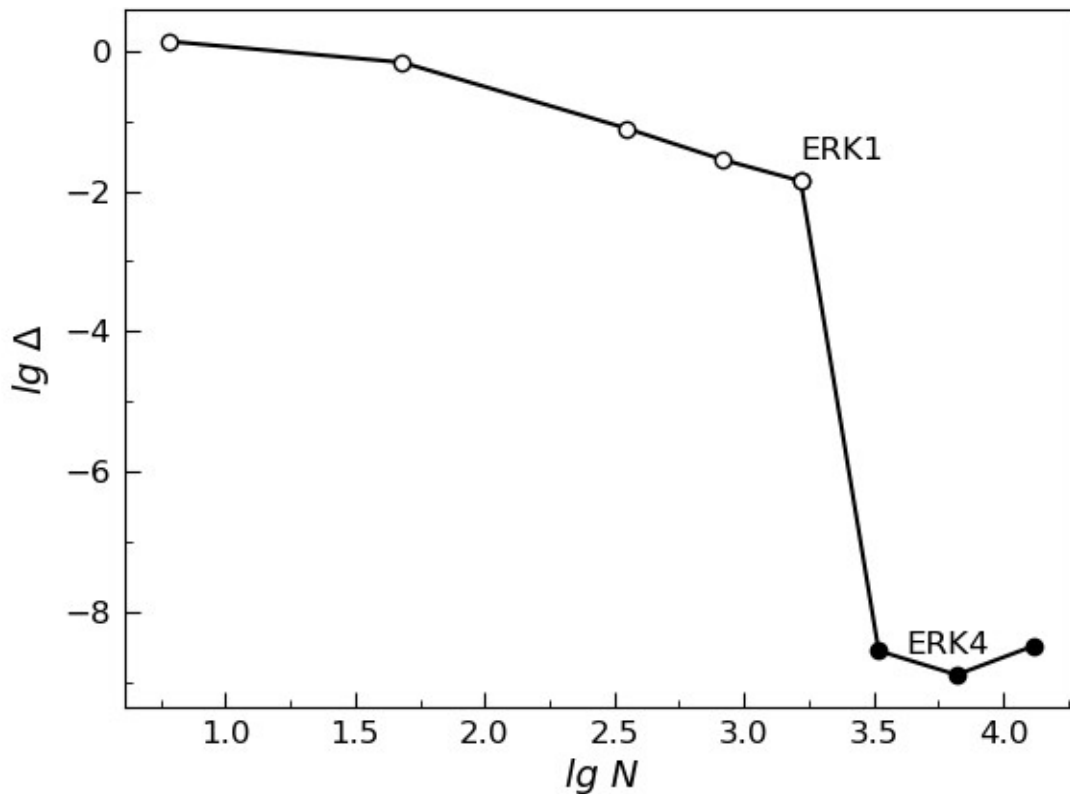


Рис. 11. Смешанная стратегия. Погрешности для  $\lambda = 10^6$ ; обозначения идентичны рис. 9.

#### 4. Заключение

Проведены численные расчеты жестких задач для ОДУ с оптимальным выбором шага по кривизне интегральной кривой. Построен тест, в котором решение выражается в элементарных функциях при использовании в качестве аргумента как времени  $t$ , так и длины дуги  $l$ . Этот тест позволил провести точное и количественное сравнение различных методов решения.

Показано, что при оптимальном выборе шага расчеты задач даже высокой жесткости удается проводить по явным схемам Рунге-Кутты. При этом схема первого порядка точности обладает наибольшей надежностью (расчет срывается только при очень высоких жесткостях), но точность ее мала. Схема второго порядка имеет заметно меньшую надежность, но более высокую точность. Схема четвертого порядка наиболее точна, но наименее надежна.

Предложена смешанная стратегия, при которой первый этап – построение удовлетворительно адаптированной к решению сетки – происходит по схеме

первого порядка; второй этап – сгущение адаптированной сетки – выполняется по схеме четвертого порядка. Это обеспечивает хорошую надежность и высокую точность численного метода.

Работа поддержана грантом РФФИ №18-01-00175.

## Библиографический список

- 1 Хайрер Э., Нерсет С., Ваннер Г. Решение обыкновенных дифференциальных уравнений. Нежесткие задачи. – М.: Мир, 1990.
- 2 Хайрер Э., Ваннер Г. Решение обыкновенных дифференциальных уравнений. Жесткие и дифференциально-алгебраические задачи. М.: Мир, 1999.
- 3 Шалашилин В.И., Кузнецов Е.Б. Метод продолжения решения по параметру и наилучшая параметризация. – М.: Эдиториал УРСС, 1999.
- 4 Белов А.А., Калиткин Н.Н., Пошивайло И.П. Геометрически-адаптивные сетки для жестких задач Коши // Доклады Академии наук. **466:3** (2016), 276–281.
- 5 Белов А.А., Калиткин Н.Н. Выбор шага по кривизне для жестких задач Коши // Математическое моделирование. **28:11** (2016), 97–112.
- 6 Белов А.А., Булатов П.Е., Калиткин Н.Н. Сравнительный анализ алгоритмов автоматического выбора шага для жестких задач Коши // Препринты ИПМ им. М.В. Келдыша. 2019. № 146, 34 с. [https://keldysh.ru/papers/2019/prep2019\\_146.pdf](https://keldysh.ru/papers/2019/prep2019_146.pdf)
76. Пошивайло И.П. Жесткие и плохо обусловленные нелинейные модели и методы их расчета. Диссертация на соискание ученой степени кандидата физико-математических наук: 05.13.18. – Москва, 2015. – 89 с.
- 8 Белов А.А., Калиткин Н.Н. Экономичные методы численного интегрирования задачи Коши для жестких систем ОДУ // Дифференциальные уравнения. **55:7** (2019), 907–918.
- 9 Richardson L.F., Gaunt J.A. The deferred approach to the limit Phil. Trans. A. 1927. Vol. 226. P. 299-349.
- 10 Марчук Г.И., Шайдуров В.В. Повышение точности решений разностных уравнений. М.: Наука, 1979.
- 11 Калиткин Н.Н. Численные методы. М.: Наука, 1978.
- 12 Калиткин Н.Н., Альшин А.Б., Альшина Е.А., Рогов Б.В. Вычисления на квазиравномерных сетках. М.: Физматлит, 2005.
- 13 Жолковский Е.К., Белов А.А., Калиткин Н.Н. Решение жестких задач Коши явными схемами с геометрически-адаптивным выбором шага // Препринты ИПМ им. М.В. Келдыша. 2018. №227, 20 с. [http://keldysh.ru/papers/2018/prep2018\\_227.pdf](http://keldysh.ru/papers/2018/prep2018_227.pdf)

## Оглавление

1. Методы численного интегрирования жестких систем .....	3
1.1. Длина дуги .....	3
1.2. Выбор шага .....	4
1.3. Оптимальный шаг .....	5
1.4. Двухэтапная стратегия.....	6
2. Апробация на тестах .....	9
2.1. Требования к тестам .....	9
2.2. Гиперболический тест .....	10
2.3. Тригонометрический тест .....	14
3. Численные расчеты .....	16
3.1. Условия тестирования.....	16
3.2. Результаты расчетов.....	18
3.3. Смешанная стратегия .....	23
4. Заключение.....	25
Библиографический список.....	26