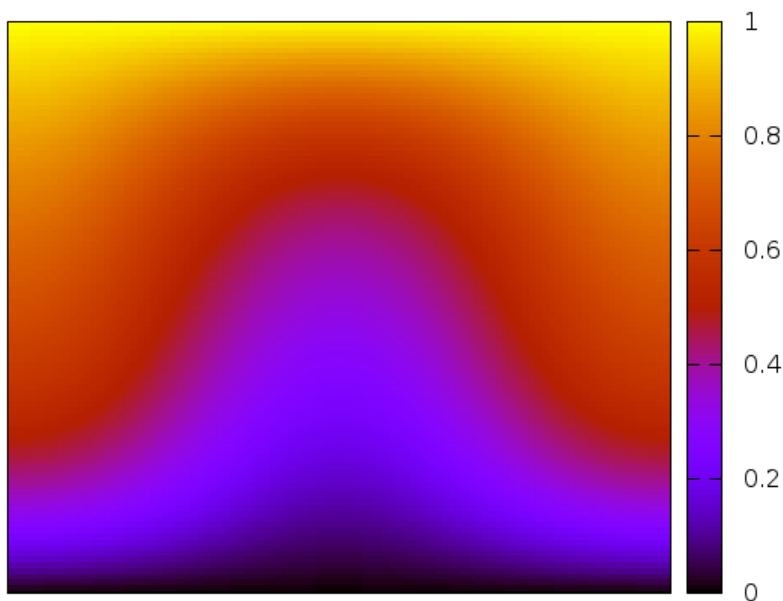

Hausarbeit zur Vorlesung

Numerische Strömungsmechanik



Name: Roland Zimmermann
Matrikelnummer: 21426901
Mail Adresse: roland.zimmermann@stud.uni-goettingen.de
Abgabetermin: 05.08.2016

Inhaltsverzeichnis

1. Einleitung	3
2. Fragestellung	3
3. Theorie	4
3.1. Die Diffusions-Advektion-Gleichung	4
3.2. Das FTCS-Schema	5
3.3. Das BTCS-Schema	6
4. Durchführung der Aufgaben	8
4.1. Aufgabe 1: Untersuchung von \vec{v}	8
4.2. Aufgabe 2: FTCS-Lösung	9
4.3. Aufgabe 3: Numerische Korrektheit der FTCS-Lösung	10
4.4. Aufgabe 4: Auswertung verschiedener FTCS-Lösungen	12
4.5. Aufgabe 5: Beurteilung der Stabilität der FTCS-Lösung	14
4.6. Aufgabe 6: BTCS-Lösung	16
5. Auswertung der Ergebnisse	18
5.1. Numerische Korrektheit der FTCS-Lösung	18
5.2. Aufgabe 5: Beurteilung der Stabilität der FTCS-Lösung	18
5.3. Aufgabe 6: Auswertung und Vergleich der BTCS-Lösung	18
Literatur	19
A. Quellcode	20
A.1. Hauptprogramme	20
A.2. Hilfsbibliotheken	28

1. Einleitung

Viele Verhaltensweisen und Effekte in der Natur können durch partielle Differenzialgleichungen beschrieben werden. Die Lösung dieser ist allerdings nur in wenigen, speziellen Fällen analytisch möglich. Für die restlichen Fälle müssen numerische Lösungsverfahren verwendet werden. Gerade diese Szenarien spielen beispielsweise in der Auto- und Luftfahrtindustrie eine relevante Rolle bei der Produktentwicklung und Optimierung. In dieser Arbeit werden, hierdurch motiviert, zwei verschiedene Lösungsverfahren für ein solches Problem implementiert und anschließend analysiert.

2. Fragestellung

In der Hausaufgabe wird die Wärmeausbreitung (siehe 3.1) in einem zweidimensionalen, rechteckigen System betrachtet. Diese soll durch ein extern erzeugtes Geschwindigkeitsfeld (beispielsweise durch einen Ventilator) verstärkt werden. Das Vektorfeld soll nicht zeitabhängig sein, und lässt sich in dimensionsloser Form durch

$$\vec{v} = (\pi \sin(2\pi x) \cos(\pi y), -2\pi \cos(2\pi x) \sin(\pi y))^t$$

beschreiben. Das zu betrachtende System ist ein Quadrat $[0, 1]^2$. Für dieses gibt es für jeden der vier Ränder je eine Randbedingung. So soll die Temperatur auf dem oberen Rand $T = 1$ und auf dem unteren $T = 0$ betragen. An dem linken und rechten Rand soll dagegen die jeweilige Normalenableitung verschwinden. Dies wird durch eine wärmeisolierende Wand in dem System bedingt. Zu Beginn soll eine Temperaturverteilung der Form $T_{start}(x, y) = y$ vorliegen.

Zur Untersuchung des System fallen mehrere Aufgaben an. So wird zuerst das Geschwindigkeitsfeld selber auf seine Divergenz untersucht (siehe 4.1). Anschließend wird eine erste numerische Lösung des Problems (im *FTCS-Schema*) durchgeführt (siehe 4.2). Daraufhin wird nun die numerische Korrektheit dieser Lösung untersucht (siehe 4.3). Anschließend werden für verschiedene Parameter Simulationen durchgeführt, um unter anderem die Stabilität zu beurteilen (siehe 4.4, 4.5). Schließlich wird das Problem noch einmal gelöst, diesmal aber im *BTCS-Schema* (siehe 4.6).

Abschließend wird noch ein kurzer Vergleich der beiden Verfahren hinsichtlich ihrer Resultate und der benötigten Rechenzeit durchgeführt.

Alle hierfür benötigten Programme und Hilfsroutinen werden in *C++* programmiert.

3. Theorie

3.1. Die Diffusions-Advektion-Gleichung

Die Ausbreitung von Flüssigkeiten, Gasen beziehungsweise von Energie kann zum einen durch die *Diffusionsgleichung* [3, S.49ff]

$$\frac{\partial}{\partial t}T = D\nabla^2T \quad (1)$$

und zum anderen durch die *Advektionsgleichung*

$$\frac{\partial}{\partial t}T = -\nabla \cdot (\vec{v} \cdot T) \quad (2)$$

beschrieben werden. Hierbei stehen T für die Verteilung der Energie (Temperatur), t für die Zeit und D für die Diffusionskonstante. Während die *Diffusionsgleichung* die zeitliche Ausbreitung der Partikel beziehungsweise der Energie, bedingt durch zufällige Bewegungen, beschreibt, beschreibt die *Advektionsgleichung* den Partikelfluss der durch ein Geschwindigkeitsfeld herbeigeführt wird.

Fasst man diese beiden Effekte aus den Gleichungen (2) und (1) zusammen [3, S.140ff], so erhält man die *Diffusions-Advektion-Gleichung*

$$\frac{\partial}{\partial t}T = D\nabla^2T - \nabla \cdot (\vec{v} \cdot T).$$

Hierbei kann nun, wenn ein gegebenes externes, divergenzfreies Geschwindigkeitsfeld gegeben ist, dieses aus der Divergenz herausgelöst werden, sodass sich

$$\frac{\partial}{\partial t}T + \vec{v} \cdot \nabla T = D\nabla^2T \quad (3)$$

ergibt.

Durch das Einführen der *Péclet-Zahl*

$$Pe = \frac{Lv\rho c_p}{\lambda}$$

ist eine Entdimensionalisierung der Gleichung (3) möglich. Hierbei stehen L für eine charakteristische Länge, v für eine Geschwindigkeit, ρ für die Dichte, c_p für die spezifische Wärmekapazität und λ für die Wärmeleitfähigkeit des jeweiligen Mediums. Somit ergibt sich

$$\frac{\partial}{\partial t}T + Pe \cdot \vec{v} \cdot \nabla T = \nabla^2T. \quad (4)$$

3.2. Das FTCS-Schema

Eine mögliche Methode zur Lösung einer solchen partiellen Differenzialgleichung ist durch die Verfahrensgruppe der finiten Differenzen gegeben. Eine hierfür geeignete Diskretisierung des Problems ist in Form eines *FTCS*-Schemas (Forward in Time, Centered in Space) gegeben. Hierbei wird der Raum in der x -Richtung in $(N_x + 1)$ und in der y -Richtung in $(N_y + 1)$ Stützstellen unterteilt, sodass sich ein zweidimensionales räumliches Gitter ergibt. Zwischen diesen Punkten herrscht ein jeweiliger Abstand von $\Delta x = X/N_x$ und $\Delta y = Y/N_y$, wobei X und Y für die Ausmaße des zu beschreibenden Problems in der x - und y -Richtung stehen.

Die Temperatur $T(x, y)$, die hiermit beschrieben werden soll, geht dadurch in die Form $T_{i,j}$ über, wobei i der Index des Punktes in der x - und j der in der y -Richtung ist. Diese Indizes laufen von 0 bis N_x beziehungsweise N_y .

Um nun noch die Zeit zu erfassen, wird auch diese in Punkte diskretisiert, die einen Abstand Δt haben, sodass sich insgesamt ein dreidimensionaler Würfel ergibt, dessen verschiedenen Ebenen für die verschiedenen Zeitpunkte stehen. Erneut geht die Temperatur $T(t, x, y) = T_{i,j}(t)$ über in $T_{i,j}^N$, wobei N für den zeitlichen Index steht. Auch hier wird von 0 an indiziert.

Diese Aufteilung erlaubt es nun, die in den Gleichung (4) vorkommenden Ableitungen zu beschreiben. Hierbei wird in dem *FTCS*-Schema die zeitliche Ableitung in der ersten und die räumlichen in der zweiten Ordnung beschrieben. Die Ordnung gibt hierbei jeweils an, wie stark der Fehler von der jeweiligen Schrittweite abhängt. Bei einem Verfahren zweiter Ordnung [1, S.18ff] hängt dieser Fehler von $\mathcal{O}(\Delta x^2)$ ab. Es ergibt sich für die Ableitungen

$$\frac{\partial}{\partial t} T(x, y) = \frac{T_{i,j}^{N+1} - T_{i,j}^N}{\Delta t}, \quad (5)$$

$$\frac{\partial}{\partial x} T(x, y) = \frac{T_{i+1,j}^N - T_{i-1,j}^N}{2\Delta x}, \quad (6)$$

$$\frac{\partial^2}{\partial x^2} T(x, y) = \frac{T_{i+1,j}^N - 2T_{i,j}^N + T_{i-1,j}^N}{\Delta x^2}, \quad (7)$$

wobei die Fehlerordnung aus Gründen der Übersichtlichkeit implizit weggelassen worden ist. Durch Einsetzen in die Gleichung (4) und Umstellen nach $T_{i,j}^{N+1}$ erhält man schließlich

$$\begin{aligned} T_{i,j}^{N+1} = T_{i,j}^N + \Delta t \left[-Pe \left(v_{i,j}^x \frac{T_{i+1,j}^N - T_{i-1,j}^N}{2\Delta x} + v_{i,j}^y \frac{T_{i,j+1}^N - T_{i,j-1}^N}{2\Delta x} \right) \right. \\ \left. + \frac{T_{i+1,j}^N - 2T_{i,j}^N + T_{i-1,j}^N}{\Delta x^2} + \frac{T_{i,j+1}^N - 2T_{i,j}^N + T_{i,j-1}^N}{\Delta y^2} \right], \end{aligned} \quad (8)$$

womit die Integration durchgeführt werden kann. Die Ausdrücke $v_{i,j}^x$ und $v_{i,j}^y$ stehen für die jeweilige x - und y -Komponente des zeitlich konstanten Geschwindigkeitsfeldes an

dem jeweiligen Ort.

Des Weiteren sind allerdings noch die Randbedingungen zu beachten. Hierbei muss die Differenzialgleichung besonders betrachtet werden, da hier nur Nachbarpunkte in eine Richtung existieren.

Bei den *Dirichlet*-Rändern können die Werte am Rand konstante gehalten werden, sodass

$$T_{Rand}^{N+1} = T_{Rand}^N \Rightarrow T_{i,0}^{N+1} = T_{i,0}^N = 0 \wedge T_{i,N_y}^{N+1} = T_{i,N_y}^N = 1$$

eine geeignete mathematische Beschreibung ist. Exemplarisch für die *Neumann*-Ränder wird der linke Rand betrachtet, auf dem die erste Ableitung in x -Richtung verschwinden soll. Durch eine Taylor-Entwicklung und einen Koeffizientenvergleich erhält man die diskretisierte Randbedingung zweiter Ordnung

$$T_{0,j}^{N+1} - \frac{4}{3}T_{1,j}^{N+1} + \frac{1}{3}T_{2,j}^{N+1} = 0.$$

In diesem Problem muss eine analoge Gleichung auch noch für den rechten Rand betrachtet werden.

3.3. Das BTCS-Schema

Ein alternativer Ansatz zur Lösung ist das *BTCS*-Schema (Backward in Time, Centered in Space), bei der räumlich erneut eine Diskretisierung zweiter Ordnung durchgeführt wird, diese allerdings zum Zeitpunkt $(N + 1)$ ausgewertet wird [3, S.170ff]. Zeitlich wird auch wieder ein Verfahren erster Ordnung verwendet. Somit ergibt sich der Ausdruck

$$T_{i,j}^N = \alpha T_{i,j}^{N+1} + \beta_{i,j}^- T_{i-1,j}^{N+1} + \beta_{i,j}^+ T_{i+1,j}^{N+1} + \gamma_{i,j}^- T_{i,j-1}^{N+1} + \gamma_{i,j}^+ T_{i,j+1}^{N+1}. \quad (9)$$

Dabei wurden die folgenden Abkürzungen eingeführt und benutzt

$$\begin{aligned} \alpha &= \Delta t \left(\frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} + \frac{1}{\Delta t} \right), \\ \beta_{i,j}^\pm &= \Delta t \left(-\frac{1}{\Delta x^2} \pm \frac{Pev_{i,j}^x}{2\Delta x} \right), \\ \gamma_{i,j}^\pm &= \Delta t \left(-\frac{1}{\Delta y^2} \pm \frac{Pev_{i,j}^y}{2\Delta y} \right). \end{aligned}$$

Dies gilt allerdings nur für die inneren Punkte des Gitters. Für den linken und rechten Rand gelten, um ebenfalls die Diskretisierung in der zweiten Ordnung umsetzen zu können, die Gleichungen

$$T_{0,j}^{N+1} - \frac{4}{3}T_{1,j}^{N+1} + \frac{1}{3}T_{2,j}^{N+1} = 0, \quad (10)$$

$$T_{N_x,j}^{N+1} - \frac{4}{3}T_{N_x-1,j}^{N+1} + \frac{1}{3}T_{N_x-2,j}^{N+1} = 0. \quad (11)$$

Um das Problem nun besser beschreiben zu können, wird das zweidimensionale Feld in einem eindimensionalen Spaltenvektor zusammengefasst. Hierfür ist es notwendig, die Indizes umzubenennen. Hierbei wird ab jetzt die folgende Konvention benutzt: durchlaufe erst alle x -Werte zu einem festen y -Wert, und erhöhe den y -Wert. Dies bewirkt, dass die Zeilen des Feldes hintereinander geschrieben werden können:

$$T_{i,j}^N = T_{i+N_y j}^N.$$

Nun kann die Gleichung (9) als eine Matrixoperation der Form

$$\vec{F} = M \cdot \vec{T}^{N+1},$$

mit der Matrix \mathbf{M} ausgedrückt werden. Hierbei geht \vec{F} aus \vec{T}^N hervor, wobei einige Einträge auf null gesetzt werden, um die Gleichungen (10) und (11) korrekt als Matrix-Gleichung darzustellen. Die Matrix \mathbf{M} hat fünf Diagonalen, die ungleich null sind. Sie kann nun als Blockmatrix

$$M = \begin{pmatrix} \mathbb{1} & 0 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ E^- & D & E^+ & 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ 0 & E^- & D & E^+ & 0 & \cdot & \cdot & \cdot & 0 \\ \cdot & & \cdot & \cdot & \cdot & & & & \cdot \\ \cdot & & & \cdot & \cdot & \cdot & & & \cdot \\ \cdot & & & & \cdot & \cdot & \cdot & & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & 0 & E^- & D & E^+ \\ 0 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \mathbb{1} \end{pmatrix},$$

mit den quadratischen Matrizen \mathbf{E}^\pm , \mathbf{D} und der Einheitsmatrix $\mathbb{1}$ geschrieben werden. Die beiden Einheitsmatrizen im oberen linken und unteren rechten Eck setzen hierbei die *Dirichlet*-Randbedingungen am oberen und rechten Rand um. Die Matrizen \mathbf{D} und \mathbf{E}^\pm

$$E^\pm = \begin{pmatrix} 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \gamma^\pm & \cdot & \cdot & \cdot & \cdot \\ \cdot & 0 & \gamma^\pm & 0 & \cdot & \cdot \\ \cdot & \cdot & 0 & \gamma^\pm & 0 & \cdot \\ \cdot & \cdot & \cdot & 0 & \gamma^\pm & \cdot \\ 0 & \cdot & \cdot & \cdot & \cdot & 0 \end{pmatrix},$$

$$D = \begin{pmatrix} 1 & -\frac{4}{3} & \frac{1}{3} & \cdot & \cdot & 0 \\ \beta^- & \alpha & \beta^+ & 0 & \cdot & \cdot \\ 0 & \beta^- & \alpha & \beta^+ & 0 & \cdot \\ \cdot & \cdot & \beta^- & \alpha & \beta^+ & 0 \\ \cdot & \cdot & 0 & \beta^- & \alpha & \beta^+ \\ 0 & \cdot & \cdot & \frac{1}{3} & -\frac{4}{3} & 1 \end{pmatrix}$$

setzen die implizite Gleichungen (9) und (10) um, sind also das Kernstück des *BTCS-Schemas*. Ihre erste und letzte Zeile weichen hierbei vom Rest ab, da sie die *Neumann-Randbedingungen* umsetzen. Um das System nun um einen Zeitschritt weiter zu simulieren, muss die Matrix \mathbf{M} invertiert werden - beziehungsweise das Gleichungssystem $\mathbf{M}\vec{x} = \vec{b}$ gelöst werden.

Hierfür hat sich das *SOR-Verfahren* (Successive Over-Relaxation), welches eine Modifikation des *Gauß-Seidel-Verfahrens* ist, als besonders geeignet herausgestellt. Dies sind beides iterative Lösungsverfahren. Um sie verwenden zu können, wird die Matrix zuerst in drei Teilmatrizen $\mathbf{M} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ zerlegt, wobei \mathbf{D} eine Diagonal-, \mathbf{L} eine untere Dreiecks- und \mathbf{U} eine obere Dreiecksmatrix ist. Zudem wird der *Residuen-Vektor* $\vec{r}_n = \mathbf{M}\vec{x}_n - \vec{b}$ eingeführt, wobei der Index den jeweiligen Iterationsschritt angibt. Dieser Vektor beschreibt die Differenz zwischen der korrekten und der iterativen Lösung der Gleichung $\mathbf{M}\vec{x} = \vec{b}$.

Durch diese Zerlegung ist es möglich eine rekursive Bestimmungsformel [2, S. 25ff]

$$\vec{x}^n = \vec{x}^{n-1} - \alpha(\mathbf{L} + \mathbf{D})^{-1}\vec{r}^{n-1}$$

aufzustellen. Hierbei ist α der so genannte *Relaxationsparameter*, welcher bei dem *Gauss-Seidel-Verfahren* gleich eins ist.

Diese Matrizen-Gleichung kann in Summenschreibweise durch Eigenschaften triagonaler Matrizen aus der linearen Algebra auch als

$$x_i^{n+1} = (1 - \alpha)x_i^n + \frac{\alpha}{m_{i,i}} \left(b_i - \sum_{j=1}^{i-1} m_{i,j}x_j^{n+1} - \sum_{j=i+1}^N m_{i,j}x_j^n \right) \quad (12)$$

aufgefasst werden [2, S.30], wobei N die Dimensionen des Problems angibt.

4. Durchführung der Aufgaben

4.1. Aufgabe 1: Untersuchung von \vec{v}

Für die numerische Betrachtung des Problems spielt der Einfluss von \vec{v} eine entscheidende Rolle. Ist dieses Geschwindigkeitsfeld nämlich divergenzfrei, so ergibt sich eine mögliche, vereinfachende Betrachtung in diskretisierter Form nach Gleichung (8). Für das externe Feld

$$\vec{v} = (\pi \sin(2\pi x) \cos(\pi y), -2\pi \cos(2\pi x) \sin(\pi y))^t$$

gilt

$$\nabla \cdot \vec{v} = \frac{\partial}{\partial x} (\pi \sin(2\pi x) \cos(\pi y)) + \frac{\partial}{\partial y} (-2\pi \cos(2\pi x) \sin(\pi y)) = 0;$$

es ist somit divergenzfrei. Es bietet sich zudem an, das Feld grafisch darzustellen. Dies kann in Abbildung (1) nachvollzogen werden.

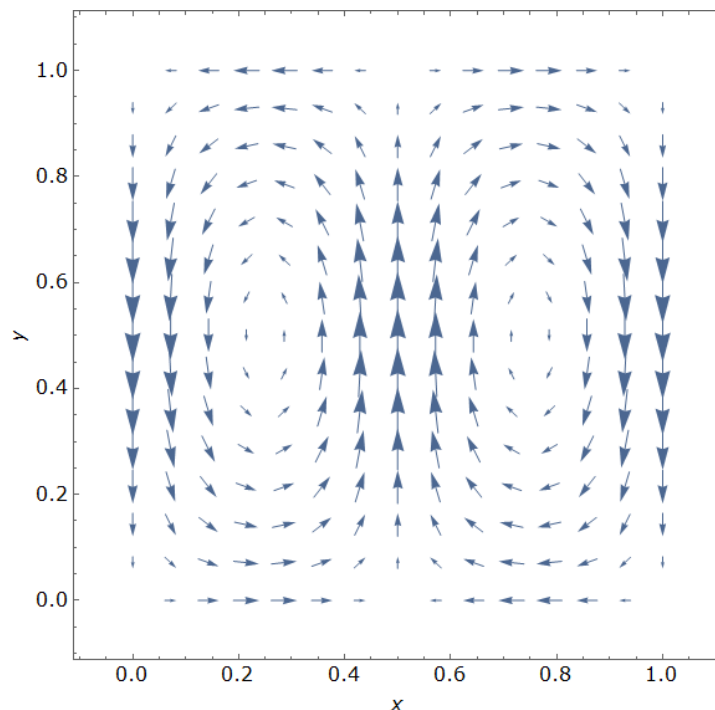


Abbildung 1: Grafische Darstellung des Geschwindigkeitsfeldes \vec{v} in Abhängigkeit von x und y .

Es ist zu erkennen, dass das Feld zwei Zentren hat, um die sich je ein Strudel bildet. Diese laufen in der Mitte zusammen, und strömen zum oberen Rand hin.

4.2. Aufgabe 2: FTCS-Lösung

Um nun eine erste numerische Lösung erzeugen zu können, wird das zuvor in (3.2) hergeleitete Verfahren benutzt. Hierfür wird direkt eine etwas allgemeinere Implementation gewählt, welche in der Integration auch mögliche Quellterme berücksichtigt. Dieser Ansatz wird gewählt, da im weiteren Verlauf dieser Arbeit (siehe 4.3) die numerische Korrektheit des Integrators überprüft werden soll. Hierfür wird eine Quelle hinzugefügt, welche eine analytische Lösung erlaubt.

Das Programm aus Quellcode (1) besteht aus zwei verschiedenen Routinen. Die MAIN-Methode liest gewünschte Eingangsparameter, wie die Schrittweite Δt , die maximale Integrationszeit t_{max} , die *Péclet-Zahl* Pe , die Gitterauflösung N_x und N_y und einen weiteren Parameter, der angibt, ob ein Quellterm erzeugt werden soll, ein (Zeile 78-119). Zudem erzeugt sie Objekte des Types `VECTOR<VECTOR<DOUBLE>>`, in welchem der

Temperaturverteilung in dem Rechteck, und ein etwaiger Quellterm abgespeichert werden. Hierfür wird ein VECTOR statt eines Feldes genutzt, um mögliche ACCESS VIOLATION EXCEPTIONS zu vermeiden und um Memoryleaks zu reduzieren (Zeile 121-163). Etwaige Leistungsverluste hierdurch sind zum einen begrenzt und steuern keinen großen Anteil zu der Gesamtlaufzeit bei.

In einer Schleife wird in den gewählten Zeitschritten bis zur Endzeit iteriert, und für jede Zeit ein Integrationsschritt durchgeführt (Zeile 131-163). Im Anschluss daran wird die benötigte Rechenzeit ausgegeben (Zeile 165-166 und 202-205). Zuletzt wird die berechnete Temperaturverteilung in die vorher angegebene Datei ausgegeben (Zeile 208-215). Der jeweilige Integrationsschritt wird von der FTCS_TIME_STEP-Methode ausgeführt. Diese legt zuerst eine Kopie des übermittelten Feldes an, um auf Grundlage dieser, die Ableitungsterme berechnen zu können. Hierfür wird in einer Schleife über das gesamte Gitter iteriert, und jeweils der $\nabla^2 T$ und der $\vec{v}_0 \cdot \nabla T$ Term einzeln berechnet und anschließend nach Gleichung (13) als Änderung von T verrechnet (Zeile 19-59). Abschließend werden noch die *Dirichlet*- und *Neumann*-Ränder beachtet und die Werte passend aktualisiert (Zeile 61-76).

4.3. Aufgabe 3: Numerische Korrektheit der FTCS-Lösung

Zur Beurteilung der Korrektheit der numerischen Lösung wird normalerweise erstmal versucht, das Ergebnis mit der analytischen Lösung zu vergleichen. Da dieses Problem allerdings nicht analytisch lösbar ist, muss hier anders verfahren werden. Als Ansatz hierbei wird eine mögliche stationäre Lösung geraten, und die Differenzialgleichung um einen Quellterm so ergänzt, dass dies die exakte Lösung des Problems ist.

Für die konstruierte stationäre Lösung wird der Ansatz

$$T^* = \cos(\pi x) \sin(\pi y) + y$$

gewählt, da sie die vier Randbedingungen erfüllt. Um den Quellterm Q zu korrigieren, wird T^* in die stationäre Gleichung

$$Pe \cdot \vec{v} \nabla T^* = \nabla^2 T^* + Q$$

eingesetzt, welche auch Quellen berücksichtigt. Es ergibt sich somit

$$Q = Pe \cdot \vec{v} \nabla T^* - \nabla^2 T^*,$$

was sich nach Einsetzen des Ansatzes und des gegebenen Geschwindigkeitsfeldes zu

$$Q = \pi^2 (2 \cos(\pi x) \sin(\pi y) - Pe \cdot \cos(\pi x)^3 \sin(2\pi y)) - 2Pe \cdot \pi \cos(2\pi x) \sin(\pi y)$$

ergibt.

Berücksichtigt man diesen Quellterm in dem *FTCS-Schema*, so folgt:

$$T_{i,j}^{N+1} = T_{i,j}^N + \Delta t \left[-Pe \left(v_{i,j}^x \frac{T_{i+1,j}^N - T_{i-1,j}^N}{2\Delta x} + v_{i,j}^y \frac{T_{i,j+1}^N - T_{i,j-1}^N}{2\Delta x} \right) + \frac{T_{i+1,j}^N - 2T_{i,j}^N + T_{i-1,j}^N}{\Delta x^2} + \frac{T_{i,j+1}^N - 2T_{i,j}^N + T_{i,j-1}^N}{\Delta y^2} + Q_{i,j} \right], \quad (13)$$

Die Berechnung und Zuweisung dieses Quellterms (Zeile 148-161) wurde im Quellcode (1) implementiert. Nun kann zu jedem Zeitpunkt die Differenz zwischen der stationären Lösung $T_{i,j}^*$ und der numerischen Lösung $T_{i,j}^N$ berechnet werden als

$$\sigma^N = \sqrt{\sum_{i,j} (T_{i,j}^N - T_{i,j}^*)^2}.$$

Diese Abweichung wird nun alle 100 Zeitschritte in der *main*-Methode (Zeile 175-196) ausgegeben.

Setzt man nun noch $N_x = N_y$, und variiert diese Werte für ein festes Δt , so können Aussagen getroffen werden, über die Abweichungen der stationären numerischen T und der tatsächlichen stationären Lösung T^* in Abhängigkeit von der Gittergröße.

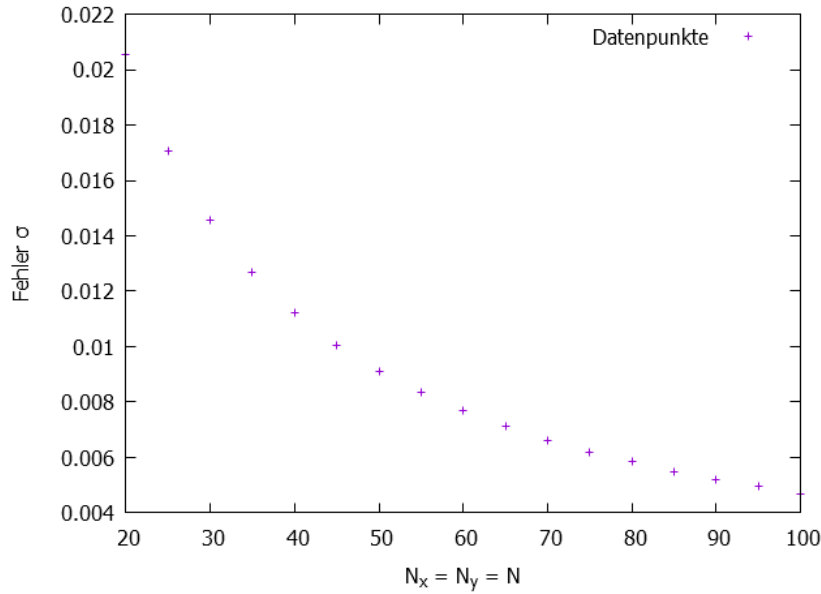


Abbildung 2: Grafische Auftragung der numerischen Abweichung σ in Abhängigkeit von der Gittergröße N .

In Abbildung (2) zeigt den Verlauf des Fehlers σ in Abhängigkeit von $N_x = N_y$ für $Pe = 2$ für einen Zeitschritt von $\Delta t = 2 \cdot 10^{-5}$ aufgetragen. Zudem wurde eine Regression

durchgeführt, welche eine Abhängigkeit der Form $\sigma = N^{-1.2}$ ergeben hat. Dazu ist noch anzumerken, dass für größere N -Werte (größer als 115) der Fehler nicht mehr weiter sinkt, sondern sogar stark anwächst, da die numerische Lösung divergent wird.

4.4. Aufgabe 4: Auswertung verschiedener FTCS-Lösungen

Im weiteren Verlauf wird der zuvor eingeführte Quellterm $Q = 0$ gesetzt. Nun kann das Temperaturfeld T für zwei verschiedene *Péclet-Zahlen* $Pe = 2$ und $Pe = 10$ simuliert, und zu den drei Zeitpunkten $t = 0.005$, $t = 0.05$ und $t = 0.5$ visualisiert werden. Hierfür wird $\Delta t = 2 \cdot 10^{-4}$ und $N = N_x = N_y = 30$ angenommen. Die Darstellungen sind in den Abbildungen (4) und (5) zu finden.

Hierbei ergibt sich unabhängig von der *Péclet-Zahl* eine ähnliche Temperaturverteilung. So ist sie im unteren Bereich deutlich geringer, als im oberen Bereich, was durch die *Dirichlet*-Randbedingungen hervorgerufen wird. Zudem ist sie in der Mitte deutlich geringer, als rechts und links daneben. Dies lässt sich gut durch die Verteilung durch das Geschwindigkeitsfeld erklären. Diese transportiert anschaulich die vorhandene Wärme aus der Mitte hinweg, und verteilt sich im oberen Bereich. Bei der größeren *Péclet-Zahl* $Pe = 10$ ergibt sich eine stärkere Unterteilung. Für beide Parameter geht das System allerdings schon nach der kurzen Zeitspanne $t = 0.05$ in eine annähernd stationäre Lösung über - die Unterschiede zwischen den Verteilung für $t = 0.05$ und $t = 0.5$ sind gering. Zur Veranschaulichung des Zusammenhanges zwischen der Verteilung und dem Geschwindigkeitsfeld wurden diese für eine Parameterwahl in Abbildung (3) dargestellt.

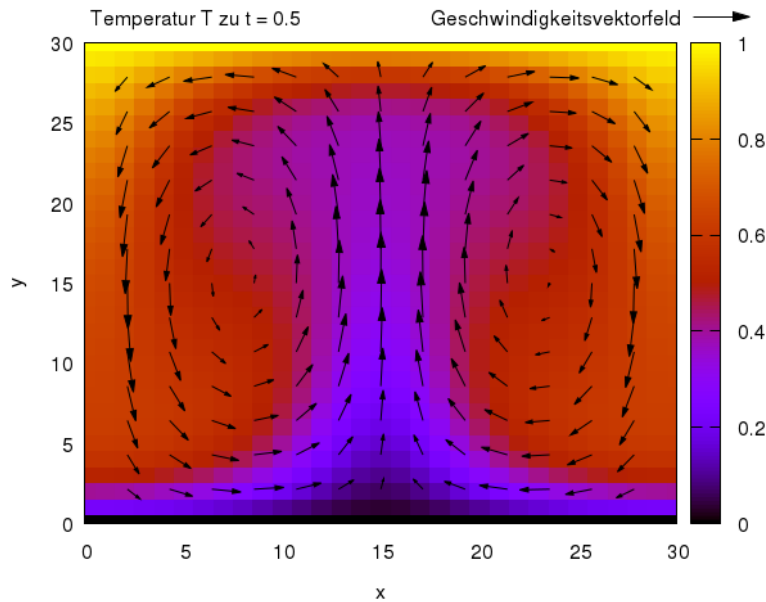


Abbildung 3: Darstellung der Temperaturverteilung und dem Geschwindigkeitsfeld \vec{v} für $Pe = 10$, $t_{max} = 0.05$, $\Delta t = 2 \cdot 10^{-4}$ und $N = N_x = N_y = 30$.

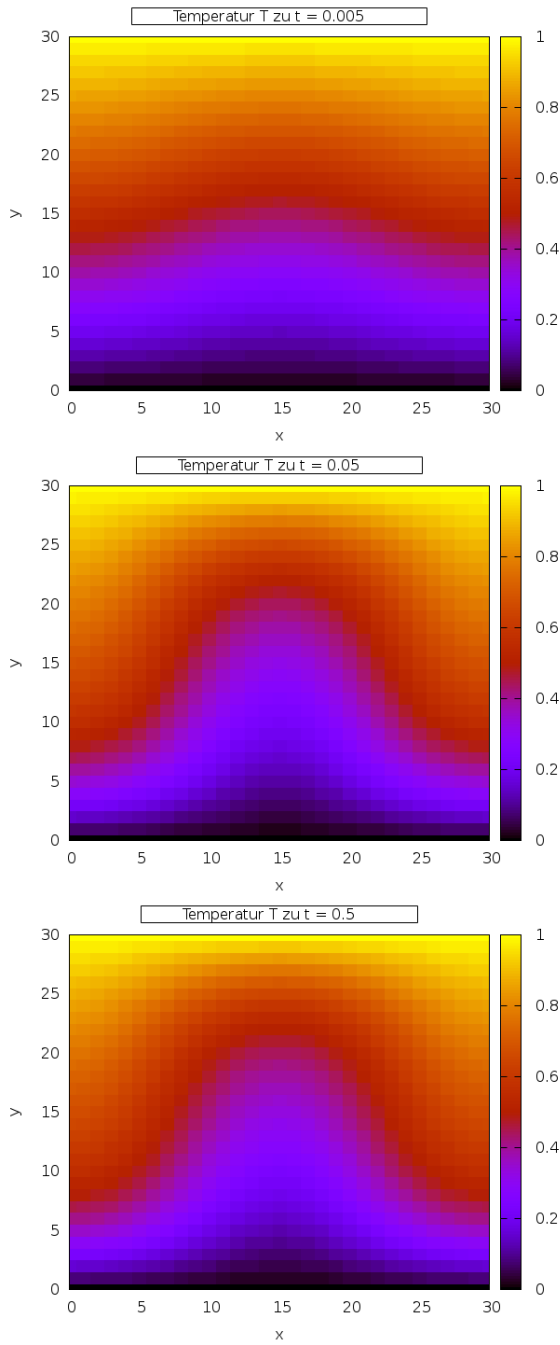


Abbildung 4: Temperatur für $Pe = 2$.

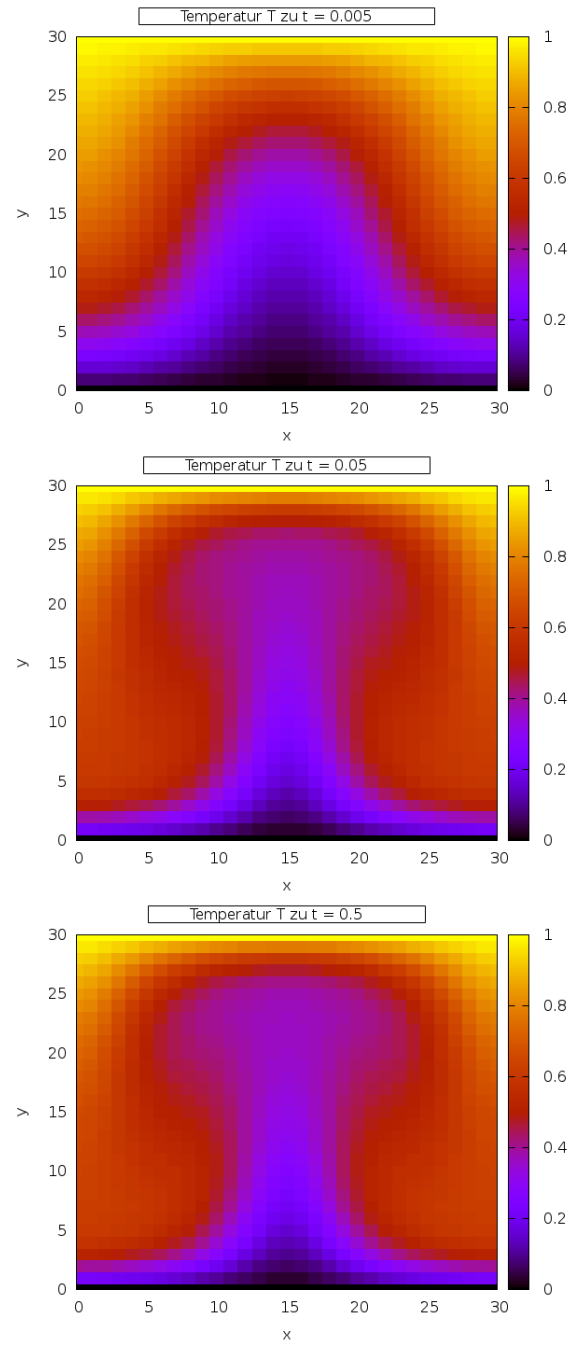


Abbildung 5: Temperatur für $Pe = 10$.

4.5. Aufgabe 5: Beurteilung der Stabilität der FTCS-Lösung

Um weitere Aussagen über die numerische Lösung treffen zu können, wird nun noch die Stabilität betrachtet. Hierfür wird zuerst für $N = N_x = N_y = 30$ die maximale Zeitschrittgröße Δt_{max}^{Pe} für die *Péclet-Zahlen* $Pe \in \{0.1, 1.0, 10\}$ empirisch ermittelt. Diese ergeben sich zu

$$\begin{aligned}\Delta t_{max}^{0.1} &= 2.788 \cdot 10^{-4}, \\ \Delta t_{max}^{1.0} &= 2.789 \cdot 10^{-4}, \\ \Delta t_{max}^{10} &= 2.795 \cdot 10^{-4}.\end{aligned}$$

Dies liefert einen ersten Eindruck: Die Stabilitätsgrenze hängt von Pe ab, und steigt für steigende Pe in gewissen, geringen Maßen in dem hier betrachteten Intervall an. Für eine genauere Aussage wird nun eine *von-Neumann-Stabilitätsanalyse* durchgeführt. Hierfür wird für die Lösung ein Ansatz der generellen Form

$$T_{i,j}^N = \xi^N e^{i(k_x \Delta x \cdot i + k_y \Delta y \cdot j)}$$

gewählt. Durch Einsetzen in Gleichung (8) und Kürzen ergibt sich der Ausdruck

$$\begin{aligned}\xi = 1 - \frac{\Delta t Pe v^x}{2\Delta x} (e^{ik_x \Delta x} - e^{-ik_x \Delta x}) - \frac{\Delta t Pe v^y}{2\Delta y} (e^{ik_y \Delta y} - e^{-ik_y \Delta y}) - \frac{2\Delta t}{\Delta x^2} - \frac{2\Delta t}{\Delta y^2} \\ + \frac{\Delta t}{\Delta x^2} (e^{ik_x \Delta x} + e^{-ik_x \Delta x}) + \frac{\Delta t}{\Delta y^2} (e^{ik_y \Delta y} + e^{-ik_y \Delta y}),\end{aligned}$$

welcher durch trigonometrische Identitäten in

$$\begin{aligned}\xi = 1 - i\Delta t Pe \left(\frac{v^x}{\Delta x} \sin(k_x \Delta x) + \frac{v^y}{\Delta y} \sin(k_y \Delta y) \right) + \frac{2\Delta t}{\Delta x^2} \cos(k_x \Delta x) + \frac{2\Delta t}{\Delta y^2} \cos(k_y \Delta y) \\ - \frac{2\Delta t}{\Delta x^2} - \frac{2\Delta t}{\Delta y^2}\end{aligned}$$

übergeht. Damit das Verfahren konvergiert, muss $|\xi|$ kleiner 1 sein. Nähert man hier den Sinus und den Cosinus mit den Maximalwerten, und schätzt die Komponenten v^x und v^y nach oben durch π und 2π ab, so ergibt sich

$$\lambda^2 (64 + (2\pi Pe \Delta x)^2) - 16\lambda < 0, \quad (14)$$

mit der *CFL-Zahl* $\lambda = \Delta t / \Delta x^2$. Eine Lösung dieser Ungleichung ergibt folgenden Bereich, sodass das *FTCS-Schema* konvergiert:

$$0 < \lambda = \frac{\Delta t}{\Delta x^2} < \frac{4}{16 + (\Delta_x Pe \pi)^2}. \quad (15)$$

Hierdurch kann der maximale Zeitschritt, für den das Verfahren konvergiert, durch

$$\Delta t_{max} = \frac{4\Delta x^2}{16 + (\pi Pe \Delta_x)^2} \quad (16)$$

ausgedrückt werden.

<i>Péclet</i> -Zahl	Δt_{max} (empirisch)	Δt_{max} (theoretisch)
0.1	$2.788 \cdot 10^{-4}$	$2.778 \cdot 10^{-4}$
1	$2.789 \cdot 10^{-4}$	$2.776 \cdot 10^{-4}$
10	$2.795 \cdot 10^{-4}$	$2.600 \cdot 10^{-4}$

Tabelle 1: Vergleich der empirischen und theoretischen abgeschätzten Werte für Δt_{max} .

Ein Vergleich zwischen den empirisch ermittelten Grenzzeitschritten und den maximalen Zeitschritten, für die nach der nach oben abgeschätzten theoretischen Betrachtung noch ein konvergentes Verhalten zu erwarten ist, bietet die Tabelle (1). Es ist zu erkennen, die theoretische Abschätzung für die kleinste *Péclet*-Zahl die beste Übereinstimmung mit dem empirisch gefunden Wert hat - dies hebt den Charakter der bei der Herleitung angenommenen Abschätzungen hervor.

4.6. Aufgabe 6: BTCS-Lösung

Nachdem zuvor die explizite Diskretisierung untersucht wurde, wird das Problem nun implizit im *BTCS-Schema* betrachtet (siehe 3.3). Dazu werden im Hauptprogramm im Quellcode (2) vier Routinen benötigt. Die MAIN-Methode liest die eingegebene Parameter erneut ein (Zeile 50-88), und befüllt die zuvor eingeführte Matrix \mathbf{M} . Dazu wird ein Objekt der Klasse `SQUARE_MATRIX` aus Quellcode (3) benutzt (Zeile 90-155). Diese erbt von der Klasse `MATRIX`, welche im Quellcode (4) und (6) definiert wird. Hierbei werden die Routinen `BETA` und `GAMMA` benutzt, welche die so benannten Variablen für den jeweiligen Gitterpunkt berechnen (Zeile 21-48). Weitergehend wird das Temperaturfeld nicht mehr in einem zwei sondern in einem eindimensionalen Objekt gespeichert. Dieses wird zunächst mit den gewünschten Startwerten initialisiert (Zeile 155-170).

Nun wird über das zu betrachtende Zeitintervall iteriert, und in jedem Integrationsschritt die zuvor aufgestellte Matrix invertiert. Dafür wurde das *SOR-Verfahren* im Quellcode (7) und (5) nach Gleichung (12) implementiert. Dabei werden in jedem einzelnen Iterationsschritt die Residuen $\vec{r} = \mathbf{M}\vec{x} - \vec{b}$ für die Lösung von $\mathbf{M}\vec{x} = \vec{b}$ berechnet, und gefordert, dass der Betrag dieser kleiner als $|\vec{r}| < 10^{-4}$ ist. Dies wird als Abbruchbedingung im dem *SOR-Verfahren* genutzt. Um die Vorteile des *SOR* wirklich nutzen zu können, muss zuerst ein geeigneter Relaxationsparameter α ermittelt werden. Hierzu wird $Pe = 10$, $N_x = N_y = 30$ und $\Delta t = 100$ gesetzt. Nun wird α mit einer Schrittweite von 0.05 ab 1.0 variiert und die benötigten Iterationen ermittelt.

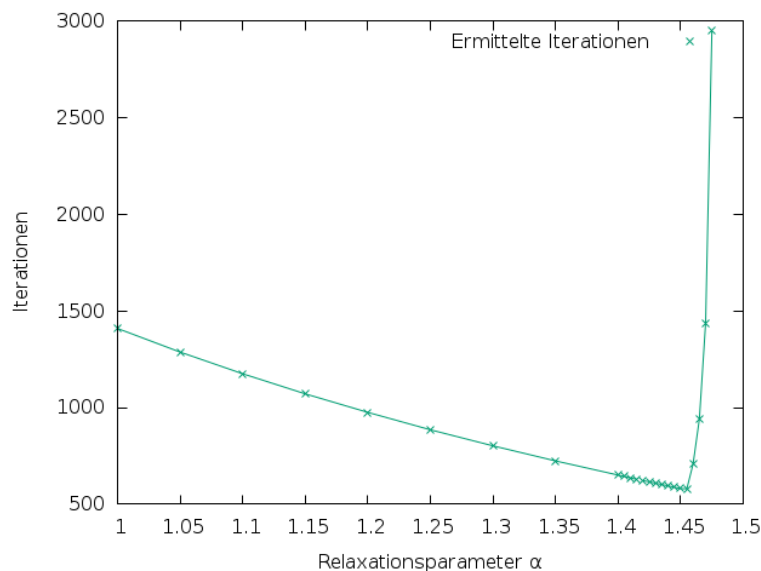


Abbildung 6: Grafische Darstellung der benötigten Iterationen zur Konvergenz in Abhängigkeit vom Relaxationsparameter α .

Die grafische Auftragung der so ermittelten Ergebnisse ist in Abbildung (6) zu finden. Es ist gut erkennbar, dass bei $\alpha = 1.455$ die geringste Anzahl an Iterationen benötigt wurde. Die Anzahl der benötigten Iterationen für höhere Relaxationsparameter α sind nicht mehr graphisch dargestellt, da hier für leicht erhöhte Werte bereits mehrere Zehntausend Iterationen benötigt wurden. Somit wird im weiteren Verlauf dieser Parameter auf 1.455 gesetzt.

Schließlich wird das zuvor implementierte und getestete Verfahren noch einmal für die Parameterwahl aus Aufgabenpunkt (4.4) angewendet, um die Ergebnisse miteinander zu vergleichen. Hierfür wurde also erneut $Pe = 10$, $\Delta t = 2 \cdot 10^{-4}$ und $t \in \{0.005, 0.05, 0.5\}$ gewählt. Die hierdurch erhaltenen Temperaturverteilungen sind erneut als Heat-map in Abbildung (7) dargestellt.

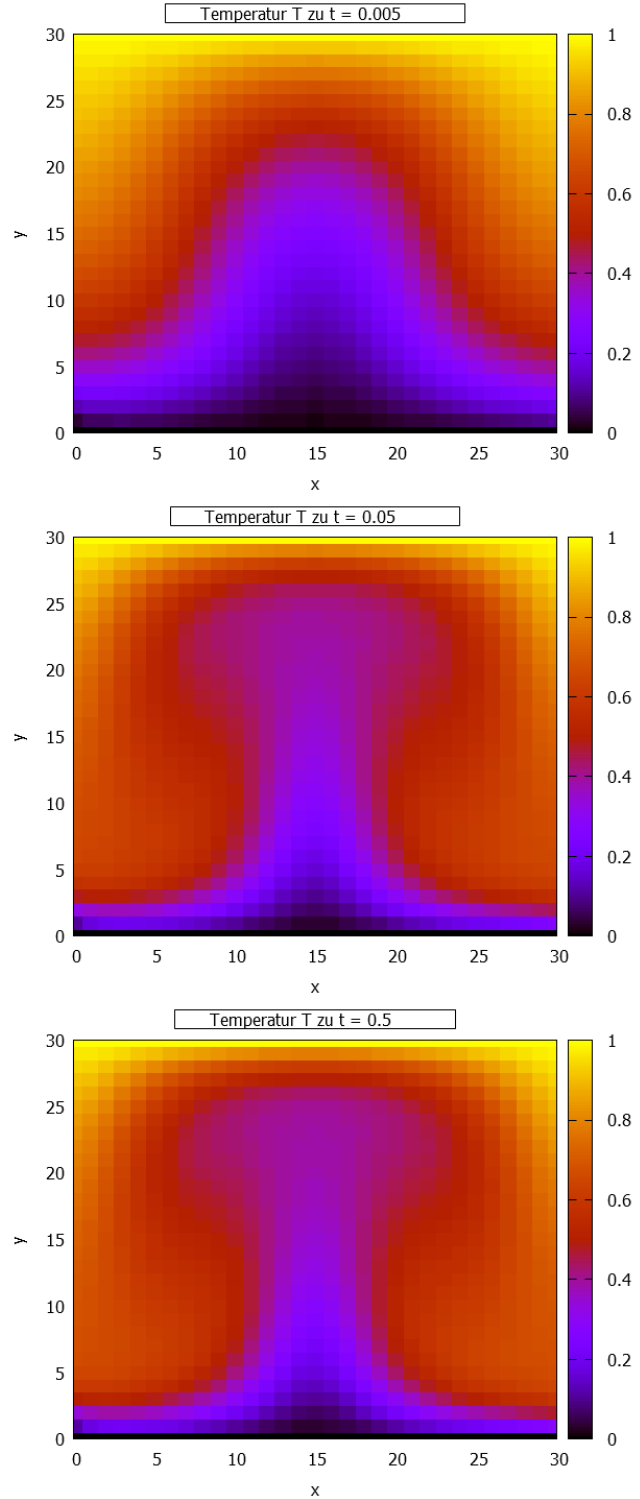


Abbildung 7: Temperatur nach dem *BTCS-Schema* für $Pe = 10$ und $\Delta t = 2 \cdot 10^{-4}$.

5. Auswertung der Ergebnisse

5.1. Numerische Korrektheit der FTCS-Lösung

Durch den Vergleich in Abschnitt (4.3) ist erkennbar, dass der numerische Fehler für größere Gitter deutlich geringer wird. Dieser Abfall hat allerdings auch eine bestimmte Grenze. Überschreitet $N = N_x = N_y$ diesen Wert, so kann die numerische Lösung nicht mehr konvergieren. Den Grund für dieses Verhalten (Abweichung der *CFL-Zahl*) ist in Abschnitt (4.5) beschrieben und erklärt. Anzumerken ist noch, dass der Fehler noch rapider gefallen wäre, wenn man den Gesamtfehler durch die Anzahl der zu berechnenden Gitterpunkte zwecks einer Normalisierung geteilt hätte.

5.2. Aufgabe 5: Beurteilung der Stabilität der FTCS-Lösung

Die zuvor in Abschnitt (4.5) hergeleitete Gleichung für den maximal gültigen Zeitschritt Δt_{max} ist durch die Abschätzungen nach oben, eine obere Schranke. Liegt λ außerhalb des Bereiches, aber nahe an der Grenze, so ist ein Konvergieren nicht sicher ausgeschlossen. Liegt λ allerdings in dem Bereich, so konvergiert das Verfahren sicher. Dieses Verhalten konnte an den zuvor empirisch ermittelten Grenzen für Δt verifiziert werden. Hierbei zeigte sich allerdings, dass die Aussagen der Gleichung (16) für geringe *Péclet-Zahlen* zu stark sind, und noch zu viele Werte ausschließen, für die das Verfahren noch konvergiert. Allerdings konnte durch exemplarische Tests bei höheren Werten ein Verhalten festgestellt werden, welches dem durch die Formel beschriebenen entspricht: je höher die *Péclet-Zahl* wird, desto feiner muss die Zeitintegration gewählt werden.

5.3. Aufgabe 6: Auswertung und Vergleich der BTCS-Lösung

Nachdem nun beide Verfahren umgesetzt und getestet worden sind, bietet sich ein Vergleich dieser an. Hierbei werden zuerst die dadurch erhaltenen Lösungen miteinander verglichen. Betrachtet man die Ergebnisse für die *Péclet-Zahl* $Pe = 10$ des *BTCS*-Verfahrens aus Abbildung (7) mit denen des *FTCS*-Verfahrens in Abbildung (5), so kann man eine hohe Übereinstimmung des Verlaufes feststellen. Hierbei ist für diese gegebene Parameterwahl kein Unterschied feststellbar.

Simulierte Zeitspanne t	$FTCS$ [sec]	$BTCS$ [sec]
1	< 1	6
10	3	7
50	15	8
100	33	9
200	68	9
400	135	10

Tabelle 2: Exemplarischer Laufzeitvergleich zwischen dem expliziten ($FTCS$) und dem impliziten ($BTCS$) Lösungsverfahren.

Da zuvor das explizite Lösungsverfahren auf numerische Korrektheit überprüft wurde, kann also auch davon ausgegangen werden, dass das implizite Lösungsverfahren korrekt umgesetzt worden ist.

Betrachtet man allerdings die Laufzeit der beiden Verfahren, so können große Unterschiede sowohl in der Größenordnung, als auch im Verlauf festgestellt werden. Diese Laufzeiten wurden exemplarisch gemessen und in Tabelle (5.3) festgehalten. Es ist anzumerken, dass die gemessene Laufzeit von dem verwendeten Rechner und anderen äußeren Einflüssen abhängen kann, sodass diese Daten nur als Anhaltspunkt für das grobe Verhalten gesehen werden können. Zudem hätten die verschiedenen Routinen auch noch weiter hinsichtlich ihrer Laufzeit optimiert werden können. Da dies den Rahmen des Themas der Arbeit überschreiten würde, wurde darauf verzichtet.

Es ist zu erkennen, dass die Laufzeit des expliziten Verfahrens in etwa linear mit der simulierten Zeitspanne t anwächst, wohingegen die des impliziten Verfahrens deutlich geringer anwächst, und näherungsweise fast konstant bleibt.

Somit zeigt sich, dass das explizite Verfahren zweiter Ordnung in diesem Fall seine Stärken nur bei sehr kurzen Zeitspannen hat. Das implizite Verfahren läuft bereits bei Simulationen für mehr als $t = 50$ schneller ab und ist zu bevorzugen.

Literatur

- [1] Andres Honecker, Thomas Pruschke, Salvatore R. Manmana,
Scriptum zur Vorlesung: Computergestütztes wissenschaftliches Rechnen.
Göttingen, Sommersemester 2016
- [2] W. Auzinger,
Lecture Notes Iterative Solution of Large Linear Systems.
Wien, 2011
- [3] F. Peric,
Computational Methods for Fluid Dynamics.
Stanford, 2002

A. Quellcode

A.1. Hauptprogramme

```
1 #include <iostream>
2 #include <fstream>
3 #include <vector>
4 #include <cstdlib>
5 #include <iomanip>
6 #include <cmath>
7
8 using namespace std;
9
10 //custom type which describes a 2d grid. use this instead of arrays for access
11 //violation safety
12 typedef vector<vector<double> > grid_t;
13
14 //integrates the ODE with the FTCS scheme and returns the error, compared with the
15 //stationary solution. expects the parameters:
16 // values: a grid_t object which contains current distribution of the temperature T
17 // v0x, v0y: a grid_t object which contains the x/y component of the velocity v0
18 // source: a grid_t object which contains the source values
19 // delta_t, delta_x, delta_y, pe: the time-step-size, the step size in x/y direction
20 // and the Péclet number
21 void ftcs_time_step(grid_t& values, grid_t v0x, grid_t v0y, grid_t source, const
22 double delta_t, const double delta_x, const double delta_y, const double pe)
23 {
24     //get the dimensions of the grid
25     size_t dimension_x = values.size();
26     size_t dimension_y = values.at(0).size();
27
28     //create copy of the grid with the same dimensions
29     //this will be used for the calculation of the derivatives
30     grid_t old_values;
31
32     for (size_t x=0; x<dimension_x; x++)
33     {
34         //add a new row
35         old_values.push_back(vector<double>());
36         for (size_t y=0; y<dimension_y; y++)
37         {
38             //copy value
39             old_values.at(x).push_back(values.at(x).at(y));
40         }
41     }
42
43     double gradient_t_times_v0;
44     double laplacian_t;
45
46     //loop over the entire inner grid (ignore the 4 outer sides of the rectangle) and
47     //calculate the new value using the FTCS-scheme
48     for (size_t x=1; x<dimension_x-1; x++)
49     {
50         for (size_t y=1; y<dimension_y-1; y++)
51         {
52             //calculates the 2nd ordner laplacian of the temp.
53             laplacian_t = 1.0/pow(delta_x,2)*(old_values.at(x+1).at(y) -
54                 2.0*old_values.at(x).at(y) + old_values.at(x-1).at(y))
55                 + 1.0/pow(delta_y,2)*(old_values.at(x).at(y+1) - 2.0*old_values.at(x).at(y)
56                 + old_values.at(x).at(y-1));
```

```

52     //calculates the inner product of the gradient of the temp. and v0
53     gradient_t_times_v0 =
54         v0x.at(x).at(y)/(2.0*delta_x)*(old_values.at(x+1).at(y)-old_values.at(x-1).at(y))
55         +
56         v0y.at(x).at(y)/(2.0*delta_y)*(old_values.at(x).at(y+1)-old_values.at(x).at(y-1));
57
58     //integration step according to eq. (8)
59     values.at(x).at(y) += delta_t * (laplacian_t - pe*gradient_t_times_v0 +
60         source.at(x).at(y));
61 }
62 }
63 //now take care of the four boundary conditions
64 //left/right boundaries (Neumann boundaries)
65 for (size_t y=0; y<dimension_y; y++)
66 {
67     values.at(0).at(y) = 1.0/3.0*(4.0*values.at(0+1).at(y) - values.at(0+2).at(y));
68     values.at(dimension_x-1).at(y) = -2.0/3.0*(-2.0*values.at(dimension_x-2).at(y) +
69         1.0/2.0*values.at(dimension_x-3).at(y));
70 }
71 //bottom/top boundaries (Dirichlet boundaries)
72 for (size_t x=0; x<dimension_x; x++)
73 {
74     values.at(x).at(0) = 0.0;
75     values.at(x).at(dimension_y-1) = 1.0;
76 }
77 }
78 int main(int argc, char* argv[])
79 {
80     //print information for the usage
81     cout << "Use this program to calculate the temperature with the FTCS-scheme in the
82         rectangle with a source, and to compare the stationary solution with the
83         numerical." << endl;
84     cout << "Expects the following set of parameters:" << endl;
85
86     cout << "\toutput file:\tThe path to the file in which the results will be stored."
87         << endl;
88     cout << "\tN_x:\t\tThe amount of grid points in the x-direction." << endl;
89     cout << "\tN_y:\t\tThe amount of grid points in the y-direction." << endl;
90     cout << "\tPe:\t\tThe Péclet-Number" << endl;
91     cout << "\tMax t\t\tThe maximum time until the system will be simulated." << endl;
92     cout << "\tDelta t\t\tThe time step size for the integration." << endl;
93     cout << "\tUse external heat source\t\tIf set to 1, then the source Q will be used."
94         << endl;
95
96     //read entered parameters
97
98     //create output stream
99     ofstream outputFile;
100     outputFile.open(argv[1]);
101     outputFile << fixed << setprecision(5);
102
103     //the amount of grid points in the x direction
104     const size_t dimension_x = (atoi)(argv[2])+1;
105     //the distance between to grid points
106     const double delta_x = 1.0/(dimension_x-1);
107
108     //the amount of grid points in the y direction
109     const size_t dimension_y = (atoi)(argv[3])+1;
110     //the distance between to grid points

```

```

107  const double delta_y = 1.0/(dimension_y-1);
108
109  //Péclet-number
110  const double pe = atof(argv[4]);
111  //maximum time until which the system will be simulated
112  const double max_t = atof(argv[5]);
113  //the time step size
114  const double delta_t = atof(argv[6]);
115
116  const bool use_source = atoi(argv[7]) == 1;
117
118  cout << "Entered parameters:" << endl;
119  cout << "\tNx: " << dimension_x << "\tNy: " << dimension_y << "\tPe: " << pe <<
    "\tMax t: " << max_t << "\tDelta t: " << delta_t << endl;
120
121  //to store the x component of the velocity v0
122  grid_t v0x;
123  //to store the y component of the velocity v0
124  grid_t v0y;
125  //to store the current temperature
126  grid_t values;
127  //to store the source of the temperature
128  grid_t sources;
129
130  //create the new grids
131  for (size_t x=0; x<dimension_x; x++)
132  {
133      //add a new row
134      v0x.push_back(vector<double>());
135      v0y.push_back(vector<double>());
136      values.push_back(vector<double>());
137      sources.push_back(vector<double>());
138
139      for (size_t y=0; y<dimension_y; y++)
140      {
141          //add the v0 values
142          v0x.at(x).push_back(M_PI*sin(2*M_PI*x*delta_x)*cos(M_PI*y*delta_y));
143          v0y.at(x).push_back(-2.0*M_PI*cos(2*M_PI*x*delta_x)*sin(M_PI*y*delta_y));
144
145          //add the start temperature of the rectangle
146          values.at(x).push_back(y*delta_y);
147
148          //add the source of the temperature of the rectangle
149          if (use_source)
150          {
151              sources.at(x).push_back(
152                  pow(M_PI,2)*(cos(M_PI*delta_x*x)*sin(M_PI*delta_y*y)*2
153                      - pe*pow(cos(x*delta_x*M_PI),3)*sin(2*M_PI*delta_y*y))
154                      -pe*2*M_PI*cos(2*M_PI*delta_x*x)*sin(M_PI*y*delta_y)
155              );
156          }
157          else
158          {
159              //no heat source is desired => set it to 0 for all points
160              sources.at(x).push_back(0.0);
161          }
162      }
163  }
164
165  //start to measure the runtime now
166  clock_t start_time = clock();
167
168  //to print just every 100th error comparison

```

```

169 size_t iteration = 0;
170 //is being used to calculate the numerical error compared to stat.solution
171 double error;
172 //loop over the time, until we have reached the maximum time
173 for (double t=0.0; t<max_t; t+=delta_t)
174 {
175     ftcs_time_step(values, v0x, v0y, sources, delta_t, delta_x, delta_y, pe);
176
177     //the heat source is beeing used => calculate the difference
178     if (use_source)
179     {
180         //now calculate the difference, compared to the stationary solution
181         //loop over all elements and compare them
182         error = 0.0;
183         for (size_t x=0; x<dimension_x; x++)
184         {
185             for (size_t y=0; y<dimension_y; y++)
186             {
187                 error += pow(values.at(x).at(y) -
188                             (cos(M_PI*delta_x*x)*sin(M_PI*delta_y*y)+y*delta_y), 2);
189             }
190         }
191
192         //print every 100th error value
193         if (iteration++ % 100 == 0)
194         {
195             cout << "time:\t" << t << "\terror:\t" << error << "\n" << flush;
196             iteration = 0;
197         }
198     }
199
200     //stop the time measurement
201     clock_t end_time = clock();
202
203     cout << "Integration finished after " << (end_time - start_time) / CLOCKS_PER_SEC <<
204          " seconds.";
205
206     //print the final grid's values in the givven output file in a matrix form
207     for (size_t y=0; y<dimension_y; y++)
208     {
209         for (size_t x=0; x<dimension_x; x++)
210         {
211             outputFile << values.at(x).at(y) << "\t";
212         }
213         outputFile << "\n";
214     }
215
216     return 1;
217 }

```

Quellcode 1: Implementation des FTCS-Schemas

```

1 #include <iostream>
2 #include <fstream>
3 #include <vector>
4 #include <cstdlib>
5 #include <iomanip>
6 #include <cmath>
7 #include "../lib/linear_system.h"
8 #include "../lib/square_matrix.h"
9
10 using namespace std;

```

```

11 using namespace numerical;
12
13 //custom type which describes a 2d grid. use this instead of arrays for access
14     violation safety
15 typedef vector<vector<double>> grid_t;
16
17 /*calculates the variable beta. expects the parameters:
18     plus: indicates whether beta plus or beta minus will be calculated
19     delta_t, delta_x, delta_y, pe: the time-step-size, the step size in x/y direction
20     and the Péclet number
21     x,y: the x/y coordinate of the point on which beta is desired
22 */
23 double beta(bool plus, double delta_t, double delta_x, double delta_y, double pe,
24     double x, double y)
25 {
26     if (plus)
27     {
28         return delta_t*(-1.0 / pow(delta_x, 2) + (pe*(M_PI*sin(2 *
29             M_PI*x*delta_x)*cos(M_PI*y*delta_y))) / (2.0*delta_x));
30     }
31     else
32     {
33         return delta_t*(-1.0 / pow(delta_x, 2) - (pe*(M_PI*sin(2 *
34             M_PI*x*delta_x)*cos(M_PI*y*delta_y))) / (2.0*delta_x));
35     }
36 }
37
38 /*calculates the variable gamma. expects the parameters:
39     plus: indicates whether gamma plus or gamma minus will be calculated
40     delta_t, delta_x, delta_y, pe: the time-step-size, the step size in x/y direction
41     and the Péclet number
42     x,y: the x/y coordinate of the point on which gamma is desired
43 */
44 double gamma(bool plus, double delta_t, double delta_x, double delta_y, double pe,
45     double x, double y)
46 {
47     if (plus)
48     {
49         return delta_t*(-1.0 / pow(delta_y, 2) + (pe*(-2.0*M_PI*cos(2 *
50             M_PI*x*delta_x)*sin(M_PI*y*delta_y))) / (2.0*delta_y));
51     }
52     else
53     {
54         return delta_t*(-1.0 / pow(delta_y, 2) - (pe*(-2.0*M_PI*cos(2 *
55             M_PI*x*delta_x)*sin(M_PI*y*delta_y))) / (2.0*delta_y));
56     }
57 }
58
59 int main(int argc, char* argv[])
60 {
61     //print information for the usage
62     cout << "Use this program to calculate the temperature with the BTCS-scheme in the
63         rectangle." << endl;
64     cout << "Expects the following set of parameters:" << endl;
65
66     cout << "\toutput file:\tThe path to the file in which the results will be stored."
67         << endl;
68     cout << "\tN_x:\t\tThe amount of grid points in the x-direction." << endl;
69     cout << "\tN_y:\t\tThe amount of grid points in the y-direction." << endl;
70     cout << "\tPe:\t\tThe Péclet-Number" << endl;
71     cout << "\tMax t\t\tThe maximum time until the system will be simulated." << endl;
72     cout << "\tDelta t\t\tThe time step size for the integration." << endl;
73 }

```

```

63 //read entered parameters
64
65 //create output stream
66 ofstream outputFile;
67 outputFile.open(argv[1]);
68 outputFile << fixed << setprecision(5);
69
70 //the amount of grid points in the x direction
71 const size_t dimension_x = (atoi)(argv[2])+1;
72 //the distance between to grid points
73 const double delta_x = 1.0/(dimension_x-1);
74
75 //the amount of grid points in the y direction
76 const size_t dimension_y = (atoi)(argv[3])+1;
77 //the distance between to grid points
78 const double delta_y = 1.0/(dimension_y-1);
79
80 //Péclet-number
81 const double pe = atof(argv[4]);
82 //maximum time until which the system will be simulated
83 const double max_t = atof(argv[5]);
84 //the time step size
85 const double delta_t = atof(argv[6]);
86
87 cout << "Entered parameters:" << endl;
88 cout << "\tNx: " << dimension_x << "\tNy: " << dimension_y << "\tPe: " << pe <<
      "\tMax t: " << max_t << "\tDelta t: " << delta_t << endl;
89
90 //the length of the value vector
91 size_t n = dimension_x*dimension_y;
92
93 //create a new square matrix (nXn) to store the matrix M
94 square_matrix coeff_matrix(n, 0);
95
96 //create const value for the variable alpha
97 const double alpha = delta_t * (2.0 / pow(delta_x, 2) + 2.0 / pow(delta_y, 2)) + 1.0;
98
99 //now set the values of the three matrices
100
101 //unit matrix
102 //main diagonal for the dirichlet boundaries again
103 for (size_t i = 0; i<dimension_x; i++)
104 {
105     coeff_matrix.set_value(i, i, 1.0);
106 }
107
108 //Matrix M
109 //contains the biggest part of the BTCS-scheme
110 for (size_t i = dimension_x+1; i<n - dimension_x-1; i++)
111 {
112     //ingore this diagonals, as they affect the boundaries
113     if (i % dimension_x == 0 || i % dimension_x == dimension_x-1)
114     {
115         continue;
116     }
117
118     //i % dimension_x gives the x coordinate of the field-value on which the current
119     //matrix's item operates
120     //i / dimension_x gives the y coordinate of the field-value on which the current
121     //matrix's item operates
122     coeff_matrix.set_value(i, i - dimension_x, gamma(false, delta_t, delta_x, delta_y,
123     pe, i % dimension_x, i/dimension_x));
124     coeff_matrix.set_value(i, i - 1, beta(false, delta_t, delta_x, delta_y, pe, i %

```

```

122     dimension_x, i/dimension_x));
123     coeff_matrix.set_value(i, i, alpha);
124     coeff_matrix.set_value(i, i + 1, beta(true, delta_t, delta_x, delta_y, pe, i %
125         dimension_x, i/dimension_x));
126     coeff_matrix.set_value(i, i + dimension_x, gamma(true, delta_t, delta_x, delta_y,
127         pe, i % dimension_x, i/dimension_x));
128 }
129
130 //now set the repeating line which take care of the neumann boundaries
131 //first left boundary
132 coeff_matrix.set_value(dimension_x, dimension_x+2, 1.0/3.0);
133 coeff_matrix.set_value(dimension_x, dimension_x+1, -4.0/3.0);
134 coeff_matrix.set_value(dimension_x, dimension_x, 1);
135 for (size_t i=2; i<dimension_x-1; i++)
136 {
137     //right
138     coeff_matrix.set_value(i*dimension_x-1, i*dimension_x-3, 1.0/3.0);
139     coeff_matrix.set_value(i*dimension_x-1, i*dimension_x-2, -4.0/3.0);
140     coeff_matrix.set_value(i*dimension_x-1, i*dimension_x-1, 1);
141
142     //left
143     coeff_matrix.set_value(i*dimension_x, i*dimension_x, 1);
144     coeff_matrix.set_value(i*dimension_x, i*dimension_x+1, -4.0/3.0);
145     coeff_matrix.set_value(i*dimension_x, i*dimension_x+2, 1.0/3.0);
146 }
147
148 //last right boundary
149 coeff_matrix.set_value((dimension_x)*(dimension_x-1)-1,
150     (dimension_x)*(dimension_x-1)-3, 1.0/3.0);
151 coeff_matrix.set_value((dimension_x)*(dimension_x-1)-1,
152     (dimension_x)*(dimension_x-1)-2, -4.0/3.0);
153 coeff_matrix.set_value((dimension_x)*(dimension_x-1)-1, (dimension_x)*(dimension_x-1)-1,
154     1);
155
156 //unit matrix
157 //main diagonal for the dirichlet boundaries again
158 for (size_t i = n - dimension_x; i<n; i++)
159 {
160     coeff_matrix.set_value(i, i, 1.0);
161 }
162
163 //create two vectors to store the temperature and to calculate the new one
164 vector<double> values;
165 vector<double> old_values;
166
167 //fill them with the starting temperature
168 for (size_t y = 0; y<dimension_y; y++)
169 {
170     for (size_t x = 0; x<dimension_x; x++)
171     {
172         old_values.push_back(y*delta_y);
173         values.push_back(y*delta_y);
174     }
175 }
176
177 //start to measure the runtime now
178 clock_t start_time = clock();
179
180 //loop over the time, until we have reached the maximum time
181 for (double t = 0.0; t<max_t; t += delta_t)
182 {
183     //set the right hand side of the equation (set some values to zero for the neuman

```

```

179     boundaries)
180     old_values.at(dimension_x) = 0;
181     for (size_t i=2; i<dimension_x-1; i++)
182     {
183         old_values.at(i*dimension_x-1) = 0;
184         old_values.at(i*dimension_x) = 0;
185     }
186     old_values.at((dimension_x)*(dimension_x-1)-1) = 0;
187
188     //invert the matrix M/solve the linear system to get the new temperature(value)
189     cout << "t: " << t << "\titerations: " << linear_system_solve_sor(coeff_matrix,
190         values, old_values, 1.455, 1e-4) << endl;
191
192     //copy the new values into the old_values vector for the next iteration
193     for (size_t i = 0; i<n; i++)
194     {
195         old_values.at(i) = values.at(i);
196     }
197
198     //measure the end time
199     clock_t end_time = clock();
200
201     //print the runtime
202     cout << "Integration finished after " << (end_time - start_time) / CLOCKS_PER_SEC <<
203         " seconds.";
204
205     //print the final grid's values in the given output file in a matrix form
206     for (size_t i = 0; i<n; i++)
207     {
208         outputFile << values.at(i) << "\t";
209         if (i > 0 && (i+1) % dimension_x == 0)
210         {
211             outputFile << endl;
212         }
213     }
214     return 1;
215 }

```

Quellcode 2: Implementation des BTCS-Schemas

A.2. Hilfsbibliotheken

```
1 #include "matrix.h"
2
3 #ifndef SQUAREMATRIX_H
4 #define SQUAREMATRIX_H
5 namespace numerical
6 {
7     class square_matrix : public matrix
8     {
9     public:
10         //creates a nXn matrix
11         square_matrix(int n) : matrix(n,n) {}
12         //creates a nXn matrix and initializes all values with defaultValue
13         square_matrix(int n, double defaultValue) : matrix(n,n, defaultValue) {};
14     };
15 }
16 #endif
```

Quellcode 3: square_matrix.h

```
1 #include <fstream>
2
3 #ifndef MATRIX_H
4 #define MATRIX_H
5 namespace numerical{
6     class matrix
7     {
8     public:
9         //creates a new nXm matrix
10         matrix(size_t n, size_t m);
11         //creates a new nXm matrix and initializes all values with defaultValue
12         matrix(size_t n, size_t m, double defaultValue);
13         //destructor
14         ~matrix();
15
16         //returns the number of rows/columns
17         size_t get_n();
18         size_t get_m();
19
20         //returns the value in row n, column m
21         double get_value(size_t n, size_t m);
22         //sets the value in row n, column m
23         void set_value(size_t n, size_t m, double value);
24     private:
25         double** data;
26         size_t n;
27         size_t m;
28     };
29 }
30 #endif
```

Quellcode 4: matrix.h

```
1 #include "square_matrix.h"
2 #include <vector>
3
4 using namespace std;
5
6 #ifndef LINEARSYSTEM_H
7 #define LINEARSYSTEM_H
8 namespace numerical
```

```

9 {
10     size_t linear_system_solve_sor(square_matrix& coeff_matrix, vector<double> &x,
        vector<double> &b, double alpha, double error_threshold);
11     double linear_system_solve_sor_step(square_matrix& coeff_matrix, vector<double> &x,
        vector<double> &b, double alpha);
12
13     void multiply_matrix_vector(matrix& matrix, vector<double> &x, vector<double>
        &result);
14 }
15 #endif

```

Quellcode 5: linear_system.h

```

1 #include "matrix.h"
2 #include <iostream>
3
4 namespace numerical{
5     //creates a new nXm matrix
6     matrix::matrix(size_t n, size_t m)
7     {
8         this->n = n;
9         this->m = m;
10        data = new double*[m];
11        for (size_t i=0;i<m;i++)
12        {
13            data[i] = new double[n];
14        }
15    }
16
17    //creates a new nXm matrix and initializes all values with defaultValue
18    matrix::matrix(size_t n, size_t m, double defaultValue)
19    {
20        this->n = n;
21        this->m = m;
22        data = new double*[m];
23        for (size_t i=0;i<m;i++)
24        {
25            data[i] = new double[n];
26            for (size_t j=0;j<n;j++)
27            {
28                data[i][j] = defaultValue;
29            }
30        }
31    }
32
33    //destructor
34    matrix::~matrix()
35    {
36        delete data;
37    }
38
39    //sets the value in the nth row, mth column
40    void matrix::set_value(size_t n, size_t m, double value)
41    {
42        data[n][m] = value;
43    }
44
45    //returns the value in the nth row, mth column
46    double matrix::get_value(size_t n, size_t m)
47    {
48        return data[n][m];
49    }
50

```

```

51 //returns the number of rows
52 size_t matrix::get_n()
53 {
54     return this->n;
55 }
56
57 //returns the number of columns
58 size_t matrix::get_m()
59 {
60     return this->m;
61 }
62 }

```

Quellcode 6: matrix.cpp

```

1 #include "linear_system.h"
2 #include "square_matrix.h"
3 #include <iostream>
4 #include <fstream>
5 #include <cmath>
6 #include <vector>
7 #include <float.h>
8
9 using namespace std;
10
11 namespace numerical
12 {
13     void multiply_matrix_vector(matrix& matrix, vector<double> &xVec, vector<double>
        &result)
14     {
15         //clear the result vector and reserve enough space
16         result.clear();
17         double val;
18         result.reserve(xVec.size());
19         for (size_t y = 0; y<matrix.get_n(); y++)
20         {
21             val = 0.0;
22             for (size_t x = 0; x<matrix.get_m(); x++)
23             {
24                 val += matrix.get_value(y, x) * xVec.at[x];
25             }
26             result.push_back(val);
27         }
28     }
29
30     //calculates one SOR step the for system coeff_matrix * x = b with the relaxation
        parameter alpha
31     double linear_system_solve_sor_step(square_matrix& coeff_matrix, vector<double> &x,
        vector<double> &b, double alpha)
32     {
33         //the new value of the current item
34         double new_value;
35         //a helper variable to store a needed sum of items in the loop
36         double sum;
37
38         //loop over all items and calculate the next elemt of x
39         for (size_t k = 0; k<coeff_matrix.get_n(); k++)
40         {
41             sum = 0.0;
42
43             //splitted the one loop into two separate loops for better performance.
44             //one of the loops is already using the new values while the other one is using
                the old one

```

```

45     for (size_t i = 0; i < k; i++)
46     {
47         sum += coeff_matrix.get_value(k, i) * x.at(i);
48     }
49     for (size_t i = k + 1; i < coeff_matrix.get_n(); i++)
50     {
51         sum += coeff_matrix.get_value(k, i) * x.at(i);
52     }
53     new_value = (1.0 - alpha) * x.at(k) + alpha / coeff_matrix.get_value(k,
54         k) * (b.at(k) - sum);
55
56     //update the new x value
57     x.at(k) = new_value;
58 }
59
60 //calculate the residue now
61 double residue = 0.0;
62 vector<double> testVector;
63 multiply_matrix_vector(coeff_matrix, x, testVector);
64 for (size_t k = 0; k < coeff_matrix.get_n(); k++)
65 {
66     residue += pow(testVector.at(k) - b.at(k), 2);
67 }
68
69 //return the length of the residue vector
70 return sqrt(residue);
71 }
72
73 //tries to solve the linear equation coeff_matrix * x = b for x iterative with the
74 //SOR method untill either:
75 //1) the residues are lower than the specified error_threshold
76 //2) the residues of 50 following iterations have been increasing => propably no
77 //convergence for this system
78 //returns the needed iterations
79
80 size_t linear_system_solve_sor(square_matrix& coeff_matrix, vector<double> &x,
81     vector<double> &b, double alpha, double error_threshold)
82 {
83     //contains the value of the residues of the current iteration. start with a high
84     //number because of the if-comparison below
85     double error = DBLMAX;
86     //contains the value of the residues of the previous iteration. start with a high
87     //number because of the if-comparison below
88     double old_error;
89     //counter of the needed iterations
90     size_t iteration = 0;
91
92     //contains how often the residues of two following iterations have been increasing
93     size_t residue_increased_counter = 0;
94
95     //loop, until a solution has been found
96     while (true)
97     {
98         //swap the old and the current error value
99         old_error = error;
100         //increase the number of needed iterations
101         iteration++;
102
103         //calculate the next SOR step
104         error = linear_system_solve_sor_step(coeff_matrix, x, b, alpha);
105
106         //check for exit condition 1)

```

```
102     if (error < error_threshold)
103     {
104         return iteration;
105     }
106     else if (error > old_error)
107     {
108         //exit condition 2
109         if (residue_increased_counter > 50){
110             return iteration;
111         }
112         residue_increased_counter++;
113     }
114     else
115     {
116         //reset the counter for the second exit condition.
117         residue_increased_counter=0;
118     }
119 }
120 }
121 }
```

Quellcode 7: linear_system.cpp