1 Übersetzung in relationale Algebra (8 Punkte)

Gegeben sei das Relationenschema aus der zweiten Übung.

```
[Personen]: {[PID: int, Name: string, Geburtsjahr: int, Wohnort: string]}
[Schüler*innen]: {[SID:(Personen→PID), Klassenstufe: int, Klassenraum: int, Schulform: string]}
[Klausuren]: {[KID: int, Fach: string, Thema: string, Dauer: int]}
[Lehrer*innen]: {[LID:(Personen→PID), Hauptfach:string, Gehalt:int, Dienstjahre:int]}
[unterrichten]: {[Schüler*in:(Schüler*innen→SID), Datum:date, Uhrzeit:time,

Lehrer*in:(Lehrer*innen→LID), Fach: string)]}
[korrigieren]: {[Klausur:(Klausuren→KID), Schüler*in:(Schüler*innen→SID),

Lehrer*in:(Lehrer*innen→LID), Note: int]}
```

Um die Lesbarkeit Ihrer Ausdrücke zu verbessern, dürfen Sie die Teilergebnisnotation aus der zweiten Übung benutzen. Benennen Sie zum Beispiel einen Ausdruck wie folgt

$$R_3 := R_1 \bowtie_{A=B} R_2$$

so kann R_3 nun in darauffolgenden Ausdrücken verwendet werden. Zudem dürfen Sie den syntaktischen Zucker für die gleichzeitige Umbenennung mehrerer Attribute aus der zweiten Übung nutzen

$$\rho_{A'\leftarrow A,...,Z'\leftarrow Z} R = \rho_{A'\leftarrow A} (... (\rho_{Z'\leftarrow Z} R)).$$

- (a) Übersetzen Sie die folgenden umgangssprachliche Anfragen in Ausdrücke der relationalen Algebra:
 - Die schlechteste (größte) Note, die in Klausuren über das Thema 'Quantenmechanik' erreicht wurde.
 - 2. Die durchschnittlichen Noten je Klausur, die von Schüler*innen der Klassenstufe 12 geschrieben wurde.
 - 3. Die Namen der Lehrer*innen, die schon mehr als 7 Schüler*innen bei einer Klausur über das Thema 'Analysis' korrigiert haben.
 - 4. Die Schulformen der Schüler*innen, die noch nie an einer Klausur teilgenommen haben, in der mehr als 10 Schüler*innen mindestens die Note 3 erreicht haben.
 - 5. Das Geburtsjahr der ältesten Schüler*innen, die genau eine Klausur geschrieben haben, in der sie von einem/einer Lehrer*in korrigiert wurden, der/die diese Schüler*innen bereits im Fach Informatik unterrichtet hat.
- (b) Übersetzen Sie folgende Anfragen aus der relationalen Algebra in natürliche Sprache:
 - 1. $R := (\sigma_{\text{Klassenstufe}} = 10 \land \text{Schulform} = \text{`Gymnasium'} \text{ Schüler*innen}) \bowtie_{\text{SID}} = \text{Schüler*in} \text{ unterrichten}$ $\pi_{\text{Gehalt}}(R \bowtie_{\text{Lehrer*in}=\text{LID}} (\sigma_{\text{Hauptfach}='\text{Sport'}} \text{ Lehrer*innen}))$
 - 2. $R_1 := \gamma_{\text{avg}(\text{Gehalt})} ((\sigma_{\text{Dienstjahre}>10} \text{ Lehrer*innen}) \bowtie_{\text{LID}=\text{Lehrer*in}} (\pi_{\text{Lehrer*in}} \text{ unterrichten}))$ $R_2 := (\text{Lehrer*innen} \bowtie_{\text{Gehalt}<\text{avg}(\text{Gehalt})} R_1) \bowtie_{\text{LID}=\text{Lehrer*in}} \text{ korrigieren}$ $\pi_{\text{Schulform}} (\text{Schüler*innen} \bowtie_{\text{SID}=\text{Schüler*in}} R_2)$
 - 3. $R_1 := (\sigma_{\text{Fach} = \text{'Physik'}} \vee \text{Fach} = \text{'Biologie'}$ Klausuren) $\bowtie_{\text{KID} = \text{Klausur}}$ korrigieren $R_2 := \text{korrigieren} \bowtie_{\text{Note} = \text{max(Note)}} \wedge \text{Klausur} = \text{KID}(\gamma_{\text{KID,max(Note)}} R_1)$ $R_3 := \sigma_{\text{count}(*)} > 5 (\gamma_{\text{Schüler*in,count}(*)} R_2)$ $\pi_{\text{Name}} (R_3 \bowtie_{\text{Schüler*in} = \text{PID}} \text{Personen})$

Lösung:

- (a) Vorgesehen ist ein Punkt pro Anfrage (je 0,5 Punkte Abzug für inkorrekte (Join-)Prädikate oder Aggregate, unterschiedliche Schemata bei Mengenoperationen, nicht-disjunkte Schemata bei Joins, uneindeutige Bezeichner, inkorrekte Bezeichner durch Verwechslung von Umbenennung und Teilergebnisnotation, usw.). Wichtig ist hierbei nur, dass die Anfragen das korrekte Ergebnis produzieren, nicht, dass sie mit der Musterlösung übereinstimmen.
 - 1. $\gamma_{\text{max(Note)}}$ (($\sigma_{\text{Thema}} = \gamma_{\text{Quantenmechanik}}$ Klausuren) $\bowtie_{\text{KID}} = \text{Klausur}$ korrigieren)
 - 2. $R := \text{korrigieren} \bowtie_{\text{Schüler*in} = \text{SID}} (\sigma_{\text{Klassenstufe}=12} \text{Schüler*innen})$ $\pi_{\text{avg(Note)}}(\gamma_{\text{Klausur,avg(Note)}} R)$
 - 3. $R_1 := \text{korrigieren} \bowtie_{\text{Klausur} = \text{KID}} (\sigma_{\text{Thema}='\text{Analysis'}} \text{Klausuren})$ $R_2 := \sigma_{\text{count}(*) > 7}(\gamma_{\text{Klausur, Lehrer*in,count}(*)} R_1)$ $\pi_{\text{Name}}(R_2 \bowtie_{\text{Lehrer*in} = \text{PID}} \text{Personen})$
 - 4. $R_1 := \rho_{\text{Klausur}} \leftarrow_{\text{Klausur}} (\sigma_{\text{count}(*)} > 10 \ (\gamma_{\text{Klausur},\text{count}(*)} (\sigma_{\text{Note} \leq 3} \ \text{korrigieren})))$ $R_2 := \text{Schüler*innen} \bowtie_{\text{SID} = \text{Schüler*in}} (\text{korrigieren} \bowtie_{\text{Klausur} = \text{Klausur}}, R_1)$ $\pi_{\text{Schulform}} (\pi_{\text{SID}, \text{Schulform}} \text{Schüler*innen} \pi_{\text{SID}, \text{Schulform}} R_2)$
 - 5. $R_1 := \rho_{\text{Lehrer*in'}\leftarrow \text{Lehrer*in}, \text{Schüler*in'}\leftarrow \text{Schüler*in}} \text{korrigieren}$ $R_2 := (\sigma_{\text{Fach = 'Informatik'}} \text{ unterrichten}) \bowtie_{\text{Lehrer*in}} = \text{Lehrer*in'} \land \text{Schüler*in} = \text{Schüler*in'}, R_1$ $R_3 := \sigma_{\text{count(*)} = 1}(\gamma_{\text{Schüler*in}, \text{count(*)}} (\pi_{\text{[korrigieren]}} R_2))$ $\gamma_{\min(\text{Geburtsjahr})}(R_3 \bowtie_{\text{Schüler*in}} = \text{PID Personen})$
- (b) Vorgesehen ist ein Punkt pro Übersetzung.
 - 1. Die Gehälter der Lehrer*innen, die Sport als Hauptfach haben und einmal ein/e Schüler*in in der 10. Klasse auf dem Gymnasium unterrichtet haben.
 - 2. Die Schulformen der Schüler*innen, deren Klausur mindestens einmal von einem/einer Lehrer*in korrigiert wurde, deren Gehalt geringer ist als das Durchschnittsgehalt der Lehrer*innen, die bereits mehr als 10 Dienstjahre haben und unterrichtet haben.
 - 3. Die Namen der Schüler*innen, die bereits mehr als fünf mal in einer Klausur in den Fächern Physik oder Biologie Klassenschlechteste*r waren, also die schlechteste Note hatten.

2 Relationale Anfragen (7 Punkte)

Betrachten Sie das vereinfachte IMDb-Schema aus der Vorlesung und den Notebooks.

Um die Lesbarkeit Ihrer Ausdrücke zu verbessern, dürfen Sie wieder die Teilergebnisnotation und den syntaktischen Zucker für die gleichzeitige Umbenennung mehrerer Attribute nutzen, die in den vorherigen Aufgaben eingeführt wurden.

Übersetzen Sie die folgenden umgangssprachlichen Anfragen in Ausdrücke der relationalen Algebra:

- 1. Die Jahre, in denen eine Komödie veröffentlicht wurde.
- 2. Die Vor- und Nachnamen der Schauspieler*innen, die schon einmal eine Rolle im Film 'Star Wars' gespielt haben.
- 3. Die Nachnamen und die Anzahl an Filmen, die Regisseur*innen mit dem gleichen Nachnamen gedreht haben.
- 4. Die Genre der Filme mit der besten Bewertung.
- 5. Die Nachnamen der Schauspieler*innen, die die Rolle des/der Superman/Superwoman in einem Actionfilm gespielt haben.
- 6. Die Anzahl an Horrorfilmen, die vom Regisseur James Wan nach 2000 gedreht wurden.
- 7. Die höchste Anzahl an Filmen die von einer/einem Regisseur*in gedreht worden sind, dessen/deren Filme vor 1995 mindestens eine Durchschnittswertung von 3 haben.

Lösung:

Vorgesehen ist ein Punkt pro Anfrage (je 0,5 Punkte Abzug für inkorrekte (Join-)Prädikate oder Aggregate, unterschiedliche Schemata bei Mengenoperationen, nicht-disjunkte Schemata bei Joins, uneindeutige Bezeichner, inkorrekte Bezeichner durch Verwechslung von Umbenennung und Teilergebnisnotation, usw.). Wichtig ist hierbei nur, dass die Anfragen das korrekte Ergebnis produzieren, nicht, dass sie mit der Musterlösung übereinstimmen.

```
    π<sub>year</sub> ((σ<sub>genre = 'Komödie'</sub> movies_genres) ⋈<sub>movie_id = id</sub> movies)
    R<sub>1</sub> := (σ<sub>name = 'Star Wars'</sub> movies) ⋈<sub>id = movie_id</sub> roles
π<sub>first_name, last_name</sub> (R<sub>1</sub> ⋈<sub>actor_id = id'</sub> (ρ<sub>id'←id</sub> actors))
    γ<sub>last_name,count(*)</sub>(π<sub>last_name, movie_id</sub>(directors ⋈<sub>id = director_id</sub> movie_directors))
    π<sub>genre</sub> (movies_genres ⋈<sub>movie_id</sub> = id (movies ⋈<sub>rank = max(rank)</sub> (γ<sub>max(rank)</sub> movies)))
    R<sub>1</sub> := ρ<sub>movie_id'←movie_id</sub>(σ<sub>genre='Action'</sub> movies_genres)
R<sub>2</sub> := R<sub>1</sub> ⋈<sub>movie_id'=movie_id</sub> (σ<sub>role='Superman'</sub> ∨ role='Superwoman'</sub> roles)
π<sub>last_name</sub> (R<sub>2</sub> ⋈<sub>actor_id=id</sub> actors)
```



6. $R_1 := \sigma_{\text{first name}='James' \land last name='Wan'}$ directors $R_2 := (\sigma_{\text{genre} = 'Horror'}, \text{ movies_genres}) \bowtie_{\text{movie_id} = \text{movie_id'}} (\rho_{\text{movie_id'} \leftarrow \text{movie_id}} (\text{movie_directors}))$ $\gamma_{\operatorname{count}(*)}((R_2 \bowtie_{\operatorname{director_id=id}} R_1) \bowtie_{\operatorname{movie_id}} = \operatorname{m_id}(\rho_{\operatorname{m_id} \leftarrow \operatorname{id}}(\sigma_{\operatorname{year} > 2000} \ \operatorname{movies})))$ 7. $R_1 := \text{movie_directors} \bowtie_{\text{movie_id=id}} (\sigma_{\text{year} < 1995} \text{ movies})$ $R_2 := \rho_{\mathrm{director_id'}\leftarrow\mathrm{director_id}}(\sigma_{\mathrm{avg(rank)}\geq 3} \ (\gamma_{\mathrm{director_id,avg(rank)}} \ R_1))$ $\gamma_{\max(\mathrm{count}(*))} \ \left(\gamma_{\mathrm{director_id},\mathrm{count}(*)} \ \left(\mathrm{movie_directors} \ \bowtie_{\mathrm{director_id}} = \mathrm{director_id}, \ R_2) \right) \right)$



3 Implementierung Co-Group-Join (5 Punkte)

Implementieren Sie den Co-Group-Join in den beiden mit *Exercise* markierten Zellen im beigefügten Notebook. Dieses finden Sie zusammen mit den anderen zur Ausführung benötigten Dateien im Anhang der Übung. Hierbei müssen die Klassen Equi_Join und Co_Group_Join vervollständigt werden.

Die Klasse Equi_Join erweitert dabei die Klasse BinaryOperator für logische, binäre Operatoren. Vervollständigen Sie hier die folgende Methode:

• get_schema: Gibt das Ausgabeschema des Operators zurück. Hier müssen Sie die Schemata der beiden Kinder konkatenieren.

Die Klasse Co_Group_Join erweitert die Klasse Equi_Join um die physische Funktionalität des Auswertens. Hierbei soll ein Equi_Join mit Hilfe einer Co-Gruppierung ausgeführt werden. Dabei soll eine Auswertung in linearer Laufzeit $\mathcal{O}(n)$ gewährleistet sein. Beachten Sie, dass diese Laufzeitbeschränkung nur gelten soll, sofern man über die Schlüssel beider Relationen joint. Ihre Implementierung soll trotzdem alle Fälle abdecken. Implementieren Sie die folgende Funktion:

• evaluate: Wertet den Operatorbaum rekursiv aus. Beachten Sie, dass ihr Join in linearer Laufzeit $\mathcal{O}(n)$ in der Anzahl an Elementen der Eingaberelationen durchgeführt werden soll, sofern über die Schlüssel beider Relationen gejoint wird. Es dürfen keine anderen Operatoren wie z.B. Theta Join als Subroutine verwendet werden. Eine Ausnahme bildet hier das Kreuzprodukt.

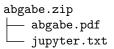
Für einen einfachen Einstieg haben wir Ihnen bereits einige Beispielanwendungen beider Klassen sowie einen Unit Test mit mehreren Testfällen im Notebook zur Verfügung gestellt, der evaluate von Co_Group_Join testet. Ihre finale Abgabe muss auf allen gültigen Eingaben das korrekte Ergebnis berechnen. Die Qualität Ihrer Implementierung wird bei der Punktevergabe berücksichtigt.



Abgabe

Lösungen sind in Teams von 2 bis 3 Studierenden bis zum 12. Mai 2022, 10:15 Uhr über Ihre persönlichen Statusseite im CMS einzureichen. Nutzen Sie hierfür die Team Groupings Funktionalität im CMS.

Ihre Abgabe muss dem folgenden Format entsprechen:



Hierbei enthält abgabe.pdf Ihre Lösungen zu Aufgabe 1 und 2 und jupyter.txt Ihre Lösung zu Aufgabe 3. Achten Sie darauf, dass Sie nur die von Ihnen zu ergänzenden Jupyter Zellen so kopieren, dass Einrückung und Formatierung korrekt sind.

Abgaben, die nicht den oben angegeben Vorgaben entsprechen, führen zu Punktabzug. Einzelabgaben werden nicht mehr korrigiert und mit 0 Punkten bewertet.