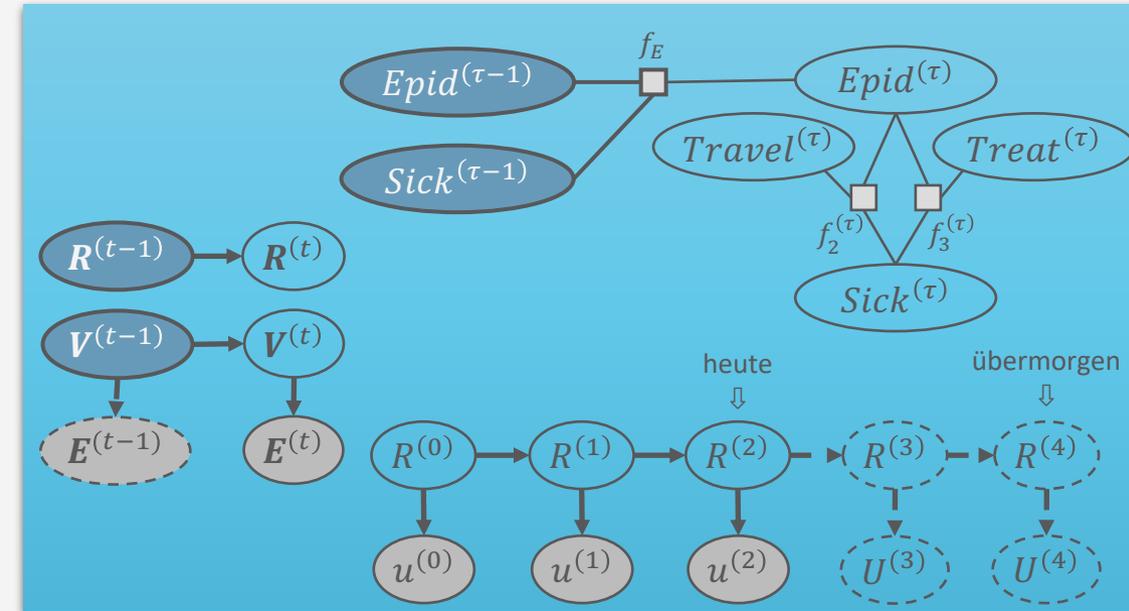


Sequentielle PGMs und Inferenz

Einführung in die
Künstliche Intelligenz



Inhalte

1. Künstliche Intelligenz & Agenten

- Agentenabstraktion, Rationalität
- Aufgabenumgebung

2. Episodische PGMs

- Gerichtetes Modell: Bayes Netze (BNs)
- Ungerichtete Modelle

3. Exakte Inferenz in episodischen PGMs

- Wahrscheinlichkeits- und Zustandsanfragen
- Direkt auf den Modellen, mittels Hilfsstrukturen

4. Approximative Inferenz in episodischen PGMs

- Wahrscheinlichkeitsanfragen
- Deterministische, stochastische Algorithmen

5. Lernalgorithmen für episodische PGMs

- Bei (nicht) vollständigen Daten, (un)bekannter Struktur

6. Sequentielle PGMs und Inferenz

- Dynamische BNs, Hidden-Markov-Modelle
- filtering / prediction / hindsight Anfragen, wahrscheinlichste Zustandssequenz
- Exakter, approximativer Algorithmus

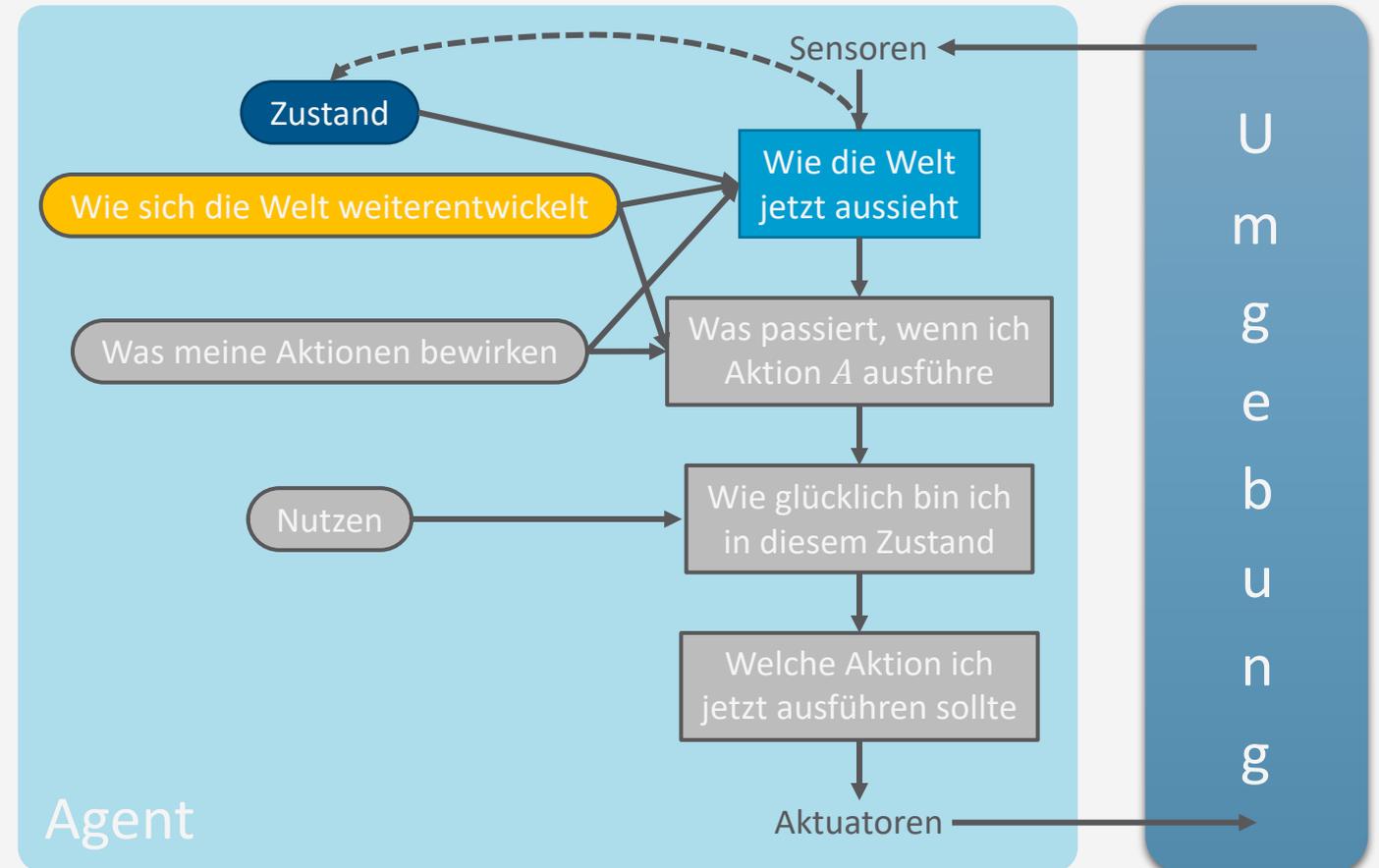
7. Entscheidungstheoretische PGMs

- Präferenzen, Nutzenprinzip
- PGMs mit Entscheidungs- und Nutzenknoten
- Berechnung der besten Aktion (Aktionssequenz)

8. Abschlussbetrachtungen

Einordnung der Vorlesung: *Modell- und nutzenbasierter Agent*

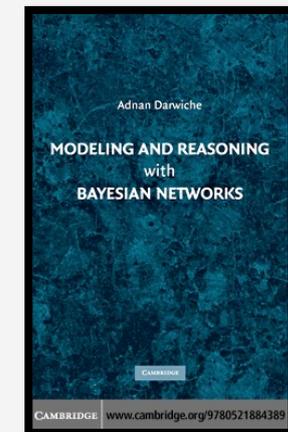
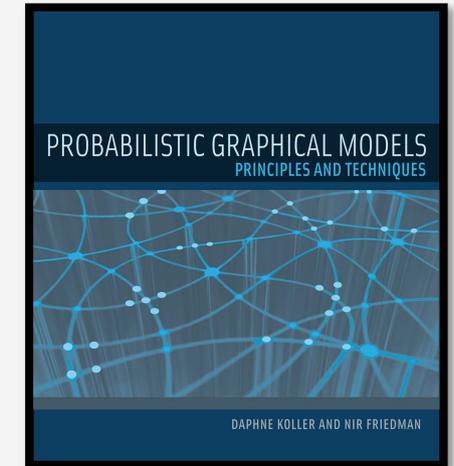
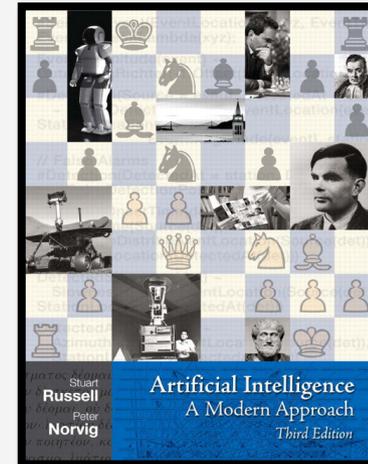
- Nachfolgende Themen der Vorlesung
 2. Episodische PGMs
 3. Exakte Inferenz in episodischen PGMs
 4. Approximative Inferenz in episodischen PGMs
 5. Lernalgorithmen für episodische PGMs
 6. **Sequentielle PGMs und Inferenz**
 7. Entscheidungstheoretische PGMs



Literaturhinweise

Inhalte dieses Themenblocks werden in den folgenden Kapiteln der Vorlesungsbücher behandelt

- AIMA(de)
 - Kap. 15.2: Inferenz in temporalen Modellen
 - Kap. 15.3: Hidden-Markov-Modelle
 - Kap. 15.5: Dynamische Bayes Netze
- PGM
 - Kap. 6.2: Zeitliche Modelle
 - Kap. 15.2: Inferenz in zeitlichen Modellen
- Wer gerne ein anderes Buch ausprobieren möchte (Fokus auf BNs):
 - Adnan Darwiche, *Modelling and Reasoning with Bayesian Networks*, 2009.



Bemerkung zu Namenskonventionen

- Gebräuchliche Namen für die gleiche Sache in PGMs
 - **Dynamisch (*dynamic*)**
 - ABER: *stationär* bezogen darauf, wie sich ein System verändert von einem Zustand zum nächsten
 - **Zeitlich (*temporal*)**
 - Änderungen zwischen Zuständen impliziert oft, dass die Zeit fortschreitet, i.e., eine zeitliche Zustandssequenz
 - Implizite Richtung der Kanten in Richtung der Zukunft
 - Vereinfachende Annahme: Diskrete Zeitschritte indiziert durch Integer (t)
 - **Sequentiell (*sequential*)**
 - Verallgemeinerung der Bezeichnung „zeitlich“, da eine Sequenz nicht nur durch fortschreitende Zeit, sondern auch durch etwas anderes entstehen kann
 - Beispiele: Räumliche Bewegung, Sequenz von Worten im Text
 - Meistens schreitet implizit dabei auch die Zeit voran

Überblick: 6. Sequentielle PGMs und Inferenz

A. *Sequentielle PGMs*

- Templates, dynamische BNs, dynamische Faktormodelle, Hidden Markov Modelle; Semantik
- Inferenzaufgaben: Wahrscheinlichkeitsanfragen (Filtering, Prediction, Hindsight), Zustandsanfragen (MPE, MAP)

B. *Sequentielle Inferenz*

- Naïve Inferenz mittels Ausrollen, Interface Algorithmus, Komplexität, Approximationen

C. *Spezialfall Hidden-Markov-Modelle*

- Viterbi-Algorithmus für MPEs
- Anfragebeantwortung durch Matrixoperationen
- Baum-Welch-Algorithmus zum Lernen

Zustand eines Systems über die Zeit

- Menge von Zufallsvariablen, um Zustand eines Systems zu beschreiben
 - $R = \{R_1, \dots, R_n\}$
 - **Template-Variablen** (Schablone, Vorlage)
- Zustand des Systems zu einem Zeitpunkt t :
 - Charakterisiert durch Belegungen von R zum Zeitpunkt t
 - Instanziierung der Template-Variablen mit t , notiert als $R^{(t)} = \{R_1^{(t)}, \dots, R_n^{(t)}\}$
 - Belegung mit Werten aus jeweils $\text{Val}(R_i)$: $r^{(t)} = \{r_1^{(t)}, \dots, r_n^{(t)}\}, r_i^{(t)} \in \text{Val}(R_i)$
- Zustand des Systems über ein diskretes Intervall $[t_1, t_2], t_1 < t_2$
 - Menge der Variablen: $R^{(t_1:t_2)} = \{R^{(t)} \mid t \in [t_1, t_2]\}$
 - Zustands-„Sequenz“: $r^{(t_1:t_2)} = \{r^{(t)} \mid t \in [t_1, t_2]\}$
 - Als Sequenz: $\boldsymbol{r}^{(t_1:t_2)} = (\boldsymbol{r}^{(t)} \mid t \in [t_1, t_2])$

Zustand eines Systems über die Zeit

- Gegeben eine Menge von Template-Variablen $\mathbf{R} = \{R_1, \dots, R_n\}$ und ein Endpunkt T
 - Menge der Variablen: $\mathbf{R}^{(0:T)}$
- Eine mögliche Welt: Belegung $\mathbf{r}^{(0:T)}$ von $\mathbf{R}^{(0:T)}$ mit Wahrscheinlichkeit belegt
 - Genannt **Trajektorie**
- Vollständige gemeinsame Wahrscheinlichkeitsverteilung über alle möglichen Trajektorien

$$P_{\mathbf{R}}^T = P(\mathbf{R}^{(0:T)}) = P(\mathbf{R}^{(0)}, \dots, \mathbf{R}^{(T)}) = P(R_1^{(0)}, \dots, R_n^{(0)}, \dots, R_1^{(T)}, \dots, R_n^{(T)})$$

- Kettenregel angewendet:

$$P_{\mathbf{R}}^T = P(\mathbf{R}^{(0:T)}) = P(\mathbf{R}^{(0)}) \prod_{t=1}^T P(\mathbf{R}^{(t)} | \mathbf{R}^{(0:(t-1))})$$

- Sehr komplexe Verteilung ohne weitere vereinfachende Annahmen

Vereinfachende Annahme 1: Markov Annahme

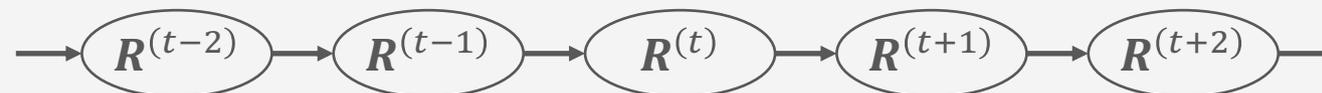
- **Markov Annahme:** Nächster Zustand hängt nur vom jetzigen Zustand ab

- Formal: $(\mathbf{R}^{(t+1)} \perp \mathbf{R}^{(0:(t-1))} | \mathbf{R}^{(t)})$

- Auswirkung auf $P(\mathbf{R}^{(0:T)})$: Faktorisierung über t

$$P_{\mathbf{R}}^T = P(\mathbf{R}^{(0:T)}) = P(\mathbf{R}^{(0)}) \prod_{t=1}^T P(\mathbf{R}^{(t)} | \mathbf{R}^{(t-1)})$$

- Auch *Markov Prozess erster Ordnung* genannt
- Aufteilung des Modells in **Zeitscheiben** (*time slices*)
- Graphische Darstellung:

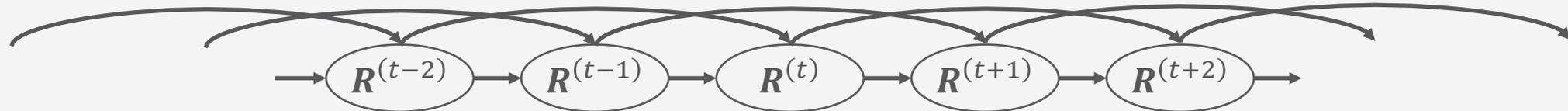


Vereinfachende Annahme 1: Markov Annahme

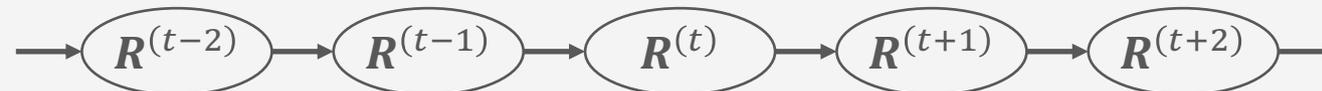
- Verallgemeinerung der Markov-Annahme: **Markov- k**
 - Der nächste Zustand hängt nur von den k vorherigen Zuständen ab

$$P_{\mathbf{R}}^T = P(\mathbf{R}^{(0:T)}) = P(\mathbf{R}^{(0)}) \prod_{t=1}^k P(\mathbf{R}^{(t)} | \mathbf{R}^{(0:(t-1))}) \prod_{t=k+1}^T P(\mathbf{R}^{(t)} | \mathbf{R}^{((t-k):(t-1))})$$

- Graphische Darstellung für $k = 2$, ergibt *Markov Prozess zweiter Ordnung*



- Markov-Annahme \triangleq Markov- k mit $k = 1$



Vereinfachende Annahme 2: Stationäres System

- Gegeben eine Menge von Template-Variablen $\mathbf{R} = \{R_1, \dots, R_n\}$ und ein Endpunkt T
- Vollständige gemeinsame Wahrscheinlichkeitsverteilung über alle möglichen Trajektorien
 - Mit Markov Annahme: Faktorisierung über t

$$P_{\mathbf{R}}^T = P(\mathbf{R}^{(0:T)}) = P(\mathbf{R}^{(0)}) \prod_{t=1}^T P(\mathbf{R}^{(t)} | \mathbf{R}^{(t-1)})$$

- Annahme **Stationarität**: $P(\mathbf{R}^{(t)} | \mathbf{R}^{(t-1)})$ identisch für jedes t
 - Formal: Es gibt ein **Übergangsmodell** (*transition model*) $P(\mathbf{R}' | \mathbf{R})$, so dass für jedes $t \geq 1$ gilt

$$P(\mathbf{R}^{(t)} = \mathbf{r}' | \mathbf{R}^{(t-1)} = \mathbf{r}) = P(\mathbf{R}' = \mathbf{r}' | \mathbf{R} = \mathbf{r})$$

Für ein stationäres System mit Markov Annahme reicht es, die folgenden Verteilungen zu spezifizieren:

$$P(\mathbf{R}^{(0)}), P(\mathbf{R}^{(\tau)} | \mathbf{R}^{(\tau-1)})$$

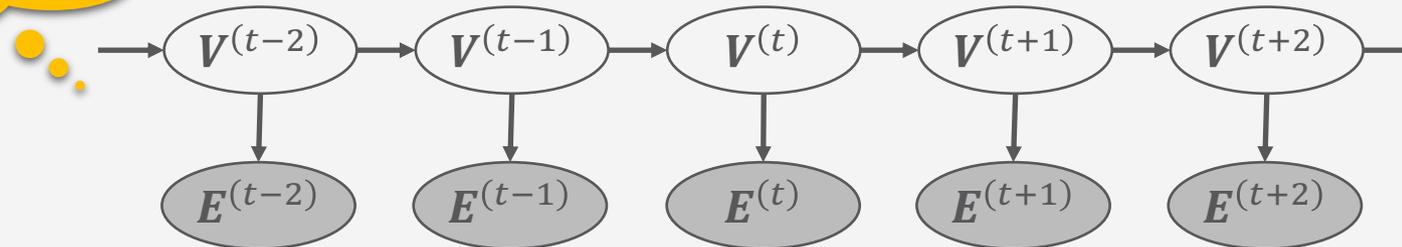
Was ist das Problem hiermit?

Beobachtbare und nicht beobachtbare Variablen

- Menge der Zufallsvariablen \mathbf{R} lässt sich in der Regel in Evidenzvariablen \mathbf{E} (jeden Zeitschritt mit Evidenz belegt) und latente Variablen \mathbf{V} aufteilen
- Verteilung ist dann, wie folgt, faktorisiert:

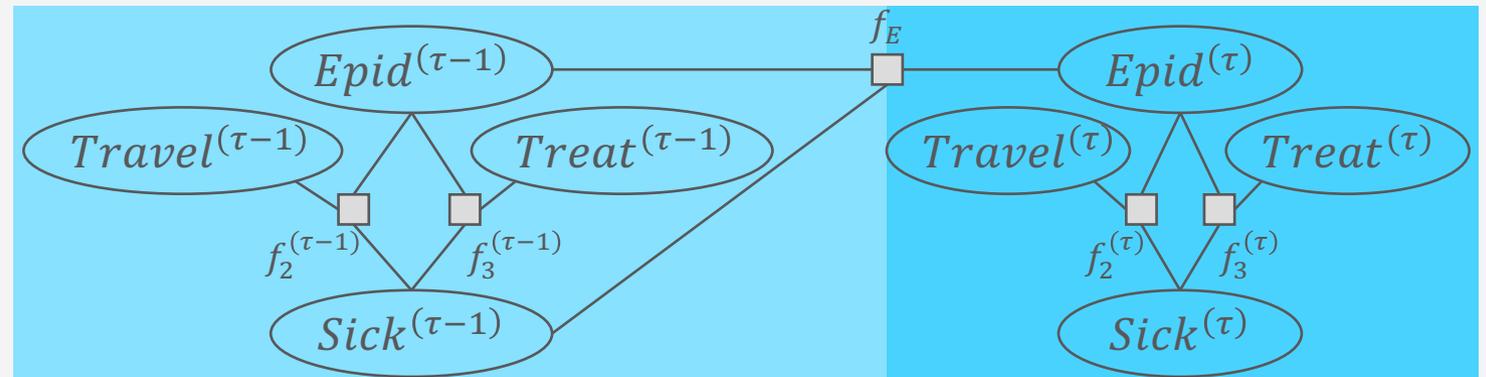
$$P_{\mathbf{R}}^T = P(\mathbf{V}^{(0:T)}, \mathbf{E}^{(0:T)}) = P(\mathbf{V}^{(0)}) \prod_{t=1}^T \underbrace{P(\mathbf{V}^{(t)} | \mathbf{V}^{(t-1)})}_{\text{Übergangmodell}} \underbrace{P(\mathbf{E}^{(t)} | \mathbf{V}^{(t)})}_{\text{Sensormodell}}$$

Was für Unabhängigkeiten herrschen hier?



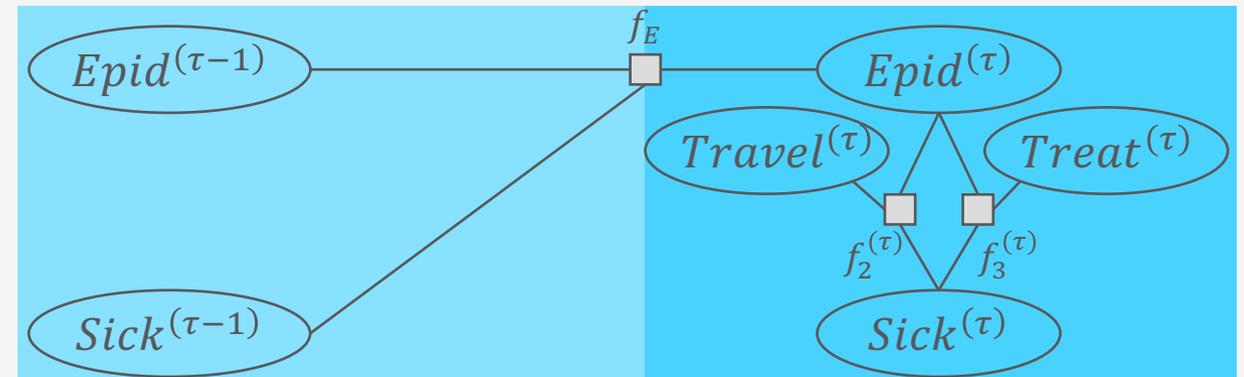
Allgemeine Faktorisierung für sequentielle Modelle

- Übergang zum nächsten Zeitschritt ausgelöst durch neue Beobachtungen / Ausführung von Aktionen
 - Markov Annahme: Effekt eines Aktion spätestens im nächsten Zeitschritt
- Modellierung
 - *Intra-Zeitscheiben*-Faktoren/CPDs: Beschreiben Vorgänge / Effekte innerhalb eines Zeitschrittes
 - Argumente nur aus Zeitschritt τ
 - Z.B. $f_2: Epid^{(\tau)}, Sick^{(\tau)}, Travel^{(\tau)}$
 - *Inter-Zeitscheiben*-Faktoren/CPDs: Beschreiben zeitliche Übergänge bzw. Vorgänge mit Effekten im nächsten Zeitschritt
 - Argumente aus $\tau - 1, \tau$
 - Z.B. $f_E: Epid^{(\tau-1)}, Sick^{(\tau-1)}, Epid^{(\tau)}$



2-Zeitscheiben-Modelle & 1.5-Zeitscheiben-Modelle

- **2-Zeitscheiben-Faktormodell** (eigentlich nur noch ein 1.5-Zeitscheiben-Faktormodell)
 - Spezifizierung eines Übergangsmodells $P(\mathbf{R}^{(\tau)} | \mathbf{R}^{(\tau-1)})$ gemäß Faktorisierung
 - Partielle Spezifizierung eines sequentiellen Modells (ohne $P(\mathbf{R}^{(0)})$)
 - Inter-Zeitscheiben-Effekte und Intra-Zeitscheiben-Effekte
 - Ausreichend: Intra-Zeitscheiben-Faktoren für τ , Inter-Zeitscheiben-Faktoren von $\tau - 1$ zu τ
- Formal: $F^{\rightarrow} = \left\{ f_i^{(\tau)} \right\}_{i=1}^n \cup \left\{ f_j \right\}_{j=1}^m$ mit
 - Intra-Zeitscheiben-Faktoren
 $f_i^{(\tau)} = \phi_i^{(\tau)} \left(R_1^{(\tau)}, \dots, R_{k_i}^{(\tau)} \right)$
 - Inter-Zeitscheiben-Faktoren
 $f_j = \phi_j \left(R_1^{(t)}, \dots, R_{k_j}^{(t)} \right), t \in \{\tau, \tau - 1\}$

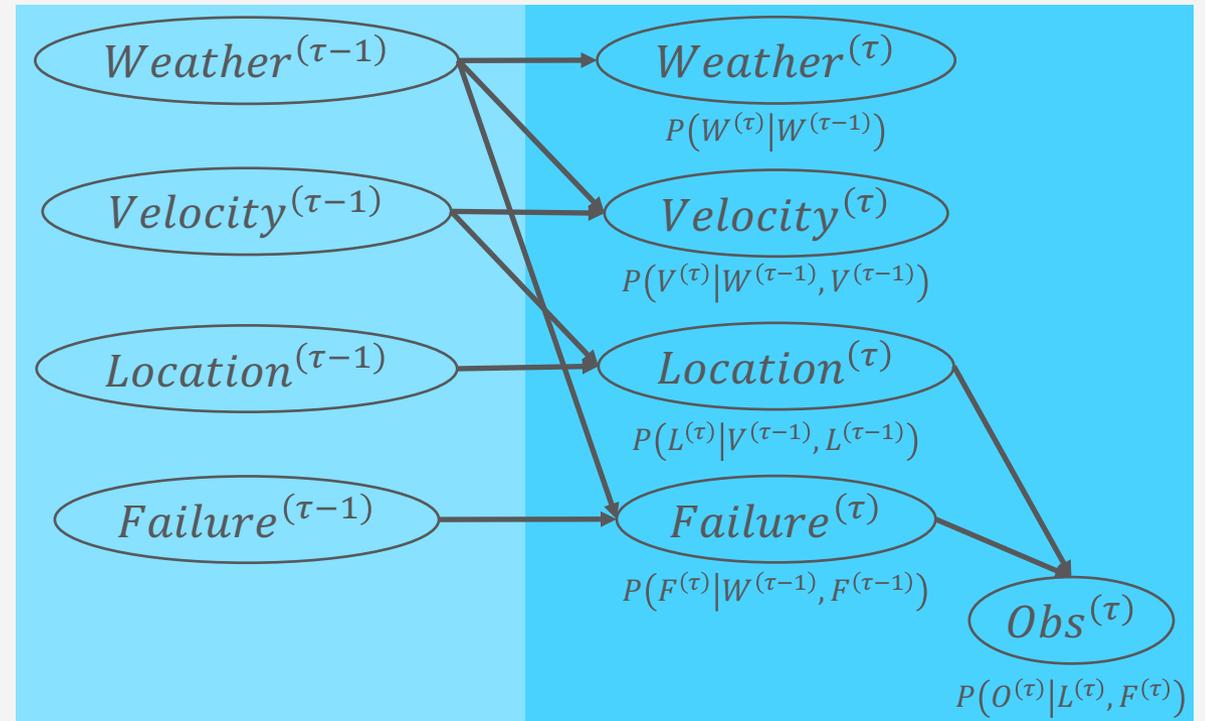


Was für Unabhängigkeiten herrschen hier?

2-Zeitscheiben-Modelle & 1.5-Zeitscheiben-Modelle

- In derselben Art können wir 2-Zeitscheiben-BNs bauen
- **2-Zeitscheiben-BN** (eigentlich nur noch ein 1.5-Zeitscheiben-BN)
 - $B^{\rightarrow} = \left\{ P \left(R_i^{(\tau)} \mid \text{Pa} \left(R_i^{(\tau)} \right) \right) \right\}_{i=1}^n$ mit
 - Für alle $R_j^{(t)} \in \text{Pa} \left(R_i^{(\tau)} \right)$ gilt $t \in \{\tau, \tau - 1\}$

Im Gegensatz zu der normalen Interpretation von BNs mit einer CPD pro Knoten haben die Zufallsvariablen aus $\tau - 1$ in B^{\rightarrow} keine CPDs assoziiert.



Dynamische Modelle

- Annahmen: Markov-1, stationärer Prozess
- Übliche Definition: Dynamische Modelle sind Tupel (M^0, M^{\rightarrow}) , wobei
 - M^0 das Verhalten des ersten Zeitschritts (Intra-Zeitscheiben-Verhalten) beschreibt
 - M^{\rightarrow} ein 2-Zeitscheiben-Modell ist, welches Intra- und Inter-Zeitscheiben-Verhalten beschreibt
 - **Template-Modell**, das instanziiert wird, indem τ mit einem Zeitpunkt t ersetzt wird
- Häufiger Zusammenhang zwischen M^0, M^{\rightarrow}
 - Intra-Zeitscheiben-Verhalten aus M^{\rightarrow} zum Zeitpunkt τ gleich dem Verhalten des ersten Zeitschritts (Intra-Zeitscheiben-Verhalten aus M^{\rightarrow} bei $\tau = 0$)
 - Inter-Zeitscheiben-Verhalten gibt es in M^0 nicht, da es keinen vorhergehenden Schritt gibt
 - Dafür wird M^0 manchmal ergänzt mit Apriori-Informationen kodiert in Verteilungen

Dynamische Faktormodelle

- **Dynamisches Faktormodell** (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_i^{(0)}\}_{i=1}^{n_0}$, $f_i^{(0)} = \phi_i^{(0)}(R_1^{(0)}, \dots, R_{k_i}^{(0)})$

- F^{\rightarrow} ein 2-Zeitscheiben-Faktormodell: $F^{\rightarrow} = \{f_i^{(\tau)}\}_{i=1}^n \cup \{f_j\}_{j=1}^m$ mit

- $f_i^{(\tau)} = \phi_i^{(\tau)}(R_1^{(\tau)}, \dots, R_{k_i}^{(\tau)})$ (Intra-Zeitscheiben-Faktoren)

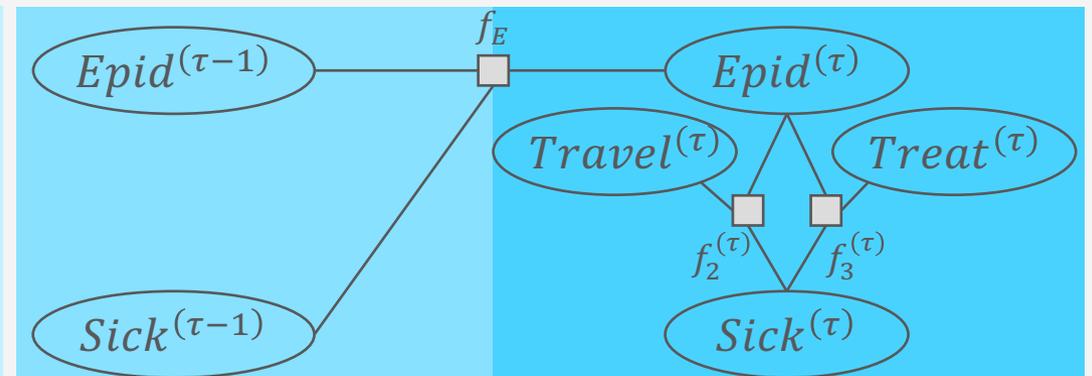
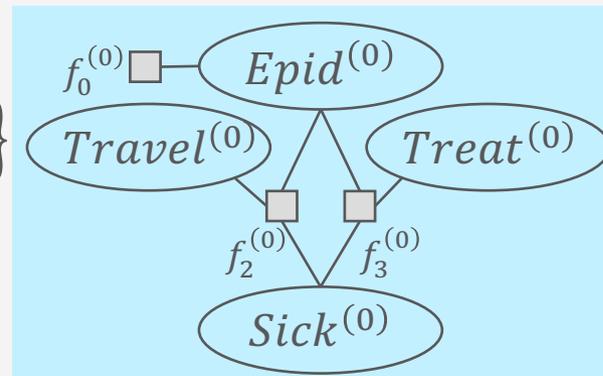
- $f_j = \phi_j(R_1^{(t)}, \dots, R_{k_j}^{(t)})$, $t \in \{\tau, \tau - 1\}$ (Inter-Zeitscheiben-Faktoren)

- **Beispiel:**

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$

- Keine Annahme zu Beobachtbarkeit



F_0 besteht i.d.R. aus den Intra-Zeitscheiben-Faktoren von F^{\rightarrow} mit $\tau = 0$, evtl. ergänzt um Quasi-Apriori-Verteilungen, i.e.,

$$\{f_i^{(0)}\}_{i=1}^{n_0} = \{f_i^{(\tau)}\}_{i=1| \tau=0}^n \cup \{f_i^{(0)}\}_{i=1}^{n'}$$

Dynamisches BN (DBN)

- Dynamisches BN (B^0, B^{\rightarrow})

- $B^0 =$

$$\left\{ P \left(R_i^{(0)} \mid \text{Pa} \left(R_i^{(0)} \right) \right) \right\}_{i=1}^{n_0}$$

- $\forall R_j^{(t)} \in \text{Pa} \left(R_i^{(0)} \right) :$
 - $t = 0$

- B^{\rightarrow} ein 2-Zeitscheiben-BN:

$$B^{\rightarrow} =$$

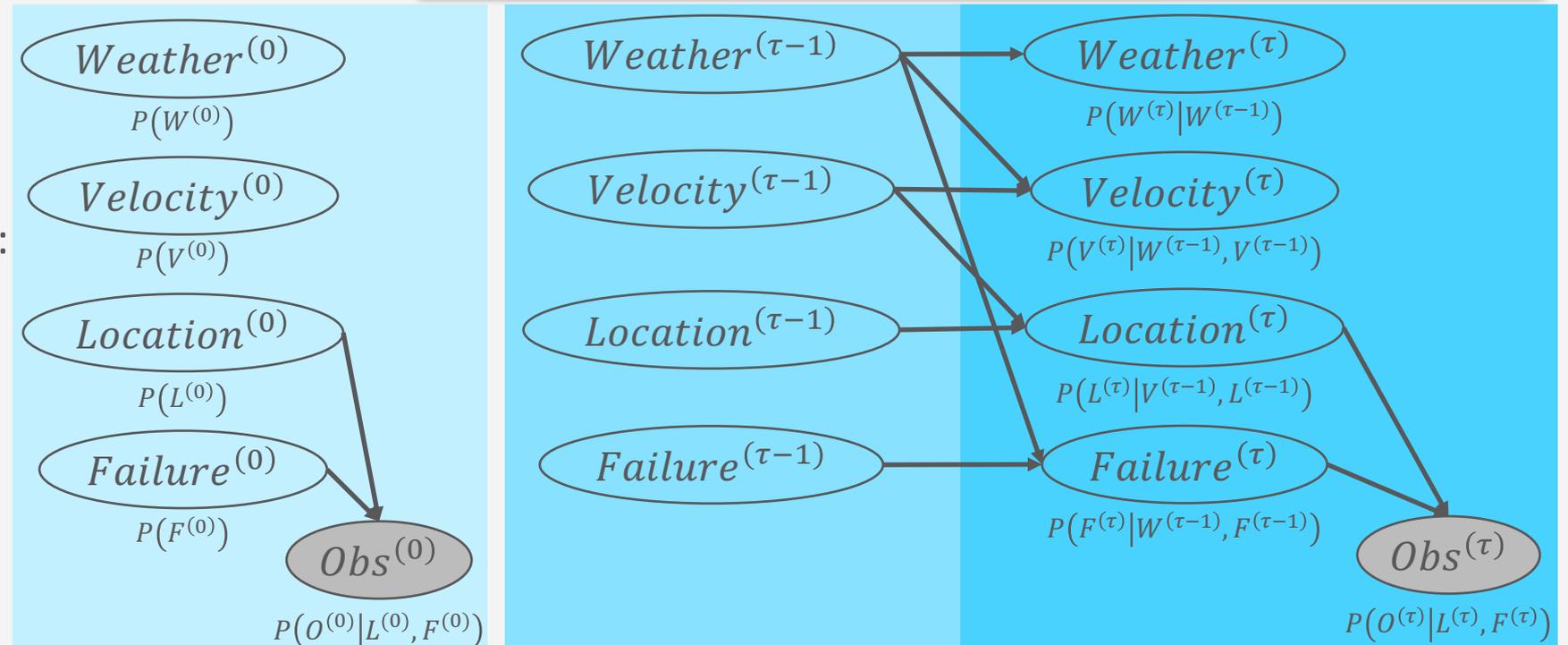
$$\left\{ P \left(R_i^{(\tau)} \mid \text{Pa} \left(R_i^{(\tau)} \right) \right) \right\}_{i=1}^n$$

- $\forall R_j^{(t)} \in \text{Pa} \left(R_i^{(\tau)} \right) :$
 - $t \in \{\tau, \tau - 1\}$

- Beispiel: Abbildung

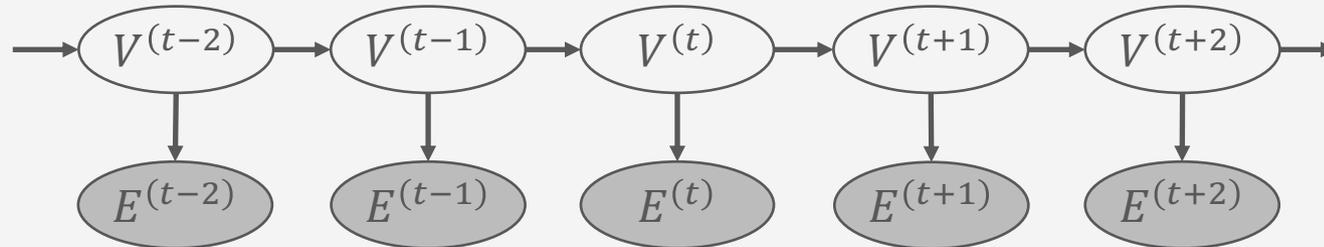
- Annahme: $\mathbf{E} = \{Obs\}$

Bei (B^0, B^{\rightarrow}) besteht B^0 i.d.R. aus Intra-Zeitscheiben-CPDs, die auch in B^{\rightarrow} gelten, und Apriori-Wahrscheinlichkeiten $P \left(R_{i'}^{(0)} \right)$ für die Wurzelknoten von B^0 , was dann die Knoten in B^{\rightarrow} mit Index $\tau - 1$ sind, welche keine CPD in B^{\rightarrow} haben.



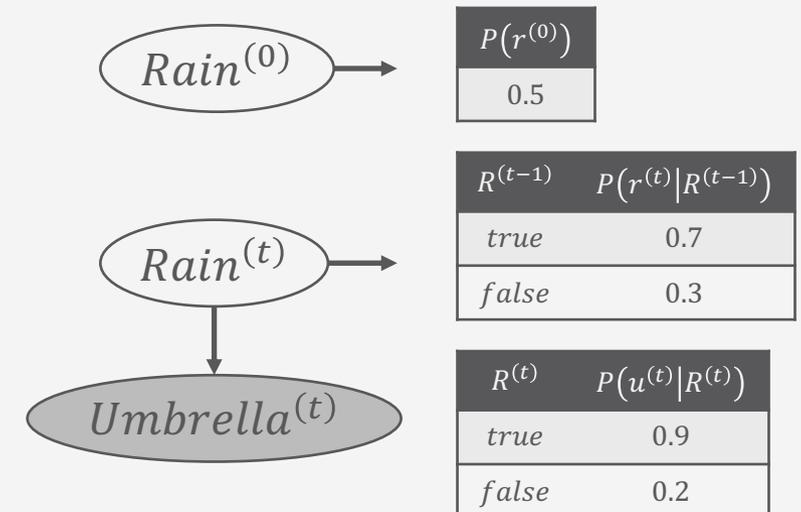
Hidden Markov Modell (HMM)

- Spezialfall des DBN mit $R = \{V, E\}$ bzw. $V = \{V\}, E = \{E\}$
 - Latente Zufallsvariable V
 - Evidenzvariable E
 - Modell:
 - $B^0 = P(V^{(0)})$
 - $B^\rightarrow = \{P(V^{(\tau)}|V^{(\tau-1)}), P(E^{(\tau)}|V^{(\tau)})\}$
 - Repräsentierte Verteilung: $P_{V,E}^T = P(V^{(0)}) \prod_{t=1}^T P(V^{(t)}|V^{(t-1)})P(E^{(t)}|V^{(t)})$



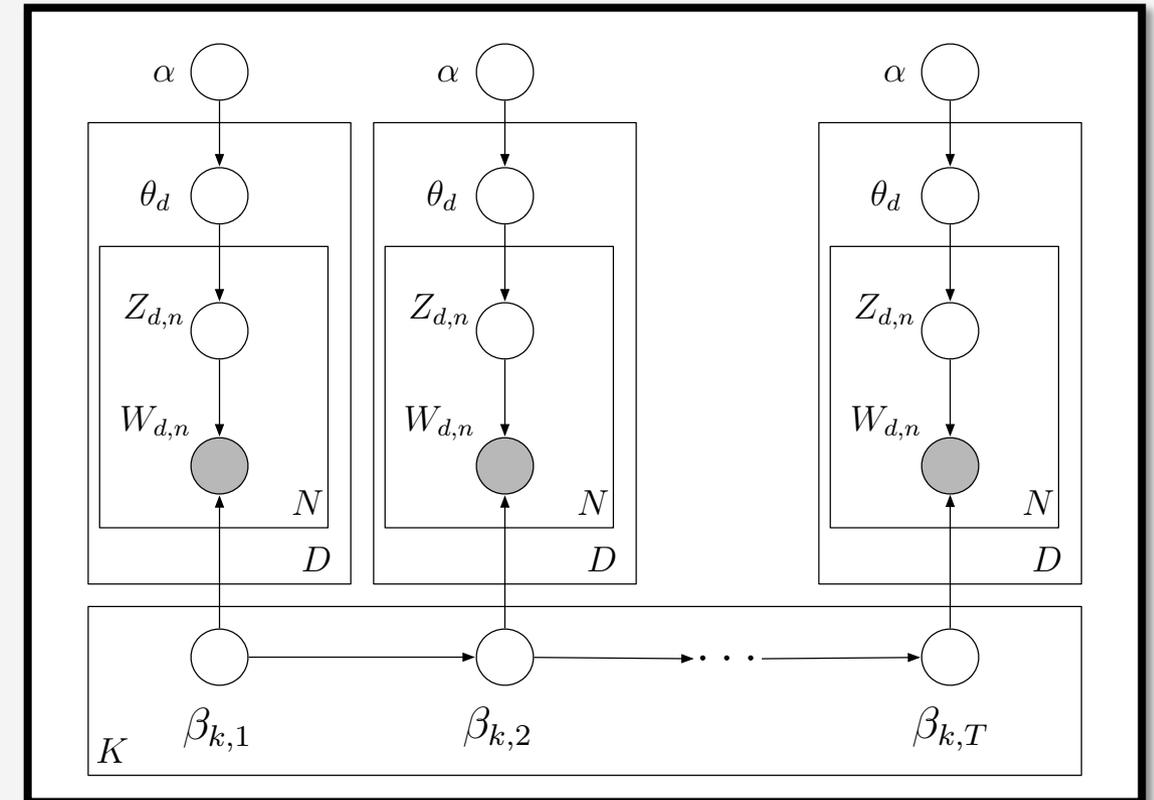
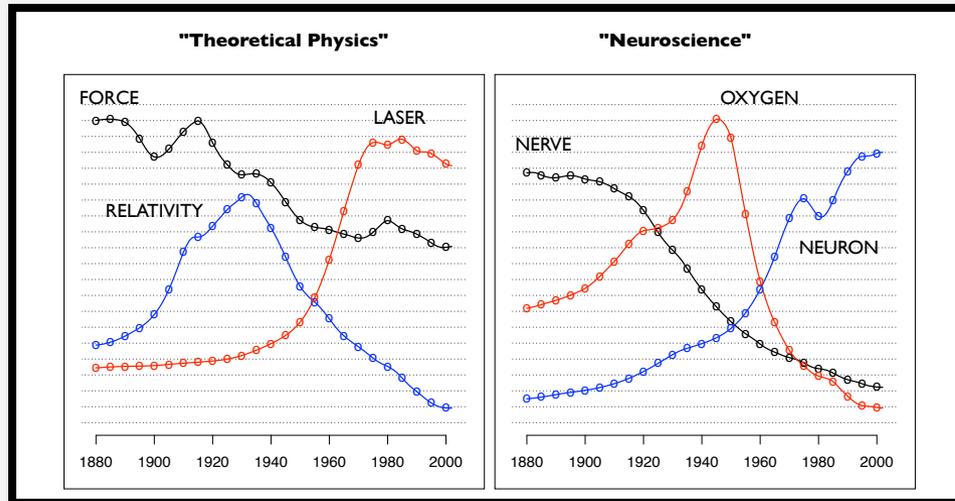
Ein DBN mit $R = V \cup E$ kann jederzeit in ein HMM umgewandelt werden:

- Mit jedem $v \in \text{Val}(V)$ als abstrakten Wert v , $\text{Val}(V) = \text{Val}(V)$
 - Mit jedem $e \in \text{Val}(E)$ als abstrakten Wert e , $\text{Val}(E) = \text{Val}(E)$
- $\text{Val}(V)$ und $\text{Val}(E)$ entsprechend groß!



Topic Modellierung über die Zeit: Dynamic Topic Modell

- Corpus in Zeitscheiben einteilen (z.B. pro Jahr)
- Annahme: Innerhalb einer Zeitscheibe, aus einem LDA Modell generiert
- Erlaube, dass sich Topic-Verteilungen von einer Zeitscheibe zur nächsten verändern



Semantik

- Semantik definiert über die vollständige gemeinsame Wahrscheinlichkeitsverteilung gegeben ein Zeitpunkt $T > 0$
 - Dynamisches Modell (M^0, M^{\rightarrow}) für T Zeitschritte **ausrollen**

$$M = M^{(0:T)} = M^0 \cup \bigcup_{t=1}^T M^{\rightarrow|\tau=t}$$

- $M^{\rightarrow|\tau=t}$: Ersetze τ mit t in M_{\rightarrow}
- Ergibt ein episodisches Modell mit bekannter Semantik, i.e., ein Modell, welches eine vollständige gemeinsame Wahrscheinlichkeitsverteilung als (normalisiertes) Produkt der Faktoren / CPDs repräsentiert

$$P_M^T = \frac{1}{Z} \prod_{\zeta \in M} \zeta$$

- ζ ein Faktor, wenn M^0, M^{\rightarrow} Faktormodelle, oder eine CPD, wenn M^0, M^{\rightarrow} BNs (dann $Z = 1$)

Ausrollen: Beispiel – $T = 2$

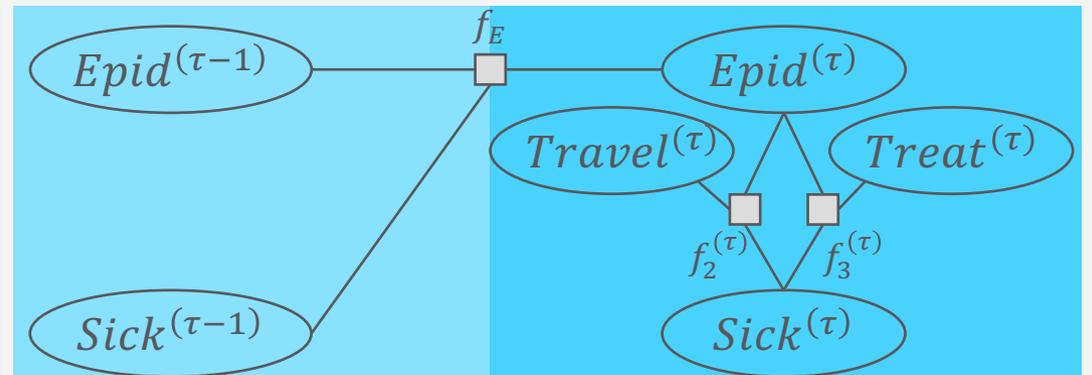
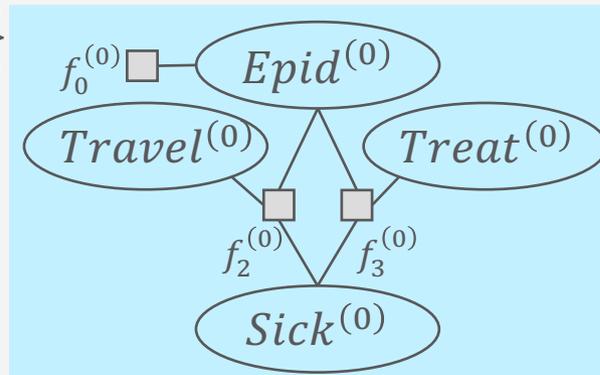
- Ausgerolltes Modell

$$F = F^0 \cup \bigcup_{t=1}^2 F^{\rightarrow|\tau=t}$$

- Dynamisches Faktormodell (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$



Ausrollen: Beispiel – $T = 2, t = 0$

- Ausgerolltes Modell

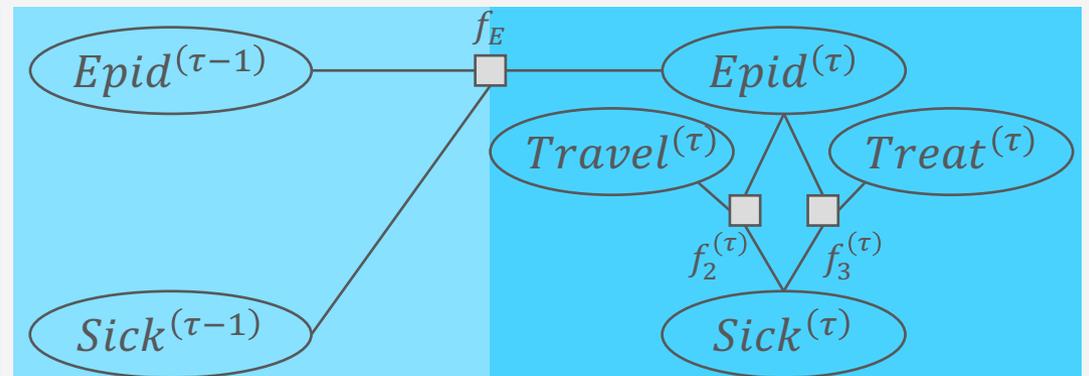
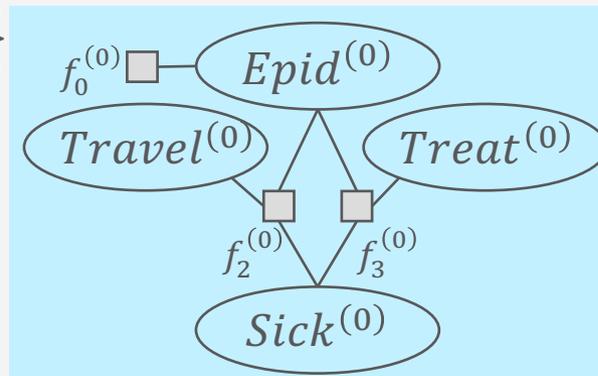
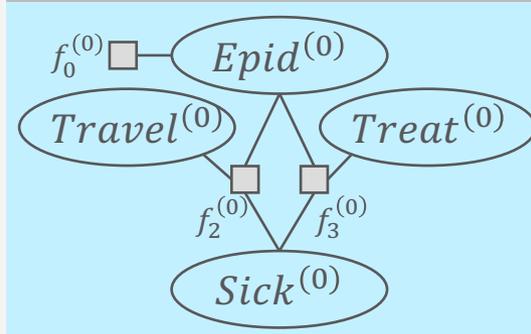
$$F = F^0 \cup \bigcup_{t=1}^2 F^{\rightarrow|\tau=t}$$

- Dynamisches Faktormodell (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$

$$F = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\} \cup \dots$$



Ausrollen: Beispiel – $T = 2, t = 1$

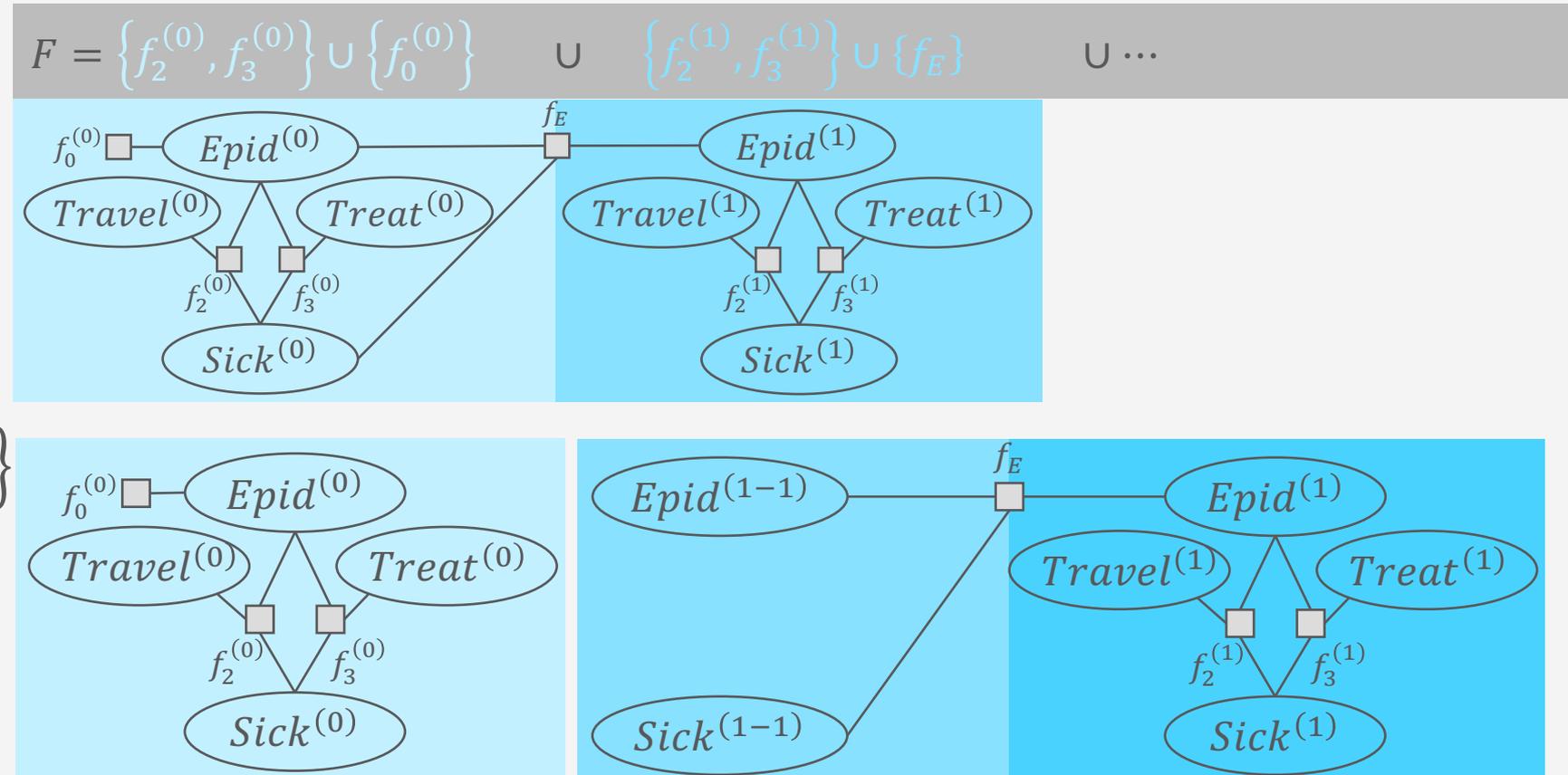
- Ausgerolltes Modell

$$F = F^0 \cup \bigcup_{t=1}^2 F^{\rightarrow|\tau=t}$$

- Dynamisches Faktormodell (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$



Ausrollen: Beispiel – $T = 2, t = 2$

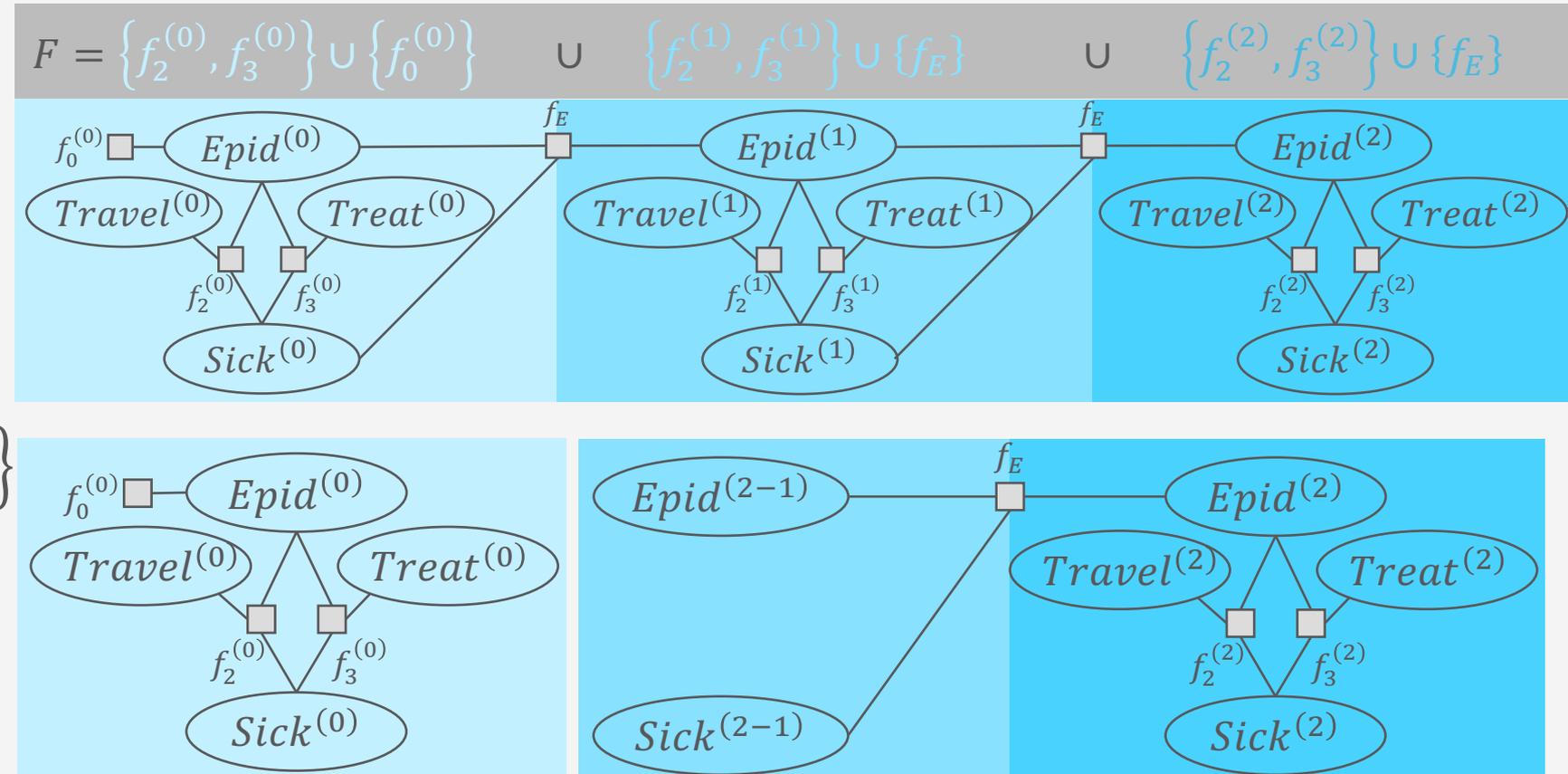
- Ausgerolltes Modell

$$F = F^0 \cup \bigcup_{t=1}^2 F^{\rightarrow|\tau=t}$$

- Dynamisches Faktormodell (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$



Ein Blick auf die
Unabhängigkeiten

Ausrollen: Beispiel – $T = 2$, Semantik

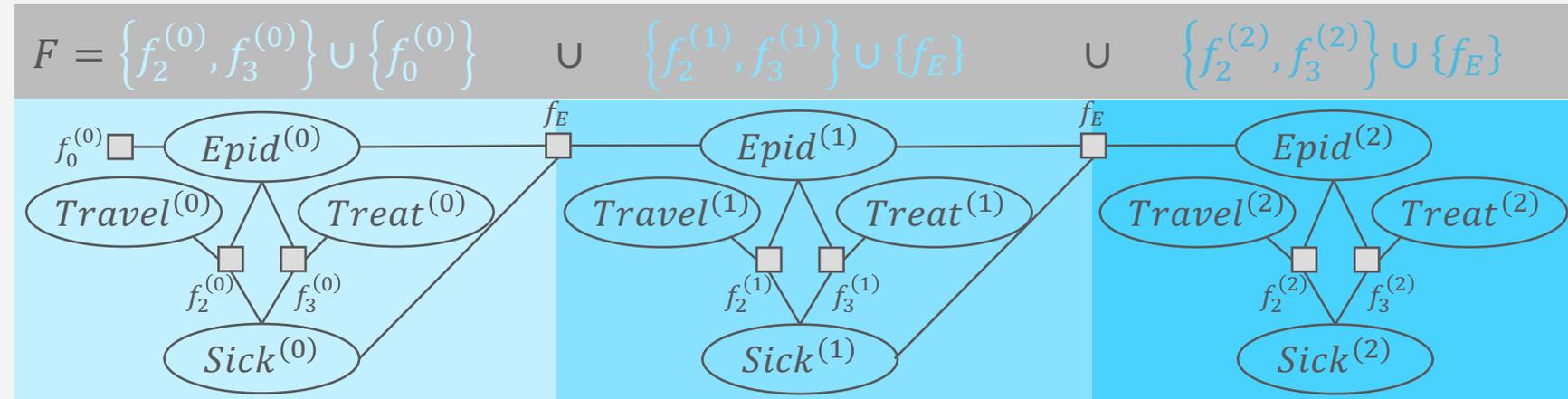
- Ausgerolltes Modell

$$F = F^0 \cup \bigcup_{t=1}^2 F^{\rightarrow|\tau=t}$$

- Dynamisches Faktormodell (F^0, F^{\rightarrow}) mit

- $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}$

- $F^{\rightarrow} = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$



$$\begin{aligned}
 P_F^T &= \frac{1}{Z} \prod_{f \in F} f = \frac{1}{Z} \phi_2^{(0)}(E^{(0)}, S^{(0)}, Tl^{(0)}) \cdot \phi_3^{(0)}(E^{(0)}, S^{(0)}, Tt^{(0)}) \cdot \phi_0^{(0)}(E^{(0)}) \\
 &\quad \cdot \phi_2^{(1)}(E^{(1)}, S^{(1)}, Tl^{(1)}) \cdot \phi_3^{(1)}(E^{(1)}, S^{(1)}, Tt^{(1)}) \cdot \phi_E(E^{(1)}) \\
 &\quad \cdot \phi_2^{(2)}(E^{(2)}, S^{(2)}, Tl^{(2)}) \cdot \phi_3^{(2)}(E^{(2)}, S^{(2)}, Tt^{(2)}) \cdot \phi_E(E^{(2)})
 \end{aligned}$$

Interfaces

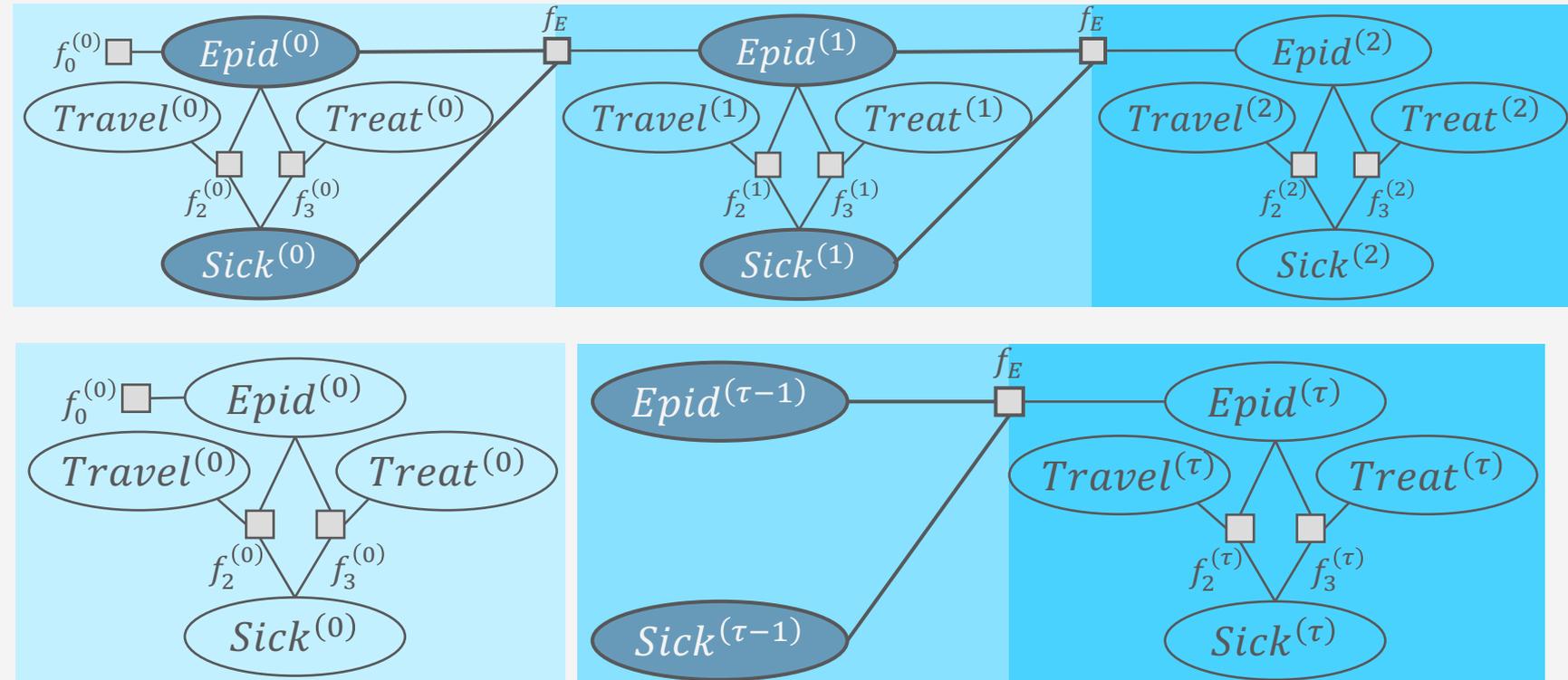
- Interface: Separierende Untermenge an Zufallsvariablen, durch die die Zeitscheiben voneinander unabhängig sind

- Beispiel:

- $Epid^{(\tau-1)}$

- $Sick^{(\tau-1)}$

- D.h., Variablen der Inter-Zeitscheiben-Faktoren f_i mit $\tau - 1$



Wie sieht das bei DBNs aus? DBN ausrollen: Beispiel – $T = 2$

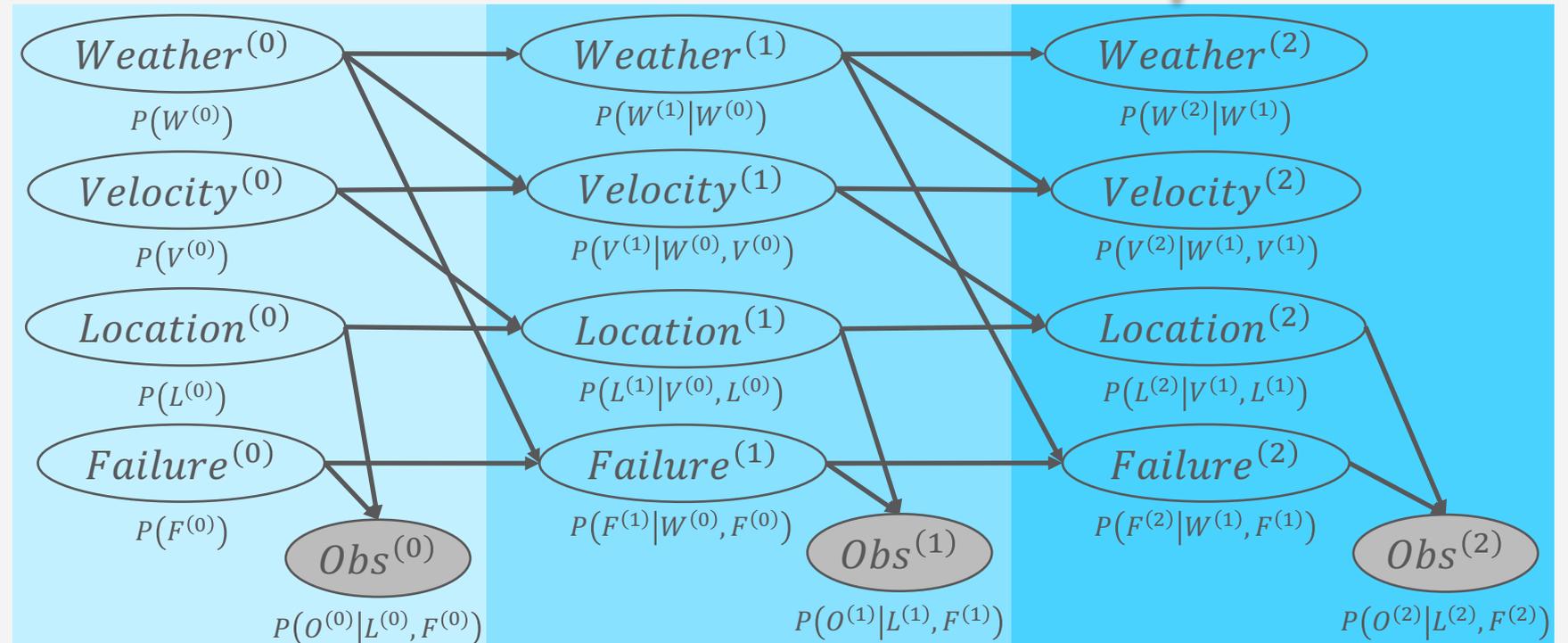
- Ausgerolltes Modell

$$B = B^0 \cup \bigcup_{t=1}^2 B^{\rightarrow|\tau=t}$$

- DBN (B^0, B^{\rightarrow}) mit

- $B^0 = \{P(O^{(0)}|L^{(0)}, F^{(0)})\} \cup \{P(W^{(0)}), P(V^{(0)}), P(L^{(0)}), P(F^{(0)})\}$

- $B^{\rightarrow} = \{P(O^{(\tau)}|L^{(\tau)}, F^{(\tau)})\} \cup \{P(W^{(\tau)}|W^{(\tau-1)}), P(V^{(\tau)}|W^{(\tau-1)}, V^{(\tau-1)}), P(L^{(\tau)}|V^{(\tau-1)}, L^{(\tau-1)}), P(F^{(\tau)}|W^{(\tau-1)}, F^{(\tau-1)})\}$



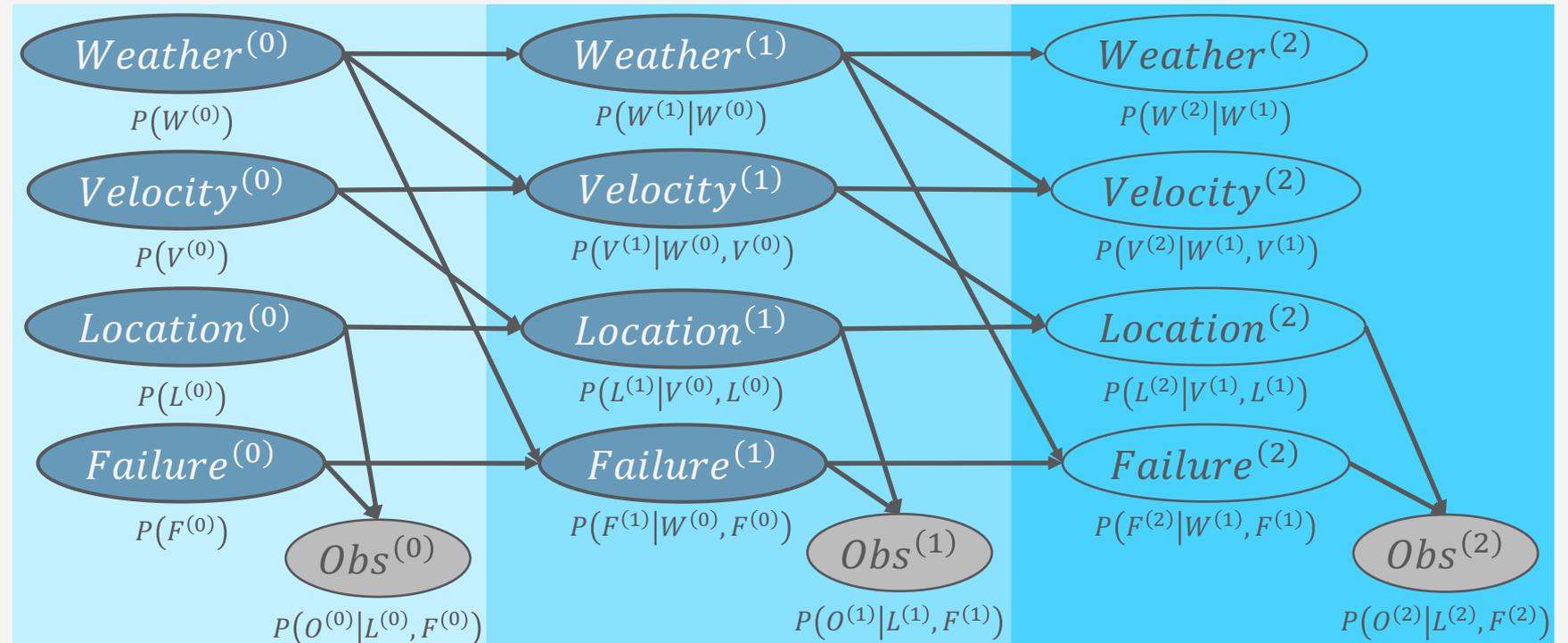
Interfaces

- Interface: Separierende Untermenge an Zufallsvariablen, durch die die Zeitscheiben voneinander unabhängig sind

- Beispiel:

- $Weather^{(\tau-1)}$
- $Velocity^{(\tau-1)}$
- $Location^{(\tau-1)}$
- $Failure^{(\tau-1)}$

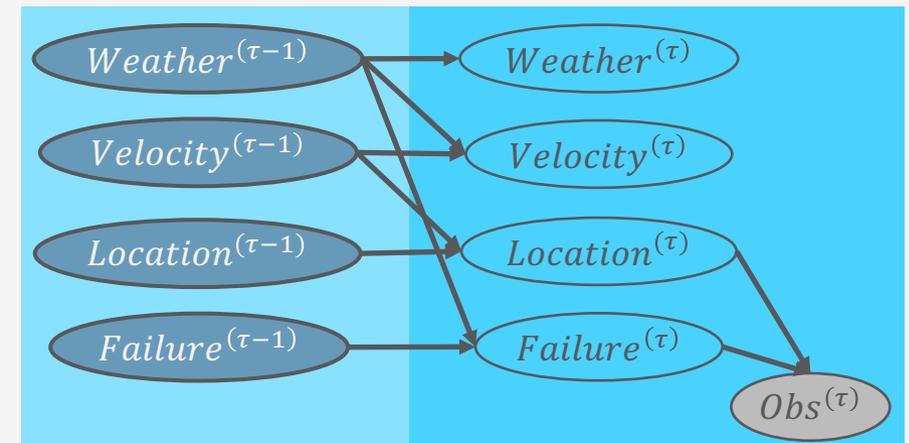
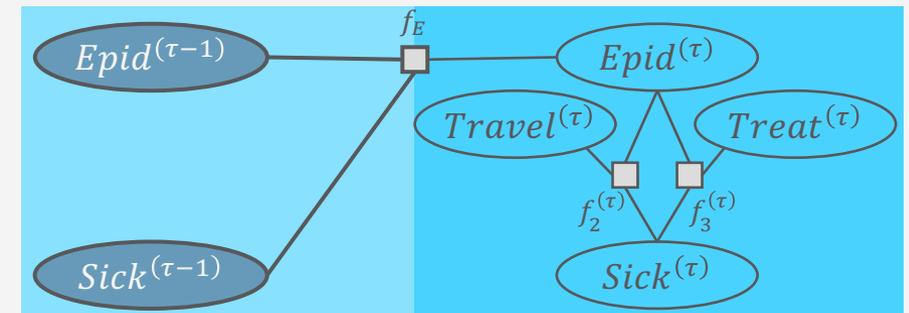
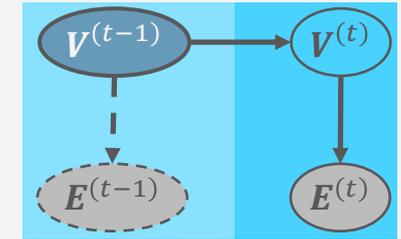
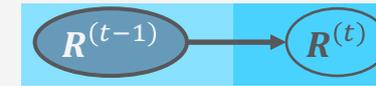
- D.h., Variablen des Übergangsmodells mit $\tau - 1$



Interfaces

- **Vorwärts-Interface $I^{(\tau-1)}$** : Separatoren der Zeitscheiben bei fortschreitender Zeit (von $M^{(\tau-1)}$ nach $M^{(\tau)}$)
 - Triviale Separatoren, in der Regel nicht minimal
 - Maximales Interface: R
 - Bei $R = V \cup E$: maximales Interface V
 - Dynamisches Faktormodell: Menge der Variablen mit Index $\tau - 1$ in F^{\rightarrow}
 - $I = \{R^{(\tau-1)} \mid R^{(\tau-1)} \in \text{rv}(F^{\rightarrow})\}$
 - Kommen nur in Inter-Zeitscheiben-Faktoren vor
 - DBN: Menge der Variablen mit Index $\tau - 1$ in B^{\rightarrow}
 - $I = \{R^{(\tau-1)} \mid R^{(\tau-1)} \in \text{rv}(B^{\rightarrow})\}$
 - Sind Elternknoten von τ -Variablen in B^{\rightarrow}

Sequentielle PGMs & Inferenz



I trennt auch von $M^{(\tau)}$ nach $M^{(\tau-1)}$; es gibt aber auch Rückwärts-Interfaces, die aus der Richtung τ nach $\tau - 1$ bestimmt werden [Näheres dazu in Kevin Murphys Dissertation, 2002]

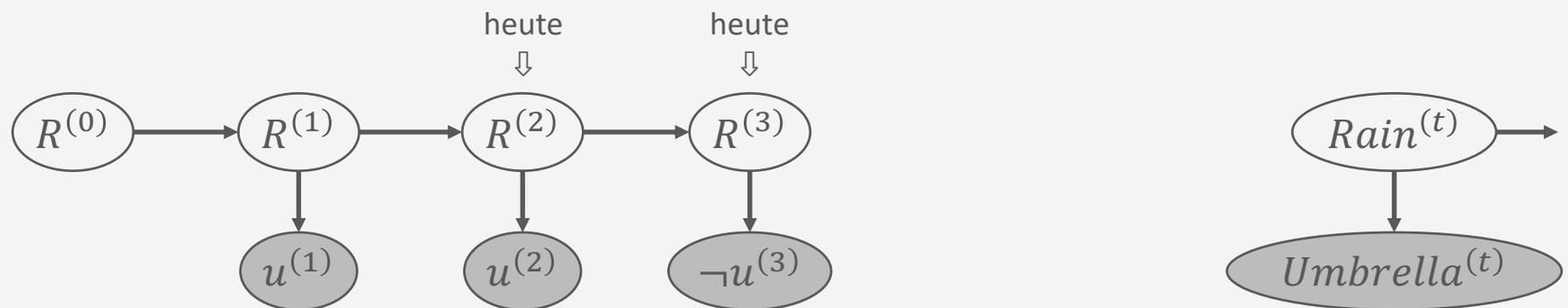
Anfragebeantwortungsproblem & Anfragetypen

- Filtern / überwachen (*filtering*)
 - Was ist die Wahrscheinlichkeit, dass es heute ($t = 2$) regnet gegeben alle Beobachtungen zu Regenschirmen bis einschließlich heute?

$$P(R^{(2)} | u^{(1)}, u^{(2)})$$

- Nächster Tag: Was ist die Wahrscheinlichkeit, dass es heute ($t = 3$) regnet gegeben alle Beobachtungen zu Regenschirmen bis einschließlich heute?

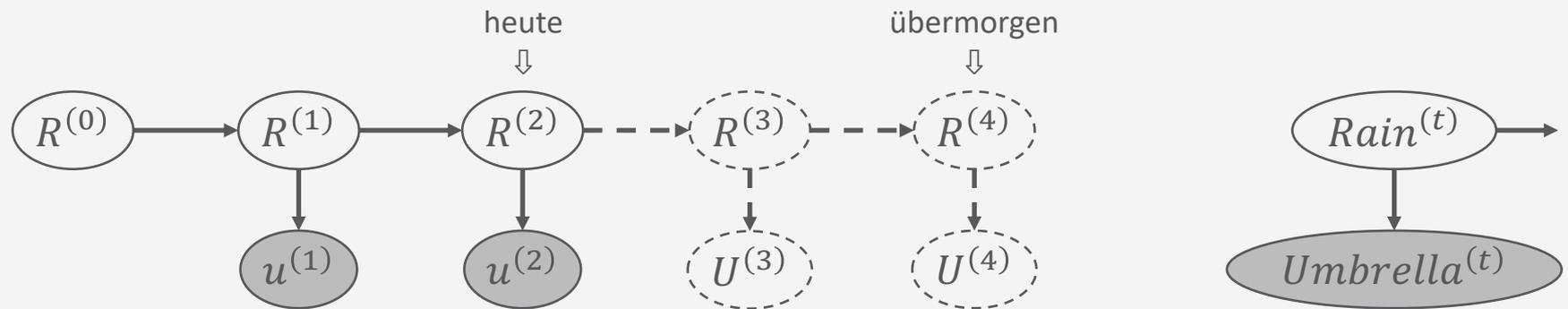
$$P(R^{(3)} | u^{(1)}, u^{(2)}, \neg u^{(3)})$$



Anfragebeantwortungsproblem & Anfragetypen

- Vorhersage (*prediction*)
 - Was ist die Wahrscheinlichkeit, dass es übermorgen regnet gegeben alle Beobachtungen zu Regenschirmen bis einschließlich heute ($t = 2$)?

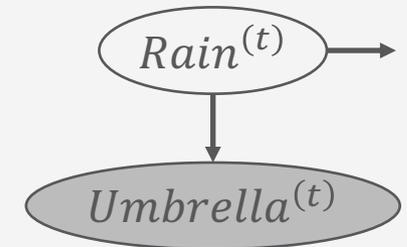
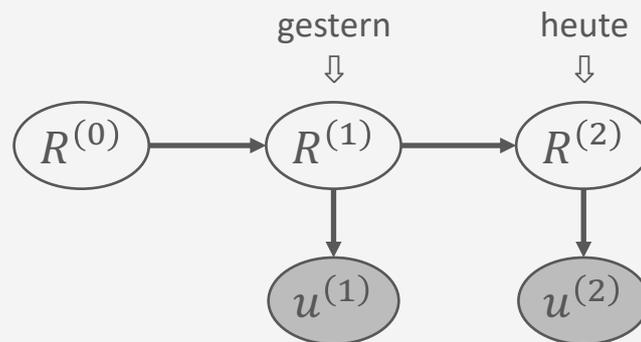
$$P(R^{(4)} | u^{(1)}, u^{(2)})$$



Anfragebeantwortungsproblem & Anfragetypen

- Rückschau (*hindsight*)
 - Auch Glättung (*smoothing*) genannt, bezieht sich aber eher auf das, was man berechnungstechnisch tut, als auf die Anfrage
 - Was ist die Wahrscheinlichkeit, dass es gestern geregnet hat gegeben alle Beobachtungen zu Regenschirmen bis einschließlich heute ($t = 2$)?

$$P(R^{(1)} | u^{(1)}, u^{(2)})$$

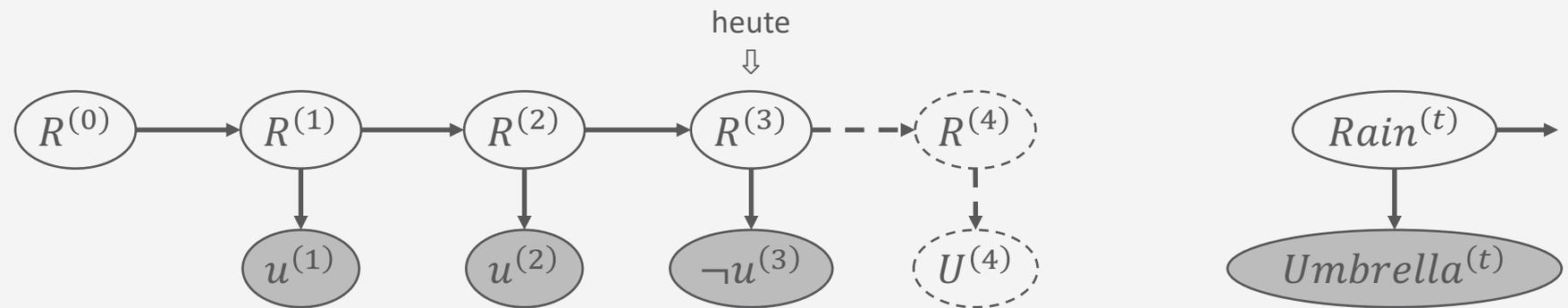


Anfragebeantwortungsproblem & Anfragetypen

- Wahrscheinlichste Erklärung (*MPE*)
 - Wenn der Regenschirm an den ersten beiden Tagen, aber nicht am dritten zu beobachten war, was ist die wahrscheinlichste Wettersequenz, die diese Regenschirm-Beobachtungen produziert hat?

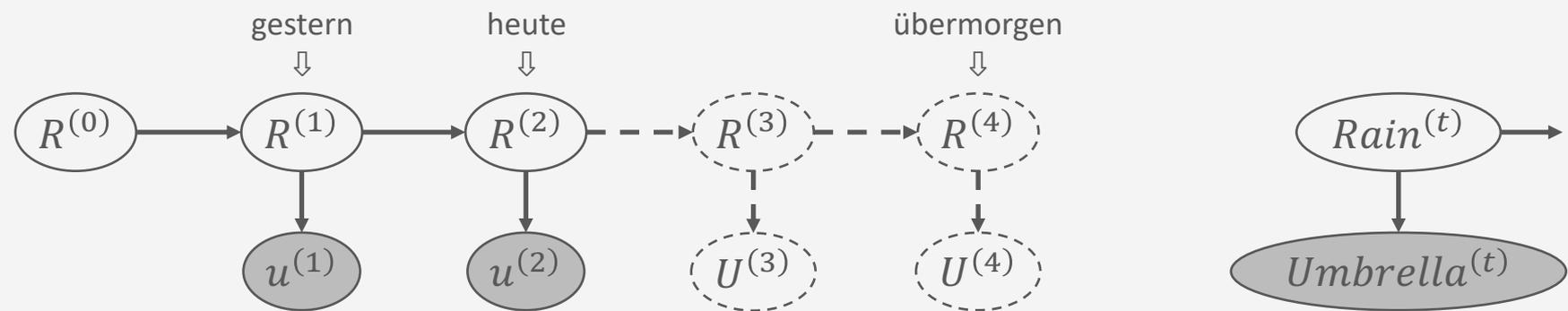
$$MPE(\mathbf{e}^{(1:3)}) = \arg \max_{\mathbf{r}^{(0:3)}} P(\mathbf{r}^{(0:3)} | \mathbf{e}^{(1:3)})$$

- Geht natürlich auch als *MAP*, wenn es mehr Variablen ohne Evidenz gibt und
 - Man nur an einer Untermenge von denen interessiert ist: $U \subseteq V$, V Variablen ohne Evidenz
 - Man nur an einigen (letzten k) Zeitschritten interessiert ist: $V^{(t-k:t)}$ bei $\mathbf{e}^{(0:t)}$



Anfragen: Formal

- Wahrscheinlichkeitsanfrage $P(\mathcal{S}^{(\pi)} | e^{(0:t)})$, t der aktuelle Zeitschritt, an ein Modell M mit Template-Zufallsvariablen R
 - Filtern / Überwachen (*filtering*): $\pi = t$
 - Vorhersage (*prediction*): $\pi > t$
 - Rückschau (*hindsight*): $\pi < t$



Anfragen: Formal

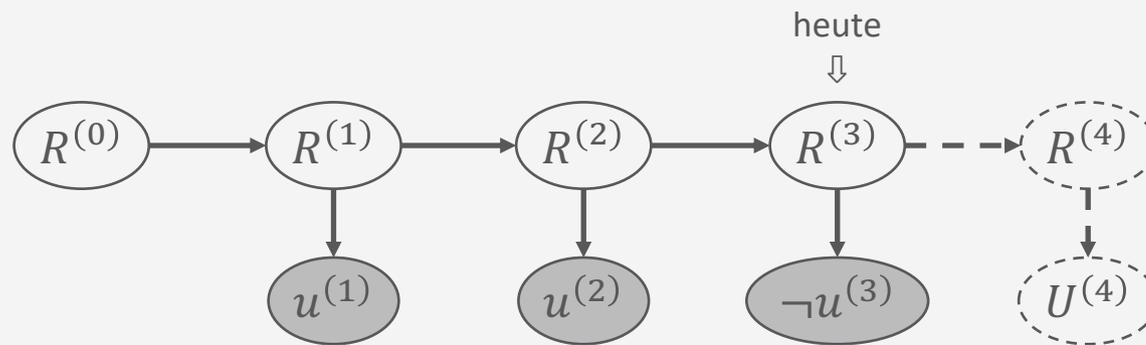
- Zustandsanfragen

$$MPE_M(\mathbf{e}^{(0:T)}) = \arg \max_{\mathbf{v}^{(0:t)} \in \text{Val}(\mathbf{V}^{(0:T)})} P(\mathbf{v}^{(0:T)} | \mathbf{e}^{(0:T)})$$

- $\mathbf{V}^{(0:T)} = \mathbf{R}^{(0:T)} \setminus \text{rv}(\mathbf{e}^{(0:T)})$

$$MAP_M(\mathbf{U}^{(t_1:t_2)} | \mathbf{e}^{(0:t)}) = \arg \max_{\mathbf{u}^{(t_1:t_2)} \in \text{Val}(\mathbf{U}^{(t_1:t_2)})} \sum_{\mathbf{v}^{(0:T)} \in \text{Val}(\mathbf{V}^{(0:T)})} P(\mathbf{u}^{(t_1:t_2)}, \mathbf{v}^{(0:T)} | \mathbf{e}^{(0:t)})$$

- $\mathbf{V}^{(0:T)} = \mathbf{R}^{(0:T)} \setminus \mathbf{U}^{(t_1:t_2)} \setminus \text{rv}(\mathbf{e}^{(0:t)}), T = \max\{t, t_1, t_2\}$



Zwischenzusammenfassung

- Sequentielle Daten modellieren
 - Annahmen: Markov-1, stationärer Prozess
 - Dynamische Modelle bestehen aus zwei episodischen Modellen
 - Eins um den ersten Schritt zu beschreiben
 - Eins um die Vorgänge innerhalb eines Schritts und den Übergang von einem Schritt zum nächsten zu beschreiben
 - Kopiervorlage
 - Semantik: Modell für T Schritte ausrollen, ergibt ein episodisches Modell mit bekannter Semantik (vollständige gemeinsame Verteilung)
- Dynamische Modelle: dynamische Faktormodelle, DBNs, HMMs
 - Folgen alle den gleichen Annahmen und der gleichen Struktur / Idee
- Inferenzaufgaben: Filtering, Prediction, Hindsight; MPE, MAP

Überblick: 6. Sequentielle PGMs und Inferenz

A. *Sequentielle PGMs*

- Templates, dynamische BNs, dynamische Faktormodelle, Hidden-Markov-Modelle; Semantik
- Inferenzaufgaben: Wahrscheinlichkeitsanfragen (Filtering, Prediction, Hindsight), Zustandsanfragen (MPE, MAP)

B. *Sequentielle Inferenz*

- Naïve Inferenz mittels Ausrollen, Interface Algorithmus, Komplexität, Approximationen

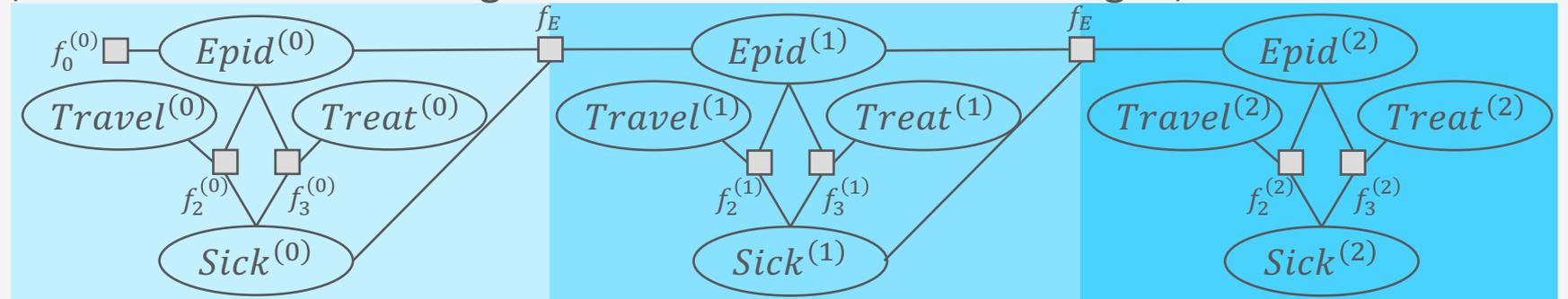
C. *Spezialfall Hidden-Markov-Modelle*

- Viterbi-Algorithmus für MPEs
- Anfragebeantwortung durch Matrixoperationen
- Baum-Welch-Algorithmus zum Lernen

Naïve Inferenz über Ausrollen des Modells

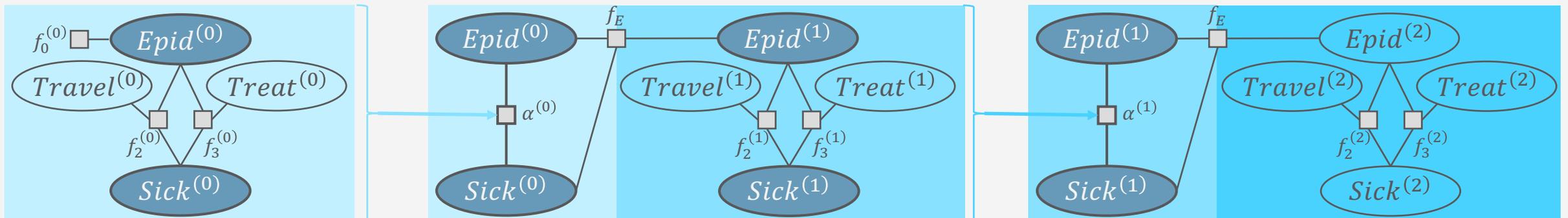
- Gegeben ein Modell $M = (M^0, M^{\rightarrow})$ und eine Anfrage $P(\mathbf{R}^{(\pi)} | \mathbf{e}^{(1:t)})$
 1. Modell für $T = \max\{t, \pi\}$ Schritte ausrollen \rightarrow episodisches Modell
 2. $P(\mathbf{R}^{(\pi)} | \mathbf{e}^{(1:t)})$ mit Algorithmus für episodische Modelle beantworten
- Probleme
 - Ausgerolltes Modell sehr groß
 - Neu ausrollen (oder anpassen), wenn sich T oder $\mathbf{e}^{(1:t)}$ ändern
 - Zeit schreitet voran: $T++$, $\mathbf{e}^{(1:t)}$ um $\mathbf{e}^{(t+1)}$ verlängert $\rightarrow M^{\rightarrow}$ für $\tau = t + 1$ anhängen, Evidenz absorbieren
 - Vor allem bei mehreren Anfragen mit sich ändernden π pro t aufwendig

Ziel: Nur mit einem aktuellen Modell $M^{(t)}$ arbeiten, welches mehrere Anfragen und fortschreitende Zeit ($t++$) effizient verarbeiten kann \rightarrow Interfaces nutzen



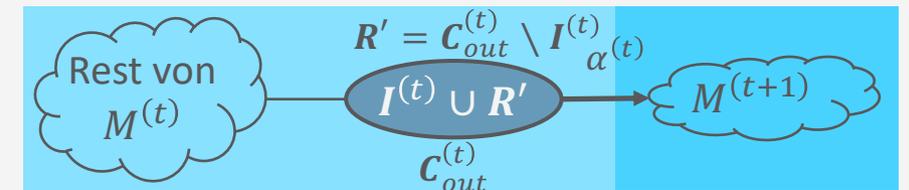
Interfaces um mit der Zeit zu gehen

- Wenn Inferenz für Zeitschritt t fertig
 - Anfrage für $I^{(t)}$ an aktuelles Modell $M^{(t)}$ stellen
 - Quasi wie im Junction Tree Algorithmus (JT): Nachricht $\alpha^{(t)}$ über Separator $I^{(t)}$ berechnen
 - Resultat der Anfrage zu Modell $M^{(t+1)}$ hinzufügen
 - $M^{(t+1)} = M^{\rightarrow|\tau=t+1} \cup \alpha^{(t)}$
 - Ohne Anbindung an $M^{(0:t)}$, da Information über $M^{(0:t)}$ in $\alpha^{(t)}$ vorhanden
- Jtrees und JT dafür nutzen



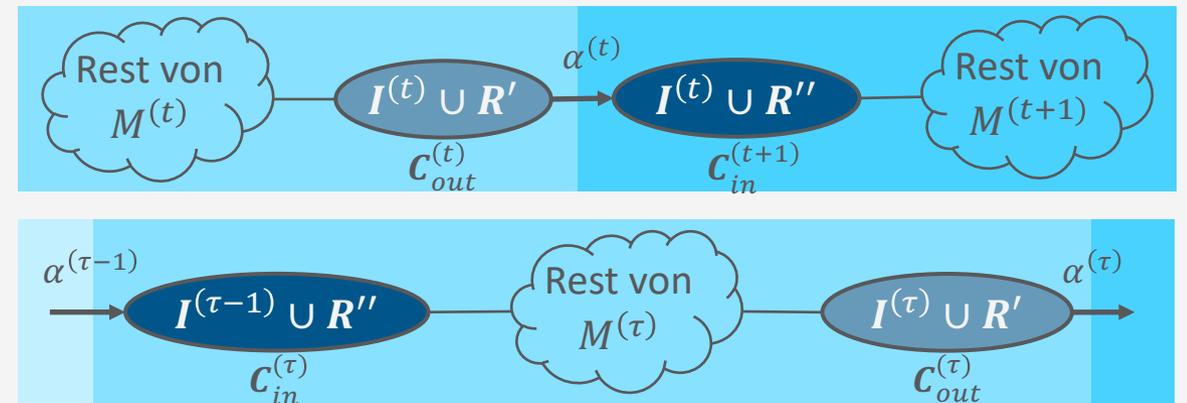
Jtrees und JT um mit der Zeit zu gehen

- Ermöglicht auch effiziente Beantwortung von mehreren Anfragen pro Zeitscheibe
- Anpassung für die Dynamik: $I^{(t)}$ in einem Cluster vorkommen lassen
 - Vermeidet wiederkehrende große Subgraph-Anfrage
 - Im Faktormodell: Faktor $\phi(I^{(\tau)})$ vor dem Bau des Jtrees zum Modell hinzufügen
 - Potentiale in ϕ irrelevant für Bau, können also unspezifiziert gelassen werden
 - ϕ kann nach dem Bau entfernt werden
 - Wenn nicht entfernt, Potentiale in ϕ uniform verteilt
 - Im BN: Beim Moralisieren des BNs, $I^{(\tau)}$ vollständig miteinander verbinden
 - Stellt sicher, dass ein Cluster \mathcal{C} $I^{(\tau)}$ enthält (Jtree Eigenschaft 2), i.e., $I^{(\tau)} \subseteq \mathcal{C}$
 - Dieses Cluster trennt die Zeitscheiben voneinander
 - So genanntes **Outcluster**, da hier Nachricht $\alpha^{(t)}$ über $I^{(\tau)}$ berechnet und losgeschickt wird (*outgoing*)



Jtrees und JT um mit der Zeit zu gehen

- Da $\alpha^{(t)}$ über $I^{(t)}$ zu einem Cluster im Jtree für $t + 1$ hinzugefügt werden muss, braucht man auch ein Cluster \mathcal{C} , welches das Interface zur vorherigen Zeitscheibe enthält
 - So genanntes **Incluster**, da hier Nachricht $\alpha^{(t)}$ über $I^{(t)}$ ankommt (*incoming*)
 - Im Faktormodell: Auch ein Faktor $\phi(I^{(\tau-1)})$ vor dem Bau des Jtrees zum Modell hinzufügen
 - Es gilt wieder, Potentiale in ϕ irrelevant für Bau, können also unspezifiziert gelassen werden
 - ϕ kann nach dem Bau entfernt werden; wenn nicht entfernt, Potentiale in ϕ uniform verteilt
- Im BN: Beim Moralisieren des BNs, $I^{(\tau-1)}$ vollständig miteinander verbinden
- Stellt sicher, dass ein Cluster $\mathcal{C} I^{(\tau-1)}$ enthält (Jtree Eigenschaft 2), i.e., $I^{(\tau-1)} \subseteq \mathcal{C}$
- In jeder $t > 0$ Zeitscheibe, In- und Outcluster vorhanden



Incluster & Outcluster: Gegenüberstellung

- Cluster, welches $I^{(\tau-1)}$ enthält \rightarrow Incluster $C_{in}^{(\tau)}$
 - Aus der Perspektive von τ : separiert Vergangenheit von der Gegenwart
 - Erhält eingehende Informationen vom Outcluster der vorherigen Zeitscheibe
- Cluster, welches $I^{(\tau)}$ enthält \rightarrow Outcluster $C_{out}^{(\tau)}$
 - Aus der Perspektive von τ : separiert Gegenwart von der Zukunft
 - Schickt Informationen raus zum Incluster der nächsten Zeitscheibe
 - Gegenwart wird zur Vergangenheit im nächster Schritt



Dynamische Jtrees

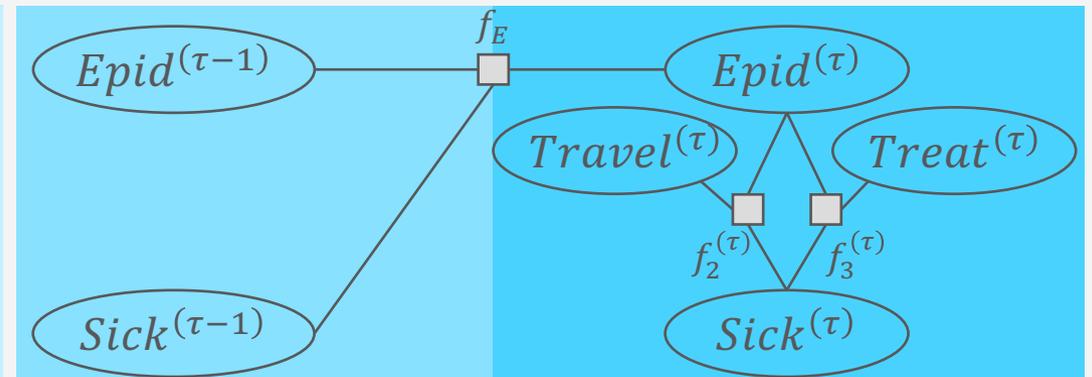
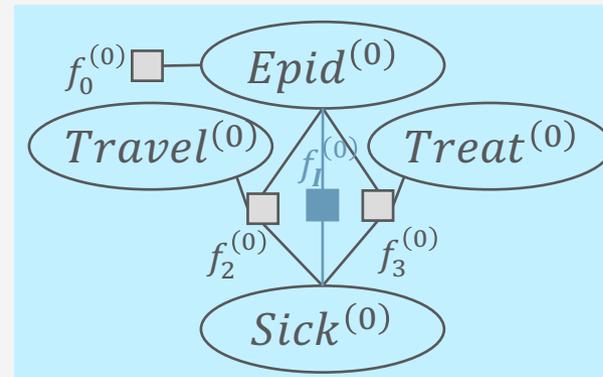
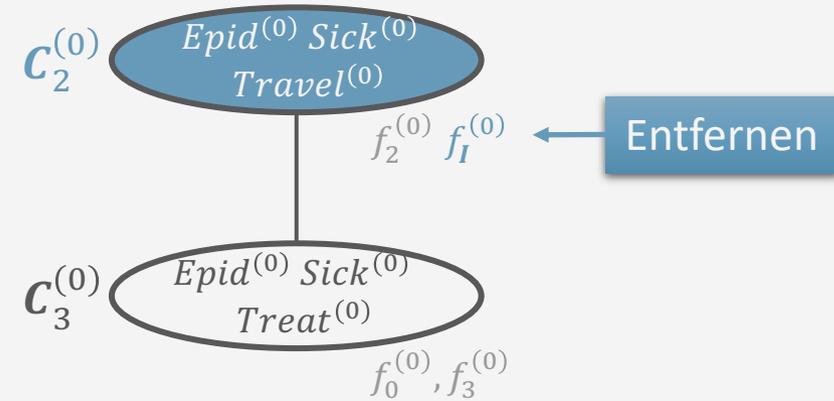
- Gegeben ein Interface $I^{(\tau)}$ für ein dynamisches Modell $M = (M^0, M^{\rightarrow})$
- Baue zwei Jtrees (J^0, J^{\rightarrow}) für M
 - J^0 für $M^0 \cup \{\phi(I^{(0)})\}$
 - $I^{(0)} = I^{(\tau)|\tau=0}$ in einem Cluster (Outcluster)
 - J^{\rightarrow} für $M^{\rightarrow} \cup \{\phi(I^{(\tau-1)}), \phi(I^{(\tau)})\}$
 - $I^{(\tau-1)} = I^{(\tau)|\tau=\tau-1}$ in einem Cluster (Incluster)
 - $I^{(\tau)}$ in einem Cluster (Outcluster)
 - **Template-Jtree**, der instanziiert wird, indem τ mit einem Zeitpunkt t ersetzt wird
- Standardverfahren zum Bau der Jtrees verwenden
 - Z.B. Dtree bauen, Cluster bestimmen, Cluster in Jtree umwandeln, minimieren

Jtree $J = (V, E)$ ein ungerichteter, azyklischer Graph mit V eine Menge von Clustern und E eine Menge von Kanten für ein Modell M ; muss drei Eigenschaften erfüllen:

1. Cluster bestehen aus Variablen aus M
2. Variablen jedes Faktors / CPD in einem Cluster enthalten
3. *Running intersection property*

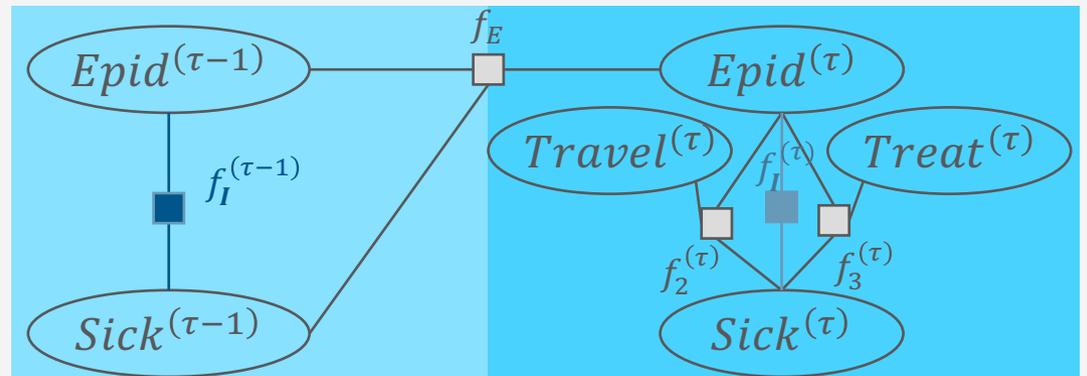
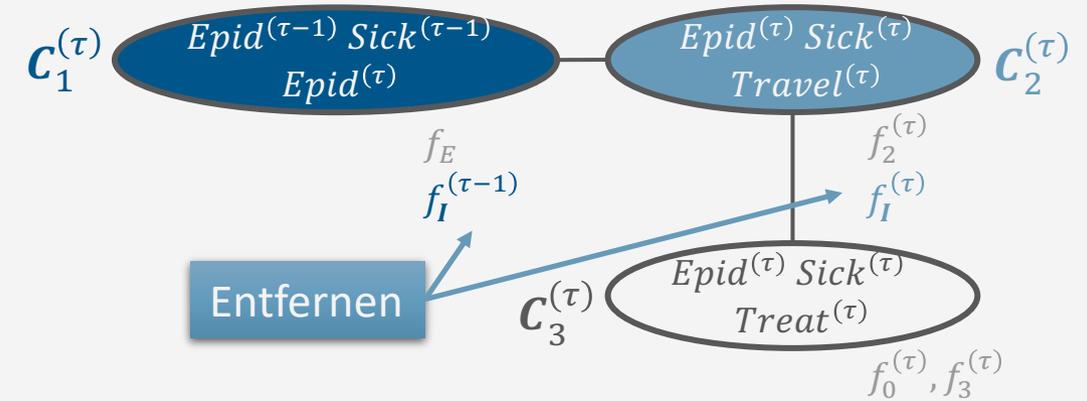
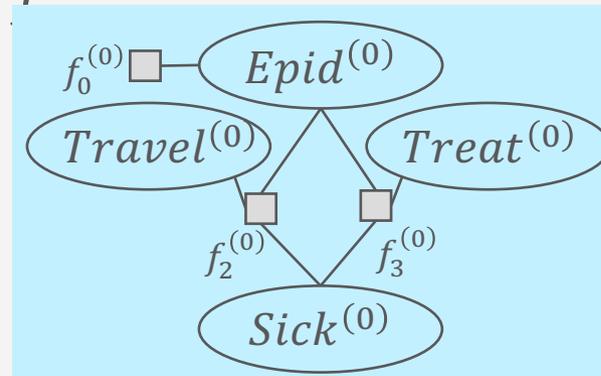
Dynamische Jtrees: Beispiel

- Dynamisches Faktormodell $F = (F^0, F^\rightarrow)$
 - $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}, F^\rightarrow = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$
 - Interface $I^{(\tau)} = \{Epid^{(\tau)}, Sick^{(\tau)}\}$
- J^0 für $F^0 \cup \{f_I^{(0)}\}$
 - $f_I^{(0)} = \phi(E^{(0)}, S^{(0)})$
 - Beide Cluster enthalten Interface $I^{(0)}$
 - In der Regel wird das gewählt, welches $f_I^{(0)}$ enthält: $\mathcal{C}_2^{(0)} = \mathcal{C}_{out}^{(0)}$



Dynamische Jtrees: Beispiel

- Dynamisches Faktormodell $F = (F^0, F^\rightarrow)$
 - $F^0 = \{f_2^{(0)}, f_3^{(0)}\} \cup \{f_0^{(0)}\}, F^\rightarrow = \{f_2^{(\tau)}, f_3^{(\tau)}\} \cup \{f_E\}$
 - Interface $I^{(\tau)} = \{Epid^{(\tau)}, Sick^{(\tau)}\}$
- J^\rightarrow für $F^\rightarrow \cup \{f_I^{(\tau-1)}, f_I^{(\tau)}\}$
 - $f_I^{(\tau-1)} = \phi(E^{(\tau-1)}, S^{(\tau-1)})$
 - $f_I^{(\tau)} = \phi(E^{(\tau)}, S^{(\tau)})$
 - Incluster: $\mathcal{C}_1^{(\tau)} = \mathcal{C}_{in}^{(\tau)}$
 - Outcluster: $\mathcal{C}_2^{(\tau)} = \mathcal{C}_{out}^{(\tau)}$



Dynamische Jtrees: Ausrollen

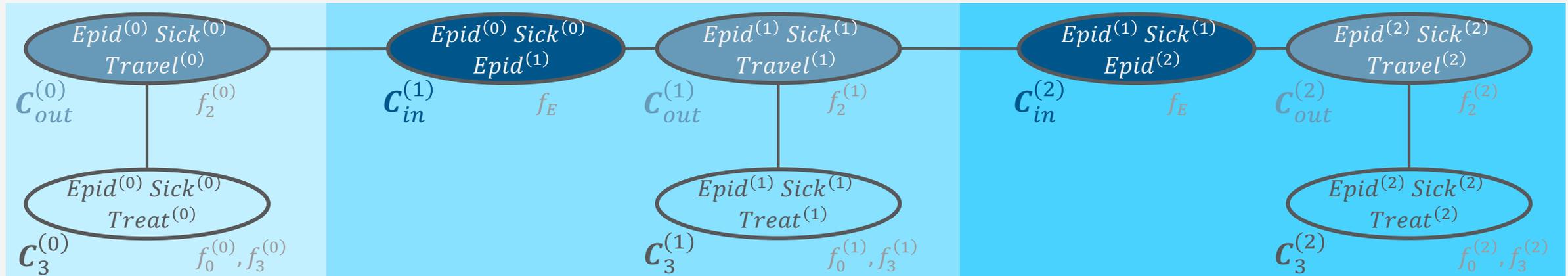
- Gegeben zwei Jtrees (J^0, J^{\rightarrow}) für ein dynamische Modell $M = (M^0, M^{\rightarrow})$ und Zahl T , rolle (J^0, J^{\rightarrow}) aus durch Bilden eines Jtrees $J^{(0:T)}$ aus Jtrees $J^{(0)}, J^{(1)}, \dots, J^{(T)}$, wobei
 - $J^{(0)} = (V^{(0)}, E^{(0)}) = J^0$
 - $J^{(t)} = (V^{(t)}, E^{(t)}) = J^{\rightarrow|_{\tau=t}}$ für alle $J^{(t)}, t > 0$
 - Für alle Zeitscheiben $t - 1, t$ für $t > 0$, eine Kante zwischen dem Outcluster von $t - 1$ und dem Incluster von t eingefügt wird
 - Formal: $J^{(0:T)} = (V, E)$ mit

$$V = V^{(0)} \cup \bigcup_{t=1}^T V^{(t)}$$

$$E = E^{(0)} \cup \bigcup_{t=1}^T E^{(t)} \cup \bigcup_{t=1}^T \left\{ \left\{ \mathbf{c}_{out}^{(t-1)}, \mathbf{c}_{in}^{(t)} \right\} \mid \mathbf{c}_{out}^{(t-1)} \wedge \mathbf{c}_{in}^{(t)} \right\}$$

Dynamische Jtrees: Ausrollen – Beispiel

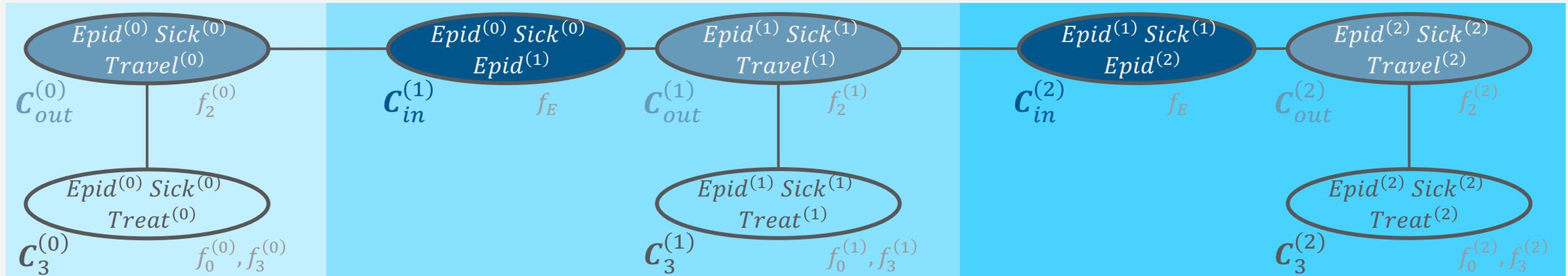
- $T = 2$



- JT für Anfragenbeantwortung nutzen
 - Evidenz behandeln, Nachrichten verschicken, Anfragen beantworten
- Probleme (wie beim Modell ausrollen)
 - Ausgerollter Jtree sehr groß
 - Neu ausrollen (oder anpassen), wenn sich $T = \max\{t, \pi\}$ oder $e^{(1:t)}$ ändern

Dynamische Jtrees: Ausrollen – Beispiel

- $T = 2$



- Ausrollen zur Anfragenbeantwortung innerhalb einer Zeitscheibe unnötig
 - Durch die Interfaces in den In- und Outclustern sind die Zeitscheiben unabhängig von einander
 - Anfragenbeantwortung innerhalb einer Zeitscheibe mittels JT
 - Übergang zur nächsten Zeitscheibe: Nachricht in Outcluster $C_{out}^{(t)}$ über $I^{(t)}$ berechnen und zum Incluster $C_{in}^{(t+1)}$ im neu instanziierten Jtree für $t + 1$ hinzufügen

Mit der Zeit gehen

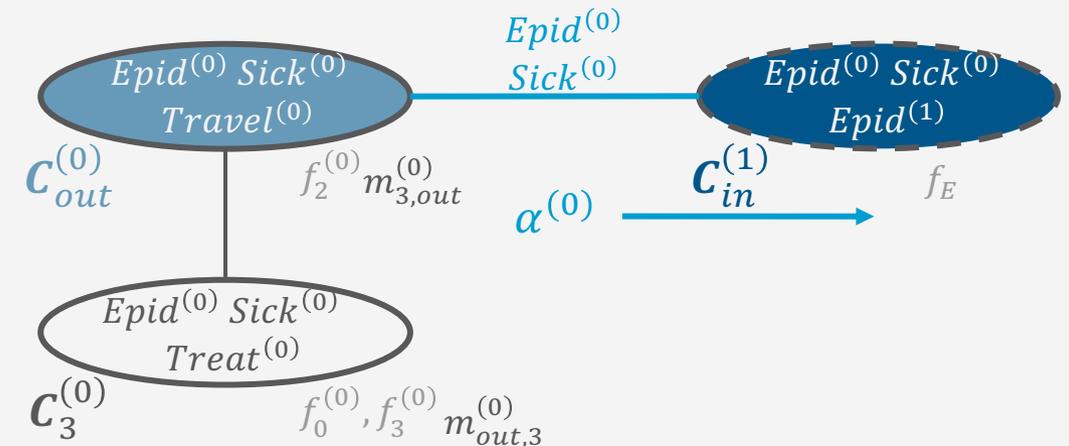
- Anfragebeantwortung innerhalb einer Zeitscheibe t folgt JT
 - Menge von Anfragen für t bei Evidenz bis einschließlich t effizient beantworten
 - Stellt auch sicher, dass alle notwendigen Informationen bei $\mathbf{C}_{out}^{(t)}$ verfügbar sind, um die **Vorwärtsnachricht** zu berechnen, die $t + 1$ unabhängig von t macht
- Übergang zur nächsten Zeitscheibe:
 - Berechne **Vorwärtsnachricht** $\alpha^{(t)}$ über den Separator zwischen $\mathbf{C}_{out}^{(t)}$ und $\mathbf{C}_{in}^{(t+1)}$ basierend auf dem lokalen Modell $F_{out}^{(t)}$ und den eingegangenen Nachrichten $m_{j,out}^{(t)}$
 - Instanziiere J^{\rightarrow} für $t + 1$
 - Füge $\alpha^{(t)}$ zum lokalen Modell von $\mathbf{C}_{in}^{(t+1)}$ hinzu
 - Lasse $J^{(t)}$ fallen
 - Mache Anfragebeantwortung in $J^{(t+1)}$ (Schritte 2-4 von JT)

Mit der Zeit gehen: Beispiel

- $t = 0$
 - aktueller Jtree
 - Evidenz $e^{(0)} = \{treat^{(0)}\}$ behandeln
 - Intra-Zeitscheiben-Nachrichten: $m_{3,out}^{(0)}, m_{out,3}^{(0)}$
 - Beantworte Anfragen $P(R_i^{(0)} | e^{(0)})$
- Inter-Zeitscheiben-Nachricht: $\alpha^{(0)}$
 - Eliminiere alle Nicht-Separator-Variablen, $Travel^{(0)}$, aus dem lokalen Modell $F_{out}^{(0)} = \{f_2^{(0)}\}$ und $m_{3,out}^{(0)}$
 - Ergebnis als Nachricht $\alpha^{(0)}$ zu $C_{in}^{(1)}$

$\alpha^{(0)}$ beinhaltet alle Informationen aus $F^{(0)}$ inklusive $e^{(0)}$, was Schritt 0 und 1 unabhängig macht.

Als nächstes: Instanziiere einen Jtree für $\tau = 1$ und füge $\alpha^{(0)}$ zum lokalen Modell von $C_{in}^{(1)}$ hinzu.

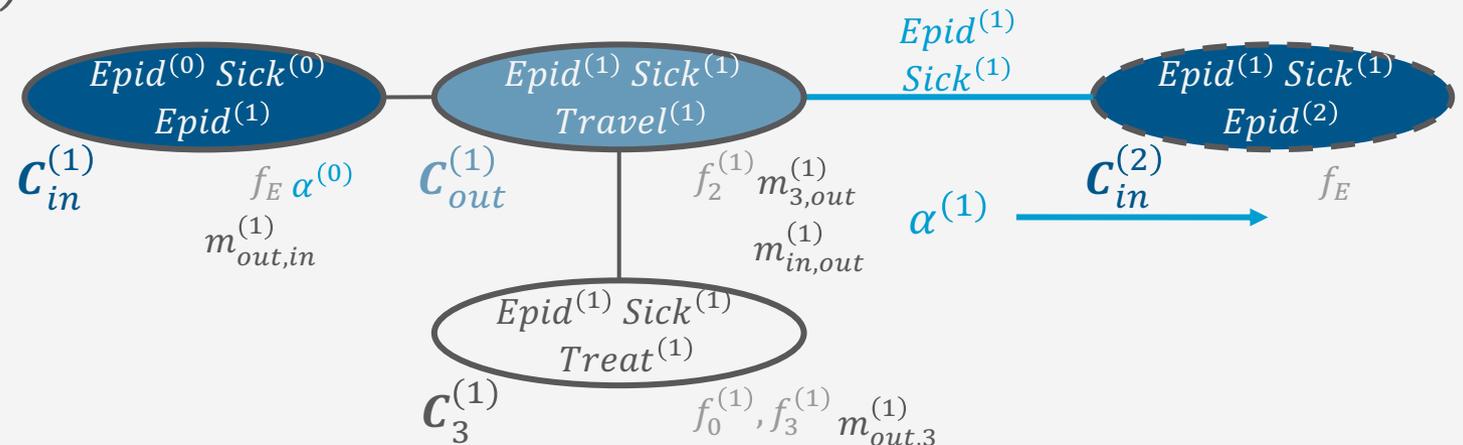


Mit der Zeit gehen: Beispiel

- $t = 1$
 - aktueller Jtree
 - Evidenz $e^{(1)} = \{treat^{(1)}\}$ behandeln
 - Intra-Z.-N.: $m_{3,out}^{(1)}, m_{out,3}^{(1)}, m_{in,out}^{(1)}, m_{out,in}^{(1)}$
 - Beantworte Anfragen $P(R_i^{(1)} | e^{(0:1)})$
- Inter-Zeitscheiben-Nachricht: $\alpha^{(1)}$
 - Eliminiere $Travel^{(1)}$ aus dem lokalen Modell $F_{out}^{(1)} = \{f_2^{(1)}\}$, $m_{in,out}^{(1)}$ und $m_{3,out}^{(1)}$
 - Ergebnis als Nachricht $\alpha^{(1)}$ zu $C_{in}^{(2)}$

Während des Nachrichtenversands werden die Informationen aus $\alpha^{(0)}$ verteilt und dadurch in $\alpha^{(1)}$ mit verrechnet. $\alpha^{(1)}$ beinhaltet nun alle Informationen aus $F^{(0:1)}$ inklusive $e^{(0:1)}$, was Schritt 1 und 2 unabhängig macht.

Als nächstes: Instanziiere einen Jtree für $\tau = 2$ und füge $\alpha^{(1)}$ zum lokalen Modell von $C_{in}^{(2)}$ hinzu.

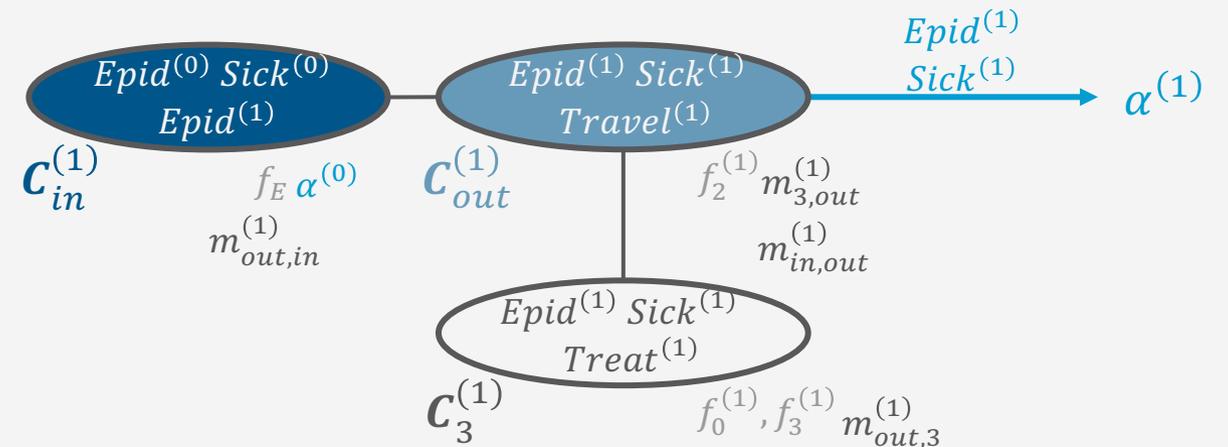


Filtern / Überwachen (*filtering*)

- Anfragen, die auf diese Art beantwortet werden: Filterungsanfragen (*filtering*)
 - $P(R^{(t)} | e^{(0:t)})$
 - Anfragen an Zufallsvariablen der aktuellen Zeitscheibe t gegeben die bisherige Evidenz $e^{(0:t)}$
- Vorteile
 - Nur ein aktueller Jtree nötig
 - Eine zusätzliche Nachricht ($\alpha^{(t)}$), um in der Zeit voranzuschreiten

- Was ist mit Vorhersagen (*prediction*) und Rückschauen (*hindsight*)?
 - $P(R^{(\pi)} | e^{(0:t)}), \pi \neq t$

Hat jemand eine Idee?



Vorhersagen (*prediction*)

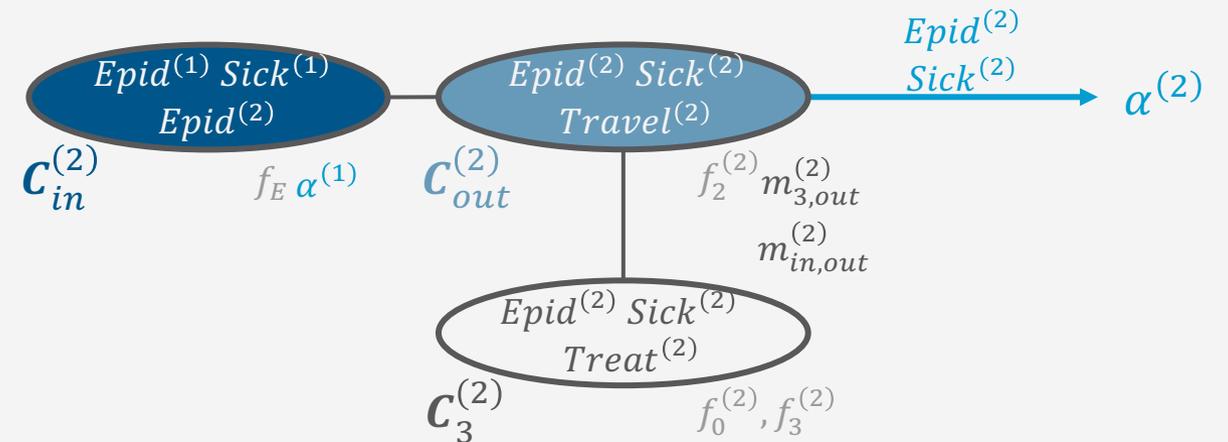
- Vorhersagen: Anfragen an zukünftige Instanzen von Zufallsvariablen $P(R^{(\pi)} | e^{(0:t)})$, $\pi > t$
 - Vorausschau in die Zukunft: sich in der Zeit bis π vorwärts bewegen ohne Evidenz zu sehen, i.e.,
 - Filterung: $P(R^{(\pi)} | e^{(0:\pi)})$ mit leerer Evidenz zwischen $t + 1$ und π in $e^{(0:\pi)}$: $e^{((t+1):\pi)} = \{\emptyset^{(t')}\}_{t'=t+1}^{\pi}$
- Anfragebeantwortung für Vorhersagen
 - Gehe vorwärts in der Zeit bis $t = \pi$
 - Eine Nachrichten-Phase von der Peripherie zum Zentrum mit *Outcluster* als Zentrum ausreichend
 - Beantworte Anfrage mit Anfragevariable $R^{(\pi)}$
 - Wenn nur eine Anfrage
 - Finde Cluster C_i^{π} , welches $R^{(\pi)}$ enthält; führe eine Nachrichten-Phase von der Peripherie zum Zentrum mit C_i^{π} als Zentrum durch; beantworte Anfrage in C_i^{π}
 - Sonst:
 - Vollständigen Nachrichtenversand (zwei Phasen) durchführen; Anfragen beantworten



Spart die Hälfte der Nachrichten

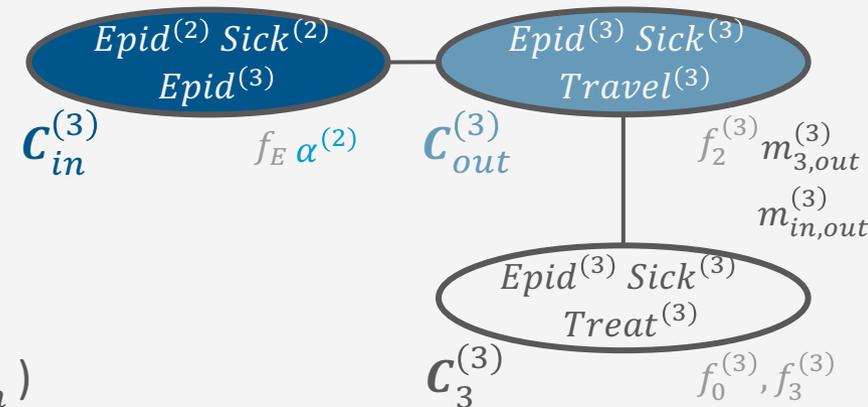
Vorhersagen (*prediction*): Beispiel

- $P(Epid^{(3)} | e^{(0:1)})$, $3 > 1$
 - Keine Evidenz für Schritte 2 und 3: $e^{(2:3)} = \{\emptyset^{(t')}\}_{t'=2}^3$
 - Gehe vorwärts in der Zeit bis $t = 3$, dann beantworte die Anfrage mit Anfragevariable $Epid^{(3)}$
 - $t = 2$, aktueller Jtree mit $\alpha^{(1)}$ im lokalen Modell von $\mathcal{C}_{in}^{(2)}$
 - Keine Evidenz zu behandeln: $\emptyset^{(2)}$
 - Versende Intra-Zeitscheiben-Nachrichten in Richtung $\mathcal{C}_{out}^{(2)}$ (alle Informationen am Outcluster sammeln)
 - Berechne Inter-Zeitscheiben-Nachricht $\alpha^{(2)}$



Vorhersagen (*prediction*): Beispiel

- $P(Epid^{(3)} | e^{(0:1)})$, $3 > 1$
 - Keine Evidenz für Schritte 2 und 3: $e^{(2:3)} = \{\emptyset^{(t')}\}_{t'=2}^3$
 - Gehe vorwärts in der Zeit bis $t = 3$, dann beantworte die Anfrage mit Anfragevariable $Epid^{(3)}$
 - $t = 3$, aktueller Jtree mit $\alpha^{(2)}$ im lokalen Modell von $C_{in}^{(3)}$
 - Keine Evidenz zu behandeln: $\emptyset^{(3)}$
 - Wenn nur Anfrage zu $Epid^{(3)}$
 - Versende Intra-Zeitscheiben-Nachrichten in Richtung $C_{out}^{(3)}$ (oder $C_3^{(3)}$)
 - Beantworte Anfrage in $C_{out}^{(3)}$
 - Sonst:
 - Alle Nachrichten versenden (auch $m_{out,3}^{(3)}$, $m_{out,in}^{(3)}$)



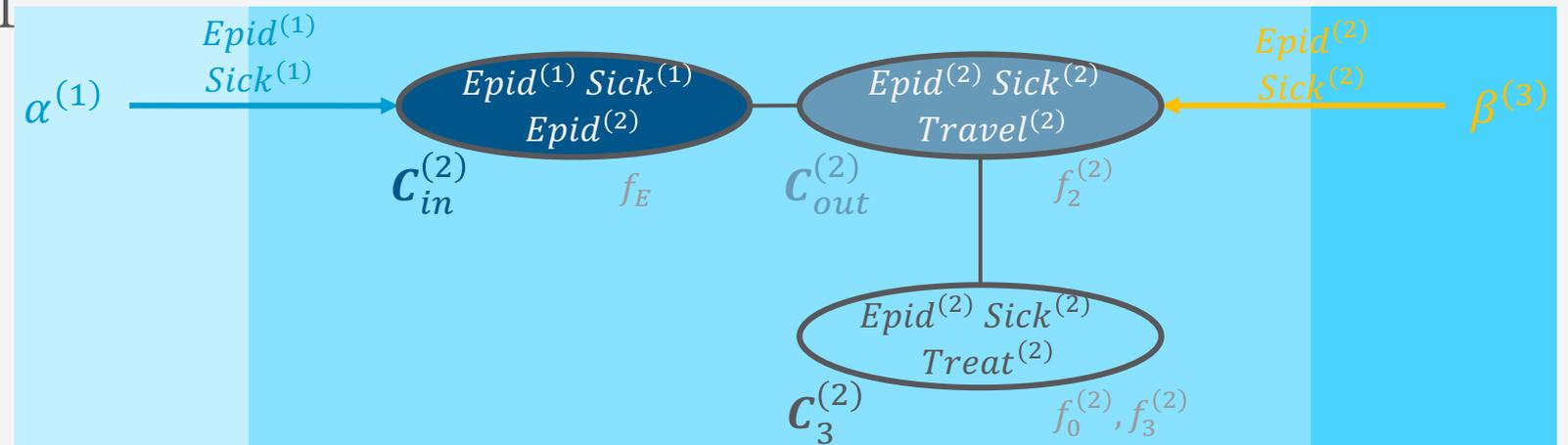
Rückschau (*hindsight*): Wieder zurück in der Zeit

- Rückschau: Anfrage an vergangene Instanzen von Zufallsvariablen $P(R^{(\pi)} | e^{(0:t)})$, $\pi < t$
 - Zurückblicken mit dem Wissen von jetzt (t) bezogen auf Evidenz, i.e.,
 - Von der jetzigen Zeitscheibe t rückwärts gehen und die Informationen, die sich zwischen π and t angesammelt haben, nach π bringen; dann eine Filterungsanfrage in π beantworten
- Anfragebeantwortung für Rückschauen
 - Bewege dich rückwärts in der Zeit bis $t = \pi$
 - Berechne dafür *Rückwärtsnachrichten* $\beta^{(t)}$ von $C_{in}^{(t)}$ zu $C_{out}^{(t-1)}$ beginnend bei t bis π
 - Eine Nachrichten-Phase von der Peripherie zum Zentrum mit *Incluster* als Zentrum ausreichend
 - Beantworte Anfrage mit Anfragevariable $R^{(\pi)}$ in J^π mit $\alpha^{(\pi-1)}$ im Incluster, $\beta^{(\pi+1)}$ im Outcluster
 - Wenn nur eine Anfrage
 - Eine Nachrichten-Phase von der Peripherie zu Cluster C_i^π , $R^{(\pi)} \in C_i^\pi$, als Zentrum; beantworte Anfrage in C_i^π
 - Sonst: Vollständigen Nachrichtenversand (zwei Phasen) durchführen; Anfragen beantworten

Rückschau (*hindsight*): Rückwärtsnachricht $\beta^{(t)}$

- Idee: Zeitscheibe unabhängig von der Zukunft machen (Vorwärtsnachricht umgekehrt)
- Aus der Perspektive von Zeitscheibe π irgendwo in der Sequenz von 0 bis t
 - Vorwärtsnachricht $\alpha^{(\pi-1)}$ beinhaltet alle Informationen zu $M^{(0:(\pi-1))}$ inklusive $e^{(0:(\pi-1))}$
 - Macht $\pi - 1$ unabhängig von π
 - Rückwärtsnachricht $\beta^{(\pi+1)}$ beinhaltet alle Informationen zu $M^{((\pi+1):t)}$ inklusive $e^{((\pi+1):t)}$
 - Macht π unabhängig von $\pi + 1$

- Immer von $C_{in}^{(t)}$ zu $C_{out}^{(t-1)}$ verschickt, bis $t = \pi$

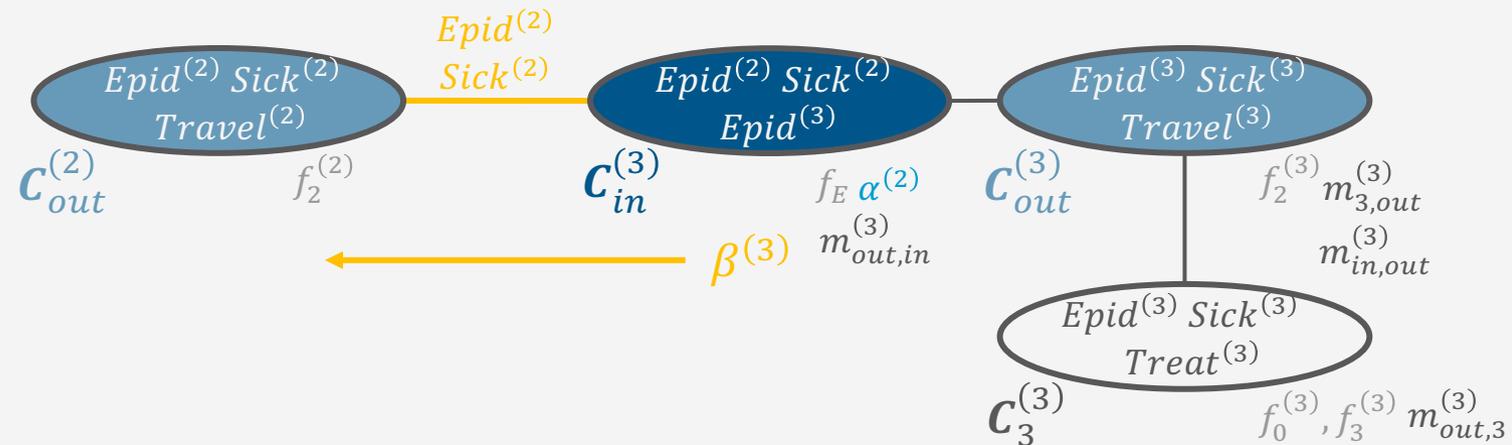


Rückschau (*hindsight*): Rückwärtsnachricht $\beta^{(t)}$

- Berechnung
 - Eliminiere alle Nichtseparator-Variablen zwischen $\mathcal{C}_{in}^{(t)}$ und $\mathcal{C}_{out}^{(t-1)}$ aus dem lokalen Modell $F_{in}^{(t)}$ und den eingehenden Nachrichten $m_{j,in}^{(t)}$
 - Für Rückwärtsnachricht nicht Vorwärtsnachricht $\alpha^{(t-1)}$ berücksichtigen, da $\alpha^{(t-1)}$ von $\mathcal{C}_{out}^{(t-1)}$ kam
 - Diese Information schon in Zeitscheibe $t - 1$ vorhanden

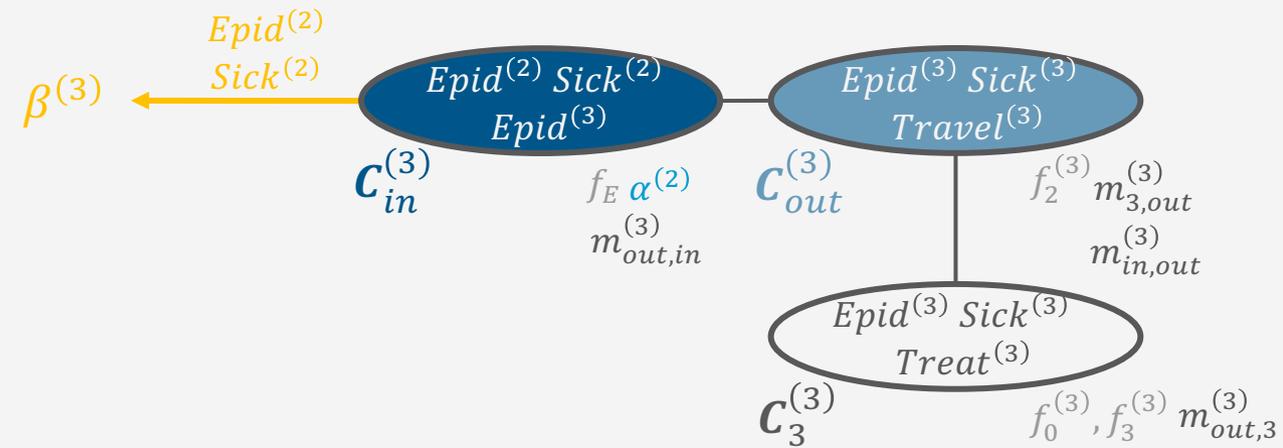
Beispiel:

- Aktuelle Zeitscheibe $t = 3$
- Berechne $\beta^{(3)}$ durch Eliminierung von $Epid^{(3)}$ aus $f_E, m_{out,in}^{(3)}$
 - Ohne $\alpha^{(2)}$



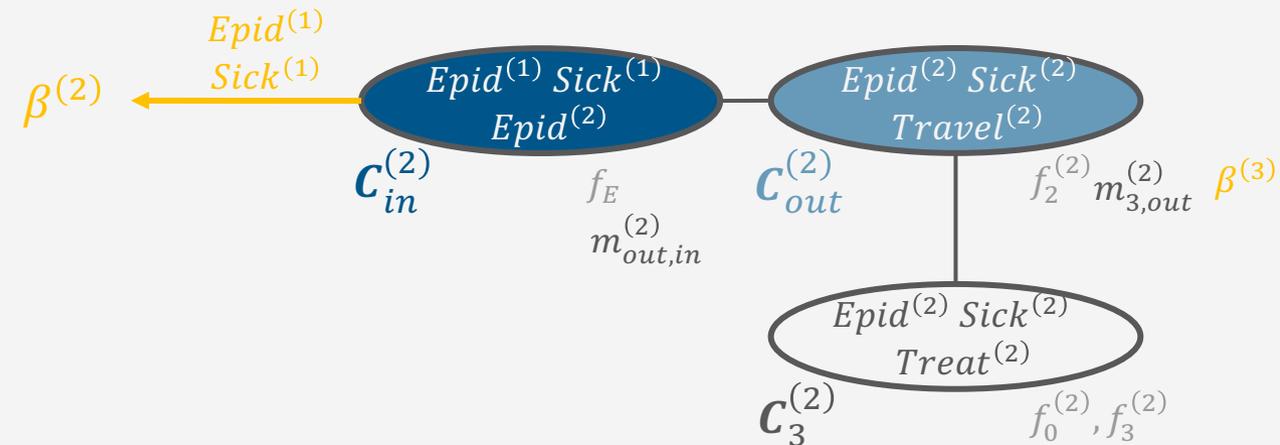
Rückschau (*hindsight*): Beispiel

- $P(Epid^{(1)} | e^{(0:3)})$, $1 < 3$
 - Gehe rückwärts in der Zeit bis $t = 1$, dann beantworte die Anfrage mit Anfragevariable $Epid^{(1)}$
 - $t = 3$, aktueller Jtree mit $\alpha^{(2)}$ im lokalen Modell von $\mathcal{C}_{in}^{(3)}$
 - Berechne Rückwärtsnachricht $\beta^{(3)}$
 - Instanziiere Jtree für $t = 2$, füge $\beta^{(3)}$ zu lokalem Modell von $\mathcal{C}_{out}^{(2)}$ hinzu



Rückschau (*hindsight*): Beispiel

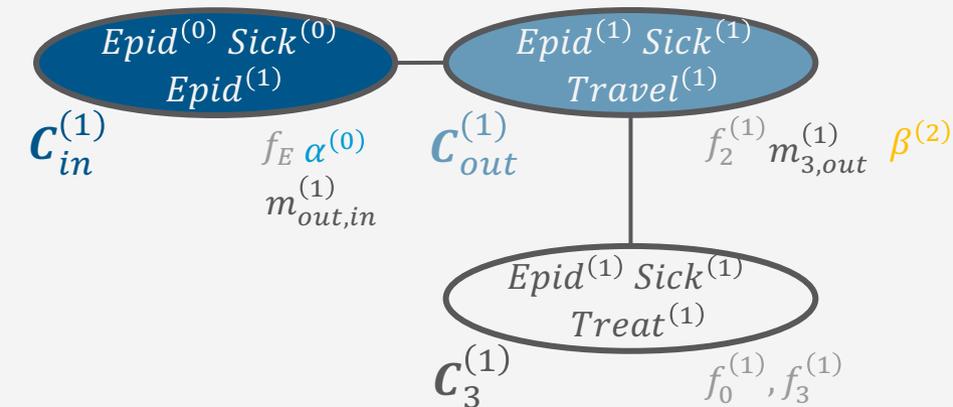
- $P(\text{Epid}^{(1)} | e^{(0:3)})$, $1 < 3$
 - Gehe rückwärts in der Zeit bis $t = 1$, dann beantworte die Anfrage mit Anfragevariable $\text{Epid}^{(1)}$
 - $t = 2$, aktueller Jtree mit $\beta^{(3)}$ im lokalen Modell von $\mathcal{C}_{out}^{(2)}$
 - Behandle Evidenz $e^{(2)}$
 - Sende Nachrichten in Richtung $\mathcal{C}_{in}^{(2)}$
 - Berechne Rückwärtsnachricht $\beta^{(2)}$
 - Instanziiere Jtree für $t = 1$, füge $\beta^{(2)}$ zu lokalem Modell von $\mathcal{C}_{out}^{(2)}$ hinzu



Rückschau (*hindsight*): Beispiel

- $P(Epid^{(1)} | e^{(0:3)})$, $1 < 3$
 - Gehe rückwärts in der Zeit bis $t = 1$, dann beantworte die Anfrage mit Anfragevariable $Epid^{(1)}$
 - $t = 1$, aktueller Jtree mit $\beta^{(2)}$ im lokalen Modell von $C_{out}^{(1)}$, $\alpha^{(0)}$ im lokalen Modell von $C_{in}^{(1)}$
 - Behandle Evidenz $e^{(1)}$
 - Wenn nur Anfrage zu $Epid^{(1)}$
 - Versende Intra-Zeitscheiben-Nachrichten in Richtung $C_{in}^{(1)}$ (oder die anderen beiden)
 - Beantworte Anfrage in $C_{in}^{(1)}$
 - Sonst: Alle Nachrichten versenden

Um Rückschau-Anfragen zu beliebigen Zeitscheiben π zu ermöglichen, müssen wir uns die Vorwärtsnachrichten merken, da $\alpha^{(\pi-1)}$ für die Anfragebeantwortung in π benötigt wird.



Allgemeine Anfragenbeantwortung

- **Vorwärts-Rückwärts-Algorithmus** für die Anfragentypen Filterung, Vorhersage, Rückschau
- Gegeben: Menge von Anfragen $\mathbf{Q}^{(0:T_q)} = \left\{ \left\{ \mathbf{Q}_i^{(\pi_i)} \right\}_{i=1}^{m_t} \right\}_{t=0}^{T_q}$
 - Bei anwachsender Evidenz $\mathbf{e}^{(0:t)}, t \in \{0, \dots, T_e\}$
 - Quasi zwei Ströme an Anfragen und Evidenz über die Zeit t
 - Aus Effizienzgründen, gehe durch $\mathbf{Q}^{(0:T_q)}$ mit ansteigendem t und durch die Anfragen $\left\{ \mathbf{Q}_i^{(\pi_i)} \right\}_{i=1}^{m_t}$ bei gegebener Evidenz $\mathbf{e}^{(0:t)}$ basierend auf dem Typ und π_i
 1. Filterungsanfragen $\mathbf{Q}_i^{(\pi_i)}, \pi_i = t$
 2. Vorhersage-Anfragen $\mathbf{Q}_i^{(\pi_i)},$ geordnet nach ansteigendem π_i ($\pi_i > t$)
 3. Rückschau-Anfragen $\mathbf{Q}_i^{(\pi_i)},$ geordnet nach absteigendem π_i ($\pi_i < t$)
 - Reihenfolge von 2. und 3. könnte man tauschen

Allgemeine Anfragenbeantwortung: Genaues Vorgehen bei Anfragen

- Gegeben t als die aktuelle Zeitscheibe, Anfragevariablen $\{Q_i^{(\pi_i)}\}_{i=1}^{m_t}$ und einem abgeschlossenen Nachrichtenversand in t
 - Für alle $Q_i^{(\pi_i)}$ mit $t = \pi_i$, beantworte $Q_i^{(\pi_i)}$ in $J^{(t)}$
 - Behalte $J^{(t)}$ und bewege dich mit $t' = t$ in der Zeit
 - Für alle $Q_i^{(\pi_i)}$ mit $t < \pi_i$, bewege dich ohne Evidenz vorwärts in der Zeit bis $t' = \max_i \pi_i$:
 - Instanziiere Jtree $J^{(t')}$ und berechne entsprechende Nachrichten
 - Wann immer $t' = \pi_i$, beantworte $Q_i^{(\pi_i)}$
 - Für alle $Q_i^{(\pi_i)}$ mit $t > \pi_i$, bewege dich rückwärts in der Zeit bis $t' = \min_i \pi_i$
 - Instanziiere Jtree $J^{(t')}$, behandle Evidenz $e^{(t')}$ und berechne entsprechende Nachrichten
 - Wann immer $t' = \pi_i$, beantworte $Q_i^{(\pi_i)}$



Welche Nachrichten müssen wir berechnen?

Dynamischer Jtree Algorithmus (DJT) – AKA *Murphy's Interface Algorithm*

```
procedure DJT( $(M^0, M^{\rightarrow}), Q^{(0:T_q)}, e^{(0:T_e)}$ )
```

```
  Baue Jtrees  $(J^0, J^{\rightarrow})$  für  $M^0, M^{\rightarrow}$ 
```

```
  for  $t$  in  $0 \dots T_q$  do
```

```
    Instanziiere  $J^{(t)}$  und lösche  $J^{(t-1)}$ 
```

```
    Füge  $\alpha^{(t-1)}$  zum Incluster von  $J^{(t)}$  hinzu
```

```
    Behandle Evidenz  $e^{(t)}$  in  $J^{(t)}$ 
```

```
    Verschicke Nachrichten in  $J^{(t)}$ 
```

```
    Beantworte Anfragen  $Q^{(t)}$ 
```

```
    Berechne  $\alpha^{(t)}$ 
```

▸ J^0 bei $t = 0$; J^{\rightarrow} sonst

▸ falls $t > 0$

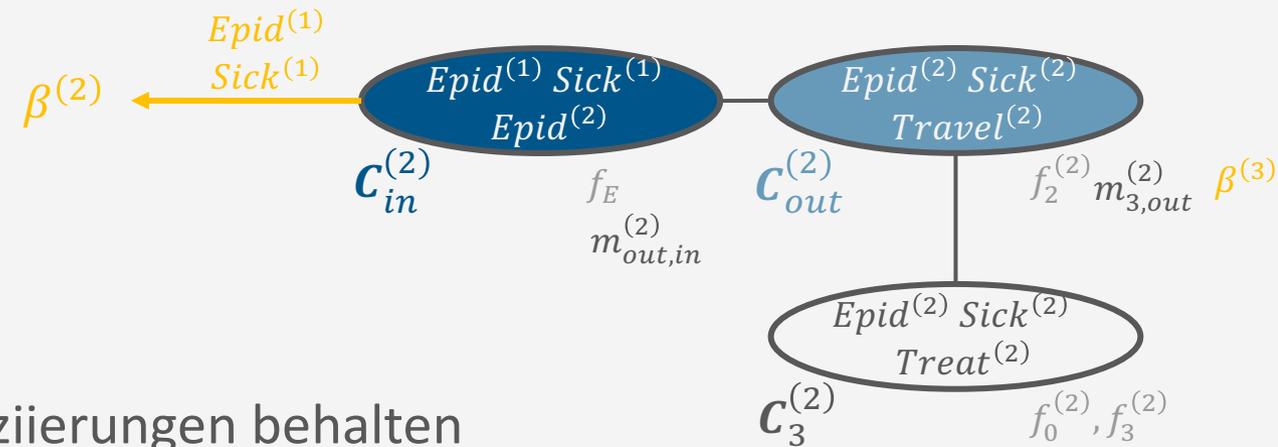
▸ Vorgehen auf vorheriger Folie

DJT

- Bei Online-Inferenz, zwei Ströme für eingehende Anfragen $Q^{(t)}$ und Evidenz $e^{(t)}$
 - Anstatt Eingaben $Q^{(0:T_q)}, e^{(0:T_e)}$

Vorgehensweisen beim Rückwärtsgehen

- Für beliebige Rückschau-Anfragen müssen Vorwärtsnachrichten gespeichert werden
 - Vorwärtsnachricht im Incluster gebraucht für Rückschau-Anfragen an diesen Zeitschritt
- Rückwärts gehen bedeutet zwischendrin Nachrichten erneut zu berechnen
 - Jeder Schritt t rückwärts ohne Anfrage
 - Berechne eine Nachrichtenphase zum Incluster mit $n - 1$ Nachrichten, n Anzahl an Cluster in J_t
 - Jeder Schritt t rückwärts mit Anfragen
 - Berechne vollständigen Nachrichtenversand mit $2(n - 1)$ Nachrichten
- Vorteil
 - Nur ein Jtree im Speicher zur Zeit
- Was, wenn DJT die Jtrees behalten hätte?
 - Bedarfsweise instanziiieren (bisher) vs. Instanziiierungen behalten

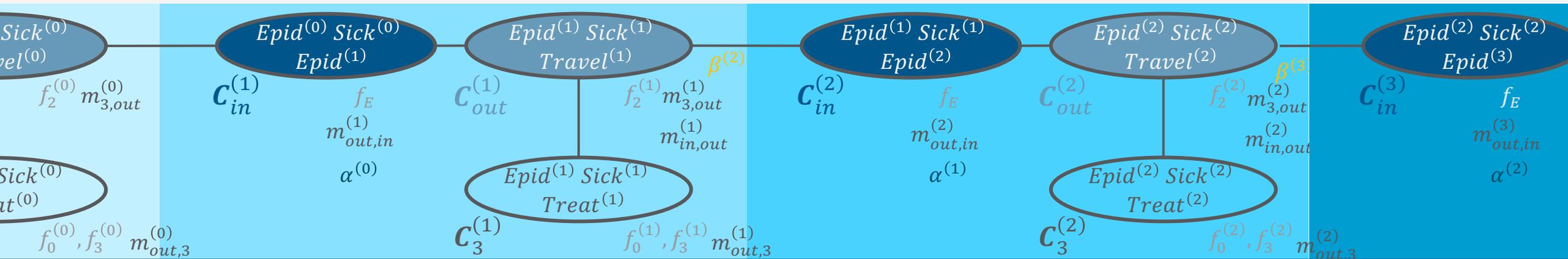


Instanziierungen behalten

- Was zahlt DJT an Speicher?
 - Lokale Modelle, Intra-Zeitscheiben-Nachrichten für jede Zeitscheibe speichern
- Wie viel Aufwand kann DJT einsparen?
 - Evidenz behandeln nicht nötig, da schon in lokalen Modellen absorbiert beim Vorwärtsgehen
 - Adaptiver Nachrichtenversand: Nur Nachrichten aktualisieren, die durch β_{t+1} geändert
 - Jeder Schritt t rückwärts ohne Anfrage
 - β_{t+1} neue Information, die beim Outcluster ankommt und zum Incluster geschafft werden muss
 - Berechne eine Nachrichten auf dem Pfad zwischen Outcluster und Incluster \rightarrow bis zu $n - 1$ Nachrichten, n Anzahl an Cluster in J_t (Worst Case: alle Cluster in Reihe mit In- und Outcluster an den jeweiligen Enden)
 - Jeder Schritt t rückwärts mit Anfragen
 - β_{t+1} neue Information, die beim Outcluster ankommt und zu allen anderen Clustern geschafft werden muss
 - Berechne eine Nachrichtenphase vom Outcluster zur Peripherie mit $n - 1$ Nachrichten, n Anzahl an Cluster in J_t

Beispiel: $P(\text{Treat}^{(1)} | e^{(0:3)})$ – Zu berechnende Nachrichten

- In $t = 2$ (keine Anfrage)
 - Instanziierungen behalten: $m_{out,in}^{(2)}$ (auf Pfad)
 - Bedarfsweise instanziiieren: $m_{3,out}^{(2)}, m_{out,in}^{(2)}$ (einwärts nach C_2^{in})
- In $t = 1$ (mit Anfrage)
 - Instanziierungen behalten: $m_{out,3}^{(1)}$ (auswärts nach $C_3^{(1)}$) + $m_{out,in}^{(1)}$ (für beliebige Anfragen)
 - Bedarfsweise instanziiieren: $m_{in,out}^{(1)}, m_{out,3}^{(1)}$ (nach $C_3^{(1)}$) + $m_{out,in}^{(1)}, m_{3,out}^{(1)}$ (für beliebige Anfragen)



Vorgehensweisen beim Rückwärtsgehen

	Instanziierungen behalten	Bedarfsweise instanziiieren
Schritte ohne Anfragen	$\leq n - 1$	$n - 1$
Schritte mit beliebigen Anfragen	$n - 1$	$2(n - 1)$
Zusätzlicher Speicher	Lokale Modelle, Intra-Zeitscheiben-Nachrichten, $\alpha^{(t)}$ Nachrichten	$\alpha^{(t)}$ Nachrichten

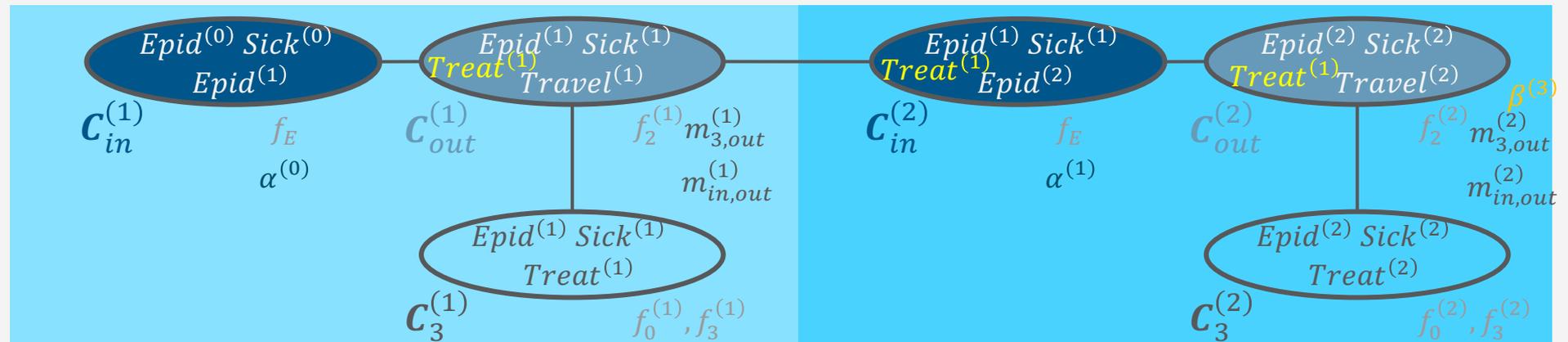
- Alle Instanziierungen im Speicher vorhalten nicht durchführbar
 - Ziel war ein speicher-effizienter Algorithmus, welcher nicht das gesamte Modell vom ersten Zeitschritt an im Speicher benötigt
- Laufzeit vs. Speicher gegeneinander abwägen
 - Letzte k Instanziierungen behalten
 - Daten / Anfragen für typischen Versatz in Anfragen analysieren
 - Bedarfsweise für Anfragen mit größerem Versatz instanziiieren

Wahrscheinlichkeitsanfragen über Zeitschritte hinweg

- Anfragevariablen aus unterschiedlichen Zeitscheiben
 - Beispiel: $P(\text{Treat}^{(1)}, \text{Travel}^{(2)} | \mathbf{e}^{(0:3)})$ bei $t = 3 \rightarrow \min = 1, \max = 2$
- Anfragebeantwortung (eine von vielen Möglichkeiten):
 - Rolle Jtrees für Zeitschritte aus, die in der Anfrage vorkommen: $J^{(\min:\max)}$
 - \min = minimale Zeitscheibe in Anfrage
 - \max = maximale Zeitscheibe in Anfrage
 - $\alpha^{(\min-1)}$ in Incluster von $J^{(\min)}$; $\beta^{(\max+1)}$ in Outcluster von $J^{(\max)}$, wenn $\max < t$
 - Füge Evidenz $\mathbf{e}^{(\min:\max)}$ ein
 - Füge alle Anfragevariablen zu Cluster \mathcal{C} mit größter Überlappung mit Anfragevariablen hinzu
 - Stelle *running intersection property* durch Erweitern von Clustern wieder her
 - Führe eine einwärts Nachrichtenphase mit \mathcal{C} als Zentrum durch und beantworte Anfrage in \mathcal{C}
 - Berechnet auch eigentliche α Nachrichten neu, da wir weniger Variablen eliminieren als sonst

Wahrscheinlichkeitsanfragen über Zeitschritte hinweg

- Anfragevariablen aus unterschiedlichen Zeitscheiben
 - Beispiel: $P(\text{Treat}^{(1)}, \text{Travel}^{(2)} | \mathbf{e}^{(0:3)})$ bei $t = 3 \rightarrow \min = 1, \max = 2$
- Anfragebeantwortung:
 - Jtrees ausrollen (1 : 2) mit $\alpha^{(0)}$ in $\mathcal{C}_{in}^{(1)}$ und $\beta^{(3)}$ in $\mathcal{C}_{out}^{(2)}$, Evidenz eingeben
 - Anfragevariablen zu $\mathcal{C}_{out}^{(2)}$ hinzufügen
 - Running intersection property herstellen
 - Nachrichten zu $\mathcal{C}_{out}^{(2)}$ senden
 - Anfrage beantw.



Sequentielles MPE

- Wahrscheinlichste Zuweisungen an alle Zufallsvariablen $V^{(0:T)}$ ohne Evidenz gegeben Evidenz $e^{(0:T)}$

$$MPE_M(e^{(0:T)}) = \arg \max_{v^{(0:t)} \in \text{Val}(V^{(0:T)})} P(v^{(0:T)} | e^{(0:T)})$$

- $V^{(0:T)} = R^{(0:T)} \setminus \text{rv}(e^{(0:T)})$
- Auch als laufende Anfrage $MPE_M(e^{(0:t)})$ möglich
- Vorgehen ähnlich zu Filterungsanfragen
 - Für jeden Zeitschritt $t = 0, \dots, T$
 - Jtree $J^{(t)}$ instanziiieren, Vorwärtsnachricht $\alpha^{(t-1)}$ zum Incluster von $J^{(t)}$ hinzufügen
 - Evidenz $e^{(t)}$ behandeln
 - Nachrichten mit **MPE-VE Operatoren** berechnen und zum Outcluster von $J^{(t)}$ schicken
 - Vorwärtsnachricht $\alpha^{(t)}$ mit **MPE-VE Operatoren** berechnen

Nur eine Nachrichtenphase zum Outcluster nötig (wie bei MPE-JT)



Sequentielles MAP (und MPE)

- Wahrscheinlichste Zuweisungen an Sequenzen von Zufallsvariablen $\mathbf{U}^{(t_1:t_2)}$ gegeben Evidenz $\mathbf{e}^{(0:t)}$

$$MAP_M(\mathbf{U}^{(t_1:t_2)} | \mathbf{e}^{(0:t)}) = \arg \max_{\mathbf{u}^{(t_1:t_2)} \in \text{Val}(\mathbf{U}^{(t_1:t_2)})} \sum_{\mathbf{v}^{(0:T)} \in \text{Val}(\mathbf{V}^{(0:T)})} P(\mathbf{u}^{(t_1:t_2)}, \mathbf{v}^{(0:T)} | \mathbf{e}^{(0:t)})$$

- $\mathbf{V}^{(0:T)} = \mathbf{R}^{(0:T)} \setminus \mathbf{U}^{(t_1:t_2)} \setminus \text{rv}(\mathbf{e}^{(0:t)}), T = \max\{t, t_1, t_2\}$
- Vorgehen wie bei MAP-JT
 - Jtrees über $t_1 : t_2$ ausrollen: $J^{(t_1:t_2)}$; $\alpha^{(t_1-1)}$ in Incluster von $J^{(t_1)}$; $\beta^{(t_2+1)}$ in Outcluster von $J^{(t_2)}$
 - Finde Subgraph über $\mathbf{U}^{(t_1:t_2)}$ in $J^{(t_1:t_2)}$
 - Schicke Nachrichten an den Rand des Subgraphen (mittels VE Operatoren berechnet)
 - Eliminiere $\mathbf{V}^{(0:T)}$ mit VE Operatoren im Modell des Subgraphen
 - Eliminiere $\mathbf{U}^{(t_1:t_2)}$ mit MPE-VE Operatoren im verbleibenden Modell

Sequentielles MAP (und MPE)

- Wahrscheinlichste Zuweisungen an Sequenzen von Zufallsvariablen $\mathbf{U}^{(t_1:t_2)}$ gegeben Evidenz $\mathbf{e}^{(0:t)}$

$$MAP_M(\mathbf{U}^{(t_1:t_2)} | \mathbf{e}^{(0:t)}) = \arg \max_{\mathbf{u}^{(t_1:t_2)} \in \text{Val}(\mathbf{U}^{(t_1:t_2)})} \sum_{\mathbf{v}^{(0:T)} \in \text{Val}(\mathbf{V}^{(0:T)})} P(\mathbf{u}^{(t_1:t_2)}, \mathbf{v}^{(0:T)} | \mathbf{e}^{(0:t)})$$

- $\mathbf{V}^{(0:T)} = \mathbf{R}^{(0:T)} \setminus \mathbf{U}^{(t_1:t_2)} \setminus \text{rv}(\mathbf{e}^{(0:t)}), T = \max\{t, t_1, t_2\}$
- MAP Anfragen über die letzten k Zeitschritte $(t - k) : t =$ MPE Anfrage über Zeitschritte $(t - k) : t$ (wie bei MAP Anfragen an vollständige Subgraphen bei MAP-JT)
 - Bis $t - k$ Nachrichten mittels VE Operatoren berechnen
 - Ab $t - k$ Vorgehen wie bei MPE Anfragen
 - Nachrichten mit **MPE-VE Operatoren** berechnen und zum Outcluster von $J^{(t)}$ schicken
 - Vorwärtsnachricht $\alpha^{(t)}$ mit **MPE-VE Operatoren** berechnen

Komplexität

- JT Komplexität für Nachrichtenversand und Anfragebeantwortung (einfache Anfragen)

$$O_{MP} = O(n_J \cdot r^w)$$

$$O_{QA} = O(r^w)$$

Größt mögliches Zwischenergebnis

- n_T, n_J Anzahl der Knoten im Dtree / Jtree
- r größte Domäne
- w Baumweite (größtes Cluster)
- DJT: JT + mit der Zeit gehen
 - Nachrichtenversand innerhalb einer Zeitscheibe $\rightarrow O_{MP}$
 - Vorwärtsnachricht eine Anfrage an Outcluster $\rightarrow O_{QA}$
 - Rückwärtsnachricht eine Anfrage an Incluster $\rightarrow O_{QA}$



Was sind die Best-Case- und Worst-Case-Fälle bzgl. Anfragen?

Komplexität

- Gegeben maximaler Zeitschritt T , der in allen Anfragen vorkommt
 - m_t Anzahl an Anfragen pro Zeitschritt t , alle Anfragen: $M = \sum_{t=0}^T m_t$, durchschnittliche Anfragen pro Zeitschritt: $m = \frac{1}{T} \sum_{t=0}^T m_t$
 - Vorwärts gehen: Informationen verteilen + Vorwärtsnachricht, i.e., $T \cdot (O_{MP} + O_{QA})$
 - **Best-Worst-Case** bzgl. Anfragen im Zeitschritt $t \in \{0, \dots, T\}$
 - Filterungsanfrage zu t , kostet je O_{QA}
 - Zusammen,

$$\begin{aligned}
 & T \cdot (O_{MP} + O_{QA}) + M \cdot O_{QA} = (T \cdot n_J + T + M) \cdot O(r^w) \\
 & = O\left((T \cdot n_J + M) \cdot r^w\right) = O\left((T \cdot n_J + T \cdot m) \cdot r^w\right)
 \end{aligned}$$

Vergleiche JT Komplexität: $O\left((n_J + m) \cdot r^w\right)$

Komplexität

- **Worst-Worst-Case** bzgl. Anfragen im Zeitschritt $t \in \{0, \dots, T\}$
 - Rückschau-Anfrage für $\pi = 0$
 - Rückwärts-Nachrichtenphase zu $\pi = 0$: $t \cdot (O_{MP} + O_{QA})$
 - Vorhersagen-Anfrage für $\pi = T$
 - Vorwärts-Nachrichtenphase zu $\pi = T$: $(T - \pi) \cdot (O_{MP} + O_{QA})$
 - Zusammen führen eine Rückschau- und eine Vorhersagen-Anfragen zu einer Komplexität von

$$\pi \cdot (O_{MP} + O_{QA}) + (T - \pi) \cdot (O_{MP} + O_{QA})$$

$$= T \cdot (O_{MP} + O_{QA})$$
 - Für m Anfragen pro Zeitscheibe haben wir

$$T \cdot (O_{MP} + O_{QA}) + m \cdot O_{QA}$$
- Für T Zeitschritte haben wir

$$T \cdot T \cdot (O_{MP} + O_{QA}) + T \cdot m \cdot O_{QA} = T^2 \cdot (O_{MP} + O_{QA}) + M \cdot O_{QA}$$

Komplexität

- **Worst-Worst-Case**

- Vorwärts gehen: Informationen verteilen + Vorwärtsnachricht, i.e., $T \cdot (O_{MP} + O_{QA})$
- Anfragen beantworten

$$T^2 \cdot (O_{MP} + O_{QA}) + M \cdot O_{QA}$$

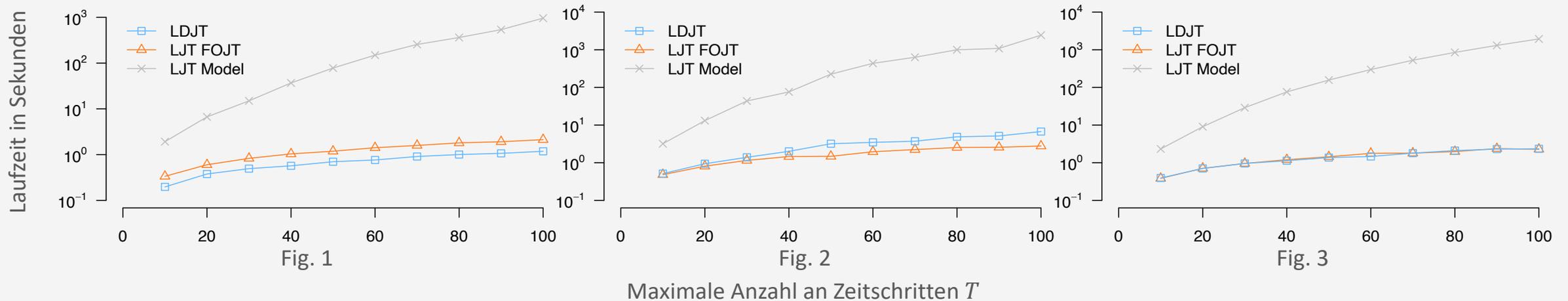
- Zusammen,

$$\begin{aligned} & T \cdot (O_{MP} + O_{QA}) + T^2 \cdot (O_{MP} + O_{QA}) + M \cdot O_{QA} \\ = & (T \cdot n_J + T + T^2 \cdot n_J + T^2 + M) \cdot O_{QA} \\ = & O \left(\left((T^2 + T) \cdot n_J + M \right) \cdot r^w \right) \\ = & O \left(\left((T^2 + T) \cdot n_J + T \cdot m \right) \cdot r^w \right) \end{aligned}$$

Vergleiche Komplexität Filterungsanfragen: $O \left((T \cdot n_J + T \cdot m) \cdot r^w \right)$

Laufzeiten

- Algorithmen im Vergleich: LDJT, LJT FOJT: Inferenz auf ausgerolltem Jtree mittels JT, LJT Modell: Inferenz auf ausgerolltem Modell mittels JT
 - Das „L“ (von Lifted) bitte ignorieren; gezeigtes Verhalten gilt so für die propositionalen Algorithmen hier
 - LJT FOJT ähnlich zu LDJT, da die gleichen Berechnungen durchgeführt werden, nur LJT FOJT hält das gesamte Modell im Speicher (klein genug, dass das geht)
- Fig. 1: Filterungsanfragen, Fig. 2: Rückschau/ Vorhersage-Anfragen zu $0, T$, Fig. 3: Rückschau-Anfragen (Vorhersage-Anfragen mit ähnlichen Laufzeiten)



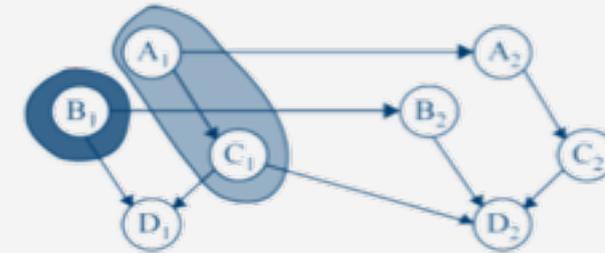
Approximationen

• Boyen-Koller Algorithmus

- Annahme: Unabhängigkeiten im Interface
 - Partitionierung des Interfaces als Eingabe gefordert
 - Beispiel: $\{\{A, C\}, \{B\}\}$ in einem Interface $\{A, B, C\}$
- Anfrage nicht mehr über das gesamte Interface, sondern als Produkt der unabhängigen Teile
 - Inter-Zeitscheiben-Nachrichten: Vereinigung der Resultate der Anfragen über die einzelnen Partitionen
 - Beispiel: VE aufrufen mit dem lokalen Modell und Nachrichten von $\mathcal{C}_{out}^{(t)}$ und der Anfrage mit Anfragevariablen $\{A, C\}$ und der Anfrage mit Anfragevariable $\{B\}$:

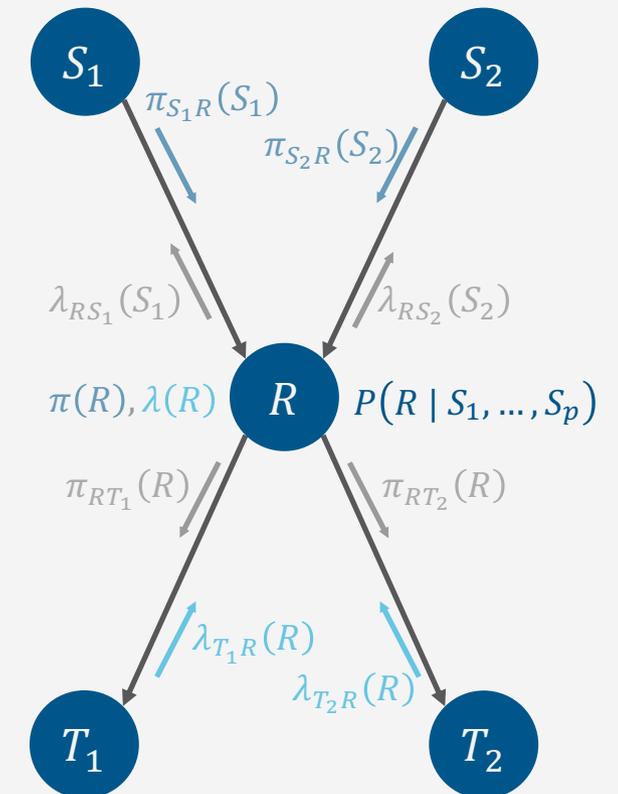
$$\alpha^{(t)} = \left\{ \text{VE} \left(F_{out}^{(t)} \cup \bigcup_{j \in \text{Nb}(\mathcal{C}_{out}^{(t)})} m_{j,out}^{(t)}, \{A, C\}, \emptyset \right), \text{VE} \left(F_{out}^{(t)} \cup \bigcup_{j \in \text{Nb}(\mathcal{C}_{out}^{(t)})} m_{j,out}^{(t)}, \{B\}, \emptyset \right) \right\}$$

- Exakter Algorithmus, wenn Unabhängigkeiten tatsächlich gelten, sonst approximativ



Approximationen

- **Factored-frontier Algorithmus** [Kevin Murphys Dissertation, 2002]
 - Zur Erinnerung:
 - Probability Propagation auf Polytree BNs: Nachrichten direkt zwischen den Knoten versandt
 - Probability Propagation in allgemeinen Graphen \rightarrow (loopy) belief propagation: Nachrichten direkt zwischen den Knoten versandt
 - Exakt bei Polytrees, sonst approximativ
 - Kerngedanke: Nutze loopy belief propagation über die Zeit



Topic Modellierung über die Zeit: Dynamic Topic Modell

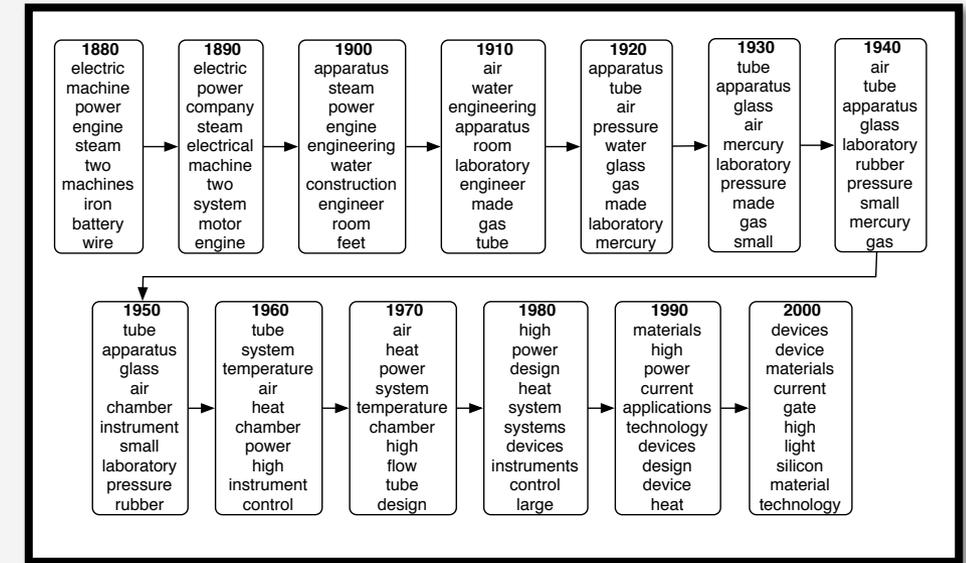
- Wie auch Standard LDA, *Sampling-basiert*
- Beispiel: Journal Science (1880-2002) JSTOR Corpus
 - Vokabular beschränkt auf die 30.000 Terme, die mehr als 10 mal vorkommen
 - Ca. 76 Mio. Worte in rund 130.000 Dokumenten

TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McMurray

Genome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encodes all of the information necessary for the life of the organism. The base sequence contains four nucleotides—adenine, thymidine, guanosine, and cytosine—which are linked together into long double-helical chains. Over the last two decades, automated DNA sequencers have made the process of obtaining the base-by-base sequence of DNA...



Zwischenzusammenfassung

- Interfaces separieren die Vergangenheit von der Gegenwart und die Gegenwart von der Zukunft
 - Inter-Zeitscheiben-Nachrichten
 - Vorwärtsnachricht transferiert alle Informationen aus der Vergangenheit bis zur aktuellen Zeitscheibe einschließlich zur nächsten Zeitscheibe
 - Rückwärtsnachricht transferiert alle Informationen aus der Zukunft bis zur aktuellen Zeitscheibe einschließlich zur vorherigen Zeitscheibe
- DJT Algorithmus
 - JT für Intra-Zeitscheiben-Anfragebeantwortung
 - Inter-Zeitscheiben-Nachrichten um sich in der Zeit zu bewegen
 - Komplexität: Abhängig von der maximalen Zeitscheibe T
 - Approximationen über Nachrichten direkt auf dem Graph oder Unabhängigkeiten im Interface

Überblick: 6. Sequentielle PGMs und Inferenz

A. *Sequentielle PGMs*

- Templates, dynamische BNs, dynamische Faktormodelle, Hidden-Markov-Modelle; Semantik
- Inferenzaufgaben: Wahrscheinlichkeitsanfragen (Filtering, Prediction, Hindsight), Zustandsanfragen (MPE, MAP)

B. *Sequentielle Inferenz*

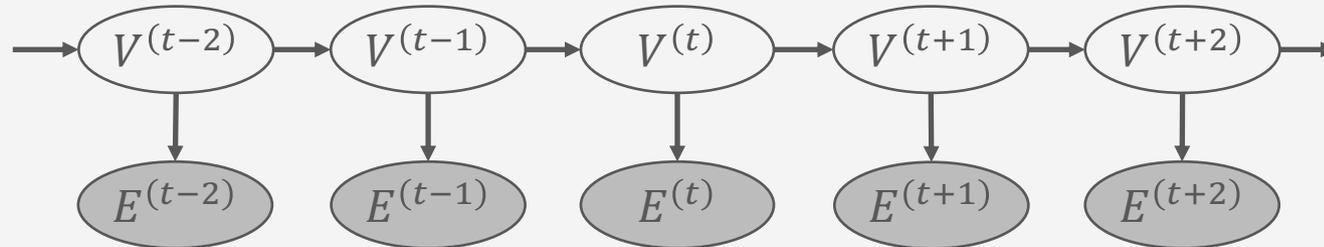
- Naïve Inferenz mittels Ausrollen, Interface Algorithmus, Komplexität, Approximationen

C. *Spezialfall Hidden-Markov-Modelle*

- Viterbi-Algorithmus für MPEs
- Anfragebeantwortung durch Matrixoperationen
- Baum-Welch-Algorithmus zum Lernen

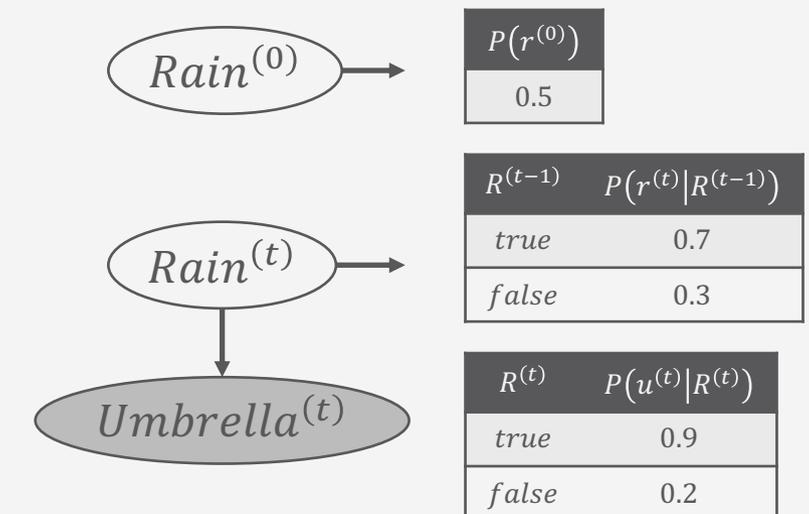
Wiederholt: Hidden Markov Modell (HMM)

- Spezialfall des DBN mit $R = \{V, E\}$ bzw. $V = \{V\}$, $E = \{E\}$
 - Latente Zufallsvariable V
 - Evidenzvariable E
 - Modell:
 - $B^0 = P(V^{(0)})$
 - $B^\rightarrow = \{P(V^{(\tau)}|V^{(\tau-1)}), P(E^{(\tau)}|V^{(\tau)})\}$
 - Repräsentierte Verteilung: $P_{V,E}^T = P(V^{(0)}) \prod_{t=1}^T P(V^{(t)}|V^{(t-1)})P(E^{(t)}|V^{(t)})$



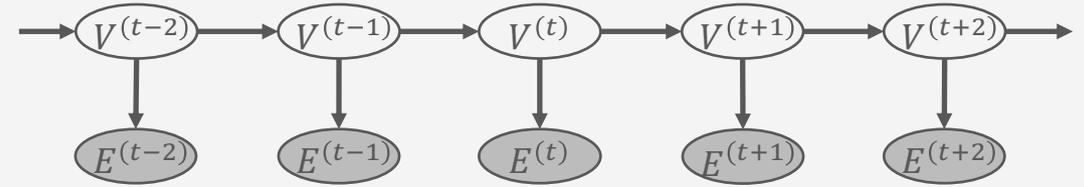
Ein DBN mit $R = V \cup E$ kann jederzeit in ein HMM umgewandelt werden:

- Mit jedem $v \in \text{Val}(V)$ als abstrakten Wert v , $\text{Val}(V) = \text{Val}(V)$
 - Mit jedem $e \in \text{Val}(E)$ als abstrakten Wert e , $\text{Val}(E) = \text{Val}(E)$
- $\text{Val}(V)$ und $\text{Val}(E)$ entsprechend groß!



Filterungsanfrage

- $P(V^{(t)} | e^{(1:t)})$
 - Verteilung über den Zustand der Umgebung → Grundlage für Agent um rationale Entscheidung zu treffen
 - Bei HMMs: $P(V^{(t)} | e^{(1:t)})$
- Was passiert bei den Berechnungen?
 - Gegeben den Ergebnissen der Filterung bis zum Zeitpunkt t , Ergebnis für Zeitpunkt $t + 1$ aus der neuen Evidenz $e^{(t+1)}$ berechnen



$$\begin{aligned}
 P(V^{(t+1)} | e^{(0:t+1)}) &= f(e^{(t+1)}, P(V^{(t)} | e^{(0:t+1)})) \\
 &= P(V^{(t+1)} | e^{(0:t)}, e^{(t+1)}) \\
 &\propto P(e^{(t+1)} | V^{(t+1)}, e^{(0:t)}) P(V^{(t+1)} | e^{(0:t)}) \\
 &\propto \underbrace{P(e^{(t+1)} | V^{(t+1)})}_{\text{Aktualisierung der Ein-Schritt-Vorhersage mit Evidenz}} \underbrace{P(V^{(t+1)} | e^{(0:t)})}_{\text{Ein-Schritt-Vorhersage}}
 \end{aligned}$$

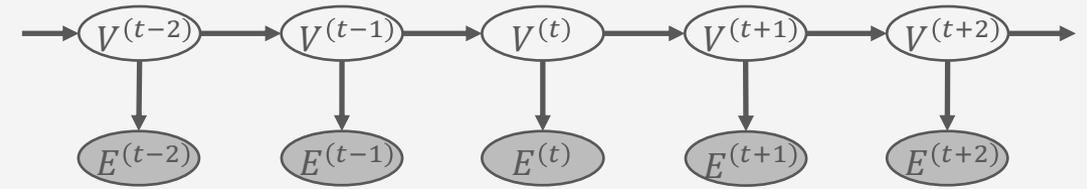
(für eine Funktion f)

(Evidenz aufteilen)

(durch Bayes Theorem)

(durch Markov Annahme)

Filterungsanfrage



- Was passiert bei den Berechnungen?
- Gegeben den Ergebnissen der Filterung bis zum Zeitpunkt t , Ergebnis für Zeitpunkt $t + 1$ aus der neuen Evidenz $e^{(t+1)}$ berechnen

$$\begin{aligned}
 \underbrace{P(\mathbf{V}^{(t+1)} | \mathbf{e}^{(0:t+1)})}_{\text{Vorwärtsnachricht } \alpha^{(t+1)}} &\propto P(\mathbf{e}^{(t+1)} | \mathbf{V}^{(t+1)}) P(\mathbf{V}^{(t+1)} | \mathbf{e}^{(0:t)}) \\
 &\propto P(\mathbf{e}^{(t+1)} | \mathbf{V}^{(t+1)}) \sum_{\mathbf{v}^{(t)} \in \text{Val}(\mathbf{V}^{(t)})} P(\mathbf{V}^{(t+1)} | \mathbf{v}^{(t)}, \mathbf{e}^{(0:t)}) P(\mathbf{v}^{(t)} | \mathbf{e}^{(0:t)}) \quad (\text{Konditionieren auf } \mathbf{V}^{(t)}) \\
 &\propto P(\mathbf{e}^{(t+1)} | \mathbf{V}^{(t+1)}) \sum_{\mathbf{v}^{(t)} \in \text{Val}(\mathbf{V}^{(t)})} \underbrace{P(\mathbf{V}^{(t+1)} | \mathbf{v}^{(t)})}_{\text{Übergangsmodell}} \underbrace{P(\mathbf{v}^{(t)} | \mathbf{e}^{(0:t)})}_{\text{Filterungsergebnis bis zum Zeitpunkt } t = \text{Vorwärtsnachricht } \alpha^{(t)}} \quad (\text{durch Markov Annahme})
 \end{aligned}$$

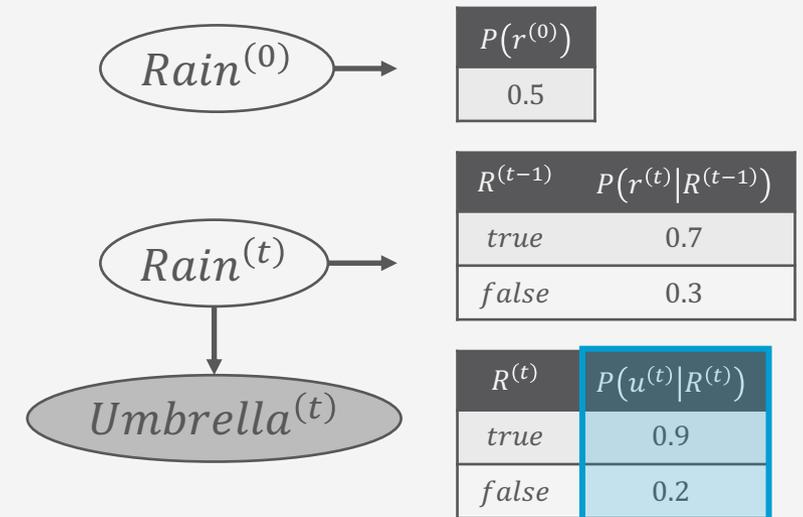
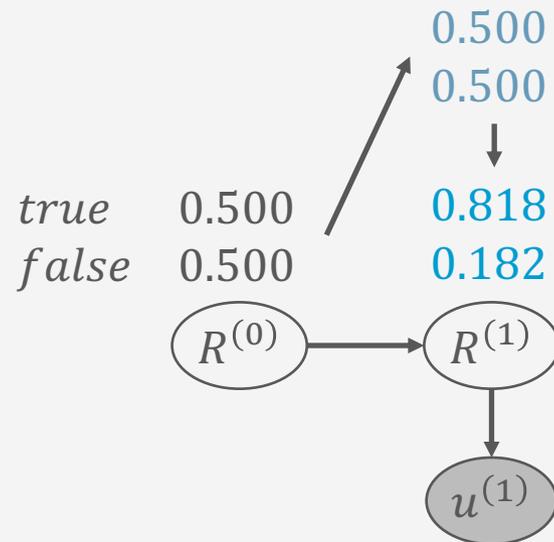
Vorwärtsnachricht $\alpha^{(t+1)}$
 = FORWARD $(\alpha^{(t)}, e^{(t+1)})$

Übergangsmodell

Filterungsergebnis bis zum Zeitpunkt t
 = Vorwärtsnachricht $\alpha^{(t)}$

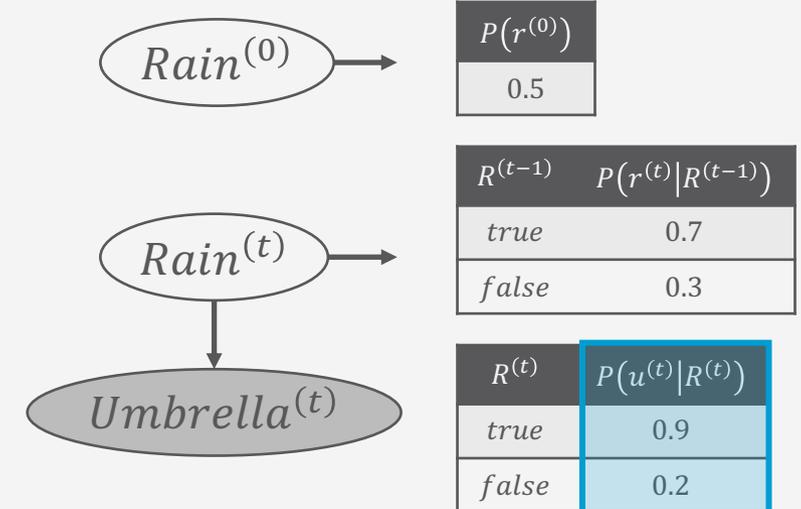
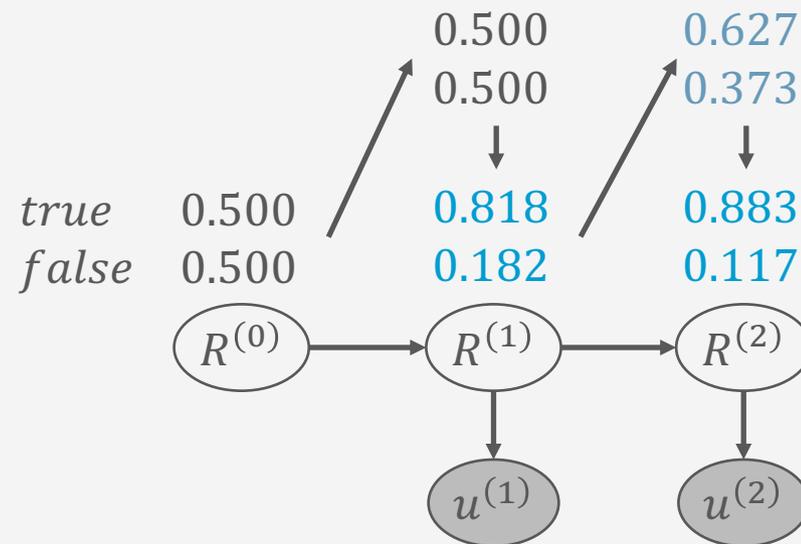
Beispiel

- Tag $t = 1$: Regenschirm beobachtet $\rightarrow u^{(1)}$
- Vorhersage von $t = 0$ nach $t = 1$ ist: $P(R^{(1)}) = \sum_{r^{(0)} \in \text{Val}(R^{(0)})} P(R^{(1)} | r^{(0)}) P(r^{(0)})$
- Aktualisieren mit der Evidenz von $t = 1$ ergibt: $P(R^{(1)} | u^{(1)}) \propto P(u^{(1)} | R^{(1)}) P(R^{(1)})$



Beispiel

- Tag $t = 2$: Regenschirm beobachtet $\rightarrow u^{(2)}$
- Vorhersage von $t = 1$ nach $t = 2$ ist: $P(R^{(2)}|u^{(1)}) = \sum_{r^{(1)} \in \text{Val}(R^{(1)})} P(R^{(2)}|r^{(1)})P(r^{(1)}|u^{(1)})$
- Aktualisieren mit der Evidenz von $t = 2$ ergibt: $P(R^{(2)}|u^{(1:2)}) \propto P(u^{(2)}|R^{(2)})P(R^{(2)}|u^{(1)})$



Vorhersagen und Rückschauen

- Vorhersage $P(\mathbf{V}^{(t+k)} | \mathbf{e}^{(0:t)})$
 - Filterungsanfrage ohne Evidenz von t bis $t + k$
 - Berechnungstechnisch passiert das gleiche wie auf den Folien zuvor
- Rückschau $P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t)})$
 - Bessere Schätzung des Zustandes zum Zeitpunkt $t - k$, da jetzt mehr Evidenz vorhanden
 - Berechnungen: Evidenz aufteilen in $\mathbf{e}^{(0:t-k)}$, $\mathbf{e}^{(t-k+1:t)}$

$$\begin{aligned}
 P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t)}) &= P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}, \mathbf{e}^{(t-k+1:t)}) && \text{(Evidenz aufteilen)} \\
 &\propto P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}) P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}, \mathbf{e}^{(0:t-k)}) && \text{(durch Bayes Theorem)} \\
 &\propto P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}) P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}) && \text{(durch Markov Annahme)} \\
 &\propto \alpha^{(t-k)} \beta^{(t-k+1)}
 \end{aligned}$$

- Vorwärtsnachricht abspeichern
- Rückwärtsnachricht mittels Rückwärtsrekursion

Vorhersagen und Rückschauen → Vorwärts-Rückwärts-Algorithmus

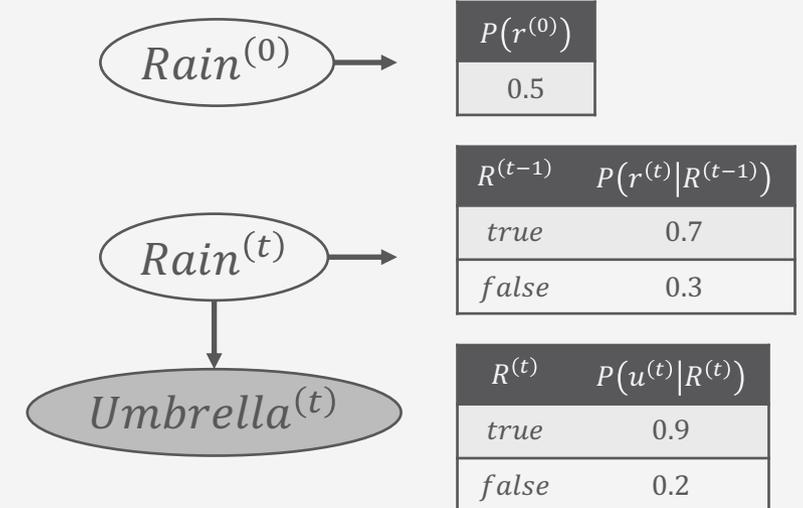
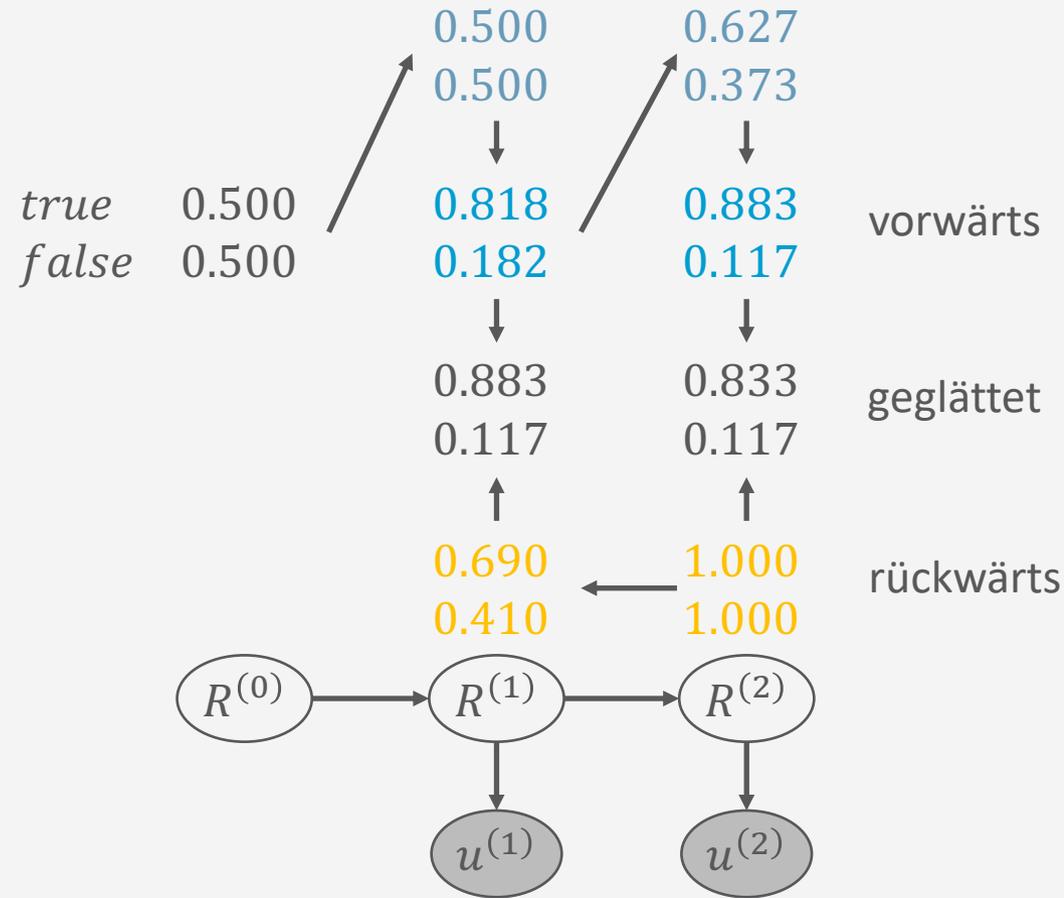
- Berechnungen: Evidenz aufteilen in $\mathbf{e}^{(0:t-k)}$, $\mathbf{e}^{(t-k+1:t)}$

$$\begin{aligned}
 P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t)}) &= P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}, \mathbf{e}^{(t-k+1:t)}) && \text{(Evidenz aufteilen)} \\
 &\propto P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}) P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}, \mathbf{e}^{(0:t-k)}) && \text{(durch Bayes Theorem)} \\
 &\propto P(\mathbf{V}^{(t-k)} | \mathbf{e}^{(0:t-k)}) P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}) && \text{(durch Markov Annahme)} \\
 &\propto \alpha^{(t-k)} \beta^{(t-k+1)}
 \end{aligned}$$

- Vorwärtsnachricht abspeichern
- Rückwärtsnachricht mittels Rückwärtsrekursion

$$\begin{aligned}
 P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}) &= \sum_{\mathbf{v}^{(t-k+1)}} P(\mathbf{e}^{(t-k+1:t)} | \mathbf{V}^{(t-k)}, \mathbf{v}^{(t-k+1)}) P(\mathbf{v}^{(t-k+1)} | \mathbf{V}^{(t-k)}) && \text{(Konditionieren auf } \mathbf{V}^{(t-k+1)}) \\
 &\propto \sum_{\mathbf{v}^{(t-k+1)}} P(\mathbf{e}^{(t-k+1:t)} | \mathbf{v}^{(t-k+1)}) P(\mathbf{v}^{(t-k+1)} | \mathbf{V}^{(t-k)}) && \text{(durch Markov Annahme)} \\
 &\propto \sum_{\mathbf{v}^{(t-k+1)}} P(\mathbf{e}^{(t-k+1)} | \mathbf{v}^{(t-k+1)}) P(\mathbf{e}^{(t-k+2:t)} | \mathbf{v}^{(t-k+1)}) P(\mathbf{v}^{(t-k+1)} | \mathbf{V}^{(t-k)}) && \text{(Kettenregel)}
 \end{aligned}$$

Beispiel



MPE

- MPE Definition

$$MPE(\mathbf{e}^{(1:t)}) = \arg \max_{\mathbf{v}^{(1:t)}} P(\mathbf{v}^{(1:t)} | \mathbf{e}^{(1:t)})$$

- Bei HMMs: $\arg \max_{\mathbf{v}^{(1:t)}} P(\mathbf{v}^{(1:t)} | \mathbf{e}^{(1:t)})$
- Was passiert bei den Berechnungen?

- Wahrscheinlichster Pfad zu jedem $\mathbf{v}^{(t+1)}$ = wahrscheinlichster Pfad zu einem $\mathbf{v}^{(t)}$ plus ein weiterer Schritt

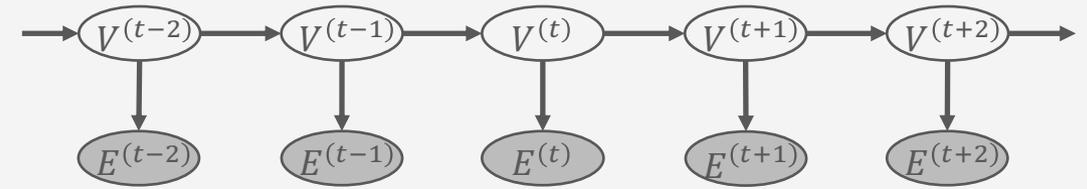
Übergangsmodell

MPE-Vorwärtsnachricht $\alpha^{(t)}$

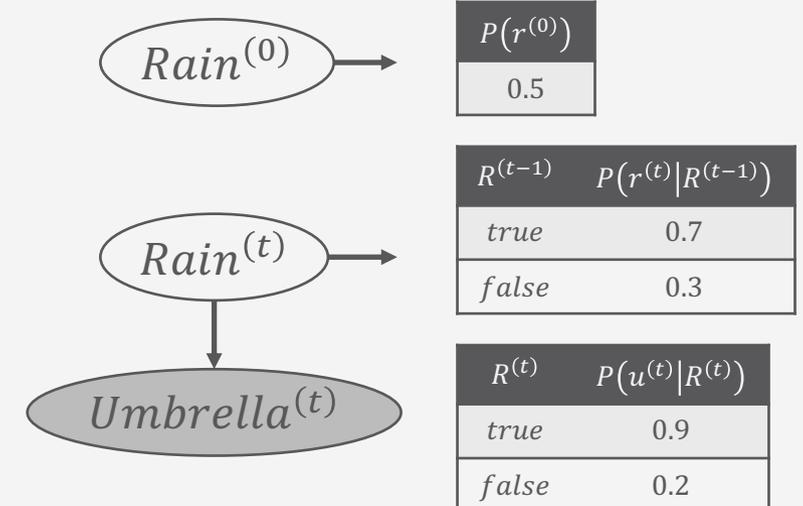
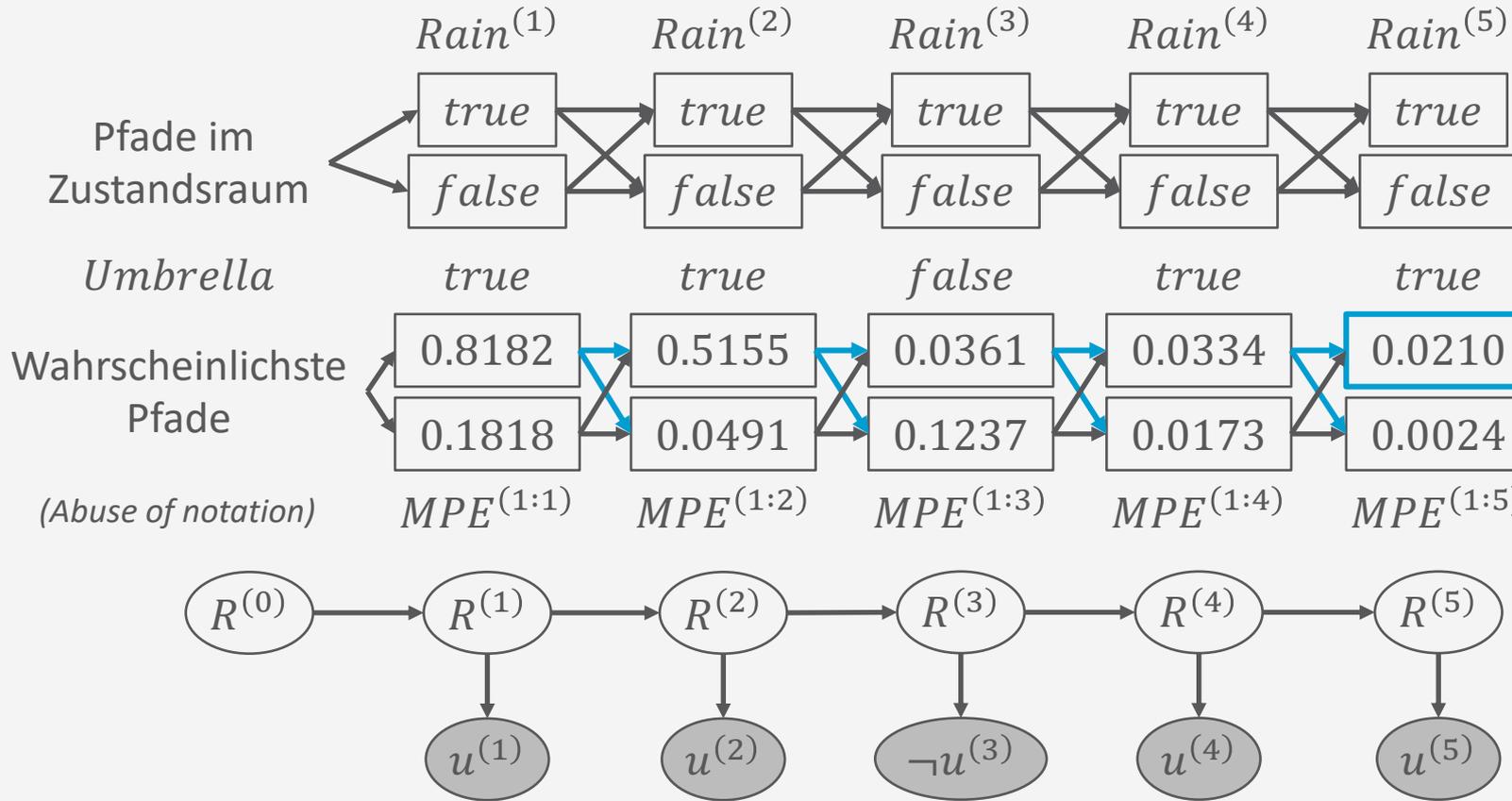
$$\begin{aligned} \max_{\mathbf{v}^{(1:t)}} P(\mathbf{v}^{(1:t)}, \mathbf{V}^{(t+1)} | \mathbf{e}^{(1:t)}) &= P(\mathbf{e}^{(t+1)} | \mathbf{V}^{(t+1)}) \max_{\mathbf{v}^{(t)}} \left(P(\mathbf{V}^{(t+1)} | \mathbf{v}^{(t)}) \max_{\mathbf{v}^{(1:t-1)}} P(\mathbf{v}^{(1:t-1)}, \mathbf{v}^{(t)} | \mathbf{e}^{(1:t)}) \right) \\ &= P(\mathbf{e}^{(t+1)} | \mathbf{V}^{(t+1)}) \max_{\mathbf{v}^{(t)}} (P(\mathbf{V}^{(t+1)} | \mathbf{v}^{(t)}) \cdot \alpha^{(t)}) \end{aligned}$$

MPE-Vorwärtsnachricht $\alpha^{(t+1)}$
= MPE-FORWARD $(\alpha^{(t)}, \mathbf{e}^{(t+1)})$

Viterbi-Algorithmus für HMMs



Beispiel



Rechnungen als Matrixoperationen
 → Effiziente Implementierung möglich

HMMs als Matrizen

- Latente Zufallsvariable V , $|\text{Val}(V)| = n$
- Evidenzvariable E
- **Übergangsmatrix** \mathfrak{I} eine $n \times n$ Matrix mit
 - $\mathfrak{I}_{ij} = P(V^{(t)} = v_i | V^{(t-1)} = v_j)$
- **Sensormatrix** $\mathfrak{D}^{(t)}$ für jeden Zeitschritt t gemäß Beobachtung $e^{(t)}$ eine $n \times n$ Matrix
 - Diagonale Einträgen $P(e^{(t)} | V^{(t)} = v_i)$

- Vorwärts- und Rückwärtsnachrichten als Spaltenvektoren

$$\alpha^{(t+1)} \propto \mathfrak{D}^{(t+1)} \mathfrak{I}^T \alpha^{(t)}$$

$$\beta^{(t-k+1)} \propto \mathfrak{I} \mathfrak{D}^{(t-k+1)} \beta^{(t-k+2)}$$

- Beispiel:

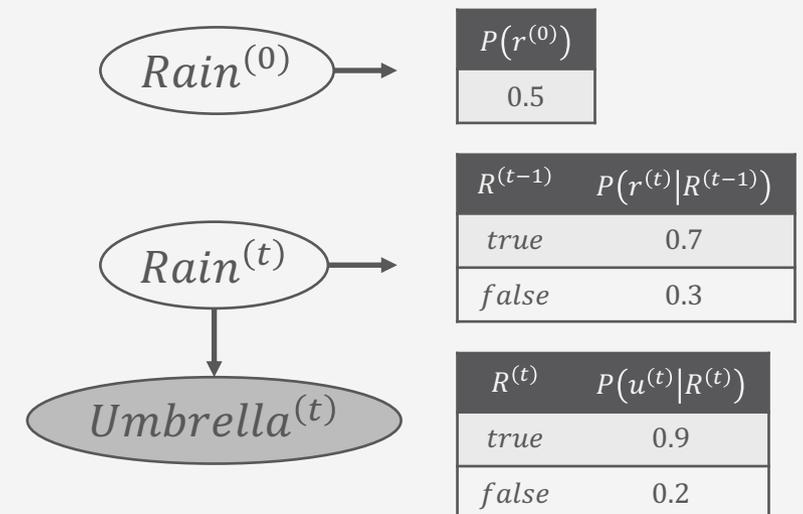
- Übergangsmatrix:

		alter Zustand v_j	
neuer Zustand v_i :	t	0.7	0.3
	f	0.3	0.7

- Sensormatrix:

$$U^{(1)} = \text{true}, \mathfrak{D}^{(1)} = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}$$

$$U^{(3)} = \text{false}, \mathfrak{D}^{(3)} = \begin{pmatrix} 0.1 & 0 \\ 0 & 0.8 \end{pmatrix}$$



Country Dance Algorithmus für HMMs

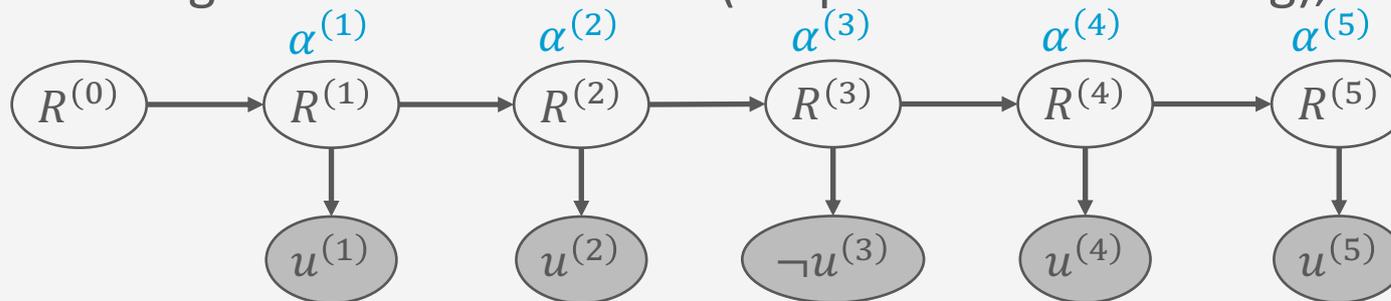
- Vorwärtsnachrichten nicht abspeichern zum Glätten, sondern Nachrichten bedarfsweise neu berechnen, indem man den Vorwärtspass rückwärts laufen lässt → Matrixinvertierung

$$\alpha^{(t+1)} = \frac{1}{Z^{(t+1)}} \mathfrak{D}^{(t+1)} \mathfrak{I}^T \alpha^{(t)}$$

$$\mathfrak{D}^{(t+1)^{-1}} \alpha^{(t+1)} = \frac{1}{Z^{(t+1)}} \mathfrak{I}^T \alpha^{(t)}$$

$$\frac{1}{Z^{(t+1)}} \mathfrak{I}^{T^{-1}} \mathfrak{D}^{(t+1)^{-1}} \alpha^{(t+1)} = \alpha^{(t)}$$

- Vorwärts gehen berechnet $\alpha^{(t)}$ (abspeichern nicht nötig), rückwärts gehen berechnet $\alpha^{(t')}, \beta^{(t')}$



Speichereffizient

$$\frac{1}{Z'} \mathfrak{I}^{\top-1} \mathfrak{D}^{(t+1)-1} \alpha^{(t+1)} = \alpha^{(t)} \rightarrow \frac{1}{Z'} \mathfrak{I}^{\top-1} \mathfrak{D}^{(2)-1} \alpha^{(2)} = \alpha^{(1)}$$

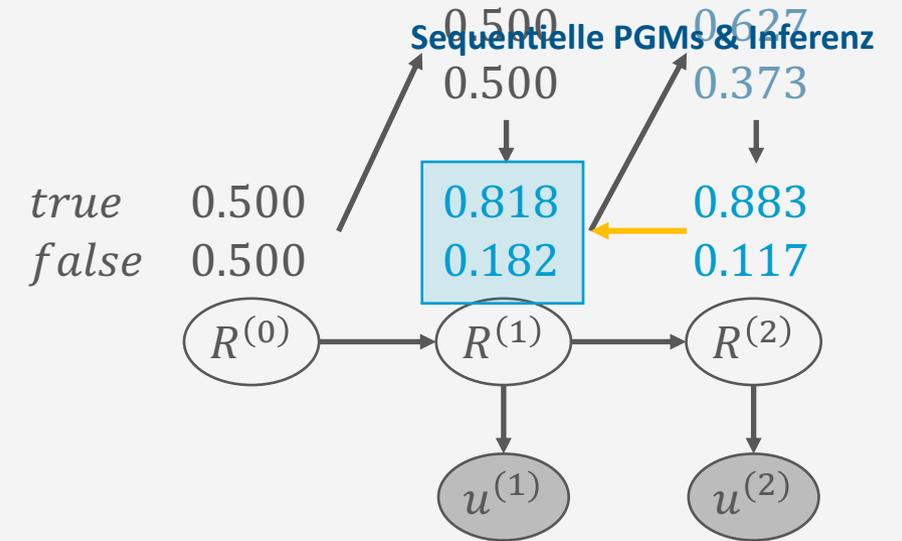
$$U^{(2)} = true, \mathfrak{D}^{(2)} = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}$$

$$\mathfrak{I} = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix} \rightarrow \mathfrak{I}^{\top} = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}$$

$$\mathfrak{D}^{(2)-1} = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}^{-1} = \frac{1}{0.9 \cdot 0.2 - 0 \cdot 0} \begin{pmatrix} 0.2 & 0 \\ 0 & 0.9 \end{pmatrix} = 5.56 \begin{pmatrix} 0.2 & 0 \\ 0 & 0.9 \end{pmatrix} = \begin{pmatrix} 1.1 & 0 \\ 0 & 5.0 \end{pmatrix}$$

$$\mathfrak{I}^{\top-1} = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}^{-1} = \frac{1}{0.49 - 0.09} \begin{pmatrix} 0.7 & -0.3 \\ -0.3 & 0.7 \end{pmatrix} = 2.5 \begin{pmatrix} 0.7 & -0.3 \\ -0.3 & 0.7 \end{pmatrix} = \begin{pmatrix} 1.75 & -0.75 \\ -0.75 & 1.75 \end{pmatrix}$$

$$\frac{1}{Z'} \begin{pmatrix} 1.75 & -0.75 \\ -0.75 & 1.75 \end{pmatrix} \begin{pmatrix} 1.1 & 0 \\ 0 & 5.0 \end{pmatrix} \begin{pmatrix} 0.883 \\ 0.117 \end{pmatrix} = \frac{1}{Z'} \begin{pmatrix} 1.944 & -3.75 \\ -0.825 & 8.75 \end{pmatrix} \begin{pmatrix} 0.883 \\ 0.117 \end{pmatrix} = \frac{1}{Z'} \begin{pmatrix} 1.278 \\ 0.295 \end{pmatrix} = \begin{pmatrix} 0.812 \\ 0.188 \end{pmatrix} = \alpha^{(1)}$$



$P(r^{(0)})$
0.5

$R^{(t-1)}$	$P(r^{(t)} R^{(t-1)})$
true	0.7
false	0.3

$R^{(t)}$	$P(u^{(t)} R^{(t)})$
true	0.9
false	0.2

$$A^{-1} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

Lernen von temporalen Modellen

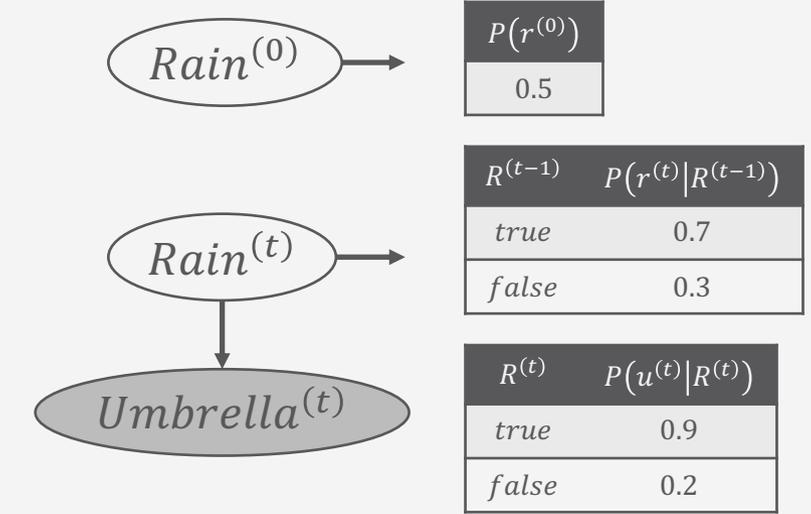
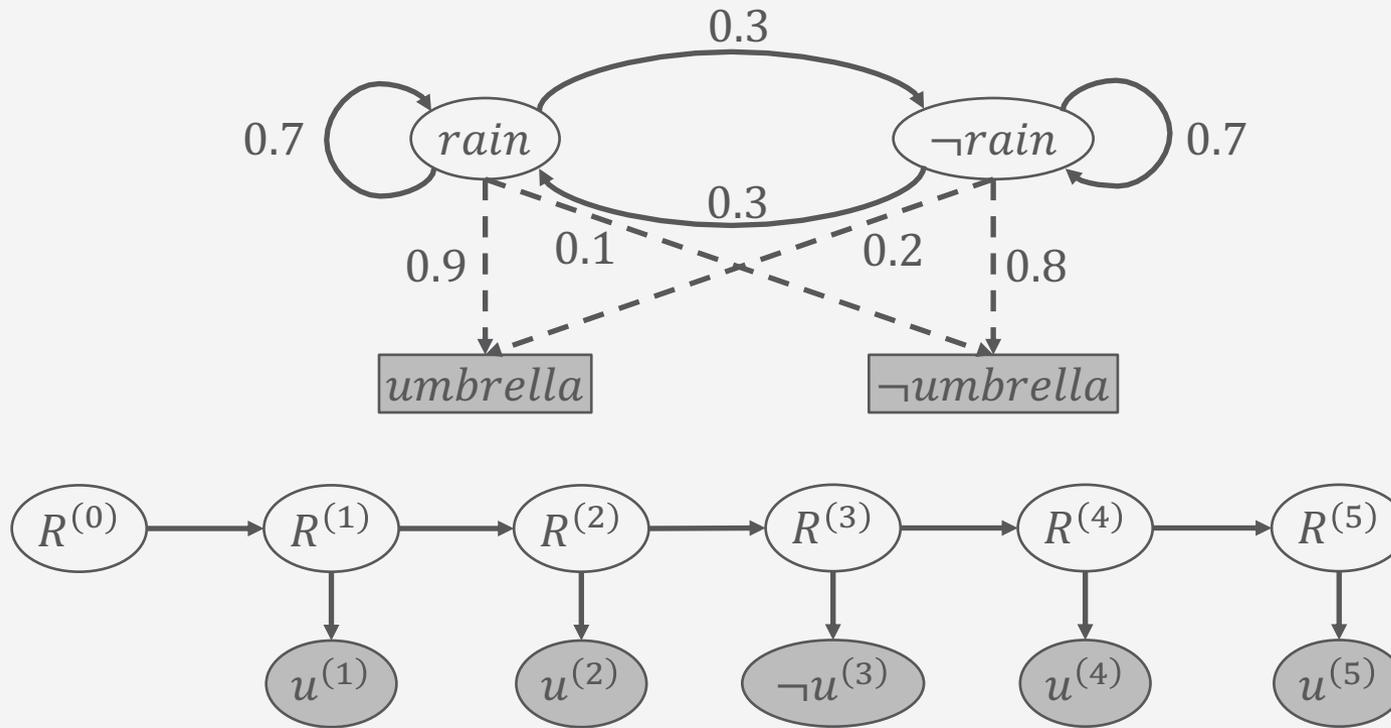
- Lernen benötigt Glättung, um bessere Schätzungen der Zustände zu bekommen
- Lernen von allgemeinen temporalen Modellen noch in den Kinderschuhen
 - Benötigt einen iterativen Ansatz im Sinne von Expectation-Maximisation (EM), um latente Zustandsvariablen zu schätzen
- Bekannter Lernalgorithmus für HMMs: [Baum-Welch-Algorithmus](#)
 - Interpretation eines HMMs als Zustandsübergangssystem

HMMs als Zustandsübergangssystem

- Menge von Zuständen $\{s_1, \dots, s_n\}$
 - Prozess bewegt sich von einem Zustand in den nächsten, was eine Zustandssequenz generiert: s_{i1}, \dots, s_{ik}
 - Markov Annahme: Nächster Zustand hängt nur vom jetzigen Zustand ab:
$$P(s_{ik} | s_{i1}, \dots, s_{ik-1}) = P(s_{ik} | s_{ik-1})$$
 - Matrix A : Übergangswahrscheinlichkeiten
- Zustände nicht beobachtbar, aber jeder Zustand generiert eine von m Beobachtungen (bzw. sichtbaren Zuständen) $\{v_1, \dots, v_m\}$
 - Matrix B : Emissionswahrscheinlichkeiten
- Für HMM:
 - $\text{Val}(V) = \{s_1, \dots, s_n\}$
 - $\text{Val}(E) = \{v_1, \dots, v_m\}$
 - Übergänge:
$$A = (a_{ij}), a_{ij} = P(s_j | s_i)$$
 - Wie in der HMM Matrixdarstellung
 - Beobachtungen:
$$B = (b_i(v)), b_i(v) = P(v | s_i)$$
 - Wie in der HMM Matrixdarstellung
 - Vektor mit initialen Wahrscheinlichkeiten:
$$\pi = (\pi_i), \pi_i = P(s_i)$$
 - HMM dargestellt durch $M = (A, B, \pi)$

Zustandsübergangssystem: Beispiel

- Latente Zufallsvariable R
- Evidenzvariable U



$P(r^{(0)})$
0.5

$R^{(t-1)}$	$P(r^{(t)} R^{(t-1)})$
true	0.7
false	0.3

$R^{(t)}$	$P(u^{(t)} R^{(t)})$
true	0.9
false	0.2

Lernproblem

- Gegeben eine Menge von Beobachtungssequenzen $O = o_1, \dots, o_k$ als Trainingsdaten und die generelle Struktur des HMMs (Anzahl an versteckten und sichtbaren Zuständen)
- Bestimme HMM Parameter $M = (A, B, \pi)$, welche am besten zu den Trainingsdaten passen, i.e., welche $P(O|M)$ maximieren
- Wenn wir Daten zu S hätten, könnten wir direkt ML-basiert die Einträge in M bestimmen

$$a_{ij} = P(s_j | s_i) = \frac{\text{Anzahl an Übergängen von } s_i \text{ nach } s_j}{\text{Anzahl an Übergängen aus } s_i \text{ raus}}$$

$$b_i(v) = P(v | s_i) = \frac{\text{Anzahl an Beobachtungen von } v \text{ in } s_i}{\text{Anzahl an Besuchen von } s_i}$$

- Andernfalls: Iterativ mittels EM Einträge als lokales Optimum bestimmen
→ **Baum-Welch-Algorithmus**

Baum-Welch-Algorithmus

- Generelle Idee: Erwartete Zähler nutzen (EM)

$$a_{ij} = P(s_j | s_i) = \frac{\text{Erwartete Anzahl an Übergängen von } s_i \text{ nach } s_j}{\text{Erwartete Anzahl an Übergängen aus } s_i \text{ raus}}$$

$$b_i(v) = P(v | s_i) = \frac{\text{Erwartete Anzahl an Beobachtungen von } v \text{ in } s_i}{\text{Erwartete Anzahl an Besuchen von } s_i}$$

$$\pi_i = P(s_i) = \text{erwartete relative Häufigkeit von } s_i \text{ zum Zeitpunkt } k = 1$$

E-Schritt

- Hilfsvariable $\xi_k(i, j)$: Wahrscheinlichkeit in Zustand s_i zum Zeitpunkt k und in Zustand s_j zum Zeitpunkt $k + 1$ zu sein, gegeben die Beobachtungssequenzen $o_1, \dots, o_T, k < T$

$$\begin{aligned}
 \xi_k(i, j) &= P(Q_k = s_i, Q_{k+1} = s_j | o_1, \dots, o_T) \\
 \xi_k(i, j) &= \frac{P(Q_k = s_i, Q_{k+1} = s_j, o_1, \dots, o_T)}{P(o_1, \dots, o_T)} \\
 &= \frac{P(Q_k = s_i, o_1, \dots, o_k) a_{ij} b_i(o_{k+1}) P(o_{k+2}, \dots, o_T | Q_{k+1} = s_j)}{P(o_1, \dots, o_T)} \\
 &= \frac{1}{Z} \text{FORWARD}_k(i) a_{ij} b_i(o_{k+1}) \text{BACKWARD}_{k+1}(j)
 \end{aligned}$$

E-Schritt

- Hilfsvariable $\gamma_k(i)$: Wahrscheinlichkeit in Zustand s_i zum Zeitpunkt k zu sein, gegeben die Beobachtungssequenzen $o_1, \dots, o_T, k < T$

$$\begin{aligned}\gamma_k(i) &= P(Q_k = s_i | o_1, \dots, o_T) \\ \gamma_k(i) &= \frac{P(Q_k = s_i, o_1, \dots, o_T)}{P(o_1, \dots, o_T)} \\ &= \frac{1}{Z} \text{FORWARD}_k(i) \text{BACKWARD}_{k+1}(j)\end{aligned}$$

E-Schritt

- Hilfsvariablen berechnen mit anfangs geratenen Parametern
 - $\xi_k(i, j) = P(Q_k = s_i, Q_{k+1} = s_j | o_1, \dots, o_T)$
 - $\gamma_k(i) = P(Q_k = s_i | o_1, \dots, o_T)$
- Parameter schätzen
 - Erwartete Anzahl an Übergängen von Zustand s_i nach Zustand $s_j = \sum_k \xi_k(i, j)$
 - Erwartete Anzahl an Übergängen aus Zustand $s_i = \sum_k \gamma_k(i)$
 - Erwartete Anzahl an Beobachtungen v in Zustand $s_i = \sum_{k, o_k=v} \gamma_k(i)$
 - Erwartete relative Häufigkeit von Zustand s_i zum Zeitpunkt ($k = 1$) = $\gamma_1(i)$

M-Schritt

- Parameter maximieren

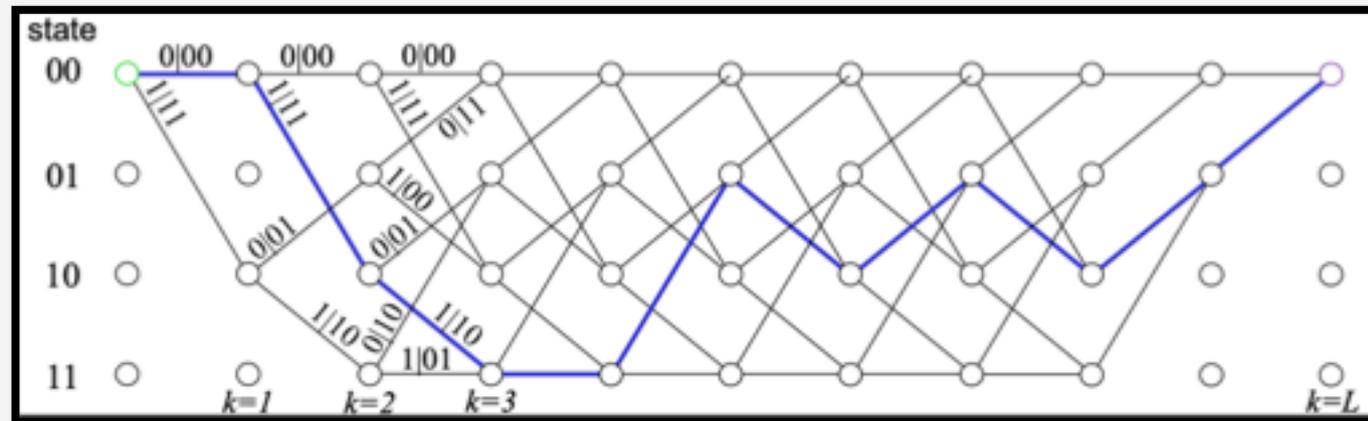
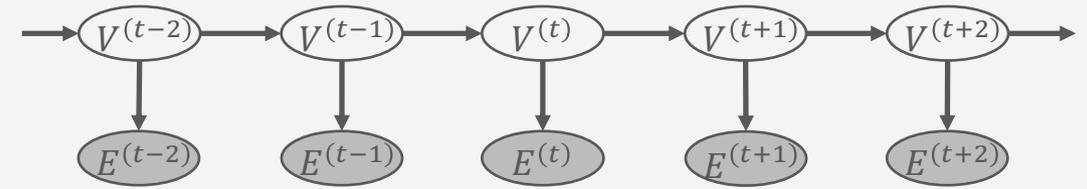
$$a_{ij} = P(s_j | s_i) = \frac{\text{Erwartete Anzahl an Übergängen von } s_i \text{ nach } s_j}{\text{Erwartete Anzahl an Übergängen aus } s_i \text{ raus}} = \frac{\sum_k \xi_k(i, j)}{\sum_k \gamma_k(i)}$$

$$b_i(v) = P(v | s_i) = \frac{\text{Erwartete Anzahl an Beobachtungen von } v \text{ in } s_i}{\text{Erwartete Anzahl an Besuchen von } s_i} = \frac{\sum_{k, o_k=v} \gamma_k(i)}{\sum_k \gamma_k(i)}$$

$$\pi_i = P(s_i) = \text{erwartete relative Häufigkeit von } s_i \text{ zum Zeitpunkt } (k = 1) = \gamma_1(i)$$

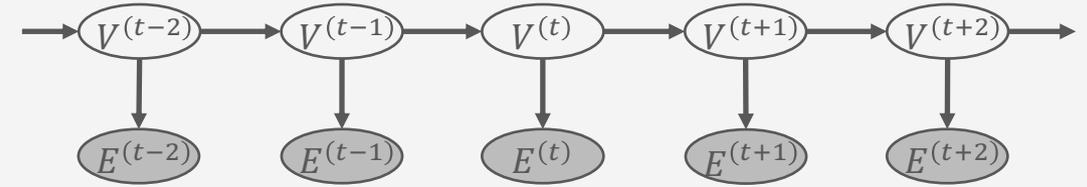
Anwendungen

- Gensequenzanalyse
 - Beispiel: MPE-Anfrage mit RNA-Strang als Evidenz E , DNA-Sequenz als latenter Zustand V
 - $\text{Val}(V) = \{A, C, G, T\}$, $\text{Val}(E) = \{A, C, G, U\}$
 - Änderungen durch: Baustein fällt weg, kommt hinzu, wird falsch übersetzt
- Signalverarbeitung
 - Rauschen herausfiltern: Was war der ursprünglich gesendete Wert $\{0,1\}$? → MPE-Anfrage



Anwendungen

- Textverarbeitung
 - Gedichte in Tamil, angereichert mit Kommentaren im Laufe der Jahrhunderte
 - MPE-Anfrage: Was im Text ist Kommentar, was ist Gedicht?



Tamil/Sanskrit text

Grammar of Old Tamil for Students

Critical Edition

Edition (1933) + Commentary

Edition (1940) + Commentary

Edition (1975) + Commentary

Edition (1990) + Commentary

Loss of original content of the manuscript over time.

Old manuscripts are of high interest these days.

Eva Wilden, *A Critical Edition and an Annotated Translation of the Akanāpūru (Part 1 - Kalīriyānirai)*, 2018

Eva Wilden, *Grammar of Old Tamil for Students*, 2018



Not visible on palm-leaf

Second row in original poem

Second row in original poem

Bold: part of original poem

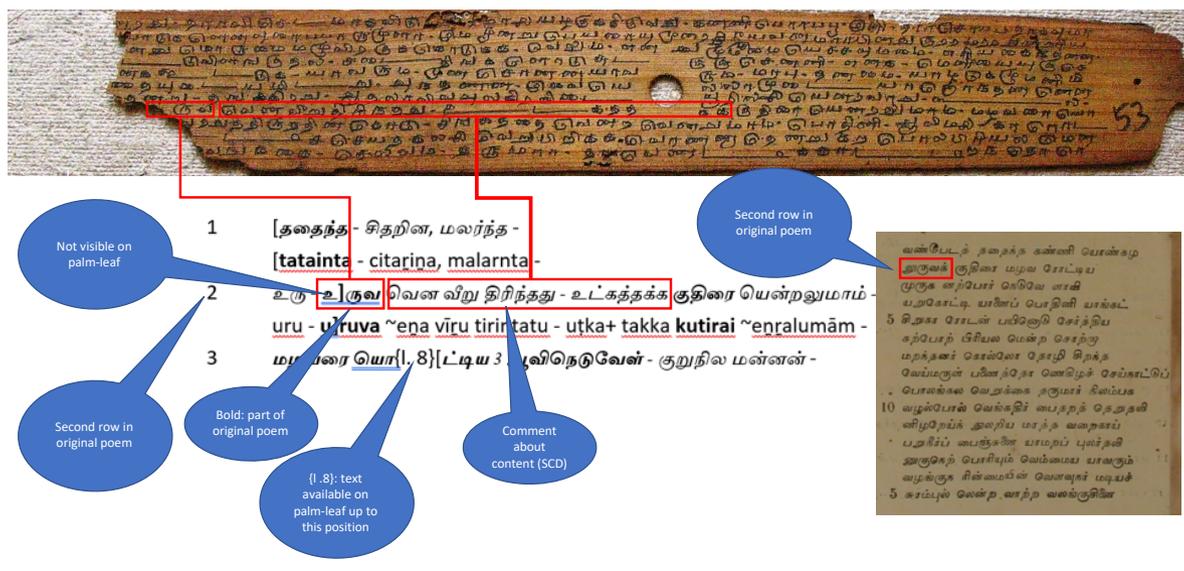
[l.8]: text available on palm-leaf up to this position

Comment about content (SCD)

1 [ததைந்த - சிதறின, மலர்ந்த -
[tatainta - citarina, malarnta -

2 உரு **உ]குவு** வென வீறு திரிந்தது - உட்கத்தக்க குதிரை யென்றலுமாம் -
uru - uṟuva ̣ena vīru tirintatu - utka+ takka kutirai ̣enralumām -

3 மயலரை யொ[.8][டடிய 3] விநெடுவேள் - குறுநில மன்னன் -



Zwischenzusammenfassung

- Direkte Interpretation der Wahrscheinlichkeiten in HMMs möglich
 - Vorwärts- und Rückwärtsnachrichten
 - Viterbi-Algorithmus für MPEs
- Berechnungen über Matrixoperationen
 - Erlauben effiziente Umsetzung
 - Country-Dance-Algorithmus um Vorwärtsnachrichten beim Rückwärtsgehen zu rekonstruieren
- Baum-Welch-Algorithmus
 - Interpretation eines HMMs als Zustandsübergangssystem
 - Lernen der Parameter eines HMMs ML- und EM-basiert

Überblick: 6. Sequentielle PGMs und Inferenz

A. *Sequentielle PGMs*

- Templates, dynamische BNs, dynamische Faktormodelle, Hidden-Markov-Modelle; Semantik
- Inferenzaufgaben: Wahrscheinlichkeitsanfragen (Filtering, Prediction, Hindsight), Zustandsanfragen (MPE, MAP)

B. *Sequentielle Inferenz*

- Naïve Inferenz mittels Ausrollen, Interface Algorithmus, Komplexität, Approximationen

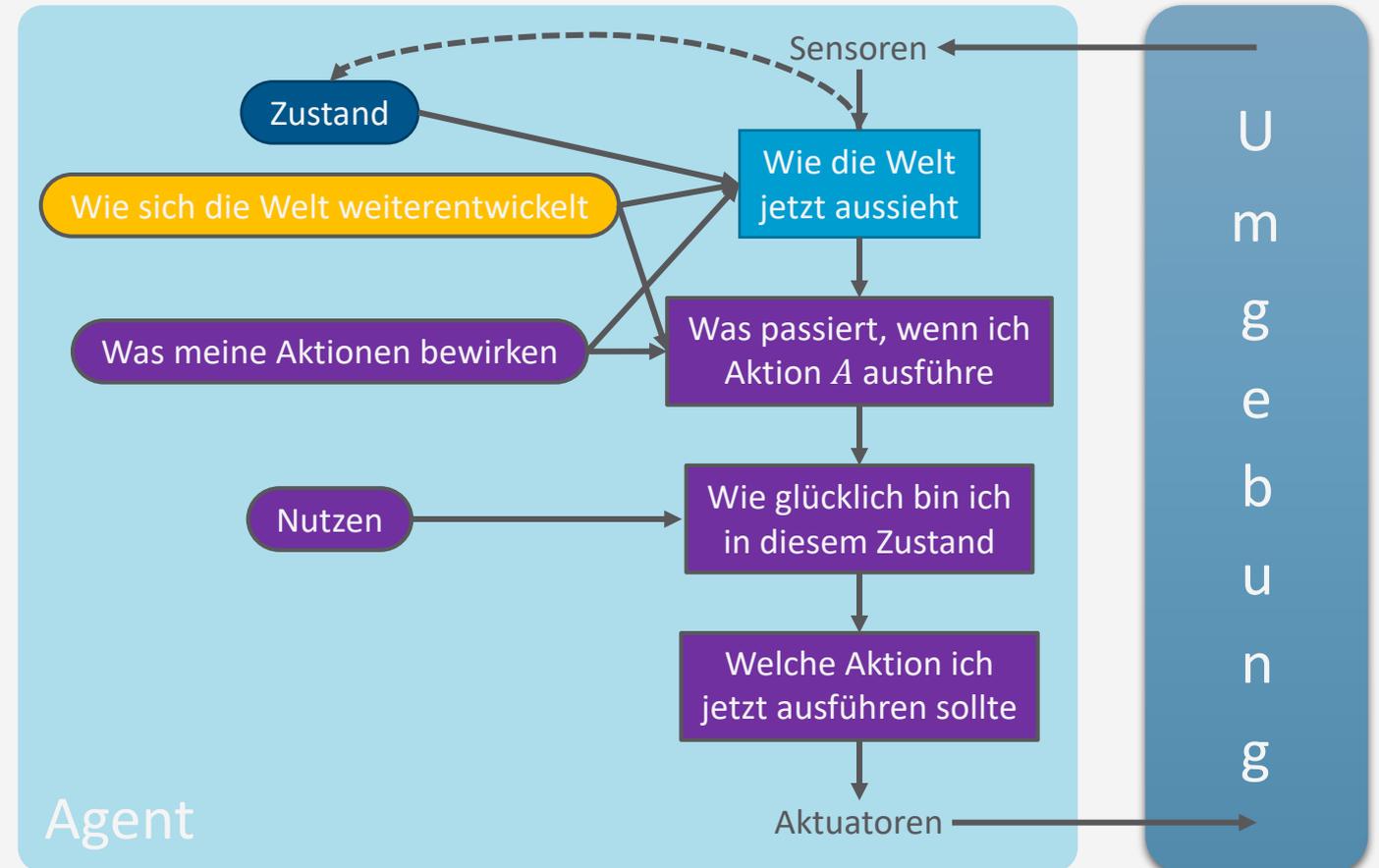
C. *Spezialfall Hidden-Markov-Modelle*

- Viterbi-Algorithmus für MPEs
- Anfragebeantwortung durch Matrixoperationen
- Baum-Welch-Algorithmus zum Lernen

→ Entscheidungstheoretische PGMs und Inferenz

Einordnung der Vorlesung: *Modell- und nutzenbasierter Agent*

- Nachfolgende Themen der Vorlesung
 2. Episodische PGMs
 3. Exakte Inferenz in episodischen PGMs
 4. Approximative Inferenz in episodischen PGMs
 5. Lernalgorithmen für episodische PGMs
 6. **Sequentielle PGMs und Inferenz**
 7. Entscheidungstheoretische PGMs



„Anhang“

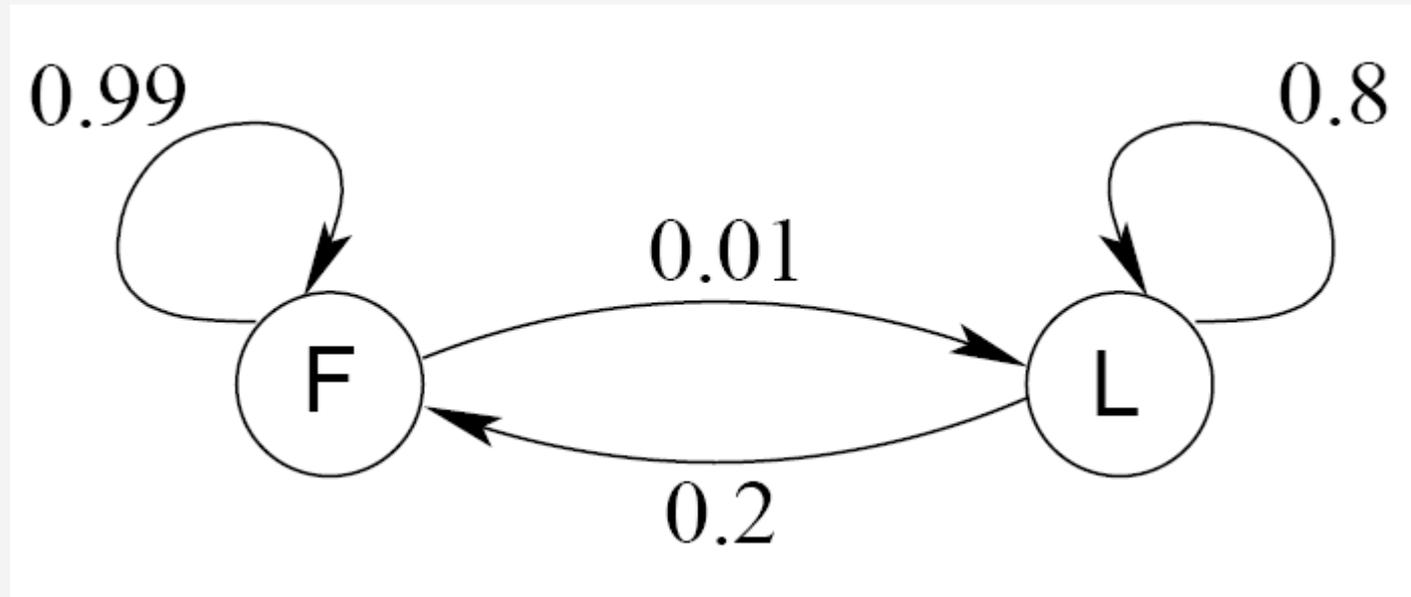
Nachfolgende Folien beinhalten ein nicht genutztes Beispiel zum Viterbi-Algorithmus

ACHTUNG: Folien sind nicht ordentlich gesetzt oder übersetzt!

The occasionally dishonest casino

- A casino uses a fair die most of the time, but occasionally switches to a loaded one
 - Fair die: $\text{Prob}(1) = \text{Prob}(2) = \dots = \text{Prob}(6) = 1/6$
 - Loaded die: $\text{Prob}(1) = \text{Prob}(2) = \dots = \text{Prob}(5) = 1/10$, $\text{Prob}(6) = 1/2$
 - These are the emission probabilities
- Transition probabilities
 - $\text{Prob}(\text{Fair} \rightarrow \text{Loaded}) = 0.01$
 - $\text{Prob}(\text{Loaded} \rightarrow \text{Fair}) = 0.2$
 - Transitions between states modeled by a Markov process

Transition model for the casino



The occasionally dishonest casino

- Known:
 - The structure of the model
 - The transition probabilities
- Hidden: What the casino did
 - FFFFFLLLLLLLLFFFF...
- Observable: The series of die tosses
 - 3415256664666153...
- What we must infer:
 - When was a fair die used?
 - When was a loaded one used?
 - The answer is a sequence
FFFFFFFFLLLLLLLLFFF...

Making the inference

- Model assigns a probability to each explanation of the observation:

$$\begin{aligned} & P(326 | FFL) \\ &= P(3 | F) \cdot P(F \rightarrow F) \cdot P(2 | F) \cdot P(F \rightarrow L) \cdot P(6 | L) \\ &= 1/6 \cdot 0.99 \cdot 1/6 \cdot 0.01 \cdot 1/2 \end{aligned}$$

Determine which explanation is most likely

- Find the path *most likely* to have produced the observed sequence

Determine probability that observed sequence was produced by the model

- Consider *all* paths that could have produced the observed sequence

Notation

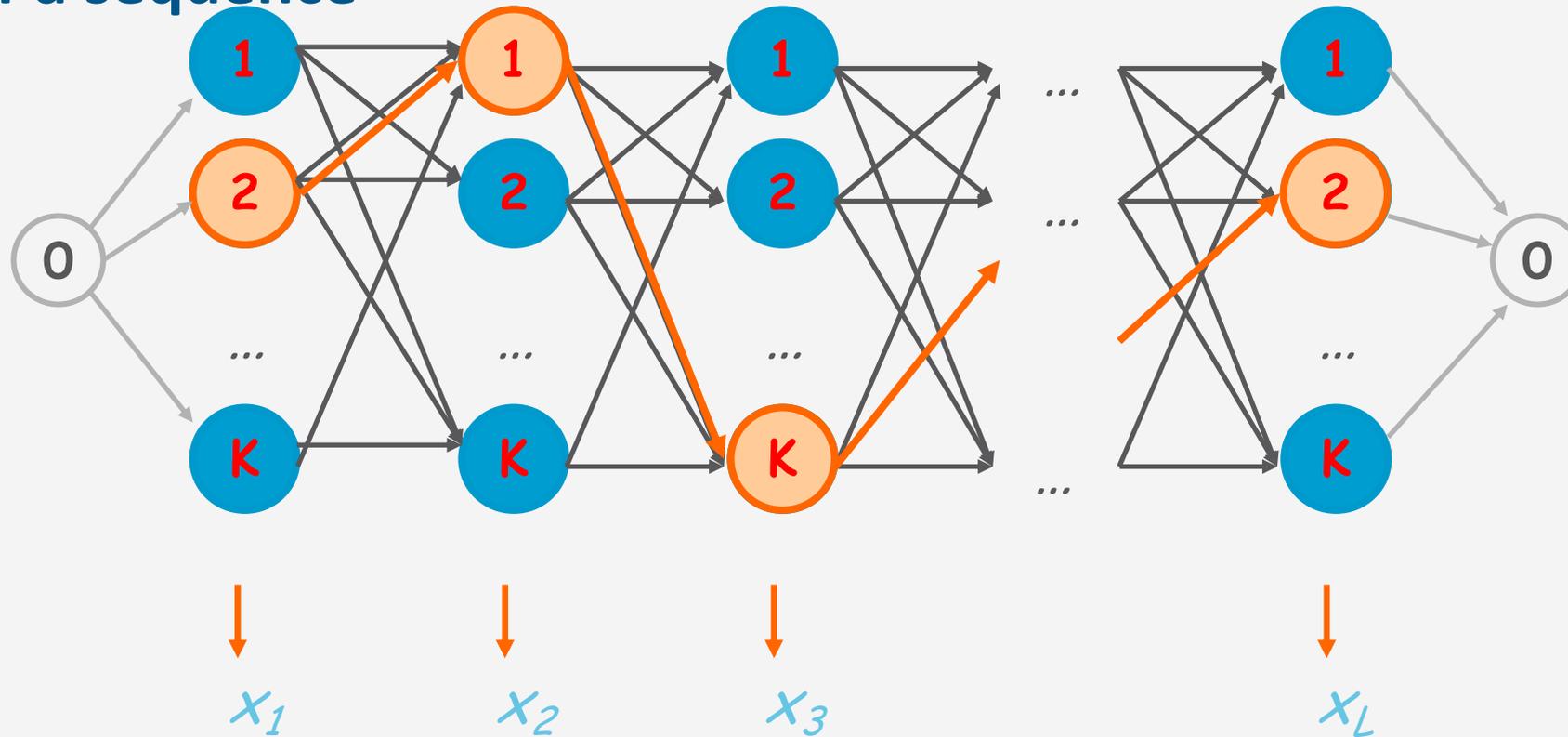
- x is the sequence of symbols emitted by model
 - x_i is the symbol emitted at time i
- A **path**, π , is a sequence of states
 - The i -th state in π is π_i
- a_{kr} is the probability of making a transition from state k to state r :

$$a_{kr} = \Pr(\pi_i = r \mid \pi_{i-1} = k)$$

- $e_k(b)$ is the probability that symbol b is emitted when in state k

$$e_k(b) = \Pr(x_i = b \mid \pi_i = k)$$

A “parse” of a sequence



$$\Pr(x, \pi) = a_{0\pi_1} \prod_{i=1}^L e_{\pi_i}(x_i) \cdot a_{\pi_i \pi_{i+1}}$$

The occasionally dishonest casino

$$x = \langle x_1, x_2, x_3 \rangle = \langle 6, 2, 6 \rangle$$

$$\pi^{(1)} = FFF$$

$$\pi^{(2)} = LLL$$

$$\begin{aligned} \Pr(x, \pi^{(2)}) &= \Pr(x, \pi^{(1)}) \Pr(x, \pi^{(2)} | x, \pi^{(1)}) \\ &= a_{FF} e_L(6) a_{FF} e_L(2) a_{FF} e_L(6) \\ &= 0.5 \times \frac{1}{6} \times 0.99 \times \frac{1}{6} \times 0.99 \times \frac{1}{6} \times 0.5 \\ &\approx 0.00227 \end{aligned}$$

$$\pi^{(3)} = LFL$$

$$\begin{aligned} \Pr(x, \pi^{(3)}) &= a_{0L} e_L(6) a_{LF} e_F(2) a_{FL} e_L(6) a_{L0} \\ &= 0.5 \times 0.5 \times 0.2 \times \frac{1}{6} \times 0.01 \times 0.5 \\ &\approx 0.0000417 \end{aligned}$$

The most probable path

The most likely path π^* satisfies

$$\pi^* = \arg \max_{\pi} \Pr(x, \pi)$$

To find π^* , consider all possible ways the last symbol of x could have been emitted

Let

$v_k(i)$ = Prob. of path $\langle \pi_1, \dots, \pi_i \rangle$ most likely
to emit $\langle x_1, \dots, x_i \rangle$ such that $\pi_i = k$

Then

$$v_k(i) = e_k(x_i) \max_r (v_r(i-1) a_{rk})$$

The Viterbi Algorithm

- Initialization ($i = 0$)

$$v_0(0) = 1, \quad v_k(0) = 0 \text{ for } k > 0$$

- Recursion ($i = 1, \dots, L$): For each state k

$$v_k(i) = e_k(x_i) \max_r (v_r(i-1) a_{rk})$$

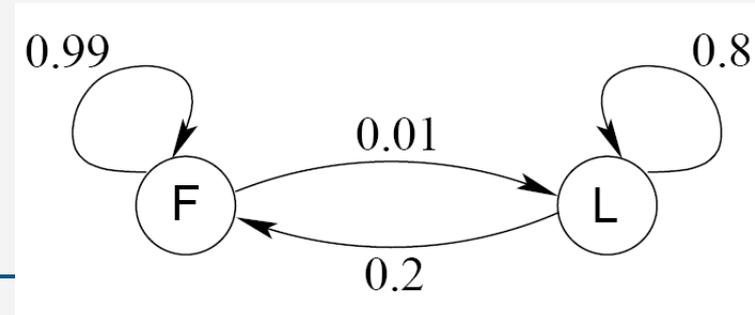
- Termination:

$$\Pr(x, \pi^*) = \max_k (v_k(\text{Length}) a_{k0})$$

To find π^ , use trace-back, as in dynamic programming*

Viterbi: Example

		ε	6	2	x	6
π	B	1	0	0		0
	F	0	$(1/6) \times (1/2) = 1/12$	$(1/6) \times \max\{(1/12) \times 0.99, (1/4) \times 0.2\} = 0.01375$		$(1/6) \times \max\{0.01375 \times 0.99, 0.02 \times 0.2\} = 0.00226875$
	L	0	$(1/2) \times (1/2) = 1/4$	$(1/10) \times \max\{(1/12) \times 0.01, (1/4) \times 0.8\} = 0.02$		$(1/2) \times \max\{0.01375 \times 0.01, 0.02 \times 0.8\} = 0.08$



Viterbi gets it right more often than not

Rolls	315116246446644245321131631164152133625144543631656626566666
Die	FFL
Viterbi	FFL
Rolls	651166453132651245636664631636663162326455235266666625151631
Die	LLLLLFFFLFF
Viterbi	LLLLLFFFLFF
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFLFF
Viterbi	FFL
Rolls	366163666466232534413661661163252562462255265252266435353336
Die	LLLLLLLLLFF
Viterbi	LLLLLLLLLFF
Rolls	233121625364414432335163243633665562466662632666612355245242
Die	FFLFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	FFLFFFFFFFFFFFFFFFFFFFFFFFF