

DAI-Labor
TU Berlin

Grundlagen der Künstlichen Intelligenz

08.12.2005: Modallogik

According to the intentional stance, an agent is assumed to decide to act and communicate based on its beliefs about its environment and its desires and intentions.
nach D.C. Dennet: The Intentional Stance

Dr.-Ing. Stefan Fricke
stefan.fricke@dai-labor.de

AIOIT
Agententechnologien in betrieblichen Anwendungen und der Telekommunikation

Gliederung

- ⇒ Einführung in Modallogik
- ⇒ Modallogiken für Zeit und Handeln
- ⇒ Intentionalität
- ⇒ Epistemische Logik mit Belief, Desire, Intention
- ⇒ BDI-Agenten
- ⇒ Zusammenfassung

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 2

Probleme klassischer Logik Einführung

Klassische Logik beschreibt Tatsachen und Zusammenhänge, aber..

- ⇒ ... kein Handeln, keine Prozesse
 - Der Effekt einer Aktion ist unvorhersagbar in nichtdeterministischen und in dynamischen Umgebungen
- ⇒ ... keine Zeit
- ⇒ ... keine Überzeugungen, Motivationen, etc
 - Ich hoffe, du glaubst, dass ich weiß, dass der Zahnarztbesuch Schmerzen bereitet. Ich beabsichtige nicht, Schmerzen zugefügt zu bekommen. (Aber ich beabsichtige zum Zahnarzt zu gehen)

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 3

Modallogik Einführung

- ⇒ Idee: Aus logischen Zusammenhängen und Inferenzen der Prädikatenlogik **sinnvolle Sachverhalte und Folgerungen** durch Einführung so genannter **Modaloperatoren** ermöglichen.
- ⇒ Modaloperatoren werden vor eine Formel geschrieben und geben ihr eine neue Interpretation. Z. B. für $p = \text{es_regnet}$

ICH_GLAUBE p	vs.	ICH_HOFFE $\neg p$
IRGENDWANN p	und	ICH_WEISS IRGENDWANN p
ICH_GLAUBE DU_HOFFST p	vs.	ICH_HOFFE DU_GLAUBST p

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 4

Possible Worlds Einführung

⇒ Verschiedene Aktionen führen in unterschiedliche Welten.

⇒ Der Effekt einer Aktion ist nicht notwendigerweise deterministisch

t=1 t=2 t=3 ...

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 5

Modallogik Einführung

- ⇒ **Erreichbarkeitsrelation** $R(W_1, W_2)$: W_2 ist von W_1 aus erreichbar.
- ⇒ **Möglichkeit**: $\diamond P$ ist folgerbar in einer Welt W genau dann, wenn P wahr ist in *mindestens* einer möglichen Welt
 - $\exists w': R(w, w') \wedge w' \models P$
- ⇒ **Notwendigkeit**: $\square P$ ist folgerbar in einer Welt W genau dann, wenn P wahr ist in *jeder* möglichen Welt
 - $\forall w': R(w, w') \Rightarrow w' \models P$

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 6

Axiomschemata für Modallogiken Einführung

K: $\Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$ (Distribution)
D: $\Box A \rightarrow \neg \Box \neg A$ (Konsistenz)
T: $\Box A \rightarrow A$ (Reflexivität)
B: $A \rightarrow \Box \Diamond A$ (Symmetrie)
4: $\Box A \rightarrow \Box \Box A$ (Transitivität)
NEC: $A \rightarrow \Box A$ (Notwendigkeit)

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 7

Possible Worlds Semantik Einführung

⇒ Ein Modell ist nicht nur eine, sondern eine Menge von Welten ...
 → ... generiert durch Aktionsauswahl,
 → ... generiert durch eine Sequenz von Aktionen,
 → ... generiert durch mehrere mögliche Ergebnisse einer Aktion,
 → ... generiert durch Wissenslücken.

⇒ Possible Worlds dienen zur Beschreibung dieses Modells...

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 8

Beispiel für Possible Worlds in der Blockswelt Einführung

⇒ $W_0: \{ \text{ontable}(C), \text{on}(B,C), \text{on}(A,B), \text{clear}(A) \}$

⇒ $\Diamond \text{clear}(B)$
 ⇒ $\Box \text{ontable}(C)$

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 9

Gliederung

- ⇒ Einführung in Modallogik
- ⇒ Modallogiken für Zeit und Handeln
 - Temporallogik
 - Dynamische Logik
- ⇒ Intentionalität
- ⇒ Epistemische Logik mit Belief, Desire, Intention
- ⇒ BDI-Agenten
- ⇒ Zusammenfassung

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 10

Modallogiken im Zusammenhang Modale Logiken

Modallogik
 Modaloperatoren: \Diamond, \Box

Wissen u. Glauben
 M-O: B, K

Temporallogik
 M-O: für Zeit (linear / branching)

Dynamische Logik
 M-O: für Aktionen

BDI Logik
 M-O: B, D, I

CTL Logik
 M-O: für Aktionen und branching time

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 11

Logik von Glauben und Wissen Modale Logiken

⇒ $K_a \phi$ Agent a weiß ϕ
 ⇒ $Bel_a \phi$ Agent a glaubt ϕ

⇒ Ein Agent weiß / glaubt eine Proposition ϕ genau dann, wenn ϕ in allen möglichen Welten gilt.

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 12

Eigenschaften von Wissen und von Glauben	Modale Logiken
(1) $K_a \varphi \wedge K_a (\varphi \rightarrow \psi) \rightarrow K_a \psi$	
(2) $K_a \varphi \rightarrow \varphi$	Wissen
(3) $K_a \varphi \rightarrow K_a (K_a \varphi)$	positive Introspektion
(4) $\neg K_a \varphi \rightarrow K_a (\neg K_a \varphi)$	negative Introspektion
\Rightarrow Für Bel _a φ analog: <ul style="list-style-type: none"> \rightarrow gilt (1) \rightarrow gilt nicht (2) \rightarrow gilt (3) \rightarrow kann (4) gelten, ist aber problematisch. 	

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 13

Beispiel	Modale Logiken
\Rightarrow 2 Personen sitzen einander gegenüber und wissen (1), dass mindestens einer von ihnen eine weiße Stirn hat (2). B sagt, er wisse nicht, ob er eine weiße Stirn habe (3). Daraus folgert A, dass er selbst eine weiße Stirn hat.	
1. $K_A(\neg \text{White}(A) \rightarrow K_B(\neg \text{White}(A)))$	(1)
2. $K_A(K_B(\text{White}(A) \vee \text{White}(B)))$	(2)
3. $K_A(\neg K_B(\text{White}(B)))$	(3)
\Rightarrow Kann A daraus schließen, dass er selbst eine weiße Stirn hat?	

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 14

Beispiel	Modale Logiken
Inferenzregeln A1: $K_A \varphi \wedge K_A (\varphi \rightarrow \psi) \rightarrow K_A \psi$; A2: $K_A \varphi \rightarrow \varphi$ R2: $\varphi \rightarrow \psi \wedge K_A \varphi \rightarrow K_A \psi$	
1. $K_A(\neg \text{White}(A) \rightarrow K_B(\neg \text{White}(A)))$	
2. $K_A(K_B(\text{White}(A) \vee \text{White}(B)))$	
3. $K_A(\neg K_B(\text{White}(B)))$	
4. $\neg \text{White}(A) \rightarrow K_B(\neg \text{White}(A))$	1, A2; $X \rightarrow Y \Leftrightarrow \neg X \vee Y$
5. $K_B(\neg \text{White}(A) \rightarrow \text{White}(B))$	2, A2
6. $K_B(\neg \text{White}(A)) \rightarrow K_B(\text{White}(B))$	5, A1
7. $\neg \text{White}(A) \rightarrow K_B(\text{White}(B))$	4, 6
8. $\neg K_B(\text{White}(B)) \rightarrow \text{White}(A)$	Umkehrschluss von 7
9. $K_A(\text{White}(A))$	3, 8, R2

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 15

Dynamische Logik	Modale Logiken
\Rightarrow Aktionssymbole a, b, ...	
\Rightarrow a;b	Sequenz
\Rightarrow a+b	nichtdeterministische Auswahl
\Rightarrow p?	deterministische Auswahl (IF-THEN) Aktion, basierend auf Wahrheitswert von p
\Rightarrow a*	0 oder mehr Wiederholungen von a
\Rightarrow <a>p	a macht p möglicherweise wahr (analog zu \diamond)
\Rightarrow [a]p	a macht p notwendigerweise wahr (analog zu \square)

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 16

Temporallogik	Modale Logiken
\Rightarrow Momente T mit einer partiellen Ordnung <, die die Vorher-Beziehung ausdrückt. \rightarrow Jeder Moment entspricht einer Possible World	
\Rightarrow p U q	p ist wahr bis q wahr wird (UNTIL)
\Rightarrow Xp	p ist im nächsten Moment wahr (NEXT)
\Rightarrow Pp	p war in einem früheren Moment wahr (PAST)
\Rightarrow Ep	p ist irgendwann in der Zukunft wahr (EVENTUALLY)
\Rightarrow Ap	p wird immer wahr sein (ALWAYS) Ap = $\neg \neg E \neg p$

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 17

Gliederung
\Rightarrow Einführung in Modallogik
\Rightarrow Modallogiken für Zeit und Handeln
\Rightarrow Intentionalität
\Rightarrow Epistemische Logik mit Belief, Desire, Intention
\Rightarrow BDI-Agenten
\Rightarrow Zusammenfassung

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 18

Fragestellungen zum Verhalten komplexer Systeme

- ⇒ Wie lässt sich das Verhalten eines Systems vorhersagen ohne Kenntnis seiner internen Struktur?
- ⇒ Welche Repräsentationen eignen sich zur Definition
 - des beobachtbaren Verhaltens?
 - des erwarteten Verhaltens?
- ⇒ Die **Intentionalität** gibt Antworten auf die Beschreibung des Verhaltens komplexer Systeme...

Zur Theorie des Intentional Stance

- ⇒ Der Philosoph **Daniel Dennett** unterscheidet zwischen drei Herangehensweisen, wenn man **Vorhersagen über ein System** treffen will:
 - physical stance,
 - design stance,
 - intentional stance...

Physical Stance

- ⇒ Der **physical stance** beschreibt Entitäten in physikalischen Begriffen.
- ⇒ Zum Beispiel geometrische Körper in der Mathematik, Mechanik im Maschinenbau, Statik in der Architektur.

Design Stance

- ⇒ Der **design stance** nimmt an, dass ein Objekt zu einem Zweck entwickelt wurde und sich entsprechend verhält.
- ⇒ Die Erklärung erfolgt anhand funktionaler Begriffe.
- ⇒ Beispiel: Das Verhalten eines Autos wird über Funktionen wie Gas geben, Kupplung treten, Lenken, Bremsen beschrieben.

Intentional Stance

- ⇒ Der **intentional stance** beschreibt das Verhalten einer Entität als das eines rationalen Agenten.
 - Es ist immer auch möglich, dasselbe Verhalten in rein physikalischen oder funktionalen Begriffen zu erklären.
- ⇒ Aus pragmatischen Gründen ist die intentionale Perspektive jedoch unverzichtbar. Was können das für pragmatische Gründe sein?
- ⇒ Die Komplexität der Entität.
 - U.U. ist es sinnvoll, einem Thermostaten die „Absicht“ zuzubilligen, die Temperatur konstant zu halten.

Theorie des Intentional Stance

- ⇒ Die Kernidee ist, das Verhalten von Individuen durch die **Zuschreibung von Wünschen und Überzeugungen** rational verständlich und damit prognostizierbar zu machen.
- ⇒ Diese Zuschreibungen beschreiben jedoch keine Tatsachen, aus denen das Verhalten aufgrund kausaler Gesetzmäßigkeiten erschlossen werden kann.
 - Vielmehr stellen intentionale Beschreibungen ein Abstraktionswerkzeug dar, das mit familiären Begriffen operiert.

Schachcomputer als Beispiel für den Intentional Stance

- ⇒ Wie setzt man sich adäquat mit einem Schachcomputer auseinander, dessen exakte innere Implementation man nicht kennt?
 - Er hat Wissen über die Figuren auf den Feldern und die Bewegungsmöglichkeiten der Figuren;
 - Er hat den Wunsch, das Spiel zu gewinnen.
- ⇒ Mit diesen Zuschreibungen kann man gute Vorhersagen oder Erklärungsansätze über sein Verhalten liefern.

Gliederung

- ⇒ Einführung in Modallogik
- ⇒ Modallogiken für Zeit und Handeln
- ⇒ Intentionalität
- ⇒ Epistemische Logik mit Belief, Desire, Intention
 - Modaloperatoren Belief, Desire, Intention
 - Axiomatisierungen für B, D, I
- ⇒ BDI-Agenten
- ⇒ Zusammenfassung

Belief, Desire, Intention (BDI) bietet dreierlei:

- ⇒ ein philosophisches Modell des menschlichen Denkens und Handelns,
 - [Bratman, 1987]
- ⇒ verschiedene Architekturimplementierungen,
 - (IRMA, PRS, JACK, JIAC)
- ⇒ eine abstrakte logische Semantik.

Epistemische Modallogiken für Agenten

Epistemische Logik nutzen, um

- ⇒ mentale Attitüden von Agenten zu modellieren:
 - Beliefs, Desires, Goals, Know How, Intentions, ...
- ⇒ das Verhalten eines Agenten zu begründen,
- ⇒ das Verhalten eines Agenten vorhersehbar zu machen,
- ⇒ sinnvolle Handlungen zu generieren.

Belief, Desire, Intention

- ⇒ $Bel_x p$: Agent x glaubt p.
- ⇒ $Des_x p$: Agent x wünscht Zustand p
 - Ziele ermöglichen in die Zukunft gerichtetes Handeln.
- ⇒ $Int_x p$ Agent x beabsichtigt, den Zustand p zu erreichen
 - Intentionen sind mit Handlungen verknüpft

Eigenschaften von Zielen (Goals, bzw. Desires)

- ⇒ Desires können nicht erfüllbar oder gar inkonsistent sein.
 - **Goals** sind Teilmengen von Desires: erreichbar und konsistent.
- ⇒ Der Agent soll glauben, dass sein Ziel erreichbar ist.
 - Das verhindert, dass der Agent Ziele annimmt, von denen er glaubt, dass sie unerreichbar sind.
 - Diese Eigenschaft wird **Realismus** genannt.
- ⇒ Ziele (und auch Intentionen) können überprüft werden.

Eigenschaften von Intentionen

- ⇒ Jede Intention soll erfüllbar sein.
- ⇒ Die Intentionen sollen gegenseitig konsistent sein.
- ⇒ Sie sollen mit den Beliefs des Agenten vereinbar sein.
 - $Int_x p$ impliziert nicht $Des_x p$
- ⇒ Ein Agent muss nicht unbedingt an die Erfüllbarkeit einer Intention glauben, ausschließen darf er sie jedoch nicht.

Possible Worlds in der Computation Tree Logic (CTL)

- ⇒ Eindeutige Vergangenheit, verzweigende Zukunft
- ⇒ Ein **Pfad** ist eine maximale Menge von Momenten, ausgehend von der Gegenwart bis in alle Momente in der Zukunft, entlang einem Zweig gemäß $<$
- ⇒ **Situation** = Welt zu einem Zeitpunkt
- ⇒ **Zustandsformeln** beziehen sich auf eine Situation
- ⇒ **Pfadformeln** beziehen sich auf einen Pfad ...

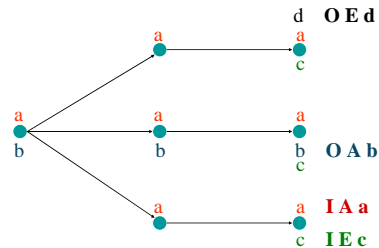
Operatoren für Possible Worlds

- ⇒ **Modale Operatoren für Situationen:**
 - $A a$ (Always) In jeder möglichen Zukunft gilt a
 - $E a$ (Eventually) In einer possible world gilt a
 - z.B. $E B p$ (irgendwann glauben, dass es regnet)
- ⇒ **Modale Operatoren für Pfade:**
 - $O f$ (optional) für mindestens einen Pfad ist f wahr
 - $I f$ (inevitable) für alle ausgehenden Pfade ist f wahr
 - z.B. $O A B p$ (in einem möglichen Leben immer glauben, dass es regnet)

Possible Worlds in der CTL

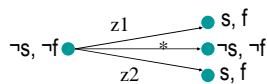
b: Cindy ist in Berlin
a: Berlin ist Hauptstadt Deutschlands

c: Cindy ist in Celle
d: es ist Herbst



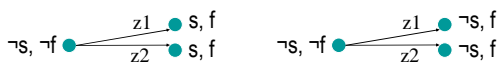
Beispiel für den Einsatz der Possible Worlds Semantik

- ⇒ Ein Agent glaubt, dass es unvermeidlich ist, dass eine Zahnfüllung (f) durch Schmerz (s) begleitet wird.



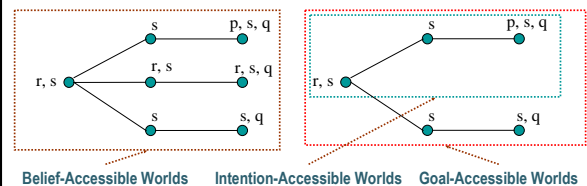
s: Schmerz;
 f: Füllung;
 z1: Zahnarzt 1;
 z2: Zahnarzt 2;
 *: andere Aktion

- ⇒ Der Agent hat das Ziel, eine Zahnfüllung zu bekommen, aber nicht notwendigerweise das Ziel, Schmerz zu erleiden.



Beziehungen zwischen B, D und I

- ⇒ Goal-Accessible Worlds sind Subwelten der Belief-Accessible Worlds eines Agenten.
- ⇒ Intention Accessible Worlds sind Subwelten der Goal-Accessible Worlds.



Gliederung	
⇒	Einführung in Modallogik
⇒	Modallogiken für Zeit und Handeln
⇒	Intentionalität
⇒	Epistemische Logik mit Belief, Desire, Intention
→	Modaloperatoren Belief, Desire, Intention
→	Axiomatisierungen für B, D, I
⇒	BDI-Agenten
⇒	Zusammenfassung

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 37

KD45: Axiomatisierung für epistemische Logiken	
⇒	$\text{Bel}(p \rightarrow q) \rightarrow (\text{Bel } p \rightarrow \text{Bel } q)$ (K)
⇒	$\text{Bel } p \rightarrow \neg \text{Bel } \neg p$ (D)
⇒	$\text{Bel } p \rightarrow \text{Bel } \text{Bel } p$ (4)
⇒	$\neg \text{Bel } p \rightarrow \text{Bel } \neg \text{Bel } p$ (5)
Inferenzregeln:	
⇒	$p \wedge p \rightarrow q \rightarrow q$ (Modus Ponens)
⇒	$p \rightarrow \text{Bel } p$ (Notwendigkeit)

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 38

Schlussfolgerungen in BDI sind nicht trivial	
Bel (clever(John) AND John=partner(Sally))	impliziert nicht:
Bel (clever(Sally))	
Des (Zahnarztbesuch) AND Bel (Zahnarztbesuch => Schmerz)	impliziert nicht:
Des (Schmerz)	

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 39

Axiomatisierungen von Agenten 1	
⇒	Wenn der Agent eine Formel beabsichtigt, dann muss er sie als Ziel haben und auch an sie glauben (starker Realismus):
→	$\text{Int}_a x \Rightarrow \text{Des}_a x$ und $\text{Des}_a x \Rightarrow \text{Bel}_a x$
→	x steht z.B. für optionally(eventually(Diplom))
→	(d.h., Agenten haben keine beliebigen Absichten)
⇒	Realismus : Agenten glauben an ihre Intentionen und Ziele
→	$\text{Int}_a x \Rightarrow \text{Bel}_a \text{Int}_a x$ und $\text{Des}_a x \Rightarrow \text{Bel}_a \text{Des}_a x$
⇒	schwacher Realismus : Nicht an die Nichterfüllbarkeit glauben
→	$\text{Int}_a x \Rightarrow \neg \text{Bel}_a \text{Int}_a \neg x$

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 40

Axiomatisierungen von Agenten 2	
⇒	Absichten werden durch Ziele gestützt
→	$\text{Int}_a x \Rightarrow \text{Des}_a \text{Int}_a x$
⇒	Bewusstsein bezüglich der Ereignisse
→	unabhängig vom Ausgang (Gelingen oder Scheitern)
→	$\text{done}(e) \Rightarrow \text{Bel}_a \text{done}(e)$
→	$\text{done}(e)$: Ereignis e hat stattgefunden (erfolgreich oder nicht)
⇒	Fallen lassen von Absichten nach endlicher Zeit
→	$\text{Int}_a x \Rightarrow \text{I E } \neg \text{Int}_a x$

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 41

Verpflichtung (commitment) gegenüber Intentionen	
⇒	Intentionen lösen Handlungen aus:
→	Aktivierung eines mit der Intention verknüpften Handlungsplans.
⇒	Entsprechend ist ein Agent gegenüber seinen Intentionen zu Handlungen verpflichtet (committed).
⇒	Agenten können Intentionen entweder aufrechterhalten oder verwerfen.
→	z.B. in Abhängigkeit davon, ob die Handlung erfolgreich durchgeführt wurde.

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 42

Drei Arten des Commitments

- ⇒ Ein **blindly committed** Agent verfolgt seine Intentionen, bis er glaubt, sie erreicht zu haben.

$$Int_a I E \varphi \Rightarrow I Int_a I E \varphi \cup Bel_a \varphi$$

- Es ist unabdingbar, dass die Intention aufrechterhalten wird, bis („until“) der Agent φ erreicht zu haben glaubt.
- Was geschieht aber, wenn der Agent $BEL(\neg\varphi)$ in seiner Wissensbasis hat?

Drei Arten des Commitments

- ⇒ Ein **single-minded** Agent behält seine Absichten, bis er glaubt, dass er sie erreicht hat oder niemals erreichen wird.

$$Int_a I E \varphi \Rightarrow I Int_a I E \varphi \cup (Bel_a \varphi \vee \neg Bel_a O E \varphi)$$

- Die Intention wird fallen gelassen, wenn der Agent nicht mehr an die Erfüllbarkeit des Zustands glaubt.
- Was bedeutet das für ein Ziel $Des_a E \varphi$?

Drei Arten des Commitments

- ⇒ Ein **open-minded** Agent bleibt seinen Intentionen solange verpflichtet, bis er sie als erreicht ansieht oder er keine **entsprechenden Ziele** mehr hat.

$$Int_a I E \varphi \Rightarrow I Int_a I E \varphi \cup (Bel_a \varphi \vee \neg Des_a O E \varphi)$$

- Die Intention wird auch dann fallen gelassen, wenn das unterstützende Ziel nicht mehr existiert.

Axiomatisierungen von Agenten 3

- ⇒ Ist eine Formel notwendigerweise gültig, dann soll der Agent nicht gezwungen sein, sie als Ziel oder Absicht aufzunehmen.

$$\rightarrow Bel_a I \varphi \not\Rightarrow Des_a I A \varphi$$

- ⇒ Ist eine Formel $(\varphi \rightarrow \psi)$ notwendigerweise gültig und hat der Agent das Ziel (oder die Absicht) φ , dann soll er nicht gezwungen sein, auch ψ als Ziel (oder Absicht) zu haben.

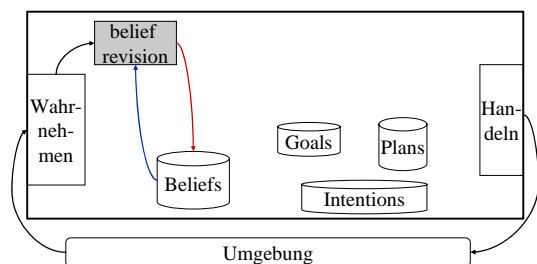
$$\rightarrow Bel_a I (\varphi \rightarrow \psi) \wedge Des_a I \varphi \not\Rightarrow Des_a I \psi$$

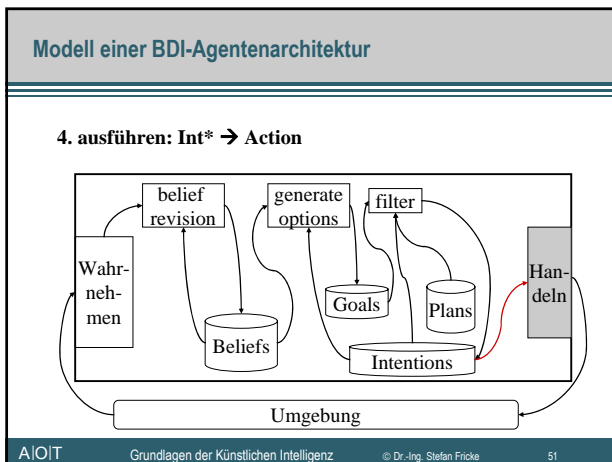
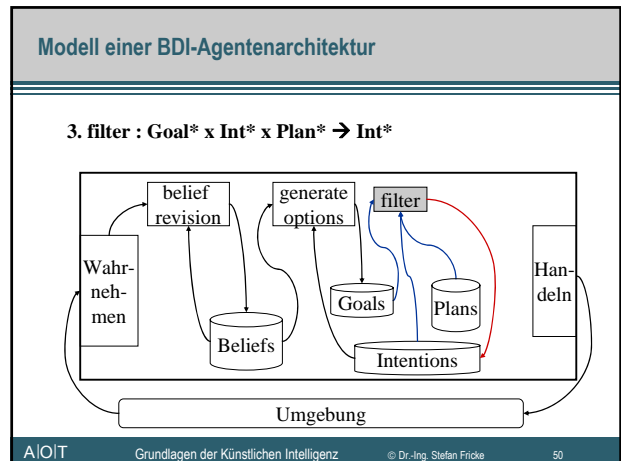
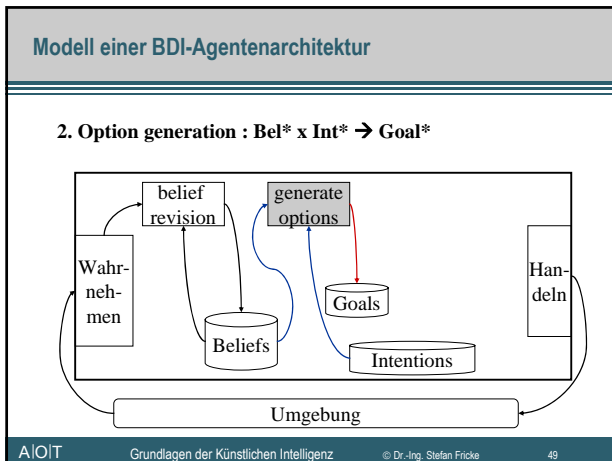
Gliederung

- ⇒ Einführung in Modallogik
- ⇒ Modallogiken für Zeit und Handeln
- ⇒ Intentionalität
- ⇒ Epistemische Logik mit Belief, Desire, Intention
- ⇒ **BDI-Agenten**
 - Allgemeine BDI-Architektur
 - Agent-0
- ⇒ Zusammenfassung

Modell einer BDI-Agentenarchitektur

1. Belief revision function: $Bel^* \times Perception \rightarrow Bel^*$





- ### Gliederung
- ⇒ Einführung in Modallogik
 - ⇒ Modallogiken für Zeit und Handeln
 - ⇒ Intentionalität
 - ⇒ Epistemische Logik mit Belief, Desire, Intention
 - ⇒ BDI-Agenten
 - Allgemeine BDI-Architektur
 - Agent-0
 - ⇒ Zusammenfassung
- AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 52

AGENT-0 [Shoham 1993]

Agent-Oriented Programming besteht nach Shoham aus:

- ⇒ logischem System zur Beschreibung mentaler Zustände,
- ⇒ interpretierter Programmiersprache für Agenten sowie
- ⇒ Agentifizierungsprozess: Kette von Übersetzungsschritten.

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 53

AGENT-0: Logisches System

AGENT-0 unterstützt folgende Sprachelemente

- ⇒ belief (mental) - BEL
- ⇒ commitment (mental) - CMT
- ⇒ capability (nicht mental) - CAN
- ⇒ Grammatik = modale Prädikatenlogik, erweitert um Zeitpunkte

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 54

mentaler Belief-Operator

BEL(<agent>, <timepoint>, <fact>)

BEL(a, t1, BEL(b, t2, like(a, b, t3)))

„Agent *a* glaubt zum Zeitpunkt *t1*, dass Agent *b* zum Zeitpunkt *t2* an *like(a, b, t3)* glaubt“

nicht-mentaler Capability-Operator

CAN(<agent>, <timepoint>, <fact>)

CAN(a, t1, open(door, t2))

Fähigkeit wird durch Zustand ausgedrückt

„Agent *a* kann zum Zeitpunkt *t1* sicherstellen, dass die Tür zum Zeitpunkt *t2* offen ist“

Mentaler Commitment-Operator

CMT(<agent>, <agent>, <timepoint>, <fact>)

CMT(a, b, t1, open(door, t2))

„Agent *a* ist Agent *b* zum Zeitpunkt *t1* zur Türöffnung verpflichtet“

Logisches System: Annahmen

⇒ Konsistenz von Beliefs und Commitments

⇒ Good faith:

CMT(a,b,t,x) ⇔ BEL(a, t, CAN(a, t, x))

⇒ Introspection:

CMT(a,b,t,x) ⇔ BEL(a, t, CMT(a, b, t, x))

¬CMT(a,b,t,x) ⇔ BEL(a, t, ¬CMT(a, b, t, x))

Die AGENT-0 -Sprache: communicative actions und private actions

⇒ **inform(t, a, fact)**

→ zum Zeitpunkt *t* dem Agenten *a* das *fact* zusenden

⇒ **request(t, a, action)**

→ Agent *a* zu Handlung *action* auffordern

→ unrequest und refrain analog.

⇒ **DO(t, p-action)**

→ zum Zeitpunkt *t* die lokale Methode *p-action* ausführen

Die AGENT-0 -Sprache: Commitment Rules

COMMIT(msgcond, mtlcond, agent, action)

msgcond = (From, Type, Content)

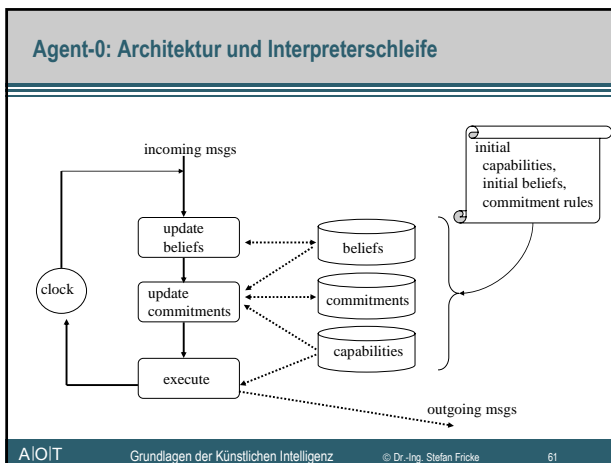
„aufgrund einer empfangenen Nachricht in einem bestimmten Zustand einem anderen Agenten zu einer Aktion verpflichtet sein.“

existenzquantifizierte Variable

**COMMIT((?a, REQUEST, ?action),
(BEL [!t, myfriend(?a)]),
?a,
?action)**

allquantifizierte Variable

„requests von Freunden werden ausgeführt“



Agent-0: Interpreterschleife: 1. Belief-Update

betrifft inform-Sprechakte:

inform(t1, a, fact(x, t2))

- Einfügen des Faktus in die Belief-Datenbasis
- Entfernen inkonsistenter Fakten

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 62

Agent-0: Interpreterschleife: 2. Commitment-Update

commit(msgcond,mtlcond,a,action)

Die Verpflichtung wird angenommen, wenn

- msgcond mappt msg,
- mtlcond ist wahr,
- action ist aktuelle capability,
- Agent ist nicht zu REFRAIN(action) verpflichtet.

unrequest(action): Entfernen der Handlung aus Commitment-DB

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 63

Agent-0: Interpreterschleife: 3. Execute

3. Ausführung von Commitments in Abhängigkeit der Zeit:

INFORM(t, b, fact)

- senden
- assert(BEL(t,a, BEL(t,b,fact)))

REQUEST, UNREQUEST

senden

DO(t, action)

- action checken,
- ggfs. ausführen

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 64

Gliederung

- ⇒ Einführung in Modallogik
- ⇒ Modallogiken für Zeit und Handeln
- ⇒ Epistemische Logiken
- ⇒ Zusammenfassung

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 65

	Syntax	Inferenz	Semantik
Modal	Möglichkeit \Diamond und Notwendigkeit \Box	Inferenzregeln für \Diamond und \Box	Possible Worlds
Dynamisch	Sequenz, Verzweigung, Test	Inferenzregeln für Ergebnisse	Possible Worlds mit Übergängen
Temporal	Zeitpunkte oder Intervalle, zeitliche Beziehungen	Inferenz für Aussagen mit zeitlicher Ergänzung	Possible Worlds mit vielfältigen Übergängen

AIOIT Grundlagen der Künstlichen Intelligenz © Dr.-Ing. Stefan Fricke 66

Zusammenfassung Intentionalität

- ⇒ Der intentional stance bietet eine pragmatische und natürliche Methode, die Regeln und Verhaltensmuster unabhängig von der tatsächlichen Implementation erkennen lässt.
- ⇒ Es geht nicht darum, ob ein Agent wirklich etwas „glaubt“, „wünscht“ oder „beabsichtigt“.
- ⇒ BDI erweist sich als besonders nützlich, wenn ein Agent über andere Agenten „nachdenken“ soll.

AIOIT

Grundlagen der Künstlichen Intelligenz

© Dr.-Ing. Stefan Fricke

67

BDI-Logiken: Zusammenfassung

- ⇒ Der Agent wird als intentionales System beschrieben
- ⇒ Modallogik wird für nicht-monotones Schließen verwendet
 - Modaloperatoren Bel, Des, Int
 - verschiedene Axiomatisierungen sind möglich
- ⇒ BDI überbrückt die Lücke zwischen Wissen und Handeln
 - Aktionen werden über Intentionen, Intentionen über Ziele und Ziele über Wissen motiviert.

AIOIT

Grundlagen der Künstlichen Intelligenz

© Dr.-Ing. Stefan Fricke

68

Vorteile einer *expliziten* Repräsentation von Zielen in Agenten

- ⇒ höhere Fehlertoleranz:
 - Schlägt eine Handlung fehl, kann der Agent besser wieder aufsetzen.
- ⇒ bessere Koordination:
 - Widersprüchliche Ziele nicht gleichzeitig verfolgen.
 - Subsumierende Ziele nur einmal verfolgen.
 - Ständiges Überdenken und Umplanen bei veränderten Randbedingungen ist möglich.

AIOIT

Grundlagen der Künstlichen Intelligenz

© Dr.-Ing. Stefan Fricke

69

BDI-Logiken: Diskussion und Ausblick

- ⇒ Welche Modaloperatoren verwenden?
 - Wissen, Know How, Kommunikation
 - Commitment, Joint Goal, Joint Commitment
- ⇒ Die Axiomatisierungen sind problematisch
 - Constraintformalismen drücken keine Aktivitäten aus
 - Wo kommen Ziele her?
- ⇒ Der Umgang mit Widersprüchen und veralteten Informationen ist unklar.

AIOIT

Grundlagen der Künstlichen Intelligenz

© Dr.-Ing. Stefan Fricke

70

Weitere Probleme klassischer Logik

Zusammenfassung

Klassische Logik ist als Repräsentationsformalismus ungeeignet für viele Aspekte der wirklichen Welt:

- ⇒ zur Repräsentation **kontinuierlicher Größen**,
- ⇒ zur Repräsentation des **Unbekannten**,
- ⇒ zur Repräsentation **unsicheren, probabilistischen Wissens**,
 - „morgen ist 80%ige Regenwahrscheinlichkeit“
- ⇒ zur Repräsentation von **Relativität und Unschärfe**.
 - „Paul ist ziemlich groß“.

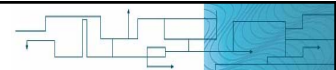
AIOIT

Grundlagen der Künstlichen Intelligenz

© Dr.-Ing. Stefan Fricke

71

DAI-Labor
TU Berlin



Grundlagen der Künstlichen Intelligenz

08.12.2005: Modallogik

Abschluss der ersten Teils GKI
„symbolische KI“

ab nächste Woche: statistische
Verfahren

Dr.-Ing. Stefan Fricke
stefan.fricke@dai-labor.de

AIOIT

Agententechnologien in
betrieblichen Anwendungen
und der Telekommunikation

Referenzen

- ⇒ **Agent-0:**
 - Y. Shoham. Agent-oriented programming. Artificial Intelligence, 60. S. 51 - 92, 1993.
 - www.agentbuilder.com: kommerzielles Agent-0
- ⇒ **BDI:**
 - Modeling Rational Agents within a BDI – Architecture (Anand. S. Rao, Michael P. Georgeff)
 - Andreas Kerlin, Anatolij Zubow, Daniel Göhring: Glauben und Absichten. Seminararbeit HU Berlin, 2002, http://www.drqoehring.de/uni/papers/Glauben_und_Absichten_062002.pdf

Anhänge

BDI-Programm

```
Beliefs = Beliefs,  
Intentions = Intentions, Goals = Goals,  
while true do  
  get next percept p  
  Beliefs = brf(Beliefs,p)  
  Intentions = options(Goals,Intentions)  
  Intentions = filter(Beliefs, Goals, Intentions)  
   $\pi$  = plan(Beliefs, Intentions)  
  while not (empty( $\pi$ ) or succeeded (Intentions, Beliefs) or  
            impossible(Intentions, Beliefs)) do  
     $\alpha$  = head( $\pi$ )  
    execute( $\alpha$ )  
     $\pi$  = tail( $\pi$ )  
    get next percept p  
    Beliefs = brf(Beliefs,p)  
    if not sound( $\pi$ , Intentions, Beliefs) then  
       $\pi$  = plan(Beliefs, Intentions)  
  end while  
end while
```

Dropping intentions that are impossible or have succeeded

Reactivity, replan

Klassische Logik und Agenten

Einleitung

Schwierige Anwendbarkeit der Prädikatenlogik für bestimmte **Umgebungstypen:**

- | | |
|-------------------------|---------------------------|
| ⇒ Nichtdeterministisch | → unzuverlässige Aktionen |
| ⇒ Partiiell beobachtbar | → unvollständiges Wissen |
| ⇒ Kontinuierlich | → Repräsentationsproblem |
| ⇒ Dynamisch | → unzuverlässiges Wissen |

Agent-0 – Zusammenfassung

- ⇒ Der Agent ist im Wesentlichen ein Regelinterpretier.
- ⇒ Ziele werden nicht modelliert.
- ⇒ Zeitpunkte sind problematisch im Zusammenhang mit Wissen.
- ⇒ Wenig Aufmerksamkeit wird dem belief-revision geschenkt.
- ⇒ AgentBuilder ist eine kommerzielle Erweiterung von Agent-0
 - mit Sprechakten und Interaktionsprotokollen
 - www.agentbuilder.com