

VITRA: Verbalisierung visueller Information

Gerd Herzog, Anselm Blocher, Klaus-Peter Gapp,
Eva Stopp und Wolfgang Wahlster
Universität des Saarlandes, Fachbereich Informatik,
Sonderforschungsbereich 314,
Postfach 15 11 50, D-66041 Saarbrücken

Zusammenfassung

Das Projekt VITRA (Visual Translator) beschäftigt sich mit Grundfragen der Beziehung zwischen Sprache und Sehen. Ziel der experimentellen Studien ist die Entwicklung wissensbasierter Systeme zur Integration von visueller Wahrnehmung und der Verarbeitung natürlicher Sprache. Hierbei konnten erstmals automatische sprachliche Beschreibungen für aus realen Bildfolgen gewonnene Trajektorien erzeugt werden. Die vorliegende Arbeit stellt den in VITRA verfolgten Ansatz zur simultanen Auswertung und natürlichsprachlichen Beschreibung zeitveränderlicher Szenen genauer vor. Bei dieser Konzeption wird die Verarbeitung auf allen Stufen in inkrementeller Weise durchgeführt, eine wichtige Voraussetzung für die langfristig angestrebte Echtzeitverarbeitung.

Diese Arbeit ist erschienen in: Informatik Forschung und Entwicklung, Band 11, Heft 1, 1996, S. 12–19.

1 Einleitung

Nachdem die beiden Teilgebiete Bild- und Sprachverstehen in der Künstlichen Intelligenz (KI) lange nahezu unabhängig voneinander behandelt wurden, gewinnt die Entwicklung wissensbasierter Systeme zur Integration von maschinellem Sehen und der Verarbeitung natürlicher Sprache zunehmend an Bedeutung [1], [7], [20], [29], [43]. Dabei ergeben sich vielfältige Anwendungsperspektiven in den unterschiedlichsten Forschungsfeldern: Robotik, Medizintechnik, Fernerkundung, Fahrerassistenzsysteme, intelligente Leitstände und Intellimediasysteme stellen einige wichtige Beispiele aus der Vielzahl potentieller Anwendungsbereiche dar.

Die Beziehung zwischen Sprache und Wahrnehmung bildet den Forschungshintergrund für das Projekt VITRA (Visual Translator), das sich mit dem Entwurf und der Realisierung von wissensbasierten Systemen für den natürlichsprachlichen Zugang zu visueller Information auseinandersetzt [19]. Langfristiges Hauptziel der in VITRA verfolgten Forschungsrichtung ist es, die komplexen Informationsverarbeitungsprozesse, welche der Interaktion von visueller Wahrnehmung und Sprachproduktion zugrunde liegen, algorithmisch zu beschreiben und zu erklären [42]. Neben diesem kognitionswissenschaftlichen Ansatz geht es aus der ingenieurwissenschaftlichen Perspektive darum, den Benutzern die Resultate eines bildverstehenden Systems besser zugänglich und leichter verständlich zu machen.

In VITRA werden unterschiedliche Diskursbereiche und verschiedenartige Kommunikationssituationen im Hinblick auf den natürlichsprachlichen Zugang zu visueller Information näher untersucht. In enger Kooperation mit dem Fraunhofer-Institut für Informations- und Datenverarbeitung (IITB), Karlsruhe, [30] sowie dem Institut für Prozeßrechentchnik und Robotik (IPR) der Universität Karlsruhe werden dabei folgende Szenarien betrachtet:

- Beantwortung von Fragen über Beobachtungen in einer Straßenverkehrsszene [3], [37]
- Generierung von Simultanbeschreibungen anhand kurzer Ausschnitte aus Videoaufnahmen eines Fußballspiels [4], [18] sowie verschiedener Straßenverkehrsszenen mit Fahrzeugen und Fußgängern [17]
- Erzeugung inkrementeller, multimodaler Wegbeschreibungen basierend auf einem 3D-Modell des Saarbrücker Campus [15], [28]
- Natürlichsprachliche Interaktion mit einem autonomen mobilen Robotersystem [27], [39]

Der vorliegende Beitrag konzentriert sich auf die automatische Auswertung und Beschreibung zeitveränderlicher Szenen ausgehend von Kamerabildfolgen.

2 Von der Bildfolge zur intentionalen Simultanbeschreibung

Beim natürlichsprachlichen Zugang zu bildverstehenden Systemen müssen drei Verarbeitungsebenen unterschieden werden [16]. Auf der sensorischen Ebene besteht die Aufgabe der Bildanalyse in der Identifizierung und Klassifizierung der beobachtbaren Objekte. Sie leistet damit den Übergang von der Bildfolge zur Szenenfolge im 3D-Raum. Die Szenenfolgenanalyse realisiert den Schritt von der rekonstruierten *geometrischen Szenenbeschreibung* [31] zu einer konzeptuellen Beschreibung des Szenengeschehens. Diese konzeptuelle Ebene dient der referenzsemantischen Verankerung sprachlicher Einheiten [6], [8], [21]. Im Hinblick auf den variierbaren Detaillierungsgrad sprachlicher Beschreibungen sind hierbei, wie in Abb. 1 angedeutet, konzeptuelle Einheiten auf unterschiedlichen Abstraktionsebenen bereitzustellen. Auf der dritten, der linguistischen Ebene werden die an der Wahrneh-

| | | |
|------------------|---|------------------------------|
| Ziel: | Verbalisierung visueller Wahrnehmungen | |
| Ausdruck: | Sprachliche Beschreibung | Erkennung von: |
| | <ul style="list-style-type: none"> • erkannter Objekte: <i>Da sind zwei PKW und ein Bus.</i> | Objekten |
| | <ul style="list-style-type: none"> • der räumlichen Lage von Objekten: <i>Der Bus steht vor der Kirche.</i> | Räumlichen Relationen |
| | <ul style="list-style-type: none"> • von Objektbewegungen: <i>Der Bus fährt an.</i> | Vorgängen |
| | <ul style="list-style-type: none"> • vermuteter Ziele und Pläne der beobachteten Agenten: <i>Der Bus wartet vor der Ampel.</i> | Plänen |

Abbildung 1: Ebenen der sprachlichen Szenenbeschreibung

mung orientierten konzeptuellen Strukturen in natürlichsprachliche Äußerungen überführt. Frühere natürlichsprachliche Zugangssysteme, wie HAM-ANS [44] und NAOS [31], beschränken sich auf eine *a posteriori* Analyse, die lediglich retrospektive Szenenbeschreibungen zulässt. Im Gegensatz dazu wird in VITRA eine *inkrementelle* Analyse durch-

geführt, bei der die Verarbeitung der visuellen Daten simultan zum Fortschreiten der Szenenfolge geschieht. Hierdurch werden Informationen zur aktuellen Szene bereitgestellt und unmittelbare Systemreaktionen, wie z.B. eine natürlichsprachliche Simultanbeschreibung, ermöglicht.

Aufgrund des eingesetzten, flexibleren Analysemodus ergibt sich in unserem Ansatz die in Abb. 2 skizzierte Verarbeitungskaskade. Die Digitisierung und Auswertung der Rohbilder erfolgt mit den durch unsere Projektpartner am IITB entwickelten Systemen [24], [25], [35], [40].

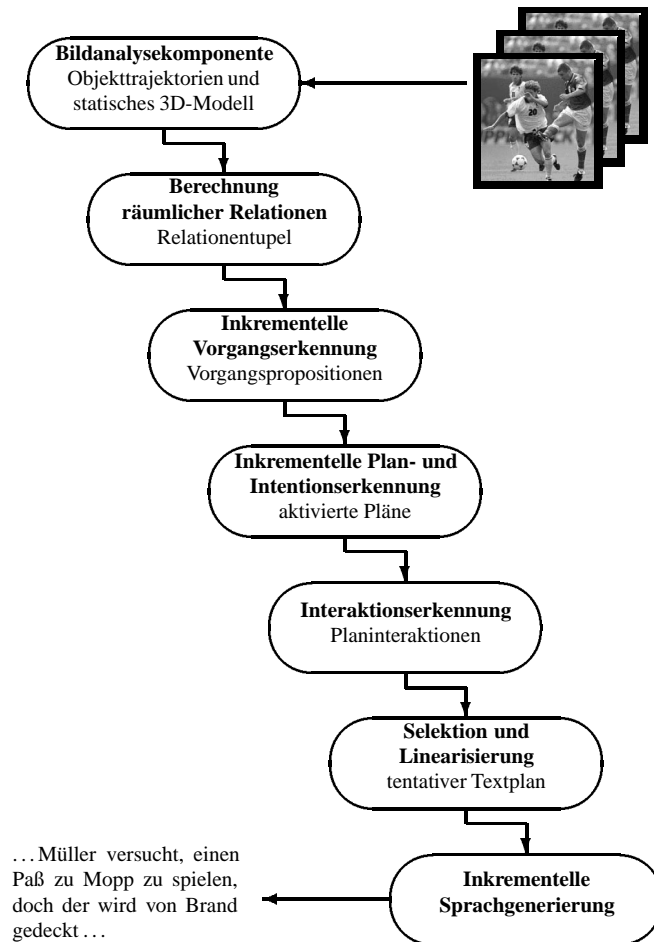


Abbildung 2: Die Verarbeitungskaskade in VITRA

Die Szenenfolgeanalyse in VITRA umfaßt neben der Extraktion räumlicher Beziehungen und interessanter Bewegungsvorgänge auch die Erkennung vermuteter Handlungsintentionen sowie von Planinteraktionen der beobachteten Agenten. Bei der Sprachproduktion müssen die mitteilungswürdigen Propositionen ausgewählt und in einem vorläufigen Textplan angeordnet werden. Unter Rückgriff auf den ständig aktualisierten Textplan werden simultan zum Ablauf der Szene natürlichsprachliche Äußerungen generiert und ausgegeben.

3 Inkrementelle Szenenfolgenanalyse

Ausgehend von der sukzessive durch das Bildanalysesystem bereitgestellten geometrischen Szenenbeschreibung erfolgt die weitergehende Interpretation der Szenenfolge, welche die für eine sprachliche Beschreibung notwendigen höheren konzeptuellen Einheiten, wie räumliche Relationen, auftretende Bewegungsvorgänge, aktivierte Handlungspläne und Planinteraktionen, zur Verfügung stellt.

3.1 Auswertung räumlicher Relationen

Die semantische Analyse räumlicher Präpositionen führt auf den Begriff *räumliche Relation* als einzelsprachunabhängige Bedeutungseinheit. Man definiert räumliche Relationen, indem man Bedingungen über räumliche Gegebenheiten einer Objektkonfiguration spezifiziert, wie z.B. Distanz zwischen Objekten (topologisch) oder ihre relative Lage bezüglich einer Orientierung (projektiv) [13], [26], [41].

Die Berechnung der Semantik räumlicher Relationen wird in VITRA in mehreren Stufen vollzogen. Die Basis bildet die *Grundbedeutung*, welche nur die geometrische Objektbeschreibung und die Gebrauchsart der Relation — intrinsisch, extrinsisch oder deiktisch — berücksichtigt [32]. Auf höherer Ebene wird mittels zusätzlichem kontextspezifischem konzeptuellem Wissen eine erweiterte Semantik räumlicher Relationen ausspezifiziert [9].

Zur Berechnung der Grundbedeutung räumlicher Relationen wird ein von der Ausdehnung des Referenzobjekts und der Gebrauchsart der Relation abhängiges *lokales Koordinatensystem* erstellt (Abb. 3). Dieses lokale Referenzsystem dient der Skalierung des Raums zur Messung von Distanz und Winkel zwischen zu lokalisierendem und Referenzobjekt. Der *Grad* der Anwendbarkeit einer räumlichen Relation wird durch geeignete Abbildung von Distanz bzw. Winkel mittels relationenspezifischer Bewertungsfunktionen auf das Anwendbarkeitsintervall $[0..1]$ ermittelt. Die Durchführung erster experimenteller Untersuchungen bezüglich der Bewertung projektiver Relationen zeigte eine lineare Korrelation zwischen Richtungsabweichung und Anwendbarkeitsgrad bei quadratischen Referenzobjekten [10, 11].

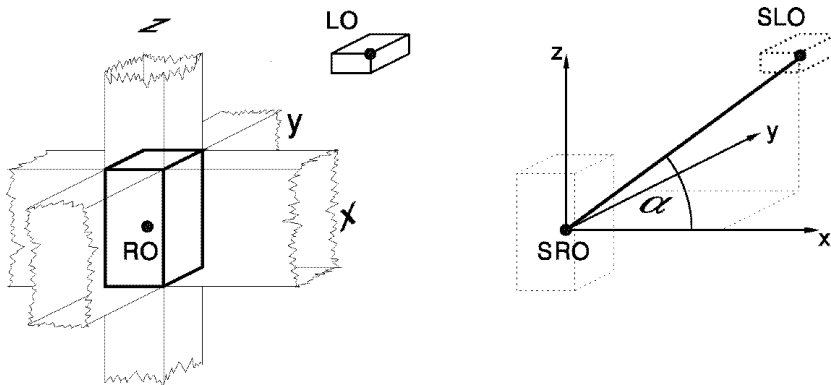


Abbildung 3: Lokales Koordinatensystem und Winkelabweichung

Die graphische Darstellung sogenannter *Anwendbarkeitsfelder* ermöglicht die Visualisierung der Anwendbarkeitsstruktur räumlicher Relationen bezüglich eines bestimmten Referenzobjekts und Kontexts. Abb. 4 gibt die graphische Darstellung zweier Anwendbarkeitsräume. Man beachte, daß hier die durch die projektive Relation bedingte Richtungsabweichung mit einem Distanzkonzept kombiniert wurde.

Auch die Kernsemantik von Relationen, welche sich nicht unmittelbar in Distanz- bzw. Richtungsrelationen kategorisieren lassen, wie z.B. *auf*, *neben* und *zwischen*, lassen sich

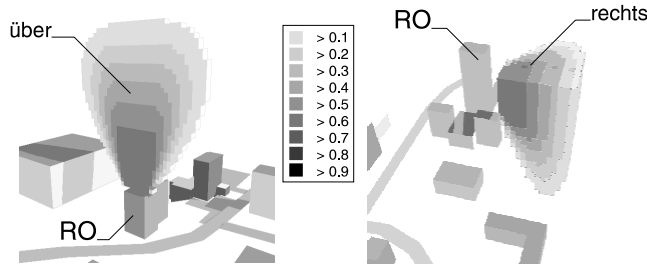


Abbildung 4: Ausschnitte der 3D Anwendbarkeitsstruktur von *über* und *rechts* zur Objektlokalisierung

im vorgestellten Modell formalisieren. So kann die Berechnung der räumlichen Relation *zwischen* beispielsweise auf eine zweimalige Berechnung der Relation *vor* zurückgeführt werden. Die Kombination elementarer Relationen durch die Bildung gewichteter Minima erlaubt die Berechnung zusammengesetzter Relationen, wie z.B. *rechts-hinter* oder *links-bei* [9].

3.2 Charakterisierung von Objektbewegungen

Neben räumlichen spielen zeitliche Aspekte bei der Auswertung einer Szenenfolge eine zentrale Rolle. Aus den erkannten Objekttrajektorien müssen konzeptuelle Beschreibungen abgeleitet werden, welche die relevanten Objektbewegungen in geeigneter Weise charakterisieren. Im Hinblick auf die natürlichsprachliche Beschreibung einer Bildsequenz dienen derartige konzeptuelle Einheiten dazu, die referentielle Bedeutung der korrespondierenden Bewegungs- und Handlungsverben zu erfassen.

Generische Vorgangsmodelle, d.h. deklarative Beschreibungen von Klassen von Bewegungsvorgängen, bilden die Grundlage für die Erkennung relevanter Objektbewegungen. Diese Vorgangskonzepte sind in einer Abstraktionsheterarchie angeordnet, die auf zeitlicher Dekomposition in Teilbewegungen und auf der Spezialisierungsrelation (*laufen* ist z.B. eine spezielle Art von *bewegen*) basiert.

Header: (BALL-TRANSFER ?p1*player ?b*ball ?p2*player)
Conditions: (same (TEAM ?p1) (TEAM ?p2))
Subconcepts: (BALLBESITZ ?p1 ?b) [I1]
(FREI-BEWEGEN ?b) [I2]
(BALLBESITZ ?p2 ?b) [I3]
Temporal-Relations: [I1] :meets [BALL-TRANSFER]
[I1] :meets [I2]
[I2] :equal [BALL-TRANSFER]
[I2] :meets [I3]

Abbildung 5: Vorgangsmodell Ball-Transfer

Abb. 5 zeigt eine vereinfachte Definition des Konzepts *ball-transfer*, bei dem auf den Ballbesitz eines Spielers eine freie Bewegung des Balles folgt, die wiederum mit dem Ballbesitz eines Mitspielers abschließt. Zeitliche Beziehungen werden hierbei mittels der qualitativen Relationen zwischen je zwei Zeitintervallen [2] ausgedrückt.

Um einen Simultanbericht für eine gerade ablaufende Szenenfolge zu ermöglichen, müssen interessante Bewegungsabläufe schritthaltend erkannt und bereits während ihres Auftretens beschrieben werden. Zu diesem Zweck werden die intervall-basierten Vorgangs-

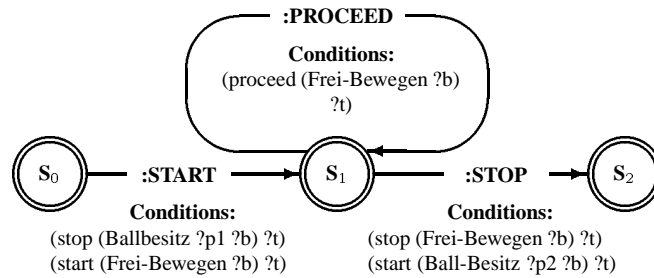


Abbildung 6: Ablaufschema für Ball-Transfer

modelle in eine für eine inkrementelle Erkennung besser geeignete Darstellungsform, sogenannte *Ablaufschemas*, transformiert. Ein Ablaufschema ist ein markierter, gerichteter Graph mit getypten Kanten, der den prototypischen Verlauf eines Vorgangs modelliert, indem die für jeden einzelnen Zeitpunkt zu erfüllenden Bedingungen formuliert werden. Die drei Verlaufsprädikate *start*, *proceed* und *stop* dienen dazu, den aktuellen Zustand eines gerade auftretenden Bewegungsvorgangs zu charakterisieren. Die inkrementelle Erkennung einer Vorgangsausprägung entspricht der schrittweisen Traversierung des zugehörigen Ablaufschemas. Die automatische Konstruktion der Ablaufschemas erfolgt mit Methoden des zeitlichen Schließens [14]. Abb. 6 zeigt das für das Konzept *Ball-Transfer* erzeugte Ablaufschema.

3.3 Intentionale Interpretation

Beim Beschreiben von Handlungen und Geschehnissen beschränkt sich ein menschlicher Beobachter im allgemeinen nicht auf visuelle Informationen, sondern bezieht auch Absichten und Ziele der Akteure in seine Schilderung mit ein. Diese intentionale Interpretation spiegelt sich in Sätzen wie „A will X tun.“, „A tat X, um Y zu erreichen.“ oder einfach „A verfolgt B.“ wider.

Ein Kriterium für die Wahl der Fußballdomäne in VITRA war die Tatsache, daß der Einfluß der bei den Agenten vermuteten Intentionen auf die sprachliche Beschreibung besonders offenkundig erscheint. Unter Berücksichtigung von Spielposition, Mannschaftszugehörigkeit und der Rollenverteilung in Standardsituationen können in jeder Situation Annahmen über stereotypische Intentionen gemacht werden.

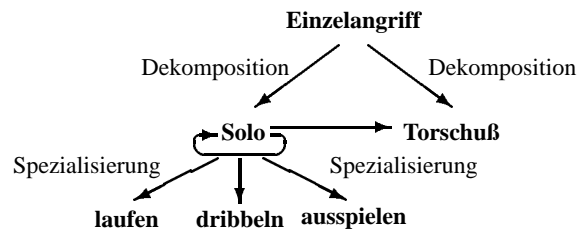


Abbildung 7: Ausschnitt einer Planhierarchie in VITRA

Für die intentionale Interpretation wird in VITRA auf Planerkennungsmethoden zurückgegriffen [33], [34]. Das Wissen über Ziele und Pläne ist in einer Planhierarchie repräsentiert, die durch Aktionen und deren hierarchische sowie zeitliche Beziehungen gebildet wird. Innere Knoten stellen abstrakte Aktionen dar. Die Blätter repräsentieren elementare, d.h. wahrnehmbare, generische Aktionen. Die Knoten enthalten zusätzlich Informationen

über notwendige Vorbedingungen und den intendierten Effekt der korrespondierenden Aktion. Jeder zusammenhängende Teilgraph, der von einer zeitlich dekomponierbaren Aktion bis zu deren Blättern reicht, entspricht einem Plan. Zeitliche Dekomposition und Spezialisierungsrelation erlauben in Plänen die Verwendung von Aufeinanderfolge, Alternativen und Wiederholungen. Abb. 7 zeigt einen beispielhaften Ausschnitt einer Planhierarchie. Die unmarkierten Zeitkanten repräsentieren dabei eine unmittelbare zeitliche Aufeinanderfolge.

Unter Einsatz von Fokussierungsheuristiken werden potentielle Pläne inkrementell erkannt, wobei drei Arten von vermuteten Intentionen abgeleitet werden können. Der intendierte Effekt des aktuellen Planschrittes kann als zustandsgerichtete Intention aufgefaßt werden („Bommer greift Mopp an, um den Ball zu bekommen“). Die aufgrund einer Planhypothese erwarteten zukünftigen Aktionen stellen handlungsgerichtete Intentionen dar („Mopp will Maier den Ball zuspiesen.“). Die Wurzel des aktuellen Plans definiert das vom Agenten verfolgte Oberziel („Mopp startet einen Einzelangriff.“).

Ausgehend von den beobachteten Plänen können in einem weiteren Verarbeitungsschritt kooperative („Während Mopp einen Einzelangriff startet, läuft sich Maier frei.“) und antagonistische Interaktionen („Als Mopp nach vorne läuft, greift Bommer ihn an.“) zwischen den Plänen mehrerer Akteure erkannt werden. Neben diesen jeweils simultanen oder sequentiellen Interaktionen gibt es die Möglichkeit, kollektive Aktionen darzustellen.

Aufbauend auf der Intentions- und Interaktionserkennung lassen sich auch Planfehlschläge interpretieren. Eine Reihe von Fehlschlagsursachen können in VITRA bereits erkannt werden [34]. Ein Plan kann fehlschlagen aufgrund falscher Annahmen des Agenten hinsichtlich der Erfülltheit von Vorbedingungen oder bedingt durch einen antagonistischen Plan eines Kontrahenten. Im Falle einer kooperativen Interaktion kann ein Versagen des Partners vorliegen.

4 Generierung natürlichsprachlicher Simultanbeschreibungen

Eine Simultanbeschreibung ist im Gegensatz zu einer retrospektiven Beschreibung dadurch geprägt, daß der vollständige Verlauf der betrachteten Szenenfolge zum jeweiligen Zeitpunkt der Textgenerierung noch nicht bekannt ist. Da sich die Beschreibung auf das aktuelle Geschehen konzentrieren soll, muß über Vorgänge und Aktionen bereits berichtet werden während diese ablaufen und eventuell noch nicht vollständig erkannt wurden. Diese Besonderheiten der Kommunikationssituation machen eine inkrementelle Verarbeitungsstrategie erforderlich, bei der mit der Verbalisierung begonnen wird, bevor der Inhalt einer Äußerung bis ins letzte Detail geplant werden kann.

Die Sprachproduktion in VITRA umfaßt Prozesse zur Selektion, Linearisierung und Einkodierung von Propositionen. Das hierbei verwendete Hörermodell stellt eine Imaginationskomponente bereit, die dazu dient, die beim Hörer vermutete visuelle Konzeptualisierung der beschriebenen Szene zu antizipieren, damit eine möglichst adequate Beschreibung erzeugt werden kann.

4.1 Selektion und Linearisierung von Propositionen

Da die betrachtete zeitveränderliche Szene schritthaltend beschrieben werden soll, unterliegt die Sprachproduktion starken zeitlichen Restriktionen. Das System kann folglich nicht über alle Vorgänge und Aktionen berichten, die es erkannt hat, sondern muß sich für diejenigen entscheiden, die es dem Hörer ermöglichen, dem Szenengeschehen zu folgen. Entsprechend dem Kooperationsprinzip von Grice [12] sollte der Hörer dabei unter Vermeidung redundanter Information über alle relevanten Fakten informiert werden.

Die Relevanz einer Proposition wird durch vielfältige Faktoren bestimmt und ändert sich dynamisch mit dem Fortschreiten der Szene. Zu den besonders wichtigen Aspekten zählen hierbei der Informationsgehalt, der mit der Komplexität des zugehörigen generischen Konzepts korreliert, die Auftrittshäufigkeit, die Aktualität und der Erkennungszustand. Um Redundanz zu vermeiden, werden solche Propositionen, die durch andere schon erwähnte Propositionen impliziert werden, nicht selektiert. Beispielsweise ist für den Hörer evident, daß nach einem Zuspiel von Mopp zu Maier letzterer in Ballbesitz ist. Die taxonomische Organisation der konzeptuellen Wissensbasis bildet eine wichtige Voraussetzung zur Ableitung derartiger Zusammenhänge.

Weitere Selektionsprozesse im Rahmen der Enkodierung betreffen die Wahl geeigneter Deskriptionen für Objekte, Lokationen und temporale Angaben, um systeminterne Repräsentationsstrukturen in natürlichsprachliche Äußerungen überführen zu können. Entscheidend ist dabei der stetige Rückgriff auf ein Textgedächtnis und das Hörermodell, deren Inhalte die Auswahlprozesse maßgeblich beeinflussen.

Der Linearisierungsprozeß bestimmt die Reihenfolge, in der die ausgewählten Propositionen im Text erwähnt werden sollen und legt damit einen vorläufigen Textplan fest. Die Linearisierung erfolgt primär unter dem Gesichtspunkt der zeitlichen Abfolge der korrespondierenden Vorgänge und Aktionen, wobei zusätzlich Fokuskriterien zur Steigerung der Textkohärenz Berücksichtigung finden.

4.2 Antizipation der Hörervorstellung

Im Hörermodell werden die geplanten Äußerungen mit angenommenen Hörervorstellungen in einer Antizipationsrückkopplungsschleife [22] ausgehend von bereits mitgeteilten Szenenabschnitten verglichen, um Verständnisprobleme bzw. redundante Informationen auszuschließen.

Das Hörermodell muß daher in der Lage sein, im aktuellen Kontext eine möglichst typische Interpretation der geplanten Äußerung zu finden [36], [38]. Dazu werden räumliche Propositionen mit Hilfe des Gradientenverfahrens auf die in Abschnitt 3.1 beschriebenen Anwendbarkeitsräume – hier interpretiert als Typikalitätsverteilungen – approximiert, so daß sie die gegebene Beschreibung im aktuellen Kontext maximaltypisch erfüllen (Abb. 8). Die Vorgehensweise bei zeitlichen Relationen ist analog.

Die zur Verbalisierung anstehende Proposition wird anhand des zugehörigen Vorgangsmodells in ihre *propositionale Elementarstruktur* überführt. Die entstehenden zeitpunktweise zusammengefaßten Mengen von räumlichen Propositionen werden mittels des Approximationsalgorithmus in die geometrische Repräsentation mentaler Bilder umgewandelt. Die durch die Abfolge der konstruierten Bildvorstellungen entstandene *imaginierte* Szenenfolge wird dann analog zum ursprünglichen Szenenfolge reanalysiert.

Zum Vergleich der intendierten und imaginierten Äußerungsgehalte werden im folgenden – wie in Abb. 9 dargestellt – verschiedene Mengen von Propositionen aus den konzeptuellen Beschreibungen von realer und imaginierten Szene gegenübergestellt [5]. Dieser Vergleich findet *nicht* auf der Ebene absoluter Koordinaten statt, da in diesem Fall der Kontext nicht berücksichtigt werden könnte. Ziel ist es, den vorläufigen Textplan hinsichtlich der folgenden drei pragmatischen Kriterien zu überprüfen und gegebenenfalls Verbesserungen vorzuschlagen:

- *Referenz:* Die verwendeten Objektbeschreibungen müssen für Sprecher und Hörer gleichermaßen verständlich sein. Redundanz soll vermieden werden.
- *Plausibilität:* Die neue Information muß in den gegebenen Kontext integriert werden können.
- *Korrektheit:* Das intendierte und das antizipierte Verständnis müssen übereinstimmen. Falsche Inferenzen durch den Hörer sollen vermieden werden.

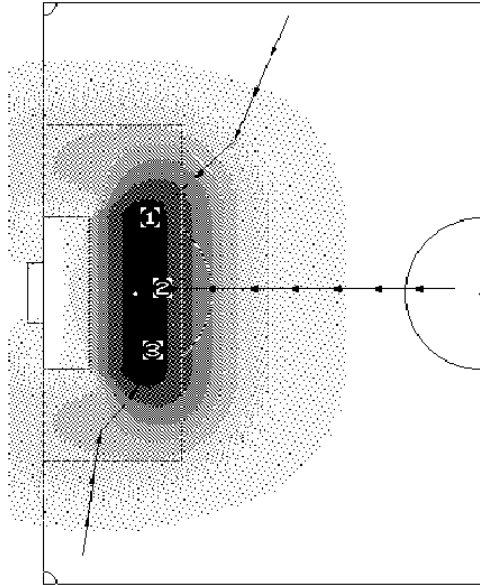


Abbildung 8: Ein Beispiel zu kontextsensitiver Approximation: (vor Spieler linker_Torraum)

Die in Abb. 10 noch einmal schematisch dargestellte Antizipationsrückkopplungsschleife resultiert nötigenfalls in einer Änderung des Textplans. Hierunter fallen insbesondere die Modifikation optionaler Tiefenkasus und die Generierung von Klärungseinschüben.

4.3 Inkrementelle Verbalisierung

Bei der Enkodierung einer ausgewählten Proposition, d.h. der Umwandlung einer konzeptuellen Beschreibung in natürlichsprachliche Äußerungen, wird zunächst ein geeignetes Verb ausgewählt und der damit assoziierte Kasusrahmen instantiiert. Die Lexikalisierung erfolgt unter Zugriff auf ein Konzeptlexikon, das die Verbindung zwischen außersprachlichen und sprachlichen Einheiten realisiert. Ergänzende Selektionsprozesse entscheiden, welche Information bezogen auf die einzelnen Kasusfüller mitgeteilt werden soll. Die selektierte Information wird in natürlichsprachliche Ausdrücke umgewandelt, die auf Zeit, Ort und Objekte referieren. In Abhängigkeit vom Textgedächtnis ist dabei gegebenenfalls auch Anaphernbildung möglich. Durch Zeitformen bzw. mittels Temporaladverbien ist es möglich, temporale Information auszudrücken; räumliche Information wird analog dazu durch geeignete Präpositionalphrasen beschrieben. Die geeignete Selektion von Attributen für interne Objektbezeichner ermöglicht dem Hörer eine eindeutige Identifikation des intendierten Referenten. Charakterisierende Attribute können im Hörermodell eingetragenes Vorwissen, die räumliche Position oder auch vorerwähnte Vorgänge und Aktionen sein, an denen das Objekt beteiligt war.

Die aus der Wortwahl und der Bestimmung morphosyntaktischer Information resultierende präverbale Struktur wird Stück für Stück an den Oberflächengenerator übergeben, der für die grammatische Enkodierung, die Linearisierung und die Flexion verantwortlich ist. Für die syntaktische Beschreibungsebene wird eine lexikalisierte Baumadjunktionsgrammatik mit Merkmalsunifikation (LTAG) verwendet [23]. Der erweiterte Lokaltätsbereich

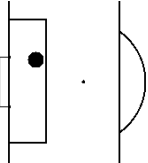
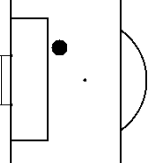
| Visuelle Perzeption | | Imaginierte Bildvorstellung | |
|---|--------|---|--------|
|  | |  | |
| Proposition | A-Grad | Proposition | A-Grad |
| (in Mopp linke_Spielhälfte) | 1.0 | (in Mopp linke_Spielhälfte) | 1.0 |
| (in Mopp linker_Torraum) | 1.0 | (vor Mopp linker_Torraum) | 0.9 |
| ... | | ... | |
| Analysierender Vergleich | | | |
| Differenz-Propositionenpaar | | $\left. \begin{array}{l} ((in \text{ Mopp linker_Torraum}) 1.0)_{vp} \\ ((vor \text{ Mopp linker_Torraum}) 0.9)_{im} \end{array} \right\} \neq$ | |
| Zu explizierende Proposition (in Mopp linker_Torraum) | | | |
| Mögliche Unterspezifikation linker_Torraum \rightarrow Torraum | | | |
| Zu generierende Proposition (in Mopp Torraum) | | | |

Abbildung 9: Ein Beispiel der Verfeinerung des Textplans durch Antizipation der Hörervorstellung

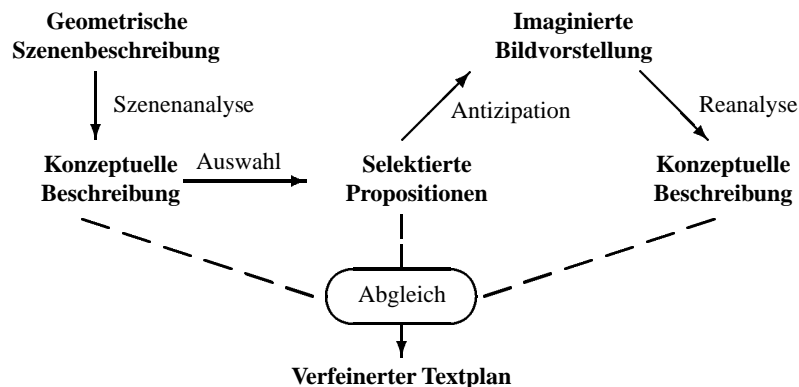


Abbildung 10: Antizipationsrückkopplungsschleife

sowie die flexible Expansion partieller Strukturen durch Substitution und Adjunktion machen LTAGs zu geeigneten Kandidaten für die inkrementelle syntaktische Generierung.

5 Zusammenfassung und Ausblick

Die sprachorientierte KI-Forschung im Bereich der Integration von Sprache und Wahrnehmung zielt auf eine (extreme) operationale Form einer Referenzsemantik ab, die bis zur sensorischen Ebene hinunterreicht. VITRA stellt das bislang einzige System dar, mit dem natürlichsprachliche Beschreibungen für Objekttrajektorien, die aus Realweltbildfolgen abgeleitet wurden, automatisch generiert werden können. Die Szenenfolgenanalyse in VITRA beschränkt sich dabei nicht allein auf rein visuelle Aspekte, sondern schließt auch die Ableitung vermuteter Handlungsintentionen der beobachteten Agenten mit ein.

In enger Kooperation mit unseren Projektpartnern am IITB wird in VITRA ein Ansatz zur Kopplung von bildverstehenden und sprachverstehenden Systemen verfolgt, bei dem die Verarbeitungsprozesse auf allen Ebenen simultan zu der betrachteten Bild- bzw.

Szenenfolge ablaufen. Langfristig wird damit eine Echtzeitverarbeitung angestrebt, die es ermöglichen soll, visuelle Information unmittelbar zu interpretieren, um Systemreaktionen, wie z.B. sprachliche Äußerungen oder motorische Aktionen eines Robotersystems, basierend auf den jeweils vorliegenden Analyseergebnissen sofort durchführen zu können.

Als Einschränkung werden zur Zeit nur Aufnahmen mit stationärer Kamera untersucht, sodaß eine Rückkoppelung von der Szenenfolgenanalyse zum Bildauswertungssystem bislang nicht realisiert werden kann. Im Hinblick auf die geplante aktive Sensorsteuerung bietet der in VITRA konsequent realisierte Ansatz der inkrementellen Verarbeitung jedoch die notwendigen Voraussetzungen. Die kontextabhängige Auswahl relevanter Objektklassen, partiell erkannte Vorgänge, aktivierte Handlungspläne und imaginierte Folgeszenen können notwendige Parameter zur Sensorsteuerung liefern.

Technische Anmerkungen

Bei den experimentellen Arbeiten zur Realisierung des natürlichsprachlichen Zugangssystems VITRA wird vorwiegend objekt-orientierte Programmierung in der KI-Programmiersprache Common Lisp und dem Common Lisp Object System (CLOS) eingesetzt. Die Realisierung der graphischen Systemoberfläche basiert auf dem Common Lisp Interface Manager (CLIM). Die Systementwicklung erfolgt auf Symbolics 36xx Lisp-Maschinen, Symbolics UX1200S Lisp-Coprozessoren und auf Hewlett Packard 9720 sowie SPARC-basierten Arbeitsplatzrechnern.

Danksagung

Die hier beschriebenen Forschungsarbeiten wurden von der Deutschen Forschungsgemeinschaft im Rahmen des Sonderforschungsbereichs 314 „Künstliche Intelligenz und wissensbasierte Systeme“, Projekt N2: VITRA, gefördert.

Literatur

- [1] Artificial Intelligence Review Journal, 8, Special Volume on the Integration of Natural Language and Vision Processing, 1994.
- [2] J. F. Allen. Towards a General Theory of Action and Time. *Artificial Intelligence*, 23(2):123–154, 1984.
- [3] E. André, G. Bosch, G. Herzog und T. Rist. Characterizing Trajectories of Moving Objects Using Natural Language Path Descriptions. In: *Proc. of the 7th ECAI*, volume 2, pp. 1–8, Brighton, UK, 1986.
- [4] E. André, G. Herzog und T. Rist. On the Simultaneous Interpretation of Real World Image Sequences and their Natural Language Description: The System SOCCER. In: *Proc. of the 8th ECAI*, pp. 449–454, Munich, 1988.
- [5] A. Blocher und J. R. J. Schirra. Optional Deep Case Filling and Focus Control with Mental Images: ANTLIMA-KOREF. In: *Proc. of the 14th IJCAI*, pp. 417–423, Montreal, Canada, 1995.
- [6] I. Carsten und T. Janson. Verfahren zur Evaluierung räumlicher Präpositionen anhand geometrischer Szenenbeschreibungen. Diplomarbeit, Fachbereich für Informatik, Univ. Hamburg, 1985.

- [7] Centre National de la Recherche Scientifique. *Images et Langages: Multimodalité et Modélisation Cognitive, Colloque Interdisciplinaire du Comité National de la Recherche Scientifique*, Paris, 1993.
- [8] M. Fürnsinn, M. N. Khenkhar und B. Ruschkowski. GEOSYS – Ein Frage-Antwort-System mit räumlichem Vorstellungsvermögen. In: C.-R. Rollinger (Hrsg.), *Probleme des (Text-) Verstehens, Ansätze der künstlichen Intelligenz*, pp. 172–184. Tübingen: Niemeyer, 1984.
- [9] K.-P. Gapp. Basic Meanings of Spatial Relations: Computation and Evaluation in 3D Space. In: *Proc. of AAAI-94*, pp. 1393–1398, Seattle, WA, 1994.
- [10] K.-P. Gapp. Angle, Distance, Shape, and their Relationship to Projective Relations. In: J. D. Moore und J. F. Lehman (Hrsg.), *Proc. of the 17th Annual Conference of the Cognitive Science Society*, pp. 112–117. Mahwah, NJ: Lawrence Erlbaum, 1995.
- [11] K.-P. Gapp. An Empirically Validated Model for Computing Spatial Relations. In: I. Wachsmuth, C.-R. Rollinger und W. Brauer (Hrsg.), *KI-95: Advances in Artificial Intelligence. 19th Annual German Conference on Artificial Intelligence*, pp. 245–256. Berlin, Heidelberg: Springer, 1995.
- [12] H. P. Grice. Logic and Conversation. In: P. Cole und J. L. Morgan (Hrsg.), *Speech Acts*, pp. 41–58. London: Academic Press, 1975.
- [13] A. Herskovits. *Language and Spatial Cognition. An Interdisciplinary Study of the Prepositions in English*. Cambridge, London: Cambridge University Press, 1986.
- [14] G. Herzog. Utilizing Interval-Based Event Representations for Incremental High-Level Scene Analysis. In: M. Aurnague, A. Borillo, M. Borillo und M. Bras (Hrsg.), *Proc. of the 4th International Workshop on Semantics of Time, Space, and Movement and Spatio-Temporal Reasoning*, pp. 425–435, Château de Bonas, France, 1992. Groupe “Langue, Raisonnement, Calcul”, Toulouse.
- [15] G. Herzog, W. Maaß und P. Wazinski. VITRA GUIDE: Utilisation du Langage Naturel et de Représentation Graphiques pour la Description d’Itinéraires. In: *Images et Langages: Multimodalité et Modélisation Cognitive, Colloque Interdisciplinaire du Comité National de la Recherche Scientifique*, pp. 243–251, Paris, 1993.
- [16] G. Herzog, T. Rist und E. André. Sprache und Raum: Natürlichsprachlicher Zugang zu visuellen Daten. In: C. Freksa und C. Habel (Hrsg.), *Repräsentation und Verarbeitung räumlichen Wissens*, pp. 207–220. Berlin, Heidelberg: Springer, 1990.
- [17] G. Herzog und K. Rohr. Integrating Vision and Language: Towards Automatic Description of Human Movements. In: I. Wachsmuth, C.-R. Rollinger und W. Brauer (Hrsg.), *KI-95: Advances in Artificial Intelligence. 19th Annual German Conference on Artificial Intelligence*, pp. 257–268. Berlin, Heidelberg: Springer, 1995.
- [18] G. Herzog, C.-K. Sung, E. André, W. Enkelmann, H.-H. Nagel, T. Rist, W. Wahlster und G. Zimmermann. Incremental Natural Language Description of Dynamic Imagery. In: C. Freksa und W. Brauer (Hrsg.), *Wissensbasierte Systeme. 3. Int. GI-Kongreß*, pp. 153–162. Berlin, Heidelberg: Springer, 1989.
- [19] G. Herzog und P. Wazinski. VIsual TRAnslator: Linking Perceptions and Natural Language Descriptions. *Artificial Intelligence Review*, 8(2/3):175–187, 1994.
- [20] B. Hildebrandt, R. Moratz, G. Rickheit und G. Sagerer. Integration von Bild- und Sprachverstehen in einer kognitiven Architektur. *Kognitionswissenschaft*, 4(3):118–128, 1995.

- [21] M. Hußmann und P. Scheffe. The Design of SWYSS, a Dialogue System for Scene Analysis. In: L. Bolc (Hrsg.), *Natural Language Communication with Pictorial Information Systems*, pp. 143–201. München: Hanser/McMillan, 1984.
- [22] A. Jameson und W. Wahlster. User Modelling in Anaphora Generation. In: *Proc. of the 5th ECAI*, pp. 222–227, Orsay, France, 1982.
- [23] A. Kilger. Using UTAGs for Incremental and Parallel Generation. *Computational Intelligence*, 10(4):591–603, 1994.
- [24] D. Koller. *Detektion, Verfolgung und Klassifikation bewegter Objekte in monokularen Bildfolgen am Beispiel von Straßenverkehrsszenen*. St. Augustin: Infix, 1992.
- [25] H. W. Kollnig. *Ermittlung von Verkehrsgeschehen durch Bildfolgenauswertung*. St. Augustin: Infix, 1995.
- [26] B. Landau und R. Jackendoff. “What” and “Where” in Spatial Language and Spatial Cognition. *Behavioral and Brain Sciences*, 16:217–265, 1993.
- [27] T. Längle, T. C. Lüth, E. Stopp, G. Herzog und G. Kamstrup. KANTRA - A Natural Language Interface for Intelligent Robots. In: U. Rembold, R. Dillman, L. O. Hertzberger und T. Kanade (Hrsg.), *Intelligent Autonomous Systems (IAS 4)*, pp. 357–364. Amsterdam: IOS, 1995.
- [28] W. Maaß. From Visual Perception to Multimodal Communication: Incremental Route Descriptions. *Artificial Intelligence Review*, 8(2/3):159–174, 1994.
- [29] P. McKeivitt (Hrsg.). *Proc. of AAAI-94 Workshop on Integration of Natural Language and Vision Processing*, Seattle, WA, 1994.
- [30] H.-H. Nagel. From Image Sequences Towards Conceptual Descriptions. *Image and Vision Computing*, 6(2):59–74, 1988.
- [31] B. Neumann und H.-J. Novak. NAOS: Ein System zur natürlichsprachlichen Beschreibung zeitveränderlicher Szenen. *Informatik Forschung und Entwicklung*, 1:83–92, 1986.
- [32] G. Retz-Schmidt. Various Views on Spatial Prepositions. *AI Magazine*, 9(2):95–105, 1988.
- [33] G. Retz-Schmidt. Recognizing Intentions, Interactions, and Causes of Plan Failures. *User Modeling and User-Adapted Interaction*, 1:173–202, 1991.
- [34] G. Retz-Schmidt. *Die Interpretation des Verhaltens mehrerer Akteure in Szenenfolgen*. Berlin, Heidelberg: Springer, 1992.
- [35] K. Rohr. Towards Model-based Recognition of Human Movements in Image Sequences. *Computer Vision, Graphics, and Image Processing (CVGIP): Image Understanding*, 59(1):94–115, 1994.
- [36] J. R. J. Schirra. *Bildbeschreibung als Verbindung von visuellem und sprachlichem Raum: Eine interdisziplinäre Untersuchung von Bildvorstellungen in einem Hörermodell*. St. Augustin: Infix, 1994.
- [37] J. R. J. Schirra, G. Bosch, C.-K. Sung und G. Zimmermann. From Image Sequences to Natural Language: A First Step Towards Automatic Perception and Description of Motions. *Applied Artificial Intelligence*, 1:287–305, 1987.
- [38] J. R. J. Schirra und E. Stopp. ANTLIMA—A Listener Model with Mental Images. In: *Proc. of the 13th IJCAI*, pp. 175–180, Chambery, France, 1993.

- [39] E. Stopp, K.-P. Gapp, G. Herzog, T. Längle und T. C. Lüth. Utilizing Spatial Relations for Natural Language Access to an Autonomous Mobile Robot. In: B. Nebel und L. Dreschler-Fischer (Hrsg.), *KI-94: Advances in Artificial Intelligence. 18th German Annual Conference on Artificial Intelligence*, pp. 39–50. Berlin, Heidelberg: Springer, 1994.
- [40] C.-K. Sung und G. Zimmermann. Detektion und Verfolgung mehrerer Objekte in Bildfolgen. In: G. Hartmann (Hrsg.), *Mustererkennung 1986; 8. DAGM-Symposium*, pp. 181–184. Berlin, Heidelberg: Springer, 1986.
- [41] L. Talmy. How Language Structures Space. In: H. Pick und L. Acredolo (Hrsg.), *Spatial Orientation: Theory, Research and Application*, pp. 225–282. New York, London: Plenum, 1983.
- [42] W. Wahlster. One Word Says More Than a Thousand Pictures. On the Automatic Verbalization of the Results of Image Sequence Analysis Systems. *Computers and Artificial Intelligence*, 8(5):479–492, 1989.
- [43] W. Wahlster. Text and Images. In: R. A. Cole, J. Mariani, H. Uszkoreit, A. Zaenen und V. Zue (Hrsg.), *Survey on Speech and Natural Language Technology*. Dordrecht: Kluwer, 1994.
- [44] W. Wahlster, H. Marburger, A. Jameson und S. Busemann. Over-answering Yes-No Questions: Extended Responses in a NL Interface to a Vision System. In: *Proc. of the 8th IJCAI*, pp. 643–646, Karlsruhe, FRG, 1983.