

Differenzbasierte Repräsentation räumlicher Relationen zur probabilistischen Szenenerkennung mittels hierarchischen Constellation Models

Bachelorarbeit
von

Joshua Enrico Link

An der Fakultät für Informatik
Institut für Anthropomatik und Robotik
Lehrstuhl Prof. Dr.-Ing. R. Dillmann

Erstgutachter: Prof. Dr.-Ing. R. Dillmann
Zweitgutachter: Prof. Dr.-Ing. Björn Hein
Betreuernder Mitarbeiter: Dipl.-Inform. Pascal Meißner

Bearbeitungszeit: 11. Juni 2017 – 10. September 2017

Hiermit erkläre ich an Eides statt, dass ich die von mir vorgelegte Arbeit selbstständig verfasst habe, dass ich die verwendeten Quellen, Internet-Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen – die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Karlsruhe, den (**Datum**)

ToDo

Joshua Enrico Link

Inhaltsverzeichnis

1 Einführung	1
2 Motivation und Problemstellung	2
2.1 Motivation	2
2.2 Fokus der Arbeit	3
3 Grundlagen	5
3.1 PSM - Probabilistic Scene Model	5
3.1.1 Aufbau	6
3.1.2 Learner - Anlernen der Daten	6
3.1.3 Model - Szenenmodell	7
3.1.4 Inference - Szenenerkennung	10
4 Konzept	12
4.1 Ansatz	12
4.2 Erkennungsalgorismus	13
4.2.1 Algorithmus: Beschreibung	13
4.2.2 Algorithmus: Pseudocode	15
4.3 Wahrscheinlichkeitsabschätzung	16
4.3.1 Objektwahrscheinlichkeit	17
4.3.2 Szenenwahrscheinlichkeit	18
5 Implementierung	21
5.1 Umbau PSM	21
5.1.1 Klassenaustausch	21
5.1.2 Datenbankeinbindung	23
5.2 Differenzbasierter Erkennungsalgorismus	24
5.2.1 Algorithmus	25
5.2.2 Einbindung PSM	27
6 Evaluation	30
6.1 Experiment 1: Büro	30
6.2 Experiment 2: Frühstück	34
6.3 Fazit	38

7 Zusammenfassung und Ausblick	39
Literaturverzeichnis	40

1. Einführung

In der Robotik ist die Servicerobotik wohl der Forschungsbereich, welcher den größten Alltagsbezug für den Menschen hat, da er sich mit der Entwicklung und Weiterentwicklung von autonomen Robotern beschäftigt, welche dem Menschen im Alltag assistieren. Man findet mittlerweile Roboter im Privaten, die das Putzen, Staubsaugen oder Rasenmähen übernehmen, in der Industrie, bei Montage und Fertigung, sowie auch in der Medizin, als Pflegehilfe, Botengänger oder Assistent.

Allerdings müssen die Roboter ihre Umwelt für komplexere Aufgaben so präzise wie möglich wahrnehmen und verstehen. Sie können Aufgaben übernehmen bei denen sie gezielt Objekte umfahren, suchen und auch aufnehmen und benutzen. Dieser Funktionsumfang kann mit dem Prinzip Programmieren durch Vormachen (PdV) ermöglicht werden, bei dem die Roboter Objekte und Tätigkeiten ihrer Umgebung kennen lernen, wieder erkennen und nachahmen können. So lässt sich die hohe Komplexität umgehen, die die manuelle Programmierung vieler Aufgaben mit sich bringen würde.

Um tatsächlich selbstständige Serviceroboter zu schaffen muss man aber noch zu einer Objekterkennung ein Kontextverständnis hinzufügen. Die Roboter müssen erkannte Objekte in einen Zusammenhang bringen, um die dadurch resultierenden Aufgaben zu verstehen. Zum Beispiel hat ein Teelöffel, welcher neben einer Tasse Tee liegt eine andere Aufgabe zu verrichten, als wenn er neben einem Becher Joghurt platziert ist. Nur am Kontext lässt sich dort entscheiden warum im einen Fall umgeführt und im anderen gelöffelt wird. Ebenso wäre ein Stück Butter verschieden zu verwenden, wenn es auf einem Frühstückstisch steht als wenn es mit anderen Zutaten neben einer Rührschüssel vorkommt.

Somit braucht man eine Szenenerkennung, welche zuverlässig die Objekte erkennen und ihren jeweiligen Kontext verstehen und einschätzen kann. Diese Erkennung ist nicht immer eindeutig, da der eben genannte Löffel ebenso zwischen einem Becher Joghurt und einer Tasse Tee liegen könnte, deshalb bietet es sich an mit Wahrscheinlichkeitsabschätzungen des vorliegenden Kontexts zu arbeiten.

2. Motivation und Problemstellung

Das Kapitel geht in Motivation auf die Relevanz des Themas und der vorliegenden Arbeit ein und in Fokus der Arbeit auf die Problemstellung die bearbeitet wird, sowie diverse Einschränkungen und Annahmen die für die Arbeit festgelegt sind.

2.1 Motivation

Wie schon in der Einführung erwähnt ist es elementar wichtig, dass man eine zuverlässige Szenenerkennung und ein gutes Kontextverständnis schafft um den Robotern die Möglichkeit zu bieten, sinnvoll mit ihrer Umwelt zu interagieren. Mit einer präzisen Wahrscheinlichkeitseinschätzung wie die momentane Umgebung beschaffen ist, lässt sich abschätzen welche Aufgaben es zu bewältigen und welche Probleme zu lösen gilt. Um dies zu gewährleisten muss man in das System möglichst viele Referenzdaten einspeisen, damit es jeden vorhanden Kontext erkennen kann. Szenen werden zu diesem Zweck aufgebaut un der Erkennung als neue Szene vorgestellt. Jede Szene die die Szenenerkennung auf diese Weise lernt, hilft das Chaos der sie umgebenen Objekte mehr und mehr zu interpretieren und einzuordnen.

Die Szenenerkennung nutzt also die Daten die sie zur Verfügung gestellt bekommt, bereits gelernte Szenen wiederzuerkennen. In dem bereits vorhandenen PSM(Probabilistic Scene Model)-System werden die Daten pro Szene zu einem Modell zusammengefasst, bei dem Auffällige Zusammenhänge berücksichtigt werden und scheinbar nicht miteinander in Verbindung stehende Objekte voneinander gelöst betrachtet werden. Zum Beispiel findet man eine Computermaus signifikant häufig vor dem Computerbildschirm und nie dahinter, allerdings kann dabei das räumliche Verhältnis der Maus zur Tastatur stark variieren. Die Maus wäre in diesem Beispiel manchmal direkt neben der Tastatur, manchmal dichter beim Bildschirm eben so oft weiter entfernt. In diesem Fall würde möglicherweise die Relation zwischen Maus und Tastatur wegfallen und nur jeweils das räumliche Verhältnis zum Bildschirm betrachtet werden. Dadurch gibt es Vorteile in der

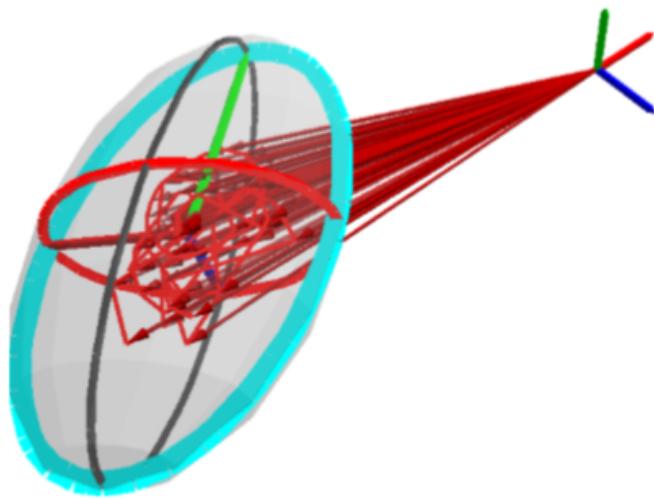


Abbildung 2.1: Beispiel: Relative Position eines Objekts zu einem anderen [Geh14]

Laufzeit und möglicherweise auch eine signifikantere Erkennung, allerdings findet natürlich auch ein Informationsverlust statt, der zu Fehlern führen kann. In dieser Arbeit wird ein Ansatz getestet, der dichter an den erhaltenen Daten arbeitet.

In Abbildung 2.1 sieht man wie das Modell eine Relation zwischen zwei Objekten aufgrund der erhaltenen Daten erstellt. Die roten Pfeile beschreiben hier die jeweiligen Relationen, die aus den Daten berechnet wurden und die Ellipsen zeigen, welchen Bereich das Modell als Basis für die Wahrscheinlichkeitsabschätzung benutzt.

2.2 Fokus der Arbeit

Ziel der vorliegenden Arbeit ist es das PSM-System zu überarbeiten und einen neuen Modus der Szenenerkennung im PSM-System zu entwickeln. Dieser Modus soll alternativ zu den bereits vorhandenen Modi auswählbar sein und das System erweitern. Außerdem soll sowohl die Positionierung als auch die Rotation der erkannten Objekten mit den bekannten Szenen verglichen werden um eine Wahrscheinlichkeitsabschätzung auszugeben, welche die Wahrscheinlichkeit aller möglichen Szenen ausgibt sowie auch die Wahrscheinlichkeit, dass es sich um keine der Szenen handelt.

Sei M das Szenenmodell, welches Szene S darstellt, D die Daten auf denen M basiert und sei O eine Menge an Objekten, die momentan von einer Objekterkennung erkannt werden. Dann sei $f(O, M, D)$ eine Funktion, die die Objekte, die Daten und das Modell annimmt und daraus die Wahrscheinlichkeit von Szene S berechnet, dass sie in O enthalten ist. Es soll also eine neue Funktion f entstehen, welche dicht an den Daten arbeitet, die von der zu testenden Szene S vorhanden sind.

Um eine interaktive Szenenerkennung zu schaffen ist es sinnvoll das System datengetrieben zu programmieren. Datengetriebene Entwicklung zeichnet sich dadurch aus, dass



Abbildung 2.2: Zwei Beispielszenen - Relative Lage der Objekte ändert die Szene [Geh14]

sich der Programmfluss ändert aufgrund der Daten die das System während der Laufzeit als Eingabe erhält. Das PSM-System wurde bereits datengetrieben programmiert und in der zu dieser Arbeit gehörigen Implementierung wurde auch am datengetriebenen Entwicklungsmodell festgehalten.

Da die Erkennung innerhalb des PSM-Systems nutzbar sein soll sind Eingabe- und Ausgabeschnittstellen sowie die Visualisierung vorgeschrieben und sollen für gute Vergleichbarkeit denen des vorhandenen Systems entsprechen.

Die Arbeit ist wie folgend strukturiert. In Kapitel 3 Grundlagen wird auf die Grundkenntnisse eingegangen die man braucht um den Rest der Arbeit gut zu verstehen. Es werden in diesem Kapitel die für die Arbeit wichtigen Komponenten des bestehenden PSM-System erklärt. In Kapitel 4 Konzept wird das theoretische Konzept thematisiert, welches zum Algorithmus geführt hat, dass den neuen Modus vom alten System unterscheidet. Außerdem wird der Algorithmus selbst erläutert. In Kapitel 5 Implementierung wird die Software beschrieben und dokumentiert die im Zuge dieser Arbeit entstanden ist sowie die Änderungen die am bestehenden PSM System vorgenommen wurden. Kapitel 6 Evaluation beschreibt die durchgeföhrten Experimente und interpretiert sie. In Kapitel 7 Zusammenfassung und Ausblick wird die Arbeit noch einmal zusammen gefasst und ein Ausblick darauf gegeben in welche Richtung man das bestehende System weiter entwickeln und auf welche Weise man möglicherweise die Szenenerkennung noch weiter verbessern kann. Am Schluss stehen die Quellen welche zur Recherche für diese Arbeit genutzt wurden.

3. Grundlagen

In diesem Kapitel wird das bestehende Probabilistic Scene Model - System vorgestellt und erklärt. Dabei wird auf das Grundprinzip, die Struktur und die einzelnen Komponenten eingegangen.

3.1 PSM - Probabilistic Scene Model

Bei der Entwicklung des Systems war das Ziel, ein System zur Erkennung von Szenen zu entwickeln. Dabei sollten viele Informationen umfasst werden. Dies umfasste die beteiligten Objekte bzw. deren Erscheinung, welche von den zur Erkennung eingesetzten Werkzeugen abhängig ist. Die Auftrittshäufigkeit der Objekte sollte ebenfalls mit einfließen. Der Hauptinformationsträger sind die Relationen der Objekte untereinander, welche in Form von relativen Objektlagen berücksichtigt wurden. Es wurde berücksichtigt, dass Relationen nicht statisch, sondern dynamisch sind. Da eine Szene auch unabhängig von ihrem Ort der Demonstration erkannt werden soll, musste die Lagebeschreibung invariant gegenüber Rotation und Translation sein.

Weiterhin sollte Robustheit gegenüber verschiedenen Störfaktoren bestehen. Fehlende Objekte sollen eine Erkennung der Szene erlauben, wenn auch mit entsprechend reduzierter Konfidenz. Da jedes Objekt fehlen kann darf es kein zentrales Referenzobjekt geben, von dem die Relationen zu den anderen Objekten der Szene ausgehen. Überzählige Objekte sollen sich nicht negativ auf das Erkennungsergebnis auswirken. Generell soll sich die Funktionsweise des entwickelten Systems im Rahmen dessen bewegen, was plausibel erscheint.

Es wird im folgenden erst auf die Struktur des Systems, dann auf die Eingabeverarbeitung, das innere Datenmodell und letztendlich auf die eigentliche Erkennung eingegangen. [Geh14]



(a) Der Teller im Kontext einer Frühstücksszene.



(b) Derselbe Teller im Rahmen eines Mittagessen-Szenario.

Abbildung 3.1: Zwei Beispielszenen - Ein Objekt in zwei Kontexten [Geh14]

3.1.1 Aufbau

Das System ist in mehrere Komponenten aufgeteilt die verschiedenen Aufgaben dienen. Für diese Arbeit wichtig sind vor allem die Komponenten Learner und Inference. Der Learner ist für das Anlernen der Daten und das Einspeichern von Szenen zuständig. Die Inference übernimmt die Erkennung der Szene indem sie in Echtzeit Daten empfängt und die Wahrscheinlichkeiten der Szenen ausgibt. Es gibt außerdem eine Visualizer-Komponente, die dafür zuständig ist, dass verschiedene programminterne Prozesse über Rviz für den Nutzer sichtbar und verständlich gemacht werden.

3.1.2 Learner - Anlernen der Daten

Die Learnerkomponente berechnet aus den Objekten, die sie bekommt, die Parameter für das Szenenmodell, dass für die Erkennung benötigt wird. Der Learner ist in mehrere Teilbereiche eingeteilt, Engine, OCM und Szenenmodell. Engine ist die Schnittstelle des Learners und kapselt die anderen beiden Architekturgliederungen voneinander ab.

Abbildung 3.2 beschreibt den Learner als Klassendiagramm. Die darin vorkommende *SceneLearningEngine* ist die Schnittstelle des Moduls zum Rest des PSM- Systems. Sie liest die in der Launch-Datei gegebenen Parameter aus und überprüft, ob alle Parameter sich mit dem passenden Datentyp auslesen lassen. Außerdem ist die Klasse für Visualisierung des Modells zuständig, welches gerade gelernt wurde.

Die gegebenen Objekte können mehrere Szenen beschreiben, da man auch Objekte aus einer Datei auslesen kann und diese jeweils eine pattern variable haben, die beschreibt zu welcher Szene sie zugehörig sind. Für jede Szene die gerade gelernt wird, gibt es einen Lerner in der *SceneModelLearner*-Klasse, welche dafür verantwortlich ist, die einzelnen Lerner zu verwalten. Die Daten der verschiedenen Szenen werden auf die dazugehörigen Lerner aufgeteilt und es gibt eine separaten Lerner für den Hintergrund, dessen Daten für den Fall wichtig sind, dass keine der Szenen in den gemessenen Objekten vertreten ist. Es ist sozusagen die Szene, die die Gegenwahrscheinlichkeit symbolisiert.

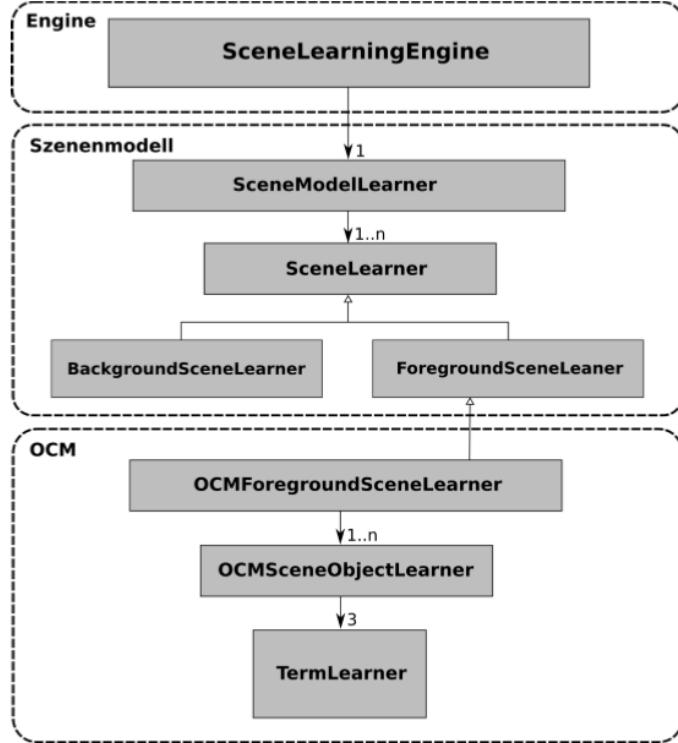


Abbildung 3.2: Klassendiagramm des Learner-Moduls [Geh14]

Alle Lerner erben vom *SceneLearner*, der eine abstrakte Basisklasse darstellt und eine Schnittstelle für das Lernen, Speichern und Visualisieren bietet. Der Lerner für den Hintergrund ist eine Instanz des *BackgroundSceneLearner*, welcher blos die Anzahl der Objekte abspeichert, um daraus einen validen Schwellenwert beziehungsweise eine sinnvolle Gegenwahrscheinlichkeit zu bestimmen. Der *ForegroundSceneLearner* ist eine abstrakte Klasse, sodass man das Modell der Berechnung leicht austauschen kann, allerdings konzentriert sich das PSM-System auf den Einsatz des OCM als Repräsentation der Szenenobjekte. Der *OCMForegroundSceneLearner* ist eine Unterklasse des *ForegroundSceneLearner* und kapselt die Lerner für den Vordergrund, welche Instanzen der Klasse *OCMSceneObjectLearner* sind. [Geh14]

3.1.3 Model - Szenenmodell

Das Szenenmodell wird als XML-Datei abgespeichert. Dadurch kann man manuell das erstellte Modell lesen, verstehen und verändern, falls dies zum testen nötig ist. Man kann auch leicht mehrere erstellte Szenenmodelle kombinieren, da man alle Parameter verändern und leicht eine zusätzliche Szene aus einem anderen Modell hinzufügen oder ersetzen kann.

Abbildung 3.3 zeigt ein Beispielmodell für eine Frühstücksszene namens *breakfast* und die immer vorhandene Hintergrundszene *background*. Jeder Szene ist ein apriori-Wert

```

<psm>
  <scenes>
    <scene type="background" name="background" priori="0.5">
      <description objects="2" volume="27"/>
    </scene>
    <scene type="ocm" name="breakfast" priori="0.5">
      <object name="CoffeeBox" type="ocm" priori="0.5">
        <slots number="2"/>
        <shape>
          <root volume="27">
            <child name="Cup">
              +<pose></pose>
            </child>
          </root>
        </shape>
      </object>
    </scene>
    <appearance>
      <mapping>
        <map id="1" name="CoffeeBox"/>
        <map id="2" name="Cup"/>
      </mapping>
    </appearance>
    <occlusion>
      <table>
        <entry values="0 1 0"/>
        <entry values="0 0 1"/>
      </table>
    </occlusion>
  </scenes>
</psm>

```

Abbildung 3.3: Beispiel eines Szenenmodells - Frühstücksszene [Geh14]

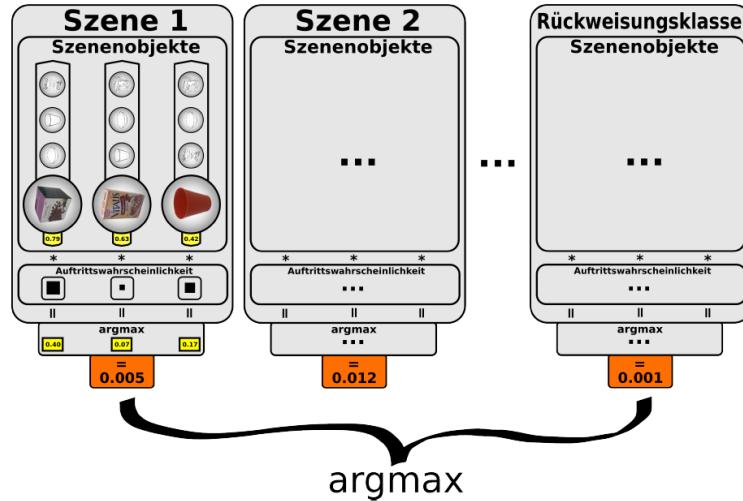


Abbildung 3.4: Grafische Darstellung mehrerer Szenenmodelle. Alle werden zur Wahrscheinlichkeitsberechnung genutzt. [Geh14]

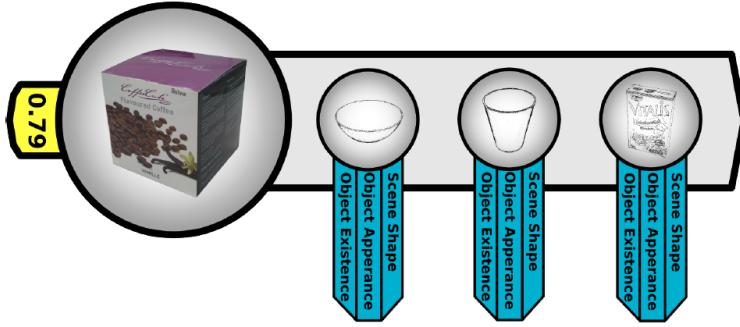


Abbildung 3.5: Das Object Constellation Model(OCM) beschreibt ein einzelnes Szenenobjekt, also ein Objekt (hier die Kaffeebox) im Kontext einer Szene. Die Slots werden durch die kleinen Kreise rechts symbolisiert und beschreiben jeweils ein Objekt durch die blau eingefärbten Parameter. In gelb sieht man die Grunfwahrscheinlichkeit des Objekts. [Geh14]

zugeordnet, welcher beschreibt, wie wahrscheinlich es grundsätzlich ist, dass die Szene auftritt, allerdings ist diese im PSM-System generisch gleichverteilt und kommt nur der Vollständigkeit halber vor. Man kann diese Werte falls nötig manuell ändern, das erstellte Szenenmodell der Learner-Komponente wird aber stets gleichverteilte Werte beinhalten. Der Teil der Hintergrundszene im Modell beinhaltet Informationen über die Anzahl unterschiedlicher Objekte und die Größe des Bereichs auf dem Szenen erkannt werden.

Die Szene namens *breakfast* beinhaltet zwei Szenenobjekte, *Cup* und *CoffeeBox*. Als Beispiel wurde eine einfache Szene mit geringer Objektanzahl gewählt um die Erklärung möglichst simpel zu halten. Außerdem ist nur das Szenenobjekt *Coffebox* ausgeklappt in der Modellabbildung, damit die Abbildung übersichtlicher ist. Für jeden Term des OCM sind die entsprechenden Parameter im Objekt gespeichert. Jedes Objekt in der Szene hat eine apriori-Wahrscheinlichkeit, welche die Wichtigkeit des Objekts in der Szene verdeutlicht. Desweiteren enthält es Informationen über die Anzahl der Slots, welche zwar aus dem kompletten Modell hergeleitet werden kann, man spart aber Ladezeit dadurch. Die Terme die im OCM vorkommen sind *Scene Shape* und *Objekt Appearance*.

Der *Scene Shape*- Term wird im Modell mit *shape* bezeichnet und beinhaltet einen Baum, welcher die Relation zwischen den Szenenobjekten beschreibt. Der Baum besteht aus einem Wurzelknoten *root* und beliebig vielen Kindsknoten *child*, welche wiederum Kindsknoten beinhalten können. Der Wurzelknoten enthält das Volumen *volume* auf dem die Erkennung arbeiten soll und die Kindsknoten jeweils Daten die die relative Lage zum direkten Elternknoten bezeichnen. Diese Daten, die unter der Bezeichnung *pose* gespeichert sind, sind eingeklappt, da es unübersichtlich wäre, und enthalten Mittelwerte, Gewichte und Kovarianzmatrizen der einzelnen Gauss-Kernel.

Der *Objekt Appearance*-Term ist in zwei Abschnitte aufgeteilt. Erst wird jedem Objektnamen ein Index zugeordnet. Anschließend gibt es eine Wahrscheinlichkeitstabelle, die für jeden Index eine Auftrittswahrscheinlichkeit definiert. Der Index 0 wird bei der Zu-

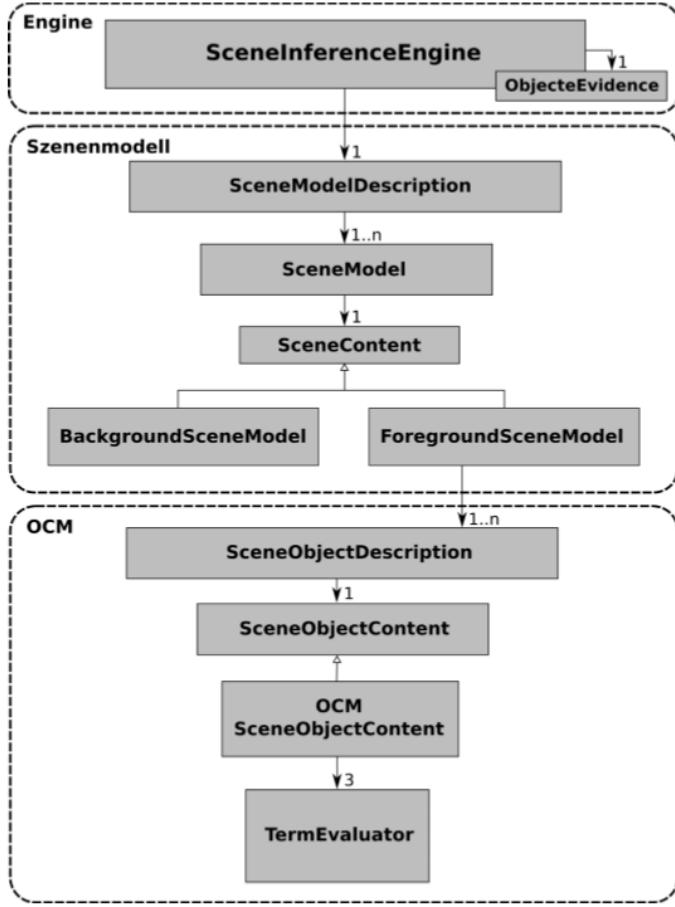


Abbildung 3.6: Klassendiagramm des Inference-Moduls [Geh14]

ordnung bewusst nicht belegt, da er ein Platzhalter für unbekannte Objekte ist. [Geh14]

3.1.4 Inference - Szenenerkennung

Die Szenenerkennung oder auch Inferenz hat zur Aufgabe, aufgrund von dem gegebenen Szenenmodell und momentan erkannten Objekten, im folgenden auch Evidenz genannt, für jede Szene im Modell einzuschätzen, wie wahrscheinlich es ist, dass diese in der Evidenz vorkommt. Die verschiedenen Szenenwahrscheinlichkeiten werden dann miteinander verrechnet und es wird die relative Wahrscheinlichkeit für jede Szene ausgegeben, dass diese vorhanden ist. Diese Erkennung passiert in Echtzeit und ist als datengetriebenes Modul zu betiteln, da sich die Ausgaben und das Verhalten des Programms, maßgeblich dadurch verändert, welche Daten die Evidenz enthält, das heißt, welche Objekte erkannt werden.

In Abbildung 3.6 sieht man das Klassendiagramm der Szenenerkennung, welches vergleichbar zu dem des Learner-Moduls strukturiert ist. Die verschiedenen Teilbereiche

sind wieder Engine, OCM und Szenenmodell. Engine ist die Schnittstelle und kapselt Szenenmodell und OCM voneinander ab. *SceneInferenceEngine* ist für das Auslesen der Parameter aus der Launch-Datei, für die Visualisierung des Erkenners und die Annahme und Weitergabe der erkannten Evidenz zuständig. Diese wird an *SceneModelDescription* weiter gegeben, welche für jede Szene im Modell eine Instanz der Klasse *SceneDescription* erstellt. Ob es sich um eine Vorder- oder Hintergrundszene handelt, kann man dadurch erkennen, die daraufhin erstellte Instanz einer Unterkategorie der Abstrakten Klasse *SceneContent* entweder vom Klassentyp *BackgroundSceneModel* oder *ForegroundSceneModel* ist.

In *ForeGroundSceneModel* sind Instanzen von *SceneObjektDescription* gekapselt, allerdings wird im PSM-System nur eine einzige Instanz benutzt, da dieses nur den Ansatz des OCM verfolgt und diese Kapselung dafür gedacht ist verschiedene Ansätze möglich zu machen. Die Unterkategorie *OCMSceneObjektContent* ist der besagte Algorithmus des PSM-Systems und beinhaltet mehrere *TermEvaluator* die für die *Object Appearance*, *Object Existance* und *SceneShape* zuständig sind.

Die jeweils berechneten Wahrscheinlichkeiten werden eingespeichert und von der *SceneInferenceEngine* über mehrere Klassen hinweg abgefragt und über das Visualisierungsmodul ausgegeben. [Geh14]

4. Konzept

Das folgende Kapitel thematisiert das Konzept, dass im Zuge der vorliegenden Arbeit entwickelt wurde, um die vorgestellte Problemstellung zu lösen. Zwischen vielen verschiedenen möglichen Ansätzen einen neuen Modus für das PSM System zu entwickeln, entschied ich mich für einen differenzbasierten Vergleich der Positions- und Rotationsrelationen. Dieser wird im ersten Abschnitt erläutert und anschließend wird der Algorithmus sprachlich, grafisch und in Pseudocode erklärt. Zum Schluss wird die Wahrscheinlichkeitsabschätzung für den Algorithmus erklärt und begründet.

4.1 Ansatz

Im vorhandenen PSM-System werden im Learner die erhaltenen Positions- und Rotationsdaten zu einem Modell zusammengefasst, dass die vorkommen aller Objekte in Relation zueinander zusammenfasst. Dieser Vorgang führt dazu, dass teilweise Zusammenhänge in den Daten betont werden, aber auch zu einem Informationsverlust da Ausreißer und mutmaßliche Fehlmessungen dadurch verloren gehen. Deshalb kann es sinnvoll sein direkt auf den gemessenen Daten zu arbeiten, um ein Ergebnis zu erhalten welches alle Daten berücksichtigt.

Wir betrachten folgendes Szenario. Ein Roboter soll seine Aufgaben aufgrund von einer Szenenerkennung einschätzen und durchführen. In der Szene "Kaffee" gibt es eine volle Kaffeetasse und einen Teelöffel und seine Aufgabe ist es mit dem Löffel den Kaffee umzurühren. In seinen Referenzdaten zu der Szene war der Löffel meist direkt neben der Tasse und nur in einem Fall ein Stück weiter entfernt. Allerdings gilt jede einzelne aufgezeichnete Referenzszene auf äquivalente Weise als Beispiel für die Szene "Kaffee". Das Parametermodell würde diese Ausreißerdaten allerdings glätten und kaum berücksichtigen, sodass der Roboter die Szene selbst mit genau dem Aufbau aus den Referenzdaten möglicherweise nicht erkennen würde. Wenn man allerdings die Erkennung direkt auf den Referenzdaten basiert, erkennt die Szenenerkennung den Ausreißer auch, da sie ja eine Instanz der Szene mit diesem vergleicht, welche diesem entspricht.



Abbildung 4.1: Beispiel: Kaffeetasse - Ausreißer in den Daten

Abbildung 4.1 verdeutlicht das genannte Szenario. Die roten Punkte stehen für die gemessenen Positionen und die grünen Pfeile stehen für die räumliche Relation zwischen den Objekten. Links sieht man, dass die meisten Messungen einen kleinen Abstand zwischen Löffel und Tasse haben, rechts ist der Ausreißer dargestellt, der möglicherweise vom alten System nicht als die gelernte Szene erkannt wird.

Im differenzbasierten Modus sollen also gemessene Objekte direkt mit den Referenzdaten verglichen werden, die das System bereits gelernt hat. Der Algorithmus betrachtet alle Objekte vollvermascht, sodass er die Szenenreferenz findet, die die maximale Ähnlichkeit zu den gemessenen Objekten hat. Darauf basierend wird die Wahrscheinlichkeit abgeschätzt, dass die gemessenen Objekte die Referenzszene enthalten oder repräsentieren. Dabei stören zusätzliche Objekte die Erkennung nicht und eine Unvollständigkeit der Szene führt zu einer kleineren Wahrscheinlichkeit aber nicht zu direkter Ablehnung, da die Szene noch durch weitere Objekterkennungen vervollständigt werden könnte.

4.2 Erkennungsalgorithmus

Der Algorithmus wurde mit den in Ansatz genannten Annahmen und Einschränkungen entwickelt und hat zur Aufgabe zu jeder Szene die auf Vorkommen geprüft wird die Instanz der Szene in den Daten zu finden, die am dichtesten an den gemessenen Daten liegt und so die höchste Wahrscheinlichkeit aufzeigt, dass die gemessenen Daten die Szene enthalten. Nachdem die Wahrscheinlichkeit einer Szene bestimmt ist wird diese mit den anderen Szenen genau wie im bestehenden PSM-System verrechnet, sodass am Ende die Relative Wahrscheinlichkeit für alle zutestenden Szenen angegeben wird.

4.2.1 Algorithmus: Beschreibung

Der Algorithmus läuft wie folgt ab. Für jedes Objekt, der zu testenden Szene, wird überprüft ob es sich um ein Objekt handelt, welches gerade von der Objekterkennung erkannt wird. Wenn dies nicht der Fall ist, wird das Objekt übersprungen. Wenn dies allerdings

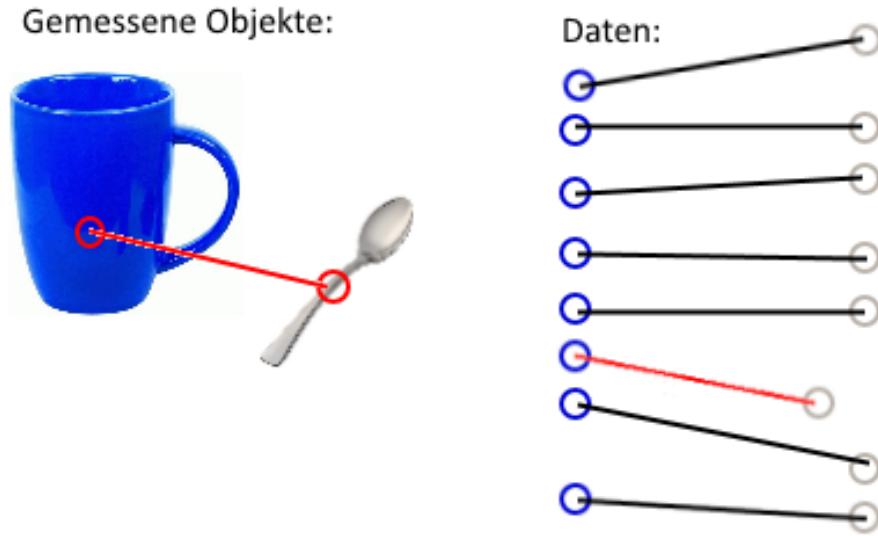


Abbildung 4.2: Beispiel: Erkennung einer Tasse und dem dazugehörigem Teelöffel. Der Algorithmus sucht die Dateninstanz, die am besten passt. Diese sind hier in rot markiert.

der Fall ist, wird das Objekt zeitweise zu unserem Referenzobjekt und der Algorithmus führt für jede Instanz der zu testenden Szene in den Daten folgendes aus:

Es wird über alle anderen Objekte der Szeneninstanz iteriert und für jedes Objekt abgefragt, ob es ein gemessenes Objekt gibt, welches die Repräsentation für das Datenobjekt sein kann. Falls das Objekt nicht in den momentan wahrgenommenen Objekten vor kommt, ist die Szene allein aus Sicht dieses Objekts betrachtet unwahrscheinlich. Deshalb wird eine Objektwahrscheinlichkeit von 0 eingespeichert, welche aussagt, dass dieses Objekt allein die Szene als nicht auffindbar beschreibt. Wenn allerdings eine Instanz der Objekts in der Iteration gefunden wird, berechnet der Algorithmus die Positions- und Rotationsrelation zwischen den beiden Objekten aus den Daten. Genauso wird die Position und Rotationsrelation der gemessenen Objekte zueinander berechnet, welche mutmaßlich den Objekten aus der Datenbank entsprechen sollten. Nach diesen Berechnungen wird die Differenz verglichen die zwischen den beiden Paaren jeweils besteht. Dabei werden Positionsrelationen und Orientierung beziehungsweise Rotation jeweils getrennt betrachtet. Aus dem Grad der Ähnlichkeit dieser Differenzen lässt sich nun die Wahrscheinlichkeit ableiten, ob das Objekt so vorkommt wie die Szene aufgezeichnet wurde. Somit hat man eine Objektwahrscheinlichkeit.

Nun werden alle Objektwahrscheinlichkeiten zu einer Szenewahrscheinlichkeit zusammengefasst. Diese wird innerhalb der Iteration über die Szeneninstanzen maximiert. Außerdem wird dieser Maximalwert wiederum innerhalb der äußersten Schleife maximiert, welche über alle Objekte iteriert, die sowohl in der zu testenden Szene sind als auch von der Objekterkennung erkannt wurden. Es wird also der absolute maximale Wert der

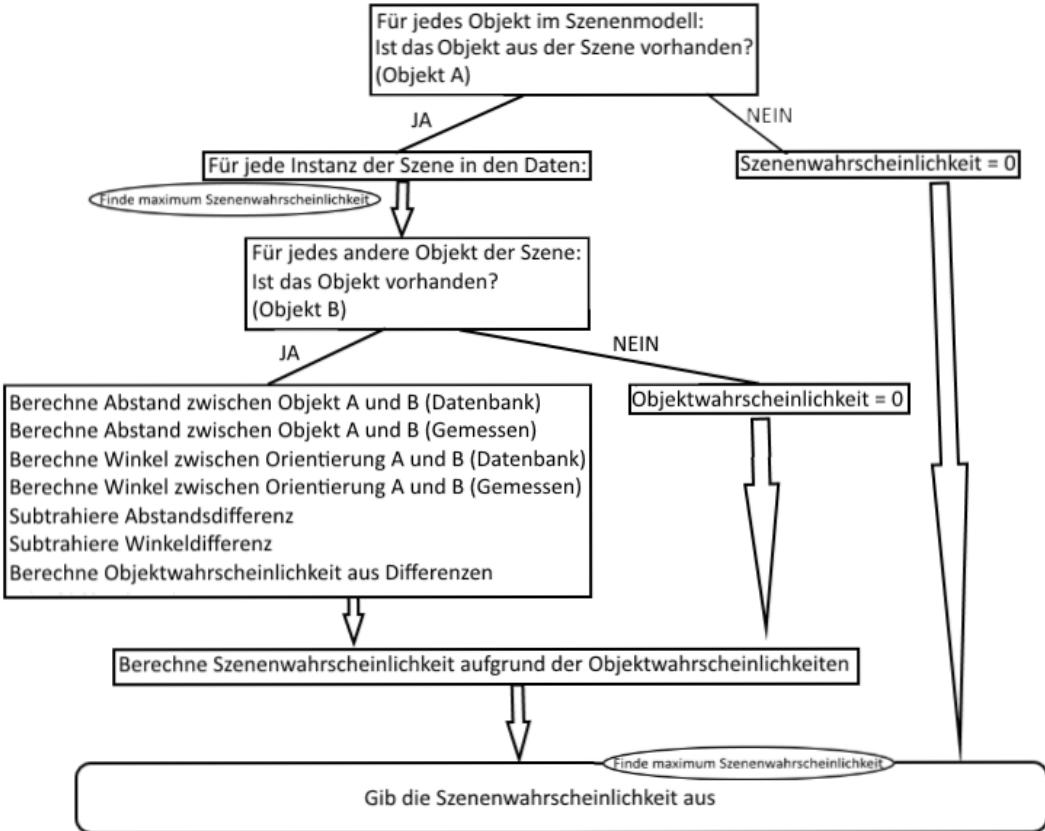


Abbildung 4.3: Algorithmus als vereinfachtes Flussdiagramm

Szenenwahrscheinlichkeit bestimmt, den man mit einer der Szeneninstanzen mit der beschriebenen Schleife mit jedwedem Referenzobjekt erzeugen kann. Dieser Maximalwert ist die Wahrscheinlichkeit, ob die zu testende Szene in den Gemessenen Objekten und potentiell weiteren unbekannten Objekten enthalten ist, welche von der Objekterkennung noch erkannt werden könnten.

Abbildung 4.3 beschreibt den Algorithmus als vereinfachtes Flussdiagramm. Einfache Linien verdeutlichen den Programmfluss in den verschiedenen Fällen und die Flussrichtung ist stets nach unten. Die Schleifen sind zusammengefasst damit das Diagramm übersichtlich bleibt.

4.2.2 Algorithmus: Pseudocode

Um das Algorithmuskonzept weiter zu verdeutlichen beschreibt Abbildung 4.4 den Algorithmus nochmals mit Pseudocode. Damit der Code verständlich wird hier eine kurze Erklärung zu der Abbildung. Die Einrückungen wurden statt den geschweiften Klammern verwendet, die in den meisten höheren Programmiersprachen vorkommen. Die Parameter sind wie folgt definiert. Szenenmodell beschreibt das Szenenmodell, dass die vorkommenden Objekte in der zu testenden Szene beinhaltet. Der Parameter Gemessen

```

BerechneSzenenWahrscheinlichkeit ( Objekt[] Szenenmodell, Objekt[] Gemessen, Szeneninstanz[] Daten )
|   SzenenWahrscheinlichkeit = 0;

|   foreach - Objekt A in Szenenmodell
|
|       if - Gemessen->EnthaeltObjektMitName ( A->Name )
|           ReferenzObjekt = Gemessen->FindeObjektMitName ( A->Name );

|           foreach - Szeneninstanz S in Daten
|               SzenenWahrscheinlichkeitsSumme = 1;
|
|               foreach - Objekt B in S
|                   Objektwahrscheinlichkeit = 0;
|
|                   if - Gemessen->EnthaeltObjektMitName ( B->Name )
|                       ZweitesObjekt = Gemessen->FindeObjektMitName ( B->Name );
|                       Objektwahrscheinlichkeit = BerechneObjektWahrscheinlichkeit( ReferenzObjekt, ZweitesObjekt, A, B);
|
|                   SzenenWahrscheinlichkeitsSumme += Objektwahrscheinlichkeit;

|               NeueSzenenWahrscheinlichkeit = SzenenWahrscheinlichkeitsSumme / Szenenmodell->Laenge;

|               if - NeueSzenenWahrscheinlichkeit > SzenenWahrscheinlichkeit
|                   SzenenWahrscheinlichkeit = NeueSzenenWahrscheinlichkeit;

```

Abbildung 4.4: Algorithmus als PseudoCode

beschreibt die Objekte die momentan erkannt werden und real oder in einer Simulation vorhanden sind. Der Parameter Daten steht für die Instanzen die zur zu testenden Szene als Referenzen gespeichert sind.

Die Variable Name, die jedes Objekt gesetzt hat ist eine Identifizierung um Objekte ihren mutmaßlichen Vorkommen in der Messung zuzuordnen. Die Funktion Enthaelt-ObjektMitName(string Name) gibt true aus, falls die aufrufende Liste von Objekten ein Objekt mit dem gegebenen Namen enthält. Ansonsten wird false ausgegeben.

Die Funktion FindeObjektMitName(string Name) gibt das Objekt aus der Liste zurück, welches den gegebenen Namen trägt. Falls kein Objekt mit dem gegebenen Namen existiert wird NULL zurückgegeben. BerechneObjektWahrscheinlichkeit(Objekt C, Objekt D, Objekt A, Objekt B) nimmt vier Objekte und berechnet die Positions- und Orientierungsunterschiede zwischen den Parametern C und D sowie zwischen A und B. Anschließend werden die Differenzen abgeglichen und basierend auf den Unterschieden eine Wahrscheinlichkeitsabschätzung zwischen 0 und 1 abgegeben, wobei 1 für "mit der Szene übereinstimmend" und 0 für "weit entfernt" steht. Auch die Komplette Funktion speichert am Ende eine Wahrscheinlichkeit zwischen 0 und 1. Sie wird nicht ausgegeben, da auf die eingespeicherte Wahrscheinlichkeit über eine andere Schnittstelle zugegriffen wird und der Algorithmus nur den Zweck erfüllt die Wahrscheinlichkeit zu berechnen.

4.3 Wahrscheinlichkeitsabschätzung

Zuerst zeigen wir, dass die Abschätzung beim einzelnen Objekt präzise ist, anfolgend wird die Berechnung der gesamten Szenenwahrscheinlichkeit erklärt. Die Wahrschein-

lichkeiten sind von gewissen Schwellenwerten abhängig, da der Maßstab und die Szene die überprüft wird unterschiedliche Anforderungen haben kann.

Bei einer Beispielszene, die den gedeckten Frühstückstisch repräsentiert, hat man sicher noch eine gewisse Tolleranz, wenn es um die räumliche Positionierung der Objekte zueinander geht, allerdings gibt es kaum Kompromisse bei der Rotation. Wenn die Teller umgekehrt wären, würde es nicht mehr dem gewohnten gedeckten Tisch entsprechen. Wenn hingegen ein Ball in einer Szene vorkommt, wird dieser mehr Tolleranz gegenüber Rotation aber möglicherweise einen kleineren Schwellenwert bei der Position haben. Man sieht also, dass die Schwellenwerte nötig sind, um in verschiedenen Szenekontexten sinnvolle Ergebnisse zu bekommen.

4.3.1 Objektwahrscheinlichkeit

Es gibt eine Instanz der zu testenden Szene. Außerdem gibt es ein Referenzobjekt mit Positions- und Rotationsdaten. Es gibt ein Testobjekt mit Positions- und Rotationsdaten. Es gibt ein gemessenes Referenzobjekt, welches dem gewählten Referenzobjekt in der aktuellen Objekterkennung entspricht. Nun wird das Objekt in den erkannten Objekten gesucht, welches die selbe Identifikation, wie unser Testobjekt hat. Dieses nennen wir gemessenes Testobjekt. Falls wir kein Objekt finden, welches die gewollten Eigenschaften hat ist die Objektwahrscheinlichkeit gleich null, da die Szene aus der Perspektive des Testobjekts nicht auffindbar ist.

Die Objektwahrscheinlichkeit beschreibt die Wahrscheinlichkeit für das Testobjekt, dass die Szene aus der es entspringt in den aktuell erkannten Objekten vorkommt. Diese wird bestimmt indem das Verhältnis zwischen Testobjekt und Referenzobjekt mit dem zwischen gemessenem Testobjekt und gemessenem Referenzobjekt abgeglichen wird. Dieser Vergleich vergleicht einerseits die Positionsrelationen sowie auch den Winkelunterschied zwischen den Orientierungen der Objekte. Die beiden Einzelvergleiche arbeiten mit Schwellenparametern, die bestimmen wie groß der Unterschied sein muss, damit die Wahrscheinlichkeit null entspricht. Beim Winkel kann dieser Wert zwischen 0 und 180 liegen, beim Positionsvergleich ist dieser Parameter eine beliebige Zahl größer oder gleich null. Ausformuliert bedeuten diese Parameter: Wie weit kann man die Szene verändern, bevor die Erkennung sie nicht mehr erkennen soll?

In Abbildung 4.5 kommen Testobjekt T, Referenzobjekt R, gemessenes Testobjekt gT, gemessenes Referenzobjekt gR und die Parameter, Positionsparameter pP und Rotationsparameter rP, vor. Alle Objekte haben geben mit `.p` ihre Position als Vektor und mit `.r` ihre Rotation als Matrix in einem Weltkoordinatensystem zurück. Die Funktion `RotationDiff` nimmt zwei Rotationsmatrizen und berechnet den Rotationsunterschied zwischen diesen. `Position(T, R, gT, gR)` berechnet den Differenzvektor zwischen den Relationen von R zu T und von gR zu gT und bestimmt anschließend dessen Länge. Die Funktion `Rotation(T, R, gT, gR)` berechnet die Rotationsmatrix, um die man Rotieren müsste um aus der Orientierungsrelation von R und T zu der vergleichbaren von gR und gT zu rotieren und gibt den Winkel dieser Rotation aus. Die Funktionen `PositionParametisiert(T, R, gT, gR, pP)` und `RotationParametisiert(T, R, gT, gR, rP)` benutzen

$$\begin{aligned}
 \text{Position}(T, R, gT, gR) &= (T.p - R.p - (gT.p - gR.p)).\text{Länge} \\
 \text{Rotation}(T, R, gT, gR) &= (\text{RotationDiff}(\text{RotationDiff}(R.r, T.r), \text{RotationDiff}(gR.r, gT.r))).\text{Winkel} \\
 \text{PositionParametisiert}(T, R, gT, gR, pP) &= \begin{cases} 0 & \text{Position}(T, R, gT, gR) > pP \\ 1 & pP = 0 \\ (pP - \text{Position}(T, R, gT, gR)) / pP & \text{sonst} \end{cases} \\
 \text{RotationParametisiert}(T, R, gT, gR, rP) &= \begin{cases} 0 & \text{Rotation}(T, R, gT, gR) > rP \\ 1 & rP = 0 \\ (rP - \text{Rotation}(T, R, gT, gR)) / rP & \text{sonst} \end{cases} \\
 \text{Objektwahrscheinlichkeit}(T, R, gT, gR, pP, rP) &= \text{PositionParametisiert}(T, R, gT, gR, pP) * \text{RotationParametisiert}(T, R, gT, gR, rP)
 \end{aligned}$$

Abbildung 4.5: Formeln zur Objektwahrscheinlichkeit

die Schwellenparameter pP und rP , um aus der Distanz oder dem Winkel die Wahrscheinlichkeitsabschätzungen zu gewinnen. Objektwahrscheinlichkeit(T, R, gT, gR, pP, rP) kombiniert die Wahrscheinlichkeit, die rein durch die Positionsrelationen begründet ist, sowie die Wahrscheinlichkeit, die nur aus der Rotation basiert, zu einer gesammten Objektwahrscheinlichkeit.

Im gesamten handelt es sich um eine Wahrscheinlichkeitsabschätzung die sich aus den Komponenten der Positions- und der Rotationswahrscheinlichkeit zusammensetzt. Die beiden einzelnen Komponenten beschreiben jeweils wie wahrscheinlich ist, dass aufgrund der gegebenen Objekte die zu testende Szene in der Objekterkennung repräsentiert wird, wenn man nur die jeweilige Komponente betrachtet.

4.3.2 Szenenwahrscheinlichkeit

Die Szenenwahrscheinlichkeit setzt sich aus den Objektwahrscheinlichkeiten der Szene zusammen. Mit jedem möglichen Referenzobjekt wird in jeder Instanz der Szene in den Daten die Objektwahrscheinlichkeit für jedes andere Objekt berechnet und anschließend in eine mögliche Szenenwahrscheinlichkeit zusammengefasst. Diese möglichen Szenenwahrscheinlichkeiten bestimmen, wie wahrscheinlich es ist, dass die Szene in den gemessenen Objekten repräsentiert ist, wenn man nur genau die gewählte Instanz der Daten und das gewählte Referenzobjekt betrachtet. Aus den möglichen Szenenwahrscheinlichkeiten wird die maximale Szenenwahrscheinlichkeit als das entgültige Ergebnis ausgewählt. Man findet also die Wahrscheinlichkeit dafür, dass die Szene vorhanden ist, ausgehend von den Referenzdaten, die am dichtesten an den gemessenen Daten sind.

Die Objektwahrscheinlichkeiten werden jeweils mit einem A-priori Wert multipliziert und ihre Wichtigkeit in der Szene zu gewichten. Objekte müssen nicht in jeder Datenaufzeichnung der Szene vorhanden sein und der A-priori Wert zeigt auf wie signifikant und wichtig das Auftreten eines Objekts in einer Szene ist. Zum Beispiel ist eine Gabel auf fast jedem gedeckten Mittagstisch, tiefe oder flache Teller werden sich jedoch teils in den Daten abwechseln und nicht beide in jeder Aufzeichnung in der Szene vorhanden sein. Somit hätte die Gabel einen höheren A-priori Wert.

Die gewichteten Objektwahrscheinlichkeiten werden nun aufaddiert. Jede Objektwahrscheinlichkeit hat einen Wert, der größer oder gleich 0 und kleiner oder gleich 1 beträgt. Außerdem wird noch 1 aufaddiert, da es ja ein Referenzobjekt gibt, bei dem man davon

```

Szene S
Instanz I
Referenzobjekt R
TestObjekt T
O(R, T, pP, rP) = Objektwahrscheinlichkeit(T, R, FindeGemessenes(T), FindeGemessenes(R), pP, pR)
Szenenwahrscheinlichkeit(S, pP, rP) = MaxForEach(S,MaxForEach(I, 1 * apriori(R) + SumForEach(I, O(R, T, pP, rP) * apriori(T))))

```

Abbildung 4.6: Formeln zur Szenenwahrscheinlichkeit

ausgeht, dass es genau an der richtigen Position ist und die gemessene Orientierung hat. Anschließend teilt man die Summe durch die Anzahl aller Objekte, um den durchschnittlichen Wert der Objektwahrscheinlichkeit zu erhalten. Dieser entspricht dem gewichteten Erwartungswert der Objektwahrscheinlichkeit von einem Objekt, dass man zufällig aus der zu testenden Instanz der Szene zieht. Dieser gewichtete Erwartungswert entspricht der Wahrscheinlichkeit, dass die Szene auf Basis der gemessenen Objekte vorhanden sein könnte. Parametisiert muss dieser Wert der Szenenwahrscheinlichkeit nicht mehr werden, da dies schon auf der Ebene der Objektwahrscheinlichkeit passiert.

In Abbildung 4.6 zeigt auf wie sich die Szenenwahrscheinlichkeit ergibt. In der Grafik die Formeln beinhaltet steht S für die zu testende Szene, I für eine Instanz aus den Daten, die die Szene beschreibt, R für eine Referenzobjekt welches variabel wählbar ist und T für das Testobjekt für das gerade die Objektwahrscheinlichkeit berechnet wird. Die Funktion $O(R, T, pP, rP)$ berechnet die Objektwahrscheinlichkeit für das gegebene Testobjekt T in Hinblick auf das Referenzobjekt R und mit den Parametern pP und rP, die die Distanztoleranz für die Position und die Winkeltoleranz für die Orientierung angeben. $\text{FindeGemessenes}(A)$ bestimmt jeweils das Objekt in den gemessenen Daten, welches dem gegebenen Objekt A entspricht, das heißt die selbe Identifikation wie A hat. $\text{Szenenwahrscheinlichkeit}(S, pP, rP)$ nimmt eine zu testende Szene S an und die beiden Schwellenparameter.

Die Funktion MaxForEach iteriert über alle gegebenen Instanzen von S und findet den Maximalwert, welcher beim zweiten Parameter errechnet werden kann. Dieser Parameter, der eine Funktion enthält, die einen Zahlenwert zurückliefert, verwendet I als die aktuelle Instanz, die gerade in der Iteration geprüft wird. MaxForEach kann aber auch eine Instanz als ersten Parameter enthalten und verhält sich dann anders. In diesem Fall wird über alle Objekte in der Szeneninstanz iteriert und jeweils das aktuelle Objekt als Referenzobjekt R ausgewählt. Der zweite Parameter muss eine Zahl zurückliefern und R kann frei in dem Ausdruck verwendet werden. Die Funktion gibt den Maximalwert aus, den sie über alle Schleifendurchläufe findet.

$\text{apriori}(A)$ gibt den gespeicherten A-priori Wert des Objekts A zurück, sodass man die Wahrscheinlichkeiten gewichten kann. Die SumForEach iteriert über alle Objekte der gegebenen Instanz I bis auf das momentane Referezobjekt R. Am Ende jedes Schleifenschritts wird der im zweiten Parameter definierte Zahlenwert aufaddiert. Die Variable T entspricht hier dem Testobjekt, welches das Iterationsobjekt der Schleife ist, also das Objekt welches aktuell in der Iteration geprüft wird.

Es wird also die Instanz in den Daten gesucht deren Objekte die höchsten Erwartungswert haben, wenn es darum geht aufgrund der Objekte und der gemessenen Daten einzuschätzen, ob die zu testende Szene in der Messung repräsentiert wird. Es ist sinnig den Maximalwert zu suchen, da jede Messung gleichermaßen vollwertiges Beispiel für die Szene gilt und damit die maximale Wahrscheinlichkeit bestimmt wird, dass die Szene, die es zu testen gilt, vorhanden ist. Man könnte natürlich auch wiederum den Erwartungswert bestimmen, wenn man eine zufällige Instanz aus allen Szenen zieht, um aus allen Szeneninstanzen ein Modell zu generieren, welches die Szene gut repräsentiert, da aber gegeben ist, dass jede einzelne Messung schon alleinstehend die Szene vollends repräsentiert, würde dieses Vorgehen die Erkennung nur unpräziser machen.

5. Implementierung

Im Kapitel Implementierung wird alles beschrieben und erklärt, was am bestehenden PSM-Projekt verändert und hinzugefügt wurde. Außerdem werden ausgewählte hinzugefügte und veränderte Klassen sowie launch-Dateien dokumentiert, sodass das Kapitel das Verständnis und die Nutzung der Neuheiten im System vereinfacht. Nachdem der Umbauprozess beschrieben wird, bei dem eine Klasse komplett aus dem PSM-Projekt ausgetauscht wurde, widmet sich das Kapitel der Umsetzung des Algorithmuskonzepts und der Einbettung in das vorhandene System.

5.1 Umbau PSM

Da das Paket *pbd_msgs* nicht konstenlos zur Verfügung gestellt wird, mussten alle Vorkommen der Klassen aus diesem Paket zu alternativen Ersatzklassen geändert werden. Teilweise konnte man dies durch simple Ersetzung erreichen, allerdings gab es nicht für jede Klasse eine Ersatzklasse mit dem selben Funktionsumfang. In den Fällen, in denen Funktionen fehlten oder geringfügig anders funktionierten, konnte man den Umbau durch kleine Anpassungen erreichen oder musste eigene Funktionen schreiben, welche die nötigen Operationen verrichten konnten. Alle auf diese Weise programmierten Funktionen wurden hinreichend auf Gleichheit mit ihren Ursprungsfunktionen in ihrer Funktionsweise getestet, indem die Ergebnisse bei gleichen Eingangsparametern abgeglichen wurden. Außerdem wurde eine Datenbankschnittstelle für das PSM-System hinzugefügt, da es Daten zum anlernen, wegen einem parallel laufenden anderen Projekt, schon in Datenbanken gibt und diese vom System noch nicht verwertet werden konnten.

5.1.1 Klassenaustausch

Der größte Arbeitsaufwand ergab sich dadurch, dass das die repräsentation der Szenenobjekte im System ausgerauscht werden musste. Im Zuge dieser Arbeit wurde im ganzen PSM-System *pbd_msgs::PbdObject* durch *ISM::Object*, die Objektrepräsentation aus dem

```
include/inference/model/foreground/ocm/shape/HierarchicalShapeModel.h
include/inference/model/foreground/ocm/shape/HierarchicalShapeModelNode.h
include/learner/SceneLearner.h
include/learner/SceneLearningEngine.h
include/learner/SceneModelLearner.h
include/learner/background/BackgroundSceneLearner.h
include/learner/foreground/ForegroundSceneLearner.h
include/learner/foreground/ocm/OcmForegroundSceneLearner.h
include/learner/foreground/ocm/SceneObjectLearner.h
include/learner/foreground/ocm/ocm/OcmSceneObjectLearner.h
include/learner/foreground/ocm/ocm/OcmTree.h
launch/learner.launch
src/inference/model/foreground/ocm/shape/HierarchicalShapeModel.cpp
src/inference/model/foreground/ocm/shape/HierarchicalShapeModelNode.cpp
src/learner.cpp
src/learner/SceneLearner.cpp
src/learner/SceneLearningEngine.cpp
src/learner/SceneModelLearner.cpp
src/learner/background/BackgroundSceneLearner.cpp
src/learner/foreground/ForegroundSceneLearner.cpp
src/learner/foreground/ocm/OcmForegroundSceneLearner.cpp
src/learner/foreground/ocm/ocm/OcmSceneObjectLearner.cpp
src/learner/foreground/ocm/ocm/OcmTree.cpp
include/visualization/psm/ProbabilisticPrimarySceneObjectVisualization.h
include/visualization/psm/ProbabilisticSecondarySceneObjectVisualization.h
include/visualization/psm/helper/CoordinateFrameVisualizer.h
include/visualization/psm/helper/GaussianKernelVisualizer.h
include/visualization/psm/helper/KinematicChainVisualizer.h
include/visualization/psm/helper/SampleVisualizer.h
src/visualization/psm/ProbabilisticPrimarySceneObjectVisualization.cpp
src/visualization/psm/ProbabilisticSecondarySceneObjectVisualization.cpp
src/visualization/psm/helper/CoordinateFrameVisualizer.cpp
src/visualization/psm/helper/GaussianKernelVisualizer.cpp
src/visualization/psm/helper/KinematicChainVisualizer.cpp
src/visualization/psm/helper/SampleVisualizer.cpp
include/trainer/PSMTrainer.h
include/trainer/TreeNode.h
include/trainer/generator/heuristic/DirectionRelationHeuristic.h
include/trainer/generator/heuristic/HeuristicalTreeGenerator.h
include/trainer/source/ObjectSetList.h
include/trainer/source/PbdSceneGraphSource.h
src/trainer/PSMTrainer.cpp
src/trainer/TreeNode.cpp
src/trainer/generator/heuristic/DirectionRelationHeuristic.cpp
src/trainer/generator/heuristic/HeuristicalTreeGenerator.cpp
src/trainer/source/ObjectSetList.cpp
src/trainer/source/PbdSceneGraphSource.cpp
```

Abbildung 5.1: Liste der Dateien die beim Umbau verändert wurden

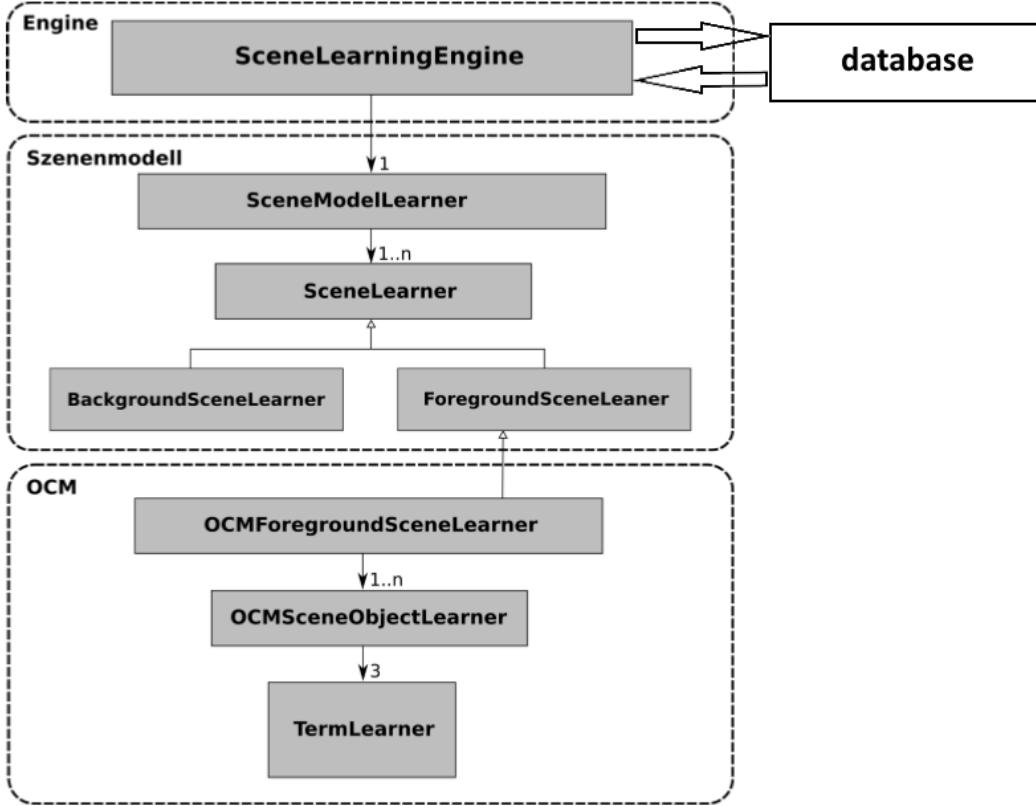


Abbildung 5.2: Datenbankanbindung des Learner-Moduls [Geh14]

ISM-Projekt ersetzt. Später wurden dann die Objekte aus dem ISM-System teils durch *AsrObject* repräsentiert und damit wurde dies auch im PSM-System eingeführt.

In Abbildung 5.1 sind Dateien aufgelistet, die im Zuge des Umbaus geändert werden mussten. Die Liste ist in drei Teile aufgeteilt, die die Änderungen verschiedenen Ordner zuordnen. Der erste Teil ist im *asr_psm*-Ordner, welcher den Kern des PSM-Projekts enthält. Der zweite Teil ist in *asr_psm_visualizations* enthalten, des Komponente die die Visualisierung des PSM-Systems übernimmt. Der dritte und letzte Part ist in *asr_relation_graph_generator* enthalten.

5.1.2 Datenbankeinbindung

Im vergleichbaren ISM-Projekt, welches eine ähnliche Zielsetzung wie das PSM-Projekt hat, allerdings einen nicht stochastischen Ansatz verfolgt, kann man die Szenen, die das System lernen soll aus einer Datenbank auslesen. Aus Gründen der Vergleichbarkeit wie auch dem Komfort ist es sinnig, diese Datenbankschnittstelle für das PSM System nachzurüsten. Um dies zu erreichen wurde ein Datenbankpfad in die Launch-Datei des Learners *learner.launch* hinzugefügt und der veraltete rosbagfiles Parameter damit ersetzt. Die Datenbank wird innerhalb der *SceneLearningEngine* Klasse ausgelesen und

```

class DifferenceForegroundInferenceAlgorithm : public ForegroundInferenceAlgorithm {
public:
    void doInference(std::vector<ISM::Object> pEvidenceList, std::ofstream& pRuntimeLogger);
    double getProbability();
    double differenceBetween(ISM::Object pRoot, ISM::Object pTar, ISM::Object pDiffRoot, ISM::Object pDiffTar);

private:
    double mMaximumDistance;
    double mMaximumRotation;
    double mProbability;
    ISM::RecordedPatternPtr mPattern;
    std::string mDataBaseName;
    std::string patternName;
    boost::shared_ptr<ISM::TableHelper> tableHelper;
    ISM::Object findObjectOfType(std::vector<ISM::Object> pList, std::string pTypeAndObservedId);
    void normalizeVector3d(Eigen::Vector3d input);
};


```

Abbildung 5.3: Header-Datei der Klasse DifferenceForegroundInferenceAlgorithm

die erhaltenen Daten, welche jeweils Vorkommen und von verschiedenen Objekten in diversen Szenen repräsentieren, werden anschließend konvertiert, sodass sie für das System sinnvolle AsrObject-Instanzen werden.

Für den Zugriff auf die Datenbank nutzt die Engine-Klasse die im ISM-System gestellte Hilfsklasse *TableHelper*, die auf die Erkennungssysteme zugeschnittene Funktionen beinhaltet, die man für den Datenbankzugriff nutzen kann.

In Abbildung 5.2 sieht man deutlich, wo der Datenbankzugriff im Learner Klassendiagramm geschieht, nämlich bei der *SceneLearningEngine*-Klasse im gekapselten Bereich der Engine. Die Pfeile symbolisieren den Austausch zwischen Programmcode und Datenbank. Der Programmcode schickt Anfragen, die Datenbank liefert die gewollten Daten, falls die Anfragen korrekt sind. [Geh14]

5.2 Differenzbasierter Erkennungsalgorithmus

Einerseits beschreibt dieses Kapitel, wie der im Konzept vorgestellte Algorithmus programmiertechnisch umgesetzt, andererseits wie dieser in das bestehende System eingebunden wurde.

5.2.1 Algorithmus

Der Algorithmus wurde umgesetzt, wie er in Kapitel 4.2 beschrieben wurde. Zu diesem Zweck wurden die Klassen *DifferenceForegroundInferenceAlgorithm* und *DifferenceBackgroundInferenceAlgorithm* hinzugefügt, welche einen Modus der Szenenerkennung mit dem beschriebenen Algorithmus implementieren, sowie die dazugehörige Hintergrundwahrscheinlichkeit bestimmen. Im Fall des Differenzbasierten Erkennungsalgorithmus wird die Hintergrundszenenwahrscheinlichkeit nicht aufgrund der Anzahl der Objekte oder dem Volumen der Szene berechnet, sondern man stellt manuell einen Schwellenwert ein, da die eben genannten Parameter keinen direkten Einfluss auf den Erkennungsalgorithmus haben. Man muss also einschätzen, wahrscheinlich eine Szene sein muss um wahrscheinlicher als jede beliebige Szene zu sein und den Parameter entsprechend setzen, da die Wahrscheinlichkeiten miteinander in Relation gesetzt werden und man die Szenenwahrscheinlichkeit aller Szenen auch in Bezug auf die Szenenwahrscheinlichkeit der Hintergrundszene berechnet.

In Abbildung 5.3 sieht man eine vereinfachte Version des Headers der Klasse *DifferenceBackgroundInferenceAlgorithm*. Darin sind alle Klassenvariablen und Methoden der Klasse aufgelistet, bis auf den Konstruktor. Dieser initialisiert verschiedene Variablen, wie zum Beispiel die Schwellenwerte der Positionierung und Rotation, welche den maximalen Bereich angeben indem Objekte überhaupt noch erkannt werden sollen. Außerdem wird im Konstruktor auf die Datenbank zugegriffen und die Szeneninstanzen, die die zu testende Szene repräsentieren, werden in der Variable *mPattern* abgespeichert.

Die Methode *doInference* ist die Hauptschnittstelle zum restlichen System und muss von *DifferenceForegroundInferenceAlgorithm* implementiert werden, da die Klasse von der abstrakten Klasse *ForegroundInferenceAlgorithm* erbt, die wiederum von der abstrakten Klasse *InferenceAlgorithm* erbt, welche dies vorschreibt. Die Methode bekommt die Objekte der Objekterkennung gegeben, sowie auch das im Learner erstellte Szenenmodell und berechnet damit die Szenenwahrscheinlichkeit. Das Szenenmodell wird nur genutzt um Objekte zu identifizieren, da sich die Erkennung nur auf die aus der Datenbank ausgelesenen Daten stützt. In der *doInference*-Methode ist der eigentliche Algorithmus mit mehreren *for*-Schleifen und unter Nutzung der anderen Methoden und Variablen der Klasse implementiert.

Die *getProbability*-Methode gibt die aktuelle Wahrscheinlichkeit aus, die in der Klasse eingespeichert ist. Auf diese wird zugegriffen, wenn die Ausgabe des PSM-Systems die Wahrscheinlichkeit der verschiedenen Szenen abfragt um diese in Relation zu setzen und auszugeben. Auch diese Methode muss in der Klasse Implementiert sein, da sie von abstrakten Klassen geerbt wird.

Die Methode *differenceBetween* implementiert berechnung der Objektwahrscheinlichkeit, die in 4.3.1 beschrieben wird. Es werden vier Objekte übergeben und erst sowohl die Positions- als auch die Rotationsrelation zwischen dem ersten und zweiten Paar berechnet und anschließend die Relationen verglichen. Bei diesem vergleich bekommt man eine gewisse Differenz, welche mit Nutzung der Schwellenwerte zu der Objektwahrscheinlichkeit, also einem *double*-Wert, umgeformt werden. Die Methode ist als public aufgeführt,

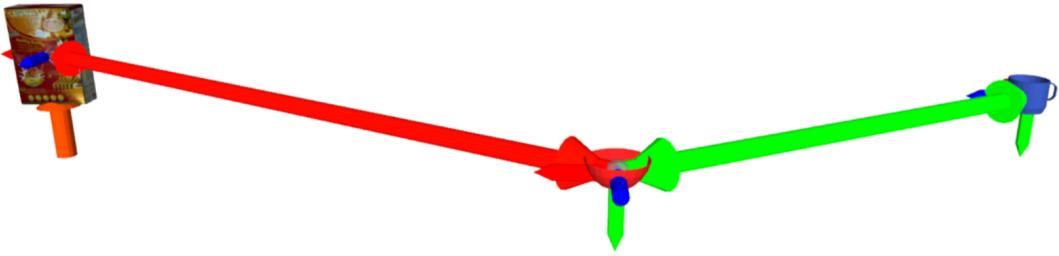


Abbildung 5.4: Inferenz des PSM-Systems [Gaf17]

damit man sich auch mit anderen Klassen an dieser bedienen kann, falls man in Zukunft eine differenzbasierte Objektwahrscheinlichkeit benötigt.

Die Variablen *mMaximumDistance* und *mMaximumRotation* sind die Schwellenwerte für Distanz und Rotation, welche zur Wahrscheinlichkeitsbestimmung notwendig sind. Für *mMaximumDistance* ist jeder Wert größer oder gleich 0 zulässig, für *mMaximumRotation* nur Werte von 0 bis 180. Wenn einer der Werte auf 0 gesetzt ist ist die Erkennung nicht auf einen Tolleranzwert von 0 gesetzt, sondern die zum Parameter gehörige Sparte ist bei der Erkennung deaktiviert. Dies hat den Sinn diese Deaktivierung ohne weitere Variablen zu implementieren. Außerdem schafft man in der Realität sowieso nicht, dass exakt die gleiche Position oder Rotation erkannt wird und hätte das Problem, dass bei Schwellenwert 0, der jeweilige Teil der Objektwahrscheinlichkeit nur 0 oder 1 sein könnte. Das ginge gegen jedweden stochastischen Ansatz.

In *mProbability* wird die Wahrscheinlichkeit eingespeichert, sobald sie berechnet ist. Diese wird dann von *getProbability* ausgegeben. Des Weiteren wird die eingespeicherte Wahrscheinlichkeit als Vergleichswert genutzt, um die maximale Wahrscheinlichkeit zu finden.

Die Variable *mPattern* speichert die Objekte die aus der Datenbank ausgelesen wurden, um sie im Laufe der Inferenz zu nutzen. Auf *mPattern* wird nur im Konstruktator schreibend zugegriffen und sonst nur lesend, damit die erhaltenen Daten niemals verfälscht werden und das System robust gegenüber Wiederholungen ist. Die Objekte sind in sogenannten Pattern gespeichert, die eine Szene repräsentieren, welche wiederum Sets enthalten. Diese Sets werden dann im Algorithmus jeweils setweise betrachtet und mit den erkannten Objekten verglichen.

Die Variable *mDataBaseName* enthält den Pfad zur Datenbank, der im Konstrukt zur Abfrage genutzt wird. Das System wurde nur mit absoluten Pfaden genutzt, aber auch mit realtiven Pfaden arbeiten können.

patternName ist der Szenenname der gerade zu testenden Szene. Dieser ist wichtig, um die richtigen Objekte aus der Datenbank zu lesen, damit die Szenenwahrscheinlichkeit der korrekten Szene bestimmt wird.

Die Variable *tableHelper* zeigt auf eine Instanz der Hilfsklasse *Tablehelper*, welche viele Funktionen beinhaltet die bei dem Zugriff auf Datenbanken nützlich sind. Im Konstruktator wird die Funktion *getRecordedPattern* genutzt, welche einen *string* annimmt, der die Bezeichnung für eine Szene ist und die Objekte dieser Szene zurückliefert. Die Objekte

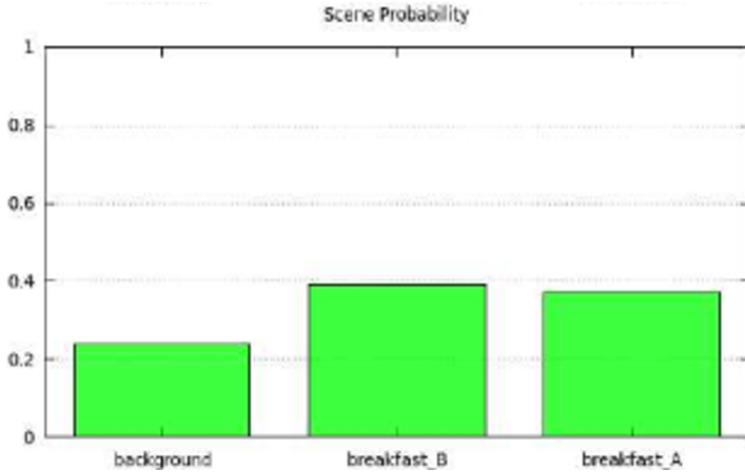


Abbildung 5.5: Wahrscheinlichkeitsausgabe der Inferenz

sind unter anderem als Sets gespeichert, sodass man setweise über die Objekte iterieren kann.

Die Methode *findObjectOfType* bekommt eine Liste von Objekten und eine eindeutige Identifikation eines Objekts übergeben und gibt das Objekt aus der Liste mit der gegebenen Identifikation zurück. Falls kein Objekt der Liste die gegebene Identifikation hat, wird ein NULL Objekt zurückgegeben.

Die Methode *normalizeVector3d* normiert einen gegebene 3D Vektor, das heißt, dieser wird auf Länge 1 gebracht. Die Funktion wurde für verschiedene Berechnungen benutzt, die während des Testens und des Forschens an dem Algorithmus und dem Modus an sich, wird allerdings momentan nicht mehr genutzt und ist veraltet. Der Vollständigkeit halber und falls noch weiter daran geforscht wird und die Funktion eventuell nochmal gebraucht wird, wurde sie in der Klasse behalten.

In Abbildung 5.6 wird der Programmcode innerhalb der *DoInference*-Methode vereinfacht dargestellt. In rot gefärbt sind dabei Kommentare, die den Programmcode erklären. Welche Objekte im Algorithmus wie benannt sind, wird im Kapitel 4.2 erklärt. Es fällt auf, dass der Algorithmus geringfügig vom Konzept abweicht, da die beiden äußeren Schleifen vertauscht sind. Allerdings macht dies für das Ergebnis keinen Unterschied und hat den Sinn, dass man niemals von einem Referenzobjekt ausgeht, welches nicht gemessen wurde, da dies ein unnötiger Schleifendurchlauf wäre.

5.2.2 Einbindung PSM

Da davon abgesehen wird die Parametrisierung des Szenenmodells zu nutzen, braucht man den Learner und das Szenenmodell nur noch zu dem Zweck, dass überliefert wird, welche Szenen potentiell erkannt werden sollen. Die einzigen Informationen die der Algorithmus nutzt sind der Name der Szene und Anzahl der enthaltenen Objekte, sowie die Identifikationen dieser. Das Learner-Modul musste für den neuen Algorithmus also nicht weiter überarbeitet werden sondern wird genauso eingesetzt, wie es auch im

```

mProbability = 0.0; Ausgabewahrscheinlichkeit
double maxProbability = 0.0; Maximale Wahrscheinlichkeit, die bis jetzt berechnet wurde
for(auto recordedObject : pEvidenceList){ Äußere Schleife über alle Objekte die gefunden wurden
    std::string currentType = recordedObject.type + recordedObject.observedId; Setzt gemessenes Referenzobjekt für diesen Schleifendurchlauf
    for(ISM::ObjectSetPtr objectSet : mPattern->objectSets){ Schleife über Szeneninstanzen aus der Datenbank
        ISM::Object testReference; Das Referenzobjekt aus der Datenbank
        for(auto object : objectSet->objects) { Schleife über alle Objekte in der Szeneninstanz des Schleifendurchlaufs
            std::string objectType = object.type + object.observedId; Bestimme die eindeutige Identifikation des Objekts
            if(objectType.compare(currentType) == 0) Überprüfe, ob das Objekt der ID des gemessenen Referenzobjekts entspricht
                testReference = *object; Falls ja, wird das Referenzobjekt gesetzt
        }
        if(testReference.type.compare("") ){ Überprüfe, ob ein Referenzobjekt gesetzt ist - Falls ja, fahre fort
            double testProbability = 0.0; Szenenwahrscheinlichkeit, die nur lokal vorkommt
            for(auto object : objectSet->objects) { Schleife über alle Objekte in der Szeneninstanz des Schleifendurchlaufs
                std::string objectType = object.type + object.observedId; Bestimme die eindeutige ID des Testobjekts
                double innerTestProbability = 0.0; Setze die Objektwahrscheinlichkeit = 0
                if(objectType.compare(currentType) != 0){ Überprüfe ob es sich nicht um das Referenzobjekt handelt - Falls ja, fahre fort
                    ISM::Object innerRecordedObject = findObjectType(pEvidenceList, objectType); Finde das gemessene Testobjekt in der Messung
                    if(innerRecordedObject.type.compare("") != 0) Überprüfe ob ein Objekt gefunden wurde - Falls ja, fahre fort
                    innerTestProbability = differenceBetween(recordedObject, innerRecordedObject, testReference, *object); Berechne die Objektwahrscheinlichkeit der Testobjekt
                }
                testProbability += innerTestProbability; Addiere die Objektwahrscheinlichkeit auf
            }
            int objectCount = objectSet->objects.size(); Bestimme die Anzahl der Objekte der Szeneninstanz
            if((testProbability + 1.0) / (objectCount * 1.0) > maxProbability) Überprüfe, ob eine größere Szenenwahrscheinlichkeit gefunden wurde
            maxProbability = (testProbability + 1.0) / (objectCount * 1.0); Falls ja, speichere diese
        }
    }
}

if(maxProbability > mProbability) Überprüfe, ob eine größere Ausgabewahrscheinlichkeit gefunden wurde
mProbability = maxProbability; Falls ja, speichere diese
}

```

Abbildung 5.6: Programmcode des Algorithmus

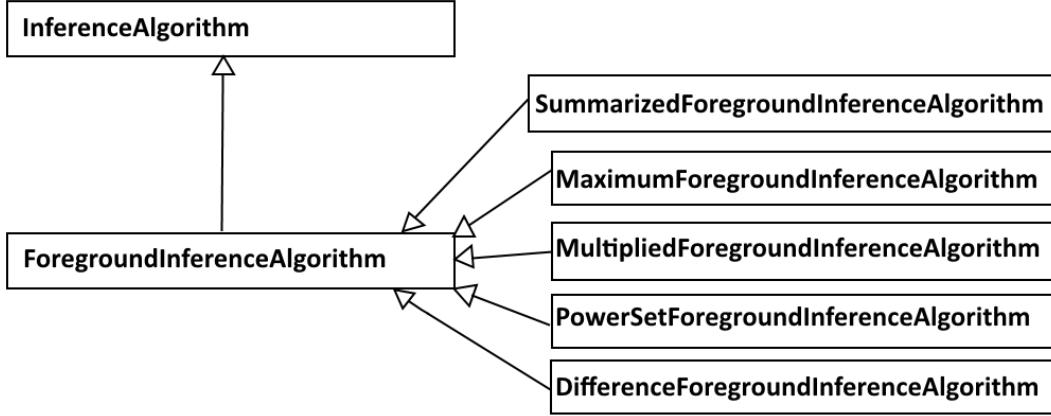


Abbildung 5.7: Klassendiagramm - InferenceAlgorithm-Klasse und Erben

restlichen PSM-System eingesetzt wird. Es wird nur eben nicht auf alle möglichen Daten zugegriffen die zur Verfügung gestellt werden. In der Launch-Datei der Erkennung *inference.launch* wurde ein neuer Parameter hinzugefügt um den Pfad der Datenbank anzugeben. Die Datenbank enthält die Referenzinstanzen von jeder zu testenden Szene. Im Projekt gab es schon verschiedene Modi mit denen man die Erkennung laufen lassen konnte, die unter dem Parameter *inference_algorithm* in der Launch-Datei eingestellt werden konnten. Zu dem neuen Modus, für den man den Parameter auf *difference* setzen muss, gibt es noch die Varianten *powerset*, *summarized*, *multiplied* und *maximum*. Von diesen wurde zuletzt nur noch der *maximum*-Modus genutzt.

Abbildung 5.7 zeigt teilweise die Vererbung der Klasse *InferenceAlgorithm*. Es wurde der Teilbaum der Erben für die Klasse *BackgroundInferenceAlgorithm* weggelassen, da dieser symmetrisch zu der anderen Hälfte der Vererbung ist. Zu Jeder Vordergrundklasse gibt es eine Hintergrundklasse die statt mit *Foreground* mit *Background* gekennzeichnet ist. Um den Algorithmus richtig einzubinden wurde demnach sowohl die Klasse *DifferenceForegroundInferenceAlgorithm* als auch die Klasse *DifferenceBackgroundAlgorithm* hinzugefügt.

6. Evaluation

Im Kapitel Evaluation wird der implementierte Modus mittels Experimenten getestet. Dies geschieht indem das Erkennungsergebnis des neuen Modus mit dem des bestehenden PSM-Systems verglichen wird. Zu diesem Zweck unterscheiden wir zwischen der differenzbasierten und der parametisierten Erkennung. Die differenzbasierte Erkennung ist der neue Modus, der im Zuge dieser Arbeit entwickelt wurde, die parametisierte Erkennung ist die schon vorhandene Erkennung die mittels des Modells und der dort angelernten Parameter Szenenwahrscheinlichkeiten einschätzt.

Experiment 1(Kapitel 6.1) hat eine typische Büroszene als Aufbau und testet die Erkennung anhand der angelernten Daten. Experiment 2(Kapitel 6.2) beinhaltet mehrere Frühstücksszenen und es werden verschiedene Objekte zu der Evidenz hinzugefügt oder weggenommen. Die Erkennungsdaten werden währenddessen jeweils von parametrisierter und differenzbasierter Erkennung gemessen, verglichen und interpretiert.

6.1 Experiment 1: Büro

In diesem Experiment wird überprüft, wie sich die Erkennung verhält, wenn genau die Daten erkannt werden, welche das System zum anlernen für die Szene erhalten hat. Die Szene enthält zwei Bildschirme, eine Computermaus und eine Tastatur. Maus und Tastatur sind nebeneinander vor den beiden Bildschirmen, welche auch nebeneinander stehen. Die Objekte und deren ungefähre Anordnung kann man in Abbildung 6.4 sehen. Die Szene soll einen gewöhnlichen Rechtshänderarbeitsplatz mit Computer beschreiben.

Zwecks dieses Experiments wurden Daten dieser Szene mit dem Learner-Modul angelernt. Dies wurde nur einmal für die differenzbasierte und die parametisierte Erkennung durchgeführt, sodass die beiden Erkennungen auf den gleichen Modelldaten ablaufen und die Erkennungen gut verglichen werden können. Die Szene *office* ist dadurch in beiden Erkennungen gleich definiert. Auch am Computer sowie an dem kompletten PSM-System mit dem getestet wurde, wurde nichts geändert während dem Experiment, um jegliche

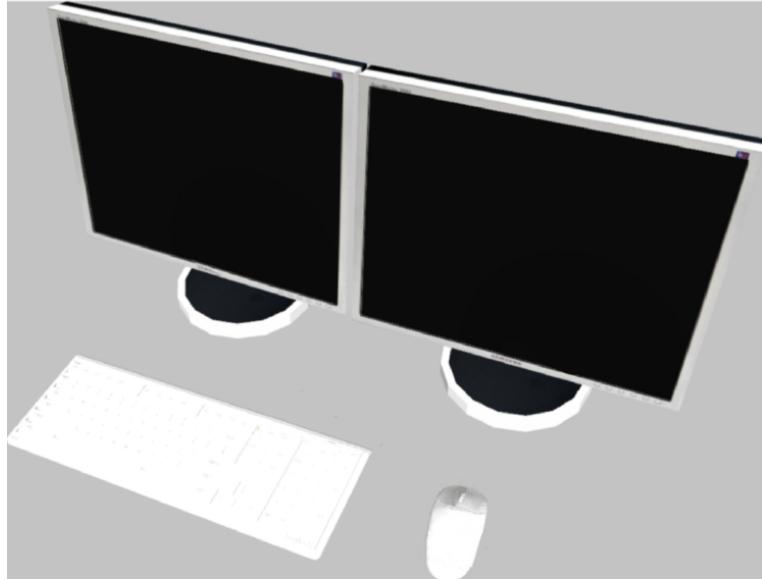


Abbildung 6.1: Alle Objekte, die bei Experiment 1 genutzt werden [Gaß17]

Unstimmigkeit zwischen den Einzelerkennungen zu vermeiden.

In Abbildung 6.2 sieht man eine Visualisierung der angelernten Daten. In der Abbildung sieht man einerseits große rote Pfeile, die jeweils die Ausrichtung des Objekts zeigen und kleinere bunte Linien, welche jeweils eine Position und Orientierung des Objekts beschreiben. Man sieht zum Beispiel, dass die Tastatur meist in einem gewissen Viereck verrückt wird, die Maus in verschiedene Richtungen verschoben wird und unter anderem, dass die beiden Bildschirme nicht verrückt werden, sondern nur in der Höhe verstellbar sind. Die Positionen sind jeweils in rot, blau und grün gehalten und zusätzlich wurden die Bereiche eingefärbt in denen sich das jeweilige Objekt am ehesten aufhält. Es wird also ein gewisses Muster zwischen den anzulernenden Objekten gesucht und eingespeichert.

Nach dem Anlernen der Daten wurden die selben Daten so im System simuliert, als würden sie gerade als Objektevidenz von der Objekterkennung erkannt werden. Die beiden Erkennungsalgorithmen wurden nacheinander laufen gelassen, allerdings wurden die jeweiligen Szeneninstanzen, das heißt, die Sets aus den zu erkennenden Daten, einander zugeordnet. Dadurch kann man die beiden Erkennungen setweise miteinander vergleichen und ein Gesamtergebnis aus den Einzelvergleichen ziehen.

In Abbildung 6.3 sind beispielhaft ein paar Vergleiche dokumentiert. Die Wahrscheinlichkeiten sind jeweils relativ zueinander berechnet, nehmen Werte zwischen 0 und 1 an und ergeben in der Summe 1. Links sieht man jeweils das Erkennungsergebnis der parametrisierten Erkennung eines Sets, rechts das jeweilige Ergebnis der differenzbasierten Erkennung. Man sieht, dass die parametrisierte Erkennung jeweils einen hohen Wahrscheinlichkeitswert für die *office*-Szene hat und einen geringen für die Hintergrundszene. Die differenzbasierte Erkennung erkennt die *office*-Szene in jedem Set mit einer beinahe

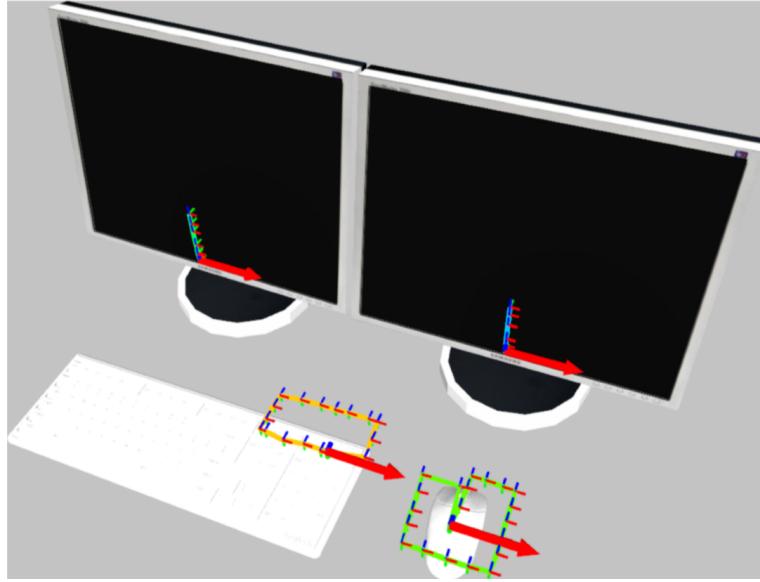


Abbildung 6.2: Angelernte Daten der Büroszene [Gaß17]

hundertprozentigen Wahrscheinlichkeit und schwangt kaum zwischen den Sets, während die parametisierte Erkennung kleine Unterschiede zwischen den einzelnen Messungen aufweist. Alle Messungen erkennen die *office*-Szene beinahe eindeutig.

Diese Beobachtungen sind wie folgt zu erklären. Dass in jeder Messung die Szene mit hoher Wahrscheinlichkeit erkannt wird ist nicht verwunderlich, da die Daten, die zur Modellerzeugung genutzt wurden, schließlich die im Verlauf des Experiments erkannten Daten enthalten. Man muss stets bedenken, dass die Wahrscheinlichkeiten, die in den Messungen vorkommen relative Wahrscheinlichkeiten sind. Das heißt, man sieht wie Wahrscheinlich es ist dass die Szene *office* auftritt im Vergleich zu der Hintergrundszene. Damit ist die Wahrscheinlichkeit, die bei der parametrisierten Erkennung für die *office*-Szene ausgegeben wird immer etwas geringer, als sie erscheint, da sie im Bezug auf eine geringe Hintergrundwahrscheinlichkeit berechnet wurde. Genauso sind die Schwankungen der eigentlichen Szenenwahrscheinlichkeit möglicherweise größer, als es auf den ersten Blick erscheint. Im gesamten sollte die Szenenwahrscheinlichkeit überall zwischen 80 und 90 Prozent liegen.

Das nicht komplett eindeutige Erkennungsergebnis lässt sich auf die Parametrisierung der Erkennung zurück führen. Es wurde ein Modell aufgrund aller Daten erstellt, welches natürlich nicht mehr jede einzelne Messung mit hundertprozentiger Wahrscheinlichkeit erkennt, da es ein größeres Modell ist, dass sich alle Einzelmessungen zur Basis genommen hat. Die Schwankungen der Messung können sich allerdings auch durch die dynamische Hintergrundwahrscheinlichkeit begründen. Wenn das Erkennungssystem unterschiedliche Daten erhält, verändert sich die errechnete Hintergrundszenenwahrscheinlichkeit, wodurch sich logischerweise auch die relative Wahrscheinlichkeit der Szenen verändert.

Die Erkennungsergebnisse für die differenzbasierte Erkennung lassen sich leicht inter-

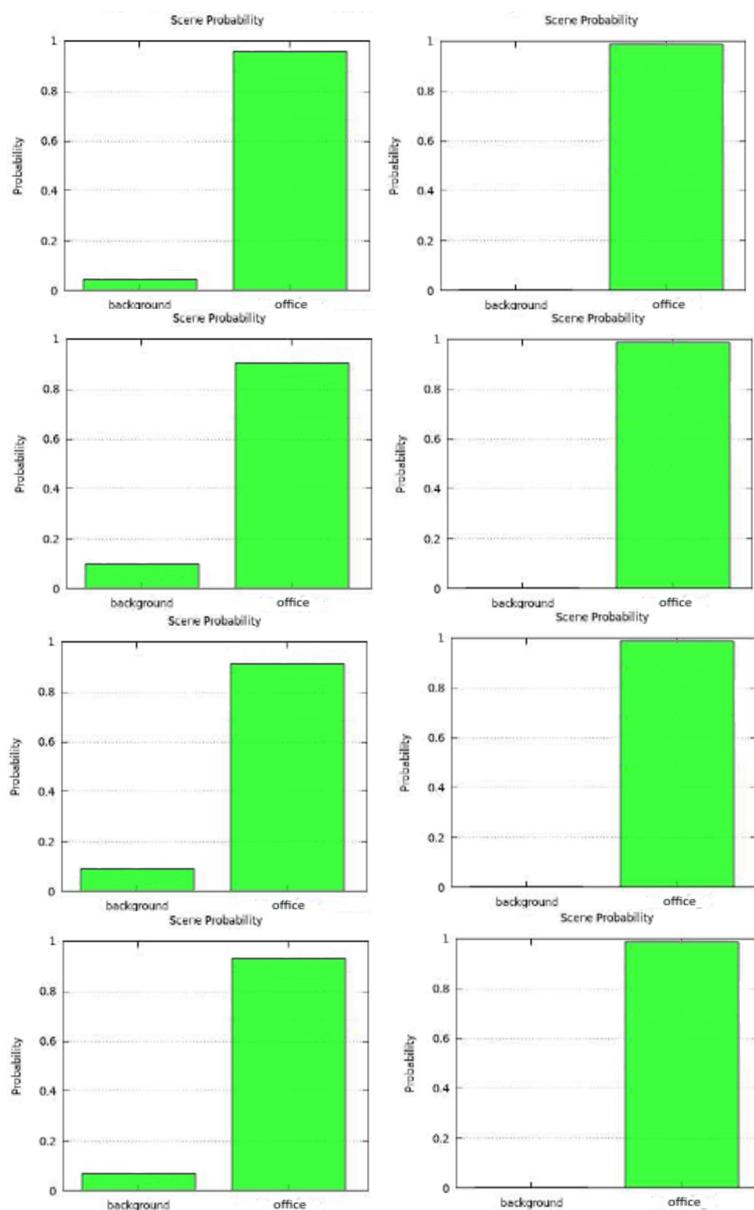


Abbildung 6.3: Erkennungswahrscheinlichkeiten der Büroszene



Abbildung 6.4: Alle Objekte, die bei Experiment 2 genutzt werden [Gaß17]

pretieren und erklären. Die Wahrscheinlichkeit ist bei jeder Messung beinahe hundertprozentig, da es immer eine Szeneninstanz in den Daten gibt, die sich von der erkannten Evidenz nicht unterscheiden lässt. Dadurch ergibt sich ein hundertprozentiges Ergebnis bei der Szenenwahrscheinlichkeit, welches dann mit der geringen Hintergrundwahrscheinlichkeit verrechnet wird. Die geringen Schwankungen lassen sich dadurch erklären, dass die Hintergrundwahrscheinlichkeit nicht statisch berechnet wird, sondern ein Schwellenwert als Parameter gesetzt wird. Dadurch schwankt die Hintergrundwahrscheinlichkeit nicht und kann auch das Ergebnis in keiner Weise beeinflussen.

Im gesamten lässt sich sagen, dass die differenzbasierte Erkennung die Szene noch sicherer erkannt hat, als die parametisierte Erkennung, da sie direkt mit den Daten arbeitet, welche sie gelernt hat. Das Experiment zeigt damit, dass es in gewissen Fällen präziser sein kann direkt mit den Beispieldaten zu arbeiten, die man von einer Szene hat. Allerdings ist es in der Realität unmöglich eine Szene perfekt aufzubauen, wie sie in den Daten vorkommt. Um die differenzbasierte Erkennung auch in einem weniger realitätsfernen Kontext zu testen folgt Experiment 2(Kapitel 6.2).

6.2 Experiment 2: Frühstück

Im Experiment Frühstück werden mehrere verschiedene Frühstücksszenen erkannt. Es werden wieder beide Erkennungsalgorithmen die gleiche Evidenz überprüfen gelassen und die daraus resultierenden Ergebnisse verglichen und interpretiert. In Abbildung 6.4 sieht man alle Objekte die im Zuge dieses Experiments vorkommen und deren Orientierung. Die Objekte werden im Versuch nur auf die Weise angeordnet sein, wie sie in der Abbildung zu erkennen sind. Aus dieser Anordnung werden dann Objekte entfernt oder hinzugefügt. Für die Verständlichkeit werden allen Objekten im folgenden Namen zugeordnet. Von links nach rechts sieht man eine Schüssel, eine Kaffeebox, eine Müslipackung, einen Becher, einen Teller, eine Cornflakespackung und einen Milchkrug.

Diese Gesamtheit der in diesem Experiment vorkommenden Objekte werden nun ihren Szenen zugeteilt. Dabei gilt keine Exklusivität eines Objekts für eine Szene, ein Objekt kann also auch in beiden Szenen vorkommen. Die zwei Frühstücksszenen die in diesem Experiment angelernt und erkannt werden sollen sieht man in Abbildung 6.5. Die linke Szene trägt im System den Titel *breakfast_A* und die rechte den Titel *breakfast_B*.

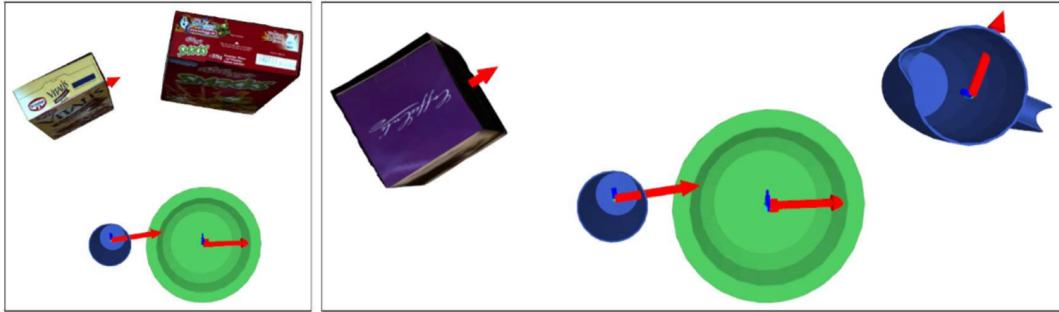


Abbildung 6.5: Beispiel beider Frühstückszenen [Gaß17]

Außerdem gibt es in jeder im folgenden vorkommenden Messung die obligatorische Hintergrundszene *background*.

Die verschiedenen Objekte die in Abbildung 6.4 gezeigt wurden, wurden nun in verschiedenen Konstellationen und Kombinationen simuliert und die entsprechenden Erkennungsergebnisse der differenzbasierten und parametisierten Erkennung festgehalten. In Abbildung 6.6 sieht man ausgewählte Erkennungsergebnisse, die bei diesem Experiment entstanden sind. Links ist jeweils das Ergebnis der parametisierten Erkennung, in der Mitte sieht man die aktuelle Evidenz die überprüft wird und rechts werden die Erkennungsergebnisse der differenzbasierten Erkennung dargestellt. Dabei gliedert sich die Abbildung zeilenweise in sieben verschiedene Messungen, auf die im folgenden als Messung 1 bis 7 referenziert wird.

Messung 1 hat einen möglichen Aufbau der Szene *breakfast_A* zur Evidenz. Wie erwartet erkennen beide Erkennungsmodi die aufgebaute Szene mit hoher Wahrscheinlichkeit und die Szene *breakfast_B* hat einen niedrigen Wert. Die differenzbasierte Erkennung hat jedoch noch eindeutigere Wahrscheinlichkeitswerte als die parametisierte Erkennung, *breakfast_A* hat also einen höheren Wert und die andere Szene einen niedrigeren. Dies lässt sich auf darauf zurückführen, dass die Szene in etwa wie sie aufgebaut wurde, auch in den Daten vorkommt. Dadurch hat die Szene eine besonders hohe Szenenwahrscheinlichkeit gegenüber der parametisierten Erkennung. Die parametisierte Erkennung hat außerdem einen höheren Wert für die Szene *breakfast_B*, da zwei Objekte aus dem Modell genau richtig platziert sind. Dies erzeugt einen etwas größeren Wert als die entsprechende Szenenwahrscheinlichkeit bei der differenzbasierten Erkennung.

Messung 2 enthält alle Objekte aus Szene *breakfast_A* bis auf die Müslipackung. Damit enthält sie drei Objekte der Szene A und zwei Objekte der Szene B. Dies sieht man auch in den Wahrscheinlichkeitsergebnissen. Die Hintergrundszene wird wahrscheinlicher, da keine Szene komplett erkannt wurde und Szene B bekommt einen höheren Wert für die relativen Wahrscheinlichkeiten. Eigentlich sinkt hauptsächlich die Szenenwahrscheinlichkeit für Szene A und dadurch steigen die beiden anderen Wahrscheinlichkeiten an, da sie im Verhältnis zur Wahrscheinlichkeit für Szene A größer werden. Die Ergebnisse der beiden Erkennungen sind ähnlich, nur die Hintergrundwahrscheinlichkeit der differenzbasierten Erkennung ist etwas höher, da die beiden anderen Szenenwahrscheinlichkeiten

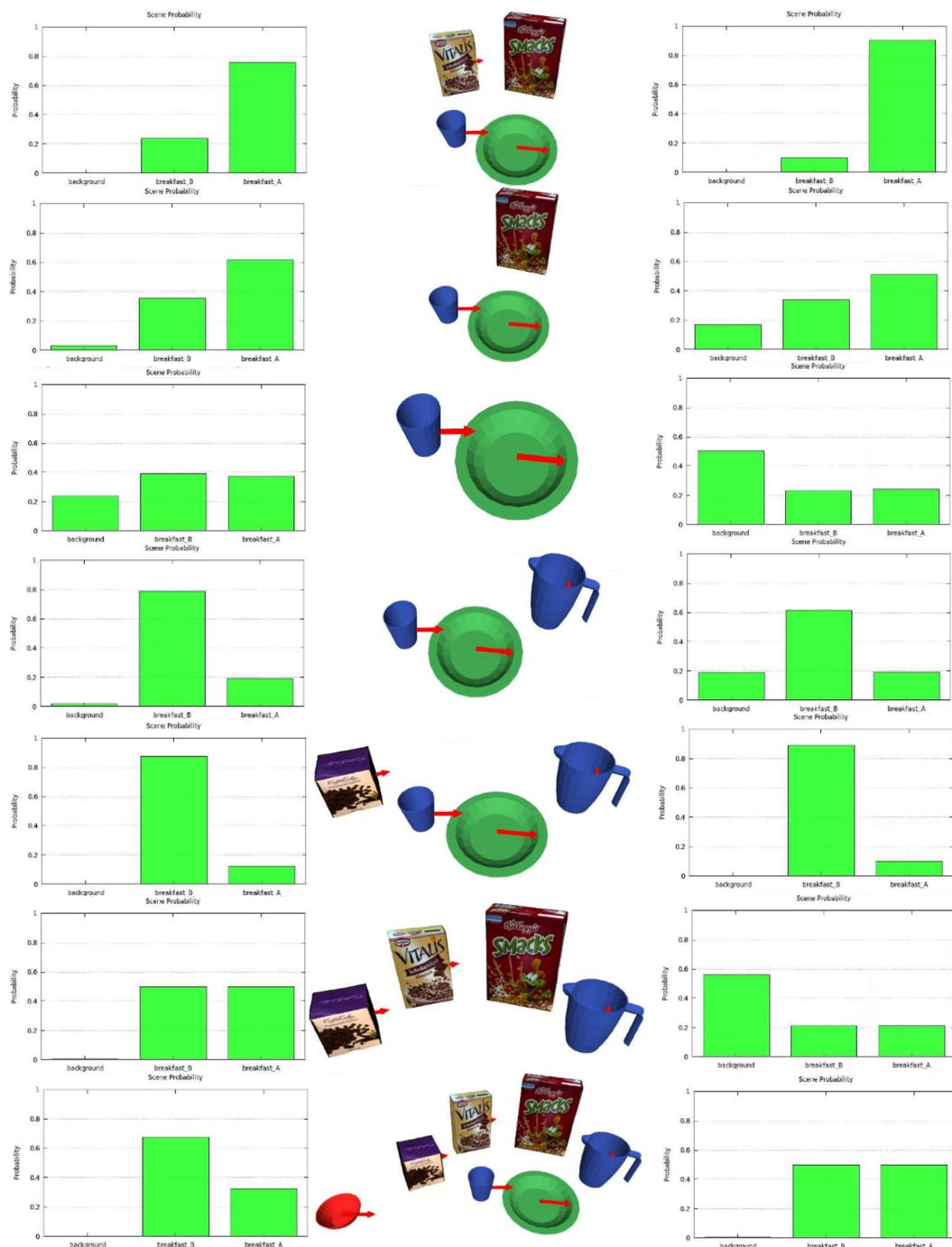


Abbildung 6.6: Erkennungswahrscheinlichkeit der Frühstückszenen [Gaß17]

etwas niedriger sind als bei der parametisierten Erkennung.

In Messung 3 sind nur noch Becher und Teller enthalten, sodass beide Szenen nur halb vorkommen. In beiden Erkennungsmodi sind jeweils die Wahrscheinlichkeiten der Szenen A und B nahezu identisch. Nur die Hintergrundwahrscheinlichkeit unterscheidet sich. Bei der parametisierten Erkennung sind sie etwas niedriger als die anderen Szenenwahrscheinlichkeiten, bei der differenzbasierten Erkennung ist sie in etwa doppelt so groß. Die differenzbasierte Erkennung ist empfindlicher bei fehlenden Objekten als die parametisierte Erkennung.

In Messung 4 wird der Milchkrug zu den vorherigen Objekten hinzugefügt. Das führt dazu, dass die Wahrscheinlichkeit für Szene B bei beiden Erkennungen erheblich ansteigt. Es fällt auf dass bei der differenzbasierten Erkennung die Hintergrundwahrscheinlichkeit erheblich höher und damit die Wahrscheinlichkeit für Szene B etwas kleiner ist. Dies zeigt wiederum, dass Szenen bei Unvollständigkeit nicht so deutlich erkannt werden wie bei der parametisierten Erkennung. Für nahezu eindeutige Erkennung braucht die differenzbasierte Erkennung in jedem Fall das Auftreten aller Objekte.

In Messung 5 entspricht die Evidenz Szene *breakfast_B* und Ergebnisse sind bei beiden Erkennungen beinahe identisch. Eine hohe Wahrscheinlichkeit bei Szene B und eine niedrige Wahrscheinlichkeit bei Szene A. Die Szene wurde ähnlich angelernt wie sie in der Evidenz vorkommt, das Ergebnis ist also nicht verwunderlich. In Messung 6 besteht die Evidenz aus der Kaffeebox, der Müslipackung, der Cornflakespackung und dem Milchkrug. Diese Messung zeigt die signifikantesten Unterschiede zwischen den Erkennungen auf. Bei der parametisierten Erkennung ist die Hintergrundwahrscheinlichkeit beinahe bei null und die Szenenwahrscheinlichkeiten für Szene A und B in etwa bei 0,5. Bei der differenzbasierten Erkennung ist die Hintergrundwahrscheinlichkeit in etwa bei 0,6 und die anderen Wahrscheinlichkeiten bei 0,2. In der Evidenz kommen jeweils zwei Objekte der Szenen vor und die Empfindlichkeit gegenüber Szenenunvollständigkeit der differenzbasierten Erkennung wird nochmals klar ersichtlich.

In Messung 7 sind alle Objekte in der Evidenz vorhanden die in Abbildung ?? gezeigt wurden. In den Wahrscheinlichkeitsergebnissen der beiden Erkennungsergebnisse gibt es unerwarteterweise offensichtliche Unterschiede. Beide Erkennungen weisen eine Hintergrundszenenwahrscheinlichkeit von nahezu null auf, während aber bei der differenzbasierten Erkennung beide Szenen gleichermaßen mit fünfzigprozentiger Auftrittswahrscheinlichkeit ausgegeben werden, wird bei der parametisierten Erkennung Szene B als in etwas doppelt so wahrscheinlich wie Szene A bestimmt.

Da das Ergebnis nicht erwartet wurde, wurde anschließend eine ähnliche Evidenz überprüft, bei der die Schüssel weggelassen wurde. Die Evidenz erzeugt mit beiden Erkennungen ein ausgeglichenes fünfzigprozentiges Ergebnis für Szene A und B. Es scheint als habe das zusätzliche Objekt, welches in keiner der Szenen vorkommt die parametisierte Erkennung beeinflusst, während die differenzbasierte Erkennung davon unbeeinflusst blieb.

6.3 Fazit

Im Zuge der vorliegenden Arbeit wurde ein neuer Modus für das PSM-System entwickelt, welcher dichter an den gegebenen Daten arbeitet. Bei den beiden Experimenten konnte man gut sehen, dass es Situationen gibt in denen der differenzbasierte Modus präziser scheint und Situationen in der gänzlich andere Ergebnisse erzeugt werden. In Experiment 1 sieht man, dass die differenzbasierte Erkennung direkt an den Daten arbeitet und diese, falls sie in der Evidenz vorkommen zuverlässig und eindeutig wiedererkennt.

In Experiment 2 war es schwieriger die Vergleichbarkeit der beiden Erkennungsmodi zu gewährleisten, da beide Erkennungen gewisse Parameter annehmen die nicht direkt vergleichbar sind. Die differenzbasierte Erkennung nimmt Schwellenwerte für die Positionierung und Rotation der Szenenobjekte, sowie für die Hintergrundwahrscheinlichkeit entgegen, während die parametisierte Erkennung alle Parameter über das Modell erhält. Alle Parameter wurden für Experiment 2 mit mehreren Testläufen feinjustiert, allerdings sind sie sicher noch nicht optimal.

Es lässt sich aber sagen, dass die differenzbasierte Erkennung mit der genutzten Parametrisierung empfindlicher auf fehlende Szenenobjekte reagiert und komplett Szenen deutlicher erkennt. Alles in allem ist der neue Modus eine Bereicherung für das PSM-System, die schon allein aus Vergleichszwecken weiter genutzt werden sollte und weiter entwickelt werden kann.

7. Zusammenfassung und Ausblick

Zwei Sätze zu jedem größeren Kapitel

[DSS93]

Literaturverzeichnis

- [DSS93] Randall Davis, Howard Shrobe und Peter Szolovits: *What is a Knowledge Representation?* AI Magazine, 14(1):17–33, 1993. <http://www.aaai.org/ojs/index.php/aimagazine/article/view/1029>.
- [Gaf17] Nikolai Andreas Gaßner: *Probabilistische Szenenerkennung durch hierarchische Constellation Models über räumliche Relationen aus demonstrierten Objekttrajektorien.* Seiten 30–85, 2017.
- [Geh14] Joachim Gehrung: *Ausgewählte Algorithmen zur kombinatorischen Optimierung der räumlichen Relationen in Probabilistic Scene Models.* Seiten 5–63, 2014.