

PRIP-TR-98

January 10, 2006

## A Graph-Based Concept for Spatiotemporal Information in Cognitive Vision<sup>1</sup>

*Adrian Ion and Yll Haxhimusa and Walter G. Kropatsch*

### **Abstract**

A concept relating story-board description of video sequences with spatio-temporal hierarchies build by local contraction processes of spatio-temporal relations is presented. Object trajectories are curves in which their ends and junctions are identified. Junction points happen when two (or more) trajectories touch or cross each other, which we interpret as the “interaction” of two objects. Trajectory connections are interpreted as the high level descriptions.

---

<sup>1</sup>Supported by the Austrian Science Fund under grant FSP-S9103-N04.

# 1 Introduction

Even though there is no generally accepted definition of cognitive vision yet, presumptions about the cognitive capabilities of a system can be made by comparing it's results with that of an entity, already 'known' and accepted to have these capabilities, the human. Also, the *Research Roadmap of Cognitive Vision* [15], presents this emerging discipline as 'a point on a spectrum of theories, models, and techniques with computer vision on one end and cognitive systems at the other'. A conclusion drawn from the previous, is that a good starting point for a representation would bring together the following:

- enable easy extraction of data for human comparison;
- bridge together high and low level abstraction data used for cognitive and computer vision processes.

After 'watching' (analyzing) a video of some complex action, one of the things, that we would expect a cognitive vision system to do, is to be able to correctly answer queries regarding the relative position of occluded objects. Let us take the video <sup>1</sup> given by a simple scenario of two black cups and a yellow ball and describe the scene in simple English words (see the description in Table 1). The description contains: **objects**: hand, cup, ball, table ; **actions**: grasp, release, move, shift etc., and **relations**: to-the-left, to-the-right, in-front-of etc.

Later, we could use this kind of description to compare the results given by the system with ones made by humans. While observing a dynamic scene, an important kind of information is that of the change of an object's location, i.e. the change of topological information. In most of the cases, this kind of change is caused by an active object (e.g. agent: hand, gravity, etc) acting on any number of passive objects (e.g. cup, ball, etc.). Queries like '*where is the ball?*' could be answered if the history of topological changes is created.

From all the work done in the domain of qualitative spatial and temporal information we would like to enumerate the following: Interval calculus [1] is used in systems that require some form of temporal reasoning capabilities. In [1] 13 interval-interval relations are defined: 'before', 'after', 'meets', 'met-by', 'overlaps', 'overlapped-by', 'started-by', 'starts', 'contains', 'during', 'ended-by', 'ends' and 'equals'. In [13], motivated by the work in [1, 7, 8],

---

<sup>1</sup><http://www.prip.tuwien.ac.at/Research/FSPCogVis/Videos/Sequence.2.DivX.avi>

an interval calculus-like formalism for the spatial domain, the so called region connection calculus (RCC) was presented. The set of 8 region-region base relations defined in [13] ( $RCC - 8$ ) are: 'is disconnected from', 'is externally connected with', 'partially overlaps', 'is a tangential proper part of', 'is non-tangential proper part of', 'has a tangential proper part', 'has non-tangential proper part', and 'equals'. A more expressive calculus can be produced with additional relations to describe regions that are either inside, partially inside, or outside other regions ( $RCC - 15$ ). Different graph based representations have been used to describe the changes/ events in a dynamic space. In [6] graphs are used to describe actions (vertices represent actions). Graphs are also used in [2], but here vertices represent objects. Balder [2] argues that arbitrary changes can be best described by state approach: the state of the world before and after the change characterizes the change completely. The Unified Modeling Language, in its state diagram, also defines a graph based representation for tracking temporal changes. The General Analysis Graph (GANAG) [14] is a hierarchical, shape-based graph that is build and used in order to recognize and verify objects. *The analysis graph can be seen as a 'recipe' for solving industrial applications, stating which kind of decisions have to be made at which stage* [14].

In Section 2 we give the spatiotemporal story-board of the video sequence. In Section 3 we describe two methods of contraction of trajectory of movements: first the spatial contraction followed by a temporal contraction (Section 3.1) and than the temporal contraction followed by a spatial contraction (Section 3.2).

## 2 Spatiotemporal Story Board of a Film

The scene history is a description of the actions and spatial changes in the scene. It should depict the spatiotemporal changes in the scene, in a way that could be used to create a human-like description (similar to the one presented in Section 4). For this we propose a graph based representation where vertices represent spatial arrangement states and edges represent actions (see Figure 1a).

Each vertex contains a topological description of the spatial arrangement of the objects in the scene, that results through a transition from a previous state, by applying the actions that link it to the current. What we refer to as objects are actually detected relevant visual entities, which in the ideal case

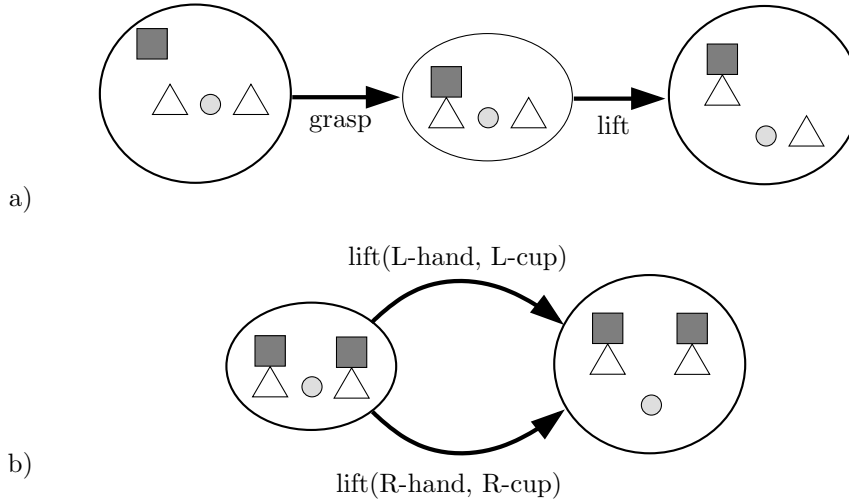


Figure 1: a) History graph. b) Parallel actions.  $\square$  Hand,  $\circ$  Ball,  $\triangle$  Cup.

would be objects or, groups of objects in a “special” physical relation e.g. occluding, containing, etc. Vertices are added when the topological description of the spatial arrangement changes. There are no vertices that contain (identify) the same topological description (scene state). If the scene enters a state, which has a topological description identical to one of the descriptions already identified by a vertex in the scene history graph (it has been in the same state in the past), then an edge/edges from the vertex identifying the previous state, to the existing vertex should be added.

Edges are associated with actions and identify the type/class of the action. Also, each edge links to the objects (from the source and destination state vertex) involved in this particular action. If an object taking part in the action cannot be identified as one of the known objects, a new instance should be created and the edge linked to it. Later on, through reasoning, the new created instance, can be identified as a previously known object or a new one (or some presumption can be made, using certain criteria). In case of simultaneous actions, more than one edge is used to connect 2 vertices. Each edge should describe the actions that happened in parallel. (Figure 1b) shows how to describe 2 hands lifting 2 cups at the same time)

The representation of the scene history as a graph allows us to create higher level abstractions. A straight forward example results from the ‘re-

usage' of vertices (disallowing multiple vertices identifying the same state). Imagine the scenario of a hand grasping and releasing the cup 10 times in a row. Besides saving space by not adding a big number of additional vertices, by identifying cycles, we can easily determine repeated actions and find the shortest way from one configuration to another. Higher level abstractions replace more complex subgraphs containing parallel actions and long sequences of actions resulting in small or unimportant changes for the objects in the system's attention.

A type of information that can be directly extracted from the spatiotemporal graph is the one of 'all known actions'. This information can be represented by a directional graph in which vertices represent unique classes of objects part in any previous action and edges represent simple actions that can involve the connected vertices (usually actions that a class of objects can perform on another class). E.g.: a hand can lift, move, grasp, release, etc. a cup.

We can observe that, in time, for a fixed set of classes of objects involved, if the actions vary enough, the graph of 'all known actions' will converge to the graph of 'all possible actions' and the presented spatiotemporal history graph, will converge to the graph 'of all possible states' (The latter is something that should be avoided, because storing/remembering everything up to the smallest details is guaranteed to sooner or later cause time and memory issues).

Another type of information, that is obtained directly (e.g. tracking) or through reasoning, is that of an object occluding or containing other objects (totally or partially, but still unrecognizable by the detection level). To store this type of information, a relabeling of the class of the occluding object should be done i.e. a cup that has been found out to contain a ball should be labeled 'cup with ball inside'.

### 3 Contraction in Spatiotemporal Space

The idea here would be to contract in 3D (2D space + time) along 'the movement trajectory'. Every frame could be represented by a region adjacency graph. In order to stretch this into time, these region adjacency graphs (region adjacency combinatorial maps) should be matched to each other, i.e. the region adjacency graph at time  $t$  is matched with the one in  $t + 1$  and so on. In this sense we could define a 'trajectory' of each region This trajec-

tory becomes a curve in 3D and with the techniques analogous with that of contraction of a 2D curve pyramid in [11], we can contract regions adjacent along this curve to produce the more abstract representation of the scene, e.g. where the movement started, where it ended etc (Figure 2).

If the analyzed scene has a structured background, then, depending on it's granularity, this is enough to detect movement using only topological information. On the other hand, this will increase the number of consecutive frames that differ with respect to topological relations. To reduce the abundance of topological states, to a set containing the most relevant ones, a set of adaptive pyramids is used. There are no constraints regarding the time intervals between 2 consecutive states. Actually, it is expected that in most of the cases where natural movement is present (not robots repeating some predefined action) these time intervals will differ quite a lot.

In subsections 3.1 and 3.2 we present two approaches, to the problem, which basically differ only in the order in which contraction in the spatial and temporal domains, is done. The first, avoids the difficult problem of graph matching by creating pyramids in the first step and then doing the matching using the pyramids. The second, while needing graph matching to be done, should have a lower memory usage. Moreover, in the ideal case, the resulting top level of the 2 approaches should be the same.

### **3.1 Spatial contraction followed by temporal contraction**

For each frame, whose topological description is different from the one of the previous frame, a space-contraction pyramid is build, that preserves only the spatial information required by the higher functionality levels (i.e reasoning) and by the time-contraction. A space-contraction pyramid is a pyramid where elements, from the same scene state, neighbored from a spatial point of view are contracted, and a time-contraction pyramid is a pyramid where elements, neighbored from a temporal point of view (consecutive scene states) are contracted.

To obtain the base level of the time-contraction pyramid from the set of space-contraction pyramids a matching step has to be performed (Figure 4). Each 2 consecutive pyramids (from a chronological perspective) have to be matched, and the vertices that represent the same object/visual entity should be linked by an edge (if it is possible i.e. if the same object/visual entity exists

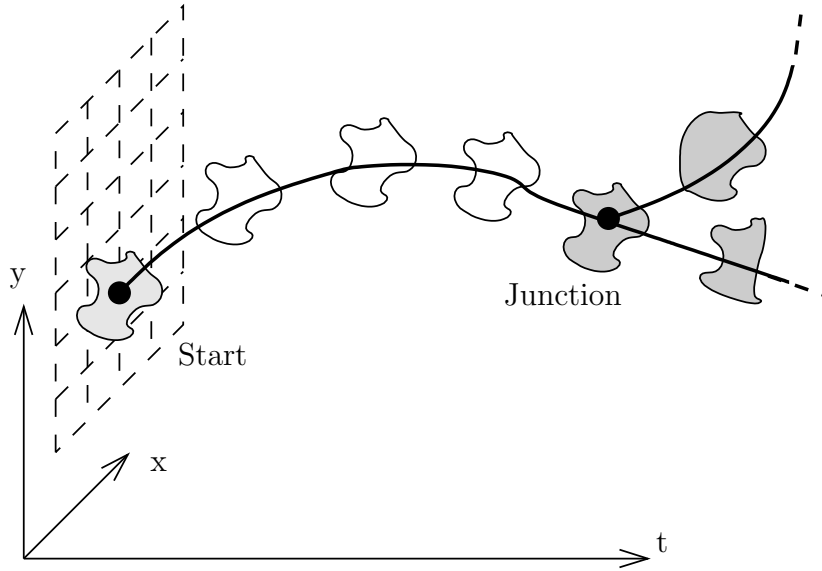


Figure 2: Trajectory of movements.

in both structures - existed in both frames). If a certain object/visual entity, that exists in one of the pyramids, does not exist in the other (occlusion, moved out of the field of view, etc.), no connecting edge can be created, thus obtaining a trajectory endpoint. If similar entities disappear and reappear at different time intervals, it will be the job of the reasoning part to decide whether it was the same instance of the same class or not.

The base level of the time-contraction pyramid contains a vertex for each of the frames in the source video, that differ in topological relations from the previous frame. Each vertex will contain the space-contraction pyramid for the region adjacency graph of the respective scene state. These vertices are linked together in a chronological manner i.e. each vertex is linked to the one of the previous and next frames. Also, as a result of the pyramid matching process mentioned before, the vertices from the consecutive space-contraction pyramids are linked together, showing the trajectories of the regions from the first through the last frame. For example: take the topological descriptions for each frame and represent them in a 3D space, where one of the dimensions is time, and the other 2 are used to represent the planar region adjacency graphs. If for every 2 consecutive graphs, the vertices representing the same object/visual entity are linked together by an edge, then following these

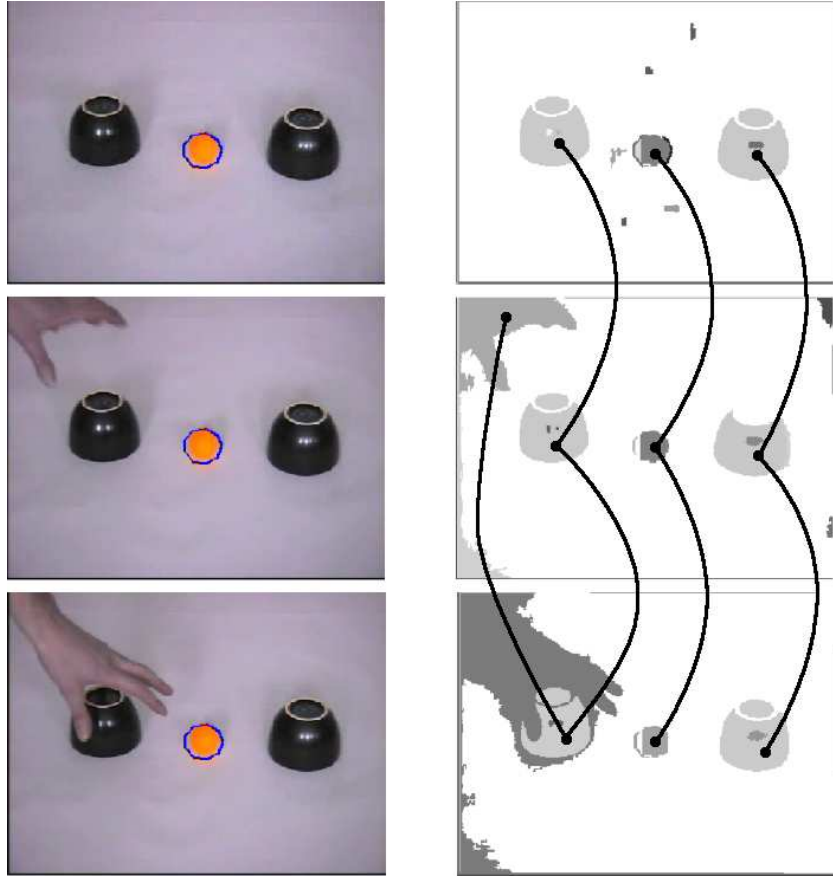


Figure 3: An example of trajectory of movements.

inter-state connection edges will produce the regions trajectory in 3D space.

Each level of the time-contraction pyramid is a chronologically ordered list of space-contraction pyramids, each element describing the topological relations of a certain scene state. The space-contraction step reduces the spatial information in areas that are not of our interest. The purpose of the time-contraction pyramid is to skip the unnecessary frames caused by the presence of the structured background (which is needed for movement detection using only topological information).



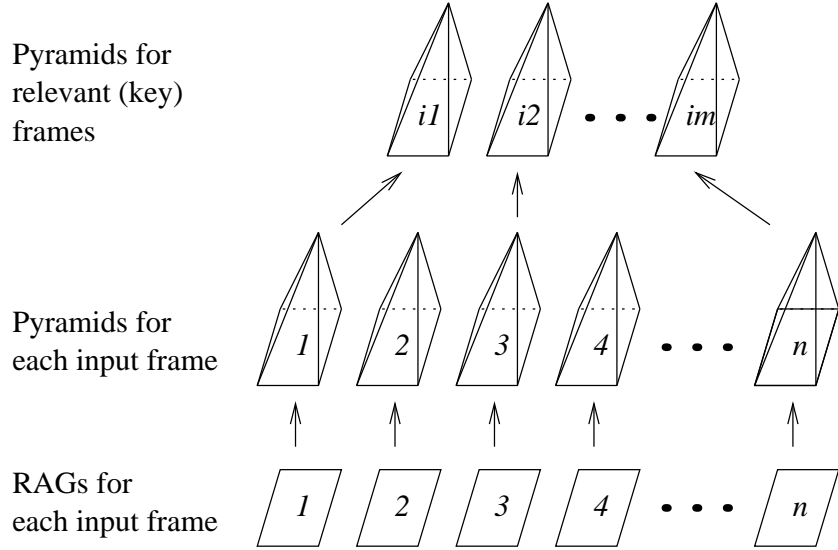


Figure 4: Space time contraction.

### 3.2 Temporal contraction followed by spatial contraction

The base level of the time-contraction pyramid contains a vertex for each of the frames in the source video, that differ in topological relations from the previous frame (Figure 5). Each of these vertices contains the region adjacency graph (RAG) for the respective frames. Through a preliminary process of matching, each vertex in a region adjacency graph should be connected with the vertex(vertices), from the two neighboring graphs, that represent the same object/visual entity (if it is possible i.e. if the same object/visual entity exists in the neighboring region adjacency graphs frame). In other words, the base level of the pyramid is the discretized evolution of the region adjacency graph of the presented scene with the exception that identical consecutive states are merged into a single state.

If we would represent the base level structure in a N dimensional space (3D for 2D state descriptions + time) we would see that we have obtained curves representing the trajectories of the different regions analyzed. A line segment parallel to the time axis, will denote a static region through the respective time interval. Each level of the pyramid is made out of a sequence

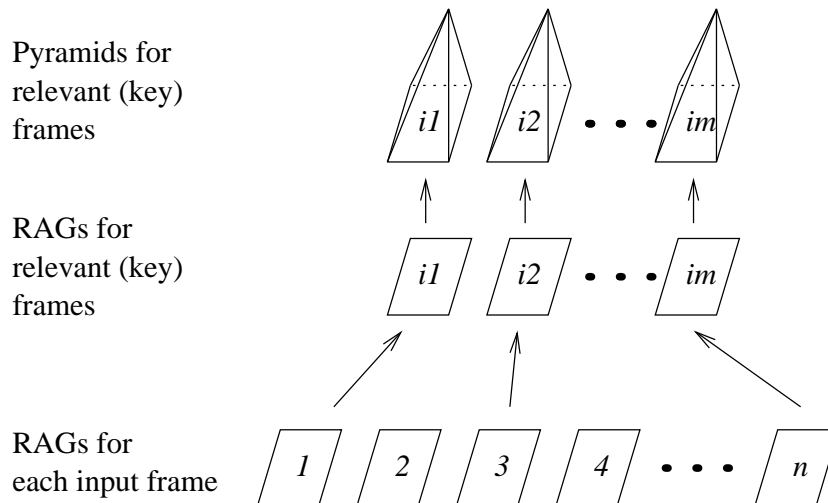


Figure 5: Time space contraction.

of region adjacency graphs. Each vertex in a region adjacency graph should be connected with the vertex(vertices), from the two neighboring graphs, that represent the same object/visual entity.

With each new level added to the time contraction pyramid, the number of topological states decreases. After reducing the number of topological states, a contraction of topological information for each state can be considered (at this level the detail regarding the background should not be important any more).

There are 2 ways that can be considered for doing this:

- contract each state independently (create a pyramid for each of the topological states at the top level of the time-contraction pyramid)
- contract all the graphs together (allow contraction kernels to span along more than one state graph)

### 3.3 Spatiotemporal Entities

The trajectories of (moving) objects (visual entities resulted from segmentation and tracked through the whole time span) represent curves connecting start, end and the junction points. Junction points happen when two (or

more) trajectories touch or cross each other, which we interpret as the ‘interaction’ of two objects.

Following the work of Kropatsch [11] the trajectory, which is a curve in 3D, and the cells, which are vertices of the graph, can be related as follows:

0-cell - an empty cell (no trajectory motion within the receptive field)

1-cell - the trajectory starts or ends in this cell (it leaves or enters the cell and intersects only *once* the boundary of the receptive field)

2-cell - the trajectory crosses the receptive field (it intersect *twice* the boundary of the receptive field).

\*-cell - a cell where more than one trajectory meet, a junction cell (the boundaries of the receptive field are intersected *more than twice*).

1-edge - trajectory intersects the connected segment boundary of the receptive field.

0-edge - no trajectory intersect the boundary of the receptive field.

It is assumed that: 1) the cells are consistent, i.e. if a trajectory crosses a boundary both cells adjacent to this boundary are in correct classes, and 2) all trajectories are well distinguishable in the base, e.g. there are no more than one single curve in one single cell of the base (except at \*-cells).

### 3.4 Selection of Contraction Kernels

Contraction should be done along the trajectory, like in curve pyramids in 2D [11, 5]. In order to undertake the contraction process, the contraction kernels must be selected. The selection rules are 1-cells and \*-cells must always survive. \*-cells are not allowed to have children. This prevents the area of unclear information<sup>2</sup> from growing. Branches of contraction kernels follow the trajectory if possible and are selected in following order: 1-cells, 2-cells, 0-cells. Receptive fields are merged as follows:

1. A 1-cell can merge with its adjacent 2-cells, then with any adjacent 0-cell and will become an 1-cell again;

---

<sup>2</sup>trajectory may intersect or may be just close to each other

2. a 2-cell can merge with both adjacent 2-cells or with any adjacent 0-cell and remains a 2-cell;
3. a 0-cell can merge with any adjacent cell and remains a 0 cell if it is merged with another 0-cell.

If the rules do not determine the contraction kernels the random selection methods [12, 10, 9] are applied. Applying these rules, the trajectory remains a simply connected curve in spatiotemporal space. At the top level (where no more contraction is possible) we find only 1-cells and \*-cells giving on overview of all movements, when and where is started, when and where the cup was grasped, and this is compact for all types.

## 4 Example

A simple, human language like description of a scene with two cups and a yellow ball is shown in Table 1. Even though the frame numbers are given, they are only for orientation purposes and can be easily eliminated from the description by putting the adverb for example 'next', 'after that', 'then' etc. The "cell type" field corresponds to the type of the scene history cell (see Section 3.3). As expected, humans tend to skip details and mention only relevant interaction of objects, this is why most of the cells have type '\*'. The previous description would be represented in the following way (see Figure 6) in the resulting top level of both approaches. The initial configuration contains 3 objects: 2 cups and 1 ball. So we initialize the objects structure with the following: *cup(1)*, *ball* and *cup(2)*. (The numerical ids in parenthesis are present to distinguish the two cups, identification could be done in many other ways. Also in the same interest, vertices are numbered to identify different positions in time.) Vertex(0) in Figure 6 depicts the initial configuration. The next vertices and edges are as follows:

1. action *move*: creates object *hand* and adds *vertex(1)*;
2. action *grasp*: links to objects *hand* and *cup(1)* and adds *vertex(2)*
3. action *lift*: links to objects *hand* and *cup(1)* and adds *vertex(3)*
4. action *move*: links to objects *hand* and *cup(1)* and adds *vertex(4)*
5. action *move*: links to objects *hand* and *cup(1)* and adds *vertex(5)*
6. action *release*: links to objects *hand* and *cup(1)* and adds *vertex(6)*

Table 1: Scene description.

cell type	Frame	Description	cell type	Frame	Description
0	16–21:	hand from left	*	92– 93:	lifts the cup, no ball visible on table
*	22:	grasps left cup	*	94–98:	shows inside, no ball inside cup
*	27–30:	moves it over ball	*	99–103:	deposit cup on its original place
*	31:	releases cup	*	104:	releases it
*	32:	grasp same cup (again)	*	105–106:	moves to the right cup
*	33–36:	shifts it to the left	*	107:	grasps it
*	37:	releases cup	*	108–113:	shifts it in front of other cup
*	38–40:	moves to right cup	*	114:	releases it
*	41:	grasps right cup	*	115–117:	moves to cup in the back
*	42–58:	shifts right cup in front of left cup (hiding left cup Fr 46–54) to the left of the original cup	*	118–121:	shifts it with fingers to the right
*	58:	releases cup	*	122:	releases cup and moves to other cup
*	59–61:	moves to the other cup	*	123:	grasps it
*	62:	grasps it	*	124–128:	shifts it to the left
*	63–65:	shifts it to the right	*	129:	releases it
*	66:	releases it	*	130:	re-grasps it
*	67–69:	moves to the left (most) cup	*	131–136:	shifts it to the right
*	70:	grasps it	*	137–139:	releases it
*	71–74:	shifts it to the right (but still to the left of the right cup)	*	140:	grasps it again
*	75:	releases it	*	141:	lifts cup
*	76–77:	moves to the right cup	*	142–150:	ball becomes visible, rolls to right until touching right
*	84:	grasps it	*	146:	cup bounces back and the hand drops cup in left position
*	85–86:	moves it to the right (but left of the right cup)	*	151:	releases cup
*	87–90:	releases it and moves up and down	*	152:	hand removes to the left
*	91:	grasps the same cup again	0	153–159:	no more movement

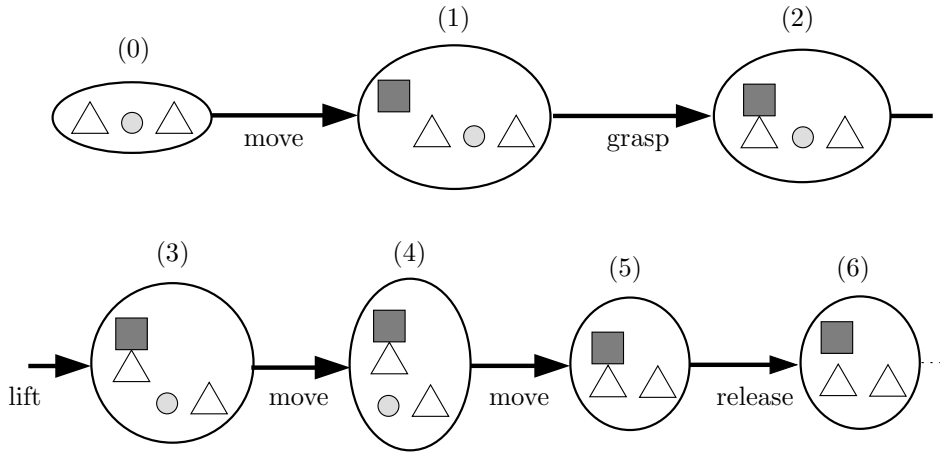


Figure 6: Example history graph.  $\square$  Hand,  $\circ$  Ball,  $\triangle$  Cup.

Although the presented approaches would work in a different way (one would first try to identify the important visual entities and then key events, while the other would start with the key events and then continue with key entities), the expected result is the same.

## 5 Human Descriptions

To motivate the research in the direction of qualitative spatial representation and reasoning, we have conducted a small set of experiments focusing on human description of videos. As a first step, 7 students (mother tongue German, descriptions made in German) were shown 2 videos (“two cups” and “yellow ball”) containing 2 identical black cups, a ball, a table (support for the cups and ball), and a hand that acts only on the cups by changing their position (on the table by pushing/shifting and in the air by picking up and holding). The 2 videos are approximately 15 and 31 seconds long.

A description of the experiments is as follows: each of the students were given a piece of paper (size A5) and told that 2 videos will be shown to them, which they should describe. After watching each video, a limited amount of time was given to describe it. No other clues were given. Of course, one can say, that seeing the hand hiding the ball using one of the cups is enough for a human (knowing the game) to focus on the ball. Which is most probably true, and can be seen on the produced descriptions. But this just enforces

the hypothesis that humans focus on a given task and do not give attention to details not related do it.

One of the first things that should be mentioned is that there were two constraints for the descriptions: one intended, which was the time allowed to write the description (2 of the descriptions are not finished), and the other one noticed, the space available on the paper for the description (more than 70% of the descriptions use up all the allocated half of the A5 paper). We can easily associate the 2 constraints with allowed processing time and available memory, and notice that humans do very well on adapting to them.

Now, getting to the descriptions themselves. Except one, all of the participants have produced narrative descriptions, with very short sentences of the form *object action direction/position*, focusing on the movement of the cups and on the position of the ball at the end of the videos. The remaining participant has used a bullet-ed list with subsections and very schematic description.

All descriptions follow a 2 section pattern<sup>3</sup>:

1. - initial configuration: contains the 3 objects initially visible, 2 cups {Tasse, Schale}, 1 ball {Ball}, and for 8 descriptions also the spatial arrangement using words like “left” {links}, “right” {rechts}, and “center” {Mitte}
2. - actions: short sentences of the form *object action direction/position* using “left cup”, “right cup”, “ball”, and “hand” {Hand} to identify objects, a whole variety of verbs for actions (e.g. “move” {bewegen}, see Table 2 for more examples) and expressions to identify positions (e.g. “between two *objects*” {zwischen}, see Table 3 for more examples).

“left” and “right” are the relational/positional words with the highest appearance (about 25 times each), followed by center/middle (less than 10 times).

One of the descriptions refers to one of the cups as “the cup with the ball” for all the time the cup is hiding the ball. On the other hand, the bullet-ed list descriptions (2, made by the same person) refer to all objects as “object” (colors are used at the initial configuration description, but only there). They contain no more than 3 actions to describe all the changes in the video and

---

<sup>3</sup>Words are given in the form English {German translation 1, German translation 2, etc.}

Table 2: Action description.

English	German
pick up	aufholen, aufheben
raise	aufheben, heben, hochheben
move	bewegen
shift	schieben
slide	schieben

Table 3: Positional description.

English	German
over	auf, über
behind	hinten
place of	auf den Platz
between two	zwischen
left	links
right	rechts
center	mitte

clearly state the final outcome. There is a description containing the wrong result, and some contain interesting hypothesis like the diameter of the ball in centimeters and the gender of the person that the hand belongs to.

After looking at the descriptions, the main observations are that:

- all the participants focused on the “implicit” problem statement (where is the ball?) and most of them basically ignored the hand;
- objects that cannot be identified easily by aspect are referred to using positions relative to the scene limits or relative to other objects;
- if the result of a position change is an interaction with another object, then this is used to describe the action, if not, then the final position is used and described relatively to the scene limits or relative to other objects using qualitative measures (left, right, front, middle, etc.);





Figure 7: The first frame of the videos shown to subjects.

- the descriptions focus on interaction between the objects, that could be considered relevant for the task.

## 6 The Experiments

The original samples of the experiments are given in the Figures 8–14. All the description are hand written, in German. The line in all the samples separates the description of the two videos (“Two-cups”<sup>4</sup> and “yellow-ball”<sup>5</sup>), since the subject knew at the beginning that there are two videos to describe. The “yellow-ball” video contains more complex events. Note that “not finished” has been put by us and the white rectangles cover a company name. A larger set of description experiments is planned and it would be interesting to make some of them in different languages.

## 7 Conclusion

This report presents a concept relating story-board description of video sequences with spatio-temporal hierarchies build by local contraction processes of spatio-temporal relations. Since object trajectories are connected curves we identify their ends and junctions and their connections as the high level descriptions. Junction points happen when two (or more) trajectories touch or cross each other, which we interpret as the ‘interaction’ of two objects.

---

<sup>4</sup><http://www.prip.tuwien.ac.at/Research/FSPCogVis/Videos/Two-Cups.mpg>

<sup>5</sup><http://www.prip.tuwien.ac.at/Research/FSPCogVis/Videos/Yellow-Ball.mpg>

We propose to derive them similar to the curve pyramid in 2D [11, 5], For the implementation we plan to use the concept of combinatorial pyramids in 3D [3, 4]. As motivation for the approach, a set of human video description experiments and their results are presented.

## References

- [1] J. Allen. An Interval-based Representation of Temporal Knowledge. In *Proc. 7th Inter. Joint Conf. on AI*, pages 221–226, 1981.
- [2] N. I. Balder. *Temporal Scene Analysis: Conceptual Descriptions of Object Movements*. PhD thesis, University of Toronto, Canada, 1975.
- [3] L. Brun and W. G. Kropatsch. The Construction of Pyramids with Combinatorial Maps. Technical Report PRIP-TR-63, Institute f. Computer Aided Automation 183/2, Pattern Recognition and Image Processing Group, TU Wien, Austria, 2000. Also available through <http://www.prip.tuwien.ac.at/ftp/pub/publications/trs/tr63.ps.gz>.
- [4] L. Brun and W. G. Kropatsch. Introduction to Combinatorial Pyramids. In G. Bertrand, A. Imiya, and R. Klette, editors, *Digital and Image Geometry*, pages 108–128. Springer, Berlin, Heidelberg, 2001.
- [5] M. Burge and W. G. Kropatsch. A Minimal Line Property Preserving Representation of Line Images. *Computing, Devoted Issue on Image Processing*, 62:355–368, 1999.
- [6] A. Chella, M. Frixione, and S. Gaglio. Understanding Dynamic Scenes. *Artificial intelligence*, 123:89–132, 2000.
- [7] B. Clarke. A Calculus of Individuals Based on Connection. *Notre Dame Journal of Formal Logic*, 23(3):204–218, 1981.
- [8] B. Clarke. Individuals and Points. *Notre Dame Journal of Formal Logic*, 26(1):61–75, 1985.
- [9] Y. Haxhimusa, R. Glantz, and W. G. Kropatsch. Constructing Stochastic Pyramids by MIDES - Maximal Independent Directed Edge Set. In E. Hancock and M. Vento, editors, *4th IAPR-TC15 Workshop on Graph-based Representation in Pattern Recognition*, volume 2726 of *Lecture*

*Notes in Computer Science*, pages 35–46, York, UK, June-July 2003. Springer, Berlin Heidelberg, New York.

- [10] Y. Haxhimusa, R. Glantz, M. Saib, G. Langs, and W. G. Kropatsch. Logarithmic Tapering Graph Pyramid. In L. van Gool, editor, *Proceedings of 24th DAGM Symposium*, pages 117–124, Swiss, 2002. Springer Verlag LNCS 2449.
- [11] W. G. Kropatsch. Property Preserving Hierarchical Graph Transformation. In C. Arricelli, L. Cordella, and G. Sanniti di Baja, editors, *Advances in Visual Form Analysis*, pages 340–349, Singapore, 1998. World Scientific.
- [12] P. Meer. Stochastic image pyramids. *Computer Vision, Graphics, and Image Processing*, 45(3):269–294, March 1989.
- [13] D. Randell, Z. Cui, and A. Cohn. A Spatial Logic Based on Regions and Connection. In *Proc. 3rd Intern. Conf. on Knowledge Representation and Reasoning*, pages 165–176. Morgan Kaufmann, 1992.
- [14] R. Sablatnig. Increasing flexibility for automatic visual inspection: the general analysis graph. *Mach. Vision Appl.*, 12(4):158–169, 2000.
- [15] D. Vernon, editor. *A Research Roadmap of Cognitive Vision (DRAFT Version 3.2)*. ECVision: The European Research Network for Cognitive Computer Vision Systems, 2004.

Video 1:

3 Objekte auf einem Hintergrund (= 4 Objekt)

3 Situationen:

- 1) ~~3~~ 3 schwarze Obj. sichtbar
- 2) 2 schw. Obj. sichtbar  
weil 1 Obj. das kleinere verdeckt
- 3) 1 schw. Obj.  
verdeckt das 2 Obj.  
selber Größe

4 Objekte:

- + 1 3 schwarze Obj.  
↳ 2 gl. Größe
- + 1 Hintergrund
- + 1 gelb
- + 1 #

---

Video 2:

4 Objekte:

- 2 schwarze gleiche Größe
- 1 kleineres gelbes
- 1 Hintergrund

Situationen:

- 1) schw. Objekt 1 verdeckt Kugel
- 2) schw. Objekt 1 tauscht Position  
bzw. Seite mit schw. Obj. 2
- 3) gelbe Objekt wieder sichtbar in ≠ Position

CA 606590 - 794 - 10 M Gedruckt auf Recyclingpapier

Figure 8: Subject 1.

Zwei schwarze Tassen auf weißem Grund. In der Mitte ein Ball. Eine Hand hebt die rechte Tasse auf den Ball, schiebt die linke in einem Halbkreis hinter der ersten Tasse vorbei. Schließlich wird die Tasse mit dem Ball hochgehoben und auf den Platz der linken gestellt → Ausgangsaufstellung.

---

Selbe Konfiguration. Linke Tasse wird über den diesmal gelben Ball gestellt, andere Tasse vor die erste geschoben. Sie tauschen die Seiten, und werden dann in Raupen ähnlichen Kriechverfahren nach und nach nach rechts verschoben. Anschließend wird die leere Tasse hochgehoben → Enttäuschung, dann wird die andere Tasse gehoben, Ball rollt heraus.

Figure 9: Subject 2.



Grauer Hintergrund

2 Tassen 1 rote Kugel (zwischen den 2 Tassen)

Paar Sekunden nichts dann kommt von rechts eine rechte Hand und nimmt rechte Tasse → über Ball, nimmt anschließend linke Tasse und verschiebt diese hinter die Tasse mit dem Ball, die 2 Tassen werden dann wieder an die ursprüngliche Position gebracht

~~Ball~~

2 Tassen 1 gelber Ball

Hand von ~~rechts~~ <sup>links</sup>, linke Tasse über Ball, rechte Tasse ~~links~~ um linke Tasse bewegt, beide Tassen sukzessive nach rechts bewegt; zunächst rechte etwas nach rechts dann linke, rechte Tasse aufgehoben und Inhalt (nichts) gezeigt linke Tasse aufgehoben Ball rollt raus und bleibt zwischen den 2 Tassen.

Figure 10: Subject 3.

2 Tassen, <sup>(schwarz)</sup> Kugel (schwarz)  
Kugel in der Mitte, links und rechts davon  
die beiden Tassen, beide umgedreht - alles  
vertikal zentriert.  
Hand von rechts kommend hebt rechte Tasse  
auf und stellt sie über Kugel, danach linke  
Tasse hinter andere, setzt in der Mitte  
befindliche, verschieben, Hand verschwindet  
nach links aus dem Bild

gelbe Kugel in der Mitte, schwarze Tassen (wie oben)  
links und rechts, Hand von links hebt  
linke Tasse auf und stellt sie über  
Kugel, zieht dann diese Tasse mit Kugel  
nach links, rechte Tasse im Halbkreis vor  
der anderen Tasse entlang bewegt <sup>ein</sup> nach links  
von Tasse mit Kugel, beide Tassen ein  
Stück nach rechts, linke Tasse (ohne Kugel)  
aufgehoben und umgedreht, beide Tassen  
stich.

(not finished!)

Figure 11: Subject 4.

Schale links, Ball in der Mitte, Schale rechts  
Rechte Schale wird über Ball gelegt  
linke Schale in Viertelkreis hinter rechte Schale  
gehoben, Hand darauf gelassen  
ebenso rechte Schale ~~in ger~~ (jetzt vordere) in gerader  
Bewegung nach rechts verschoben  
  
~~linke Schale hintere Schale nach vorne geschoben~~

---

Schale links, Ball Mitte, Schale rechts  
linke Schale (LS) über Ball gelegt  
~~re~~ rechte Schale (RS) in Halbkreis vor LS  
nach links bewegt  
in mehreren Schritten abwechselnd beide Schalen nach  
rechts verschoben, ~~linke LS anheben~~ ~~hier~~  
RS anheben, darunter schauen  
Schalen ~~sch~~ schiebend vertauschen  
LS anheben, Ball ist zu sehen, wieder darüber  
legen

CA 606590 - 794 - 10 M Gedruckt auf Recyclingpapier

Figure 12: Subject 5.



• Zu sehen sind 2 schwarze Keramikgefäße und eine graue Kugel vor hellgrünem Grund. Die rechte Tasse wird durch eine Hand angehoben und auf / über die Kugel gestülpt. Anschließend werden die Gefäße verschoben.

---

• Zu sehen ist die Szene aus Video 1, jedoch aus anderer Perspektive, mit Farbe und Sound.

Diesmal wird die linke Tasse über die (jetzt) gelbe Kugel gestülpt und verschoben. Zuletzt wird die Kugel wieder aufgedeckt.

Figure 13: Subject 6.

2 schwarze Tassen auf jenem Kinderpunkt  
 Tassen mit Öffnung nach unten auf Ebene  
 Aus den Tassen 1 roten Ball ca 3cm im Durchmesser  
 Hand benutzt im Bild, rote Tasse wird über  
 Ball gehalten, linke Tasse in 2-Richtung hinter  
 die 1. Tasse durch andere Tasse weiter  
 an linken Platz. Vorder-Tasse werden unter  
 Platz

2 Tasse, 1 gelber (Bild) (Bild)  
 linke Tasse durch Hand auf linken Ball.  
 linke Tasse werden von Tische, rote Tasse  
 von der linken Tasse (2-Richtung). ~~Hand~~  
 1. Tasse wird auf die linke Seite umgedreht.  
 beide Tassen in die Tische, ~~die Tasse~~ werden  
 Tasse, Ball, Tasse weiter im Hunger -  
 (not finished!)

Figure 14: Subject 7.