

# Real-Time Range Acquisition by Adaptive Structured Light

Thomas P. Koninckx and Luc Van Gool

**Abstract**—The goal of this paper is to provide a “self-adaptive” system for real-time range acquisition. Reconstructions are based on a single frame structured light illumination. Instead of using generic, static coding that is supposed to work under all circumstances, system adaptation is proposed. This occurs on-the-fly and renders the system more robust against instant scene variability and creates suitable patterns at startup. A continuous trade-off between speed and quality is made. A weighted combination of different coding cues—based upon pattern color, geometry, and tracking—yields a robust way to solve the correspondence problem. The individual coding cues are automatically adapted within a considered family of patterns. The weights to combine them are based on the average consistency with the result within a small time-window. The integration itself is done by reformulating the problem as a graph cut. Also, the camera-projector configuration is taken into account for generating the projection patterns. The correctness of the range maps is not guaranteed, but an estimation of the uncertainty is provided for each part of the reconstruction. Our prototype is implemented using unmodified consumer hardware only and, therefore, is cheap. Frame rates vary between 10 and 25 fps, dependent on scene complexity.

**Index Terms**—Imaging geometry, depth cues, range data, shape, real-time systems.

## 1 INTRODUCTION

### 1.1 Rationale

THE last decade has seen tremendous progress on the ease of noncontact optical surface digitization [1]. Traditional techniques such as stereo rigs (e.g., [2], [3]) and laser scanners (e.g., [4], [5]) got competition from more flexible and scalable alternatives [6]. Nevertheless, several challenges remain. First, most processes are still quite slow. This makes the production of 3D models still relatively expensive. For the same reason, the *detailed* capturing of dynamic scenes is largely an open issue. Second, 3D capturing techniques tend to fail on certain types of surfaces. As a matter of fact, one method could fail where another would be particularly successful. For instance, a specific surface texture may interfere with certain types of structured light patterns but would leave other patterns quite unaffected and would even support stereo. But typically, one method is used at a time, with the net result that for quite a few objects no model can be produced automatically. All too often still, one has to resort to adapting the scene to the system, e.g., by powdering the objects to be scanned.

We propose a scanner which inverts this paradigm in that it adapts itself to the scene. It seeks to generate better patterns online by taking the properties of scene and setup into account. Most classical approaches, in contrast, use static patterns which are designed beforehand. The adaptivity will improve the robustness of both pattern detection and (de)coding. In addition to online adaptations, the camera-projector configuration—determined during an off-line calibration—is explicitly taken into account. In case of a

quasistatic scene, the scanner behaves as if it were equipped with a static yet dedicated projection pattern, as the adaptations will only occur at startup now. The proposed approach allows operation in real-time and each reconstruction is based on only a single input image. As such, the system can deal with moving and deforming objects. As a matter of fact, a high speed of 3D extraction is an important ally in making the online adaptations possible, as the latest 3D data will contain strong information about what is to be expected in the next snapshot.

Next to the internal feedback (i.e., adapted patterns), high-speed operation also enables direct feedback to the user. By providing an almost immediate visualization of the 3D results, the user can judge the partially reconstructed surface and steer the scanning process. This is especially useful for the detection of “holes” and for filling them in by manipulating the object into a more appropriate pose. As such, human insight solves the otherwise difficult problem of view-planning. Such feedback is already successfully applied in the scanner developed by Rusinkiewicz et al. [7]. In contradistinction to that scanner, the system proposed here is “one-shot” and adds the pattern adaptivity, the importance of which we will demonstrate. The one-shot character has the advantage of allowing for faster object motions, and the adaptivity enables the system to deal with more strongly textured surfaces.

### 1.2 Related Work

Structured light range acquisition comes in many flavors. For an overview, see, e.g., [6] and the concise taxonomy of Table 1. We classify approaches based on their capabilities for real-time operation—which allows for direct (visual) feedback—and whether they are one-shot or multishot based. Only the former methods can deal with deforming or fast-moving objects. A system can use dense codes—each individual pixel, up to the resolution used, is directly related to its position in the projection pattern—or a sparse set of landmarks. Pattern coding also typically implies the

• The authors are with the Katholieke Universiteit Leuven, ESAT-PSI, Kasteelpark Arenberg 10, B-3001 Leuven-Heverlee, Belgium.  
E-mail: {tkoninck, luc.vangoool}@esat.kuleuven.be.

Manuscript received 30 Nov. 2004; revised 4 Apr. 2005; accepted 11 July 2005; published online 13 Jan. 2006.

Recommended for acceptance by H. Sawhney.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-0639-1104.

TABLE 1  
A Taxonomy of Related Active Structured Light Methods with Their Primary Capabilities and Codification Strategies Listed

	real-time	one-shot	dense coding	color coding	tracking. & time	spatial connect.
stripe boundary codes [7]	y	n	y	n	y	0
time sequential [8]	n	n	y	n	n	0
color coded slots [9]	n	y	y	y	n	1D
coded grid [10]	n	y	n	n	n	2D
adapt. color coding [11]	n	n	y	y	n	1D
1-shot color code [12]	y	y	y	y	n	1D
slit coded lines [13]	n	y	n	n	n	2D
coded chessboard [14]	n	y	y	n	n	2D
color coded dots [15]	n	y	y	y	n	2D
optimal time seq. [16]	n	n	y	n	n	0
dyn.prog.+space time [17]	n	y	y	y	y	1D
relative labeling [18]	n	y	n	n	n	2D
real-time adapt. struct. light	y	y	n	y	y	2D

use of color, markers, pattern geometry, and so forth, and can integrate information over a 1D or 2D spatial neighborhood. Integration over time, in case of moving objects, asks for pattern tracking. Strictly time coded systems use a series of subsequent projections, with different but predefined patterns. See, e.g., the early work by Altschuler et al. [8]. The acquisition process is typically speeded up by limiting the number of necessary input images. The limit case is one-shot acquisition [9], [10]. In order to still solve for the correspondence between projected pattern elements and their image projections, the pattern elements should be coded using only this single input image. Most researchers use a spatial neighborhood that yields a unique identifier for the relative position in the pattern. Examples are the use of colors [12], [17], or a set of markers with a unique shape or configuration [13], [14], or a combination of both [15]. In the seminal work of Caspi et al. [11], noise margins are taken into account for deciding which colors can be used for coding. Related to this, Horn and Kiryati [16] proposed a technique to design an optimal sequential structured light pattern of length K. The main difference with our work is that adaptations no longer happen offline but also online, and that 3D is acquired from a single projection rather than a time series.

A weak assumption of temporal continuity underlies our adaptation, as we will continually optimize a series of parameters for subsequent frames. The higher the frame rate, the more valid this assumption becomes. In contrast to most other approaches, we won't use very dense coding where every pattern element would have a unique identifier. For one-shot acquisition, dense coding almost implies the use of color coding (e.g. [9], [17]), which is fragile in colored scenes, or the extraction of rather extensive feature configurations [14], which is time consuming. On the other hand, very sparse coding together with a relative labeling approach that relies on global smoothness and strict monotonicity [18] cannot robustly deal with occlusions. Instead of relying on a single "coding cue,"

we exploit a weighted combination of different such cues. These multiple sparse encodings result in a rather dense, overall encoding. This is further combined with relative labeling. A graph cut algorithm leads to optimal correspondences, together with local confidence measures. In contrast to the work by Zhang et al. [17], where the pattern is decoded line by line, our strategy goes directly after its 2D interconnections. Mutual influences between lines and the omission of image rectification also render the algorithm considerably faster and less sensitive to noise. We explicitly focus on the niche of one-shot acquisition in real-time of moving scenery. Note that in case the scene is quasistatic and/or extended processing can be afforded, offline time sequential or space time approaches (see, e.g., the challenging results in [19]) might be considered instead. We can achieve a typical framerate of 20 hz on a P4 at 2.26 Ghz, processing an input video of  $640 \times 480$  pixels. In this work, focus is on solving the correspondence problem, based on a more extended version of the codification technique initially outlined in [22]. The solution is computed based on the graph-cut solver (compare) developed by Boykov and Kolmogorov [20]. Pattern detection and refinement is described in earlier work. (See [21]). In [23], our system was used for interactive modeling. For extended (implementation) details on this work, consult [24].

The rest of the paper is organized as follows: Section 2 describes the labeling problem. In Section 3, the pattern coding and graph cut-based integration of coding cues is outlined. Results are shown in Section 4, and Section 5 concludes the paper.

## 2 A LABELING PROBLEM

Range computations are based on optical ray-plane triangulation (see Fig. 1). The basic pattern we use for triangulation is a set of black parallel stripes on a white background, in the remainder of the paper referred to as the "base pattern." Every stripe boundary delimits a planar

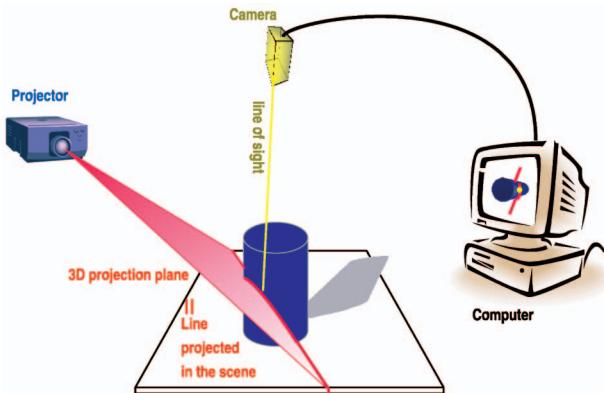


Fig. 1. Basic optical triangulation: After pattern recognition, the projection planes delimited by the stripe boundaries and the “projector center” are intersected with the corresponding lines of sight.

surface in space, a so-called “triangulation plane.” The reflection on the object of this projection pattern is identified in the camera image. In what follows, we will refer to such a detected boundary simply as a “stripe”  $S$ . Given a recognized stripe and the knowledge of the plane it originates from, the intersection with the line of sight for a pixel yields the corresponding 3D position.

The purpose of the projection pattern in a single frame acquisition process is twofold. First, it should be easy to detect precisely and as such support triangulation. Second, it should make it possible to establish correspondences between the camera image and the projector pattern in order to distinguish the different triangulation planes. Here, we focus on the latter issue, called “pattern coding” or “pattern codification”; the former is described elsewhere [21].

Suffice it to say that the black-and-white pattern yields a maximal contrast ratio and by adapting pattern resolution and orientation (see below), its detectability can be optimized even further.

The correspondence problem boils down to determining a labeling at every time  $t$ :

$$f^t : S_j^t \rightarrow L^t \text{ with } L^t = [0, n_{\max}], j \in \{0, m\} \quad (1)$$

in which  $L$  is the set of all possible labels  $n$  for the current projection pattern. The correct localization of a stripe  $S$  in the image together with its label  $n$  in the pattern are needed for its 3D reconstruction. The integer  $j$  is the identifier of a stripe  $S$  within the set of  $m$  detected stripes in the image. Discontinuities, textures, etc., let a single triangulation plane cause several detected stripes  $S$ . Thus,  $f^t$  is a *many to one* relationship and possibly  $m \geq n_{\max}$ . As a consequence, the complexity of the scene is reflected in the runtime of the algorithm and, strictly speaking, the timing of the algorithm is indeterministic. Note that the superscript  $t$  will be omitted at times, when the current frame is referred to.

Establishing the correct image-projection pattern correspondences is complicated by (see also Fig. 2):

- A. **Ghost boundaries:** Scene texture and geometry cause spurious stripes (false boundaries, which do not result from the pattern).
- B. **Horizontal occlusions:** A depth discontinuity parallel to the stripes of the pattern may cause stripes to

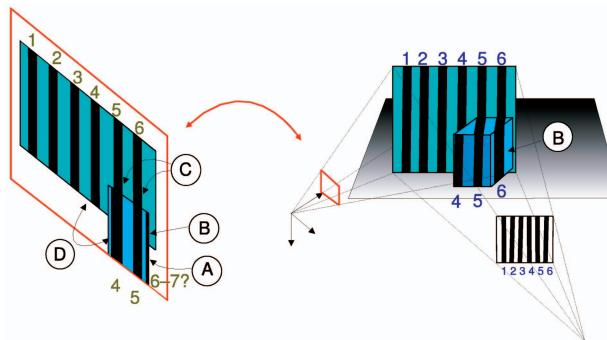


Fig. 2. Complications in labeling-based pattern projection techniques: (A) “False labels” are introduced at the cube’s edge. (B) Not all of the pattern is visible. (C) Pieces of the pattern are wrongly interconnected. (D) Monotonicity is not preserved.

disappear in the image. A labeling based on a strictly incremental neighbor relationship that crosses such a discontinuity will be inconsistent.

- C. **Vertical occlusions:** A depth discontinuity that is not parallel with the pattern can let two different stripes  $S_i$  and  $S_j$  appear as a single stripe  $S$  in the image.
- D. **Monotonicity is not preserved:** The ranking of the stripes in the projected pattern may not be preserved in their image in the proximity of depth discontinuities.

### 3 DYNAMICALLY ADAPTED PATTERNS

An overview of the processing pipeline is given in Fig. 3. The first section explains this pipeline and subsequent sections describe its different units or a combination thereof, in more detail.

#### 3.1 Online Processing Steps: Overview

The acquisition starts with projecting a pattern on the object (“project structured light”) and with taking images thereof with a consumer video camera (“video capture”). The “segmentation and pattern identification” unit extracts the positions and identity of the stripe edges in the video images. The “pattern identification” exploits different coding cues that are included in the pattern. These coding cues are fed into a graph cut algorithm yielding “optimal correspondences” between camera and projector. This enables the system to produce a “3D reconstruction,” which is presented to the user on a monitor (“visualize”). This immediate visualization could be a goal in its own right, e.g., in tele-immersive and enhanced reality systems. It can also ease the job of producing complete object models interactively (“integrate and visualize”).

This pipeline reacts upon its own results, thereby forming an autonomous closed loop system. This functionality is provided by the “adaptation and control system” unit. The quality of the stripe extraction and identification is monitored continuously. The goal is to adapt the pattern in such a way that interference by the scene is minimized, and as a consequence, the pattern detection circumstances are optimized. The “tracking” unit exploits the temporal continuity that typically exists between consecutive frames, without making the system critically dependent on it. It has a dual function. First, it enables speeding up the algorithm by limiting the search for stripes to regions of interest

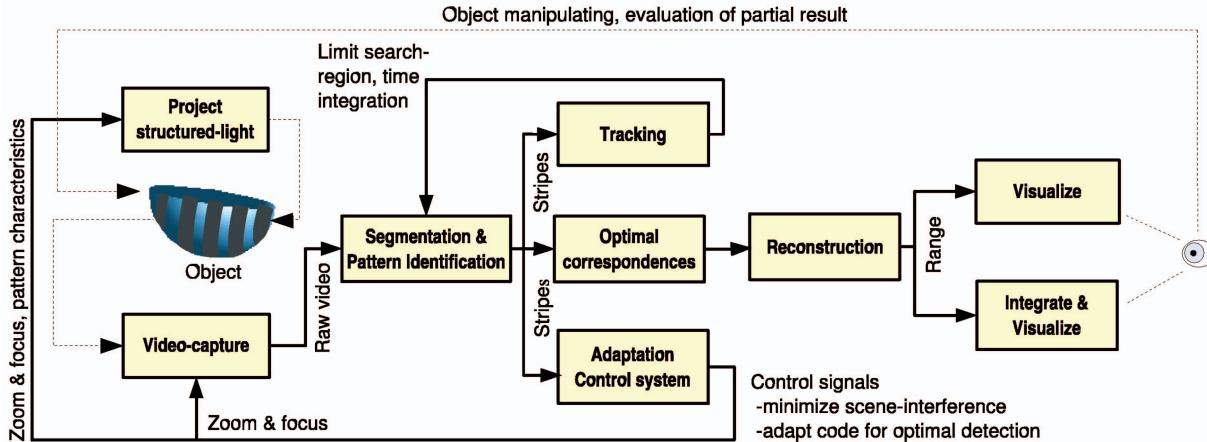


Fig. 3. Overview of the real-time adaptively coded reconstruction pipeline. The topmost recurrent arrow is the feedback from the user to the system, who can evaluate the scanning quality by the immediate visualization, the bottom feedback loop is the adaptation of the projected pattern to the scene by the algorithm.

(ROIs). Second, it allows for denser coding. In highly dynamic situations, tracking won't be possible anymore. The scanner responds in this situation by lowering the resolution of both the pattern coding and the base pattern to a level where tracking is not necessary to keep the same frame rate. In this way, resolution is exchanged for reaction time, which is visually quite acceptable, as a limited resolution is less observable for a quickly moving object than a delayed reaction or incorrect model. A correct low-resolution model, which follows the dynamic content of the scene is to be preferred over a (partially) incorrect high-resolution model, with lower responsiveness. Hence, temporal continuity affects the algorithm in a positive way, but is not necessary for it to function properly.

### 3.2 Pattern Projection, Acquisition, and Identification

Pattern projection and acquisition is based on consumer-grade hardware (see Section 4). As we work with online adapted patterns, it is important to retain a strict synchronization between camera and projector at all times. Lens distortions of both devices and color aberrations are provided for at this stage of the pipeline. The remainder of this section gives an overview of our coding strategies, which will yield the cues for pattern identification. Every cue yields labels  $n$  (1), which will be denoted  $n_{code,i,j}$ ,  $n_{color,j}$ ,  $n_{track,i,j}$ , and  $n_{trackS,j}$  for, respectively, geometric coding, color coding, tracking of code points, and tracking of stripes. The subscript indicates the "coding type," the coding line involved (if applicable, see below), and the identity of the stripe affected. An optional superscript  $t$  indicates the timestamp, and a superscript \* refers to the unknown correct label that is sought. The hat (^) notation will be used to indicate a predicted value.

#### 3.2.1 Geometric Coding: Code Points ( $x_{code}$ )

We submerge colored "coding lines" in the pattern "below" the basic stripe pattern (Fig. 4), referred to in the following as the "base pattern." If a coding line is not coincident with an epipolar line, the intersections of this coding line and the base pattern will all lie on different epipolar lines in the camera image. As a consequence, points of such a coding line that are recognized in the camera image can be

transferred to a unique point of the projection pattern by use of the epipolar geometry. This is schematically illustrated for a horizontal epipolar geometry in Fig. 4. The more this coding line is tilted with respect to the epipolar line, the lower the noise on the transfer to the projection pattern will be. However, a lower inclination allows more coding lines below each other, thereby increasing the density of such intersections or "code points" as geometric coding cues—this while still respecting the constraint that a single epipolar line should intersect only a single coding line. Fig. 5 illustrates this trade-off. This repetition of coding lines below each other is a first design parameter of the geometric code and will be called the vertical repetition  $R_v$ . It determines the inclination of the coding lines.  $R_v$  corresponds to the number of individual coding lines encountered during a vertical transition of the pattern.

We will refer to a code point in the image as  $x_{code,i,j}$ , with  $j$  the identifier of the stripe boundary on which the point resides (see (1)), and  $i \in [0, R_v]$  an identifier for its originating coding line. Another design parameter for the coding lines is the color  $C_i$  which is used. How  $R_v$  and  $C_i$  are selected and adapted, given the current scene, is explained in Section 3.5.

The decoding  $x_{code,i,j} \rightarrow n_{code,i,j}$  of a code point to the corresponding label yields a first set of cues. Each stripe  $S_j$  can receive zero, one, or multiple such labels  $n_{code}$ . Such a label solves the correspondence problem locally as it relates

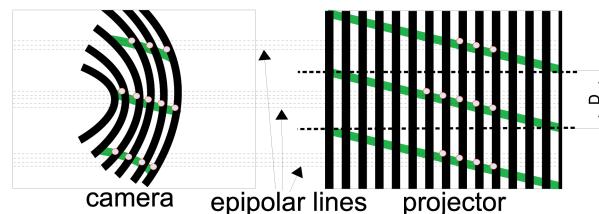


Fig. 4. Geometric coding: Colored "coding lines" underly the base stripe pattern (vertical on the right). Points of this coding line that are recognized in the camera image can be uniquely related to their corresponding point in the projection pattern based on the epipolar geometry. This resolves the ambiguities due to the periodic stripe pattern. In this schematic representation, the correspondences are indicated by dots.

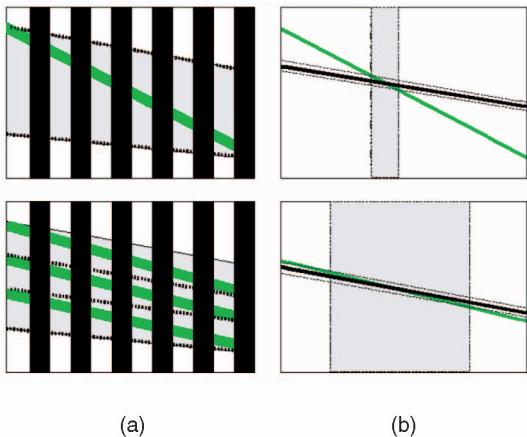


Fig. 5. (a) More coding lines result in higher coding density. (b) Uncertainty interval on the decoding for the same level of noise on the detection. The intersection becomes worse conditioned.

the location of the code point to the corresponding location in the projection pattern. It takes the form of a ranking number (the number of the originating stripe boundary) in the projection pattern.

### 3.2.2 Sparse Color Coding

We use the high-intensity black-white and white-black transitions of the base pattern as the triangulation planes. We now change the color of a limited set of slots between the stripes of the base pattern from white to yellow, magenta, or cyan. These saturated and high-intensity colors let the stripe boundary still stand out clearly. Which color  $C_s$  will be used is decided on-the-fly based on the colors in the current scene (see below). The orientation parallel to the base pattern makes it easy to discriminate between this sparse color coding and the aforementioned geometric coding. At the intersection point between the geometric coding and the sparse color coding, the latter occludes the former as the resulting coding cue is of higher quality. Allowing the geometric code to divide the sparse color code in several individual “pieces” would seriously degrade its robustness for only a limited increase in coding density.

If  $S_j$  and  $S_k$  are consecutive edges with opposite gradient values, we check the interstripe color. If above the detection threshold for  $C_s$  and isolated, we assign  $S_j$  and  $S_k$  labels  $n_{color,j}$  and  $n_{color,j} + 1 = n_{color,k}$ . Each  $S_j$  has zero or one label  $n_{color}$ . Whereas the  $n_{code}$ s are dense but quite noisy, the  $n_{color}$ s are rather sparse but mostly correct.

The detection of the colored stripes is based on a classification of the average color of the stripe in h,s,v-space. A constraint is that such a colored stripe should be locally unique. The erroneous interpretation of the color of an object in the scene as a code should be avoided. We use the same color for all coded stripes; the discussion about how to choose these colors is postponed to Section 3.5. We disambiguate between multiple color-coded stripes by taking the one which is closest to the predicted position based on the tracking, as will be explained in Section 3.4.

### 3.2.3 Tracking of Code Points ( $x_{code}$ )

After the initial detection of the code points  $x_{code,i,j}^t$ , we run the following predictor-corrector style algorithm to keep

track of these points over time, going from their positions in the previous frame to those in the current one:

- A predictor line  $l_{predictor,i}$  is fitted to each set of code points  $x_{i,j}^{t-1}$  in the previous frame that belong to the same code line  $i$  (i.e., regression over all intersections of one code line with the different stripes). Grouping of code points into such sets is a by-product of the stripe identification algorithm in [22]. As coding lines are not very oblique with respect to the epipolar lines, fitting a piecewise straight line is sufficiently accurate for our purposes. (Note: If the coding line would coincide with the epipolar line, it would show straight in the image. Given the multiple coding lines, diagonal between consecutive epipoles, the piecewise line approximation suffices. Tests with splines did not yield better results.)
- The intersections  $l_{predictor,i} \cap S_j, \forall j$  yield predicted code points  $\hat{x}_{i,j}^t$  for the current frame based on the predictor lines from the previous frame. Their positions are refined to positions  $x_{i,j}^t$  by convolving the image data around the position of  $S_j$  with the  $[5 \times 7]$  environment of the projection pattern around the code point. This matched filter approach yields the new  $x_{i,j}^t$ , where the maximum response is found. The search for this refinement is confined to a local neighborhood around the predicted positions. The latter minimizes the disturbing influence of interfering textures. Points which are “spatially isolated” are removed.
- For every  $x_{i,j}^t$ , we search the closest match  $x_{i,j}^{t-1}$  obeying the epipolar geometry between the current and the previous frame. As the same code point does not change position in the projection pattern, it lies on the same epipolar plane at all times. This means that the corresponding epipolar line connects a point that is tracked between consecutive camera frames. The motion of a stripe  $S_j$  is limited given the frame rate (between 10 and 20 fps dependent on scene complexity). This again allows the system to restrict the size of the search region. As the features that are tracked are projected, they obey the epipolar geometry between camera and projector at all times, irrespective of object deformations.

The relationship  $x_{i,j}^{t-1} \leftrightarrow x_{i,j}^t \Rightarrow n_{i,j}^{t-1} \rightarrow n_{track,i,j}^t$  allows transferring labels on time  $t - 1$  to the tracked code points on time  $t$ . This process is illustrated in Fig. 6.

### 3.2.4 Tracking of Pattern Stripes ( $S$ )

The same algorithm as was used for tracking the code points is used to track the stripes  $S_j$ . We now use the midpoint in the vertical direction of the stripe as one additional feature point for every  $S_j$ . Midpoints are checked for obeying interframe epipolar geometry and such can be tracked. This assumes that stripes in consecutive frames have similar lengths or symmetrical growth. Time continuity due to the rather high frame rates makes this assumption hold for most of the stripes during execution of the algorithm. This yields a new set of labels  $n_{trackS,j}^t$ . Although noisy on average, the label will never be far from the real label. The higher the processing speed, the more valid this becomes.

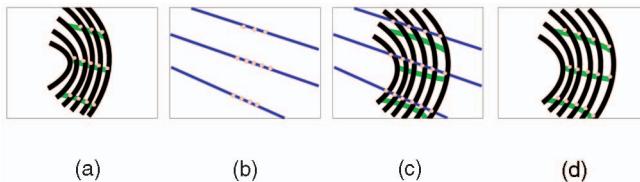


Fig. 6. Tracking code points over time: detection, curve fit (line approximate), prediction, and trimming of the predicted points.

### 3.2.5 Neighbor Relationship

After the detection of all stripes, a neighbor relationship is constructed. To do so, we create an ordered set of stripes  $Q$  based on ranking them by their position in the image. The leftmost stripe comes first; the rightmost, last. A top-down ordering is used as a secondary criterion to compare stripes of different heights.

A stripe  $S_i$  is represented as a lookup table that, for every vertical pixel position  $y$  in the image, stores the corresponding horizontal pixel position  $x$ . A stripe, however, does not necessarily exist (contain entries) for all heights  $y$ , but only in a range  $[y_{i,\min}, y_{i,\max}]$ . This let us define the  $i$ th stripe more formally as

$$S_i(y) = x, \forall y \in [y_{i,\min}, y_{i,\max}].$$

The vertical range of pixels shared between the  $i$ th and  $k$ th stripe will be called

$$\text{overlap}(S_i, S_k) = [\max(y_{i,\min}, y_{k,\min}), \min(y_{i,\max}, y_{k,\max})].$$

Any two stripes  $S_i$  and  $S_k$ , can now be compared as follows:

$$S_i \leq S_k \Rightarrow \{y_{i,\max} \leq y_{k,\min} \text{ OR } S_i(y) \leq S_k(y), \forall y \in \text{overlap}(S_i, S_k)\}.$$

During the construction of  $Q$ , it would be too expensive to compare every pixel of all stripes to each other. Rather, a limited number of samples is used. Afterward, we check for degenerate cases between consecutive entries in  $Q$ . If  $S_i(y) \leq S_k(y)$  AND  $S_i(y') \geq S_k(y')$ , both stripes are crossing or partially coincident. This is impossible, and the problem should be resolved. In such a case, the stripe with the orientation closest to the local average within a window of five stripes in  $Q$  is retained.

Given this ordered set of stripes  $Q$ , we can search for  $\text{neighbors}(S_i)$ , the *neighboring stripes* of  $S_i$ . Starting from the  $i$ th position in  $Q$ , all stripes  $S_k$  for which  $\text{overlap}(S_i, S_k) \neq 0$  AND  $\text{overlap}(S_k, S_l) = 0$ , with  $S_l \in \text{neighbors}(S_i)$ , will be added to  $\text{neighbors}(S_i)$ . We continue this search until all neighbors cover the y-range of the original stripe for 90 percent, or have checked more than 200 entries in  $Q$ . This is done for every stripe. Given the order in  $Q$  and the fact that many stripes have similar lengths, this step is executed very fast. Also during this transition, the average spacing in horizontal direction between the pixels on the stripes is computed (called  $dist$ ), together with the standard deviation  $\sigma$  of this distribution.

The number of pixels on two neighboring stripes, of which the horizontal distance resides within the  $dist + 2\sigma$  tolerance interval, is a quality measure of their neighbor relationship that expresses how “parallel” they

are. The latter will be indicated by  $\text{par}(S_i, S_k)$ . Finally, let  $\text{length}(S_i) = (y_{i,\max} - y_{i,\min})$  be the vertical length.

This enables weighting of every neighbor relationship by a quality factor, which is designed to stress stripes that are quasiparallel over most of their lengths and which favors the longer pair of stripes:

$$Nb(S_i, S_k) = \frac{\text{par}(S_i, S_k)}{\text{overlap}(S_i, S_k)} \times \frac{(\text{length}(S_i) + \text{length}(S_k))/2}{\text{imageheight}}. \quad (2)$$

The first factor of (2) is the number of pixels that lie within the tolerance interval, relative to the maximum that would be possible. This expresses the quality of the overlap (perfectly parallel stripes are likely to result from consecutive boundaries in the pattern). The next factor renders long overlapping stripes relatively more important.

This relationship gives us an additional cue for assigning labels, and makes it possible to “spread out” the other cues over regions that went uncoded so far:  $S_k \in \text{neighbors}(S_i) \Rightarrow n_k = n_i + 1$ .

## 3.3 Optimal Correspondences

### 3.3.1 Problem Statement

The labeling from the different cues of the previous section will almost never be fully consistent right away. Increasing the quality of the coding in order to render the problem better conditioned can only be achieved at the cost of a decrease in coding density. Furthermore, relative labeling of neighbors will always be noisy, but remains a prerequisite to enable single frame range acquisition. Combining all cues into a single consistent labeling poses an extra workload but is necessary to achieve robustness.

We propose to search for the *weighted least squares* solution to this problem. This can be considered as an optimal approximation of all labeling cues, and as such will yield optimal camera-projector correspondences. This solution is not limited to the cues in this paper. Any additional set of cues, regardless of density or noise, can be integrated in a straightforward way.

We seek to solve the following minimization problem:

$$\begin{aligned} \text{Min} \left[ \sum_j \left\{ \sum_i \alpha |n_j^* - n_{\text{code},i,j}|^2 + \beta |n_j^* - n_{\text{color},j}|^2 + \right. \right. \\ \left. \left. \sum_i \gamma |n_j^* - n_{\text{track},i,j}|^2 + \delta |n_j^* - n_{\text{trackS},j}|^2 + \right. \right. \\ \left. \left. \epsilon \sum_k (Nb(S_j, S_k) \times (n_j^* + 1 - n_k^*)^2) \right\} \right], \end{aligned} \quad (3)$$

in which  $j$  and  $k$  run over all stripes  $S$ , and  $i$  over all coding lines. The labels  $n_j^*$  form the solution we look for. The first four terms express that the weighted and squared differences between the resulting labeling  $n_j^*$  and the corresponding cue  $n_{x,j}$ —respectively, geometric coding, sparse color coding, tracking of geometric code points, and stripe midpoints—should be minimized. The last term introduces the neighbor relationship between the stripes.

### 3.3.2 Cast into a Graph Cut

Given we aim at an algorithm with real-time capabilities, we need a time and memory efficient way to solve (3). To do

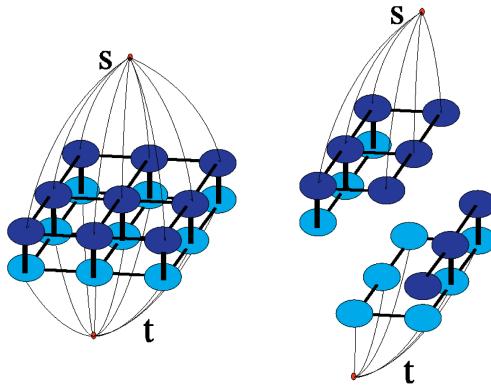


Fig. 7. A graph cut divides a set of nodes into two disjoint sets  $S$  and  $T$ .

so, we cast the problem into the search for a minimal cut on a graph.

A weighted graph  $G = \langle V, E \rangle$  consists of a set of nodes  $V$  and a set of weighted edges  $E$  to connect them. Our graph also contains some additional special nodes called terminals, of which one is the source  $s$  and one is the sink  $t$ . The nodes  $V$  form chains that connect  $s$  and  $t$ . If the nodes of a chain are also sideways connected we get a structure as visualized in Fig. 7. A binary cut on this graph divides the set of nodes  $V$  into two subsets of which one contains  $s$  and one  $t$ . We will generate such a cut, so that the cutting surface minimizes the total sum of the weights of all edges which are broken. This cut will then be seen to provide the solution of (3). If the two endpoints of every chain are connected to, respectively,  $s$  and  $t$ , we know that every chain will be cut. The problem of searching this binary cut can be solved very efficiently (see Boykov and Kolmogorov [20]).

Now, we turn back to the problem of (3). We will associate with every stripe  $S_i$  a chain of nodes (vertical in Fig. 7) that interconnects  $s$  and  $t$ . The length of every chain equals  $n_{max}$ , the number of all possible labels  $L$  in (1). Every edge of this chain represents a possible label for the corresponding stripe. The edges are initialized with unity weights. At the position of the label predicted by a cue  $n_x$  of (3), the edge is weakened by the corresponding  $\alpha, \beta, \gamma$ , or  $\delta$  cost. As a matter of fact, instead of only weakening an edge locally, we place a Gaussian around each weakened edge and let it influence a local environment on the chain. The size of the Gaussian corresponds to the uncertainty on the specific cue (see below). As every chain is connected to  $s$  and  $t$ , each stripe will be given a label, as the stripe's chain certainly will be cut. See also Fig. 8.

The  $\epsilon$ -term in (3), corresponding to the neighbor relationship, can be taken into account by a sideways connection of the nodes of chains that correspond to neighbors. These edges run diagonally through the node structure and encourage neighbors to have consecutive labels. See Figs. 8 and 9.

We want our solution to be unique, which means that we want every chain to be cut only once. We can achieve this by avoiding the possibility that the cutting surface can "fold back" in the node structure. See Fig. 9. This is achieved by putting a constraint on the minimal edge weight along a chain versus the cost of breaking an edge which connects two chains. Define  $\min(A)$  as  $(1 - \alpha - \beta - \gamma - \delta)$ , which corresponds to the minimal edge cost that can occur along a

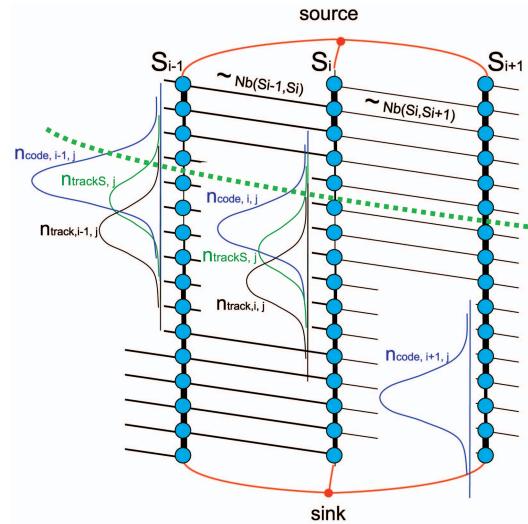


Fig. 8. Three neighboring chains. The different coding cues influence a common region on the different chains, which will attract the cutting surface. Individual outliers are easily regularized as the cost for the cutting surface to reach this point is too high due to the sideways interconnections.

chain.  $B$  corresponds to the cost of cutting the edge between two chains  $j$  and  $k$ :  $(\epsilon \times Nb(S_j, S_k)(n_j^* + 1 - n_k^*)^2)$ . As can be seen in Fig. 9, a cutting surface could turn back to avoid an expensive edge. The largest gain would occur when an edge with cost  $\max(B)$  could be avoided, at the minimum cost of  $2\min(A) + \min(B)$  (with  $\min(B)$  and  $\max(B)$ , respective of the global min and max of all "B-costs"). This leads to the following constraint on the edge weights if this is not to happen (i.e., no folding back):

$$2\min(A) > \max(B) - \min(B). \quad (4)$$

The graph cut is speeded up considerably by limiting the length of each chain to a window around the solution predicted by a much faster prelabeling. For this prelabeling, we assign labels recursively based on the neighbor relationship and the geometric coding only. At this stage, no efforts are made to overcome inconsistencies in the labeling. Every stripe is labeled based upon the closest cue in terms of neighbors or the median of the coding cues available on the

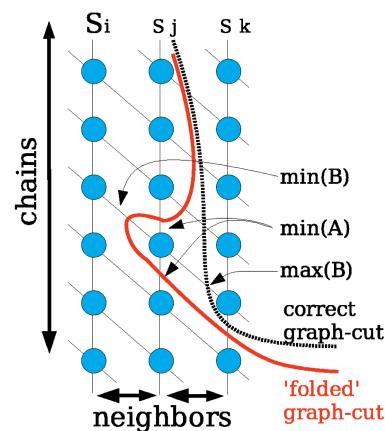


Fig. 9. A graph cut can fold back at a minimal cost of  $2\min(A) + \min(B)$  to avoid an edge which has a maximum cost of  $\max(B)$ . In order to have a unique solution for every stripe  $S$  this folding should be avoided.

stripe itself. To avoid assigning labels based on outliers of the set of cues, only these cues are used that respect the neighbor relationships up to an acceptable deviation. The latter is based on the noise level measured on the prelabeling—before versus after the graph cut—in the previous frame. The search window around this predicted labeling is adapted to be always larger than twice the average mislabeling on this prediction over the last 10 frames. A minimum window of 20 labels takes into account a possibly large and unpredicted error during this prelabeling.

### 3.3.3 Inconsistent Solutions—Uncertainty on the Solution

The solution for the labeling problem solves most but not all of the complications listed in the introduction of Section 2: False stripes (A), missing stripes (B), and nonmonotonicity (D) are dealt with appropriately. If enough cues point to a labeling that is different from the one predicted by the neighbor relationship, the latter will be automatically ignored. In the case of false stripes, one ends up with several stripes on the same scan-line with the same label. This is recognized and the stripe which deviates most from the local average orientation is removed.

However, no solution for the wrongly interconnected stripes which can occur around vertical occlusions has been proposed so far.

To this end, we will interpret the correspondence between the different coding cues and the labeling assigned by the graph cut as a local *consistency* of the solution. Define the residual for stripe  $S_j$  on cue  $k$ , in which  $k \in [0, 4]$  refers to the cues described in Section 3.2, as  $r_{j,k} = n_{k,j} - n_j^*$ . With  $r_j$  we refer to the total residual for  $S_j$ , weighted by the corresponding  $\alpha, \beta, \gamma, \delta$  and  $\epsilon$  factor:  $r_j = \sum_k (\zeta_k * |r_{j,k}|)$ ,  $\zeta \in [\alpha, \beta, \gamma, \delta, \epsilon]$ . Now, we design the consistency measure for  $S_j$  so that it rapidly decreases with the size of the residual, and with the sparseness of the stripe's coding:

$$\text{consistency}(S_j) = \frac{e^{-r_j^2}}{e^{\max Cues - \#cues}}. \quad (5)$$

$\max Cues$  corresponds to the maximum number of cues a stripe can receive:  $\max Cues = [2 \times R_v + 3]$  (geometric coding + code tracking, color coding, stripe tracking, neighbor relationship),  $\#cues$  reflects the actual number of cues received. This consistency is a measure of the local uncertainty on the labeling—and, thus, the 3D reconstruction—assigned by the graph cut. The more a label of a stripe respects the cues available, the more the consistency will approach 1. The denominator ensures that a stripe which is more heavily coded receives a higher weight.

We can now decide when to switch to a conflict resolution step by placing a threshold on the consistency. The resolution mechanism itself is based on a secondary graph which is constructed. Every stripe is no longer assigned a single chain of nodes, but each stripe is divided in multiple small segments (at distances equal to the mean distance between consecutive stripes) and is assigned a set of nodes. The solution of this graph effectively allows us to split stripes in the longitudinal direction, as multiple nodes now can originate from the same stripe. This way, it can resolve the above error (situation (C) in Fig. 2). The chain

lengths can be reduced to the range of labels assigned during the primary graph. By combining costs of edges into longitudinal direction, one can come to a graph with chainlength two, which only expresses the binary decision whether or no to split a chain locally. The short chain lengths, and limited size render this graph solvable in very limited time. This two-step, iterated approach with a global graph ignoring a possible problem, and a smaller but more performant secondary graph to fix this, keeps the problem tractable, given the tough time constraints.

### 3.4 Tracking

Tracking of code points was used earlier as a coding cue. Now, we will build further upon this result to speed up the “segmentation and pattern identification” of the next pattern. Note that the outcome is a mere speedup, and provides no additional coding information; feedback goes to the detection stage only (see middle recurrent arrow in Fig. 3) The tracked coding points define interframe translation vectors over time. These multiple translations are combined into a single average translation, which is fed continuously into a Kalman filter. The state vector comprises the position and speed of this average vector (linear motion model for the center of mass of the ROI). The predicted translation is used to update an ROI for every frame, in which we will look after the new pattern. For the ROI, an extended version is used of the part of the image for which we did find the pattern in the previous frame. This extended ROI allows us to use a very limited motion model. The result is a net speedup of the algorithm. Note that this mechanism should be disabled—or extended—in the presence of multiple, independently moving objects. A possible extension is to initialize different trackers for each region, defined by motion segmentation. In this work, this multiobject tracking however is not implemented yet. The Kalman tracker also allows one to distinguish static from dynamic scenes.

The labeling of the previous frame together with the overall translation vector in the image yield a prediction of the label of every stripe in the current image. This allows for denser geometric coding, without an increase in the decoding noise. Referring again to Fig. 4: If a single epipolar line is allowed to intersect multiple coding lines of the projected pattern, the code density can be increased. This is done by moving the coding lines closer to each other in the vertical direction, without changing the tilt. In the example of Fig. 4 this means making the vertical distance between coding lines smaller than  $D$ . As said, the ambiguity because of the multiple intersections can be resolved based on the label predicted for this position by the tracking unit. The intersection closest to the one predicted is chosen. We call the number of intersections a single epipolar line can have with coding lines the horizontal repetition  $R_h$ . This is a second design parameter of the geometrical coding. We can summarize that an increase of  $R_v$  results in a *denser but noisier code*, where an increase of  $R_h$  yields a *denser code* with identical noise level on the decoding but *requires tracking* of the pattern. The higher  $R_h$ , the less noise on the tracking can be tolerated. See Section 4 for an evaluation of code density versus noise level.

### 3.5 Pattern Adaptations and Cue Weighting

The ability to decode and identify the pattern is of primary importance to our system. As we achieve rather high frame rates (10-20 fps) we exploit time continuity to gradually change pattern settings to make the codes and pattern as detectable as possible.

The geometric code has  $R_v$  and  $R_h$  as free design parameters.  $R_h$  is chosen as high as possible as long as the average residual on the stripe tracking satisfies  $\frac{|r_0|}{2 \times R_h} < n_{max}$ . In this case, it will be normally no problem to resolve the ambiguity of multiple coding line intersections.  $R_v$ , the number of times coding lines are repeated below each other, is chosen so as to achieve high coding densities. The constraint here is that the residual on the coding remains limited:  $|r_0| < T$ , in which  $T$  is an arbitrary threshold. (In our tests, we took for  $T$  half the inconsistency threshold.) In our experiments, we use a low  $R_v$  to bootstrap the system and gradually increase  $R_v$ . The color of the coding lines is chosen to minimize interference from scene colors that are similar. The implementation is based on binning scene colors, and such selecting an appropriate color for the coding. The width of the base pattern is adapted in order to always have a stripe width after detection of approximately three pixels. Note that the base pattern, however, is always uniform and equidistant. This assures visibility and keeps the detector matched to the pattern. For details on the implementation of the adaptations, see [22].

For the location of the color-coded stripes, the total residual  $r_j$  for each stripe is used. As this coding normally gives a high-quality coding cue, it is used to overcome problems in a region which is highly uncertain. As the detector also splits stripes based on a sudden color gradient in the longitudinal direction, this also is a preemptive resolution for the next frame of a possible conflict in this area (see the previous section). If this stripe crosses a depth discontinuity and is likely to be wrongly connected to another stripe, it won't happen anymore if one of them is assigned a color. The number of color-coded slots is chosen  $\leq R_h$ , as this defines the minimal distance in the image needed to discriminate between different coding lines. It is implicitly determined by the discriminant power of the tracking. The color  $C_s$  is chosen in a similar way as for the geometric code. By binning the colors in the scene into groups on the color-disc ( $h,s,v$ -space), we can decide which colors are more likely to interfere.

A last set of parameters which should be determined are the relative weights  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ , and  $\epsilon$  for the different cues in expression (3). The cues for  $k \in [0, 3]$  are noisy measurements of the unknown correct label  $n_j^*$ . Noise on these measurements is modeled as a normal distribution with mean 0 and standard deviation  $\sigma_k$ :

$$n_{j,k} \sim n_{j,k}^* + N_k, \quad N_k \sim \mathcal{N}(0, \sigma_k^2).$$

The optimal estimation of  $n_j^*$  given the measurements  $n_{j,k}$  is the one that maximizes the posterior probability:

$$\prod_{jk} P(n_j^*, \sigma_k | n_{j,k}) = \prod_{jk} \frac{P(n_{j,k} | n_j^*, \sigma_k) P(n_j^*, \sigma_k)}{P(n_{j,k})} \quad (6)$$

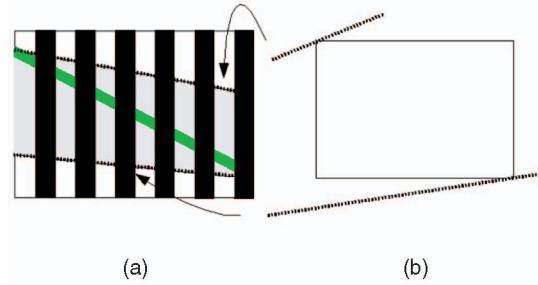


Fig. 10. A code with a single line, which will fall within the image.

in which the denominator can be seen as a normalization factor. (Summation of the numerator over all  $(n_j^*, \sigma_k)$ ). Given the noise model on  $n_{j,k}$ , the data likelihood becomes:

$$P(n_{j,k} | n_j^*, \sigma_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2} \frac{|n_{j,k} - n_j^*|^2}{\sigma_k^2}\right).$$

The prior  $P(n_j^*, \sigma_k)$  in (6) is taken to be  $P(n_j^*) \times \text{Const}$  as we have no prior knowledge about  $\sigma_k$ .  $P(n_j^*)$  can be modeled as an exponential distribution  $\exp(-R(n_j^*)/\lambda)$ .  $\lambda$  defines the width of this distribution and  $R(n_j^*)$  is a regularizer which expresses that the labeling should be continuous over the stripe set  $Q$ . This is exactly what is expressed by the expression  $(Nb(S_j, S_k) \times (n_j^* + 1 - n_k^*)^2)$  from (3).

Instead of maximizing the probability of (6), we minimize its negative logarithm. This leads to

$$E = \sum_k -\log(\sqrt{2\pi\sigma_k^2}) + \frac{1}{2} \sum_{jk} \frac{\|n_{j,k} - n_j^*\|^2}{\sigma_k^2} + \frac{1}{\lambda} \sum_{jk} Nb(S_j, S_k) \times (n_j^* + 1 - n_k^*)^2. \quad (7)$$

If we compare (7) with (3), we see that for the relative weights of the cues  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ , the variation  $\sigma_k$  on the corresponding residuals  $r_k$  can be used as a valid steering parameter. The weight  $\epsilon$  or  $1/\lambda$  is left as a parameter to the user and expresses her/his belief in the correctness of the relative labeling. For continuous surfaces it can be taken higher than for complex irregular shapes. In practical experiments, we limited the adaptations of the weights to vary within a window around an initialization value, to avoid drifting away to trivial solutions.

### 3.6 Reconstruction, Visualization, and Integration

A dual-head GPU (nVidia Geforce4-TI) is used for visualization and pattern projection. The primary head gives the visual feedback to the user about the acquired 3D data via the point-based or textured surface model, which it sends to a monitor. The second head is steering the projector. The base and code patterns are composed and rendered on the video board. The CPU only sends a vector description of the pattern down to the rendering pipeline. The GPU is also used for correcting distortions and can do part of the preprocessing [21], [22].

A reconstruction algorithm, working with the labeled set of stripes [21] and an algorithm integrating the different

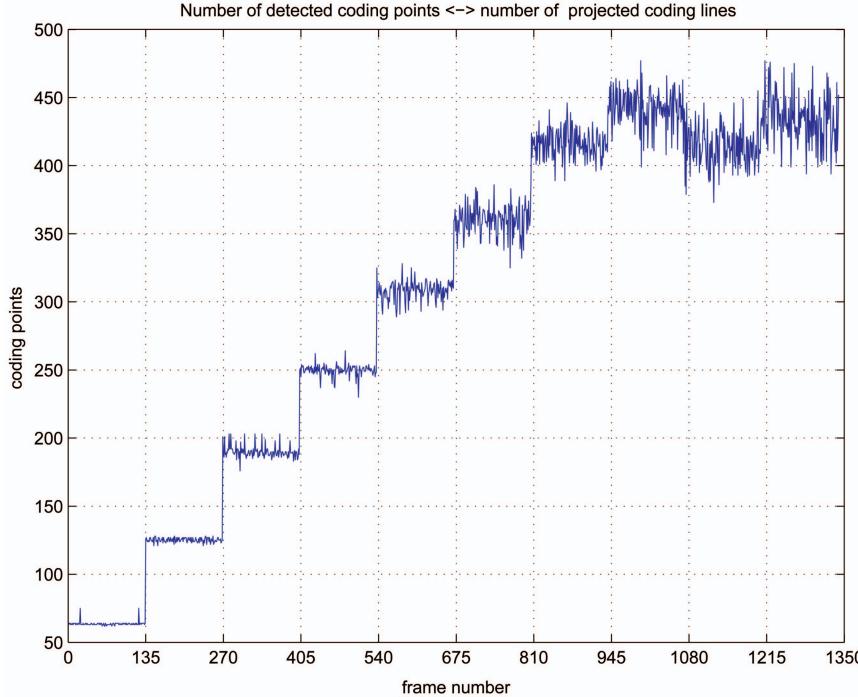


Fig. 11. Number of detected code points increases with the number of code lines. A saturation effect starts around  $R_v = 8$ .

partial reconstructions (i.e., online registration) [23] run simultaneously on the same machine or can be executed by different computers. In this case, a small server at the side of the scanner provides multicast broadcasting of the range data over the network. At the side of the integration, a client application gets its data from the network. A multi-threaded implementation enables benefit from multiple CPUs as well.

Working distances vary from approximately 75 cm to 5 m, and once a suitable distance was chosen, the working volume is determined by the depth of field of the projector.

The camera and projector are mounted together on a small transportable jig. This makes it possible to move the device around without losing the calibration. The projection axis and viewing axis are only about 25 cm apart and make an angle of less than 10 degrees. This short baseline

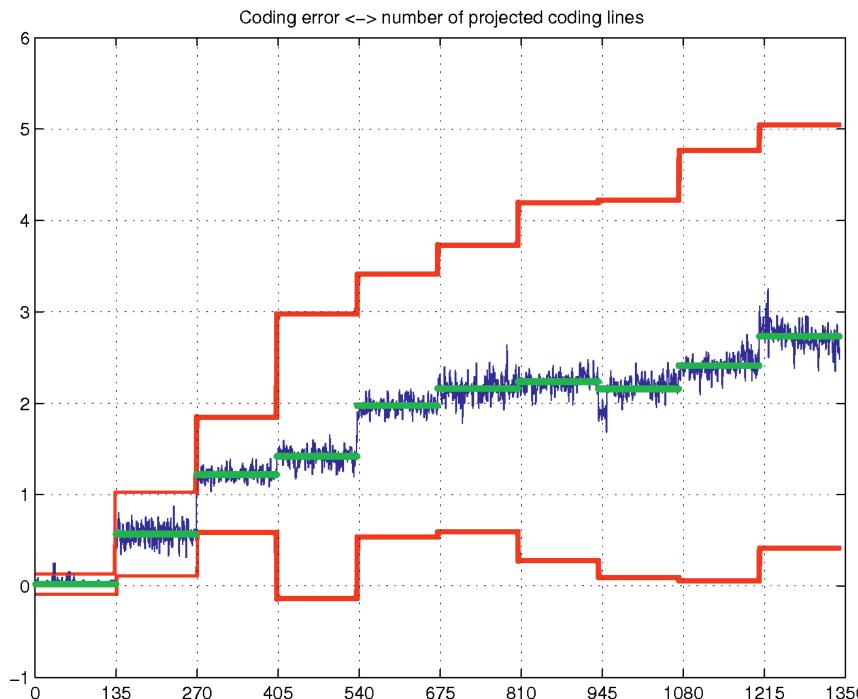


Fig. 12. Denser codes ( $R_v$  adapted) are more affected by decoding noise. Average deviations and their  $2\sigma$  interval are shown.

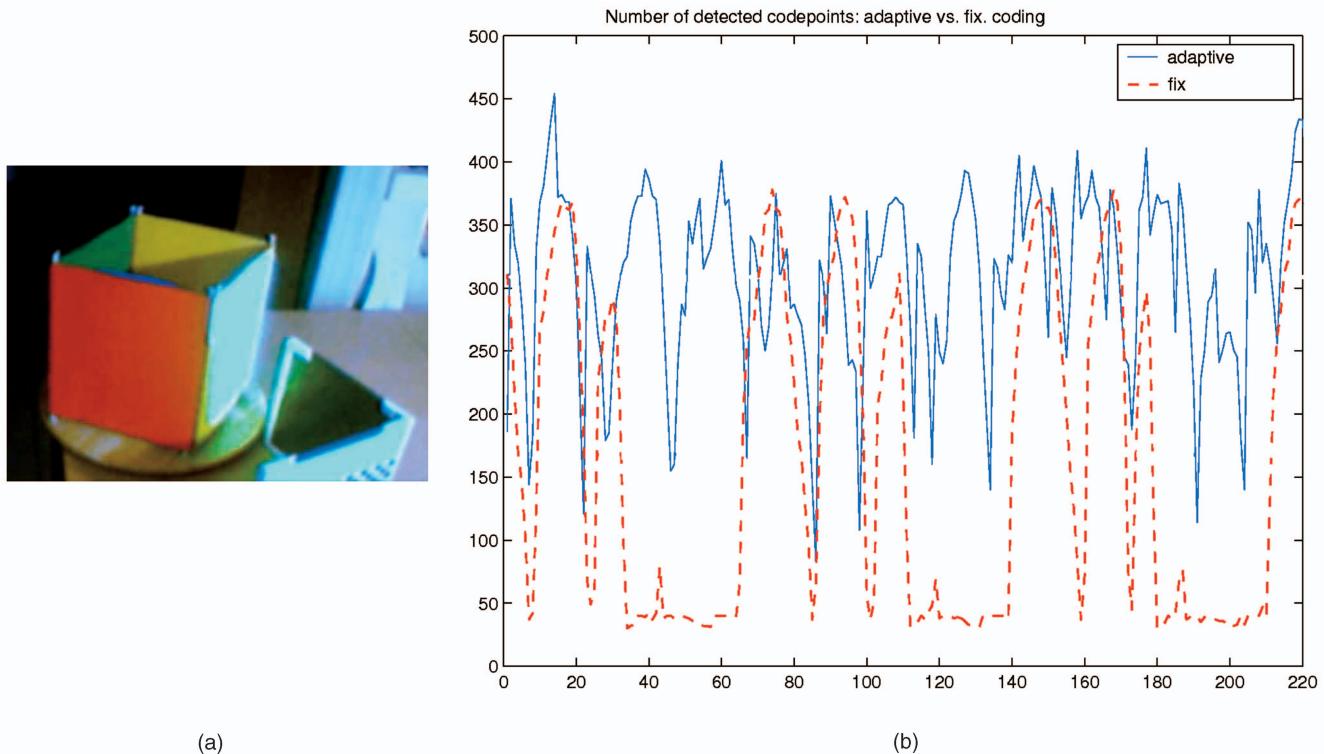


Fig. 13. (a) The scene is composed of a rotating cube with colored sides. (b) Instant number of reconstructed 3D points for a sequence of 220 frames with and without adaptation.

renders the setup compact and minimizes the part of the surface that is not visible to the projector and the camera simultaneously.

### 3.7 Pattern Outline and Setup Calibration

As the setup is assumed to remain fixed, the influence we want it to have on the pattern can be taken care of in advance and offline. Mainly, the geometric coding can benefit from this. By choosing the code line(s) diagonal in the pattern of Fig. 4, we probably put part of the code-line outside of the image. A solution is to generate the code the other way around. This is shown in Fig. 10. The two epipolar lines that correspond to the opposite corners of the image, define a cutout of the projected pattern that is actually visible (see Fig. 10a). A code line that diagonally spans this cutout would be completely visible in the image and would therefore make a better choice. This is also the code line that will have the biggest tilt with respect to the corresponding epipolar line, and as such the best conditioned intersection or decoding. For higher coding densities, this still holds and the tilt will be adapted as described earlier.

Also, the system calibration is done only once and offline. The internal and external calibration matrices of both the camera and projector are needed to decode and reconstruct. As mentioned before, radial distortions of both devices are provided for. The delays between the devices and in the graphical subsystem of the computer are also estimated beforehand and should be taken into account to keep camera and projector synchronized. All of this is integrated in an easy calibration tool, which enables setup of the system in less than a minute. That is important to make reconfiguration of the setup easy and fast. A radiometric color calibration of both devices and an estimate of

the color-crosstalk between both devices was added during recent development. It turned out to make the decoding based on sparse color coding more robust. It is, however, not used for generating the results in this paper.

## 4 RESULTS

This section will first demonstrate the relation between coding density, decoding noise and processing speed, crucial for the pattern adaptations. Next, we demonstrate some of the adaptations. Afterward, substeps, results, and limitations of the algorithm are shown for a complex object. Finally more challenging scenes, and fast movements are captured.

Fig. 11 shows how an increase of code density ( $R_v + 1$  after every 135th frame) yields more detected code points in the input image. At  $R_v = 8$ , saturation sets in because points lie too close to each other for the prediction and trim

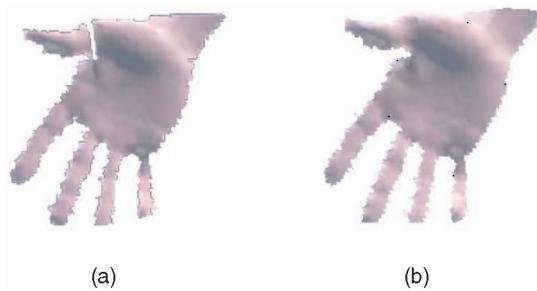


Fig. 14. (a) Wrongly detected base pattern. A stripe crosses an occlusion. A major error in the reconstruction results. (b) Conflict resolution.

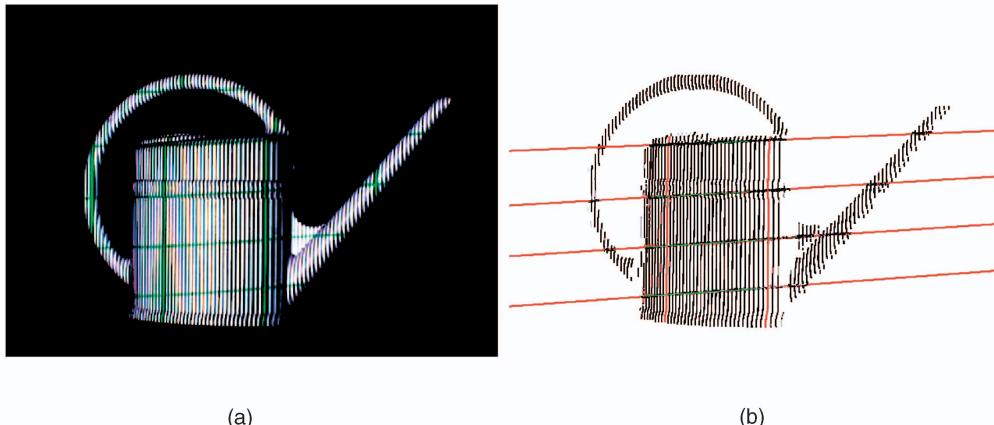


Fig. 15. (a) Input image with structured illumination. (b) Regression lines for code point tracking and trimmed predictions.

mechanism to work properly. As more code points also result in a more computational intensive operation for each frame, the overall frame rate will steadily drop. (In this experiment, the frame rate decreases from approximately 14 fps to 18 fps). As was discussed earlier, an increase in code density (higher  $R_v$  for equal  $R_h$ ) results in a higher noise level on the decoding. Fig. 12 shows the average deviation of all code points for a single frame, and the average for every set of 135 frames (central thicker line) from the label that was assigned by the graph cut. The

$2\sigma$  uncertainty interval (top and bottom “staircase” profile) around this average deviation shows that a high code density based only on  $R_h$  will make it hard to distinguish true discontinuities from noise. This illustrates the need for a combined  $R_v$  and  $R_h$  increase together with tracking.

The positive influences outlined above, on the decoding noise, density, and speed already indirectly motivate code adaptations. To demonstrate the effect more clearly, an experiment with and without color-and-resolution adaptations is performed. The setup was chosen to be dynamic and strongly colored, making it a highly unsupportive scene for color coding, let alone for a sequential approach. To make the results visually clearer, the setup was chosen so that, for quite some frames, the codes cannot be retrieved at all and the reconstruction disappears completely. In a more mixed scene, the result would be less extreme, and the relative labeling would have to fill in more of the missing data. This typically results in more noise and possible errors around occlusions (see below). Fig. 13 shows the setup and the total number of reconstructed points for every frame with and without pattern adaptations of a sequence recorded with the cube rotating. The plot shows that more 3D data

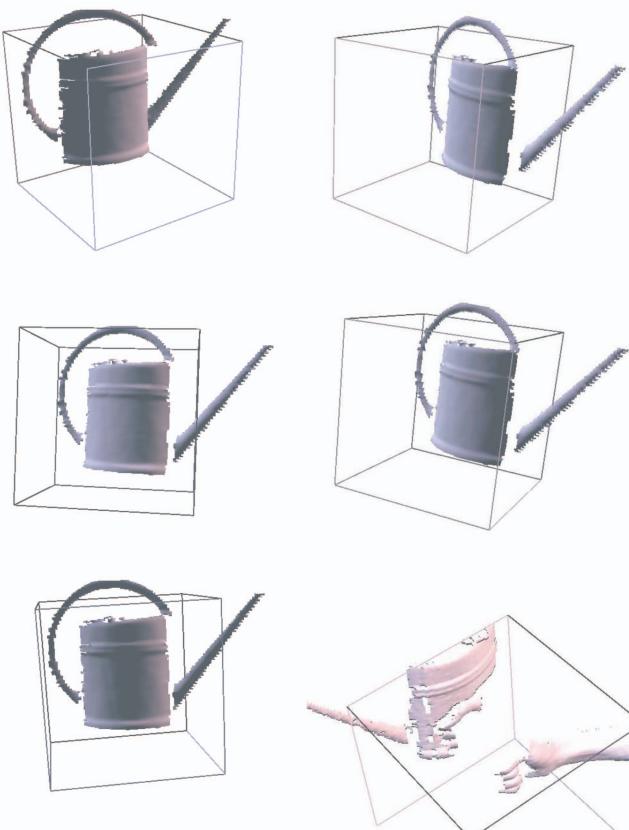


Fig. 16. Rendering of the watering can from a set of new camera viewpoints. Bottom right: The graph cut can correctly handle the discontinuity, but coding is not dense enough to decide everywhere on background versus foreground (see, e.g., the thumb).

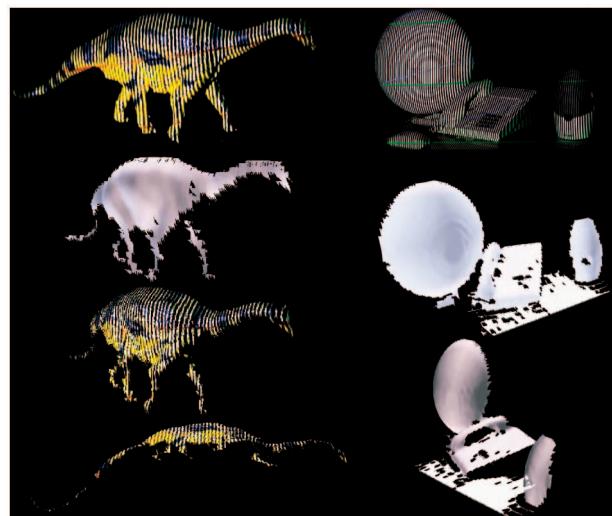


Fig. 17. Left: Input and reconstruction of a strongly textured plastic dinosaur figurine. Right: A more complex scene with more occlusions.

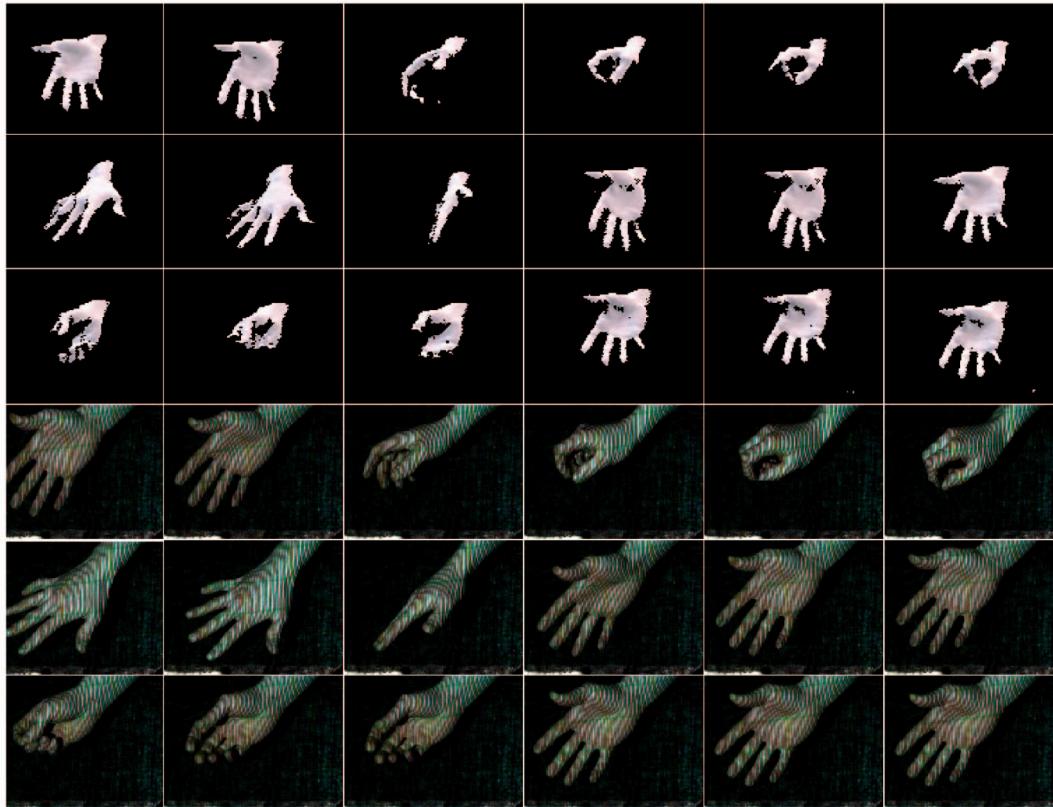


Fig. 18. Top: A series of reconstructions for a moving and deforming hand. Every fifth frame is shown. Bottom: The corresponding input images. Note that the camera viewpoint for the reconstruction is moved slightly toward the bottom left, in order to visualize jitter and/or flaws.

could be retrieved from the same scene, both due to the resolution adaptations of the base pattern and the more complete coding. Also, the number of code points retrieved on a per frame basis is now higher, resulting in a denser overall coding which is less affected by noise.

Fig. 14 shows how a labeling inconsistency results in a major flaw in the 3D reconstruction, if no action is taken. A single stripe is wrongly detected and “crosses” the occlusion around the thumb. Fig. 14a shows an intermediate reconstruction without the resolution mechanism for the inconsistency applied. Fig. 14b shows the result after the secondary graph, which split the stripe crossing the occlusion and so solved the error.

Fig. 15 shows how the regression lines resulting from the previous frame predict the location of the geometric code in the current frame. Regions of the object that are more densely coded will result in a lower uncertainty.

Fig. 16 shows the resulting reconstruction from different camera viewpoints. The bottom right subfigure demonstrates one of the limitations of the algorithm. An arm (the object in the front of the working volume) is now inserted between the watering can and the scanner and a reconstruction is made. The depth discontinuity between the foreground and the background object is picked up correctly. However, some parts of the foreground (e.g., the thumb of the hand) are wrongly assigned to the background, and as such receive an incorrect depth value. This is caused by the sparse coding: the thumb is not coded explicitly and receives its label via the neighbor relationship and graph cut. When no local coding information is available, the graph cut splits foreground and background

mainly based on edge information (the graph edges are weighted by the corresponding gradient magnitude in the images). Although this will still result in the correct solution in most cases, it is not guaranteed.

Fig. 17 shows two additional examples. The dinosaur figurine (approximately 12 cm in length), is highly textured. The code line switched color to gain robustness. Given the strong texture and high degree of shininess of the material, the reconstruction seems quite acceptable. In the bottom view, the relatively high depth range within this statue can be seen (see, e.g., the curved tail). In the same figure, a model of a more complex desk scene is shown as well. The reconstruction shows a desk with a ball next to a phone and a speaker. Fig. 18 shows a sequence of shots taken of a moving and deforming hand. The camera viewpoint is moved toward the bottom left to show the differences with the input images. Our pipeline was able to process this sequence at approximately 15 fps.

The test on the hand sequence (see movies, which can be found at <http://computer.org/tpami/archives.htm>, and Fig. 18) was run on a single desktop computer with a Pentium 4 CPU at 2.26 Ghz. An average frame rate of 14.5 fps resulted, with approximately 10,000 points/frame. Typical frame rates for the system vary between 10 and 25 fps, with a typical output between 100,000 and 500,000 points/sec.

## 5 CONCLUSION AND FUTURE WORK

We propose an approach for high speed 3D acquisition based on consumer grade projection and video devices. To solve the correspondence problem between camera and projector, a combination of geometric coding, color coding

and tracking over time is used. The correspondences are distributed over the detected pattern based on relative labeling.

In order to combine the different not fully consistent cues, a least squares solution is implemented by using a graph cut.

The design parameters for the pattern are adapted online to make it more robust against current noise levels on the decoding and to ease the detection. This results in a strategy for online adaptively coded structured light. Relative weights of the different cues are data determined, and adjusted over time. Based on the residual on the data approximation, we can assign a local consistency measure to the reconstruction.

Future work will focus on more elaborate pattern adaptations and further optimizations for speed.

## ACKNOWLEDGMENTS

The authors gratefully acknowledge support by the KULeuven Research Council (GOA-project "Marvel"), and the Flemish Institute for the Advancement of Science in Industry (IWT).

## REFERENCES

- [1] F. Blais, "Review of 20 Years of Range Sensor Development," *J. Electronic Imaging*, vol. 13, no. 1, pp. 231-240, Jan. 2004.
- [2] D. Scharstein and R. Szeliski, "High-Accuracy Stereo Depth Maps Using Structured Light," *Proc. IEEE CS Conf. Computer Vision*, pp. 195-202, 2003.
- [3] Z. Zhang, R. Deriche, O. Faugeras, and Q. Luong, "A Robust Technique for Matching Two Uncalibrated Images through the Recovery of the Unknown Epipolar Geometry," *Artificial Intelligence J.*, vol. 78, nos. 1-2, pp. 87-119, Oct. 1995.
- [4] J. Forest and J. Salvi, "A Review of Laser Scanning Three-Dimensional Digitisers," *Proc. IEEE/RJS Int'l Conf. Intelligent Robots and Systems*, pp. 73-78, 2002.
- [5] F. Blais, M. Picard, and G. Godin, "Recursive Model Optimization Using ICP and Free Moving 3D Data Acquisition," *Proc. Fourth Int'l Conf. 3-D Digital Imaging and Modeling*, pp. 251-259, 2003.
- [6] J. Batlle, E. Mouaddib, and J. Salvi, "Recent Progress in Coded Structured Light as a Technique to Solve the Correspondence Problem: A Survey," *Pattern Recognition*, vol. 31, no. 7, pp. 963-982, 1998.
- [7] S. Rusinkiewicz, O. Hall-Holt, and M. Levoy, "Real-Time 3D Model Acquisition," *Proc. SIGGRAPH*, pp. 438-446, 2002.
- [8] M.D. Altschuler, B.R. Altschuler, J. Dijaki, L.A. Tamburino, and B. Woolford, "Robot Vision by Encoded Light Beams," *Three-Dimensional Machine Vision*, T. Kanade, ed., vol 87, pp. 97-149. Kluwer Academic, 1987.
- [9] K. Boyer and A. Kak, "Color-Encoded Structured Light for Rapid Active Ranging," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 9, no. 1, pp. 14-28, Jan. 1987.
- [10] M. Proesmans and L. Van Gool, "One-Shot Active 3D Image Capture," *Proc. SPIE*, vol. 3023, Three-Dimensional Image Capture, pp. 50-61, 1997.
- [11] D. Caspi, N. Kiryati, and J. Shamir, "Range Imaging With Adaptive Color Structured Light," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 5, pp. 470-480, May 1998.
- [12] C. Je, S. Lee, and R. Park, "High-Contrast Color Stripe Pattern for Rapid Structured-Light Range Imaging," *Proc. Eighth European Conf. Comp. Vision*, pp. 95-107, 2004.
- [13] M. Maruyama and S. Abe, "Range Sensing by Projecting Multiple Slits with Random Cuts," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 15, no. 6, pp. 647-651, June 1993.
- [14] P. Vuylsteke and A. Oosterlinck, "Range Image Acquisition with a Single Binary-Encoded Light Pattern," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 12, no. 2, pp. 148-164, Feb. 1990.
- [15] P. Griffin, L. Narasimhan, and S. Yee, "Generation of Uniquely Encoded Light Patterns for Range Data Acquisition," *Pattern Recognition*, vol. 25, no. 6, pp. 609-616, 1992.
- [16] E. Horn and N. Kiryati, "Toward Optimal Structured Light Patterns," *Image and Vision Computing*, vol. 17, pp. 87-97, 1999.
- [17] L. Zhang, B. Curless, and S. Seitz, "Rapid Shape Acquisition Using Color Structured Light and Multi-Pass Dynamic Programming," *Proc. First Int'l Symp. 3D Data Processing Visualization and Transmission*, pp. 24-36, 2002.
- [18] A. Strat and M. Oliveira, "A Point-and-Shoot Color 3D Camera," *Proc. Fourth Int'l Conf. 3-D Digital Imaging and Modeling*, pp. 483-490, 2003.
- [19] L. Zhang, N. Snavely, B. Curless, and S. Seitz, "Spacetime Faces: High-Resolution Capture for Modeling and Animation," *Proc. ACM Ann. Conf. Computer Graphics*, pp. 548-558, 2004.
- [20] Y. Boykov and V. Kolmogorov, "An Experimental Comparison of Min-Cut/Max-Flow Algorithms for Energy Minimization in Vision," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, no. 9, pp. 1138-1154, Sept. 2004.
- [21] T. Koninckx and L. Van Gool, "High-Speed Active 3D Acquisition Based in a Pattern-Specific Mesh," *Proc. SPIE: EI Photometrics*, pp. 26-37, 2003.
- [22] T. Koninckx, A. Griesser, and L. Van Gool, "Real-Time Range Scanning of Deformable Surfaces by Adaptively Coded Structured Light," *Proc. Fourth Int'l Conf. 3-D Digital Imaging and Modeling*, pp. 293-302, 2003.
- [23] T. Jaeggli, T.P. Koninckx, and L. Van Gool, "Online 3D Acquisition and Model Integration," *Proc. IEEE Int'l Workshop Projector-Camera Systems*, 2003.
- [24] T.P. Koninckx, "Adaptive Structured Light," PhD thesis, Catholic Univ. Leuven, 2005, <http://homes.esat.kuleuven.be/~tkoninck/publications.html>.



**Thomas P. Koninckx** received the MS degree in electrical engineering from the Katholieke Universiteit Leuven in 2001. Since then, he has been a research assistant in the Department of Electrical Engineering at the Katholieke Universiteit Leuven. He received the PhD degree in 2005. His research interests include active 3D acquisition, real-time vision, and the use of computer vision for biomedical applications.



**Luc Van Gool** received the PhD degree in electrical engineering from the Katholieke Universiteit Leuven. He is a professor of computer vision at the Katholieke Universiteit Leuven and the Swiss Federal Institute of Technology in Zurich (ETH). His research interests include 3D reconstruction, motion capture, animation, texture, object recognition, and computer vision for archaeology.

▷ For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).