

# Information Extraction from Complex Circular Charts

A. Antonacopoulos<sup>1</sup> and D.P. Kennedy

Department of Computer Science, University of Liverpool,  
Peach Street, Liverpool, L69 7ZF, United Kingdom

## Abstract

*This paper is concerned with the recognition of circular charts (disks). The application domain is tachograph disks. These disks exhibit typical properties of circular charts but are more complex as different types of information are present at the same time and there are serious artefacts that hinder the recognition process. The method starts by analysing the characteristics of the greyscale image to separate the background from the recorded information. The position of the disk and its centre are then established. Finally, the orientation of the disk is determined and the information (recorded in a circular manner) is read and converted into a linear representation. Experimental results of testing against professionally prepared ground-truth show an average accuracy of 94% in the recognition of the target information.*

## 1. Introduction

The recognition of charts in general, is a problem that has attracted relatively little attention [1]. Charts are very frequently embedded in textual documents as illustrations and routinely ignored by OCR systems. Other charts are in their own right significant sources of information. Examples of the latter include humidity/temperature charts and tachograph discs. These are all paper-based records of some activity in a graphical form.

This paper is concerned with the recognition of circular charts (disks). The application domain is tachograph disks, which exhibit typical characteristics of the circular chart recognition problem. In fact, their complexity, in terms of the different types of information recorded at the same time and the presence of serious artefacts, makes this a relatively difficult application.

Tachograph disks record information about a driver's activity (driving, duty, break), the distance travelled and the vehicle's speed during the course of a 24-hour period. Each disk is a legal document and constitutes a formal record, which can be inspected by the police (for breach of driving regulations) and used as evidence in court. It is

<sup>1</sup>Contact author (e-mail: aa@csc.liv.ac.uk)

common practice for haulage and travel companies to collect the tachograph disks from their vehicles in regular intervals and have them analysed by service bureaux who produce transcripts. The analysis that takes place is in the vast majority of cases manual (rotating a disk under a device with a magnifying lens and entering the information in a computer system) [2]. A semi-automated system using a turntable under a fixed linear CCD has been developed to read relatively coarse information. It essentially detects the change of line thickness around a circular arc (see mode trace below). The main problem with this approach is that it is very limited in the information it can read and requires the manual placing of the disc on the turntable.

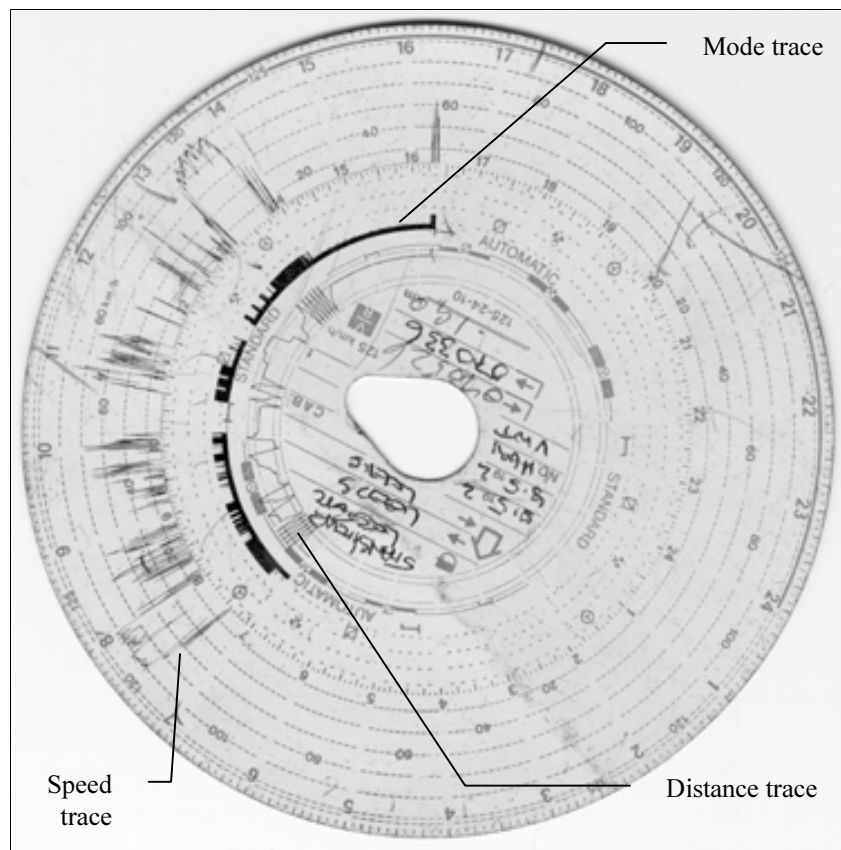
The proposed recognition approach is the first to use scanned tachograph disks. These images can be obtained in a fully automated way using a high-speed commercial flatbed scanner with a sheet feeder. This is a cost-effective solution and requires minimum setting-up effort. A tachograph disk image is shown in Figure 1.

## 2. Layout and Image Characteristics

A tachograph disk is a circular thin card with one side partitioned into 24 notional segments representing the hours of the day. The card itself is impregnated with special ink and covered in a thin layer of wax. The disk is inserted into the tachograph device which rotates it with time, completing a full rotation in 24 hours. Three styli in the tachograph device scratch the wax and, as a result, the normally invisible ink is exposed to air and will darken to become black.

Information is recorded in four different parts of the disk. The first part, the handwriting in the centre, includes details of the driver, the vehicle and the overall trip (place of departure, place of arrival and total mileage). Each of the other three parts is a trace created by a separate stylus in the tachograph device.

The trace closest to the centre is the distance recording. This shows the distance covered at any time period during the 24 hours. When the vehicle is moving, this trace oscillates within a specific range. The distance from peak to peak is 10 Kilometres, the distance from peak to trough,



**Figure 1. A scanned tachograph disk.**

5 Kilometres. When the distance trace is flat, the vehicle is stationary. The distance trace can be used to determine whether the manually entered distance information is correct.

The second trace recording is the mode trace. This trace is of primary importance in determining whether the driver has exceeded the legal driving time limit. In this type of tachograph, the mode trace shows the time spent in one of four activities: *driving*, *other working* (work-related tasks other than driving), *stand-by* (waiting times, sleeping cabin times during the trip) and *resting* periods. On the chart, each of these activities is denoted by a different thickness of the trace. The thickest of the bands in the mode trace is the driving time. The next thickest in the mode trace is the 'other working' time, followed by the stand-by time. Finally, the thinnest band in the mode trace denoted the resting time. On the chart in Figure 1, the mode trace shows recordings for driving, stand-by and resting periods. The thickness for the 'other working' time would lie between driving, the thickest, and stand-by, the narrow trace.

The third trace is the speed trace, recording the speed for the vehicle at any given moment. This trace can be checked to determine whether the vehicle was speeding.

It should be mentioned, that while the proposed system is designed for European Union standard tachograph charts, there are little differences in the layout of between those and tachograph charts used in the US, Canada, Australia and other parts of the world.

The nature of the tachograph disk and its (frequently rough) handling give rise to a number of artefacts in the image. First, a number of permanent marks may be present on the surface, interfering with the traces. A slight scratch will leave a dark mark (see top-right of Figure 1), while a crease will result in a wider grey area (see bottom-right of Figure 1), possibly with discontinuity, which may cause loss of data. Second, careless insertion and removal of the disk from the tachograph device may cause damage to the edge of the disk, causing the edge to become frayed and darken unevenly, leading to an indistinct edge in the image after scanning. Finally, the thickness of the card disk causes a shadow to be formed during scanning (see top of disk in Figure 1). The above artefacts are unique to this type of document (c.f. general documents) and present unique challenges in the extraction and recognition of the required information.

### 3. The Method

The method starts by locating the disk in the image in terms of the centre and radius of its circular outline form. The orientation of the disk is then established in order to determine the actual time in the 24-hour day that any part of a trace corresponds to. Finally, the mode trace is extracted and read to accurately determine the different types of driver activity at each time of the day.

The main goal is to determine whether a driver has exceeded the legal driving time limit. That is why in this paper only the techniques devised to extract and recognise the mode trace are presented. However, the techniques described below to locate the disk and determine its orientation are the same as when one of the other traces is extracted.

The techniques used to locate the disk are described in the next section. In Section 3.2, the determination of the disk orientation is described. Finally, in Section 4, experimental results are presented and discussed.

#### 3.1. Location of the disk

The location of the disk in the image involves the determination of the circular outline of the disk and the precise location of its centre (key for the accuracy of later stages). The location of the disk is necessary since it is not possible to pre-determine exactly where the disk will be placed on the scanner.

The disk is scanned at 400 dpi in order to preserve the thin line denoting rest periods in the mode trace. The image representation is chosen to be 8-bit greyscale in order to distinguish between the traces (darker) and other printed information on the disk (lighter).

For the location of the circular outline, the image is subsampled, resulting to one third of its original size. This step reduces the image access overhead and the reduction of resolution in this step only is not critical to the overall accuracy.

The histogram of the whole image is calculated. There are distinct regions in the histogram that correspond to different parts of the image. The region in the lighter end of the histogram corresponds to the non-disk background (image of the scanner cover). The next slightly darker distinct region corresponds to the disk background (the unscratched and unprinted surface area). Toward the middle there is a smooth hump corresponding to the pre-printed information on the disk, while at the very beginning of the histogram there is a small sharp peak that denotes the presence of the dark traces. Unfortunately, the same black peak also corresponds to the shadow and, if any, the occasional scratches.

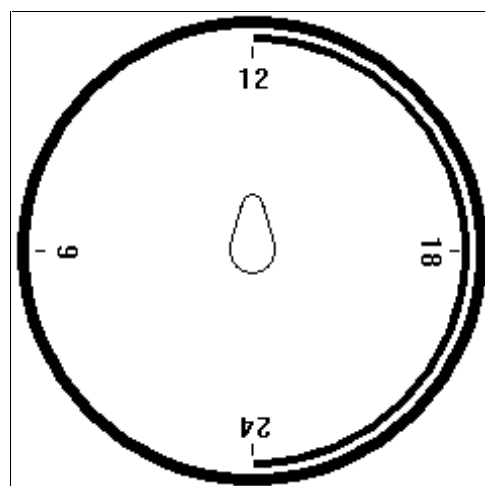


Figure 2. The main 'landmarks' on a tachograph disk.

Using the information from the histogram, the outside edge of the circular disk is traced. This is not a straightforward process as the shadow of the card distorts the circular form. It should be noted that tracing the inside edge does not give good results as it is frayed and also has the tick marks. The approach devised here is that, while tracing the outside edge, an attempt to move in towards the centre is made at the same time. The edge pixel only moves inward if the inward pixels are darker and the overall shape of the edge does not become distorted (according to an elasticity criterion).

The outside edge trace is used to estimate the parameters of the circular form (centre and radius). These parameters are obtained as a result of a Hough [3] transform. The processing and searching involved is minimised by using only the pixels of the identified edge and by restricting the accumulator to a range of expected (given the identified edge) circle parameters. It should be noted that small variations do exist between scanned tachograph disks due to different scanning parameters and due to the actual disk being slightly different (due to different manufacturers).

The initial Hough transform is performed on the reduced image and produces a coarse estimate of the centre and radius. A second Hough transform is performed on the full-size image data (only within the narrow range identified from the first transform). The result is an accurate estimate of the centre and radius of the circle. The accurate determination of the centre is crucial to the success of any subsequent step.

#### 3.2. Orientation determination

It is very important to correctly establish the orientation of the disk so that the traces can be read starting from

00:01 and ending at 24:00. The purpose of this step is to identify the axis from the '24' to the '12' mark (see Figure 2). The technique devised here exploits the fact that the longest axis of the elongated hole of the chart is aligned with the desired 24-12 axis. The vast majority (around 90%) of tachograph disks have this elongated hole, while those that do not, have other similarly identifiable features that align with the 24-12 axis.

The process starts by tracing the inside edge of the hole and identifying the edge point furthest away from the centre of the disk. The axis formed by these two points is a good first estimate of the 24-12 with a worst case accuracy of  $\pm 7.5$  degrees.

Using the estimate resulting from the processing of the hole as a starting basis, the method proceeds to find the 24-12 axis more accurately by identifying the actual '24' mark. A strip in the form of an arc inside the outer edge of the disk is examined. The strip extends either side of the axis estimate in the direction of the '24' mark. By processing the projection profile (smoothed) of the image data within the strip, tick marks, noise and digits are distinguished. In the profile, the '24' mark has the width of two digits (or of two single digits close together) and is preceded by single digit width peaks. Having identified the '24' mark, the axis through the centre that connects it to the '12' mark (by extending in the opposite direction) is accurately determined.

### 3.3. Reading the mode trace

Having established the centre of the tachograph disk and the start-of-the-day radius (the radius that crosses the '24' mark) the next step is to extract and recognise the mode trace. The mode trace is expected to lie within a certain circular band of the chart (extending above and below the half radius point). The task is to identify whether there is a trace at all and, if there is, to classify it according to its relative thickness.

The reading process involves recording the dark pixel runs (within the band limits) along consecutive radii covering the whole of the disk. In this way the circular trace band becomes a linear strip. Within this linear strip, distinct areas of similar height are classified as 'driving', 'other working', 'stand-by' and 'rest' according to their height. The start and the end points (in the horizontal direction) of each area correspond to specific times in the 24-hour day.

It is worth mentioning at this point that it is sufficient to distinguish between relative differences in thickness of mode trace. Therefore, even in the presence of distortion resulting from unwinding the circular trace (mainly dependent on the correct identification of the centre of the chart), the method is able to separate the different modes.

At the end of the reading process, the result is a sequence of time intervals, each one associated with an indication of the corresponding mode.

## 4. Results and Conclusions

The method was tested against a representative sample of tachograph disks. The sample did include average cases but also a higher-than-average number of exceptions (almost no presence of trace and badly creased and scratched disks). The location of the disk and determination of the orientation were tested using a sample of about 20 scanned disks. The accuracy of the mode trace was tested using a sample of 10 scanned disks with ground-truth prepared by professionals.

The accuracy of recognition of the driving time only (most important goal) is, on average, 94%. This includes the rare cases with recognition rate of around 59%, whilst in the majority of the cases the recognition rate is over 98%. The average recognition accuracy of all the modes excluding rest time is 83%, with similar variance as above. The rest time is the most difficult to recognise correctly because it is very thin (1 pixel in a 400dpi image). The average recognition accuracy of the whole mode trace, including rest, is 66%.

Further development is currently being carried out to fine-tune the method using more extensive data sets and associated ground-truth.

In conclusion, this paper has presented a method for extracting information from complex circular charts, which include different types of data in the presence of noise and other artefacts. The method recognises the mode trace on tachograph disks but the underlying techniques can be applied to the recognition of other information recorded in a circular form.

## Acknowledgements

The authors are grateful to Tachograph Analysis Consultants Limited of Liverpool, UK, for providing them with tachograph disks and ground-truth data.

## References

- [1] N. Yokokura and T. Watanabe, "Layout-based Approach for Extracting Constructive Elements of Bar Charts", *Proc. 2<sup>nd</sup> IAPR Workshop on Graphics Recognition (GREC'97)*, Nancy, France, August 22-23, 1997, pp. 119-126.
- [2] VDO Kienzle. Manual for the evaluation and use of the original Kienzle tachograph chart.
- [3] P.V.C. Hough, A Method and Means for Recognizing Complex Patterns, US Patent 3,069, 654, 1962.