# Audio–Visual Person Tracking
## A Practical Approach

Fotios Talantzis, Aristodemos Pnevmatikakis
and Anthony G Constantinides

Imperial College Press

# Audio-Visual Person Tracking
## A Practical Approach

Fotios Talantzis
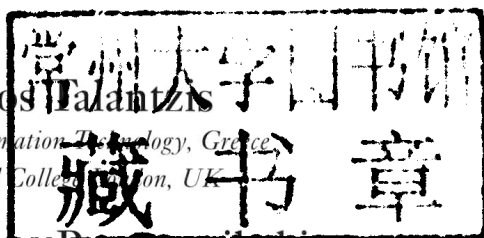*Athens Information Technology, Greece,*
*Imperial College London, UK*

Aristodemos Pnevmatikakis
*Athens Information Technology, Greece*

Anthony G Constantinides
*Imperial College London, UK*

**Communications and Signal Processing — Vol. 4**
**AUDIO VISUAL PERSON TRACKING**
**A Practical Approach**

# Audio-Visual Person Tracking
## A Practical Approach

# Communications and Signal Processing

Editors:   Prof. A. Manikas & Prof. A. G. Constantinides
*(Imperial College London, UK)*

*To my mother*
F. Talantzis

*To Efi, Athena and Katerina*
A. Pnevmatikakis

*To my students, past, present and future*
A.G. Constantinides

# Preface

Computing systems that are aware of human presence in order to provide heterogeneous services are gaining importance in living and working spaces, in entertainment, security and retail. A central role to such systems is the ability to sense humans and often track them in space across time. Tracking has become a mature topic in radar applications but requires a different set of sensors and algorithms when it involves humans. People generally do not like carrying tracking devices, a fact that inhibits service provision greatly. Instead, person tracking in this book is discussed in one of its unobtrusive flavours i.e. with the use of visual and audio modalities.

This book is about tracking humans using cameras and microphones, focusing on particle filtering algorithms. There are a few excellent texts on tracking [Blackman (1986); Blackman and Popoli (1999)], some of which focus on particle filters [Ristic *et al.* (2004)]. All these texts though focus on radar or sonar tracking. Audio-visual tracking needs different types of measurements on different types of signals: image [Gonzalez and Woods (2007)], video [Tekalp (1995); Forsyth and Ponce (2002); Shapiro and Stockman (2001)] and audio [Brandstein and Ward (2001)] signal processing elements need to be cast into the tracking frameworks. Two early works [Blake and Isard (1998); MacCormick (2002)] paved the way for visual tracking, but audio tracking still lacks a comprehensive text. A recent work covers audio-visual tracking [Zhu and Huang (2007)], mostly from the sensors and applications point of view.

Our aspiration is to fill in the gap between traditional tracking texts and signal processing texts. It is meant to be a solid introduction for the researcher starting in the field but also a good reference for people already working in it. It equips the reader with all the tools to measure the presence of humans in audio and visual signals and convert these measurements into

likelihood functions. These likelihood functions are suitable for driving many types of tracking algorithms, but the emphasis is on particle filtering. This became an obvious choice after inspecting the evolution of the relevant literature in the past decade that slowly moved away from deterministic and Kalman versions to the more versatile framework of particle filters.

We believe that the coverage of the material is end-to-end, in the sense that the theoretical foundation of particle filtering and the necessary image, video, audio and array signal processing elements are first established, followed by working examples and MATLAB [Mathworks (2010); Gilat (2004)] implementations. The MATLAB implementations aim to serve as skeletons for the employment of novel systems. We felt that the book would not be complete without a chapter discussing applications and real-world systems. This allowed us to give a more meaningful aspect to an otherwise abstract scientific problem.

This book is written by two generations of authors: Two former students and their PhD supervisor. Thus, in numerous ways the completion of this book would not have been possible without the contribution of the supervisor in terms of guidance, inspiration and patience for over a decade. Additionally, the other two authors hope to have brought into the book a hands-on approach to a rather modern signal processing topic. Either way, throughout the composition, we all felt the obligation to create a book that will familiarise researchers with the topic and quickly enable them to further advance the algorithms.

This book project has not been an easy task. It is the consolidation of years of research in the field, some of it funded by European research projects such as Computers in the Human Interaction Loop [CHIL (2007)], Cognitive Care and Guidance for Active Ageing [HERMES (2010)] and Real-Time Context-Aware and Personalised Media Streaming Environments for Large Scale Broadcasting Applications [e Director 2012 (2010)]. It has been the return-on-investment after the supervision of numerous students, discussions with valuable colleagues and most importantly spending hours away from our families.

<div align="right">

October 15, 2010
Fotios Talantzis,
Aristodemos Pnevmatikakis and
Antony G. Constantinides

</div>

# Acknowledgments

The experience from many tracking systems is depicted in this book. All these systems have been implemented in the Autonomic and Grid Computing Lab of Athens Information Technology with the help of many people, either students (former or current) of the authors, or researchers collaborating with them. The authors wish to thank:

- Andreas Stergiou for the discussions held during endless implementations of visual trackers.
- Nikos Katsarakis for researching and implementing face detectors and particle filter visual and audio-visual trackers.
- Martin and Rasmus Andersen for their MSc theses resulting into a real-time multi-camera and multi-cue particle filter visual 3D tracker. Many thanks also to Prof. Zheng Hua Tan for co-supervising the two theses.
- Panos Papageorgiou for helping with the implementation of our very first visual particle filter tracker.
- Panos Kasianidis and Vasilis Mylonakis for their help in audio direction of arrival estimation systems.
- Ghassan Karame, Genadios Genaro and George Miggos for their implementations of 3D and audio-visual trackers during their MSc theses.
- Elias Rentzeperis, Thodoris Petsatodis and Christos Boukis for sharing their research in voice activity detection.
- Achilleas Anagnostopoulos for his thesis resulting in an augmented reality system with integrated palm tracking.
- Eric Lehmann for his valuable contribution in the early stages of the Audio Tracker.

# Contents

# List of Figures