# Fixed-Point Signal Processing

# Synthesis Lectures on Signal Processing

**Editor**
**José Moura,** *Carnegie Mellon University*

**Fixed-Point Signal Processing**
Wayne T. Padgett and David V. Anderson
2009

**Advanced Radar Detection Schemes Under Mismatched Signal Models**
Francesco Bandiera, Danilo Orlando, Giuseppe Ricci
2009

**DSP for MATLAB™ and LabVIEW™ IV: LMS Adaptive Filtering**
Forester W. Isen
2009

**DSP for MATLAB™ and LabVIEW™ III: Digital Filter Design**
Forester W. Isen
2008

**DSP for MATLAB™ and LabVIEW™ II: Discrete Frequency Transforms**
Forester W. Isen
2008

**DSP for MATLAB™ and LabVIEW™ I: Fundamentals of Discrete Signal Processing**
Forester W. Isen
2008

**The Theory of Linear Prediction**
P. P. Vaidyanathan
2007

**Nonlinear Source Separation**
Luis B. Almeida
2006

**Spectral Analysis of Signals: The Missing Data Case**
Yanwei Wang, Jian Li, Petre Stoica
2006

# Fixed-Point Signal Processing

Wayne T. Padgett
Rose-Hulman Institute of Technology

David V. Anderson
Georgia Institute of Technology

*SYNTHESIS LECTURES ON SIGNAL PROCESSING #9*

## ABSTRACT

This book is intended to fill the gap between the "ideal precision" digital signal processing (DSP) that is widely taught, and the limited precision implementation skills that are commonly required in fixed–point processors and field programmable gate arrays (FPGAs). These skills are often neglected at the university level, particularly for undergraduates. We have attempted to create a resource both for a DSP elective course, and for the practicing engineer with a need to understand fixed–point implementation. Although we assume a background in DSP, Chapter 2 contains a review of basic theory, and Chapter 3 reviews random processes to support the noise model of quantization error. Chapter 4 details the binary arithmetic that underlies fixed–point processors, and then introduces fractional format for binary numbers. Chapter 5 covers the noise model for quantization error and the effects of coefficient quantization in filters. Because of the numerical sensitivity of IIR filters, they are used extensively as an example system in both Chapters 5 and 6. Fortunately, the principles of dealing with limited precision can be applied to a wide variety of numerically sensitive systems, not just IIR filters. Chapter 6 discusses the problems of product roundoff error, and various methods of scaling to avoid overflow. Chapter 7 discusses limit cycle effects and a few common methods for minimizing them.

There are a number of simple exercises integrated into the text to allow you to test your understanding. Answers to the exercises are included in the footnotes. A number of Matlab examples are provided in the text. They generally assume access to the Fixed–Point Toolbox. If you lack access to this software, consider either purchasing or requesting an evaluation license from The Mathworks. The code listed in the text and other helpful Matlab code is also available at `http://www.morganclaypool.com/page/padgett` and `http://www.rose-hulman.edu/~padgett/fpsp`. You will also find Matlab exercises designed to demonstrate each of the four types of error discussed in Chapters 5 and 6. Simulink examples are also provided on the web site.

# Contents

# Notes

## ACKNOWLEDGMENTS

We would like to gratefully acknowledge the support of Texas Instruments and The Mathworks in the development of this material.

## PERMISSIONS

Figure 5.8 is adapted from [26]. Permission is pending.
Figures 5.14 through 5.19 are adapted from [23]. Permission is pending.

## PRELIMINARY

This is a preliminary version of this book.


Wayne T. Padgett and David V. Anderson
September 2009