# Assurance and Security for AI-enabled Systems

**Joshua D. Harguess**
**Nathaniel D. Bastian**
**Teresa L. Pace**
*Editors*

**22–23 April 2024**
**National Harbor, Maryland, United States**

*Sponsored and Published by*
SPIE

**Volume 13054**

# SPIE. DIGITAL LIBRARY

SPIEDigitalLibrary.org

# Contents