# High Performance Python

*Micha Gorelick and Ian Ozsvald*

# Table of Contents