# Contents

*Contents*

*Contents*

**A Warning Note:**
A big computer, a complex algorithm and a long time does not equal science.
[Robert Gentleman, SSC 2003, Halifax, June 2003]