# Contents

**Chapter 2. A Qualitative Approach to Many-core Architecture** . . . .    27
Benoît DUPONT DE DINECHIN

**Chapter 3. The Plural Many-core Architecture – High Performance
at Low Power** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .    53
Ran GINOSAR