

MACHINE TRANSLATION

An Introductory Guide

Douglas Arnold

Lorna Balkan

Siety Meijer

R. Lee Humphreys

Louisa Sadler

This book was originally published by NCC Blackwell Ltd. as *Machine Translation: an Introductory Guide*, NCC Blackwell, London, 1994, ISBN: 1855542-17x. However, it is now out of print, so we are making this copy available. Please continue to cite the original publication details.

© Arnold, Balkan, Humphreys, Meijer, Sadler, 1994, 1996.

The right of Arnold, Balkan, Humphreys, Meijer, and Sadler to be identified as the author of this work has been asserted in accordance with the Copyright, Designs and Patents Act 1988.

First published 1994

First published in the USA 1994

Originally published as D.J. Arnold, Lorna Balkan, Siety Meijer, R. Lee. Humphreys and Louisa Sadler, 1994, *Machine Translation: an Introductory Guide*, ISBN: 1855542-17x. Originally published in the UK by NCC Blackwell Ltd., 108 Cowley Rd, Oxford OX4 1JF, and in the USA by Blackwell Publishers, 238 Main St. Cambridge, Mass. 02142. Blackwells-NCC, London.

This copy differs from the published version only in that some cartoons are omitted (there is blank space where the cartoons appear), and it has been reduced to print two pages on one A4 page (the original is slightly larger).

It remains copyright © the authors, and may not be reproduced in whole or part without their permission. This can be obtained by writing to: Doug Arnold, Department of Language & Linguistics, University of Essex, Colchester, CO4 3SQ, UK, doug@essex.ac.uk.

March 6, 2001

Preface

Automatic translation between human languages (‘Machine Translation’) is a Science Fiction staple, and a long-term scientific dream of enormous social, political, and scientific importance. It was one of the earliest applications suggested for digital computers, but turning this dream into reality has turned out to be a much harder, and in many ways a much more interesting task than at first appeared. Nevertheless, though there remain many outstanding problems, some degree of automatic translation is now a daily reality, and it is likely that during the next decade the bulk of routine technical and business translation will be done with some kind of automatic translation tool, from humble databases containing canned translations of technical terms to genuine Machine Translation Systems that can produce reasonable draft translations (provided the input observes certain restrictions on subject matter, style, and vocabulary).

Unfortunately, how this is possible or what it really means is hard to appreciate for those without the time, patience, or training to read the relevant academic research papers, which in any case do not give a very good picture of what is involved in practice. It was for this reason that we decided to try to write a book which would be genuinely introductory (in the sense of not presupposing a background in any relevant discipline), but which would look at all aspects of Machine Translation: covering questions of what it is like to use a modern Machine Translation system, through questions about how it is done, to questions of evaluating systems, and what developments can be foreseen in the near to medium future.

We would like to express our thanks to various people. First, we would like to thank each other. The process of writing this book has been slower than we originally hoped (five authors is five pairs of hands, but also five sets of opinions). However, we think that our extensive discussions and revisions have in the end produced a better book in terms of content, style, presentation, and so on. We think we deserve no little credit for maintaining a pleasant working atmosphere while expending this level of effort and commitment while under pressure caused by other academic responsibilities.

We would also like to thank our colleagues at the Computational Linguistics and Machine Translation (CL/MT) group at the University of Essex for suggestions and practical support, especially Lisa Hamilton, Kerry Maxwell, Dave Moffat, Tim Nicholas, Melissa Parker, Martin Rondell and Andy Way.

ii Preface

For proofreading and constructive criticism we would like to thank John Roberts of the Department of Language and Linguistics at the University of Essex, and John Roberts and Karen Woods of NCC Blackwell. We are also grateful to those people who have helped us by checking the examples which are in languages other than English and Dutch, especially Laurence Danlos (French), and Nicola Jörn (German).

Of course, none of them is responsible for the errors of content, style or presentation that remain.

D.J. Arnold
L. Balkan
R. Lee Humphreys
S. Meijer
L. Sadler

Colchester, August 1993.

Contents

| | |
|---|-----------|
| Preface | i |
| 1 Introduction and Overview | 1 |
| 1.1 Introduction | 1 |
| 1.2 Why MT Matters | 4 |
| 1.3 Popular Conceptions and Misconceptions | 6 |
| 1.4 A Bit of History | 13 |
| 1.5 Summary | 16 |
| 1.6 Further Reading | 16 |
| 2 Machine Translation in Practice | 19 |
| 2.1 Introduction | 19 |
| 2.2 The Scenario | 20 |
| 2.3 Document Preparation: Authoring and Pre-Editing | 28 |
| 2.4 The Translation Process | 30 |
| 2.5 Document Revision | 34 |
| 2.6 Summary | 35 |
| 2.7 Further Reading | 36 |

| | | |
|----------|--|------------|
| 3 | Representation and Processing | 37 |
| 3.1 | Introduction | 37 |
| 3.2 | Representing Linguistic Knowledge | 39 |
| 3.3 | Processing | 52 |
| 3.4 | Summary | 60 |
| 3.5 | Further Reading | 60 |
| 4 | Machine Translation Engines | 63 |
| 4.1 | Introduction | 63 |
| 4.2 | Transformer Architectures | 63 |
| 4.3 | Linguistic Knowledge Architectures | 71 |
| 4.4 | Summary | 85 |
| 4.5 | Further Reading | 85 |
| 5 | Dictionaries | 87 |
| 5.1 | Introduction | 87 |
| 5.2 | Paper Dictionaries | 89 |
| 5.3 | Types of Word Information | 91 |
| 5.4 | Dictionaries and Morphology | 98 |
| 5.5 | Terminology | 106 |
| 5.6 | Summary | 109 |
| 5.7 | Further Reading | 109 |
| 6 | Translation Problems | 111 |
| 6.1 | Introduction | 111 |

| | | |
|----------|---|------------|
| 6.2 | Ambiguity | 111 |
| 6.3 | Lexical and Structural Mismatches | 115 |
| 6.4 | Multiword units: Idioms and Collocations | 121 |
| 6.5 | Summary | 127 |
| 6.6 | Further Reading | 128 |
| 7 | Representation and Processing Revisited: Meaning | 129 |
| 7.1 | Introduction | 129 |
| 7.2 | Semantics | 130 |
| 7.3 | Pragmatics | 136 |
| 7.4 | Real World Knowledge | 139 |
| 7.5 | Summary | 144 |
| 7.6 | Further Reading | 144 |
| 8 | Input | 147 |
| 8.1 | Introduction | 147 |
| 8.2 | The Electronic Document | 147 |
| 8.3 | Controlled Languages | 156 |
| 8.4 | Sublanguage MT | 159 |
| 8.5 | Summary | 163 |
| 8.6 | Further Reading | 164 |
| 9 | Evaluating MT Systems | 165 |
| 9.1 | Introduction | 165 |
| 9.2 | Some Central Issues | 165 |

vi CONTENTS

| | | |
|-----------|--|------------|
| 9.3 | Evaluation of Engine Performance | 168 |
| 9.4 | Operational Evaluation | 178 |
| 9.5 | Summary | 180 |
| 9.6 | Further Reading | 180 |
| 10 | New Directions in MT | 183 |
| 10.1 | Introduction | 183 |
| 10.2 | Rule-Based MT | 185 |
| 10.3 | Resources for MT | 195 |
| 10.4 | Empirical Approaches to MT | 198 |
| 10.5 | Summary | 205 |
| 10.6 | Further Reading | 205 |
| | Useful Addresses | 207 |
| | Glossary | 209 |