# Modern Information Retrieval

## the concepts and technology behind search

### Second edition

document
query
search

Ricardo Baeza-Yates
Berthier Ribeiro-Neto

# Modern Information Retrieval
## The Concepts and Technology behind Search

**Ricardo Baeza-Yates**
**Berthier Ribeiro-Neto**

**Second edition**

**Addison-Wesley**

To be filled by Pearson

*To Helena, Rosa, and our children*


Amo los libros exploradores,
libros con bosque o nieve,
profundidad o cielo

Un libro, un libro lleno
de contactos humanos, de camisas,
un libro sin soledad,
con hombres y herramientas,
un libro es la victoria.

de "Oda al Libro" (I) y (II),
en *Odas Elementales*, 1954.

Pablo Neruda

território de homens livres
que será nosso país
e será pátria de todos.
Irmãos, cantai esse mundo
que não verei, mas virá
um dia, dentro de mil anos,
talvez mais. . . não tenho pressa.

de "Cidade Prevista" no livro
*A Rosa do Povo*, 1945.

Carlos Drummond de Andrade

I love books that explore,
books with a forest or snow,
depth or sky

A book, a book full
of human contacts, of shirts,
a book without solitude,
with people and tools,
a book is the victory.

from "Ode to the Book" (I) and (II),
in *Elemental Odes*, 1954.

Pablo Neruda

territory of free men
that will be our country
and will be the nation of all
Brothers, sing that world
which I'll not see, but which will
come
one day, in a thousand years,
maybe more. . . no hurry.

from "Prevised City" in the book
*The Rose of the People*, 1945.

Carlos Drummond de Andrade

# Contents

## 6  Documents: Languages & Properties

*with Gonzalo Navarro and Nivio Ziviani*