# SPARSER

version 2.3b

## USER DOCUMENTATION

# Preliminary edition

January 1995

David D. McDonald, Ph.D.
davidmcdonald@alum.mit.edu

# 1. Introduction

## 1.1 Capabilities

Sparser is a software system written in the language Common Lisp, consisting of a directory of files for loading into the user's Lisp image, a pre-loaded image that may be launched on its own, or some combination of the two.

At its heart, Sparser is a bottom-up, phrase-structure-based chart parser, optimized for semantic grammars and partial parsing. It has a highly efficient implementation based on resources of reused data structures and the option of an indelible, deterministic parsing algorithm.

Included with Sparser is a notation for constructing a grammar of context free and limited context-sensitive rewrite rules (§4). Phrase structure rules may incorporate user-defined objects and operations over them to be used as the referents of the phrases that the rules construct (§4.2).

Since real texts can pose problems for analysis that are not readily or best analyzed using phrase structure techniques, Sparser provides a high-level, task-oriented programming language for defining grammars using a whole range of techniques. A Lisp application programmer's interface (API) for user-defined parsing networks over terminals and initial non-terminals is provided (§7). User-defined code may be executed upon the completion of rules or words (§7.1). The interpretation of input characters is table-driven (§6.2.4). The definition of what constitutes a terminal can be modified (§6.4.1). Special routines may be executed at newlines (§6.4.3).

Sparser takes a string or file of ascii characters as input, and supplies a populated chart as output. An API for examining this chart is supplied (§5), as are trace routines (§8).

## 1.2 Overview

Before Sparser can be used you must load its files into Lisp or activate a stored Lisp image that already includes them. You must also load the grammar you intend to use, which will typically be a set of files of rule definition forms that you created earlier. At that point you can have Sparser analyze files of text or individual strings by calling the appropriate function (§2.5). You can also embed calls to run Sparser in your own functions, along with code for examining the results of the its analysis.

To write a phrase structure grammar for Sparser you have to understand the conventions of the rule notation (§4); you will also want to look at the trace routines that have been provided (§8). To write routines of your own for examining  Sparser's analysis you have to understand the data structures that make up Sparser's representations (§3, §5). To write more advanced routines to take advantage of the user-code hooks within Sparser you also need to understand

something of Sparser's algorithms so that you will know when the hooks are called and what is available for your routines to examine at those points (§6).