

# Discovering Functional Sites of Amino Acid Sequences Using Sorted Variable Generalization

Takashi Ishikawa<sup>1</sup>      Shigeki Mitaku<sup>2</sup>      Takao Terano<sup>3</sup>  
takashi@j.kisarazu.ac.jp    mitaku@cc.tuat.ac.jp    terano@gssm.otsuka.tsukuba.ac.jp

Makiko Suwa<sup>2</sup>      Takatsugu Hirokawa<sup>2</sup>  
suwa@cc.tuat.ac.jp    hirokawa@cc.tuat.ac.jp

<sup>1</sup> Kisarazu National College of Technology  
2-11-1 Kiyomidai-higashi, Kisarazu, Chiba 292, Japan

<sup>2</sup> Tokyo University of Agriculture and Technology  
2-24-16 Naka-cho, Koganei-shi, Tokyo 184, Japan

<sup>3</sup> University of Tsukuba  
3-29-1 Otsuka, Bunkyo-ku, Tokyo 112, Japan

## Abstract

*This research develops a method for discovering functional sites of amino acid sequences using an Inductive Logic Programming (ILP) method with sorted variable generalization. Functional sites provide clues to building a knowledge base for prediction of protein functions from amino acid sequences. The proposed method generates hypotheses of functional sites directly from aligned amino acid sequences using an ILP method extended with sorted variable generalization. The proposed method is shown to be useful for discovering functional sites by an example application to the case of bacteriorhodopsin-like proteins.*

## 1 Introduction

This research develops a method for discovering functional sites of amino acid sequences using an *Inductive Logic Programming* (ILP) method with *sorted variable generalization*. Functional sites provide clues to building a knowledge base for prediction of protein functions from amino acid sequences. Our approach [3] [4] [5] is based on the following assumption: *If there exist any functional sites, then we are able to predict specific functions of a protein from its amino acid sequence.* In order to discover functional sites of amino acid sequences, we use a machine learning technique with a framework of *Inductive Logic Programming* (ILP) [8]. The proposed method generates hypotheses of functional sites directly from aligned amino acid sequences using an ILP method extended with *sorted variable generalization*. *Sorted variable generalization* is a generalization operator of induction to generalize a constant term by replacing it with a *sort* symbol to which the constant belongs [6]. A *sort* is defined as a subset of constants. Therefore we require only to prepare *sorts* for the domain of interest. The proposed method is shown to be useful for discovering functional sites by an example application to the case of bacteriorhodopsin-like proteins. Using the propose method we are able to discover functional sites of amino acid sequences seeming to related to these specific functions.

## 2 Discovering functional sites

The problem addressed is to discover functional sites of amino acid sequences of proteins with a common specific function. A functional site is a subsequence of an amino acid sequence that only exists in sequences of a certain function. The subsequence has length of one or more and may contain any groups of amino acids. The use of information about secondary structure increases the validity of functional sites with short length. Therefore the problem is reduced to discovering a combination of functional sites existing only in given amino acid sequences with a specific protein function.

The problem is specified as follows:

