# machine learning approaches to bioinformatics
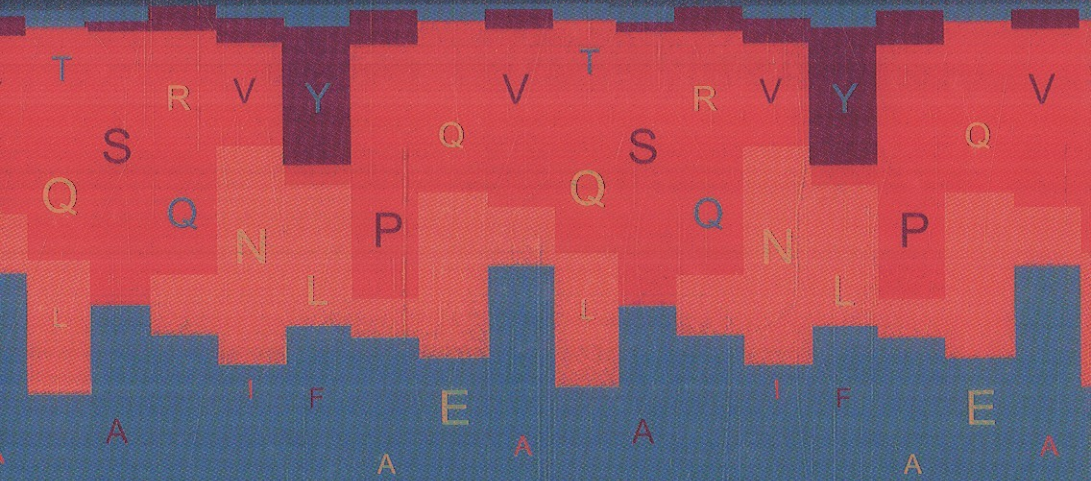
zheng rong yang

# Contents