# Syllabus of CS6222 (2025 Spring)

**Course Name:**

CS6222: Advanced Topics in Computational Biology

**Course Description:**

This course introduces fundamental concepts and methodologies of bioinformatics and computational biology. Topics covered include sequence, structure and function databases of nucleic acid and protein molecules; advanced sequence and structure alignment methods; molecular dynamics and Monte Carlo simulations; AI and deep learning methods; and protein and RNA folding and structure prediction (homologous modeling, threading, *ab initio* folding, and deep learning).

Special emphasis is placed on the latest breakthroughs in deep learning biomolecular structure predictions brought about by Google DeepMind and other leading teams of the field, with a focus on introducing the core technologies behind the breakthroughs, including but not limited to convolutional neural networks, graph neural networks, multimodal networks, transformers, language models, and diffusion models. The classes emphasize understanding computational biology concepts and their practical application, aiming to equip students with the skills to utilize cutting-edge bioinformatics tools/methods to solve problems in their own research projects.

**Instructor:**

Prof Yang Zhang (Email: zhang@nus.edu.sg; Phone: 6601-1241)

**Orgainzation:**

Department of Computer Science
School of Computing
National University of Singapore

**Schedule and location:**

Monday 4-6pm from 2025/1/13 to 2025/4/8, at COM1-VCRM

**Textbook:**

No textbook is required for this course. All teaching materials will be posted on the course website.

**Presentation, homework, & grades:**

There will be homework assignments, including code writing and literature reading and presentation. Final grade consists of literature presentation (20%), homework (30%), and exam (50%).

# Table of content