

CS4220: Knowledge Discovery Methods in Bioinformatics

Course Briefing

Wong Limsoon



Recommended “Pre-requisites”



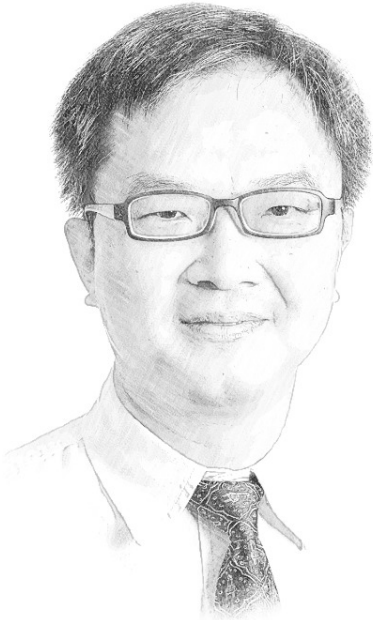
- **Completed modules on**
 - Programming
 - Algorithms
 - Basic molecular biology
 - **ST2334 Probability & Statistics**
 - **CS2220 Introduction to Computational Biology**

Objectives

- **Exposure to knowledge-discovery techniques**
 - **Enhance flexible & logical problem-solving skill**
 - **Understand bioinformatics problems and their solution in depth**
 - A modern network-based perspective
- **To achieve goals above, we expose students to case studies spanning gene expression and proteomic analysis, protein functional prediction, epistatic interaction analysis, etc.**

Professors

- **Wong Limsoon**



- **Anders Skanderup**



Contents of Course Overview



- **Time Table**
- **Course Syllabus**
- **Course Homepage**
- **Teaching Style**
- **Project, Assignments, Exams**
- **Readings**
- **Assessment**

- **Quick Overview of Themes and Applications of Bioinformatics**

Time Table

- **Lecture**
 - Tuesday 12nn-3pm at COM1-204
- **“Tutorial” (it is integrated into lecture)**
- **Emails**
 - wongls@comp.nus.edu.sg
 - skanderupamj@gis.a-star.edu.sg
- **Consultation**
 - Any time; just make appt

Course Syllabus

- **Essence of Biostatistics**
 - Statistical estimation
 - Hypothesis testing
 - Principle component analysis
- **Essence of Data Mining**
 - Clustering
 - Association rules
 - Classification
 - Class-imbalance learning
- **Gene Expression Analysis**
 - Basic gene expression analysis
 - Improving reproducibility
 - Dealing with small sample
- **Proteomic Profile Analysis**
 - Basic proteomic profile analysis
 - Improving consistency
 - Improving coverage
- **Batch effects**
 - Visualization
 - Normalization
 - PC1 removal
 - Batch effect-resistant feature selection
 - Batch effect-resistant classifiers
- **Network Perturbations in Disease Context**
- **Other classic or hot topics to be determined**

Teaching Style

- **Bioinformatics is a broad area**
 - **Need to learn a lot of material by yourself**
 - Reading papers
 - Try the exercises
 - Practise on the w
- And do this before each lecture*
- **Don't expect to be told everything**

Assignments, Project, & Exam



- **Assignments (30-40% of marks)**
 - 3 to 4 assignments
 - Some are simple programming assignments
- **Project (20-30% of marks)**
 - Based on a case study in the class
 - 8-10 pages of report / ppt slides expected
- **Exam (40% of marks)**
 - 1 final open-book exam

Be Honest

- **Exam**
 - Absence w/o good cause results in ZERO mark
 - Cheating results in ZERO mark
- **Discussion on assignments & project is allowed**
- **Blatant plagiarism is not allowed**
 - Offender gets ZERO mark for assignment or exam
 - Penalty applies to those who copied AND those who allowed their assignments to be copied

Background Readings



- **Every lecture will be accompanied by a small set of “must-read” and “good-to-read” articles**
- **The “must-read” articles are considered lecture notes and are examinable**

Related Courses

- **CS2220 Introduction to Computational Biology**
 - Understand bioinformatics problems; interpretational skills
- **CS4330 Combinatorial Methods in Bioinformatics**
- **CS4220 Knowledge Discovery Methods in Bioinformatics**
 - Gene expression, proteomic profiling, protein interaction, transcription factor interaction, pathway perturbation
- **CS5238 Advanced Combinatorial Methods in Bioinformatics**
 - Seq alignment, whole-genome alignment, suffix tree, seq indexing, motif finding, RNA sec struct prediction, phylogeny reconstruction
- **CS6222 Computational frontier in precision medicine**
- **Etc ...**

Any questions?



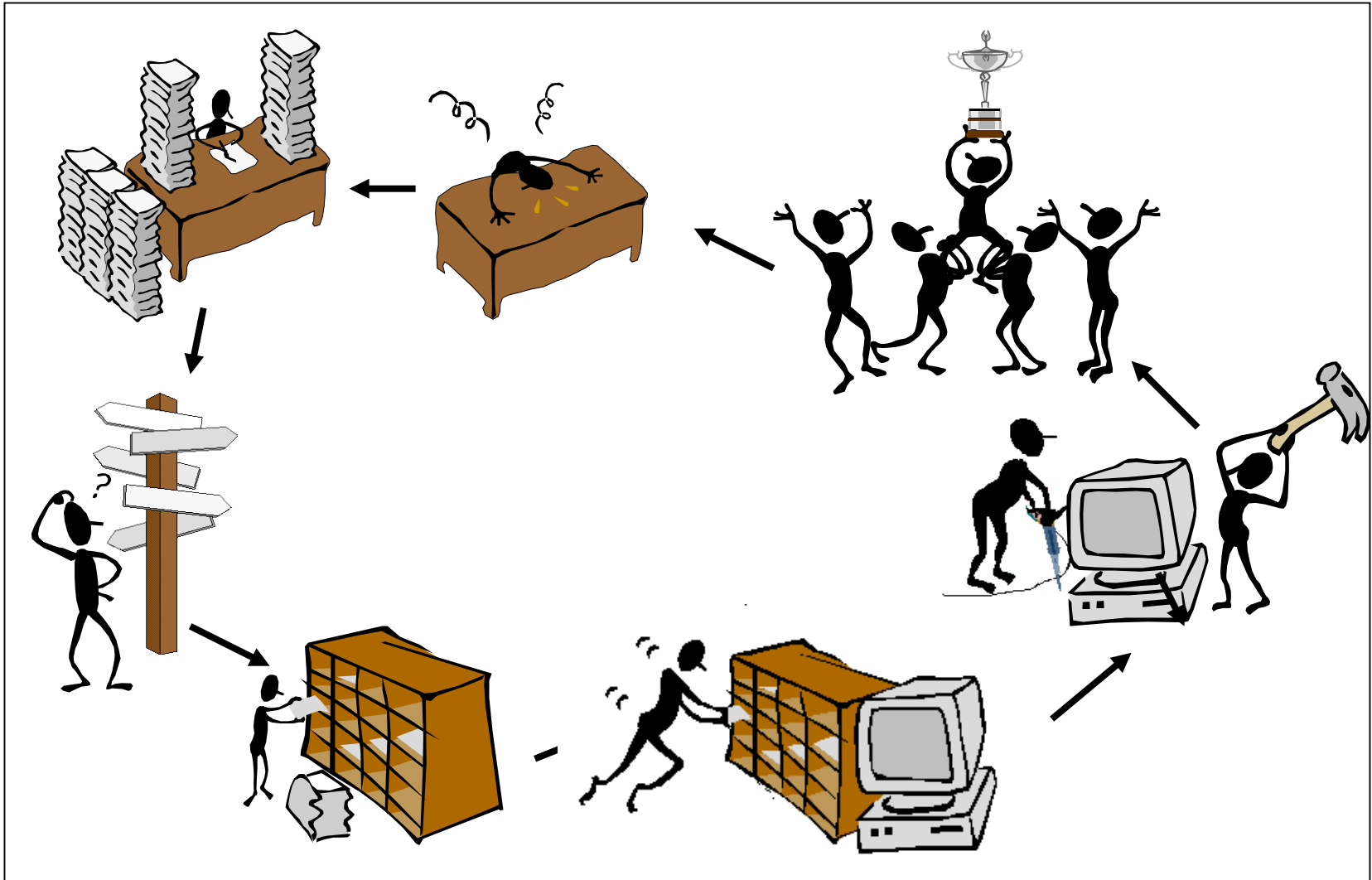
I hope you will enjoy this class 😊

Themes and Applications of Bioinformatics

**These slides are for those who have not
taken CS2220 to read at your own leisure**



What is Bioinformatics?



Themes of Bioinformatics

Themes of This Course

Bioinformatics involves

Data Mgmt +

Knowledge Discovery +

Sequence Analysis +

Physical Modeling + ...

Knowledge Discovery =

Statistics + Algorithms + Databases

The Promises of Bioinformatics



To the patient:

Better drug, better treatment

To the pharma:

Save time, save cost, make more \$

To the scientist:

Better science

Fulfilling the Promise via Drugs

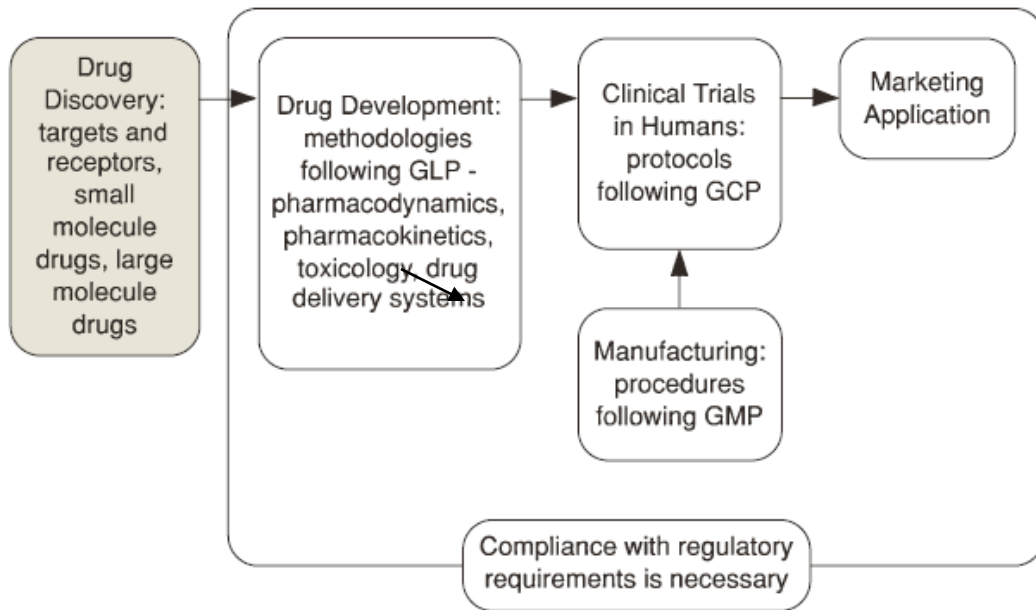


Figure from Rick Ng, *Drugs: From Discovery to Approval*

- **Bioinformatics is applicable to drug development**
- **Drug discovery: Design small molecules that bind target proteins**
 - Which proteins?
 - What should binding accomplish?
- **Biomarkers**

Pervasiveness of Bioinformatics



- **Bioinformatics is mandatory for large-scale biology**
 - e.g., High-throughput, massively-parallel measurements, or “lab on a chip” miniaturization
- **Computational data analysis is mandatory for indirect experimental methods**
 - e.g., reconstruction based on phase contrast or wave diffraction.
- **What about the rest of biology (and medicine) ?**
- **Limitless opportunities!**

Some Bioinformatics Problems

- **Biological Data Searching**
- **Biological Data Integration**
- **Gene/Promoter finding**
- **Cis-regulatory DNA**
- **Gene/Protein Network**
- **Protein/RNA Structure Prediction**
- **Evolutionary Tree reconstruction**
- **Infer Protein Function**
- **Disease Diagnosis**
- **Disease Prognosis**
- **Disease Treatment Optimization, ...**

Commonly Used Data Sources

These slides are for those who have not taken CS2220 to read at your own leisure



Introductory References

- **S.K. Ng, “Molecular Biology for the Practical Bioinformatician”, *The Practical Bioinformatician*, Chapter 1, pages 1-30, WSPC, 2004**
- **Lots of useful videos,**
http://www.as.wvu.edu/~dray/Bio_219.html
- **Materials from CS2220,**
<http://www.comp.nus.edu.sg/~wongls/courses/cs2220/2019>