# 课程大纲
# COURSE SYLLABUS

| 1. | 课程代码/名称<br>**Course Code/Title** | 基因组和组学数据分析/Genomics Data Analysis |
|---|---|---|
| 2. | 课程性质<br>**Compulsory/Elective** | 专业-选修/Elective |
| 3. | 课程学分/学时<br>**Course Credit/Hours** | 3 |
| 4. | 授课语言<br>**Teaching Language** | 英文为主，必要时辅以少量中文解释；教材、课件、考试为英文<br>English, with a few Chinese. Textbooks, ppts and examinations are in English |
| 5. | 授课教师<br>**Instructor(s)** | 靳文菲,生物系<br>Dr. Wenei JIN, Department of Biology, SUSTech<br>jinwf@sustech.eud.cn |
| 6. | 先修要求<br>**Pre-requisites** | Prerequisites include a college level mathematics, statistics and molecular biology |
| 7. | 教学目标<br>**Course Objectives** | |

**Course Objectives**

Genomics is an interdisciplinary field of biology focusing on the structure, function, evolution, mapping, and editing of genomes. This subject will help students to understand life in a whole picture --- functional genomics, comparative genomics, evolutionary genomics, transcriptomics, 3D genomics, their interrelations and influence on the organism. Furthermore, this course emphasize on computational analyses of the genomics. Various existing methods will be critically described and the strengths and limitations of each will be discussed, with practical assignments utilizing the tools. It is to train students' vigorous Scientific Spirit and inspire their scientific curiosity.

**Learning Outcomes**

With the completion of this course, The student could

1) Be familiar with the major genomic database and database searching

2) Be familiar with Linux and master at least one programming language

3) Conduct various genomic analysis

4) Analyze next generation sequencing data including DNA-seq, RNA-seq, ChIP-seq, single cell sequencing data.

| 8. | 教学方法<br>**Teaching Methods** | |
|---|---|---|
| | PPT presentation, class discussion, written assignments, computational practice and quizzes | |
| 9. | 教学内容<br>**Course Contents** | |
| | **Section 1** | I Introduction of Genomics and Basic computational skills (Linux/shell+ python/R) |

| | | Hours: 10 |
| --- | --- | --- |
| | | **1. Past, Present and Future of Genomics and Course Introduction**<br>1.1 What is genomics?<br>1.2 The origin and development of genomics<br>1.3 Present Genomics<br>1.4 Challenges and future of Genomics<br>1.5 Course Introduction: Goals, outline, evaluation/examination and learning guidelines<br><br>**2. Linux and Linux commands**<br>2.1 Server and operating systems<br>2.2 Linux operating system and Open Source Software<br>2.3 Terminal and basic Linux commands<br>2.4 File system and server management<br>2.5 Personal setting<br><br>**3. Programing language and shell**<br>3.1 Principles of programming languages<br>3.2 Script languages and bash shell<br>3.3 Basic shell functions<br>3.4 I/O Redirection and file descriptors<br>3.5 Pattern matching in shell<br>3.6 Biological data analysis: Modularization and pipeline<br><br>**4. Programming language Python**<br>4.1 The features of Python<br>4.2 Data types and variable<br>4.3 Control structures<br>4.4 Functions and procedures<br>4.5 Classes & instances<br>4.6 Modules & packages<br><br>**5. R Language Statistics and Drawing**<br>5.1 Quick start R<br>5.2 Basic principles and concepts<br>5.3 Data operation in R (Vectors, matrices, arrays, data frames)<br>5.4 Plot figures<br>5.5 Statistical Analysis of R<br>5.6 Function definition and programing<br>5.7 packages |
| **Section 2** | | II Basic sequence analysis<br>Hours: 6<br><br>**6 Pairwise sequence alignments**<br>6.1 Sequence change over time<br>6.2 Pairwise sequence comparisons<br>6.3 Dynamic programming alignment<br>6.3.1 Global alignment (Needleman-Wunsch)<br>6.3.2 Local alignment (Smith-Waterman)<br>6.4 Sequence Similarity Searching |

| | | 6.4.1 FASTA Algorithm |
|---|---|---|
| | | 6.4.2 BLAST Algorithm |
| | | |
| | | 7. Multiple Sequence Alignment and Phylogenetics |
| | | 7.1 Significance of multiple sequence alignment |
| | | 7.2 Progressive Alignment (ClustalW) |
| | | 7.3 Basics of phylogeny: Characters, traits, nodes, branches, lineages |
| | | 7.4 Molecular clock and modeling sequence evolution |
| | | 7.5 Distances and clustering algorithm: UPGMA and Neighbor Joining (NJ) |
| | | 7.6 From sequence alignments to trees: Parsimony methods |
| | | 7.7 Probability based approach: Maximum likelihood methods |
| **Section 3** | | III Next Generation Sequencing (NGS) and cancer genomics |
| | | Hours: 8 |
| | | 1. NGS and Short reads mapping |
| | | Introduction to Genomic Technologies |
| | | From Sanger sequencing to NGS |
| | | Principles of NGS: Massive parallel sequencing |
| | | Features of NGS data: Short reads |
| | | Uses Trie structure (Trie and Suffix Array) to search a reference genome |
| | | Burrows–Wheeler transform（BWT） |
| | | |
| | | 2. Variant calling and output |
| | | Genetic variants: structure variants, SNV, CNV |
| | | SAM format for mapped reads |
| | | Approaches for variants calling |
| | | VCF format for saving called variants |
| | | |
| | | 3. Cancer genomics and single cell cancer genomics |
| | | Calling variants in cancer genomics |
| | | Single cell cancer genomes |
| | | Tumor microevolution |
| **Section 4** | | IV Transcriptomic and epigenomic analysis |
| | | Hours: 10 |
| | | |
| | | 1.Gene expression profiling and RNA-seq |
| | | What's the advantage of RNA-seq compared with microarray? |
| | | What factors should we consider for RNA-seq data normalization? |
| | | What's the advantage of single cell sequencing over bulk cells? |
| | | 2. Single cell RNA-seq |
| | | Cellular heterogeneity |
| | | Single cell RNA-seq technologies |
| | | Distinct cell populations |
| | | Pseudo-time inference |
| | | 3. Epigenome and data anlysis |
| | | Definition of epigenetics? |
| | | How to detect genome-wide DNA methylation? |
| | | How to detect genome-wide nucleosome positioning and chromatin accessibility? |
| | | How to identify genome-wide TF binging sites? How to do the peak calling? |

| | | What is Hi-C? How to identify the significant interaction？<br>4. Single cell epigenomics<br>challenges<br>scDNAse-seq<br>scMNase-seq<br>scATAC-seq<br>multipe-omics<br><br>5. Gene Ontology and enrichment analysis<br>Gene ontology (GO) program<br>Structure of GO<br>Gene annotation in GO<br>GO/pathway enrichment analysis<br>Gene set enrichment analysis (GSEA) |
|---|---|---|
| **Section 5** | | V Population Genomics and association study<br>Hours: 12<br><br>1. Haplotype and linkage disequilibrium<br>What is Haplotype？<br>What is linkage disequilibrium?<br>Calculation of linkage disequilibrium<br>Complete LD and perfect LD<br>Recombination rate and LD block<br><br>2. Population genomics<br>Effective population size (Ne)<br>The major forces shaping population<br>Population substructure<br>Measure population structure (F-statistics)<br>Approaches for analysis of population structure<br>Analysis of molecular variance (AMOVA)<br>Dimensionality reduction<br>Model based approaches<br><br>3. Approaches for natural selection detection<br>Divergence rate and phylogenetic shadowing<br>Changed function-altering mutation, e.g., dN/dS or KN/KS<br>Polymorphism deviating from interspecies divergence e.g. Hudson-Kreitman-Aguade (HKA) test and McDonald-Kreitman (MK) test<br>Changed allele frequency spectrum e.g., Tajima's D<br>Increased derived allele frequencies<br>Extended haplotype homozygosity (EHH), e.g., iHS<br>Locus-specific population differentiation, e.g., FST<br>Biased ancestry contribution in admixed population.<br>Composite strategies. e.g. combine multiple factors and Likelihood-ratio test<br><br>4. Genomics and evolution theory<br>Evolution is a unifying theme in biology<br>History of "evolutionary thought"<br>Darwin's Four Postulates<br>1) Individuals within species vary. |

| | | |
|---|---|---|
| **10.** | | 2) Some variations are heritable.<br>3) More offspring are produced than can survive<br>4) Survival and reproduction are nonrandom<br>Modern evolutionary theory<br><br>5. Genomics and human evolution<br>Classic approaches for study human evolution<br>Human origin models<br>Mitochondrial and Y-chromosome detailed "Out of African" theory<br>Genomic approach revolutionized our understanding of human evolution<br>Human origin model based on genomic data<br>Human migration and natural selection<br><br>6.Gene mapping for identifying disease associated variants<br>Linkage analysis for rare disease/Mendelian diseases<br>Genetic model for complex disease:<br>1) Common disease common variant (CDCV),<br>2) Common disease rare variant (CDRV)<br>Methods for identifying disease associated variants<br>1) Family based association study<br>2) Case control based association study<br>3) Association study based on next generation sequencing (NGS)<br>Challenge in identifying disease associated variants |
| | **Section 6** | VI Genomics and Big data<br>Hours: 2<br><br>Accumulation of genomic data<br>Clustering algorithms<br>1) Hierarchical agglomerative clustering<br>2) Partitioning methods<br>Two approaches for dimensionality reduction (Feature Selection and feature extraction)<br>Linear reduction<br>1) Principal component analysis (PCA)<br>2) Singular Value Decomposition (SVD)<br>3) Multi-Dimensional Scaling (MDS)<br>Non-linear reduction<br>1) t-distributed stochastic neighbor embedding （t-SNE）<br>2) Uniform Manifold Approximation and Projection (UMAP) |
| | **Section 7** | |
| | **Section 8** | |
| | **Section 9** | |
| | **Section 10** | |
| | **…………** | |
| **10.** | 课程考核<br>**Course Assessment** | |

请再此注明：①考查/考试；②分数构成。
```
Total score 100
Attendance 10
Class Performance 20
Assignments 20
Mid-term Exam 20
Final Presentation/Exam 30
```

I encourage you to ask questions during the class, and you will get credit of Class Performance.

| 11. | 教材及其它参考资料<br>**Textbook and Supplementary Readings** |
|---|---|

Reference books:
Introduction to Genomics. Arthur M. Lesk. Oxford University Press; 3 edition. ISBN-10: 0198754833.
Bioinformatics and Functional Genomics. Jonathan Pevsner. Wiley-Blackwell; 3 edition. ISBN-10: 1118581784
Medical Genetics And Genomics, Csaba Szalai, ISBN: 9789632791876

Will have some docs at the beginning of the course