

# Supplemental Material

*CBE—Life Sciences Education*

Rosenwald *et al.*

**Supplemental Table 1: The CourseSource Bioinformatics Framework**

Topic	Learning Goals	Sample Learning Objectives
<b>Computation in the life sciences</b>	What is the role of computation in hypothesis-driven discovery processes within the life sciences?	<ul style="list-style-type: none"> <li>• Describe the role of bioinformatics in the scientific research method.</li> <li>• Explain the necessity for computation in life sciences research.</li> <li>• Explain the role of wet-bench techniques in verifying computational results in life science research.</li> <li>• Compare and contrast computer-based research with wet-lab research.</li> <li>• Read a scientific article and evaluate how bioinformatics methods were employed by the authors to explore a particular hypothesis.</li> <li>• Given a scientific question, develop a hypothesis and define computational approaches that could be used to explore the hypothesis.</li> <li>• Evaluate the social, legal, and ethical implications of computational approaches to understanding biology.</li> </ul>
	What computational concepts are important in bioinformatics?	<ul style="list-style-type: none"> <li>• Define the term <i>algorithm</i>.</li> <li>• Explain the difference between a <i>heuristic</i> (approximate) algorithm and an <i>exact</i> algorithm.</li> <li>• Describe the three basic programming structures: sequential, repetition (e.g., <b>while</b>, <b>for</b>) and selection (e.g., <b>if</b>).</li> <li>• Use variables and data structures (e.g., lists, arrays, scalars, hash functions).</li> <li>• Describe what a regular expression is.</li> <li>• Explain the concept of cloud computing.</li> <li>• Describe the importance of “big data” in bioinformatics.</li> <li>• Describe the means by which “big data” are managed and stored (e.g. <a href="http://dmptool.org">dmptool.org</a>).</li> </ul>
	What statistical concepts are important in bioinformatics?	<ul style="list-style-type: none"> <li>• Calculate average, median, mode, range, standard deviations for a given data set.</li> <li>• Calculate p-values using a t-test for discrete data.</li> <li>• Calculate p-values using a z-test for continuous data.</li> <li>• Calculate an e-value statistic.</li> <li>• Describe the importance of statistical analysis of big data sets.</li> <li>• Create a network to illustrate gene interactions.</li> </ul>
<b>DNA - Information Storage [GENOMICS]</b>	Where are data about the genome found (e.g., nucleotide sequence, epigenomics) and how are they stored and	<ul style="list-style-type: none"> <li>• Describe how nucleotide sequence data are represented (FASTA, FASTQ, GenBank).</li> <li>• Describe the nucleotide databases available at NCBI.</li> </ul>

	accessed?	<ul style="list-style-type: none"> <li>Describe how the NCBI nucleotide databases intersect with other nucleotide databases (EBI, DDBJ, UniProt, etc.).</li> <li>Compare and contrast the data contained in different nucleotide databases.</li> <li>Search for a sequence record in a nucleotide database with a given accession number.</li> <li>Create a collection of nucleotide sequence records that meet a specified criterion (e.g., gene name or symbol).</li> <li>Determine the DNA methylation state of a particular region of a genome.</li> <li>Describe the types of metadata that accompany sequence data to make for useful biological interpretation (e.g. biological source, accession number, GeneID, journal articles, etc.).</li> </ul>
	How can bioinformatics tools be employed to analyze genetic information?	<ul style="list-style-type: none"> <li>Calculate the alignment score between two DNA sequences using a provided scoring matrix.</li> <li>Perform a BLASTN search and interpret the results.</li> <li>Explain the BLASTN algorithm for nucleotide sequence information.</li> <li>Interpret the biological significance of an e-value.</li> <li>Annotate a prokaryotic gene (derive a model).</li> <li>Annotate a eukaryotic gene (derive a model).</li> <li>Create and interpret a multiple sequence alignment (e.g. T-COFFEE, MUSCLE, etc.).</li> <li>For a genomic region of interest (e.g., the neighborhood of a particular gene), use a genome browser to view nearby genes, transcription factor binding regions, epigenetic information, etc.</li> <li>Describe Hidden Markov Models and how they can be used to assess motifs.</li> </ul>
<b>RNA - Information Transfer</b> <b>[TRANSCRIPTOMICS]</b>	Where are data about the transcriptome found (e.g., expression, epigenomics and structure) and how are they stored and accessed?	<ul style="list-style-type: none"> <li>Identify the euchromatin/heterochromatin boundaries, histone states in a given sequence, and the nucleosome modifications</li> <li>Compare and contrast DNA structure at telomeres and centromeres.</li> <li>Describe the RNA databases available at NCBI (e.g., ESTs, UniGene).</li> <li>Describe the types of metadata that accompany sequence data to make for useful biological interpretation (e.g. biological source, accession number, GeneID, journal articles, etc.).</li> </ul>
	How can bioinformatics tools be employed to examine <i>transfer</i> of genetic	<ul style="list-style-type: none"> <li>Given a microarray or RNA-seq data file, find the set of significantly differentially expressed genes.</li> <li>Perform motif discovery on the promoter</li> </ul>

	information?	<p>regions of a set of genes identified by aChIP-seq experiment.</p> <ul style="list-style-type: none"> <li>• Use RNA structure prediction programs (e.g., RNAsoft, RNAfold, RNAstructure) to evaluate possible structures for an RNA sequence.</li> <li>• Identify the possible different splice isoforms possible from a given gene sequence.</li> </ul>
<b>Protein - Information in Action</b> <b>[PROTEOMICS]</b>	Where are data about the proteome found (e.g., amino acid sequence and structure) and how are they stored and accessed?	<ul style="list-style-type: none"> <li>• Describe how protein sequence data are represented (e.g., FASTA, GenBank, etc.)</li> <li>• Describe the different protein databases available at NCBI (sequence, structure, function).</li> <li>• Describe how the NCBI databases intersect with other databases (e.g., EBI, DDBJ, UniProt, etc.).</li> <li>• Compare and contrast data contained in different databases.</li> <li>• Search for a protein record in a database with a given accession number.</li> <li>• Create a collection of records that meet a specified criterion (e.g., gene name or symbol).</li> <li>• Describe the types of metadata that accompany sequence, structure, and function data to make for useful biological interpretation (e.g. biological source, accession number, UniProt number, journal articles, etc.).</li> </ul>
	How can bioinformatics tools be employed to examine protein structure and function?	<ul style="list-style-type: none"> <li>• Explain the BLASTP, BLASTX, tBLASTn, tBLASTx algorithms for protein sequence information</li> <li>• Interpret the biological significance of an e-value.</li> <li>• Describe Hidden Markov Models and how they can be used to assess motifs.</li> <li>• Query a dataset with a specific protein sequence to learn about potential functions (e.g. Pfam, CDD, SwissProt, UniProt, etc.).</li> <li>• View and interpret the structure output from Protein Data Bank (e.g. Cn3D, Jmol, etc.).</li> <li>• Propose potential functions for a given protein structure.</li> <li>• Explain the outputs from protein-folding algorithms to predict structure from sequence.</li> <li>• Understand how protein structures are determined (e.g. NMR, crystallography).</li> <li>• Compare and contrast the output from 2-D gel experiments.</li> <li>• Analyze the output from mass spectrometry analysis (e.g., use the MASCOT package).</li> </ul>
<b>Small Molecules - Cellular Homeostasis</b> <b>[METABOLOMICS &amp;]</b>	Where are data about metabolomics and systems biology found and how are	<ul style="list-style-type: none"> <li>• Describe how metabolomics data are represented (e.g. Human Metabolome Database, METLIN</li> </ul>

<b>SYSTEMS BIOLOGY]</b>	they stored and accessed?	<p>Database, etc.)</p> <ul style="list-style-type: none"> <li>Describe the different metabolomics databases that are available.</li> <li>Describe the types of metadata that accompany metabolomics data to make for useful biological interpretation.</li> </ul>
	How can bioinformatics tools be employed to examine flow of molecules within pathways?	<ul style="list-style-type: none"> <li>Perform a GO analysis to identify the pathways relevant to a set of genes (e.g., identified by a transcriptomic study or a proteomic experiment).</li> <li>Use the KEGG pathway database to look up the interaction network of a pathway.</li> <li>Interpret the data from experiments (e.g., mass spectrometry, nuclear magnetic resonance, etc.) to determine levels of small molecule metabolites.</li> </ul>
<b>Ecology and Evolution [METAGENOMICS]</b>	How can bioinformatics tools be employed to examine ecological niches?	<ul style="list-style-type: none"> <li>Create and interpret a multiple sequence alignment (e.g., T-COFFEE, MUSCLE, etc.).</li> <li>Describe the components of a phylogenetic tree (e.g., root, node, leaf).</li> <li>Explain the various types of phylogenetic trees (e.g., rooted, unrooted).</li> <li>Interpret a phylogenetic tree (e.g., which organism is most closely related to a given organism in the tree)</li> <li>Sketch a phylogenetic tree from its Newick representation.</li> <li>Use bootstrapping to assess the quality of a phylogenetic tree.</li> <li>Create a phylogenetic tree for a set of related sequences (nucleotide or amino acid) (e.g. MEGA).</li> <li>Use pre-existing tools to analyze a metagenomic data set to determine the set of organisms present in a metagenomic sample (e.g., 16srRNA, Greengenes, mothur, etc.)</li> </ul>
<b>Computational Skills</b>	How do biologists employ software development as part of the scientific discovery process?	<ul style="list-style-type: none"> <li>Write a script to calculate the reverse complement of a nucleotide sequence</li> <li>Write a script to determine reading frames of a nucleotide sequence.</li> <li>Write a script to calculate melting point of double-stranded DNA.</li> <li>Write a script to retrieve the promoter regions for a set of related genes.</li> <li>Write a script to find the longest open reading frame in a given nucleotide sequence</li> <li>Write a script to calculate the reverse complement of a nucleotide sequence.</li> <li>Write a script to convert an RNA sequence to cDNA and to amino acids</li> <li>Write a script to calculate molecular weight and isoelectric point.</li> <li>Write a script to count amino acid frequency.</li> <li>Write a script that compares the relative hydrophilicity/hydrophobicity of two</li> </ul>

		protein sequences.
	<p>What higher-level computational skills can be used in bioinformatics research?</p>	<ul style="list-style-type: none"> <li>• Use a spreadsheet to perform simple data analysis.</li> <li>• Use a spreadsheet to open, read, parse, modify and output comma-separate (.csv) files that will be ready to use in subsequent tools.</li> <li>• Perform elementary statistical analysis on an “omics” dataset (e.g. using Excel or Weka).</li> <li>• Perform Input/Output with data files.</li> <li>• Interact with remote servers.</li> <li>• Construct a bioinformatics pipeline.</li> <li>• Use open source libraries and packages (e.g., BioPerl, Biopython, R, BioConductor).</li> <li>• Use programs at the Unix/Linux command line to analyze bioinformatics data.</li> <li>• Use graph theory to represent data networks.</li> </ul>