

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ  
Львівський національний університет імені Івана Франка  
Біологічний факультет  
Кафедра генетики та біотехнології

Затверджено  
на засіданні кафедри генетики та біотехнології  
(протокол № 6 від 15 березня 2023 р.)

Завідувач кафедри. \_\_\_\_\_

  
Віктор ФЕДОРЕНКО

Силабус з навчальної дисципліни «Біоінформатика»  
що викладається в межах ОПП «Біохімія», «Біофізика», «Ботаніка»,  
«Генетика», «Зоологія», «Мікробіологія», «Фізіологія людини і тварин»,  
«Фізіологія рослин» другого (магістерського) рівня вищої освіти  
для здобувачів зі спеціальності 091 Біологія та біохімія

Львів 2023

Назва дисципліни	Біоінформатика
Статус дисципліни	Нормативна
Обсяг дисципліни	Лекції – 32 год, лабораторні – 16 год, самостійна робота – 102 год; разом – 150 год (5 кредитів ECTS)
Викладається для	Магістрів біології та біохімії, денна форма), 2 (літній) семестр навчання
Форма контролю	Іспит
Галузь знань	Галузь знань – 9 Біологія, спеціальність – 091 Біологія та біохімія
Адреса викладання	Біологічний факультет, вул. Грушевського 4, Львів 79005
Розклад занять	<a href="https://bioweb.lnu.edu.ua/students/rozklad-ispytiv">https://bioweb.lnu.edu.ua/students/rozklad-ispytiv</a>
Викладач (-і)	Богдан Омелянович Осташ (лекції) Юрій Васильович Ребець (лабораторні заняття)
Профайли викладачів	<a href="http://bioweb.lnu.edu.ua/employee/ostash-b-o">http://bioweb.lnu.edu.ua/employee/ostash-b-o</a> <a href="https://www.researchgate.net/profile/Yuriy-Rebets">https://www.researchgate.net/profile/Yuriy-Rebets</a>
Контактний тел.	032 2394407
E-mail:	<a href="mailto:bohdan.ostash@lnu.edu.ua">bohdan.ostash@lnu.edu.ua</a> <a href="mailto:yurko.rebets@gmail.com">yurko.rebets@gmail.com</a>
Сторінка курсу на сайті кафедри	<a href="https://bioweb.lnu.edu.ua/course/bioinformatyka">https://bioweb.lnu.edu.ua/course/bioinformatyka</a>
Консультації	<i>Очні консультації:</i> II семестр (2023 р), щовівторка, 11:30-13:00 <i>Онлайн-консультації:</i> у форматі “питання-відповідь” через електронну пошту, в робочі дні тижня, з 10:00-16:00; очікуйте на відповідь не пізніше ніж за три доби з моменту надходження питання

### 1. Коротка анотація до курсу

Станом на літо 2023 року в біомедичних базах даних є інформація про гени більше 400 000 організмів, або більш ніж 21 трильйон пар нуклеотидів (п.н.). Це колосальна кількість інформації про генетичну основу життя, і її кількість невпинно зростає. Нуклеотидні послідовності не мають змісту доти, доки дослідник його не визначить, як от здатність кодувати білок, бути промотором тощо. Систематичний аналіз усього масиву даних експериментальними методами нереалістичний як економічно так і технологічно. Розширення таких баз даних стимулювало розвиток методів комп'ютерного аналізу нуклеотидних і амінокислотних послідовностей (НАП). Ці методи допомогли упорядкувати дані у межах різних спеціалізованих веб-сервісів і класифікувати їх відповідно до різних критеріїв і потреб дослідників. Однак найважливішим результатом синтезу молекулярної біології, геноміки, інформатики, комп'ютерних технологій і статистики – біоінформатики – стало те, що ця нова галузь перетворилася у самостійне знаряддя наукового відкриття. Курс “Біоінформатика” включає розділи, присвячені структурам баз даних; попарному і множинному вирівнюванню НАП; моделям, що узагальнюють множинні вирівнювання; основам філогенетичної реконструкції на основі; передбаченню функцій генів і пошукові мотивів у НАП; передбаченню тривимірних структур білків; аналізу молекул РНК.

*Ключові слова:* математичні моделі в біології, множинні і попарні вирівнювання, моделі Маркова, філогенетичні реконструкції, бази НАП.

### 2. Мета та завдання курсу

**Мета:** сформувати у слухачів курсу систему знань про основні бази даних нуклеотидних й амінокислотних послідовностей (НАП) та їхнього комп'ютерного аналізу, зокрема

попарного і множинного вирівнювання, моделей на основі вирівнювань, філогенетичної реконструкції на основі НАП.

**Завдання:** а) викласти концептуальні математичні та молекулярно-біологічні засади, на яких ґрунтується біоінформатика; б) ознайомити студентів з наявним арсеналом методів біоінформатики, їхніх можливостей і обмежень, останніх тенденцій розвитку дисципліни; в) сформувати у студентів арсенал активних знань у галузі цієї дисципліни, тобто навчити їх формулювати наукові питання, відповіді на які можна шукати із залученням методів біоінформатики.

### 3. Формат курсу – очний або дистанційний

### 4. Результати навчання

Після курсу студент буде: а) знати основні бази даних НАП, принципи функціонування алгоритмів вирівнювання НАП, і спектр питань щодо структури, функції та еволюції НАП, на які можна шукати відповіді із застосуванням біоінформатичних веб-ресурсів і програм; б) мати базове розуміння методів філогенетичної реконструкції на основі НАП; в) вміти користуватись базами даних, де зберігаються нуклеотидні та амінокислотні послідовності, інформація про структури геномів, їхню експресію (транскриптомні дані – генні чіпи і RNAseq) та відповідні протеоми; г) вміти порівнювати НАП з гомологічними послідовностями за допомогою методів попарного і множинного вирівнювань; оперувати паттернами і профілями; д) вміти виявляти регуляторні послідовності у геномах – повтори, паліндроми, консервативні мотиви; передбачати функціональність гена (білка) на основі аналізу консервативних доменів і каталітичних центрів; е) вміти користуватись методами передбачення і моделювання структури білків; здійснювати молекулярно-філогенетичний аналіз.

Курс розроблено так, щоб сформувати у студентів загальні і фахові компетентності:

ЗК01. Здатність працювати у міжнародному контексті.

ЗК02. Здатність використовувати інформаційні та комунікаційні технології.

ЗК03. Здатність генерувати нові ідеї (креативність).

ФК01. Здатність користуватися новітніми досягненнями біології, необхідними для професійної, дослідницької та/або інноваційної діяльності.

ФК02. Здатність формулювати задачі моделювання, створювати моделі об'єктів і процесів на прикладі різних рівнів організації живого із використанням математичних методів й інформаційних технологій.

ФК03. Здатність користуватися сучасними інформаційними технологіями та аналізувати інформацію в галузі біології і на межі предметних галузей.

ФК06. Здатність прогнозувати напрямки розвитку сучасної біології на основі загального аналізу розвитку науки і технологій.

ФК07. Здатність діагностувати стан біологічних систем за результатами дослідження організмів різних рівнів організації

ФК10. Здатність використовувати результати наукового пошуку в практичній діяльності.

#### Програмні результати навчання:

- ПР1. Володіти державною та іноземною мовами на рівні, достатньому для спілкування з професійних питань та презентації результатів власних досліджень.
- ПР2. Використовувати бібліотеки, інформаційні бази даних, інтернет ресурси для пошуку необхідної інформації.
- ПР4. Розв'язувати складні задачі в галузі біології, генерувати та оцінювати ідеї.
- ПР5. Аналізувати та оцінювати вплив досягнень біології на розвиток суспільства.
- ПР6. Аналізувати біологічні явища та процеси на молекулярному, клітинному, організменному, популяційно-видовому та біосферному рівнях з точки зору фундаментальних загальнонаукових знань, а також за використання спеціальних сучасних методів досліджень.
- ПР11. Проводити статистичну обробку, аналіз та узагальнення отриманих експериментальних даних із використанням програмних засобів та сучасних інформаційних технологій.
- ПР14. Дотримуватись норм академічної доброчесності під час навчання та провадження наукової діяльності, знати основні правові норми щодо захисту інтелектуальної власності.
- ПР16. Критично осмислювати теорії, принципи, методи з різних галузей біології для вирішення практичних задач і проблем.

#### **5. Пререквізити та необхідне обладнання для вивчення курсу.**

Знання англійської мови на рівні, достатньому для перекладу наукових статей; необхідні знання з основ генетики, біохімії, зоології та ботаніки. Розуміння базових математичних понять (логарифм, частка значень, експонента, відсоткові величини) та теорії імовірностей та статистичного аналізу даних. Базові навички роботи з комп'ютером. Наявність комп'ютера/смартфона з підключенням до інтернету (для лабораторних занять у випадку дистанційного формату навчання).

#### **6. Політики курсу.**

Відвідування лекційної частини курсу вільне. Матеріали лекційного курсу (PowerPoint-презентації) та електронний підручник буде надано електронною поштою усім студентам. Усі статті і матеріали, або гіперпосилання до них, що згадано нижче у схемі курсу (п. 7) – буде надано. Перша частина курсу (включно з лекцією про моделі на основі множинних вирівнювань – див. нижче) закінчується письмовим модулем. Написання модуля у визначений час обов'язкове, відсутність можлива лише за умови поважної причини, що має бути задокументовано (довідка про хворобу тощо). Відвідування лабораторних занять обов'язкове, під час яких студенти отримують бали за виконання контрольних завдань. Більше про систему оцінювання – див. нижче розділ 8. Очікується, що студенти дотримуватимуться правил Академічної доброчесності – див. [http://www.lnu.edu.ua/wp-content/uploads/2019/06/reg\\_academic\\_virtue.pdf](http://www.lnu.edu.ua/wp-content/uploads/2019/06/reg_academic_virtue.pdf). Нульова

толерантність (у вигляді недопуску до іспиту) до плагіату, списування, хабарництва. Зниження оцінки при виявленні фактів несамостійного підготовлення завдань до практичних занять (нерозуміння підготовленої презентації, механічне використання перекладів, згенерованих автоматичними перекладачами тексту).

## 7. Схема курсу

### Тиждень 1

**Лекція 1 (2 год). Вступ.** Що таке ДНК і білок. Центральна догма молекулярної біології ХХ століття, її сучасне тлумачення з точки зору епігенетики й теорії інформації. Біоінформатика як синтез методів молекулярної біології, генетики, інформатики і статистики. Маргарет О. Дейгоф і перші моделі еволюції НАП. Нуклеотид, кодон, амінокислотний залишок – елементарні одиниці інформації, якими оперує біоінформатика. Типи даних, що генерують геномні, транскриптомні і протеомні методи досліджень. Інтерактом. Системний аналіз. Роль біоінформатичних методів у біологічних дослідженнях. Журнал Nucleic Acids Research – провідник у світі біоінформатики. *Матеріали* – презентація лекції bioinf-L10.pdf. *Література*: [1, 2, 4] (див. список наприкінці схеми курсу).

**Лабораторна 1 (2 год).** Біоінформатичні сервіси на веб-порталі NCBI – PubMed, GenBank, Genome, Taxonomy, GEO datasets. Пошук інформації в PubMed. Ідентифікатори статей в PubMed – doi, PMID. Підрозділ GenBank – архів нуклеотидних даних і продуктів їхньої трансляції. Структура файлів GenBank. Змінні і постійні ідентифікатори НАП в GenBank. Депонування нових НАП у GenBank. Підрозділ Genome – структура бази даних і її використання. Поняття геномного переглядача.

**Самостійна робота (24 год).** Ознайомлення з відкритим ресурсом для вивчення алгоритмів, що застосовуються в біоінформатиці - <http://rosalind.info/problems/locations/>. Основний фокус – на алгоритмах, що дають змогу аналізувати НАП – див. <http://rosalind.info/problems/list-view/>.

### Тиждень 2

**Лекція 2 (2 год). Математичні моделі НАП – концептуальні засади.** Біологічна модель – на прикладі абетки і мови. Що таке інформація? Символьне повідомлення. Що таке частота, імовірність та вірогідність події? Імовірність (частота) трапляння підпоследовності (слова) у последовності (тексті) – моделі Бернуллі і Маркова. Поняття Байєзової статистики стосовно аналізу НАП. Окремі випадки використання елементів Байєзової статистики, вірогідності і різноманітних розподілів імовірності до розв’язання біологічних питань. *Матеріали* – презентація лекції bioinf-L20.pdf. *Література*: [1, 2,4].

**Самостійна робота (10 год).** Ознайомлення з можливостями бази геномних даних PATRIC <https://www.patricbrc.org/>

### Тиждень 3

**Лекція 3 (2 год). Математичні моделі еволюції** Математичні моделі еволюції нуклеотидних последовностей. Моделі еволюції нуклеотидних последовностей як приклад параметризованих моделей. Модель Джакса-Кімури JC69, її параметри. Теорія молекулярного годинника, її практичне застосування. Типи матриць заміщення – одиничні, емпіричні, параметризовані. Райони низької складності в НАП та повтори.

Повтори – кількісно домінуюча форма організації генетичного матеріалу. Неструктуровані білки як приклад послідовностей з низькою складністю. *Матеріали* – презентація лекції bioinf-L30.pdf. *Література*: [1, 4, 6].

**Лабораторна робота 2** (2 год). Спеціалізовані веб-сервіси аналізу НАП. StrepDB – інтегрована база даних про структуру геномів модельних бактерій класу Actinobacteria. Типи завдань, які можна виконувати на основі StrepDB. FlyBase. MaizeGDB. TBDB. ExPaSy – перелік основних біоінформатичних веб-сервісів і баз даних.

**Самостійна робота** (5 год). Завантаження програми для аналізу нуклеотидних і амінокислотних послідовностей UGENE <https://unipro-ugene.software.informer.com/1.9/>

#### Тижень 4

**Лекція 4** (2 год). **Порівняння НАП – концептуальні засади.** Еволюційна спорідненість (гомологія) як концептуальна основа порівняння НАП. Гомологічність, подібність, ідентичність. Локальне і глобальне вирівнювання. Підпослідовності, прогалини, штрафи, рахунок вирівнювання. Еволюція НАП як процес Маркова. Матриці мутаційних даних PAM. Матриці BLOSUM. Емпіричні матриці кодонних заміщень і їхнє застосування в оцінці еволюції НАП. *Матеріали* – презентація лекції bioinf-L40.pdf. *Література*: [1, 2, 7].

**Самостійна робота** (5 год). Продовження вивчення можливостей пакета UGENE. Дотплот-аналіз.

#### Тижень 5

**Лекція 5** (2 год). **Попарне вирівнювання НАП.** Принцип графічного ілюстрування попарного вирівнювання НАП. Типи перебудов НАП, які можна виявляти за допомогою дотплот-аналізу – повтори, повні і часткові інверсії. Поняття “вікна” вирівнювання. Приклади програм відкритого типу для дотплот-аналізу на рівні окремих генів і геномів. Методи динамічного програмування у вирівнюванні НАП. Алгоритм локального вирівнювання Сміта-Уотермана з використанням унітарної матриці заміщень. Алгоритм глобального вирівнювання Нідельмана-Ванча. Порівняння рахунків вирівнювання НАП на основі унітарної матриці та BLOSUM62. *Матеріали* – презентація лекції bioinf-L50.pdf. *Література*: [1, 2, 3].

**Лабораторна робота 3** (2 год). Попарне вирівнювання – наявні сервіси. Сторінка уведення даних для програми BLAST. Результати програми BLASTP.

**Самостійна робота** (5 год). Ознайомлення з пакетом програм попарного вирівнювання що базуються на алгоритмах динамічного програмування глобального і локального вирівнювання: [https://fasta.bioch.virginia.edu/fasta\\_www2/fasta\\_list2.shtml](https://fasta.bioch.virginia.edu/fasta_www2/fasta_list2.shtml) . Як запит для білок-білкових вирівнювань використайте послідовність за № доступу GAW38200.

#### Тижень 6

**Лекція 6** (2 год). **Веб-сервіс BLAST.** Евристичні модифікації алгоритму локального попарного вирівнювання, що лежать в основі BLAST (Basic Local Alignment Search Tool) – “засівні слова порівняння”, афінні штрафи, пороги подібності. Статистична оцінка результатів BLAST –  $e$ ,  $p$ ,  $bits$ ,  $gaps$ . Родина програм BLAST – blastn, blastp, blastx, tblastn. PSI-BLAST – метод порівняння “профілів” білків. Структура початкової сторінки BLAST, її параметри за замовчуванням і можливості налаштування відповідно до мети

дослідження. Структура сторінки результатів BLAST. Приклади вирівнювання високоподібних і віддалених НАП. *Матеріали* – презентація лекції bioinf-L60.pdf. *Література*: [1, 2].

**Самостійна робота** (5 год). Робота з алгоритмами blastx, blastn на веб-порталі NCBI

### Тиждень 7

**Лекція 7 (2 год). Множинне вирівнювання НАП.** Концепція множинних вирівнювань НАП. Прогресивний принцип множинного вирівнювання. Інформація, яку надає множинне вирівнювання НАП. Глобальні і локальні множинні вирівнювання. Веб-сервіси, що надають послугу множинного вирівнювання – CLUSTAL W2/Ω, MUSCLE, T-COFFEE. Ілюстрування множинних вирівнювань. *Матеріали* – презентація лекції bioinf-L70.pdf. *Література*: [1, 2].

**Лабораторна робота 4** (2 год). Алгоритм blastx. Статистика попарних вирівнювань. Налаштування сторінки пошуку для виявлення малоподібних послідовностей.

**Самостійна робота** (5 год). Поглиблене ознайомлення з алгоритмічними основами множинних вирівнювань [https://www.ncbi.nlm.nih.gov/CBBresearch/Przytycka/download/lectures/PCB\\_Lect05\\_Multip\\_Align.pdf](https://www.ncbi.nlm.nih.gov/CBBresearch/Przytycka/download/lectures/PCB_Lect05_Multip_Align.pdf). Референтні бази множинних вирівнювань <http://www.lbgi.fr/wscoperr?Balibase&FileMoi&macsimHtml&BB11033>

### Тиждень 8

**Лекція 8 (2 год). Узагальнюючі моделі множинних вирівнювань** – консенсусний рядок, паттерни. Синтаксис паттернів. PROSITE. Прості профілі, паттерни і позиційно-специфічні матриці (PSSM/PSWM). Поняття зваженого рахунку позиції вирівнювання і псевдорахунку. Бази PSSM – CDD. Алгоритм PSI-BLAST. *Матеріали* – презентація лекції bioinf-L80.pdf. *Література*: [1, 2].

**Самостійна робота** (5 год). Ознайомлення з базою PROSITE, сервісом Conserved Domain Database <https://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>

### Тиждень 9

**Лекція 9 (2 год). Приховані моделі Маркова.** Генералізовані профілі. Концепція стану ознаки. Видимий шлях символів і прихований шлях станів. Принцип побудови й функціонування прихованої моделі Маркова (HMM) на прикладі аналізу 5'-ділянки екзон-інтронного переходу. Сервіси на основі HMM – HHPred, TMHMM, GeneMark, Pfam тощо. *Матеріали* – презентація лекції bioinf-L90.pdf. *Література*: [1, 3, 4].

**Лабораторна 5** (2 год). Сервіси множинного вирівнювання – MUSCLE, T-COFFEE.

**Самостійна робота** (5 год). Ознайомлення з програмами попарного вирівнювання, що функціонують на основі HMM – HMMER3 (<https://www.ebi.ac.uk/Tools/hmmer/>), HHPred (<https://toolkit.tuebingen.mpg.de/tools/hhpred>). Як запит для білок-білкових вирівнювань використайте амінокислотну послідовність за номером доступу ADL32277.

### Тиждень 10

**Лабораторна 5** (2 год). Вступ до філогенетичного аналізу. Вибір даних і моделі еволюції. Наявні онлайн-сервіси для вибору моделі еволюції (IQ-Tree) Письмовий контроль (модуль) за змістом перших 9 лекцій курсу.

**Самостійна робота** (10 год). Ознайомлення antiSMASH – HMM-опосередкованим сервісом виявлення генів вторинного метаболізму у геномах бактерій: <https://antismash.secondarymetabolites.org/#!/start>. Як запит для пошуку використайте номер доступу до нуклеотидної послідовності генома *S. albidoflavus* J1074 – NC\_020990.

### Тиждень 11

**Лекція 10** (2 год). **Молекулярна філогенетики – засади.** Концепція філогенетичного дерева, її біологічний зміст. Основні терміни – клада, нода, корінь, аутгруп, шкала дивергенції. Філогенетичний сигнал. Матеріал для аналізу – нуклеотидні, кодонні чи амінокислотні послідовності? Стратегії вибору масиву даних для філогенетичного аналізу й тлумачення результатів. Гомологи, паралоги, ортологи. Еволюційна модель у філогенетиці. *Матеріали* – презентація лекції bioinf-L100.pdf. *Література*: [1, 2, 6].

**Самостійна робота** (13 год). Як еволюція і філогенетика стосується мого повсякденного життя? – кожен зі студентів групи має вибрати і прочитати по одній статті зі списку, що є на веб-ресурсі університету Берклі: [https://evolution.berkeley.edu/evolibrary/search/topics.php?topic\\_id=15](https://evolution.berkeley.edu/evolibrary/search/topics.php?topic_id=15).

### Тиждень 12

**Лекція 11** (2 год). **Молекулярна філогенетика і філогеноміка.** Дистанційні і позиційні методи філогенетичного аналізу. Метод “з’єднання сусідів” (NJ). Метод максимальної вірогідності (ML). Статистична оцінка достовірності отриманих філогенетичних дерев – метод бутстрап-аналізу для методу NJ і aLRT – для ML. Філогеномний аналіз і систематика життя. Значення філогенетичних підходів у популяційній генетиці і судовій практиці. Аналіз 16S рРНК. Філогенетичний веб-сервер Phylogeny.fr. *Матеріали* – презентація лекції bioinf-L100.pdf. *Література*: [1, 2, 6, 7].

**Самостійна робота** (10 год). Ознайомлення з базою PROSITE, сервісом Conserved Domain Database <https://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>

### Тиждень 13

**Лекція 12** (2 год). Філогенія у межах одного виду/популяції – концептуальні відмінності від філогенії видів. Коалесцентна теорія. Фіксовані мутації між видами і поліморфізм у межах виду. Філогенетична реконструкція у вірусних популяціях, на прикладі вірусу імунодефіциту людини (HIV). Особливості біології HIV. Маркерні гени HIV. Філогенетична реконструкція HIV – глобальний рівень, між популяціями, у межах популяції, в одній особі. Про що свідчить топологія і довжина гілок дерева HIV? Практичне застосування філогенії HIV. *Матеріали* – презентація лекції bioinf-L100.pdf. *Література*: [1, 2, 6].

**Лабораторна 6** (2 год). Філогенетична реконструкція на сервері phylogeny.fr. Дерево RuBisCo на основі кількох масивів – окремо для покритонасінних та для нижчих рослин, з включенням імовірних білків RuBisCo бактерійного походження.

**Самостійна робота** (5 год). Ознайомлення з життям і науковим доробком українсько-американського вченого-генетика Теодозія Добжанського –

<https://www.youtube.com/watch?v=TH2AC8fu34M>. Сумісність біблійного та еволюційного вчень – погляд Т. Добжанського [https://www.pbs.org/wgbh/evolution/library/10/2/text\\_pop/1\\_102\\_01.html](https://www.pbs.org/wgbh/evolution/library/10/2/text_pop/1_102_01.html)

### Тиждень 14

**Лекція 13 (2 год). Ідентифікація кодувальних і операторних послідовностей.** Моделі прокариотичного і еукариотичного гена – і біологічна дійсність. Ген, відкрита рамка зчитування (orf), кодує послідовність, кодон. Виявлення кодуєчих послідовностей за гомологією – BLAST. Виявлення кодуєчих послідовностей *ab initio* – за рахунок порівняння частот вживання кодонів у досліджуваному гені і певному референтному геномі; за рахунок аналізу вживання нуклеотидів у третій позиції кодона. Врахування даних транскриптоміки у виявленні кодуєчих послідовностей. Програми GeneMark, PRODIGAL, GLIMMER. Пошук операторних послідовностей – програми RegPredict, MEME. Бази даних операторних послідовностей – TransFac тощо. *Матеріали* – презентація лекції bioinf-L110.pdf. *Література*: [1, 2, 7].

### Тиждень 15

**Лекція 14 (2 год). Аналіз білкових структур.** Класифікація білків. Поняття родини і фолду. Бази даних Pfam, SCOP. Тривимірні моделі білків – яку інформацію вони містять? PDB. Програма пошуку структурної гомології – HHpred. Веб-сервер ExPaSy для визначення основних параметрів білкових послідовностей та імовірних ділянок їхнього протеазного розщеплення і посттрансляційної модифікації. Програми для моделювання третинної структури білків і докінгу малих молекул. Веб-сервер STRING для аналізу функції гена у всій сукупності зв'язків з сусідніми генами і спорідненими геномами. KEGG, AlphaFold. *Матеріали* – презентація лекції bioinf-L120.pdf. *Література*: [6, 7].

**Лабораторна 7 (2 год).** Пошук операторів на сервісі MEME <https://meme-suite.org/meme/>

### Тиждень 16

**Лекція 15 (2 год). Аналіз РНК.** Виявлення рРНК й тРНК у геномах. Аналіз даних RNAseq. Бази даних тРНК. Передбачення вторинної структури РНК та оцінка її стабільності. Бази даних рРНК для потреб молекулярної таксономії. Бази даних некодуєчих РНК. Бази даних виявлення CRISPR-елементів у геномах бактерій. *Матеріали* – презентація лекції bioinf-L130.pdf. *Література*: [1, 2, 7].

**Лабораторна 8 (2 год).** Бази тривимірних структур білків. PDB. Phyre2 (<http://www.sbg.bio.ic.ac.uk/~phyre2/html/page.cgi?id=index>) – сервер передбачення тривимірних структур білків. Сервіси передбачення вторинних структур РНК - <http://rna.tbi.univie.ac.at/cgi-bin/RNAWebSuite/RNAfold.cgi>

### Література

1. Осташ Б.О. Біоінформатика: аналіз генетичних послідовностей. Електронний підручник. Видавничий центр ЛНУ ім. Івана Франка, 2022, 232 стор. ISBN 978-617-10-0729-1. Доступ онлайн: <http://dspace.lnlibrary.lviv.ua/handle/123456789/169>
2. Higgs PG, Attwood TK. Bioinformatics and Molecular Evolution. Blackwell Publishing, Oxford, 2005. 398 p. ISBN 1-4051-0683-2.

3. Durbin R, Eddy S, Krogh A, Mitchison G. Biological Sequence Analysis. Probabilistic Models of Proteins and Nucleic Acids. Cambridge University Press, Cambridge, 1998. 371 p. ISBN-13 978-0-521-62971-3
4. Borodovsky M, Ekisheva S. Problems and Solutions in Biological Sequence Analysis. Cambridge University Press, Cambridge, 2006. 362 p. ISBN-13 978-0-521-61230-2
5. Bioinformatics: a practical guide to the analysis of genes and proteins, 2<sup>nd</sup> Ed / AD Baxevanis, BFF Ouellette. – New York: John Wiley & Sons, 2001. – 455 p.
6. Allman ES, Rhodes JA. Mathematical Models in Biology. An Introduction. Cambridge University Press, Cambridge, 2003. 386 p.
7. Pevsner J. Bioinformatics and functional genomics. 3<sup>rd</sup> edition. Wiley Blackwell, London. – 2015- 1116 p. ISBN 978-1-118-58178-0.

## 8. Система оцінювання та вимоги

<b>Загальна система оцінювання курсу</b>	виконання завдань упродовж семестру/іспит – 50/50
<b>Вимоги до письмової роботи (модуль)</b>	За змістом перших дев'яти лекцій буде виконано поточний контроль знань у вигляді написання модуля. В модуль входять: визначення термінів (10 балів), два питання (по 6 б. кожне), одна схема чи таблиця, яку треба заповнити/зобразити (8 б.). Максимальна оцінка за модуль – 20 балів. Написання модуля обов'язкове.
<b>Лабораторні заняття</b>	Ще 30 балів студент може набрати упродовж семестру за виконання шести контрольних завдань упродовж лабораторних занять (Бази даних, BLAST, MSA, WebLogo, Phylo1, Phylo2).
<b>Умови допуску до підсумкового контролю</b>	Сумарний бал за модуль і за виконання контрольних завдань під час лабораторних занять має становити не менше 25, у такому співвідношенні: не менше 5 за модуль і не менше 20 за лабораторні.
<b>Іспит</b>	50 балів. Набір питань аналогічно до модуля; до термінів, питань і схем додаються тести. Письмова підготовка на протязі не більше 30 хв, далі усна відповідь. На іспит виноситься весь матеріал курсу

## 9. Навчальні методи

- Словесні (лекції, дискусії, пояснення)
- Лабораторні заняття
- За типом пізнавальної діяльності – проблемно пошукові, репродуктивні.

## 10. Перелік питань і типових задач, що виноситься на іспит

1. Моделі, що описують частоту зустрічності “слів” у генетичних послідовностях – Бернуллі і ланцюги Маркова.
2. Консервативні і неконсервативні заміщення у попарних вирівнюваннях, способи їхнього оцінювання; підходи до уведення у вирівнювання прогалин

3. PubMed – стратегії пошуку бібліографічних даних, значення ідентифікаторів статей
4. Структура флет-файлу GenBank
5. Національний центр біотехнологічної інформації США (NCBI) – структура і функції
6. Теорія прийнятних точкових мутацій (PAM) М. Дейгоф
7. Опишіть процес побудови позиційно-специфічних рахункових матриць (PSSM).
8. Алгоритми динамічного програмування у попарному вирівнюванні.
9. Основні терміни і теоретичні засади попарного вирівнювання
10. Опишіть етапи прогресивного принципу множинного вирівнювання
11. Моделі еволюції нуклеотидних послідовностей
12. Моделі Маркова в аналізі генетичних послідовностей
13. Способи опису складності НАП та районів з низькою складністю. Біологічне значення НАП низької складності.
14. Моделі заміщення в кодонних послідовностях
15. Опишіть основні елементи сторінки результатів BLAST
16. Підходи до уведення прогалин у попарні й множинні вирівнювання
17. Принципи побудови і використання прихованих моделей Маркова в аналізі НАП
18. Статистична оцінка попарних вирівнювань. Теорема Карліна-Альтшуля
19. Як розуміти унітарні й емпіричні системи обчислення рахунків вирівнювань з точки зору теорії інформації? Наведіть приклади. Теорема Альтшуля
20. У рівнянні з теореми Карліна-Альтшуля число очікування  $E$  прямо пропорційне простору пошуку (розміру бази даних) і обернено рахунку вирівнювання  $S$ . Тобто, у більших базах імовірність виявити випадкові збіги мала б зростати. Однак на практиці більші бази даних приводять до вирівнювань з дуже низьким  $E$ . Чому?
21. Які підходи використовують до оцінювання розривів (прогалин) у попарних вирівнюваннях?
22. Основні параметри опису попарного вирівнювання:  $S$ ,  $S'$ ,  $E$ ,  $P$ , ідентичність  $ID$ , подібність  $SI$
23. НАП з точки зору інформації: інформаційний вміст, символічне повідомлення, частотний розподіл символів, імовірність зустрічності символів/слів у послідовності.
24. Параметри налаштування програми BLAST
25. PAM-матриці і спосіб їхньої побудови
26. Опишіть принципи побудови механістичних та емпіричних матриць заміщень; що відмінного і спільного між ними.
27. Принцип побудови матриць BLOSUM, їхні відмітні властивості (порівняно з PAM).
28. Відмінні і спільні риси паттернів, позиційно-специфічних матриць і прихованих моделей Маркова
29. Чому не існує не одна, а серія матриць PAM або BLOSUM? Як обчислюють матрицю PAM250? Для яких потреб слід використовувати матрицю PAM20?
30. Основні властивості моделей еволюції нуклеотидних послідовностей.
31. Який біологічний зміст несе уведення розривів (прогалин) у попарні вирівнювання? Чи можна переставляти місцями позиції вирівнювання, і якщо так – то який біологічний процес відображає така маніпуляція НАП
32. Принципи дотплот-аналізу НАП
33. Що спонукало дослідників до розробки емпіричних підходів до оцінки вирівнювань амінокислотних послідовностей? Чому емпіричні підходи не набули поширення для нуклеотидних послідовностей?
34. Спільне і відмінне в попарному і множинному вирівнюваннях
35. Стисло опишіть моделі, які створюють на основі множинного вирівнювання НАП. Що в них відмінного?
36. Еволюційні засади попарного вирівнювання. Основні терміни.

37. Поясніть, що таке позиційно-незалежні та позиційно-специфічні рахунки вирівнювання, наведіть приклади їхнього використання в аналізі НАП
38. Концепція псевдорахунків у біоінформатиці – приклади її використання
39. Опишіть відомі вам способи множинного вирівнювання НАП
40. Принцип функціонування алгоритму пошуку оптимального глобального вирівнювання (Нідельмана-Ванча)
41. Відмінність між алгоритмами глобального і локального вирівнювання НАП
42. За якими ознаками певну модель чи процес можна віднести до Марковського?
43. Яку роль відіграють константи  $\lambda$  й  $K$  в рівнянні для обчислення числа  $E$ ?
44. Етапи побудови позиційно-специфічної рахункової матриці
45. Моделі оцінки частот символів/слів у НАП
46. Вичерпні та евристичні підходи до попарного вирівнювання НАП.
47. Порівняйте позиційно-специфічні матриці й приховані моделі Маркова як методи опису множинних вирівнювань.
48. ДНК- і білкові логотипи – принцип побудови, та інформація, яку він містить
49. У чому полягає складність побудови множинних вирівнювань? Які є способи оцінки їхньої якості?
50. Принцип функціонування алгоритму BLAST.
51. Опишіть, як зміна налаштувань програми BLAST впливає на результат пошуку гомологів?
52. Прогресивний та ітеративний принципи множинного вирівнювання
53. Які фактори ускладнюють безпосереднє використання даних про частоти появи мутацій у ДНК для коректного опису її еволюції?
54. Основні етапи побудови і використання позиційно-специфічних вагових матриць
55. Можливості й обмеження методів попарного і множинного вирівнювання для виявлення гомологічних НАП
56. Які дані містить сторінка результатів BLASTP?
57. Основні етапи пошуку оптимального локального вирівнювання за алгоритмом Сміта-Уотермана
58. Еволюційні засади філогенії
59. Основні поняття у галузі філогенетичної реконструкції
60. Яку інформацію містить філогенетичне дерево, а яку – не містить?
61. Основні елементи філогенетичного дерева, взаємозв'язок між ними
62. Вибір даних для філогенетичної реконструкції
63. Що таке модель еволюції у філогенетичній реконструкції і яке її значення?
64. Опишіть основні етапи філогенетичної реконструкції
65. Методи статистичної оцінки філогенетичних дерев
66. Підхід до філогенетичної реконструкції: з'єднання сусідів (NJ)
67. Підхід до філогенетичної реконструкції: максимальної ошадності (MP)
68. Підхід до філогенетичної реконструкції: максимальної вірогідності (ML)
69. Чому дерево-провідник з прогресивних підходів до множинного вирівнювання не є філогенетичним деревом?
70. У чому полягає суть курування вихідного множинного вирівнювання, що передусе філогенетичній реконструкції?
71. На чому ґрунтується суть пошуку мотивів ДНК?
72. Принцип функціонування алгоритму максимізації очікування (EM)?
73. Принцип функціонування програми MEME
74. Поясніть поняття родина, фолд, клас білка (за системою SCOP)
75. Поясніть принципи пошуку гомологів за допомогою програм BLAST й HHPred
76. Чому порівняння первинних і третинних структур генетичних послідовностей дає різний результат?

77. Практичне використання філогенетичного аналізу
78. В чому полягає принцип зважування генетичних послідовностей у моделях множинних вирівнювань?
79. Моделі заміщення в нуклеотидних і кодонних послідовностях, їхні особливості порівняно з моделями амінокислотних заміщень

### Задачі

1. Прочитайте такий паттерн: [HTS]-C-x-{P}-C-x(2)-C-{CP}-x(2)-C-[PEG]  
Що можна сказати про частоти вживання амінокислотних залишків у першій, другій та останній позиціях вирівнювання, з якого побудовано паттерн?
2. Обчисліть імовірність виявлення у геномі з GC-складом 60% такої послідовності: GGAATTCC
3. Як часто в геномі кишкової палички (GC склад – 50%) можна знайти послідовність GGGCCC? Послідовність TTGCAA?
4. На основі моделі Маркова 1-го порядку обчисліть імовірність виявлення у геномі з GC-складом 72% тетрануклеотиду AAAC (частоти динуклеотидів: AA – 0.016; AC – 0.057)
5. На основі моделі Маркова 1-го порядку обчисліть імовірність виявлення у геномі з GC-складом 70% тетрануклеотиду CCCG (частоти динуклеотидів: CC – 0.114; CG – 0.148)
6. Імовірність зустрічі тетрануклеотиду TCGG у геномі А становить 0.4%, а в геномі Б – 0.23%. Котрий із геномів має вищий GC-склад? Як можна дізнатися, який саме GC-склад вищезгаданих геномів (числове значення не потрібне, запропонуйте шлях розв'язку)?
7. Ендонуклеаза рестрикції EcoRI розпізнає паліндромну гексануклеотидну послідовність GAATTC. Яка імовірність виявити таку послідовність у геномі з GC-складом 72%? На скільки фрагментів буде гідролізовано такий геном, за умови що його розмір становить 6 млн п н?
8. Обчисліть імовірність зустріти у позиції  $n$  генома *Streptomyces coelicolor* залишок цитозину, за умови що у позиції  $n-1$  розташовано залишок аденіну. Для довідки: частоти динуклеотидів  $P(AA)$ ,  $P(AC)$ ,  $P(AG)$ ,  $P(AT)$  становлять 0.016, 0.057, 0.047 та 0.018, відповідно.
9. Обчисліть імовірність зустріти у позиції  $n$  генома *Streptomyces coelicolor* залишок тиміну, за умови що у позиції  $n-1$  розташовано залишок аденіну. Для довідки: частоти динуклеотидів  $P(AA)$ ,  $P(AC)$ ,  $P(AG)$ ,  $P(AT)$  становлять 0.016, 0.057, 0.047 та 0.018, відповідно.
10. Який інформаційний вміст (ентропію Шеннона) має одна позиція ДНК в геномі, що містить 80 % АТ-пар? Який максимально можливий інформаційний вміст позиції ДНК? Білка?
11. Внизу наведено множинне вирівнювання. Створіть паттерн на його основі

CUTD_DESAG/57-68	<b>CvHCGaCVpVCP</b>
CUTD_DESAG/89-100	<b>CpGCRrCEdVCP</b>
DGKE_HUMAN/135-146	<b>CfYCMvCKqVCG</b>
DGKE_MOUSE/132-143	<b>CsYCVfCRqVCG</b>
DMRX_METJA/159-170	<b>CkLCLkCInVCP</b>

12. Рахунок 4 для амінокислотних залишків валіну (V) та ізолейцину (I) у матриці PAM250 відображає співвідношення цільової частоти  $M_{ij}$  (імовірність зустріти V й I в одній позиції вирівнювання гомологічних білків) до фонові частоти  $V f_i$  ( =

імовірність, що ці два залишки будуть в одній позиції вирівнювання випадкових білків). Skorиставшись рівнянням М. Дейгоф для обчислення  $S$  (№ 25 у лекціях і № 19 у конспекті), обрахуйте, чому дорівнює це співвідношення для випадку V-I. Напишіть, у скільки разів частіше пара V-I зустрічається у вирівнюваннях гомологічних білків ніж випадкових.

13. Обчисліть рахунок пари амінокислот К/Е на основі підходу С. Хенікоф та Х. Хенікоф (BLOSUM, рівняння 27 у лекціях) за таких умов: цільова частота  $p_{KE} = 0.0041$ , фонові частоти  $p_K = 0.058$ ,  $p_E = 0.054$ .
14. У масиві даних, який використовували для побудови матриці BLOSUM62, пари лейцину L/L у попарних вирівнюваннях зустрічалися частіше ніж пари триптофану W/W ( $p_{LL} = 0.0371$ ,  $p_{WW} = 0.0065$ ). Водночас, триптофан – набагато рідкісніший ніж лейцин ( $f_L = 0.099$ ,  $f_W = 0.013$ ). Використавши ці дані а також рівняння з теореми Альтшуля (і прийнявши, що коефіцієнт  $\lambda$  у рівнянні становить 0.33), обчисліть рахунки  $S$  для пар W/W й L/L.
15. У вирівнюваннях гомологічних послідовностей пари A/L зустрічаються частіше ніж пари K/E ( $p_{AL} = 0.0044$ ,  $p_{KE} = 0.0041$ ), але A й L –поширеніші амінокислоти ( $p_A = 0.074$ ,  $p_L = 0.0099$ ,  $p_K = 0.058$ ,  $p_E = 0.054$ ). Обчисліть рахунки вирівнювання для цих двох пар на основі підходу С. Хенікоф та Х. Хенікоф (BLOSUM, рівняння 27 у лекції).
16. На еволюційній відстані 1PAM цільова частота заміщення аспартату (D) на залишок глутамату (E),  $M_{DE}$ , становить 0.0056, а на відстані 250PAM – 0.11. Обчисліть рахунок заміщення  $S_{DE}$  на відстанях 1PAM і 250PAM (фонові частоти аспартату – 0.047). Виходячи з отриманих даних, вкажіть, на якій еволюційній відстані заміщення аспартату на глутамат можна вважати прийнятною мутацією.
17. Цільові частоти заміщення  $M_{ij}$   $A \rightarrow E = 0.09$ , й  $E \rightarrow A = 0.05$ . Фонові частоти цих амінокислот становлять:  $f_A = 0.09$   $f_E = 0.047$ . Обчисліть рахунок пари  $S_{A,E}$  на основі підходу М. Дейгоф (рівняння №25, лекція).
18. У матрицях PAM і BLOSUM збігові двох ак залишків приписують максимальні значення, які однак відрізняються для різних пар амінокислот (тобто,  $S_{W,W} = 11$ ,  $S_{L,L} = 4$  тощо). За яких умов усі можливі ідентичні пари амінокислотних залишків мали б однаковий рахунок вирівнювання? Наведіть відповідь у математичній формі.
19. На рисунку наведено заповнену матрицю для алгоритма Нідельмана–Ванча. Відстежте оптимальний шлях вирівнювання, зобразіть відповідне попарне вирівнювання і визначте його рахунок  $S$ .

	G	T	A	C	G	T	C	G	G	
0	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
A	-3	-5	-8	2	-1	-4	-7	-10	-13	-16
T	-6	-8	3	0	-3	-6	4	1	-2	-5
A	-9	-11	0	11	8	5	2	-1	-4	-7
C	-12	-14	-3	8	19	16	13	10	7	4
A	-15	-17	-6	5	16	14	11	8	5	2
T	-18	-20	-9	2	13	11	22	19	16	13
G	-21	-10	-12	-1	10	21	19	17	27	24
T	-24	-13	-2	-4	7	18	29	26	24	22
C	-27	-16	-5	-7	4	15	26	37	34	31
T	-30	-19	-8	-10	1	12	23	34	32	29

20. Задано дві амінокислотні послідовності, ABCDEFG та ABCDEDEFG. Побудуйте дотплот-графік цих послідовностей. Яку генетичну перебудову ілюструє отриманий графік?
21. Задано дві амінокислотні послідовності: ALITTLE й ALITLLE. Виконайте їхнє локальне попарне вирівнювання за двох режимів обчислення рахунків. Перший: збіг +1, незбіг 0, нема штрафів за розрив (прогалина = 0). Другий: збіг +1, незбіг 0, штраф за відкриття прогалини -3, за продовження 0. (Тут прогалина на кінці

послідовності не рахується як розрив, а є незбігом – тому це локальне вирівнювання). Запишіть рахунки отриманих вами вирівнювань.

22. Розгляньте логотип НАП. Що він підсумовує? Які позиції абсолютно консервативні? Які містять приблизно однакові кількості двох різних амінокислот? Що означає вісь ординат?



23. Внизу наведено попарне вирівнювання. Обчисліть його рахунок  $S$  за такою схемою: збіги +4, незбіги -3, відкриття прогалини -5, продовження -0.1. Чи можна ці послідовності вирівняти краще (отримати вищий рахунок)?

AGCTTCGAC-C  
ACCTTCGACAC

24. Який мав би бути мінімальний біт-рахунок  $S'$  для пари послідовностей завдовжки 156 та 182 ак залишки, аби гарантувати відсутність випадкових вирівнювань з рахунком рівним або більшим  $S'$  ?
25. Порівнюють дві амінокислотні послідовності завдовжки 200 ак. Яке очікуване число випадкових вирівнювань послідовностей такого ж розміру можна отримати з біт-рахунком  $S' = 15$  ?
26. Гомолог заданої послідовності  $Y$  міститься у двох базах даних – SwissProt та PDB. Після виконання BLAST-пошуку у базі SwissProt вирівнювання  $Y$  з гомологом отримало значення  $E = 1.5$ . Таке саме вирівнювання при BLAST-аналізі бази PDB мало значення  $E = 0.075$ . Яке співвідношення розмірів вищезгаданих баз даних? Котра більша?
27. Порівнюють дві амінокислотні послідовності завдовжки 150 ак. Яким має бути біт-рахунок  $S'$  вирівнювання цих двох послідовностей, аби очікуване число випадкових вирівнювань послідовностей такого ж розміру з рахунком  $\geq S'$  було 0.001 ?
28. Задану послідовність Б використано для BLAST-пошуку гомологів у базі GenBank. Отримано хіт, що вирівнюється з Б, і це вирівнювання описується такими значеннями: біт-рахунок  $S' = 160$  число очікування  $E = 1.88e-35$ . Який розмір мав простір пошуку? Приймавши, що розмір послідовності Б становить 300 ак залишків, визначте розмір бази GenBank.
29. Задану послідовність Б розміром 400 ак залишків використано для BLAST-пошуку гомологів у геномах ссавців. Отримано три хіти, з такими рахунками вирівнювання  $S$  й асоційованими з ними значеннями  $E$ :  $S = 52, E = 0.011$ ;  $S = 48, E = 0.055$ ;  $S = 52, E = 0.011$ ;  $S = 40, E = 1.352$ . Які з вирівнювань можна вважати не випадковими? Відповідь обґрунтуйте. Обчисліть розмір бази даних, проти якої порівнювали послідовність Б, враховуючи що  $K = 0.15, \lambda = 0.4$ .
30. Яким має бути біт-рахунок  $S'$  вирівнювання заданої послідовності завдовжки 200 ак залишків із знайденою, виявленою у базі GenBank ( $96 \times 10^9$  залишків), аби це вирівнювання вважалося не випадковим?

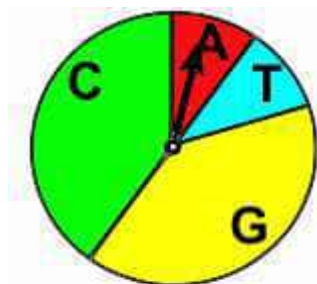
31. Послідовність  $W$  використано для BLAST-пошуку гомологів у базі GenBank ( $98 \times 10^9$  ак залишків). Отримано хіт, що вирівнюється з  $W$ , з біт-рахунком  $S'$  140 і числом очікування  $2.8e-29$ . Який розмір послідовності  $W$  ?
32. Якщо попарне вирівнювання, отримане при аналізі бази даних, має число очікування 0.0, то його слід розцінювати як випадкове чи не випадкове? Відповідь обґрунтуйте.
33. Задану послідовність  $R$  200 ак залишків завдовжки використано для BLAST-пошуку гомологів у базі розміром ( $80 \times 10^9$  ак залишків). Отримано хіт, що вирівнюється з  $R$ , і це вирівнювання має число очікування  $E = 0.0035$ . Який біт-рахунок  $S'$  має це вирівнювання?
34. Внизу наведено множинне вирівнювання. Побудуйте матрицю частот на його основі

GGATGCTAG  
 GTATCGTAG  
 CGATGGTAC  
 GCATACAAG  
 CGAAGGTAG  
 GCATTCTAT

35. Дослідникові треба гідролізувати кільцеву хромосому *Mycoplasma pneumoniae* на якомога більшу кількість фрагментів. Розмір хромосоми – 500 т п н, GC-склад – 42 %. У розпорядженні дослідника – дві ендонуклеази рестрикції (EP): *ApaI*, що розпізнає послідовність GGGCCC, й *EcoRI* – розпізнає GAATTC. Яку EP дослідникові слід вибрати? На яку кількість фрагментів очікується “розрізати” хромосому обраною EP?
36. Рисунок внизу зображує ланцюг Маркова, що генерує нуклеотидну послідовність із певним частотним розподілом символів (= нуклеотидів). Нехай задано дві послідовності:

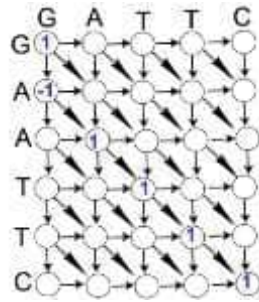
1: GATCGAATTCATTAATCTTA  
 2: GGGCACCTCCATGTTGGCCT

- Яку з цих послідовностей з більшою вірогідністю генеруватиме наведена модель? Обґрунтуйте відповідь, визначивши імовірності генерування послідовностей 1 і 2 за умови цієї моделі



$p_A=0.1, p_C=0.41, p_G=0.39, p_T=0.1$

37. Рисунок внизу показує оптимальний шлях вирівнювання двох нуклеотидних послідовностей. Зобразіть це вирівнювання у текстовому (попарному) форматі. Обчисліть рахунок цього попарного вирівнювання (збіг: +1, незбіг – 0, прогалина - 1)



38. На рисунку внизу зображено два попарні вирівнювання. Котре з них глобальне, а котре – локальне? Обчисліть рахунок вирівнювання обидвох за такою схемою: збіг +2, незбіг 0, штраф за відкриття прогалини -1, за продовження -0.1 (в локальному вирівнюванні визначте рахунок лише ділянок, що перекриваються: t-c --- g-c).

```
--T--CC-C-AGT--TATGT-CAGGGGACACG-A-GCATGCAGA-GAC
| | | | | | | | | | | | | | | | | | | | | | | | |
AATTGCCGCC-GTCGT-T-TTCAG-----CA-GTTATG-T-CAGAT--C
```

```
tccCAGTTATGTCAGggggacacgagcatgcagagac
| | | | | | | | | | | | | | | | | | | | | | | | |
```

```
aattgccgccgctcgtttttcagCAGTTATGTCAGatc
```

39. Для двох 40-нт послідовностей  $S_0$  та  $S_1$  внизу у таблиці підсумовано частоти  $S_1 = j$  та  $S_0 = i$ . Виведіть таблицю кондиційних імовірностей заміщення нуклеотидного залишка  $i$  у послідовності  $S_0$  на залишок  $j$  в  $S_1$  –  $P(S_1 = j | S_0 = i)$  – на основі даних вищенаведеної таблиці.

$S_1 \backslash S_0$	A	G	C	T
A	7	0	1	1
G	1	9	2	0
C	0	2	7	2
T	1	0	1	6

40. Внизу наведено дві вирівняні послідовності,  $S_0$  та  $S_1$ . Обчисліть загальну імовірність появи нових залишків тиміну у  $S_1$ , а також кондиційну –  $P(S_1 = T | S_0 = A)$ .

```
S0: ACGCCACTAGCTACAATCG
S1: ACCGCACTTGCTTCAATCG
```

41. Послідовність  $S_0$  складається із 40 п. н.  $S_1$ , що походить з  $S_0$ , містить 11 заміщень порівняно з предковою послідовністю. Обчисліть еволюційну відстань  $d$  між  $S_1$  й  $S_0$  за моделлю Джакса-Кентора (JC69).
42. Для множинного вирівнювання вибрано 5 послідовностей. На першому етапі програма CLUSTAL W виконала попарне вирівнювання усіх заданих послідовностей; ці результати наведено внизу у вигляді матриці відсотків ідентичності пар послідовностей. Виходячи з цієї матриці, намалюйте дерево-провідник, або опишіть словами, в якому порядку програма буде будувати множинне вирівнювання.

	1	2	3	4	5
1: Micromonospora	100.00				
2: Streptomyces	47.29	100.00			

3: Mycobacterium	46.06	49.37	100.00		
4: marine	40.49	40.49	49.37	100.00	
5: Microbacterium	40.12	40.00	50.95	70.03	100.00

43. Розгляньте два білкові логотипи. Які позиції абсолютно консервативні? Чому у першій позиції залишок валіну вищий (в обидвох логотипах) за залишок ізолейцину? Яка причина різної висоти літер у двох логотипах? Що означає вісь ординат?



44. Послідовність  $S_0$  складається із 100 п. н.  $S_1$ , що походить з  $S_0$ , містить 11 транзицій і 9 трансверсій порівняно з предковою послідовністю. Обчисліть еволюційну відстань  $d$  між  $S_1$  й  $S_0$  за моделлю Джакса-Кентора (JC69).
45. Внизу ліворуч наведено попарне вирівнювання білкових послідовностей, яке виконано на основі такої системи обчислення рахунку вирівнювання  $S$ : збіги і незбіги – BLOSUM62, штраф за відкриття прогалини -5. Обчисліть  $S$  цього вирівнювання. За яких умов було б вигіднішим вирівнювання цих двох послідовностей у спосіб, що зображено внизу праворуч?

```

ALDGTWSP
ALEGPTSP

```

```

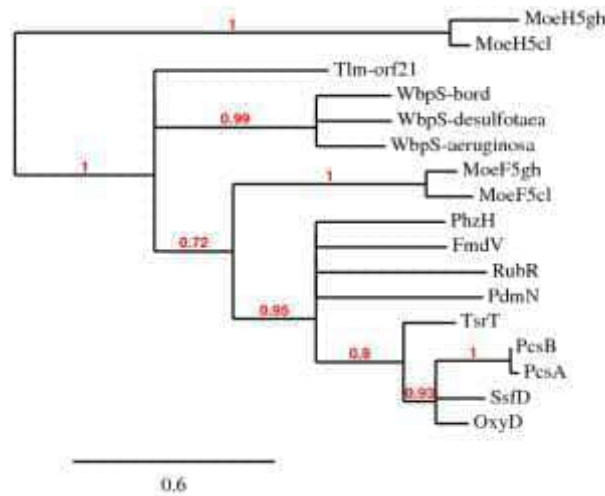
ALDGW-TSP
ALEG-PTSP

```

46. Уявімо геном, який в кожному раунді реплікації зазнає мутації з імовірністю 3%. Якою є імовірність не виявити мутацію через два раунди реплікації? Чи відрізнятиметься імовірність не виявити мутацію через два покоління від імовірності того, що мутацій на протязі двох раундів реплікації не було? Останню відповідь обґрунтуйте.
47. Уявімо, що розроблено медичний тест для скринінгу інфекційного захворювання. Якщо особу інфіковано, то тест у 99% випадків це виявляє, а якщо особу не інфіковано, то у 2% таких людей тест теж помилково виявляє хворобу. Тепер

припустимо, що 0.1% населення має інфекцію. Обчисліть імовірність того, що особа, яка пройшла тест, насправді має інфекційне захворювання.

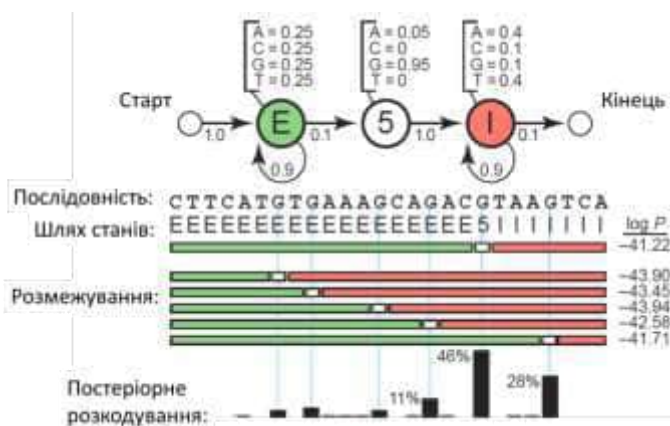
48. Дайте визначення основним елементам філогенетичного дерева, зображеного на малюнку



49. Внизу наведено типовий результат програми BLASTP. Поясніть зміст усіх термінів і позначень, наведених на ньому



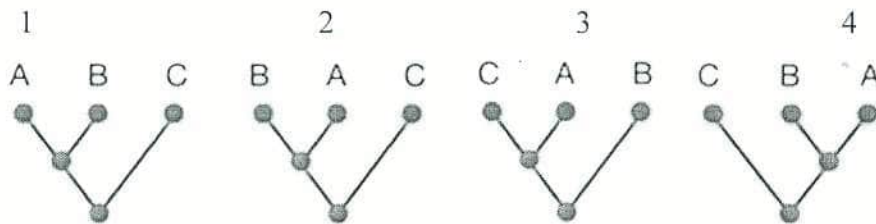
50. Поясніть, що зображено на малюнку, який наведено внизу. Як використовується математична модель, яку відображено на цьому малюнку?



51. Знайдіть оптимальний шлях вирівнювання за методом Нідельмана-Ванча для таких двох послідовностей

	C	O	E	L	A	C	A	N	T	H
P										
E										
L										
I										
C										
A										
N										

52. Внизу наведено 4 філогенетичні дерева, 1-4. Виберіть з них одне, що має топологію відмінну від решти трьох. Вибір поясніть



### 11. Опитування

Буде виконано в кінці курсу

Автор

Богдан ОСТАШ

"Погоджено"

Голова методичної ради  
біологічного факультету  
Віталій ГОНЧАРЕНКО

" 15 " березня 2023 р.

Гарант ОПП «Біохімія»  
Наталія СИБІРНА


15

березня 2023 р.

Гарант ОПП «Біофізика»  
Марта БУРА

2023 р.

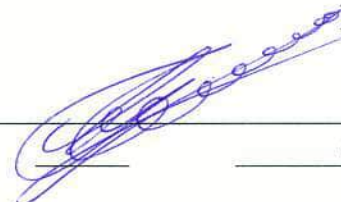
Гарант ОПП «Ботаніка»

  
\_\_\_\_\_ Анастасія ОДИНЦОВА  
\_\_\_\_\_ 15 \_\_\_\_\_ 03 \_\_\_\_\_ 2023 р.


Гарант ОПП «Генетика»

  
\_\_\_\_\_ Наталія ГОЛУБ  
\_\_\_\_\_ 15 \_\_\_\_\_ 03 \_\_\_\_\_ 2023 р.

Гарант ОПП «Зоологія»

  
\_\_\_\_\_ Андрій БОКОТЕЙ  
\_\_\_\_\_ 15/03 / \_\_\_\_\_ 2023 р.

Гарант ОПП «Мікробіологія»

  
\_\_\_\_\_ Світлана ГНАТУШ  
\_\_\_\_\_ 15.03 \_\_\_\_\_ 2023 р.

Гарант ОПП «Фізіологія людини і тварин»

  
\_\_\_\_\_ Оксана ІККЕРТ  
\_\_\_\_\_ \_\_\_\_\_ 2023 р.

Гарант ОПП «Фізіологія рослин»

  
\_\_\_\_\_ Наталія РОМАНІУК  
\_\_\_\_\_ 15 \_\_\_\_\_ Березня \_\_\_\_\_ 2023 р.