



PRESENTACIÓN

Análisis e interpretación de datos de alto rendimiento

Vivimos en una época de transformación sin precedentes en el ámbito de la genómica, lo que ha propiciado avances significativos en nuestra comprensión de la biología, la medicina de precisión y la investigación de enfermedades a nivel molecular. La reducción de los costos de secuenciación masiva (NGS) y el desarrollo de la biología computacional han sido pilares fundamentales en este progreso, permitiendo el análisis de grandes conjuntos de datos genómicos de manera eficiente y precisa. En respuesta a esta demanda creciente, este curso se enfoca en proporcionar los conocimientos esenciales para el procesamiento y análisis de datos de secuenciación masiva NGS, adoptando un enfoque práctico y aplicado.

Comenzaremos explorando la evolución y los principios básicos de las tecnologías de secuenciación, así como sus diversas aplicaciones en la investigación biomédica. Además, realizaremos visitas a instalaciones de secuenciación de última generación y laboratorios de preparación de muestras, incluyendo el Centro de Secuenciación y el clúster Urederra en Nasertic, que es uno de los nodos líderes en secuenciación y supercomputación (HPC) en Navarra.

La mayor parte de la computación práctica en este curso se llevará a cabo en un entorno de nube simulado, que emula un entorno HPC, donde todas las necesidades computacionales serán atendidas mediante el uso de entornos Docker y Conda. Además, aprovecharemos herramientas de inteligencia artificial como ChatGPT y GitHub Copilot para mejorar el aprendizaje y la creación de código asociado a la biología computacional.

Durante el curso, también adquiriremos habilidades esenciales en el manejo del entorno UNIX, utilizando comandos de una sola línea (one liners) para manipular y procesar datos de manera rápida y efectiva. Exploraremos el uso de entornos de gestión de paquetes como Conda y Mamba para configurar y mantener nuestras herramientas y librerías bioinformáticas. Además, nos familiarizaremos con la creación, configuración y despliegue de contenedores Docker, lo que nos permitirá construir ambientes de desarrollo reproducibles y portátiles. Asimismo, aprenderemos a manejar puertos y a generar salidas en formatos como HTML, facilitando la presentación y visualización de nuestros resultados.

Todo este conjunto de habilidades estará orientado a abordar grandes interrogantes biológicas mediante el análisis de datos masivos genómicos y transcriptómicos. Nos sumergiremos en el análisis de expresión génica diferencial y exploraremos los últimos avances en la cuantificación y variación de transcritos. Concluiremos el curso adquiriendo habilidades en lenguajes de flujo de trabajo bioinformático como Snakemake y Nextflow, fundamentales para la creación de pipelines y workflows en entornos HPC.

Al finalizar este módulo, los alumnos estarán capacitados para diseñar y analizar tecnologías modernas de NGS, trabajar eficientemente en entornos HPC-UNIX y preparar algoritmos para llevar a cabo todos los procesos necesarios en el análisis de datos genómicos. Esto les permitirá contribuir significativamente en la investigación biomédica y la medicina de precisión, generando conocimiento y respuestas a importantes preguntas en el campo de la biología molecular.



- **Carácter:** Optativa
- **ECTS:** 3
- **Curso y semestre:** 1º curso 2º semestre
- **Idioma:** Español. Se requieren conocimientos de inglés.
- **Título:** Máster en Métodos Computacionales en Ciencias.
- **Módulo 4** Optativo y **materia 4.1** Optatividad
- **Profesor responsable de la asignatura:** Igor Ruiz de los Mozos
- **Horario:** [Calendario del Máster](#)
- **Aula:** 1 edificio Los Castaños

COMPETENCIAS

	Análisis e interpretación de datos de alto rendimiento
	COMPETENCIAS BÁSICAS
CB6	Poseer y comprender conocimientos que aporten una base u oportunidad de ser originales en el desarrollo o aplicación de ideas, a menudo en un contexto de investigación.
CB7	Que los estudiantes sepan aplicar los conocimientos adquiridos y su capacidad de resolución de problemas en entornos nuevos o poco conocidos dentro de contextos más amplios (o multidisciplinares) relacionados con su área de estudio.
CB10	Que los estudiantes posean las habilidades de aprendizaje que les permitan continuar estudiando de un modo que habrá de ser en gran medida autodirigido o autónomo.
	COMPETENCIAS GENERALES
CG3	Conocer los principales problemas que se presentan en la adquisición y tratamiento de datos experimentales y cómo darles respuesta.
CG4	Comunicar tanto de manera oral como escrita un tema o datos de investigación en el área de las ciencias experimentales.



	COMPETENCIAS ESPECÍFICAS
CE5	Aplicar los métodos computacionales de procesamiento de datos a un problema científico particular de la disciplina de interés para el estudiante.
CE6	Diseñar un experimento científico para que sea rico en información, recogiendo gran cantidad de datos de manera estructurada que faciliten su procesamiento posterior.
CE8	Adquirir datos (bien en el laboratorio, o bien mediante minería on-line), organizarlos, filtrarlos, procesarlos, representarlos y refinarlos.
CE9	Extraer información de los datos con técnicas computacionales siguiendo un método científico.
CE10	Presentar los datos experimentales y la información científica de manera que se comuniquen de manera eficiente y fidedigna.
	COMPETENCIAS ESPECÍFICAS DE OPTATIVIDAD
CEOP5	Diseñar y analizar experimentos en el ámbito de las tecnologías ómicas, especialmente en el campo de la genómica y de la transcriptómica, con el fin de buscar nuevos biomarcadores, dianas terapéuticas o estudiar mecanismos moleculares.

PROGRAMA

Teoría de la secuenciación

- Evolución de las plataformas de secuenciación
- Desde las muestras hasta las librerías de secuenciación
- Desde la librería hasta las lecturas (reads)
- Estructura de los NGS (next-generation sequencing), control de calidad del pre-procesado



Universidad de Navarra

- Visión general de los flujos de trabajo bioinformáticos

Secuenciación práctica

- Visita al los centros de secuenciación masiva de Nasertic
- Experimentos RNA-Seq
- Secuenciación del exoma y genoma completo (Whole Exome/Genome)
- Secuenciación Single-cell

Unix y clusters de alto rendimiento

- Visita al centro de Super Computación Urederra de Nasertic
- Instalación de una máquina virtual en Linux en los computadores de los alumnos
- Curso intensivo de UNIX/bash
- Curso intensivo en Ordenadores de Alto Rendimiento (High-Performance Computers)

Alineamiento NGS con computación avanzada

- Teoría del alineamiento NGS, desde datos crudos multiplexados a lecturas alineadas.
- Mejorar los parámetros de alineamiento y estado del arte de los análisis NGS.

Detección de diferencias NGS (variant calling) y secuenciación genómica

- Traducción de las lecturas de secuenciación para obtener el genoma completo de interés
- Identificar mutaciones de nucleótidos puntuales que pueden convertirse en diferencias fenotípicas o virulentas

Comparación computacional del genoma

- Desde la secuencia consenso del DNA hasta la exploración genómica para inferir ganancia o pérdida de funciones.
- Especialización genómica y la evolución genómica.
- Análisis genómico en un caso de pandemia.

Análisis de RNA-seq

- Tecnologías transcriptómicas de secuenciación
- Cuantificación de la expresión génica
- Clusterización y caracterización de muestras, hierarchica y k-means
- Reducción de dimensiones PCA
- Análisis de expresión diferencial: teoría y práctica
- Integración y anotación de datos
- Análisis funcional: Gene Set Enrichment Analysis
- Corrección de desviación por lotes
- Cuantificación a nivel de transcrito, nuevos algoritmos sin mapeo genómico

Computación de alto rendimiento con workflows en Snakemake

- Iniciación a lenguaje de *workflows* en Snakemake.
- Análisis NGS en Snakemake.
- Análisis de reproducibilidad, trazabilidad y aprovechamiento de recursos en entornos HPC.



Universidad
de Navarra

ACTIVIDADES FORMATIVAS

Cada uno de los módulos tendrá un formato diferente. El primer módulo será más teórico y se realizarán visitas guiadas a los laboratorios. El segundo y tercer módulo serán 90% prácticos. En el trabajo final se espera que se apliquen los conceptos utilizados en los tres módulos.

EVALUACIÓN

Cada uno de los módulos valdrá un 20% de la nota (con 3 módulos 60%) y el último trabajo valdrá 40%.

HORARIOS DE ATENCIÓN

No hay horario específico para ello. Concertar cita previa por e-mail y se acuerda la fecha y hora de la tutoría:

Igor Ruiz de los Mozos: igor.ruiz@unav.es

BIBLIOGRAFÍA

[Bioinformatics Data Skills](#)

by [Vince Buffalo](#), July 2015, Publisher(s): O'Reilly Media, Inc. ISBN: 9781449367503

<https://github.com/Kur1sutaru/bioinformatics-data-skills/blob/main/book-bioinformatics-data-skills.pdf>

Next-generation Sequencing and Sequence Data Analysis

By: Ping Chiu, Kuo. Sharjah : Bentham Science Publishers. 2016. [Localízalo en la Biblioteca](#) [Recurso electrónico]

[Bioinformatics and functional genomics](#)

Pevsner, Jonathan. John Wiley & Sons, 2015. [Localízalo en la Biblioteca](#)

Genome Data Analysis

By: Kim, Ju Han. Singapore : Springer Singapore, 2019. XVI, 367 p. : 645 il., 236 il. col. Language: English, Base de datos: Catálogo de la Biblioteca de la Universidad de Navarra. [Localízalo en la Biblioteca](#) [Recurso electrónico]

Computational methods for next generation sequencing data analysis

Edited by Ion Măndoiu, Alexander Zelikovsky. [Localízalo en la Biblioteca](#) [Recurso electrónico]

Biostar Handbook Student Edition



Universidad
de Navarra

Dr. István Albert, An introduction to Bioinformatics as a scientific field of study. Covers all aspects of bioinformatics. <https://biostar.myshopify.com/>