

# From Prompt to Pipeline: Large Language Models for Scientific Workflow Development in Bioinformatics

Khairul Alam and Banani Roy

Department of Computer Science, University of Saskatchewan, 105 Administration  
Pl, Saskatoon, S7N 5A2, Saskatchewan, Canada.

\*Corresponding author(s). E-mail(s): [kha060@usask.ca](mailto:kha060@usask.ca);  
Contributing authors: [banani.roy@usask.ca](mailto:banani.roy@usask.ca);

## Abstract

**Purpose:** The increasing complexity of bioinformatics data analysis necessitates the use of Scientific Workflow Systems (SWSs) such as Galaxy and Nextflow to support scalable, reproducible, and automated scientific workflows. However, developing and understanding these workflows remains challenging for domain scientists (e.g., Bioinformaticians, Biologists, or Geneticists), particularly those without programming experience. This difficulty arises from the reliance on numerous black-box tools/modules and the complex infrastructure required to work with SWSs. This study investigates whether modern Large Language Models (LLMs) can assist in generating accurate, complete, and usable bioinformatics workflows and identifies the prompting strategies that maximize their effectiveness. The goal is to assess their potential to lower the barrier for domain specialists, enhance reproducibility, and improve the accessibility, usability, and maintainability of bioinformatics workflows.

**Methods:** We evaluate three state-of-the-art LLMs, GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3, across three research questions: (1) Can LLMs assist in scientific workflow development? (2) How do LLM-generated workflows compare in terms of completeness, correctness, and usability? (3) What prompting strategies best support workflow generation? We select a set of representative workflows across Galaxy (graphical-based) and Nextflow (script-based), covering diverse bioinformatics tasks such as SNP analysis, RNA-seq, DNA methylation, and data retrieval. The generated workflows are manually evaluated by domain experts for correctness, completeness, tool appropriateness, and executability, using established baselines from the Galaxy Training Network (GTN) and nf-core repositories. To ensure consistency and reduce complexity, all models are used with default settings, and prompting strategies, including instruction-only, role-based, and chain-of-thought, are applied uniformly across tasks to maintain comparability while preserving workflow quality.

**Results:** LLMs demonstrated varying levels of proficiency across tasks. Gemini 2.5 Flash produced the most accurate and user-friendly workflows for Galaxy SWS, while DeepSeek-V3 performed well for Nextflow. GPT-4o performed best with structured prompts that provided explicit context, while DeepSeek-V3 offered rich technical outputs, including command-line instructions and directory structures, albeit with verbosity and occasional inconsistencies. Prompting strategies significantly influenced output quality, with role-based and chain-of-thought prompts improving completeness and correctness. Overall, LLMs demonstrated promising capabilities in lowering the barrier to scientific workflow development, with prompting strategy and model selection playing critical roles in effectiveness.

**Conclusion:** This study demonstrates that LLMs can assist scientific workflow development, particularly when guided by well-crafted prompts. By translating natural language into workflow, these models can reduce barriers for both novice and expert users, enhance reproducibility, and accelerate scientific discovery. Future research should expand to other domains, explore adaptive prompting strategies, and develop integrated tooling that leverages LLMs more effectively in real-world scientific environments.

**Keywords:** Scientific Workflows, Large Language Models, Bioinformatics, GPT-4o, Gemini 2.5 Flash, DeepSeek-V3, Prompt, Galaxy, Nextflow

# 1 Introduction

The exponential growth of biological data and the increasing complexity of computational analyses have made scientific workflows [1] indispensable in modern bioinformatics [2]. These workflows (workflow means scientific workflow for the rest of the paper) are essential for enabling scalable, reproducible, and dependable large-scale data analysis, particularly when deployed on distributed and cloud infrastructures [3, 4]. In bioinformatics, workflows routinely process terabyte-scale datasets generated by advanced DNA, RNA, NGS sequencing, and other technologies [4] across diverse experimental contexts [5]. Typically, these workflows comprise multiple computational steps, including data preprocessing, quality control, sequence aggregation, machine learning for classification or clustering, statistical analysis, and result visualization. Each step is carried out by specific tools/modules, often developed by third parties, facilitated by scientific workflow systems (SWSs)’ no/low-code environments [6]. SWSs have become essential platforms for constructing, executing, and managing bioinformatics workflows by integrating diverse analytical and data processing tools/modules into reproducible pipelines (workflows)[7]. SWSs such as Galaxy [8], Nextflow [9], Snakemake [10], Airflow [11], and others [2, 4, 12, 13] provide structured environments that support the automation, customization, and reuse of complex analyses. In response to the growing demand for scalable and transparent data processing, hundreds of SWSs [14] have emerged, offering access to thousands of bioinformatics tools/modules. For instance, Galaxy’s ToolShed alone hosts over 10,500 tools [15], Nextflow has over 1500 modules, enabling researchers to assemble rich and versatile workflows tailored to their specific needs. These platforms have significantly advanced reproducibility and efficiency in computational biology by facilitating workflow sharing, versioning, and collaboration [16].

Despite the widespread adoption of SWSs, workflow development remains a significant challenge for both novice and expert users [17–20]. Practitioners, particularly those without programming backgrounds, face steep learning curves due to the need for scripting, tool configuration, dependency management, and infrastructure provisioning. Workflow developers are usually domain scientists with deep subject-matter expertise but limited experience in software engineering or distributed/cloud computing. While SWSs offer a broad array of services, such as data cleaning, statistical analysis, modeling, and high-performance computing, these capabilities are typically leveraged effectively only by a minority of users who are proficient in workflow composition, programming, and infrastructure management [21]. Consequently, many researchers struggle to design workflows to manage complex data and operate on these SWSs. These difficulties stem from the diversity of tools/modules, scripts, and data formats essential to ensure data integrity, reproducibility, and scalability. The resulting complexity hampers workflow development, adaptation, and interpretation, ultimately impeding data exploration and scientific discovery [20, 22].

Recent advancements in LLMs such as GPT-4o [23], Gemini 2.5 Flash [24], and DeepSeek-v3 [25] have demonstrated their exceptional ability to comprehend complex instructions and generate accurate, context-aware code, making them increasingly relevant in software engineering, particularly in code generation [26]. With their enhanced language comprehension and generative capabilities, they are being extensively explored to automate various software design and development aspects [27]. Despite these promising developments, research exploring the application of LLMs in workflows, especially through practical, real-world scenarios in bioinformatics, remains sparse. This gap highlights the need for studies investigating how LLMs can support practitioners in workflow development. Given their ability to generate code, we posit that LLMs can similarly generate accurate workflows with effective prompts. Thus, we plan to investigate whether LLMs can assist in workflow development. We choose GPT-4o, Gemini 2.5 Flash, and DeepSeek-v3 as they are among the most advanced, publicly accessible LLMs as of 2025. GPT-4o is known for its strong general-purpose capabilities, broad tool knowledge, and superior performance in code generation tasks. Gemini 2.5 Flash has demonstrated strengths in structured reasoning, contextual understanding, and system integration, making it a promising candidate for guiding workflow development. DeepSeek-v3, a newer open-source model trained with a code-heavy corpus, shows competitive performance in scientific and software engineering tasks, and offers an opportunity to evaluate performance beyond proprietary models.

To evaluate the capabilities of LLMs in assisting workflow development, we focus on two widely adopted and most popular SWSs for the bioinformatics domain: *Galaxy* [8] and *Nextflow* [9]. These platforms represent complementary paradigms; Galaxy offers a user-friendly, graphical interface tailored for domain scientists with limited programming experience, while Nextflow provides a script-based environment optimized for developers seeking scalability and reproducibility. Galaxy is known for its accessibility, extensive tool repository via ToolShed, and emphasis on transparency and reproducibility through automatic logging of analysis steps. It supports deployment across diverse infrastructures, including local servers, institutional clusters, and cloud platforms. In contrast, Nextflow enables the composition of complex, modular pipelines inspired by Unix principles, with built-in support for containerization (Docker, Singularity) and distributed and cloud computing environments (e.g., Kubernetes, AWS Batch, Slurm). Its portability and integration with version control and software packaging tools ensure the principle of *write once, run anywhere*. Together, Galaxy and Nextflow provide a robust foundation for assessing the generalizability, usability, and automation potential of LLMs in workflow development across both GUI-driven and script-driven environments.

Integrating LLMs into bioinformatics workflow development presents a compelling opportunity to reduce the cognitive and technical burden on workflow developers. With appropriate prompting, LLMs could potentially help users generate workflows, configure parameters, select tools/modules, and even interpret error messages. However, concerns persist regarding the reliability of generated outputs, especially in scientific contexts where precision and reproducibility are critical. To ensure the effectiveness of our approach, we have meticulously applied state-of-the-art prompt engineering strategies, drawing on insights from prominent studies in the field [28–34].

In this study, we conduct a comprehensive evaluation of state-of-the-art LLMs (*GPT-4o*, *Gemini 2.5 Flash*, and *DeepSeek-V3*) in their ability to support scientific workflow development in bioinformatics. We begin by assessing the capacity of these models to explain foundational concepts related to scientific workflows and SWSs, followed by an in-depth evaluation of their ability to develop workflows across two widely used platforms: *Galaxy* and *Nextflow*. We then analyze the generated workflows for correctness, completeness, and usability, comparing them against community-curated baselines (Galaxy Training Network [35] and nf-core [36]). Furthermore, we investigate how different prompting strategies (e.g., instruction-only, role-based, and chain-of-thought) affect output quality and workflow fidelity. By identifying recurring challenges and model-specific limitations, we propose directions for improving LLM-supported workflow development, including integrating domain-specific knowledge and refining prompt engineering techniques. This work contributes to the emerging intersection of natural language processing and computational bioinformatics by offering practical insights into how LLMs can bridge the gap between domain scientists and automated workflow generation. Specifically, we aim to answer the following three research questions (RQs):

- **RQ1:** To what extent can LLMs (e.g., GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3) assist in developing scientific workflows?
- **RQ2:** How do workflows generated by different LLMs compare regarding completeness, correctness, and usability?
- **RQ3:** What prompting strategies should a workflow developer follow?

**Paper Organization:** The remainder of this paper is outlined as follows: Section 2 provides background on Scientific Workflows, Scientific Workflow Systems (SWSs), Large Language Models (LLMs), and prompting techniques, establishing the foundation for our investigation. Section 3 outlines our study design, including the selection of workflow platforms, LLMs, and evaluation criteria. In Section 4, we present the findings of our analysis, followed by an in-depth discussion and interpretation in Section 5. Section 6 examines potential threats to the validity of our results, considering both methodological and contextual limitations. Section 7 situates our research within the broader landscape of related efforts. Finally, Section 8 concludes the paper, highlighting directions for future research.

## 2 Background

Scientific workflows combine multiple software artifacts on distributed stacks for advanced data analysis [3]. We ground the reader by providing a brief overview of scientific workflows and workflow systems and introducing LLMs, including their utility to facilitate the creation of workflows.

## 2.1 Scientific Workflows

A scientific workflow is a structured and reproducible sequence of data processing tasks, typically modeled as a directed acyclic graph (DAG), where nodes represent computational operations and edges indicate data dependencies [1, 37]. It automates complex scientific procedures, such as data acquisition, transformation, and analysis, thereby accelerating discovery through efficient and reliable execution [38]. Analogous to a cooking recipe, where ingredients, preparation steps, and instructions combine to produce a consistent outcome, scientific workflows ensure reproducibility and scalability in computational experiments. Lin et al. [39] defined workflows as the computerized automation of scientific processes that streamline both tasks and dataflows. In this paper, we use the terms *workflow*, *scientific workflow*, and *pipeline* interchangeably to refer to such structured processes. A workflow can be graphic-based and script-based. Figure 1 illustrates an example of an *RNA sequencing* workflow steps developed using *Nextflow* SWS, obtained from [36]. This pipeline takes FASTQ files as input, performs quality control and trimming, aligns or pseudo-aligns reads to the reference genome, quantifies gene and transcript expression, and generates a comprehensive expression matrix alongside detailed QC reports.

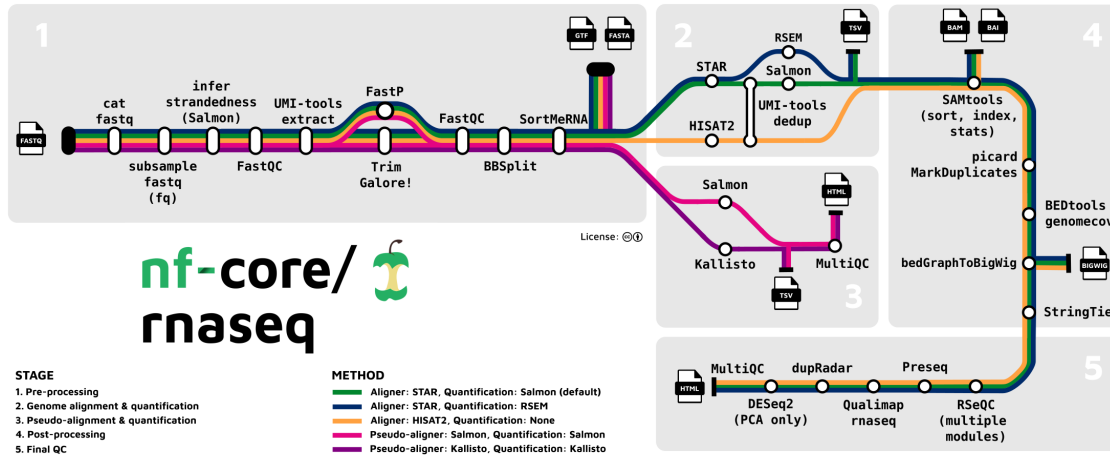


Fig. 1: RNA sequencing analysis pipeline (workflow)

## 2.2 Scientific Workflow Systems

SWSs are specialized software platforms designed to automate, manage, and execute complex sequences of computational tasks, where the execution order is controlled by a formal representation of workflow logic [39]. In bioinformatics, SWSs play a critical role in handling large-scale, heterogeneous data by enabling reproducible, scalable, and modular analysis pipelines. Prominent systems such as Galaxy [8], Nextflow [9], Taverna [40], Snakemake [10], Apache Airflow [11], and Pegasus [12] have been developed to support domain scientists in composing and executing workflows for a wide range of bioinformatics tasks, including genome assembly, variant calling, RNA-seq analysis, and metagenomics. The increasing reliance on workflow-based analysis in computational biology has led to the emergence of hundreds of SWSs [14], each offering varying degrees of abstraction and control. Core features commonly supported by these platforms include graphical or script-based workflow design, task scheduling and execution, real-time monitoring, provenance tracking for reproducibility, error handling, and integration with high-performance or cloud computing infrastructures [41]. SWSs thus provide an essential foundation for managing the complexity and scale of modern bioinformatics research.

## 2.3 Large Language Models

LLMs are advanced AI systems capable of processing and generating human-like text using vast training data. Their development has been revolutionized by the transformer architecture introduced by Vaswani et al. in *Attention Is All You Need* [42]. Leveraging neural networks and transformer

variants, LLMs like GPT-3, GPT-4, Codex (OpenAI), DeepSeek-V3 (DeepSeek), Gemini, PaLM (Google), LLaMA (Meta), and Claude (Anthropic) and others [43] have become cornerstones of natural language processing. Tools like Codex and ChatGPT assist developers by suggesting code snippets, completing functions, and enhancing productivity [44]. They aid in code review, bug detection, and quality assurance, promoting reliable software [45]. Additionally, they generate and update documentation, user guides, and API references, simplifying software maintenance [46, 47]. LLMs even support cross-platform development by translating code between languages, improving collaboration among diverse teams [48, 49]. LLMs are widely acknowledged for their ability to generate commonsense knowledge, leveraging the vast training data [3]. This capability makes them practical and accessible knowledge repositories that can be queried through natural language prompts. A key limitation of LLMs is their tendency to produce *hallucinations* [50, 51].

## 2.4 Usages of LLMs

Recent advances in LLMs have sparked significant interest in automating various aspects of software development, including code generation, refactoring, and debugging. Numerous studies have evaluated LLMs for their ability to generate syntactically correct and semantically relevant code from natural language prompts [52–54]. Beyond code generation, LLMs have also shown promise in interactive programming tasks such as modifying code based on user intent, understanding incomplete problem statements, and suggesting appropriate development patterns [55]. LLMs like Codex, ChatGPT, and GitHub Copilot have demonstrated practical utility in supporting programmers by reducing the need for frequent online searches and offering useful scaffolds for common tasks. However, users often encounter challenges related to debugging, editing, and verifying model-generated code, particularly for more complex programming tasks [56]. Despite these limitations, studies show LLMs can be beneficial in educational settings, improving student learning outcomes and reducing cognitive load during programming exercises [57]. In comparative evaluations, ChatGPT has outperformed other code generation tools in terms of correctness and technical debt, highlighting the growing maturity of general-purpose LLMs in software development [58].

While LLMs have been widely studied in the context of general software engineering, their application in scientific data analysis and workflow automation, particularly in domains like bioinformatics, remains relatively nascent. Initial efforts by Liu et al. [59] explored how code generators like Codex could assist non-programmers in performing data analysis tasks. Similarly, Sanger et al. [3] conducted a qualitative assessment of ChatGPT’s effectiveness in supporting scientific workflow design. More recently, Xu et al. [6] introduced LLM4Workflow, a framework that integrates LLMs into a workflow model. Building on these foundational studies, our research investigates the underexplored potential of LLMs in the domain of bioinformatics workflow development. Specifically, we assess how state-of-the-art LLMs can generate, explain, and refine workflows for scientific tasks using two widely adopted SWSs: Galaxy and Nextflow. By systematically analyzing the influence of different prompting strategies, such as instruction-only, role-based, and chain-of-thought, we aim to unlock more effective and reliable uses of LLMs in enabling non-programmers and domain scientists to construct robust, reproducible, and scalable workflows. This work contributes to bridging the gap between advanced language models and practical scientific computing applications.

## 2.5 Prompt Engineering

Prompt engineering refers to the process of crafting input instructions that effectively guide an LLM to generate desired outputs. A prompt typically consists of natural language text that conveys a task, includes relevant context, and may specify output format or style [33]. Since LLMs are sensitive to how information is framed, the structure, clarity, and specificity of a prompt significantly influence the relevance, accuracy, and usefulness of the model’s response [60]. Prompt engineering has therefore emerged as a critical practice for aligning LLM behavior with user goals, especially in complex domains such as scientific workflow development.

In the context of this study, prompt engineering plays a pivotal role in enabling LLMs to understand domain-specific requirements in bioinformatics and translate them into structured workflows. Drawing from recent research that emphasizes prompt design as a programming paradigm in itself [61], we explore strategies [28–34, 62–73] for adapting prompts to scientific tasks, such as generating analysis pipelines or interpreting bioinformatics procedures. These strategies include varying prompt templates, combining task descriptions with expected outcomes, and incorporating domain-relevant terminology to improve performance across different use cases [74].



To systematically investigate prompt effectiveness, we evaluate three widely used prompting techniques: (1) *instruction-only* prompting, where the task is explicitly described in a natural language instruction; (2) *role-based* prompting, where the model is asked to assume the perspective of a domain expert; and (3) *chain-of-thought* prompting, which encourages step-by-step reasoning. These three techniques represent diverse and widely adopted paradigms in prompt engineering, each targeting a different axis of LLM interaction: clarity of task (instruction-only), context emulation (role-based), and reasoning process (chain-of-thought). Instruction-only prompts are the most accessible and commonly used format in practical applications, offering a strong baseline for comparison. Role-based prompts introduce a layer of contextual alignment by asking the model to simulate expert knowledge, which is particularly relevant in specialized domains like bioinformatics. Meanwhile, chain-of-thought prompting has been shown to enhance performance on tasks that require multi-step reasoning and logical coherence, traits that mirror the sequential and conditional structure of scientific workflows. Together, these three approaches span a broad spectrum of prompting complexity and cognitive alignment, allowing us to meaningfully assess the impact of prompt design on workflow generation quality. By focusing on these strategies, we aim to provide actionable insights that are both theoretically grounded and practically relevant for researchers and developers using LLMs in bioinformatics and beyond. As scientific workflows often involve multi-step reasoning, domain-specific knowledge, and structured output, prompt engineering becomes essential in bridging the gap between high-level user intent and executable computational pipelines. By comparing prompt outcomes across different LLM platforms, our study contributes empirical insights into the design of effective prompts for scientific workflow generation. This reinforces the notion that prompt engineering is not merely a preliminary step but a central component in harnessing LLMs for domain-specific scientific computing.

### 3 Study Design

This study investigates the capabilities of modern LLMs in generating scientific workflows for bioinformatics through natural language prompts. We evaluate the performance of three state-of-the-art LLMs, GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3, in translating high-level user instructions into workflows within two widely used SWSs: Galaxy and Nextflow. Our experimental design focuses on three main research questions: **(RQ1)** *How effectively can LLMs support workflow development?*, **(RQ2)** *How do their outputs compare in terms of quality?* **(RQ3)** *Which prompting strategies yield the best results?*. Through this investigation, we aim to assess the technical capabilities of LLMs and identify best practices in prompt formulation that can support novice and expert users in constructing robust, domain-relevant bioinformatics workflows.

We select *Galaxy* and *Nextflow* as representative SWSs due to their widespread adoption within the bioinformatics community [4, 7] and complementary design paradigms. Galaxy provides a user-friendly, graphical web interface that enables researchers, particularly those without programming backgrounds, to design and execute workflows through intuitive drag-and-drop functionality. On the other hand, Nextflow adopts a script-based approach, emphasizing high reproducibility, scalability, and seamless integration with distributed and cloud computing environments. Its flexibility and support for modular, version-controlled workflows make it a preferred choice for expert users engaged in large-scale, automated analyses. For model selection, we consider three advanced LLMs with diverse strengths. *GPT-4o* is chosen for its demonstrated superiority in code generation and reasoning tasks across multiple domains. *Gemini 2.5 Flash* is included for its efficiency in producing structured outputs and handling multi-step instructions, traits particularly relevant for workflow synthesis. *DeepSeek-V3* is selected for its strong performance in open-source benchmarks and its alignment with code-centric modeling, making it a competitive and transparent alternative to proprietary models. We evaluate each model using its default configuration settings, unless explicitly stated otherwise. This decision reflects the perspective of domain scientists, who may not possess expertise in adjusting model parameters. Requiring manual configuration could introduce unnecessary complexity and create additional barriers to practical use in real-world scientific settings.

We select ten bioinformatics workflows that represent a diversity of real-world data analysis tasks and domain-specific complexities (Table 1). These workflows are sourced from two highly regarded and widely used community platforms: the *Galaxy Training Network (GTN)* and the *nf-core* initiative [36, 85]. GTN workflows offer clearly documented, step-by-step tutorials that cover core bioinformatics applications such as *RNA-seq*, *variant calling*, and *quality control*. These workflows are ideal for evaluating how well LLMs can assist practitioners in reproducing structured analyses from natural language descriptions. In contrast, *nf-core* workflows are community-curated, production-ready

**Table 1:** Overview of the used workflows. (W  $\rightarrow$  Workflow)

SWS	WID.	Workflow	Link
Galaxy	W1	Which coding exon has the highest number of single-nucleotide polymorphisms on human chromosome 22?	[75]
	W2	How to get from peak regions to a list of gene names?	[76]
	W3	Do genes on opposite strands ever overlap? If so, how often?	[77]
	W4	How can the quality of NGS raw data be assessed, what parameters should be checked, and how can the quality of the dataset be improved?	[78]
	W5	Which genes are on a draft bacterial genome, and which other genomic components can be found on a draft bacterial genome?	[79]
Nextflow	W1	nf-core/demo: This bioinformatics pipeline performs quality control and reporting on sequencing data.	[80]
	W2	nf-core/fetchngs: It is a bioinformatics pipeline to fetch metadata and raw FastQ files from databases.	[81]
	W3	nf-core/methylseq: This workflow is used for Methylation (Bisulfite) sequencing data.	[82]
	W4	nf-core/rnaseq: This workflow is used to analyse RNA sequencing data with a reference genome and annotation.	[83]
	W5	nf-core/phyloplace: This bioinformatics performs phylogenetic placement with EPA-NG.	[84]

pipelines built using Nextflow. They target high-throughput, reproducible research in environments that demand robustness, modularity, and scalability. Each nf-core pipeline adheres to strict guidelines, including automated testing, documentation standards, and best practices in computational reproducibility. Selecting workflows from nf-core allows us to evaluate LLM performance in replicating more sophisticated, script-based analyses essential for bioinformaticians. By incorporating workflows from both GTN and nf-core, we ensure our evaluation captures a broad spectrum of user expertise, computational complexity, and bioinformatics domains. This dual-source approach enables a comprehensive assessment of LLMs’ capabilities across graphical and code-based environments.

To address **RQ1**, we investigate the extent to which state-of-the-art LLMs can support the understanding and development of workflows in the bioinformatics domain. Our evaluation targets two essential capabilities: (1) *foundational understanding of workflows and workflow systems*, and (2) *construction of complete, domain-appropriate workflows*. To assess this, we begin by posing a set of fundamental questions (Table 2) that probe LLMs’ conceptual grasp of workflows, SWSs, and their core components, including tasks, tools, and data flow. This forms a baseline for determining whether models can meaningfully reason about workflows before generating them. We extend this analysis by evaluating each model’s comprehension of *platform-specific background knowledge* for our selected SWSs, using targeted questions (Table 3) tailored to each platform’s functionality, benefits, and architectural components. Additionally, we design *workflow-specific background prompts* (Tables 4 and 5) that correspond to representative bioinformatics workflows (Table 1). These background questions reflect the type of biological and computational knowledge a domain scientist needs to understand before workflow creation and are essential for evaluating the contextual readiness of LLMs to generate accurate workflows. This step served two purposes: (1) *to confirm the LLM’s baseline understanding of the domain context*, and (2) *to prime the model with relevant knowledge before executing the workflow prompt*. These background questions are curated based on domain expertise and established tutorials. While it is possible to pose many more questions to comprehensively assess LLM knowledge, our chosen set strikes a balance between breadth and depth. These questions are carefully curated to reflect key conceptual, platform-specific, and workflow-specific knowledge areas that a typical bioinformatics practitioner would need to grasp prior to workflow development. By focusing on essential topics, such as core workflow components, the principles of Galaxy and Nextflow, and the biological and computational context of each selected workflow, we ensure that the evaluation remains both focused and meaningful. Our goal is not to exhaustively test all possible knowledge domains, but to assess whether LLMs can understand and reason through the types of background information most relevant for supporting real-world workflow construction. This strategic sampling allows us to derive practical insights while keeping the evaluation tractable and grounded in realistic user scenarios.

Following the comprehension phase, we evaluate the ability of each LLM to generate end-to-end workflows using three widely adopted prompting strategies (Details are in RQ3): (1) **Instruction-only**, which provides a direct task specification (e.g., Generate a Nextflow pipeline to analyze bisulfite sequencing data.); (2) **Role-based**, which frames the model as a domain expert (e.g., You are a bioinformatics workflow expert. Write a Galaxy workflow to determine gene overlap on opposite strands....); and (3) **Chain-of-thought**, which encourages step-by-step reasoning and decomposition of the task (e.g., Let’s break this task down step-by-step to fetch metadata and raw FastQ files....), which guides the model through step-by-step planning and justification. These strategies emulate realistic human-LLM interactions and provide insight into how prompting affects output quality and completeness. Each prompt is designed to simulate a domain scientist seeking assistance in workflow design. The evaluation process begins with an *instruction-only* prompt. If the resulting workflow is valid, complete, and aligns with community baselines (e.g., GTN or nf-core [35, 36]), we do not proceed to more elaborate prompts. Otherwise, we escalate to role-based and CoT strategies to improve output quality. Across all tasks, we assess LLM-generated responses for *syntactic correctness*, *executability*, *completeness*, and *domain relevance*.

In summary, our approach to RQ1 assesses whether LLMs can accurately understand the scientific context, correctly answer foundational and background questions, and construct usable workflows across a spectrum of bioinformatics tasks. This dual focus on *workflow comprehension and generation* allows us to systematically measure LLM capabilities in supporting domain scientists, particularly those without programming expertise, in developing robust and reproducible workflows. Detailed results are described in Section 4, including prompt-response summaries, workflow validation details, and model-specific performance insights.

To address **RQ2**, we systematically evaluate the workflows generated by each LLM along three core dimensions that define the practical utility of scientific workflows in bioinformatics: *completeness*, *correctness*, and *usability*. Our objective is not merely to determine whether an LLM can produce a workflow, but to assess how well the generated workflows align with real-world scientific and technical standards. For this purpose, we select a diverse set of bioinformatics workflows drawn from two widely used scientific workflow systems: Galaxy (graphical) and Nextflow (script-based). These workflows span a range of analysis tasks, including RNA-seq, bacterial genome annotation, quality control of NGS data, and phylogenetic placement, ensuring broad coverage across typical bioinformatics use cases. Each task is benchmarked against a curated, community-reviewed, and publicly documented implementation from authoritative sources such as the Galaxy Training Network and nf-core [36, 85]. To simulate realistic usage scenarios, workflows are generated through conversational prompting using three strategies: *instruction-only*, *role-based*, and *chain-of-thought*. Prompts are applied incrementally, starting with instruction-only and escalating to more structured prompting only when necessary, reflecting typical user behavior where minimal input is preferred.

The first two authors, with over five and nine years of experience in scientific workflows and workflow systems, respectively, manually evaluate each generated workflow. Assessments are based on public documentation, workflow specifications, and verified execution examples. In addition to technical accuracy, we examine narrative clarity, such as explanatory comments and the rationale behind tool selection, which enhances usability, particularly for novice users.

Overall, this evaluation framework is designed to go beyond superficial correctness by measuring how closely LLM-generated workflows approximate expert-designed standards in terms of reproducibility, executability, and scientific relevance. The results offer critical insights into the practical viability of using LLMs for real-world workflow development across diverse tasks and user expertise levels.

Understanding the prompting strategy for interacting with LLMs is essential for achieving accurate, relevant, and reliable responses. Existing research [28, 28, 33, 86] has shown that different prompting techniques, such as *zero-shot*, *few-shot*, *role-based*, and *chain-of-thought*, lead to varying levels of performance across tasks, particularly in domains requiring structured and sequential reasoning like bioinformatics workflows. Without a deliberate prompting strategy, user queries may yield vague, incomplete, or even incorrect workflows, ultimately limiting the practical utility of LLMs and increasing inefficiency. Therefore, it is essential to evaluate which prompting strategies produce the most accurate, complete, and executable workflows. To address the **(RQ3)**, we conduct an empirical investigation into prompting methods grounded in recent literature on prompt engineering and LLM behavior [e.g., [28–34, 86, 87]]. Our study first surveys the theoretical foundations of these strategies and their prior success in domains involving procedural generation, planning, and code synthesis. We then experimentally evaluate different prompting techniques across diverse workflow scenarios,



refining them through iterative testing. In Section 4, we summarize our findings, providing actionable insights and best practices to help workflow developers optimize their prompting approaches for workflow generation. We provide a practical guideline to support workflow developers who may not possess expertise in prompt engineering or fine-tuning large language models. Recognizing that these users often lack formal training in AI or programming, our guideline is designed to be accessible and actionable, enabling them to effectively adopt our prompting approach for workflow development without requiring deep technical knowledge.

## 4 Results

In this section, we present the results of our evaluation of LLMs across three core research questions. We assess their ability to understand fundamental workflow concepts, generate complete and correct scientific workflows, and respond effectively to various prompting strategies. Our findings are based on qualitative and execution-based analyses using workflows from 1.

### 4.1 RQ1: To what extent can LLMs assist in the development of scientific workflows?

**Motivation:** Workflows in bioinformatics are essential for organizing, automating, and scaling complex data analyses. However, designing such workflows typically requires a combination of domain knowledge, technical proficiency, and iterative refinement skills that are often distributed across interdisciplinary teams. For many domain scientists, particularly those from non-computational backgrounds, understanding and constructing workflows using systems like Galaxy or Nextflow presents a significant challenge. SWSs involve abstract representations of tasks, parameter configurations, and tools/modules interactions that may be difficult to grasp without prior experience. Recent advancements in LLMs present an opportunity to bridge this gap. LLMs have demonstrated strong capabilities in natural language understanding, code generation, and step-by-step reasoning. Workflows in bioinformatics are essential for organizing, automating, and scaling complex data analyses. However, designing such workflows typically requires a combination of domain knowledge, technical proficiency, and iterative refinement skills that are often distributed across interdisciplinary teams. For many domain scientists, particularly those from non-computational backgrounds, understanding and constructing workflows using systems like Galaxy or Nextflow presents a significant challenge. Recent advancements in LLMs present an opportunity to bridge this gap.

**Approach:** To address this research question, we implement a structured three-phase evaluation designed to capture both the conceptual understanding and practical construction capabilities of LLMs in scientific workflow development. The first phase focuses on assessing whether the models can accurately interpret and articulate foundational concepts related to scientific workflows and SWSs, which are essential for meaningful and context-aware workflow generation. We curate a set of nine fundamental prompts (P1–P9), listed in Table 2, reflecting common conceptual inquiries that a domain scientist might pose when engaging with workflows.

**Table 2:** Prompts used to evaluate LLMs’ understanding of core concepts in scientific workflows and workflow systems.

Prompt ID	Question (Prompt)
P1	What is a scientific workflow?
P2	Why are scientific workflows used in research?
P3	What are the main steps in a typical scientific workflow?
P4	What is a Scientific Workflow System (SWS)?
P5	How does an SWS help automate data analysis?
P6	What is a task or process in a workflow?
P7	What is input and output data in the context of a workflow?
P8	What is a workflow tool, and what does it do?
P9	What is meant by data flow or data dependency in workflows?

To further ensure contextual grounding, we supplement this phase with a second set of platform-specific background questions targeting Galaxy and Nextflow (Table 3).

**Table 3:** Prompts used to assess LLMs’ understanding of Galaxy and Nextflow Workflow Systems.

Prompt ID	Galaxy (Prompt)	Nextflow (Prompt)
P1	What is the Galaxy platform, and how is it used in bioinformatics?	What is the Nextflow workflow system, and how is it applied in bioinformatics?
P2	How does Galaxy support reproducibility and transparency in workflows?	How does Nextflow enable reproducibility and portability in bioinformatics pipelines?
P3	What types of analyses can be performed using Galaxy (e.g., RNA-seq, variant calling)?	What types of bioinformatics analyses are commonly implemented using Nextflow (e.g., RNA-seq, metagenomics)?
P4	What are the benefits of Galaxy for bioinformatics workflow?	What are the benefits of using Nextflow for building and scaling bioinformatics workflows?
P5	What are the main components of the Galaxy Scientific Workflow System?	What are Process, Channel, Module, and Operator in Nextflow?
P6	What is Galaxy ToolShed?	What is nf-core?

These are intended to evaluate each LLM’s familiarity with the architecture, features, and design paradigms of the two selected SWSs. Importantly, these questions are submitted to the models prior to issuing any workflow generation prompts, thereby allowing the LLMs to align their reasoning with our experimental intent. All background questions are posed to *GPT-4o*, *Gemini 2.5 Flash*, and *DeepSeek-V3*. Each model is instructed to provide concise yet informative answers. The responses are independently reviewed by two authors with substantial domain expertise (over five and nine years of experience in scientific workflows, respectively) and are cross-referenced against authoritative literature and official documentation to verify their accuracy and relevance. The complete set of responses is made available in the supplementary materials [88].

In the second phase of our evaluation, we focus on assessing LLMs’ contextual understanding of specific bioinformatics workflows before prompting them to generate complete pipeline implementations. To formulate relevant background questions, we first implement each of the selected workflows by meticulously following the step-by-step instructions provided in the Galaxy Training Network (GTN) and nf-core documentation. These tutorials are typically enriched with biological and computational context, which we carefully analyze to extract the foundational concepts required for informed workflow generation. This preliminary implementation process serves two purposes: (1) *it ensures that we, as evaluators, develop a deep understanding of each workflow’s intent and computational logic*; and (2) *it allows us to derive meaningful, targeted questions that reflect the type of background knowledge a domain scientist would need prior to developing a workflow*. The two authors independently review the documentation and execute each workflow, and then collaboratively synthesize key contextual elements into background prompts designed to test the LLMs’ readiness for the task. The complete sets of background questions for workflows based on Galaxy and Nextflow are provided in Table 4 and Table 5, respectively.

In the final phase of our study addressing RQ1, we task the LLMs with generating bioinformatics workflows using Galaxy or Nextflow, depending on the platform associated with each selected use case. This stage is designed to evaluate the extent to which LLMs can synthesize and operationalize workflow knowledge into executable pipelines. Beyond verifying syntactic validity, we focus on whether the generated workflows exhibit appropriate tool selection, logical task sequencing, and adherence to community-established conventions for reproducibility and usability.

To systematically guide this generation process, we adopt a three-tiered prompting framework that varies in the degree of contextual scaffolding: *instruction-only*, *role-based*, and *chain-of-thought (CoT)* prompting. The instruction-only strategy presents a concise directive without any added context, aiming to simulate a general user query. The role-based approach enhances this by assigning the LLM a specific professional identity, typically that of a bioinformatics expert, to encourage more domain-informed reasoning. Finally, the CoT prompting strategy explicitly encourages step-by-step workflow planning, justification of tool/module selection, description of input/output formats, and commentary on task dependencies and execution logic. For instance, in the Galaxy-based workflow *W2: Peaks to Genes*, the CoT prompt elicits structured responses in which the LLM summarizes

**Table 4:** Background questions designed for Galaxy workflows to assess LLMs’ foundational understanding before workflow generation. Each set corresponds to a specific use case.

WID	Prompt ID	Background Question
W1	P1	What are nucleotides, chromosomes, exons, and SNPs?
	P2	What is a sequencing read, and how is its quality measured?
	P3	What is the difference between raw and processed sequencing data?
	P4	What is RNA splicing?
	P5	What is the difference between Exon and Intron?
W2	P1	What is ChIP-seq, and what biological information does it provide?
	P2	What is a peak in genomic data, and how is it represented in a BED file?
	P3	What is gene annotation, and how is it used to associate peaks with genes?
	P4	What is the transcription start site (TSS), and why is it important in gene regulation?
	P5	What are the key file formats in peak annotation, and how do they differ?
	P6	What is the promoter region of genes?
	P7	Why is it important to identify target genes from genomic regions in research?
W3	P1	What is the structure of double-stranded DNA, and what are the forward and reverse strands?
	P2	What are genome and chromosome?
	P3	What are reference genome and GENCODE?
	P4	How is strand orientation recorded in genome annotation files?
	P5	How do forward/reverse strands affect read mapping and quantification?
	P6	What does it mean for genes on opposite strands to overlap, and why is this an interesting biological question?
	P7	What are genomic intervals, and how is the overlap between features defined?
	P8	How can understanding gene strand orientation impact the interpretation of genetic information?
	P9	What’s the difference between sequence and annotation?
	P10	What are FASTA, BED, GTF, GFF3, and VCF?
	P11	What are GRCh37, GRCh38, hg19, and hg38?
W4	P1	What is Sanger sequencing? Provide some examples.
	P2	What are primers in sequencing?
	P3	What is the AB1 sequence file?
	P4	What are CHD8 and AOPEP?
	P5	What are sense, antisense, and consensus sequences?
	P6	What is the reverse-complement sequence, and why compute it for the antisense strand?
	P7	Why is converting from AB1 sequence file to FASTQ essential?
W5	P1	What is genome annotation, and why is it critical for bacterial genomics?
	P2	What is Bakta for bacterial genome annotation?
	P3	What are coding sequences (CDS), and what roles do rRNAs and tRNAs play?
	P4	What are plasmids?
	P5	What are integrons?
	P6	What are GFF3 and GenBank formats, and how are they used in annotation pipelines?
	P7	How do you evaluate the quality of an annotated genome?

the biological objective, selects appropriate annotation tools (e.g., ChIPseeker), and provides justifications for each component. In the case of Nextflow, the *W4: nf-core/rnaseq* workflow, the CoT strategy effectively prompts the LLM to outline the modular design, describe input channels, and specify configuration parameters aligned with Nextflow syntax.

We apply this prompting strategy uniformly across all considered workflows and LLMs, ensuring methodological consistency for fair comparison. This design allows us to assess how varying degrees of prompt specificity influence the LLMs’ ability to generate accurate, complete, and usable workflows in both graphical (Galaxy) and script-based (Nextflow) Scientific Workflow Systems. Ultimately,

**Table 5:** Background questions designed for nf-core workflows to assess LLMs’ foundational understanding before workflow generation. Each set corresponds to a specific pipeline.

WID	Prompt ID	Background Question
W1	P1	What is the purpose of performing quality control on sequencing data?
	P2	What are sequencing adapters, and why is it necessary to remove them from raw reads?
	P3	How does read quality affect downstream bioinformatics analyses such as alignment or quantification?
	P4	What is the difference between single-end and paired-end reads, and how does it influence preprocessing?
	P5	Why is it important to visualize QC metrics before proceeding with downstream analysis?
	P6	What is FastQC, and what metrics does it generate for assessing sequencing read quality?
W2	P1	What is NGS?
	P2	What are SRA, ENA, DDBJ, and GEO and why are they important for NGS data retrieval?
	P3	What is a FASTQ file, and what type of biological data does it contain?
W3	P1	What is DNA methylation and why is bisulfite sequencing used to measure it?
	P2	What is genome alignment?
	P3	What are the challenges in aligning bisulfite-converted reads to a reference genome?
	P4	What types of input data (e.g., single-end or paired-end FASTQ) are required for methylation analysis?
	P5	How are methylation levels quantified and visualized across genomic regions?
	P6	Briefly describe Bismark and BWA-meth?
W4	P1	What is RNA-seq and how does it help quantify gene expression?
	P2	What is a reference genome and annotation?
	P3	What is a gene expression matrix?
	P4	What preprocessing steps (e.g., adapter trimming, QC) are necessary before alignment?
	P5	What are transcript-level vs. gene-level quantifications and when are they used?
W5	P1	What is phylogenetic placement, and how does it differ from reconstructing a full phylogenetic tree?
	P2	Why is phylogenetic placement useful for short or fragmentary sequences?
	P3	What is EPA-NG?
	P4	What are query and reference sequences?

this phase serves as a bridge between conceptual understanding and practical construction, offering insight into the generative capabilities of LLMs within real-world bioinformatics contexts.

Drawing from existing literature on effective prompting strategies (discussed in RQ3), we design and apply structured prompts to guide LLMs in workflow generation. Tables 6 and 7 present representative prompts used to guide the LLMs in generating our workflows. These examples illustrate how varying levels of contextual scaffolding, ranging from instruction-only to role-based and chain-of-thought prompts, are employed to elicit structured and contextually appropriate responses. While the prompts shown are specific to selected case studies, a consistent prompting strategy is applied across all workflows evaluated in this study, as detailed in the results section. The complete set of prompts and corresponding LLM outputs are included in the supplementary material to promote transparency and support reproducibility. All LLMs are queried using default model settings to simulate a realistic usage scenario, where domain experts in bioinformatics may rely on general prompting capabilities without advanced expertise in prompt engineering.

To validate the generated workflows, we utilize the Galaxy web interface available at <https://usegalaxy.org/>, which offers a fully hosted environment for executing and testing Galaxy-based workflows. For Nextflow workflows, we follow the official installation procedures provided in the documentation at [nextflow.io/docs/latest/install.html](https://nextflow.io/docs/latest/install.html). Additionally, for users seeking a lightweight or cloud-based setup, workflows can be executed through GitHub Codespaces using the pre-configured Nextflow training environment available at [github.com/codespaces/nextflow-io/training](https://github.com/codespaces/nextflow-io/training). This setup

**Table 6:** Prompts used for Galaxy workflow generation across different use cases

Workflow ID	Prompt Type and Content
W1	<b>Instruction-Only:</b> Create a Galaxy workflow that identifies the coding exon with the highest number of single-nucleotide polymorphisms (SNPs) on human chromosome 22. Include all necessary tools, input data types, and expected outputs.
W2	<b>Instruction-Only:</b> Create a Galaxy workflow to identify target genes from ChIP-seq peak regions. Use the dataset <code>GSE37268.mof3.out.hpeak.txt.gz</code> and a mouse gene annotation file from UCSC. <b>Role-Based:</b> You are a bioinformatics workflow developer. Create a Galaxy workflow to identify target genes from ChIP-seq peak regions. The peak file contains numeric chromosome identifiers (e.g., 1–20 for X, 21 for Y) instead of UCSC-style labels. Preprocess the chromosome field accordingly to ensure compatibility with standard annotations.
W3	<b>Instruction-Only:</b> Create a Galaxy workflow to determine whether genes on opposite DNA strands overlap, using a gene annotation dataset in BED format.
W4	<b>Instruction-Only:</b> Create a Galaxy workflow to perform quality control on raw next-generation sequencing (NGS) data, ensuring the dataset is assessed and improved for high-quality downstream analysis.
W5	<b>Instruction-Only:</b> Create a Galaxy workflow to annotate an assembled bacterial genome provided in FASTA format. Use tools to predict genes and functional elements, identify plasmid replicons, detect integron structures, and find insertion sequences. Convert tool-specific outputs into GFF3 format and load all annotations into a genome browser for interactive visualization. <b>Role-Based:</b> You are a bioinformatics workflow developer working within the Galaxy platform. Create a workflow to annotate a draft bacterial genome, evaluate the annotation, format the outputs for visualization, and visualize the results. <b>Chain-of-Thought:</b> We have a draft bacterial genome in FASTA format. First, identify genomic components such as CDS, rRNAs, and tRNAs with structural and functional annotations. Save outputs in standard formats for downstream use (e.g., feature tables, protein/nucleotide sequences). Check annotation completeness and summarize statistics. Finally, organize outputs for compatibility with genome browsers. Based on this reasoning, generate a Galaxy workflow that completes these tasks.

offers a streamlined alternative to local installation and ensures reproducibility in a controlled execution environment.

We adopt a stepwise prompting strategy designed to minimize unnecessary complexity while ensuring the generation of accurate and complete workflows. Each workflow generation task begins with an instruction-only prompt. If the output from this initial prompt aligns with the expected workflow structure, based on baselines such as GTN or nf-core documentation, we do not proceed with further prompting. However, if the workflow is incomplete, incorrect, or misaligned with the baseline, we escalate to a role-based prompt to encourage more domain-specific reasoning. If this also fails to yield a satisfactory result, we employ a chain-of-thought (CoT) prompt, which guides the model through a step-by-step reasoning process. This tiered prompting strategy enables a principled trade-off between prompting complexity and output quality, while also allowing us to assess the minimal level of guidance required for successful workflow generation.

**Results of RQ1:** To evaluate the LLMs’ foundational understanding of scientific workflows and workflow systems, we present a set of core conceptual prompts (Table 2). All three models, GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3, are able to generate responses that are technically accurate and aligned with the intended meaning of each question. To further assess their effectiveness, we evaluate the responses using six qualitative criteria: *correctness*, *clarity*, *completeness*, *terminology*, *structure*, and *usability*. A comparison of their performance is summarized in Table 8. While all models satisfied the correctness and terminology criteria, notable differences emerged across the remaining dimensions. GPT-4o demonstrated strong clarity and accessibility, offering responses that are well-suited for foundational learning and novice users. Gemini 2.5 Flash provided the most comprehensive and academically precise answers, making it ideal for users with intermediate experience or those seeking formal educational material. DeepSeek-V3 delivered contextually accurate responses enriched with practical examples; however, its outputs are occasionally verbose and informally structured, limiting their immediate instructional usability without refinement.

We next evaluate the responses to the platform-specific background prompts designed to assess LLMs’ contextual understanding of Galaxy and Nextflow (Table 3). The evaluation results, summarized in Table 9, reveal distinct strengths and limitations across the three models. GPT-4o produces accurate, well-structured, and beginner-friendly responses characterized by clear language and logical organization. However, it occasionally omits advanced implementation details such as containerization strategies and execution environments, which are important for reproducibility and deployment.



Table 7: Prompts used for Nextflow workflow generation across different use cases

Workflow ID	Prompt Type and Content
W1	<b>Instruction-Only:</b> We want to create a workflow using Nextflow that performs quality assessment of FASTQ files, preprocesses sequencing data, and generates an aggregated summary of the results. The input data should be provided via the sample sheet available at <a href="#">samplesheet_test_illumina_amplicon.csv</a> . Please outline the step-by-step modules required to build this workflow.
W2	<b>Instruction-Only:</b> We want to create a Nextflow workflow that takes a CSV file containing sample database IDs (e.g., SRA, ENA, GEO, or DDBJ) as input, retrieves associated metadata, and downloads raw FASTQ files. Please outline the step-by-step modules to create the workflow.
W3	<p><b>Instruction-Only:</b> We want to create a Nextflow workflow for DNA methylation analysis using bisulfite-converted sequencing (BS-seq) data. The workflow should begin by pre-processing raw reads provided in a CSV file containing FASTQ input information, followed by alignment to a reference genome and comprehensive quality control assessments. Could you please provide a step-by-step outline of the modules required to construct this workflow?</p> <p><b>Role-Based:</b> You are a bioinformatics workflow developer with expertise in building reproducible pipelines using Nextflow and nf-core modules. Your task is to design a Nextflow workflow for DNA methylation analysis using bisulfite-converted sequencing (BS-seq) data. The workflow should begin by pre-processing raw reads specified in a CSV file (containing sample IDs and FASTQ paths), then proceed to align reads to a reference genome, extract methylation information, and generate comprehensive quality control reports. Please provide a step-by-step outline of the modules required to build this workflow, using only officially available nf-core or Nextflow modules, and include their exact module names where appropriate.</p>
W4	<b>Instruction-Only:</b> We want to create a Nextflow workflow for RNA-seq analysis that performs quality control, adapter trimming, alignment, quantification, and differential expression analysis. The workflow should accept a CSV sample sheet containing FASTQ file links and genome references as input. It must include quality checks, reference genome indexing if needed, and gene/transcript quantification, and produce a gene expression matrix and extensive QC report using appropriate tools. Please provide a step-by-step outline of the modules to construct this workflow.
W5	<b>Instruction-Only:</b> We want to create a Nextflow workflow for the phylogenetic placement of query sequences onto a reference tree, supporting two operational modes. In placement-only mode, the workflow uses a query FASTA file, a reference alignment, a rooted phylogenetic tree in Newick format, an evolutionary model, and optionally a taxonomy file for classification. In search-plus-placement mode, it begins with a large unfiltered FASTA file and HMM profiles to identify candidate sequences, followed by alignment, placement, classification, and reporting using the same reference inputs. The workflow aligns selected sequences, places them onto the tree, classifies their positions, and generates summary reports and visualizations. Could you please provide a step-by-step outline of the modules required to construct this workflow?

Table 8: Comparison of LLM Responses across Evaluation Criteria for Fundamental Concepts (Table 2 )

Criterion	GPT-4o	Gemini 2.5 Flash	DeepSeek-V3
Correctness	✓ Accurate but basic definitions	✓ Most precise with academic rigor	✓ Technically sound with practical examples
Clarity	✓ Very clear and straightforward	✓ Clear but occasionally formal	⚠ Clear but Verbose
Completeness	⚠ Covers basics but lacks depth i.e., omits advanced concepts like validation	✓ Comprehensive answers with good coverage of key concepts	✓ Covers most concepts explicitly, often using enumeration
Terminology	✓ Uses familiar and readable technical terms.	✓ Uses formal and precise academic terminology	✓ Appropriate with real-world SWSs (e.g., Galaxy, KNIME)
Structure	✓ Consistent and clean	✓ Well-organized	⚠ Often uses lists and informal bullets
Usability	✓ Excellent for beginner tutorials	✓ Good for intermediate audiences	⚠ Needs cleanup for instructional material

Gemini 2.5 Flash delivers the most academically rigorous and comprehensive responses, offering precise explanations that cover tool integration, ecosystem capabilities, and platform-specific components. Its slightly formal tone and structured presentation make it particularly well-suited for intermediate users or instructional contexts. DeepSeek-V3 performs well in technical accuracy and real-world applicability, often referencing specific tools and platforms such as ToolShed and WorkflowHub. However, its responses are frequently verbose, rely heavily on bullet-point formatting, and lack the narrative flow necessary for effective instructional use without additional editing. Overall, Gemini provides the most complete and balanced output, GPT-4o excels in clarity and accessibility for novice users, and DeepSeek-V3 offers valuable practical insight but would benefit from refinement to improve instructional usability.

**Table 9:** Comparison of LLM Responses across Evaluation Criteria for Background Questions about Galaxy and Nextflow (Table 3)

Criterion	GPT-4o	Gemini 2.5 Flash	DeepSeek-V3
<b>Correctness</b>	✓ Accurate but concise; omits container/cloud details	✓ Most precise with deep tool/system understanding	✓ Correct with real examples from SWSs
<b>Clarity</b>	✓ Clear and beginner-friendly	✓ Clear with slight formal tone	⚠ Verbose and uses list-heavy structure
<b>Completeness</b>	⚠ Covers basics; lacks advanced implementation details	✓ Thorough with ecosystem, portability, and usage depth	✓ Covers broad scope and tool support explicitly
<b>Terminology</b>	✓ Accessible language for broad users	✓ Formal academic terminology well-used	✓ Includes domain-specific terms and tools (e.g., ToolShed, WorkflowHub)
<b>Structure</b>	✓ Narrative, well-flowing format	✓ Well-structured with natural coherence	⚠ Checklists over narrative, uneven depth
<b>Usability</b>	✓ Suitable for introductory tutorials	✓ Great for intermediate bioinformatics learners	⚠ Needs formatting edits for teaching or documentation

To evaluate the ability of LLMs to address workflow-specific background questions, we provide a set of foundational biological and technical prompts (Table 4) related to representative Galaxy workflows. The results of this evaluation (Table 10) indicate that Gemini 2.5 Flash consistently provides the most accurate and comprehensive responses. For example, in response to *W2-P3 (What is gene annotation, and how is it used to associate peaks with genes?)*, Gemini offers a detailed explanation encompassing both functional annotation and genomic coordinate mapping, whereas GPT-4o and DeepSeek-V3 provide less insight into the mechanisms linking peaks to nearby genes. Similarly, for *W5-P2 (What is Bakta for bacterial genome annotation?)*, Gemini elaborates on database usage and output formats, while other models only describe its general purpose. GPT-4o, although less comprehensive, stands out for its clarity and ease of understanding, making it particularly suitable for beginners. For example, its explanation of *W1-P1 (What are nucleotides, chromosomes, exons, and SNPs?)* was concise, well-structured, and accessible to users with minimal background in genomics. On the other hand, DeepSeek-V3, while generally accurate, shows inconsistencies in depth and completeness. In *W3-P6 (What does it mean for genes on opposite strands to overlap?)*, DeepSeek fails to address key implications of antisense transcription and regulatory complexity that are captured by Gemini. These findings indicate that Gemini 2.5 Flash is best suited for users seeking technical rigor, GPT-4o excels in communication for broader audiences, and DeepSeek-V3 may require improvement in consistency and contextual depth for bioinformatics education and support tasks.

To assess how well LLMs explain workflow-specific background concepts in Nextflow pipelines, we analyze their responses across a curated set of prompts spanning five real-world Nextflow workflows. These prompts covered foundational biological and computational topics relevant to workflows such as *nf-core/rnaseq*, *fetchngs*, and *phyloplace*. Table 11 summarizes the comparative performance across.

**Table 10:** Comparison of LLM Responses for Background Prompts in Galaxy Workflows

Criterion	GPT-4o	Gemini 2.5 Flash	DeepSeek-V3
<b>Correctness</b>	✓ Generally accurate; some minor conflation in technical distinctions	✓ Highly accurate and aligned with academic definitions	✓ Mostly correct but occasionally vague or incomplete
<b>Clarity</b>	✓ Very clear and accessible to beginners	⚠ Clear but often formal or verbose	⚠ Mixed clarity; some explanations too shallow or ambiguous
<b>Completeness</b>	⚠ Basic coverage; misses deeper insights	✓ Comprehensive, includes biological relevance and technical context	⚠ Omits some intermediate concepts or format distinctions
<b>Conciseness</b>	✓ Concise and to the point	⚠ Informative but often lengthy	✓ Generally concise but may lack depth
<b>Terminology</b>	✓ Uses common and understandable terms	✓ Uses precise and domain-specific vocabulary	⚠ Uses general terms; sometimes lacks specificity
<b>Mistakes/Errors</b>	Minor blending of concepts in definitions	Occasionally overuses jargon without explanation	Misses important context (e.g., gene-peak associations)
<b>Overall Verdict</b>	<b>Most user-friendly</b>	<b>Best academic rigor and depth</b>	<b>Needs improvement in consistency and detail</b>

**Table 11:** Comparison of LLM Responses for Background Prompts in Nextflow Workflows

Criterion	GPT-4o	Gemini 2.5 Flash	DeepSeek-V3
<b>Correctness</b>	✓ Generally correct with some minor oversights in tool-specific context	✓ Highly accurate and comprehensive across workflows	✓ Mostly accurate with detailed examples
<b>Clarity</b>	✓ Clear and easy to follow for most users	⚠ Slightly formal tone but readable	⚠ Occasionally dense and list-heavy
<b>Completeness</b>	⚠ Omits details like parameter flags and metadata resolution	✓ Covers technical tools, configurations, and context thoroughly	⚠ Strong examples, but some answers are tool-biased or partial
<b>Terminology</b>	✓ Balanced use of domain-relevant terms	✓ Advanced bioinformatics vocabulary used appropriately	⚠ Inconsistent: sometimes general, sometimes overly specific
<b>Conciseness</b>	✓ Well-balanced for instructional use	⚠ Verbose in longer answers	✓ Compact but requires interpretation
<b>Mistakes/Errors</b>	Minor: occasionally omits integration aspects or setup steps	Few: may overwhelm novices with detail	Minor: tool links are correct but lack rationale
<b>Overall Verdict</b>	<b>Best for clarity and beginners</b>	<b>Best for completeness and experts</b>	<b>Good for technical detail, less suited for guided explanation</b>

Gemini 2.5 Flash addresses complex topics such as metadata retrieval in *fetchngs*, reference genome handling in *phyloplace*, and bisulfite alignment nuances in *methyseq*, often incorporating tool-specific details, flags, and file structure references. For example, Gemini correctly explained the importance of ENA API calls for resolving sample metadata and provided detailed descriptions of output artifacts. Its responses are comprehensive, though sometimes overly verbose or formal, which could be challenging for novice users. GPT-4o, by contrast, prioritizes clarity and accessibility. It provides succinct, beginner-friendly explanations that are logically structured and easy to follow. For example, its explanation of the FASTQ format and the rationale for adapter trimming in RNA-seq workflows is precise and instructional. However, GPT-4o occasionally lacks depth in areas such as the resolution of experiment-level identifiers or the rationale behind multiple HMM profiles in *phyloplace*. DeepSeek-V3 delivered technically correct responses with a strong practical orientation. It excels in referencing command-line tools, directory structures, and usage examples. In the context of *fetchngs*, for instance, it correctly described MD5 checksum validation and discussed the role of sample mapping fields. However, its answers often leaned toward being list-heavy and missed some of the conceptual explanations, such as the biological rationale behind bin refinement in MAG analysis or how GAPPA supports phylogenetic visualization. In summary, while all three LLMs perform reasonably well in covering core concepts, Gemini 2.5 Flash is best suited for users seeking comprehensive, GPT-4o excels in readability and instructional clarity, and DeepSeek-V3 offers detailed practical insights but requires refinement for conceptual completeness and consistency.

#### 4.1.1 Workflow Development Using LLMs for Galaxy:

For Workflow W1: *Identifying the exon with the highest number of SNPs on human chromosome 22* within the Galaxy SWS, we first manually construct and execute the workflow using the Galaxy platform. This was done by following the step-by-step guidance provided in the Galaxy Training Network (GTN) tutorial [75]. Utilizing the sample dataset provided by GTN, we determine that the exon *ENST00000253255.7\_cds\_0\_0\_chr22\_46256561\_r* contains the highest number of SNPs (27 in total) on human chromosome 22. Subsequently, we apply the instruction-only prompt (W1, Table 6) to GPT-4o, asking it to generate the corresponding workflow steps. The model successfully produces a logically ordered Galaxy workflow that correctly identifies the same exon. Notably, while the GTN tutorial assumes that the input files are already restricted to chromosome 22, GPT-4o explicitly includes filtering steps to isolate chromosome 22 from genome-wide datasets. This reflects a more realistic real-world scenario and demonstrates the model’s capacity to generalize beyond the assumptions of curated tutorials. In terms of tool selection, we observe minor differences in the implementation details. The GTN workflow uses the *bedtools intersect intervals* tool to compute overlaps between SNPs and exons, whereas GPT-4o suggests the use of the *Join* tool to perform the same operation. Additionally, while the GTN workflow employs the *Datamash* tool to count SNPs per exon, GPT-4o proposes using the *Group* tool. Despite these differences in tool choices, the overall logic and structure of the workflow remain correct and effective.

To evaluate the performance of Gemini 2.5 Flash, we use the same instruction-only prompt used previously. The model generates a clear and well-structured workflow, offering step-by-step guidance that includes both data acquisition and filtering procedures. Similar to the GTN tutorial, Gemini 2.5 Flash recommends using the *bedtools intersect intervals* tool to compute overlaps between SNPs and exons. However, for aggregating SNP counts per exon, it opted for the *Group* tool instead of the *Datamash* tool used in GTN. The workflow also includes the use of the *Sort* tool, a common step across all three implementations: GTN, GPT-4o, and Gemini 2.5 Flash. Upon executing the workflow generated by Gemini 2.5 Flash using the same GTN dataset, we obtain the same correct result, identifying the exon *ENST00000253255.7\_cds\_0\_0\_chr22\_46256561\_r* with 27 SNPs, as observed in the baseline and GPT-4o workflows.

For DeepSeek-V3, we again apply the same instruction-only prompt and are successfully able to generate the steps to develop correct and executable Galaxy workflow for identifying the coding exon with the highest number of SNPs on chromosome 22. Similar to GPT-4o, DeepSeek-V3 recommends the use of the *Intersect* tool to compute overlaps between SNPs and exon regions. Additionally, the model offers flexibility by suggesting alternative tools for certain tasks. For example, to count the number of SNPs per exon, it proposes using either the *Group* or *Count* tool, both valid and appropriate choices. For identifying the exon with the highest SNP count, it recommends the combination of *Sort* followed by *Select*, which aligns with standard Galaxy practices. Notably, DeepSeek-V3 also provides supplementary information beyond the core task. It includes guidance on how to export,

import, and execute the workflow within the Galaxy platform, along with suggestions for potential public data sources for exons and SNPs. While such details go beyond the minimal workflow requirements, they enhance usability for novice users. Overall, DeepSeek-V3 successfully produces a functionally correct workflow from the instruction-only prompt, demonstrating its capability to support bioinformatics workflow design.

Given that the instruction-only prompt produce valid and complete workflows across all three LLMs, we do not proceed with the more elaborate Role-based or CoT prompting strategies for this task. This outcome indicates that for well-defined and well-scoped problems like Workflow W1 (Galaxy), minimal prompting is sufficient to guide advanced LLMs in generating accurate and executable workflows.

*While all three LLMs, successfully generate correct and executable workflows from the instruction-only prompt, Gemini 2.5 Flash stands out for its combination of technical accuracy and instructional clarity. Its response closely aligns with the GTN tutorial, uses appropriate tools, and provides a structured, step-by-step guide that balances detail with usability. Although GPT-4o produces a concise and accurate workflow with thoughtful generalizations, such as including chromosome filtering, it lacks some of the explanatory depth seen in Gemini’s response. DeepSeek-V3 offers additional contextual guidance and tool flexibility, which may benefit novice users, but it leans toward verbosity and lacks the structured clarity of Gemini. Overall, Gemini 2.5 Flash offers the most intuitive and comprehensive support for this task, making it the preferred choice for users seeking both correctness and clear instructional guidance. The LLMs generated workflows are shared in the supplementary material.*

We next evaluate Workflow W2: *How to get from peak regions to a list of gene names?*. The primary objective is to identify genes that overlap with experimentally derived peak regions and to summarize their distribution across the genome. Following the GTN tutorial, we first upload the dataset *GSE37268\_mof3.out.hpeak.txt.gz*, obtained from the *Gene Expression Omnibus (GEO)*. To annotate the peaks, we retrieve a gene annotation file for the mouse genome from the UCSC Genome Browser and compute overlaps between the peak regions and annotated genes. Finally, we quantify the number of overlapping genes per chromosome to identify regions with enriched regulatory activity. Our analysis shows that chromosome 11 (chr11) contains the highest number of overlapping genes (2164), suggesting notable regulatory significance in this region.

Then, we begin the workflow generation by applying an instruction-only prompt (W2, Table 6) to GPT-4o. The model returns a logically ordered and mostly accurate workflow consistent with common analysis practices. However, it omits a critical preprocessing step: the chromosome identifiers in the input file are represented numerically (e.g., 1, 20, 21), which are incompatible with standard UCSC annotations that use labels such as chr1, chrX, and chrY. Without normalization of these identifiers, such as prefixing with chr and mapping 20 and 21 to X and Y, the *BEDTools Intersect* operation fails to correctly associate peaks with genes, compromising result accuracy. To address this limitation, we escalate to a role-based prompt (W2, Table 6), instructing the model to act as a bioinformatics workflow developer. This approach provides the necessary context to guide GPT-4o in including the missing preprocessing step, thereby producing a more accurate and executable workflow. Following the steps, we are able to obtain the correct result from the workflow.

To evaluate Gemini 2.5 Flash’s ability to generate Workflow W2, we apply the same instruction-only prompt. Gemini responds with a comprehensive and intuitively organized step-by-step guide that demonstrates a strong understanding of both the biological context and the technical requirements. Notably, the model includes critical preprocessing stages, such as acquiring ChIP-seq peak data from public repositories (e.g., *GSE37268\_mof3.out.hpeak.txt.gz*) and retrieving gene annotation files from the UCSC Main Table Browser with appropriate parameterization. It accurately outlines core analytical steps using standard Galaxy tools, including *BEDTools Intersect* for overlapping genomic intervals, followed by *Cut columns* and *Select unique lines* to generate a non-redundant list of gene names. The clarity, depth, and correctness of the response make it particularly suitable for novice users and confirm the model’s capability to produce a valid workflow using minimal prompting.

Similarly, DeepSeek-V3 performs well under the instruction-only prompt, generating a clear and well-structured workflow for identifying target genes from ChIP-seq peak regions. Although its tool choices differ slightly, the methodology remains sound. For instance, DeepSeek-V3 recommends using *ChIPseeker* for annotating peaks with the nearest genes, while also suggesting *BEDTools closest* as a viable option. Importantly, unlike GPT-4o, DeepSeek-V3 explicitly acknowledges the need to normalize chromosome identifiers (e.g., converting numeric values to UCSC-style labels), ensuring compatibility with downstream gene annotation. Furthermore, the model enhances usability by including a visual workflow diagram, which helps users conceptualize the overall analysis pipeline.



*In comparing the three LLMs on this task, clear differences emerge in completeness, clarity, and sensitivity to data formatting nuances. GPT-4o generates a mostly accurate and executable workflow but omits the crucial preprocessing step for chromosome label conversion. This omission necessitated an additional role-based prompt to achieve a fully functional pipeline. In contrast, Gemini 2.5 Flash offers the most technically complete and contextually aware solution, implicitly aligning data sources and addressing format compatibility through precise instructions. DeepSeek-V3 also performs commendably, offering methodologically valid alternatives, recognizing critical preprocessing needs, and improving accessibility through visual aids. Overall, Gemini delivers the most robust and reliable output, DeepSeek enhances interpretability and ease of use, while GPT-4o benefits from additional contextual scaffolding to meet expert-level expectations.*

Our third workflow, *W3: Do genes on opposite strands ever overlap? If so, how often?*, explores the prevalence of gene overlaps between forward (+) and reverse (−) DNA strands, offering insights into the complexity of genomic organization. Following the GTN tutorial, we construct and execute a workflow that utilizes gene annotation data from the UCSC Table Browser. The analysis reveals that approximately 38.46% of the genes exhibit strand overlap, indicating a substantial level of strand-specific gene co-localization.

To generate this workflow using GPT-4o, we apply an instruction-only prompt (W3, Table 6). The model produces a concise and logically ordered set of instructions outlining the core steps required for workflow construction in Galaxy. Based on this guidance, we are able to successfully recreate the workflow and replicate the results, confirming the correctness and usability of the generated steps. However, the output exhibits limitations in procedural specificity. For instance, while the model accurately notes the need to separate genes by strand into two datasets (forward and reverse), it fails to detail how to perform this operation. It does not indicate that strand information is typically located in the sixth column of BED-formatted files, nor does it reference where to obtain the gene annotation dataset, an omission that could hinder users unfamiliar with genomic data sources or file conventions. These limitations underscore a broader observation: although GPT-4o’s high-level instructions are sufficient for experienced users who understand data structures and workflows, novice users may struggle without additional domain-specific guidance. Despite this, the generated workflow was complete and executable, so we did not proceed with alternative prompting strategies such as role-based or chain-of-thought prompting for this task.

To evaluate Gemini 2.5 Flash’s performance on Workflow W3, we issue the same instruction-only prompt used previously. The model delivers a comprehensive, step-by-step guide that not only identifies the appropriate analytical tools, such as conditional filtering for strand separation and *BEDTools Intersect* for overlap analysis, but also includes explicit parameter configurations (e.g., using `c6 == '+'` to isolate forward-strand genes). In addition to the analytical workflow, Gemini provides detailed instructions on how to construct the workflow within the Galaxy interface, including drag-and-drop operations in the Galaxy workflow editor and guidance on saving and executing the pipeline. This level of procedural depth and platform-specific instruction makes the output highly accessible to users with limited experience in Galaxy. DeepSeek-V3 also performs well under the same instruction-only prompt, producing a correct and well-structured workflow for identifying overlapping genes on opposite DNA strands. The model suggests a methodologically sound approach that includes optional sorting steps, strand separation via filtered expressions, and intersection analysis. It specifies key technical details, such as enforcing a minimum one bp overlap, and outlines optional post-processing steps, including counting overlaps or isolating unique gene pairs for further analysis.

*Across all three models, we observe the successful generation of valid Galaxy workflows from the instruction-only prompt. GPT-4o provides a concise and logically structured outline but lacks critical implementation details such as the specific column (typically column 6 in BED format) used to represent strand information and the strandedness options required during intersection. These omissions could hinder less experienced users. In contrast, Gemini 2.5 Flash delivers the most comprehensive and accessible solution, with clear explanations of strand-specific filtering, intersection logic, and platform-level guidance for workflow execution. DeepSeek-V3 offers a robust and technically accurate workflow, enriched with optional enhancements that improve usability and extend analytical depth. Overall, Gemini 2.5 Flash emerges as the most effective model for supporting both novice and intermediate users, balancing accuracy, clarity, and practical usability in workflow construction.*

We next examine Workflow *W4: How can the quality of NGS raw data be assessed, what parameters should be checked, and how can the quality of the dataset be improved?*. This workflow addresses a fundamental preprocessing step in next-generation sequencing (NGS) analysis: assessing and improving the quality of raw *FASTQ* files before downstream applications. Following the GTN tutorial, we

construct a workflow that begins by inspecting raw sequencing reads and evaluating quality metrics using tools such as *FastQC* and *FASTQE* (for short Illumina reads), or *NanoPlot* and *PycoQC* (for long Nanopore reads). These tools report per-base quality scores, GC content, adapter contamination, sequence duplication levels, read length distributions, and run-level metrics. The workflow then incorporates trimming and filtering steps using tools like *Cutadapt* and *Fastp* to remove low-quality bases, adapter sequences, ambiguous reads, and short fragments. Paired-end data handling is explicitly supported, and the final results are summarized using *MultiQC*. The tutorial emphasizes the iterative nature of quality control: *assess*  $\rightarrow$  *clean*  $\rightarrow$  *re-assess*.

To generate the workflow using LLMs, we apply the instruction-only prompt (W4, Table 6) to GPT-4o. The resulting workflow closely mirrors the core structure of the GTN pipeline, adhering to the *assess*  $\rightarrow$  *clean*  $\rightarrow$  *re-assess*. GPT-4o outlines the primary steps effectively, recommending *FastQC* for initial assessment, *Cutadapt* for trimming, and *MultiQC* for summarizing results. However, its workflow is streamlined and focuses primarily on short-read Illumina data, lacking the platform-specific enhancements present in the GTN tutorial. However, GPT-4o produces a minimal yet functional workflow suitable for general Illumina data, prioritizing simplicity and automation.

The workflow generated by Gemini 2.5 Flash closely mirrors the core structure of the GTN Quality Control tutorial, adhering to the standard *assess*  $\rightarrow$  *clean*  $\rightarrow$  *re-assess* paradigm. It utilizes essential tools such as *FastQC* for both pre- and post-trimming quality assessment, *MultiQC* for summarizing quality reports, and trimming tools like *Trimmomatic* or *fastp* to remove low-quality bases and adapter sequences. While well-suited for short-read (Illumina) data and effective for general QC workflows, Gemini’s output remains focused on commonly used tools and does not incorporate platform-specific enhancements. In contrast, the GTN tutorial extends its coverage to include *NanoPlot* and *PycoQC* for long-read Oxford Nanopore data, and offers additional options like *Cutadapt*, making it more inclusive for a broader range of sequencing technologies.

DeepSeek-V3 also aligns closely with the conceptual structure of the GTN Quality Control tutorial, particularly in its adoption of the *assess*  $\rightarrow$  *clean*  $\rightarrow$  *re-assess* paradigm and its use of core tools such as *FastQC*, *MultiQC*, and trimming utilities like *Cutadapt*, *Trim Galore*, *Trimmomatic*, and *fastp*. The model introduces optional filtering tools such as *PRINSEQ* and *BBDuk*, and allows users to choose between *Cutadapt* and *Trim Galore!* for adapter removal. Additionally, it provides detailed implementation guidance, including parameter recommendations, memory usage optimization, paired-end read handling instructions, and version tracking to promote reproducibility. The inclusion of a workflow diagram further enhances usability and instructional value.

*When comparing all three LLMs on W4, each successfully generates a valid Galaxy-based QC workflow, but they vary in depth and adaptability. GPT-4o provides a simplified yet accurate representation of the QC process using a minimal toolset, making it suitable for rapid assessment of Illumina datasets but lacking in tool diversity and parameter specificity. Gemini 2.5 Flash improves on this by offering alternative tools, implementation context, and better instructional clarity, which benefits novice users. DeepSeek-V3, however, produces the most comprehensive and customizable workflow, capturing both the conceptual rigor and practical nuances of multi-platform NGS data quality control. Its extended toolset and configuration guidance make it particularly well-suited for diverse sequencing protocols and robust real-world usage.*

Our final workflow, W5: *Which genes are on a draft bacterial genome, and which other genomic components can be found on a draft bacterial genome?*, focuses on systematically annotating an assembled bacterial genome to identify and characterize its genetic content. The workflow begins by accepting a draft genome assembly in FASTA (contigs) format and uses *Bakta* to predict and annotate core genomic features, including coding sequences, tRNAs, rRNAs, and other functional elements. It proceeds to detect plasmid replicons using *PlasmidFinder*, insertion sequences with *ISEScan*, and integron structures with *IntegronFinder*. The resulting annotations are reformatted into *GFF3* format, and the final step involves visualizing the annotated genome interactively using *JBrowse*. Following the GTN tutorial, we develop the workflow and transform raw genome assemblies into richly annotated and visually interpretable bacterial genomes following the GTN instructions.

To evaluate the ability of LLMs to generate this workflow, we provide an instruction-only prompt (W5, Table 6). When applied to GPT-4o, the model generates a logically structured workflow that addresses the core objectives, including gene prediction, detection of plasmids, integrons, and insertion sequences, and genome visualization. However, several deviations from the GTN tutorial are evident. Notably, GPT-4o recommends *Prokka* for gene annotation instead of *Bakta*, which is the preferred tool in the GTN due to its active maintenance and more comprehensive annotation features. While the inclusion of *PlasmidFinder* and *IntegronFinder* aligns with the GTN protocol, the model’s

suggestion to use tools like *ISMMapper* or *MobileElementFinder* for insertion sequence identification introduces challenges as these tools are not readily available within standard Galaxy instances or the Galaxy ToolShed, limiting practical reproducibility. Additionally, GPT-4o recommends *Trackster* for visualization, whereas the GTN workflow uses *JBrowse*, which offers better integration for multi-track genomic visualization in Galaxy. Another limitation of the GPT-4o-generated workflow is the omission of key preprocessing steps emphasized in the GTN tutorial, such as renaming FASTA headers to conform with tool input requirements, and the absence of instructions for combining and formatting output files prior to visualization. These gaps underscore the challenge of translating high-level workflow synthesis into fully executable, platform-aware implementations. While GPT-4o captures the overarching intent of the annotation pipeline, it lacks the implementation-specific knowledge needed to produce a robust and directly usable Galaxy workflow.

Since the workflow generated by GPT-4o in response to the instruction-only prompt included tools not available within the Galaxy platform, such as *ISMMapper*, *MobileElementFinder*, and *Trackster*, we escalate to a role-based prompt (W5, Table 6). The goal is to assess whether providing a more context-aware prompt can guide the model to recommend Galaxy-compatible tools and generate a workflow that adheres more closely to the platform’s actual capabilities. We observe that role-based workflow uses the *Bakta* tool for comprehensive genome annotation. It begins with uploading an assembled FASTA file and utilizes *Bakta* to predict genes and functional elements, generating standardized outputs such as GFF3, GenBank, and FASTA files. However, the GPT-4o-generated workflow incorporates an evaluation step using QUAST (and optionally BUSCO) to assess annotation quality and completeness. It also introduces post-processing steps such as GFF3 sorting and genome visualization via *Trackster*. Despite this conceptual extension, we notice that certain tools proposed here, such as *Trackster* for interactive visualization and *GFF3sort* for annotation formatting, are currently not available in Galaxy instances. Thus, we are unable to develop the workflow.

We then proceed with the chain-of-thought (CoT) prompting strategy to further guide GPT-4o in generating a bacterial genome annotation workflow (W5, Table 6). This prompt encourages the model to reason through the task step by step, with the aim of producing a more complete and contextually grounded workflow. While the CoT-generated workflow successfully captures the overarching goals of bacterial genome annotation, namely, structural and functional annotation, quality assessment, and preparation of outputs for visualization, it diverges in several important ways from community-curated standards such as those outlined in the GTN tutorial. Both workflows correctly employ *Bakta* as the core annotation engine. However, the CoT-generated workflow introduces *BUSCO* for completeness evaluation, a widely used tool in general genome annotation but not included in the GTN protocol for bacterial genome workflows. More significantly, the GTN tutorial emphasizes a broader approach to structural annotation by incorporating specialized tools such as *PlasmidFinder*, *IntegronFinder*, and *ISEScan*, which detect plasmids, integrons, and insertion sequences, respectively, elements that are essential for accurately characterizing mobile genetic elements in bacterial genomes. The omission of these tools from the generated workflow leads to a structurally incomplete annotation. Additionally, the workflow proposes the use of formatting and indexing tools such as *GFF3sort*, *GFF3 to GenBank*, and *SAMtools faidx* for organizing annotation outputs and enabling visualization. These tools, however, are not available within the Galaxy platform.

In summary, while the generated workflow aligns with the conceptual goals of bacterial genome annotation, it contains key deviations, specifically, the exclusion of critical structural annotation tools and the suggestion of tools not available in Galaxy. As such, the GPT-4o-generated workflow is partially correct but cannot be considered fully usable without further refinement.

We then move to Gemini 2.5 Flash, where we provide the same instruction-only prompt used in the previous test. We observe that Gemini 2.5 Flash is able to return a well-structured, comprehensive set of steps along with relevant parameter details. It correctly identifies and integrates all the necessary tools for structural and functional annotation, including additional conversion and visualization steps. Notably, all suggested tools are available within the Galaxy platform, making the workflow immediately executable without requiring external scripting or unsupported components. For Gemini, the instruction-only prompt is sufficient to generate a complete, accurate, and Galaxy-compliant bacterial genome annotation workflow.

The workflow generated by DeepSeek-V3 using an instruction-only prompt demonstrates a generally accurate understanding of the bacterial genome annotation process, yet diverges from the GTN standard in several key areas. For gene prediction, DeepSeek-V3 recommends *Prodigal*, a well-known tool for prokaryotic CDS identification. For functional annotation, it suggests *EggNOG-mapper* or *BLAST+*, both of which are valid in general but require manual output conversion to GFF3, unlike

the GTN tutorial’s use of *Bakta*, which integrates both gene prediction and functional annotation while producing standardized outputs natively compatible with downstream tools like *JBrowse*. On the other hand, the tools suggested by DeepSeek-V3 for structural annotation *PlasmidFinder*, *IntegronFinder*, and *ISEScan*, are consistent with those used in the GTN tutorial, indicating a correct understanding of the critical elements required to capture plasmid-borne features and mobile genetic elements. However, for non-coding RNA prediction, DeepSeek-V3 proposes the use of *Infernal*, a powerful RNA homology search tool, which is unfortunately not available on Galaxy, again highlighting a case of tool hallucination. Despite these issues, DeepSeek-V3 correctly outlines the merging of GFF3 outputs and the use of *JBrowse* for genome visualization, aligning well with the GTN approach in terms of output integration and browser-based annotation exploration.

*In comparison to GPT-4o and Gemini 2.5 Flash, DeepSeek-V3 performs moderately well. Like GPT-4o, it hallucinates the availability of tools not supported in Galaxy, reducing practical executability. Gemini, by contrast, stands out for producing a structurally complete and Galaxy-native workflow using only the instruction prompt. While DeepSeek-V3 captures the correct conceptual flow and covers all major analytical components, its output is hindered by limited platform awareness and reliance on tools outside the Galaxy environment. In summary, DeepSeek-V3 offers strong logical grounding but lacks the operational precision demonstrated by Gemini.*

Among the three LLMs, Gemini 2.5 Flash consistently performs best, producing complete, accurate, and Galaxy-compliant workflows without requiring additional prompting. Its outputs demonstrate strong contextual understanding, correct tool selection, detailed parameter guidance, and seamless alignment with Galaxy’s capabilities, making it particularly effective for both novice and intermediate users. DeepSeek-V3 shows a solid grasp of workflow logic and offers methodologically valid suggestions, including flexible tool alternatives and usability enhancements, but is limited by occasional tool hallucinations, recommending software not available within Galaxy, which undermines executability. GPT-4o, while concise and easy to follow, often omits critical steps or recommends unsupported tools, requiring further prompting to reach usability. Overall, Gemini stands out for its balance of technical completeness, clarity, and platform awareness, making it the most reliable LLM for automated Galaxy workflow development.

#### 4.1.2 Workflow Development using LLMs for Nextflow:

To evaluate the extent to which LLMs can assist in developing scientific workflows in bioinformatics using Nextflow, we curate a diverse set of representative workflows from the nf-core framework (Table 1). The selected workflows (e.g., *fetchngs*, *methyseq*, *rnaseq*, *phyloplace*) span multiple domains, including genomics, transcriptomics, epigenomics, and regulatory genomics, thereby enabling a comprehensive and domain-balanced assessment. The *nf-core/demo* pipeline serves as a foundational example, involving straightforward quality control and adapter trimming of single- or paired-end FASTQ reads, making it useful for testing basic workflow generation capabilities. *Fetchngs* introduces more advanced functionality by automating the retrieval of sequencing data from public repositories, challenging models to reason about metadata management, remote data access, and parsing of structured sample sheets. *Rnaseq* and *methyseq* represent complex, end-to-end pipelines that integrate numerous modules, manage intricate dependencies, and involve advanced data transformations, ideal for evaluating the models’ ability to recommend appropriate toolchains, configurations, and parameters. The inclusion of *phyloplace* extends the evaluation into the domain of phylogenetic analysis, specifically testing the models’ ability to handle workflows that map query sequences onto a reference phylogenetic tree using specialized inputs and classification tools. These workflows offer a robust and multifaceted benchmark suite to assess how effectively LLMs understand and generate Nextflow-based scientific workflows.

The *nf-core/demo* (W1 Table 1) workflow is designed to perform a basic bioinformatics data processing task using a minimal dataset. Its core objective is to demonstrate essential steps commonly used in sequencing data preprocessing. First, it conducts quality control of raw FASTQ files using the module *FastQC* to assess sequence quality metrics. Next, it applies module *Seqtk* to trim adapter sequences and low-quality bases, improving the quality of the reads. Finally, it aggregates the quality control results using *MultiQC*, generating a summary report that provides an overview of the data quality before and after trimming. While the workflow is simplified and intended for



demonstration, it effectively showcases a typical read preprocessing pipeline involving QC, trimming, and result summarization. We evaluate the workflow using a sample sheet that includes both single-end and paired-end data, available at [https://raw.githubusercontent.com/nf-core/test-datasets/viralrecon/samplesheet/samplesheet\\_test\\_illumina\\_amplicon.csv](https://raw.githubusercontent.com/nf-core/test-datasets/viralrecon/samplesheet/samplesheet_test_illumina_amplicon.csv). The workflow is executed using the following command:

```
nextflow run nf-core/demo --profile docker, test --outdir 'demo-results'
```

We do not delve into the syntax details here, as our primary objective is not to teach Nextflow but to assess whether LLMs can effectively assist in workflow development. Interested readers may consult the official documentation for further guidance. For our purpose, we simply run the workflow on the provided sample dataset and inspect the results. We obtain several output reports from the workflow execution, including FastQC and MultiQC quality control summaries, the execution report, and the pipeline DAG visualization. After verifying these outputs, we proceed to test the same workflow development task using LLMs.

We begin by applying the instruction-only prompt (W1, Table 7) to GPT-4o. The model generates a logically structured set of workflow steps, starting with initial quality control using the *FastQC* module, followed by adapter and quality trimming using tools such as *Trim Galore!*, *fastp*, or *Trim-momatic*. It then suggests performing a second round of quality control to validate the trimming process and concludes with summarizing the results using *MultiQC*. These steps are largely consistent with those in the official nf-core/demo workflow, which also performs quality control, trimming, and reporting. The primary divergence lies in the trimming step: while the nf-core/demo pipeline uses *seqtk*, GPT-4o proposes alternative tools that are widely accepted and functionally equivalent within the context of short-read preprocessing. Overall, GPT-4o demonstrates a strong ability to generate a correct and coherent workflow based solely on an instruction-only prompt, indicating a solid understanding of standard practices in NGS data preprocessing with Nextflow.

We then provide the same instruction-only prompt to *Gemini 2.5 Flash*. The model suggests a detailed and modular workflow design consistent with nf-core best practices. The proposed steps begin with parsing the sample sheet to extract FASTQ file paths using Nextflow channels. It then performs an initial quality assessment using the *FastQC* module, followed by adapter and quality trimming with either *seqtk* or *Trim Galore!*. This is followed by a second round of quality control using *FastQC* on the trimmed reads, and the final step involves aggregating all quality control outputs into a comprehensive report using *MultiQC*. While this plan largely mirrors the nf-core/demo workflow, Gemini extends it by emphasizing modularity and offering both *seqtk* (used in nf-core/demo) and *Trim Galore!* (which wraps Cutadapt) as valid trimming options. Overall, Gemini 2.5 Flash demonstrates a clear understanding of the workflow structure and appropriately maps tools to each processing stage, showing strong competence in workflow planning from an instruction-only prompt.

We also evaluate the response from DeepSeek-V3 using the same instruction-only prompt. The model provides a well-structured and modular workflow outline for FASTQ quality assessment and preprocessing that is closely aligned with nf-core conventions. The workflow begins with parsing the input sample sheet, followed by initial quality control using *FastQC*, adapter trimming via *Cutadapt*, and then a separate quality filtering step using *FastP*. A second round of quality control is performed post-filtering using *FastQC*, and finally, *MultiQC* is used to generate an aggregated summary report. The model explicitly maps each step to specific nf-core modules and provides a clear implementation structure in *main.nf*, along with a sample configuration in *nextflow.config*.

*Among the three models, DeepSeek-V3 provides the most comprehensive and implementation-ready response, outlining a modular Nextflow workflow with nf-core modules, configuration details, and code structure. Gemini 2.5 Flash also performs well, offering accurate steps with moderate implementation guidance and alignment with nf-core practices. In contrast, GPT-4o delivers a concise and correct high-level outline but lacks the depth needed for direct workflow deployment. Overall, DeepSeek-V3 stands out as the most effective model for building fully functional, production-ready Nextflow pipelines based on instruction-only prompts.*

Our second workflow, *nf-core/fetchngs*, is designed to automate the retrieval and preparation of sequencing data from major public repositories including *ENA*, *SRA*, *DDBJ*, and *GEO*. Given a list of accession identifiers, the pipeline retrieves associated metadata and downloads the corresponding raw FASTQ files via *FTP*, *Aspera*, or *SRA tools*. It then compiles a standardized samplesheet and organizes the output in a format that is immediately compatible with downstream nf-core pipelines such as *RNA-seq* or *ATAC-seq*. This automation streamlines the data acquisition process, enhances reproducibility, and reduces the manual burden of managing large-scale sequencing datasets.



To assess LLMs’ ability to support the development of such workflows, we first execute the *nf-core/fetchngs* pipeline using a test CSV file containing sample IDs, confirming that it successfully fetches both FASTQ files and metadata. We then provide GPT-4o with an instruction-only prompt (W2, Table 7) to generate the workflow. The output from GPT-4o closely aligns with the core logic of the *fetchngs* pipeline: parsing sample IDs from a CSV, retrieving metadata, downloading sequencing files, and organizing outputs into structured directories. However, GPT-4o demonstrates greater flexibility by suggesting the integration of additional tools, such as *pysradb* utilities, which are particularly useful for SRA or GEO-specific workflows. GPT-4o also suggests the quality checking of the FAST files. This flexibility, while advantageous, also necessitates more careful implementation, especially in resolving ID types, managing errors, and ensuring reproducible execution through proper containerization.

We then provide the same instruction-only prompt to *Gemini 2.5 Flash*. The workflow generated by Gemini 2.5 Flash for fetching NGS data is notably detailed and architecturally robust, closely aligning with the *nf-core fetchngs* pipeline in both structure and objective. Like *fetchngs*, it accepts a CSV of public sample IDs (e.g., SRA, ENA, GEO, DDBJ), retrieves metadata using the ENA API, and downloads FASTQ files using both direct FTP methods and SRA tools as a fallback. Both workflows include metadata collation and output organization steps. However, Gemini’s generated design adds granularity by modularizing each logical step, e.g., separating *download\_fastq\_direct* and *download\_fastq\_sra*, validating MD5 checksums, and producing a final *samplesheet.csv*, and assigning clear containerized environments with resource constraints for each module.

One advantage of the Gemini-generated workflow is its explicit fallback mechanism: if direct FASTQ downloads fail, it reroutes to SRA tools (*prefetch* + *fastq-dump*) selectively. Additionally, the inclusion of detailed validation, resource configuration, and Python-based logic for metadata parsing shows an advanced awareness of real-world edge cases (e.g., batch API querying, empty downloads, ENA metadata field inconsistencies). Compared to *fetchngs*, which abstracts many steps behind a high-level interface, Gemini’s approach makes the internals transparent, facilitating educational or research-focused exploration. In the context of evaluating LLM capabilities, Gemini 2.5 Flash demonstrates strong potential in guiding sophisticated scientific workflow design, though its correctness and maintainability would benefit from cross-validation against established pipelines like *fetchngs*.

We then give the same prompt to the DeepSeek-V3. The generated workflow shares strong structural alignment with the *nf-core fetchngs* pipeline in retrieving metadata and FASTQ files from public repositories using a list of accessions. DeepSeek-V3 begins by parsing a CSV input, fetching metadata using ENA or SRA APIs, downloading FASTQ files (via *enaDataGet* or *fasterq-dump*), and organizing outputs. However, it diverges in its emphasis on user-defined metadata (project, condition) and its incorporation of flexible tuple-based channels to carry metadata across modules. This offers a lightweight yet extensible design. Its use of tuple-passing enables downstream modules (like *organize\_output*) to rename files and structure directories using custom project metadata, a flexibility *fetchngs* does not natively offer. It also outlines future-proof enhancements like MD5 validation, cloud output integration, and notifications, suggesting practical awareness of real-world deployment. However, compared to the rigorously tested *fetchngs*, the DeepSeek-V3 workflow requires more manual scripting and error handling, particularly around metadata inconsistencies and download failures. It also lacks built-in support for containerized environments, dependency management, and profiling, which are strengths of *nf-core* workflows. Overall, DeepSeek-V3 demonstrates strong conceptual competence in building modular, extensible workflows and offers a solid foundation for assessing LLM capabilities in scientific workflow construction.

*The workflow generated by GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3 for W2: nf-core/fetchngs, all demonstrate the capability of LLMs to assist in workflow development, particularly in automating metadata retrieval and FASTQ file downloads from public sequencing databases. GPT-4o provides a concise, modular breakdown that aligns well with nf-core’s fetchngs, while also offering flexibility through tool alternatives like pysradb and Entrez utilities. Gemini 2.5 Flash delivers the most production-ready solution, featuring a fully containerized, resource-configured, and modular DSL2 workflow that includes fallback strategies, MD5 validation, and samplesheet generation, closely mirroring the robustness of nf-core standards. DeepSeek-V3 emphasizes conceptual modularity with user-defined metadata propagation, conditional tool logic, and practical extensibility for cloud and notification integration, though it lacks some execution-level rigor. While all three LLMs show potential, Gemini 2.5 Flash outperforms the others in completeness, reproducibility, and implementation detail, making it particularly well-suited for generating executable, real-world scientific workflows.*

We then move to our next workflow *nf-core/methylseq*, a robust, containerized pipeline designed for comprehensive analysis of bisulfite-converted sequencing (BS-seq) data to facilitate DNA methylation profiling. It encompasses all critical stages of methylation analysis, including data pre-processing, genome alignment, post-alignment processing, and quality control reporting. The workflow automates key tasks such as raw read quality assessment with *FastQC*, adapter trimming using *Trim Galore!*, and alignment to a reference genome via either *Bismark* (with *Bowtie2* or *HISAT2*) or the *bwa-meth* and *MethylDackel* toolchain, with optional GPU acceleration supported through *Parabricks*. In our implementation, we utilize the Bismark-based workflow. The pipeline further includes deduplication of aligned reads, extraction of cytosine methylation metrics, and comprehensive quality evaluations using *Qualimap*, *Preseq*, and *HS Metrics*. Outputs include deduplicated BAM files, detailed methylation calls, and bias reports, all of which are consolidated into a single MultiQC report for streamlined visualization. Additionally, the pipeline generates extensive provenance metadata, including execution logs, software versions, and configuration parameters, ensuring full reproducibility across computing environments. Its modular structure and support for high-performance computing platforms make it well-suited for scalable, large-scale epigenomic studies.

To evaluate LLM support for constructing this workflow, we first test GPT-4o using an instruction-only prompt (W3, Table 7). We observe that the workflow generated by GPT-4o exhibits several inaccuracies while suggesting modules, even though the overall sequence of analytical steps was conceptually sound. For example, the model suggests a module called *fastqc.raw* for raw read quality control, which does not exist in the official Nextflow or nf-core module repositories; the correct module name is simply *fastqc*. Similarly, for adapter and quality trimming, the model proposes *trim.reads*, whereas the appropriate module is *trimgalore*. These naming inconsistencies are found throughout the generated workflow, suggesting that the model’s responses are not grounded in the actual set of available Nextflow modules. Consequently, we adopt a role-based prompting (W3, Table 7) approach to improve the specificity and correctness of the module recommendations.

We evaluate the workflow produced by GPT-4o using the role-based prompt. The GPT-4o-generated workflow correctly captures the high-level stages of DNA methylation analysis using bisulfite-converted sequencing (BS-seq) data, beginning with FASTQ quality control and adapter trimming, followed by bisulfite-aware alignment, deduplication, methylation extraction, and MultiQC-based reporting. Importantly, the model successfully identified the correct nf-core module names such as *fastqc*, *trim.galore*, *bismark/genomeprepare*, *bismark/align*, *bismark/deduplicate*, and *bismark/methylationextractor*, indicating that role-based prompting can enhance grounding in actual Nextflow components. However, the generated workflow also exhibits several limitations. First, it does not include configurable options for alternative aligners (e.g., *bismark\_hisat*) or GPU-accelerated execution. Second, while the model included key optional QC modules like *qualimap*, *preseq*, and *picard/collectHsmetrics*, it does not explain when or why to use them, missing opportunities to contextualize their role in specific experimental designs (e.g., capture-based methylation assays). Third, the model omitted support for protocol-specific presets such as RRBS (–rrbs) and PBAT (–pbat), which are natively handled in the nf-core implementation. Finally, although the core logic was structurally correct, some processes like CSV parsing are abstracted as custom steps rather than referencing existing solutions.

To assess the capabilities of Gemini 2.5 Flash in constructing complex scientific workflows, we evaluate its response to an instruction-only prompt (W3, Table 7) for developing a DNA methylation analysis pipeline using bisulfite-converted sequencing (BS-seq) data. The model successfully outlines all critical stages of the workflow, including CSV-based input parsing, raw read quality control, adapter trimming, bisulfite-aware alignment, PCR duplicate removal, methylation extraction, and comprehensive quality control reporting. It recommends appropriate and widely used tools such as *FastQC*, *Trim Galore!*, *Bismark*, *MethylDackel*, *Qualimap*, and *MultiQC*, and its modular structure reflects a good understanding of Nextflow DSL2 practices, including containerization, channel dependencies, and resource management. Furthermore, the workflow is logically organized and aligns well with the stages implemented in the official nf-core/methylseq pipeline. However, Gemini shows some limitations. It does not fully specify the conditional logic or flag handling needed to toggle between these options, nor does it elaborate on the inclusion of protocol-specific parameters. Despite these minor gaps, Gemini 2.5 Flash demonstrates strong proficiency in designing structurally sound and tool-appropriate BS-seq workflows, making it a capable model for assisting in the development of complex, real-world bioinformatics pipelines.

We then, evaluate DeepSeek-V3 for its ability to generate a Nextflow workflow for bisulfite-converted sequencing (BS-seq) data analysis using the same instruction-only prompt. The model

successfully outlines the key stages of a BS-seq workflow, including input parsing, quality control with FastQC, adapter, and base trimming using Trim Galore!, bisulfite-aware alignment with Bismark or bwa-meth, deduplication, methylation extraction, and downstream quality reporting with MultiQC. The response is structured, includes example CSV formats, command-line syntax for each tool, and even proposed a skeleton main.nf script and directory structure. This indicates that DeepSeek-V3 has a solid conceptual grasp of both BS-seq data analysis and the Nextflow DSL2 design pattern. However, the workflow generated by DeepSeek-V3 also exhibits critical shortcomings. While it named commonly used tools, workflow steps like FASTQC, TRIM\_GALORE, or BISMARK are presented as abstract module calls rather than as references to actual modules in the nf-core/modules repository, limiting reproducibility. Additionally, the pipeline lacks any reference to containerization profiles, resource management, or error handling, all critical for robust, production-ready deployment. However, the generated steps are correct, and we can develop the workflow following the instructions.

*In comparing GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3 for generating a Nextflow workflow for DNA methylation analysis using bisulfite-converted sequencing (BS-seq) data, we uncover notable differences in how each model responds to different prompting strategies. Under instruction-only prompting, Gemini 2.5 Flash produces the most accessible and beginner-friendly output. It delivers logically structured steps, accurate tool selection, and detailed descriptions suitable for users with limited workflow development experience. DeepSeek-V3 also performed well under instruction-only prompting, offering a clear conceptual outline along with practical elements such as directory structure, example CSV input, and bash commands. This makes it particularly valuable for intermediate users who seek both clarity and low-level control. In contrast, GPT-4o initially underperforms under instruction-only prompting, often hallucinating incorrect module names (e.g., fastqc\_raw, trim\_reads) and tool references. However, when guided with a role-based prompt, GPT-4o produces the most technically precise and nf-core-compliant workflow, correctly referencing official modules such as fastqc, trim\_galore, and bismark/align. In summary, Gemini is best positioned for novice users, followed by DeepSeek-V3 for intermediate users, while GPT-4o offers the highest accuracy and implementation fidelity for advanced users through prompt specialization.*

Our next workflow, *nf-core/rnaseq*, processes raw RNA-seq data into high-quality, analysis-ready outputs. Its main objective is to perform comprehensive quality control, adapter trimming, read alignment, transcript quantification, and reports in a reproducible and scalable manner. The workflow begins with raw FASTQ files and a metadata samplesheet, and takes a reference genome (FASTA) and annotation file (GTF). It includes modules for inferring strand specificity, performing quality checks (e.g., FastQC, MultiQC), trimming adapters (Trim Galore!), filtering contaminants (BBSplit), removing rRNA (SortMeRNA), and aligning reads using STAR or HISAT2. Quantification is performed using Salmon, RSEM, or StringTie, with optional pseudoalignment via Salmon or Kallisto. The pipeline also supports UMI handling, duplicate marking, transcript assembly, and generation of coverage files. Extensive quality metrics are collected using tools such as *RSeQC*, *Qualimap*, *Preseq*, and *dupRad*. All results are summarized in an interactive MultiQC report. Designed with portability and reproducibility in mind, the pipeline can be run using Docker, Singularity, or Conda, making it ideal for large-scale or collaborative transcriptomic studies. We execute the workflow using the sample dataset and obtain the desired results. To evaluate the ability of LLMs to generate this complex workflow, we provide an instruction-only prompt (W4, Table 7) and assess whether the models can accurately develop *rnaseq* pipeline.

We observe that the GPT-4o-generated workflow effectively covers the core analytical stages of RNA-seq data processing, including quality control, adapter trimming, genome alignment, gene/transcript quantification, and differential expression analysis. It accepts a CSV sample sheet with FASTQ links and reference genome inputs, making it flexible and suitable for standard analyses. However, the GPT-4o response has several limitations. The *nf-core* workflow includes a wider range of features such as pseudoalignment options (e.g., Salmon and Kallisto), UMI handling, strandness autodetection, biotype-level QC, spike-in normalization, and automated GTF/GFF validation, all of which enhance robustness and reproducibility. Moreover, *nf-core* pipelines benefit from strict input schema validation, containerization, continuous integration testing, and broad community support, making them more reliable for large-scale or production-grade use. The absence of these advanced features in the custom workflow may not affect small or focused studies but could introduce reproducibility challenges and limit flexibility in more complex experimental designs. However, the steps we obtain are fully functional and easygoing.

We evaluate the output of Gemini 2.5 Flash for constructing a Nextflow-based RNA-seq workflow using an instruction-only prompt. The generated workflow successfully captures all key analytical

stages, ranging from sample sheet parsing and raw read quality control to adapter trimming, reference genome indexing, alignment or pseudo-alignment, transcript quantification, differential expression analysis, and quality reporting, demonstrating strong structural alignment with the nf-core/rnaseq pipeline. Gemini provides flexibility in tool selection, supporting both traditional aligners (e.g., STAR) and pseudo-aligners (e.g., Salmon, Kallisto), and adopts modular design principles using channels and containers, consistent with best practices in DSL2-based Nextflow development. However, the Gemini-generated workflow lacks some infrastructure features, such as automated strandedness detection, UMI processing, GTF validation, biotype-level QC summaries, and reproducibility guarantees via CI testing and schema validation. Moreover, while it mentions statistical tools for differential expression (e.g., DESeq2), it does not detail complex contrast designs or robust metadata handling, which are integral to nf-core. As a result, although Gemini provides a conceptually complete and interpretable scaffold suitable for developing the workflow, it can be improved.

We further evaluate the performance of DeepSeek-V3 using the same instruction-only prompt. The RNA-seq workflow generated by DeepSeek-V3 presents a clear, methodical structure that encompasses all essential stages of RNA-seq data processing, starting from sample sheet parsing and reference genome indexing to quality control, adapter trimming, alignment, quantification, differential expression analysis, and final reporting. Its design closely mirrors the nf-core/rnaseq workflow in terms of both modularity and core functionality. DeepSeek appropriately suggests widely adopted tools such as STAR and HISAT2 for alignment, Salmon for quantification, and FastQC, Qualimap, and MultiQC for quality assessment. In addition, the model provides practical implementation notes, including the use of *main.nf*, *modules.config*, and *nextflow.config*, which enhance its utility for users with experience in pipeline development. However, it lacks some capabilities integrated into nf-core, such as automated strandness detection, UMI handling, TPM/FPKM normalization options, comprehensive biotype-specific QC metrics, and rigorous input schema validation. DeepSeek also includes practical implementation notes using *main.nf*, *modules.config*, and *nextflow.config*. DeepSeek’s generated steps are methodologically sound and capture key conceptual elements of a robust RNA-seq workflow, making it suitable for users familiar with pipeline customization.

*The RNA-seq workflow outputs generated by GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3 all cover the essential analytical stages, quality control, trimming, alignment, quantification, and differential expression analysis but differ significantly in depth, infrastructure awareness, and production readiness. GPT-4o provides a concise and methodologically correct outline but lacks some features such as pseudo-alignment options, detailed quality control tools, and workflow execution infrastructure, making it the least complete among the three. Gemini 2.5 Flash offers broader tool coverage, supporting both alignment-based and alignment-free methods, and presents a modular structure that is conceptually clear and well-suited for educational use; however, it stops short of providing implementation details like Nextflow process logic or configuration scaffolding. In contrast, DeepSeek-V3 delivers the most comprehensive and technically mature output, it not only includes accurate tool suggestions and optional enhancements (e.g., rMATS, ComBat) but also outlines concrete Nextflow components (main.nf, modules.config, nextflow.config), container usage, and execution profiles, closely resembling the structure and flexibility of the nf-core/rnaseq pipeline. Thus, DeepSeek-V3 outperforms the others in terms of completeness and deployability, followed by Gemini for clarity and breadth, with GPT-4o being the most minimal in practical usability.*

We conclude our evaluation with the *nf-core/phyloplace* workflow, which is designed to facilitate efficient and reproducible phylogenetic placement of query sequences onto a reference tree. The workflow operates in two distinct modes: (i) placement-only, where input sequences are directly aligned to a reference multiple sequence alignment and placed onto a fixed phylogenetic tree using likelihood-based methods; and (ii) search-plus-placement, where input sequences from a larger, unfiltered dataset are first screened using profile-based searches to identify candidates for subsequent alignment and placement. After placement, the workflow grafts the query sequences onto the reference tree and performs taxonomic classification, generating annotated phylogenies, placement summaries, and visualizations. Developed using Nextflow DSL2 and fully containerized modules, *nf-core/phyloplace* ensures portability, scalability, and reproducibility across diverse computing environments. It is particularly well-suited for high-throughput analyses in environmental and metagenomic research, where rapid and accurate phylogenetic integration is critical.

We execute the workflow following the instruction of *nf-core/phyloplace*. The workflow generates a comprehensive set of output files that summarize the results of the phylogenetic placement and classification process. The workflow produces several key output files organized per sample. These

include: (i) placement files containing likelihood-based placements of query sequences on the reference tree; (ii) grafted phylogenetic trees in Newick format with query sequences inserted into the reference topology; (iii) classification tables summarizing the taxonomic or user-defined annotation of each placed sequence; (iv) sequence alignments showing how queries are aligned to the reference multiple sequence alignment; and (v) visualizations and summary reports, such as heat trees and placement statistics, generated via automated reporting tools. A final MultiQC report aggregates log and quality-control metrics across all samples. These outputs facilitate interpretation, visualization, and downstream analysis of the phylogenetic context of newly placed sequences.

To evaluate the ability of LLMs to assist in developing the *nf-core/phyloplace* workflow, we begin by providing an instruction-only prompt (W5, Table 7) to GPT-4o. The resulting workflow aligns closely with the official *nf-core/phyloplace* implementation in terms of structure, logic, and tool selection. GPT-4o successfully captures both operational modes, placement-only and search-plus-placement, and accurately identifies the required input components, such as query FASTA files, reference alignments, rooted phylogenetic trees, evolutionary models, HMM profiles, and optional taxonomy files. The model outlines key steps including candidate filtering using *hmmsearch*, sequence alignment via *mafft*, phylogenetic placement using *epa-ng*, tree grafting with *gappa*, optional taxonomic classification, and summarization with *MultiQC*. In addition, GPT-4o includes implementation-aware suggestions such as input validation, modular output organization, and execution metadata logging, practices consistent with *nf-core* standards.

We then provide the same instruction-only prompt to Gemini 2.5 Flash. We observe that it accurately captures both operational modes and correctly identifies the necessary input components, including query and reference sequences, phylogenetic tree, evolutionary model, HMM profiles, and optional taxonomy. It organizes the workflow into well-defined functional modules such as sequence alignment, phylogenetic placement, classification, and report generation and correctly associates standard tools like MAFFT, HMMER (*hmmsearch*), EPA-NG, and GAPPA with each step. Though, it omits some *nf-core*-specific features such as execution metadata (e.g., DAGs and trace files), the overall structure is technically sound and implementation-ready.

We finally give the same prompt to DeepSeek-V3. It also captures both placement-only and search-plus-placement modes and structures the pipeline into logically ordered modules. It specifies necessary input parameters, including reference alignment, rooted phylogenetic tree, evolutionary model, taxonomy file, and HMM profiles, and introduces validation steps to ensure input integrity. The outlined modules, covering HMM-based sequence filtering (*hmmsearch*), alignment (MAFFT or HMMER), placement (EPA-NG or *pplacer*), and classification (*gappa* or *taxit*), are consistent with tools used in the official pipeline. DeepSeek-V3 additionally incorporates optional visualization using tools such as ETE3, R, and Python, and provides a pseudocode-level implementation using Nextflow syntax that demonstrates channel logic and parameter control.

*All three LLMs, successfully reconstruct a modular Nextflow workflow for phylogenetic placement, capturing both operational modes and correctly identifying core steps such as input validation, sequence alignment, phylogenetic placement, classification, and reporting. GPT-4o provides a technically accurate and semantically faithful outline closely aligned with the nf-core implementation. Gemini 2.5 Flash presents a well-organized modular architecture using intuitive labels and standard tools, though it omits some nf-core-specific features like execution metadata. DeepSeek-V3 stands out by offering the most implementation-ready response, including detailed parameter definitions, alternative tools (e.g., pplacer, taxit), and a fully structured Nextflow pseudocode, making it especially valuable for advanced users or developers seeking practical scaffolding. While all three are valid, DeepSeek-V3 provides the most comprehensive and adaptable solution, GPT-4o ensures the highest fidelity to the reference workflow, and Gemini excels in clarity for novices. Therefore, DeepSeek-V3 is best suited for expert-driven implementation, GPT-4o for faithful documentation reproduction, and Gemini for accessible modular reasoning.*



In summary, the results demonstrate that LLMs can substantially assist in developing scientific workflows for both Galaxy and Nextflow platforms. Gemini 2.5 Flash consistently outperforms other models for Galaxy workflows by generating accurate, tool-aware, and pedagogically rich workflows, offering the best balance between technical detail and usability. Its outputs are both complete and executable within the Galaxy environment, making it the most effective for supporting workflow design in this system. For Nextflow, DeepSeek-V3 emerges as the most capable model. It generates structurally sound, implementation-ready workflows that align closely with nf-core standards, including modular Nextflow DSL-2 structures, accurate tool mapping, container usage, and configuration files. While GPT-4o and Gemini provide reasonable outputs, they either lack the depth of implementation (GPT-4o) or fall short on infrastructure-level detail and flexibility (Gemini) compared to DeepSeek. Overall, the findings confirm that LLMs can play a transformative role in supporting workflow development, with model performance varying by system and use-case complexity.

## 4.2 RQ2: How do workflows generated by different LLMs compare regarding completeness, correctness, and usability?

**Motivation:** In scientific workflow development, particularly in fields like bioinformatics, the effectiveness of a workflow depends not only on its syntactic validity but also on its *completeness, correctness, and usability*. *Completeness* refers to the inclusion of all essential analytical steps required for the task. *Correctness* ensures the logical coherence of the workflow, including the use of appropriate tools, parameters, and configurations. *Usability* evaluates how easily the workflow can be understood, executed, and maintained, especially in line with community standards and platform-specific conventions. Assessing these dimensions is crucial for determining the extent to which LLMs can assist novice and expert users in designing robust, reproducible workflows for real-world applications.

**Approach:** We evaluate ten representative bioinformatics workflows across two major Scientific Workflow Systems (SWSs): Galaxy (graphical interface) and Nextflow (script-based). For each workflow task, we provide a natural language prompt (for a few cases, more prompts) to each LLM and collect the generated workflow. To ensure fair assessment, prompts are designed consistently, escalating from instruction-only to role-based and chain-of-thought when initial results lack adequacy.

The generated workflows are then manually assessed by two domain experts with over five and nine years of experience and compared to established community baselines from the Galaxy Training Network (GTN) and nf-core repositories. Evaluators rate each workflow across three dimensions:

- **Completeness:** Does the workflow include all essential steps from input handling to result generation?
- **Correctness:** Are the tools correctly chosen, configured, and logically sequenced?
- **Usability:** Can the workflow be executed with minimal modification, and does it adhere to standard practices?

Assessments considered both the logical structure and contextual relevance of the generated workflows, with special attention given to domain appropriateness and the presence of clear documentation or commentary.

**Results:** The comparative evaluation reveals that no single LLM is universally superior across all scenarios. Rather, performance varies by platform and task complexity.

- **Galaxy workflows:** Gemini 2.5 Flash consistently outperforms the other models in the Galaxy context. Its responses are well-structured, domain-aware, and aligned closely with GTN tutorials. Gemini captures platform-specific toolchains, provides a clear rationale for each step, and produces workflows that are both accurate and accessible. These traits made Gemini particularly effective for Galaxy, where usability and clarity are critical for non-programming users. This finding aligns with the results from **RQ1**, where Gemini demonstrated the strongest understanding of Galaxy’s architecture and tool ecosystem.
- **Nextflow workflows:** DeepSeek-V3 emerges as the most effective model for generating Nextflow pipelines. It produces workflows with detailed implementation logic, often including configuration blocks, directory structures, and explicit references to nf-core modules. These outputs are technically rich and closer in structure to production-ready scripts. While verbose at times, DeepSeek’s

workflows required less modification to be operational and demonstrated an in-depth understanding of modular, containerized environments. This strength in script-based reasoning is also evident. DeepSeek-V3 shows strong familiarity with Nextflow concepts and execution patterns.

GPT-4o, while generally well-performing, shows mixed results. Its workflows are often logically sound and clearly presented, making them useful in instructional contexts. However, it occasionally omits setup steps, assumes preprocessed inputs, or lacks sufficient detail for immediate execution. Its performance improves significantly under role-based or chain-of-thought prompting, suggesting that GPT-4o is best suited for users who can guide it with more context-aware instructions.

Across all evaluated workflows, the three qualitative dimensions, completeness, correctness, and usability, reveal complementary strengths and limitations among the LLMs. In terms of completeness, Gemini consistently includes all major processing steps and data transformations, especially in Galaxy workflows, often mirroring the structure of official tutorials. DeepSeek also performs well in completeness, particularly for Nextflow, where its outputs capture not only the high-level steps but also underlying configuration details such as profiles and container usage. Correctness is generally high across all models, though some models (especially GPT-4o) occasionally select reasonable but suboptimal tools, or miss platform-specific conventions. Gemini demonstrated strong correctness by adhering to community standards (e.g., ToolShed in Galaxy, nf-core in Nextflow), while DeepSeek shows tool-level precision, particularly in script-based workflows. The greatest variance is observed in usability. Gemini’s outputs are highly readable, modular, and well-suited for instructional or real-world use, especially by users unfamiliar with the workflow system. GPT-4o excels in usability when prompts are structured, though it sometimes lacks operational context. DeepSeek, despite its technical accuracy, often produces verbose and densely structured workflows that could overwhelm novice users, requiring post-processing or simplification for practical adoption. These dimension-wise differences underline the importance of aligning LLM choice with user expertise and platform characteristics.

These results highlight that the effectiveness of an LLM for scientific workflow generation is strongly dependent on the characteristics of the workflow system in use. For graphical environments like Galaxy, where clarity and usability are paramount, Gemini is most effective. For script-based environments like Nextflow, which demand technical completeness and execution logic, DeepSeek performs best. Prompting strategy and model-specific tuning play a critical role in maximizing output quality.

#### 4.2.1 RQ3: What prompting strategies should a workflow developer follow?

**Motivation:** LLMs are sensitive to the way users phrase their requests, and the same task may yield vastly different outcomes depending on the prompting strategy. While LLMs have shown promise in generating scientific workflows, the quality, structure, and reliability of the outputs are not solely functions of the model’s architecture, they are also critically shaped by how the task is presented. For workflow developers, especially those unfamiliar with prompt engineering, this variability introduces uncertainty and inefficiency in practice. Therefore, understanding which prompting strategies, such as instruction-only, role-based, or chain-of-thought, produce more accurate and usable outputs is essential. This RQ investigates how different prompt formulations impact the generation of workflows in terms of correctness, completeness, and usability across models. By identifying effective strategies, we aim to empower developers with practical guidance for interacting with LLMs to support reproducible and efficient scientific workflow development.

**Approach:** Professionals can achieve various objectives through effective prompts, such as question answering, complex or sequential analysis, and code generation. However, the potential of effective prompting has yet to be fully realized within the workflow community. We recognize the need to demonstrate how effective prompting strategies can empower workflow developers by enhancing automation, improving data handling, and streamlining complex workflows. To address this, we aim to explore existing prompting strategies and identify the most effective approaches for our analysis by exploring the prompting strategies [28–34, 62–73]. Effective prompting indicates a strong likelihood of receiving the desired recommendations [28, 29, 31–34]. In the following, we summarize our strategies to create effective prompting for our work.

##### **Results:**

**Clear and Specific Instructions:** Providing clear and specific instructions in LLMs is essential for generating the desired output. Ambiguity in the prompt can lead to unexpected responses. For each workflow, we provide precise instructions to get the responses, and most of the time, we receive

the desired responses. For example, we ask *What is the purpose of the Galaxy ToolShed?*. This background question is direct and scoped to a specific component of the Galaxy ecosystem. Models like Gemini 2.5 Flash respond with clear definitions and references to community-driven tool sharing, which makes comparison across LLMs reliable.

**Experimenting with context and examples:** Incorporating well-defined context and relevant examples in prompts can significantly improve the ability of LLMs to generate accurate and domain-specific workflows. Contextualizing prompts ensures that LLMs understand the nuances of the task, while examples serve as guides to steer the model toward producing precise and actionable outputs. This approach is essential when working with specialized domains where generic responses may lack relevance. By providing clear context and concrete examples, we are able to design complex workflows tailored to specific scientific applications. For instance, our prompt, *Create a Galaxy workflow that identifies the exon with the highest number of SNPs on human chromosome 22 using VCF and GTF files* is clear about the biological target (exon with most SNPs), the input file types (VCF, GTF), and the chromosome of interest, helping LLMs produce a complete and executable workflow.

**Leveraging System 1 and System 2 questions:** To improve the LLMs’ responses, understanding the concept of System 1 and System 2 questions is essential. System 1 questions typically involve quick, intuitive, or pattern recognition-based answers. For example, *What does a DAG (Directed Acyclic Graph) represent in scientific workflows?*. On the contrary, System 2 questions involve more deliberate, analytical, complex problem-solving and effortful thinking. For example, *How does the design of a scientific workflow system affect its scalability and performance in large-scale data processing?*

**Contextual Role Assignment:** Embedding the model in a specific expert role improves technical relevance and response structure. For instance, the prompt *You are a bioinformatics workflow developer with expertise in building reproducible pipelines using Nextflow and nf-core modules. Your task is to design a Nextflow workflow for DNA methylation analysis using bisulfite-converted sequencing (BS-seq) data*, assigns the role so that LLMs can respond accordingly.

**Controlling Output Verbosity:** Adjusting LLM verbosity tailors responses to the desired detail level. Users can specify response length for concise summaries or detailed explanations, aligning outputs with their needs. Controlling verbosity enhances response relevance and improves workflow design efficiency by matching outputs to the task’s context and purpose.

**User Intent:** Understanding the user’s goal and desired output is fundamental to formulating effective prompts. A clear awareness of the intent behind the interaction, whether it involves information retrieval, content generation, or problem-solving, ensures that the prompt aligns with the user’s expectations. Users can guide LLMs to produce accurate, relevant, and actionable responses by tailoring the prompt to address the specific objective. This focus on intent-driven prompt design is essential in workflow development, where precision and relevance directly impact workflows quality.

**Model Understanding:** Understanding the model’s knowledge and limitations is essential for designing prompts that maximize its strengths and mitigate weaknesses. Even advanced models like ChatGPT may struggle with certain tasks or produce errors. Awareness of these constraints allows users to compose prompts that align with the model’s capabilities. Additionally, incorporating strategies such as providing explicit instructions, clarifying context, or breaking complex tasks into smaller components can help address potential shortcomings. This awareness is vital in workflow development, where precision and domain-specific accuracy are crucial.

**Domain Specificity:** When working within specialized domains like scientific workflows, providing additional context and carefully chosen examples can significantly improve the robustness and relevance of model-generated responses. Context helps the model understand the specific requirements and constraints of the domain, while examples act as guides to shape the response. This approach ensures that the generated outputs align more closely with domain-specific expectations, enabling the development of precise and actionable scientific workflows.

**System-Aware Prompting:** Including system-specific terms (e.g., Galaxy or Nextflow) helps models tailor responses to the correct platform. For example, we ask *In the Galaxy platform, how are histories used to manage data?*. The inclusion of *Galaxy* helped models focus their responses specifically on GUI-based data lineage, rather than confusing it with script-based logging mechanisms in other SWSs.

**Iterative Testing and Refining:** Iterative testing and refinement are key to improving prompt responses. By analyzing outputs and adjusting prompts, users can fine-tune the model for greater accuracy and relevance. This process is especially valuable in workflow development as it enables the gradual optimization of prompts to achieve precise and actionable workflows.

**Stepwise Reasoning:** Asking models to break down the process logically before generation improves completeness and correctness. Example from Workflow Prompt (Galaxy W4): *First, list the required input files and tools. Then describe each step of the workflow from quality control to mapping. Finally, convert this plan into a Galaxy workflow.* This prompt led GPT-4o to reason about tool chaining (*assess*  $\rightarrow$  *clean*  $\rightarrow$  *re-assess*) before emitting the final workflow, resulting in better modularity and explanation.

**Balancing User Intent and LLMs Creativity:** Balancing user intent with LLM creativity is crucial for meaningful responses. While LLMs excel at innovation, prompts must guide outputs to align with user goals. This balance is vital in workflow development, where precision and problem-solving are essential for accurate results.

**Ensuring ethical usage and avoiding biases:** To ensure ethical use, addressing biases in LLM responses is vital. Clear guidelines, inclusive language, critical evaluation, and content filtering help mitigate issues and prevent harmful stereotypes. Maintaining ethical standards in workflow development ensures that outputs remain unbiased, inclusive, and suitable for diverse audiences, fostering trust and reliability in the model’s applications.

**Prompt Chaining and Multi-turn Conversations:** Prompt chaining connects multiple prompts sequentially, enabling dynamic, context-aware interactions with LLMs. Breaking complex queries into steps facilitates deeper exploration and refined discussions. In workflow development, this approach ensures step-by-step accuracy and cohesive integration of processes.

**Handling ambiguous or contradictory inputs:** LLMs may occasionally receive ambiguous or contradictory inputs. In such cases, prompt design with clarification is essential. Providing additional context or explicitly stating assumptions can further mitigate confusion. In workflow development, where precision is paramount, designing prompts to handle ambiguity ensures consistent and reliable outputs, even in complex or unclear scenarios.

We also employ several additional prompting techniques to enhance workflow accuracy and completeness. These include feedback loops for error evaluation, progressive prompting to break tasks into steps, backward prompting to verify prior outputs and multi-modal prompts. Together, these strategies refine workflow design and ensure robust results.

Following these strategies, we finalize the prompts for our study, as shown in Table 2, 3, 4, 5. This table includes prompts designed to gather background information and prior knowledge relevant to our selected workflows. The final workflow-specific prompts are presented in Tables 6 and 7, ensuring a clearer objective of each workflow while minimizing redundancy.

To standardize the prompt design process and facilitate reproducible evaluation, we identify generalizable prompt patterns for constructing effective workflow-generation queries in both Galaxy and Nextflow environments. These patterns emerge through a close examination of prompt structures used in our experiments, particularly those presented in Tables 6 and 7.

For Galaxy, which features a graphical user interface and emphasizes accessibility for non-programmers, prompts benefit from highlighting the biological objective, the input file types, and the desired output. A typical prompt for Galaxy follows this pattern: *Create a Galaxy workflow that performs [biological goal or analysis task] using [input file formats, such as VCF, GTF, or FASTQ]. The workflow should include steps for [intermediate processing, such as quality control, filtering, or mapping] and produce [expected output, such as variant calls or expression matrices].* Prompts constructed using this pattern ensure that the model captures both the domain intent and Galaxy-specific processing logic. For instance, a clear example from our study is:

*Create a Galaxy workflow that identifies the exon with the highest number of SNPs on human chromosome 22 using VCF and GTF files,* which allowed LLMs to generate workflows with tool-specific operations aligned to GTN tutorials.

In contrast, Nextflow is a script-based system designed for high reproducibility and modular pipeline development. Prompts targeting Nextflow must emphasize the computational modules, input-output relationships, and optional configuration logic. An effective prompt pattern in this case is: *Create a Nextflow workflow that performs [bioinformatics task] using [tools, such as fastp, STAR, featureCounts]. The workflow should take [data types] as input and output [result types], including steps for [processes like alignment, quantification, or filtering].* An example that reflects this pattern is: *Create a Nextflow workflow to process raw FASTQ files using fastp, STAR, and featureCounts, and perform differential expression analysis,* which guides the LLM to structure the script according to Nextflow syntax and pipeline conventions. These base patterns can be enhanced through prompting strategies such as role-based framing (e.g., *You are a bioinformatics workflow developer using Nextflow*), stepwise reasoning (*First list the inputs and tools, then describe each step before generating the workflow*), and ambiguity handling (*Assume the input data is raw and unfiltered*). Such augmentations help LLMs disambiguate task requirements, infer missing steps, and align their responses more closely with real-world workflows.

By articulating these prompt patterns, we provide a foundation for both practitioners and researchers to systematically guide LLMs in generating domain-relevant workflows across different scientific workflow systems.

## 5 Discussion

This study systematically evaluates how state-of-the-art LLMs perform in the development of scientific workflows, with a particular focus on bioinformatics scenarios using Galaxy and Nextflow. Through a task-oriented analysis across three research questions, we examine the capabilities and limitations of LLMs in (i) addressing fundamental and background knowledge of scientific workflow systems, (ii) generating complete and executable workflows, and (iii) responding to different prompting strategies. Our results collectively offer a grounded understanding of the practical value and current constraints of LLMs in computational research workflows.

One of the most compelling findings is that LLMs, particularly Gemini 2.5 Flash and GPT-4o, demonstrate promising performance in addressing fundamental questions about workflow systems. These models can provide syntactically sound and semantically relevant responses to common background queries on workflow components and domain-specific tasks. However, performance variability exists across tools, with some models omitting contextual steps or offering generalized responses, especially when minimal prompting is used. This observation underscores the need for careful prompt design even for tasks that may appear routine to human experts.

When evaluating LLM-generated workflows (RQ2), we uncover significant differences in quality, particularly in the dimensions of completeness, correctness, and usability. Gemini 2.5 Flash consistently produces high-quality workflows that are logically structured, well-annotated, and accessible for novice users. DeepSeek-V3 also performs reasonably well, especially in generating technically detailed responses including command-line and directory structures. However, GPT-4o exhibits inconsistencies when tasked with minimal instruction, often omitting critical setup steps or default configurations. These results suggest that while LLMs can be helpful in automating parts of the workflow development process, their outputs must be critically evaluated, especially in high-stakes scientific analyses.

Our investigation into prompting strategies (RQ3) further illuminates the dependencies between user input structure and model performance. The tiered evaluation, using instruction-only, role-based, and chain-of-thought prompts, revealed that prompt formulation plays a pivotal role in shaping output quality. For simpler tasks, most models respond adequately to concise instructions. However, as the complexity of the task increased, instruction-only prompts are often insufficient. GPT-4o, in particular, demonstrates the most improvement when guided by role-based or step-wise chain-of-thought prompts. In contrast, Gemini generally maintains robust performance even with minimal instruction, suggesting that some models are better tuned for general usability, while others require more deliberate contextual scaffolding to achieve optimal outcomes.

Importantly, the study highlights a practical implication: prompt engineering is not merely an academic curiosity but a crucial skill for practitioners leveraging LLMs in scientific workflow development. While the models exhibit impressive generative capacity, their effectiveness is shaped by the



user’s ability to structure requests, assess outputs, and iteratively refine instructions. For research groups aiming to integrate LLMs into their scientific pipeline design process, this emphasizes the need to develop domain-specific prompting strategies and validation routines to ensure robustness and reproducibility.

Nevertheless, the study also reveals limitations. LLMs still occasionally produce hallucinated tools or incompatible parameter settings, and they rarely reflect the full diversity of domain-specific conventions. Furthermore, workflow narratives are often limited in their ability to explain the rationale behind tool selection or parameterization, an essential aspect of reproducibility and scientific communication. These issues point to an important area for future work: the integration of LLMs with domain-aware validation engines and curated knowledge bases to improve fidelity and interpretability.

In summary, while LLMs exhibit growing competence in the generation and explanation of scientific workflows, they are not yet fully autonomous solutions. Their utility depends critically on model choice, prompt design, and human oversight. By understanding and addressing these dependencies, workflow developers can better harness the capabilities of LLMs for accelerating computational research while maintaining scientific rigor and transparency.

## 6 Threats to Validity

Despite the rigorous evaluation conducted in this study, we acknowledge several potential threats to validity to contextualize the scope and applicability of our findings.

**Internal Validity:** refers to the extent to which a causal relationship can be established between the applied treatment and the observed outcomes [89]. A central threat to the internal validity of our investigation lies in the subjectivity of manual workflow evaluation. While the assessment of correctness, completeness, and usability was conducted by two domain experts with significant experience in scientific workflows, human interpretation is inherently susceptible to bias. To mitigate this, we relied on well-defined rubrics and cross-validation against established baselines from the Galaxy Training Network (GTN) and nf-core repositories. Nevertheless, the absence of quantitative metrics for some usability aspects introduces potential inconsistencies in evaluation outcomes. Additionally, while we adopted a stepwise prompting approach, progressing from instruction-only to role-based and chain-of-thought prompts, our criteria for escalation may still carry implicit judgment calls. In rare cases, the decision not to escalate a prompt may overlook latent model capabilities that might emerge under more nuanced scenarios.

**External Validity:** refers to the extent to which the research findings can be generalized beyond the specific context of the study to broader populations, settings, or situations [89]. The generalizability of our findings may be limited by the specific selection of LLMs and workflows. Our study evaluated three LLMs *GPT-4o*, *Gemini 2.5 Flash*, and *DeepSeek-V3*, all chosen based on their availability and relevance as of early 2025. As the LLM landscape evolves rapidly, newer models may outperform or behave differently under the same experimental settings. Furthermore, although the selected workflows span a wide range of bioinformatics tasks and platforms (Galaxy and Nextflow), they may not capture the full diversity of scientific workflows used in other domains or emerging subfields within bioinformatics. Our workflow tasks also emphasized well-documented use cases from GTN and nf-core. While this ensures reproducibility and standardization, it may not reflect more exploratory or unconventional workflows encountered in cutting-edge research, where documentation is sparse or evolving.

**Construct Validity:** concerns the extent to which the measurements and observations accurately reflect the theoretical concepts they are intended to represent [89]. While *completeness*, *correctness*, and *usability* are standard in workflow assessment, their practical interpretation may vary across evaluators and contexts. For example, a workflow deemed usable for an expert might not be equally accessible to a novice. Although we attempted to account for varying user expertise levels by noting model behavior across different prompting strategies, our construct definitions could still be refined further in future studies using more granular user studies or empirical execution logs. Furthermore, the assumption that community-curated workflows represent ground truth may overlook acceptable alternative implementations. Some LLM-generated workflows deviated from reference workflows in structure or tool usage but remained valid and executable. While we considered such variation acceptable when logically justified, this introduces a gray area in evaluation that may affect reproducibility across studies.

**Conclusion Validity:** refers to the degree to which the conclusions drawn from a study are based on appropriate and reliable analysis of the data. Our conclusions regarding model effectiveness

and prompting strategy are based on a relatively small sample size of workflows and limited prompting techniques. While this setup was chosen to balance depth and manageability, the limited scope may affect the statistical power of our inferences. Additionally, since all models were evaluated using their default configurations, we did not explore the impact of parameter tuning or prompt fine-tuning, which might significantly influence output quality. By recognizing these limitations, we aim to foster transparency and encourage further research to replicate, validate, and extend our findings across broader contexts, user bases, and scientific domains.

By recognizing these limitations, we aim to foster transparency and encourage further research to replicate, validate, and extend our findings across broader contexts, user bases, and scientific domains.

## 7 Related Work

Scientific workflow development plays a central role in modern bioinformatics, supporting complex analyses that involve multiple tools/modules, datasets, and computing environments. However, constructing these workflows remains a significant challenge, particularly for researchers without extensive programming expertise. To simplify workflow construction and scientific analysis, researchers have developed several approaches. For instance, Palmblad et al. [90] used EDAM and semantic tool annotations to enable automated workflow generation in mass spectrometry-based proteomics. Similarly, Di Bernardo et al. [91] utilized data types to generate workflows automatically. Koop et al. introduced *VisComplete*, a system designed to help users build visualization workflows based on prior workflows in the VisTrails framework [92]. However, this approach did not verify the correctness of the reused workflows. These early methods, while innovative, were limited in scope, applicable to a narrow range of bioinformatics analyses, and prone to errors as tools evolved. Addressing these challenges, Kumar et al. [16] developed a deep learning-based model for recommending tools during workflow construction by analyzing workflows hosted on the European Galaxy server. However, their approach overlooked tool compatibility, failed to account for outdated tools and was restricted to workflows from the European Galaxy server. Building on this, Kumar et al. [93] developed a transformer-based tool recommendation system, but it lacked generalizability beyond the European Galaxy platform. Recent advances in LLMs have opened new opportunities for simplifying workflow development through natural language interaction and automated code generation. Several studies have explored the potential of LLMs to support scientific workflow development, highlighting both their potential and limitations. Sanger et al. [3] conducted a qualitative assessment of ChatGPT’s capabilities in assisting with workflow design across general scientific domains. While their study demonstrated ChatGPT’s ability to structure high-level workflows and reduce development complexity, it was limited to a single model and did not address domain-specific accuracy or tool-level fidelity. In contrast, our work conducts a systematic, domain-specific evaluation of three advanced LLMs, GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3, focusing on bioinformatics workflows across two widely adopted scientific workflow systems, Galaxy and Nextflow. Beyond workflow generation, we explicitly assess the role of prompting strategies (instruction-only, role-based, and chain-of-thought), offering a deeper analysis of LLM behavior in structured scientific contexts. Our evaluation is grounded in community-validated baselines (e.g., GTN and nf-core), enabling reproducible comparisons of workflow quality in terms of completeness, correctness, and usability. Pickard et al. [94] proposed *BRAD*, a retrieval-augmented agent that integrates LLMs with external databases to support tasks such as question answering and gene enrichment analysis. While BRAD focuses on multi-modal task execution through an agentic interface, our study centers on the generation of end-to-end executable workflows. Similarly, Riffle et al. [95] introduced *OLAF*, a conversational platform that enables users to execute single-cell analyses via modular LLM-driven interaction. Unlike these agent-based systems, our approach provides a model-agnostic, benchmark-driven comparison of LLMs, emphasizing prompt effectiveness and workflow synthesis across platforms rather than execution or interaction design. Tang et al. [96] presented *BioCoder*, a benchmark suite for evaluating LLM performance in generating bioinformatics code snippets in Python and Java. Their work focuses on isolated function-level code generation and uses fuzz-testing to measure correctness. While informative for programming capabilities, BioCoder does not consider workflow-level composition or usability in scientific contexts. Our study moves beyond code synthesis by generating and validating complete workflows that mirror real-world bioinformatics pipelines, thereby addressing aspects of tool chaining, data flow, and reproducibility.

Other notable efforts include the Playbook Workflow Builder (PWB) [97], which offers a GUI-based interface for assembling workflows through semantically annotated building blocks, and

PROTEUS [98], which applies LLMs for hierarchical planning and hypothesis generation in proteomics. These systems support accessible workflow construction or autonomous task planning, yet they do not evaluate the quality of LLM-generated workflows against executable community standards, nor do they systematically assess the effect of prompt design. Despite these advancements, the integration of LLMs into practical, domain-specific scientific workflow development remains under-explored. Prior work tends to focus either on high-level planning or isolated code generation, lacking a comprehensive analysis that connects prompt design, model behavior, and platform-specific execution. Our study fills this gap by offering a unified evaluation of LLMs for bioinformatics workflow generation, grounded in real platforms, community standards, and empirical metrics.

## 8 Conclusion

This study presents a systematic evaluation of LLMs for generating scientific workflows in bioinformatics, focusing on two widely adopted workflow systems, Galaxy and Nextflow. By analyzing the performance of GPT-4o, Gemini 2.5 Flash, and DeepSeek-V3 across a curated set of bioinformatics tasks, we examined the models’ capabilities in addressing fundamental workflow concepts, producing end-to-end executable pipelines, and responding to various prompting strategies. Our results highlight that LLMs can serve as effective co-developers for scientific workflows, especially when guided by well-structured prompts. Gemini 2.5 Flash consistently performed well with minimal prompting, offering clear, complete, and intuitive outputs, making it particularly suitable for novice users. GPT-4o, while occasionally omitting details in instruction-only prompts, excelled when provided with role-based context, indicating that prompt engineering plays a critical role in leveraging its full potential. DeepSeek-V3 also demonstrated strong technical depth, especially in producing scripts with directory structures and command-line instructions, although its verbosity and inconsistency may hinder usability for some users. Importantly, our study underscores the need for deliberate prompting strategies. Simple prompts often yield incomplete or ambiguous outputs, whereas thoughtfully layered approaches, incorporating context, roles, and logical decomposition, significantly improve workflow generation. These findings are relevant to developers using LLMs for workflow authoring and designing future LLM-integrated tools for scientific research.

In summary, while LLMs have not yet achieved full autonomy in scientific workflow construction, they already offer valuable support for both novice and expert users. Their integration into scientific toolchains can democratize workflow development, reduce entry barriers, and accelerate research reproducibility, provided their use is guided by informed prompting and critical validation. Future work should explore scalable evaluation frameworks, broader domain coverage, and user-centered studies to further solidify the role of LLMs in scientific workflows.

**Supplementary information.** All supplementary files can be obtained using [88]

**Acknowledgements.** This research is supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC), and by the industry-stream NSERC CREATE in Software Analytics Research (SOAR).

## Declarations

- Funding: This research is supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC), and by the industry-stream NSERC CREATE in Software Analytics Research (SOAR)
- Conflict of interest/Competing interests: We have no conflict of interest.
- Ethics approval and consent to participate: Not Applicable
- Consent for publication: Not Applicable
- Data availability: Our data can be obtained using [88]
- Materials availability: Our data can be obtained using [88]
- Code availability Not Applicable
- Author contribution: Khairul Alam designed and conducted the study, developed the evaluation framework, constructed prompts, carried out the analysis, and evaluated the LLM-generated workflows. He also led the drafting and writing of the manuscript. Banani Roy contributed to the evaluation of the workflows, provided critical insights on prompt engineering and evaluation criteria, reviewed and edited the manuscript, and offered guidance throughout the study.

## References

- [1] Liu, J., Pacitti, E., Valduriez, P., Mattoso, M.: A survey of data-intensive scientific workflow management. *Journal of Grid Computing* **13**(4), 457–493 (2015)
- [2] Cohen-Boulakia, S., Belhajjame, K., *et al.*: Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities. *FGCS* **75**, 284–298 (2017)
- [3] Sanger, M., De Mecquenem, N., *et al.*: A qualitative assessment of using chatgpt as large language model for scientific workflow development. *GigaScience* **13**, 030 (2024)
- [4] Wratten, L., Wilm, A., Goke, J.: Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers. *Nature methods* **18**(10), 1161–1168 (2021)
- [5] Goodwin, S., McPherson, J.D., McCombie, W.R.: Coming of age: ten years of next-generation sequencing technologies. *Nature reviews genetics* **17**(6), 333–351 (2016)
- [6] Xu, J., Du, W., *et al.*: Llm4workflow: An llm-based automated workflow model generation tool. In: 2024 39th IEEE/ACM International Conference on ASE, pp. 2394–2398 (2024). IEEE
- [7] Langer, B.E., Amaral, A., *et al.*: Empowering bioinformatics communities with nextflow and nf-core. *bioRxiv*, 2024–05 (2024)
- [8] Goecks, J., Nekrutenko, A., *et al.*: Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome biology* **11**, 1–13 (2010)
- [9] Di Tommaso, P., Chatzou, M., *et al.*: Nextflow enables reproducible computational workflows. *Nature biotechnology* **35**(4), 316–319 (2017)
- [10] Koster, J., Rahmann, S.: Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* **28**(19), 2520–2522 (2012)
- [11] Harenslak, P., *et al.*: Data Pipelines with Apache Airflow. Simon and Schuster, ??? (2021)
- [12] Deelman, E., Vahi, K., *et al.*: Pegasus, a workflow management system for science automation. *FGCS* **46**, 17–35 (2015)
- [13] Goodstadt, L.: Ruffus: a lightweight python library for computational pipelines. *Bioinformatics* **26**(21), 2778–2779 (2010)
- [14] Khodak, A., Chapman, B., *et al.*: Existing-Workflow-systems. [Online; accessed Thursday May 29 2025] (2024). <https://github.com/common-workflow-language/common-workflow-language/wiki/Existing-Workflow-systems>
- [15] Galaxy: Galaxy Tool Shed. [Online; accessed Thursday May 29 2025] (2025). <https://toolshed.g2.bx.psu.edu/>
- [16] Kumar, A., Rasche, H., Gruning, B., *et al.*: Tool recommender system in galaxy using deep learning. *GigaScience* **10**(1), 152 (2021)
- [17] Alam, K., Roy, B., Serebrenik, A.: Reusability challenges of scientific workflows: A case study for galaxy. In: 2023 30th APSEC, pp. 289–298 (2023). IEEE
- [18] Silva, R., Bard, D., Chard, K., Foster, I., Gibbs, T., Goble, C., Godoy, W., Gustafsson, J., Hudson, S., Jha, S., *et al.*: Workflows community summit 2024: future trends and challenges in scientific workflows (2024)
- [19] Menaka, M., Kumar, K.S.: Workflow scheduling in cloud environment—challenges, tools, limitations & methodologies: A review. *Measurement: Sensors* **24**, 100436 (2022)

- [20] Deelman, E., Peterka, T., *et al.*: The future of scientific workflows. *The International Journal of HPC Applications* **32**(1), 159–175 (2018)
- [21] De Roure, D., Goble, C.: *myexperiment—a web 2.0 virtual research environment* (2007)
- [22] Cohen-Boulakia, S., Leser, U.: Search, adapt, and reuse: the future of scientific workflows. *ACM SIGMOD Record* **40**(2), 6–16 (2011)
- [23] Hurst, A., Lerer, A., Goucher, A.P., Perelman, A., Ramesh, A., Clark, A., Ostrow, A., Welihinda, A., Hayes, A., Radford, A., *et al.*: Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024)
- [24] Gemini: Google Deepmind Gemini 2.5. [Online; accessed Thursday May 29 2025] (2024). <https://deepmind.google/models/gemini/>
- [25] Liu, A., Feng, B., Xue, B., Wang, *et al.*: Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024)
- [26] Ouyang, S., Zhang, J.M., Harman, M., Wang, M.: Llm is like a box of chocolates: the non-determinism of chatgpt in code generation. *arXiv preprint arXiv:2308.02828* (2023)
- [27] Hou, X., Zhao, Y., *et al.*: Large language models for software engineering: A systematic literature review. *ACM TOSEM* (2023)
- [28] Ekin, S.: Prompt engineering for chatgpt: a quick guide to techniques, tips, and best practices. *Authorea Preprints* (2023)
- [29] Arora, S., Narayan, A., *et al.*: Ask me anything: A simple strategy for prompting language models. In: *The Eleventh International Conference on Learning Representations* (2022)
- [30] Reynolds, L., McDonell, K.: Prompt programming for large language models: Beyond the few-shot paradigm. In: *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pp. 1–7 (2021)
- [31] Wei, J., Wang, X., *et al.*: Chain-of-thought prompting elicits reasoning in large language models. *Advances in NIPS* **35**, 24824–24837 (2022)
- [32] Zhou, Y., Muresanu, A.I., *et al.*: Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910* (2022)
- [33] Marvin, G., Hellen, N., *et al.*: Prompt engineering in large language models. In: *International Conference on Data Intelligence and Cognitive Informatics*, pp. 387–402 (2023). Springer
- [34] Lin, Z.: How to write effective prompts for large language models. *Nature Human Behaviour* **8**(4), 611–615 (2024)
- [35] Hiltmann, S., *et al.*: Galaxy training: A powerful framework for teaching! *PLoS computational biology* **19**(1), 1010752 (2023)
- [36] Ewels, P.A., *et al.*: The nf-core framework for community-curated bioinformatics pipelines. *Nature biotechnology* **38**(3), 276–278 (2020)
- [37] Barker, A., Van Hemert, J.: Scientific workflow: a survey and research directions. In: *International Conference on Parallel Processing and Applied Mathematics*, pp. 746–753 (2007). Springer
- [38] Ludäscher, B., Weske, M., McPhillips, T., Bowers, S.: Scientific workflows: Business as usual? In: *International Conference on Business Process Management*, pp. 31–47 (2009). Springer
- [39] Lin, C., Lu, S., *et al.*: A reference architecture for scientific workflow management systems and the view soa solution. *IEEE TSC* **2**(1), 79–92 (2009)



- [40] Oinn, T., Addis, M., *et al.*: Taverna: a tool for the composition and enactment of bioinformatics workflows. *Bioinformatics* **20**(17), 3045–3054 (2004)
- [41] Liew, C.S., Atkinson, M.P., *et al.*: Scientific workflows: moving across paradigms. *ACM CSUR* **49**(4), 1–39 (2016)
- [42] Vaswani, A.: Attention is all you need. *Advances in Neural Information Processing Systems* (2017)
- [43] Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., *et al.*: A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023)
- [44] Liu, J., Xia, C.S., *et al.*: Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in NIPS* **36** (2024)
- [45] Yu, J., Liang, P., *et al.*: Security code review by llms: A deep dive into responses. *arXiv preprint arXiv:2401.16310* (2024)
- [46] Nam, D., Macvean, A., Hellendoorn, V., *et al.*: Using an llm to help with code understanding. In: *Proceedings of the IEEE/ACM 46th ICSE*, pp. 1–13 (2024)
- [47] Su, Y., Wan, C., *et al.*: Hotgpt: How to make software documentation more useful with a large language model? In: *Proceedings of the 19th Workshop on Hot Topics in OS*, pp. 87–93 (2023)
- [48] Pan, R., Ibrahimzada, A.R., *et al.*: Understanding the effectiveness of large language models in code translation. *arXiv preprint arXiv:2308.03109* (2023)
- [49] Yang, Z., Liu, F., *et al.*: Exploring and unleashing the power of large language models in automated code translation. *Proceedings of the ACM on Software Engineering* **1**(FSE), 1585–1608 (2024)
- [50] Manakul, P., Liusie, A., Gales, M.J.: Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896* (2023)
- [51] Peng, B., *et al.*: Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813* (2023)
- [52] Ugare, S., *et al.*: Improving llm code generation with grammar augmentation. *arXiv e-prints*, 2403 (2024)
- [53] Gu, Q.: Llm-based code generation method for golang compiler testing. In: *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the FSE*, pp. 2201–2203 (2023)
- [54] Le, H., Wang, Y., *et al.*: Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in NIPS* **35**, 21314–21328 (2022)
- [55] Jiang, E., Toh, E., *et al.*: Discovering the syntax and strategies of natural language programming with generative language models. In: *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pp. 1–19 (2022)
- [56] Vaithilingam, P., *et al.*: Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In: *Chi Conference on Human Factors in Computing Systems Extended Abstracts*, pp. 1–7 (2022)
- [57] Kazemitabaar, M., Chow, J., *et al.*: Studying the effect of ai code generators on supporting novice learners in introductory programming. In: *Proceedings of the 2023 CHI on Human Factors in Computing Systems*, pp. 1–23 (2023)
- [58] Yetiştirilen, B., *et al.*: Evaluating the code quality of ai-assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and chatgpt. *arXiv preprint*

arXiv:2304.10778 (2023)

- [59] Liu, M.X., Sarkar, A., *et al.*: “what it wants me to say”: Bridging the abstraction gap between end-user programmers and code-generating large language models. In: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, pp. 1–31 (2023)
- [60] Oppenlaender, J.: A taxonomy of prompt modifiers for text-to-image generation. Behaviour & Information Technology, 1–14 (2023)
- [61] White, J., *et al.*: A prompt pattern catalog to enhance prompt engineering with chatgpt. arXiv preprint arXiv:2302.11382 (2023)
- [62] Liu, Y., Deng, G., Xu, Z., *et al.*: Jailbreaking chatgpt via prompt engineering: An empirical study. arXiv preprint arXiv:2305.13860 (2023)
- [63] Wang, L., Chen, X., Deng, X., *et al.*: Prompt engineering in consistency and reliability with the evidence-based guideline for llms. NPJ digital medicine **7**(1), 41 (2024)
- [64] Sahoo, P., Singh, A.K., *et al.*: A systematic survey of prompt engineering in large language models: Techniques and applications. arXiv preprint arXiv:2402.07927 (2024)
- [65] Chen, B., Zhang, Z., Langrené, N., Zhu, S.: Unleashing the potential of prompt engineering in large language models: a comprehensive review. arXiv preprint arXiv:2310.14735 (2023)
- [66] Parameswaran, A.G., Shankar, S., Asawa, P., Jain, N., Wang, Y.: Revisiting prompt engineering via declarative crowdsourcing. arXiv preprint arXiv:2308.03854 (2023)
- [67] Vatsal, S., Dubey, H.: A survey of prompt engineering methods in large language models for different nlp tasks. arXiv preprint arXiv:2407.12994 (2024)
- [68] Shah, C.: From prompt engineering to prompt science with humans in the loop. Communications of the ACM **68**(6), 54–61 (2025)
- [69] Smith, R., Fries, J.A., Hancock, B., Bach, S.H.: Language models in the loop: Incorporating prompting into weak supervision. ACM/JMS Journal of Data Science **1**(2), 1–30 (2024)
- [70] Cohn, C., Hutchins, N., Biswas, G.: Towards a formative feedback generation agent: Leveraging a human-in-the-loop, chain-of-thought prompting approach with llms to evaluate formative assessment responses in k-12 science. (2023)
- [71] Razdaibiedina, A., Mao, Y., Hou, R., Khabsa, M., Lewis, M., Almahairi, A.: Progressive prompts: Continual learning for language models. arXiv preprint arXiv:2301.12314 (2023)
- [72] Knoth, N., Tolzin, A., Janson, A., Leimeister, J.M.: Ai literacy and its implications for prompt engineering strategies. Computers and Education: Artificial Intelligence **6**, 100225 (2024)
- [73] Lo, L.S.: The art and science of prompt engineering: a new literacy in the information age. Internet Reference Services Quarterly **27**(4), 203–210 (2023)
- [74] Lo, L.S.: The clear path: A framework for enhancing information literacy through prompt engineering. The Journal of Academic Librarianship **49**(4), 102720 (2023)
- [75] Hiltemann, S., *et al.*: Galaxy basics for genomics (2024)
- [76] Pajon, A., Blank, C., *et al.*: From peaks to genes (2024)
- [77] Clements, D., Gallardo, C.: Introduction to genomics and galaxy (2024)
- [78] Batut, B., Doyle, M., *et al.*: Quality Control (Galaxy Training Materials). [Online; accessed 2025-06-25] (2025). <https://training.galaxyproject.org/training-material/topics/sequence-analysis/tutorials/quality-control/tutorial.html>

- [79] Batut, B.: Bacterial genome annotation (2024)
- [80] Hakkaart, C., Hörtenhuber, M., bot: Nf-core/demo: Nf-core/demo 1.0.1 - Roasted Sweet Potato. <https://doi.org/10.5281/zenodo.13951181> . <https://doi.org/10.5281/zenodo.13951181>
- [81] Patel, H., Garcia, M.U., et al.: Nf-core/fetchngs: Nf-core/fetchngs V1.12.0 - Titanium Platypus. <https://doi.org/10.5281/zenodo.10728509> . <https://doi.org/10.5281/zenodo.10728509>
- [82] Ewels, P., Sateesh\_Peri, Hüther, P., et al.: Nf-core/methylseq: Endless Tofu. <https://doi.org/10.5281/zenodo.14502249> . <https://doi.org/10.5281/zenodo.14502249>
- [83] Patel, H., Ewels, P., Manning, J., et al.: Nf-core/rnaseq: Nf-core/rnaseq V3.18.0 - Lithium Lynx. <https://doi.org/10.5281/zenodo.14537300> . <https://doi.org/10.5281/zenodo.14537300>
- [84] Lundin, D., bot, Yates, J.A.F.: Nf-core/phyloplace: Phylosearch Release. <https://doi.org/10.5281/zenodo.14906186> . <https://doi.org/10.5281/zenodo.14906186>
- [85] Hiltmann, S., Rasche, H., Gladman, S., and: Galaxy training: A powerful framework for teaching! PLoS Comput Biol **19**(1), 1010752 (2023) <https://doi.org/10.1371/journal.pcbi.1010752>
- [86] Chang, K., Xu, S., Wang, C., Luo, Y., Xiao, T., Zhu, J.: Efficient prompting methods for large language models: A survey. arXiv preprint arXiv:2404.01077 (2024)
- [87] Liu, P., *et al.*: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM Computing Surveys **55**(9), 1–35 (2023)
- [88] Alam, K.: From Prompt to Pipeline: Large Language Models for Scientific Workflow Development in Bioinformatics. <https://doi.org/10.5281/zenodo.16416384> . <https://doi.org/10.5281/zenodo.16416384>
- [89] Wohlin, C., Runeson, P., Höst, M., Ohlsson, M.C., Regnell, B., Wesslén, A.: Experimentation in software engineering. Springer Science & Business Media (2012)
- [90] Palmblad, M., *et al.*: Automated workflow composition in mass spectrometry-based proteomics. Bioinformatics **35**(4), 656–664 (2019)
- [91] DiBernardo, M., *et al.*: Semi-automatic web service composition for the life sciences using the biomoby semantic web framework. Journal of biomedical informatics **41**(5), 837–847 (2008)
- [92] Koop, D., *et al.*: Viscomplete: Automating suggestions for visualization pipelines. IEEE Transactions on Visualization and Computer Graphics **14**(6), 1691–1698 (2008)
- [93] Kumar, A., Grüning, B., *et al.*: Transformer-based tool recommendation system in galaxy. BMC bioinformatics **24**(1), 446 (2023)
- [94] Pickard, J., Prakash, R., et al.: Language model powered digital biology with brad. arXiv preprint arXiv:2409.02864 (2024)
- [95] Riffle, D., Shirooni, N., He, C., Murali, M., Nayak, S., Gopalan, R., Lopez, D.G.: Olaf: An open life science analysis framework for conversational bioinformatics powered by large language models. arXiv preprint arXiv:2504.03976 (2025)
- [96] Tang, X., Qian, B., Gao, R., *et al.*: Biocoder: a benchmark for bioinformatics code generation with large language models. Bioinformatics **40**(Supplement\_1), 266–276 (2024)
- [97] Clarke, D.J., Evangelista, J.E., *et al.*: Playbook workflow builder: Interactive construction of bioinformatics workflows. PLOS Computational Biology **21**(4), 1012901 (2025)
- [98] Ding, N., Qu, S., et al.: Automating exploratory proteomics research via language models. arXiv preprint arXiv:2411.03743 (2024)