

# Automated Bioinformatics Analysis via AutoBA

Juexiao Zhou<sup>1,2,†</sup>, Bin Zhang<sup>1,2,†</sup>, Xiuying Chen<sup>1,2</sup>, Haoyang Li<sup>1,2</sup>, Xiaopeng Xu<sup>1,2</sup>, Siyuan Chen<sup>1,2</sup>,  
Xin Gao<sup>1,2,\*</sup>

**Abstract**—Here we introduce Automated Bioinformatics Analysis (AutoBA), the first autonomous AI agent designed explicitly for regular omics data analysis based on large language models. AutoBA simplifies the analytical process by requiring minimal user input while delivering detailed step-by-step plans for various bioinformatics tasks. Through rigorous validation by expert bioinformaticians, AutoBA’s robustness and adaptability are affirmed across a diverse range of omics analysis cases, including whole genome/exome sequencing (WGS/WES), chromatin immunoprecipitation assays with sequencing (ChIP-seq), RNA sequencing (RNA-seq), single-cell RNA-seq, spatial transcriptomics and so on. AutoBA’s unique capacity to self-design analysis processes based on input data variations further underscores its versatility. Compared with online bioinformatic services, AutoBA offers five LLM backends, with options for both online and local usage, prioritizing data security and user privacy. Moreover, different from the predefined pipeline, AutoBA has adaptability in sync with emerging bioinformatics tools. Overall, AutoBA represents an advanced and convenient tool, offering robustness and adaptability for conventional bioinformatics analysis.

**Index Terms**—Bioinformatics, Omics analysis, Large language model, Agent.



## 1 INTRODUCTION

BIOINFORMATICS is an interdisciplinary field that encompasses computational, statistical, and biological approaches to analyze, understand and interpret complex biological data [1], [2], [3]. With the rapid growth of gigabyte-sized biological data generated from various high-throughput technologies, bioinformatics has become an essential tool for researchers to make sense of these massive datasets and extract meaningful biological insights. The applications of bioinformatics typically cover diverse fields such as genome analysis [4], [5], [6], structural bioinformatics [7], [8], [9], systems biology [10], data and text mining [11], [12], [13], phylogenetics [14], [15], [16], and population analysis [17], [18], which has further enabled significant advances in personalized medicine [19] and drug discovery

[5].

In broad terms, bioinformatics could be categorized into two primary domains: the development of innovative algorithms to address various biological challenges [20], [21], [22], [23], [24], and the application of established tools to analyze extensive biological datasets [25], [26], especially high-throughput sequencing data. Developing new bioinformatics software requires a substantial grasp of biology and programming expertise. Alongside the development of novel computational methods, one of the most prevalent applications of bioinformatics is the investigation of biological data using the existing tools and pipelines [27], [28], which typically involves a sequential, flow-based analysis of omics data, encompassing variety types of datasets like whole genome sequencing (WGS) [29], whole exome sequencing (WES), RNA sequencing (RNA-seq) [30], single-cell RNA-seq (scRNA-Seq) [31], transposase-accessible chromatin with sequencing (ATAC-Seq) [32], ChIP-seq [33], and spatial transcriptomics [34].

For example, the conventional analytical framework for bulk RNA-seq involves a meticulously structured sequence

<sup>1</sup>Computer Science Program, Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Kingdom of Saudi Arabia

<sup>2</sup>Computational Bioscience Research Center, King Abdullah University of Science and Technology, Thuwal 23955-6900, Kingdom of Saudi Arabia

<sup>†</sup>These authors contributed equally to this work.

\*To whom correspondence should be addressed; E-mail: xin.gao@kaust.edu.sa.

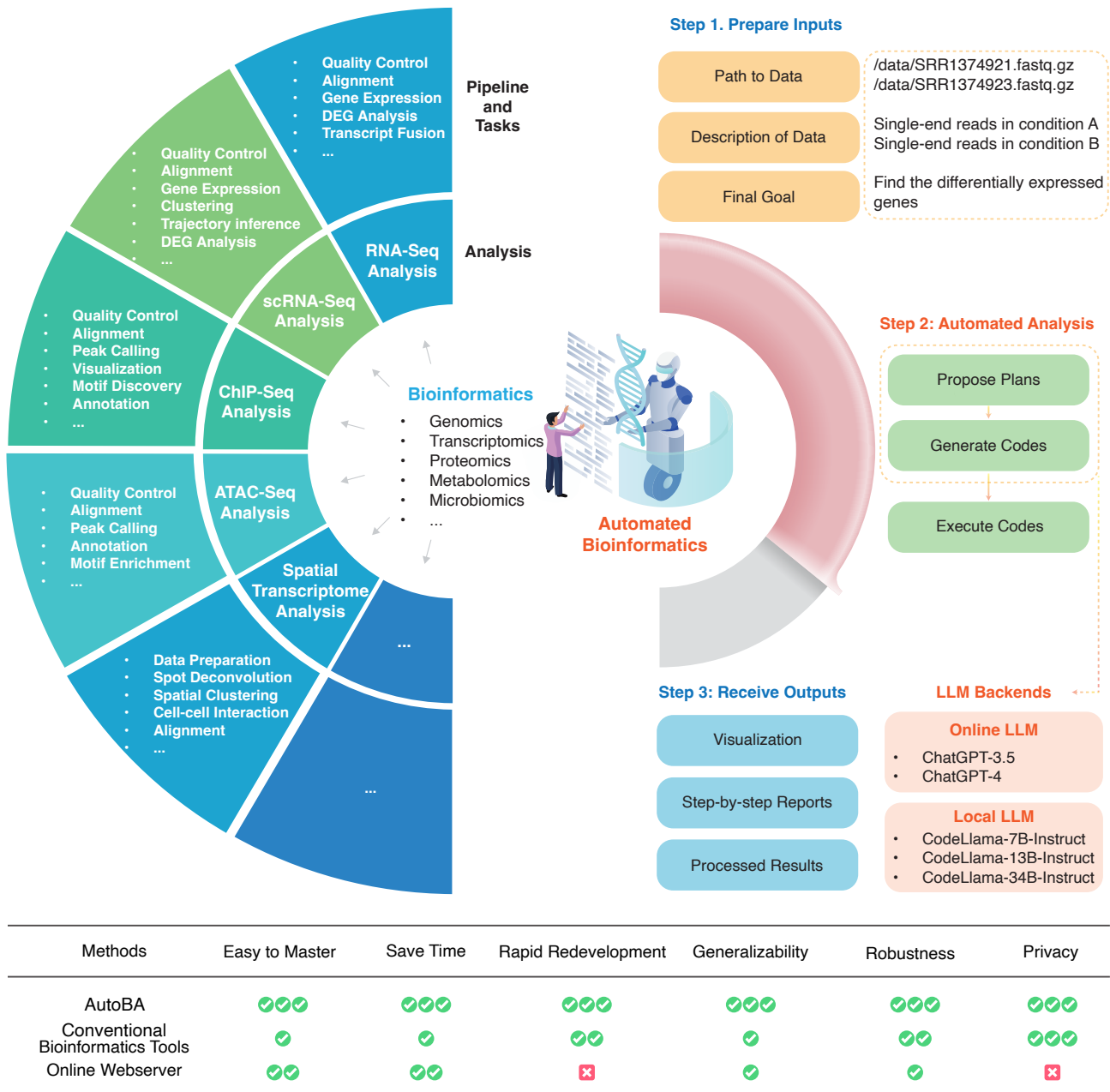


Fig. 1. **Design of AutoBA.** AutoBA stands as the first autonomous AI agent meticulously crafted for conventional bioinformatics analysis. Remarkably user-friendly, AutoBA simplifies the analytical process by requiring minimal user input, including data path, data description, and the final objective, while delivering detailed step-by-step plans for various bioinformatics tasks. With these inputs, it autonomously proposes analysis plans, generates code, executes codes, and conducts subsequent data analysis by using our well-designed prompts. AutoBA was implemented as open-source software that offers five LLM backends, with options for both online and local deployment, prioritizing data security and user privacy and offering a streamlined and efficient solution for bioinformatics tasks. Step 1 and Step 3 require human involvement, while Step 2 requires no human involvement.

of computational steps [35]. This intricate pipeline reveals its complexity through a series of carefully orchestrated stages. It begins with quality control [36], progresses to tasks such as adapter trimming [37] and the removal of low-quality reads, and then moves on to critical steps like genome or

transcriptome alignment [38]. Furthermore, it extends to some advanced tasks, including the identification of splice junctions [39], quantification through read counting [40], and the rigorous examination of differential gene expression [41]. Moreover, the pipeline delves into the intricate domain

TABLE 1

**Summary of AutoBA application scenarios in bioinformatics multi-omics analysis.** The table displays a comprehensive list of real-world cases utilized to assess AutoBA, providing information on the class of the cases, the respective task name, and the corresponding case ID.

Bioinformatics Pipelines	Tasks	Types of Omics	Case ID
WGS data analysis	Genome assembly	Genomics	1.1
WGS/WES data analysis	Somatic SNV+indel calling	Genomics	2.1
WGS/WES data analysis	Somatic SNV+indel calling and annotation	Genomics	2.2
WGS/WES data analysis	Structure variation identification with normal	Genomics	2.3
WGS/WES data analysis	Structure variation identification without normal	Genomics	2.4
ChIP-seq data analysis	Peak calling	Genomics	3.1
ChIP-seq data analysis	Motif discovery for binding sites	Genomics	3.2
ChIP-seq data analysis	Functional enrichment of target gene	Genomics	3.3
Bisulfite-Seq data analysis	Identifying DNA methylation	Genomics	4.1
ATAC-seq data analysis	Identifying open chromatin regions	Genomics	5.1
DNase-seq data analysis	Identifying DnaseI hypersensitive site	Genomics	6.1
4C-seq data analysis	Find genomics interactions	Genomics	7.1
Nanopore DNA sequencing data analysis	Genome assembly	Genomics	8.1
Nanopore DNA sequencing data analysis	Tandem repeats variation identification	Genomics	8.2
PacBio DNA sequencing data analysis	Genome assembly	Genomics	9.1
RNA-Seq data analysis	Find Differentially expressed genes	Transcriptomics	10.1
RNA-Seq data analysis	Identify the top5 downregulated genes	Transcriptomics	10.2
RNA-Seq data analysis	Predict Fusion gene with annotation	Transcriptomics	10.3
RNA-Seq data analysis	Isoform expression	Transcriptomics	10.4
RNA-Seq data analysis	Splicing analysis	Transcriptomics	10.5
RNA-Seq data analysis	APA analysis	Transcriptomics	10.6
RNA-Seq data analysis	RNA editing	Transcriptomics	10.7
RNA-Seq data analysis	Circular RNA identification	Transcriptomics	10.8
Small RNA sequencing data analysis	microRNA quantification	Transcriptomics	11.1
Small RNA sequencing data analysis	microRNA prediction	Transcriptomics	11.2
CAGE-seq data analysis	TSS identification	Transcriptomics	12.1
3' end-seq data analysis	PAS (polyadenylation site) identification	Transcriptomics	13.1
Nanopore RNA sequencing data analysis	Isoform expression	Transcriptomics	14.1
PacBio RNA sequencing data analysis	Isoform expression	Transcriptomics	15.1
CLIP-seq data analysis	Identify protein-RNA crosslink sites	Transcriptomics	16.1
RIP-seq data analysis	Find enriched genes bounded by RBP	Transcriptomics	16.2
Ribo-seq data analysis	Identify translated ORFs	Transcriptomics	17.1
single-cell RNA-seq data analysis	Cell clustering from fastq data	Transcriptomics	18.1
single-cell RNA-seq data analysis	Find differentially expressed genes based on count matrix	Transcriptomics	18.2
single-cell RNA-seq data analysis	Find marker genes based on count matrix	Transcriptomics	18.3
single-cell RNA-seq data analysis	Cell clustering and visualization	Transcriptomics	18.4
Spatial transcriptomics	Neighborhood enrichment analysis	Transcriptomics	19.1
Spatial transcriptomics	Single-cell mapping	Transcriptomics	19.2
Mass spectrometry data analysis	Protein expression quantification	Proteomics	20.1
Mass spectrometry data analysis	Metabolites quantification	Metabolomics	21.1

of alternative splicing [42] and isoform analysis [43]. This progressive journey ultimately ends in downstream tasks like the exploration of functional enrichment [44], providing a comprehensive range of analytical pursuits. Compared to bulk RNA-seq, ChIP-seq involves distinct downstream tasks, such as peak calling [45], motif discovery [46], peak annotation [47] and so on. In summary, the analysis of different types of omics data requires professional skills and

an understanding of the corresponding field. Moreover, the methods and pipelines might vary across different bioinformaticians and they even may evolve with the development of more advanced algorithms.

Meanwhile, online bioinformatics analysis platforms are currently in vogue, such as iDEP [48], ICARUS [49] and STellaris [50]. However, they often necessitate the uploading of either raw data or pre-processed statistics by users, which

could potentially give rise to additional privacy concerns and data leakage risks [51].

In the context described above, the bioinformatics community grapples with essential concerns regarding the standardization, portability, and reproducibility of analysis pipelines [52], [53], [54]. Moreover, achieving proficiency in utilizing these pipelines for data analysis demands additional training, posing challenges for many wet lab researchers due to its potential complexity and time-consuming nature. Even dry-lab researchers may find the repetitive process of running and debugging these pipelines to be quite tedious [55]. Consequently, there is a growing anticipation within the community for the development of a more user-friendly, low-code, multi-functional, automated, and natural language-driven intelligent tool tailored for bioinformatics analysis. Such a tool has the potential to generate significant excitement and benefit researchers across the field.

Over the past few months, the rapid advancement of Large Language Models (LLMs) [56] has raised substantial expectations for the enhancement of scientific research, particularly in the field of biology [57], [58], [59]. These advancements hold promise for applications such as disease diagnosis [60], [61], [62], [63], drug discovery [64], and all. In the realm of bioinformatics, LLMs, such as ChatGPT, also demonstrate immense potential in tasks related to bioinformatics education [65] and code generation [66]. While researchers have found ChatGPT to be a valuable tool in facilitating bioinformatics research, such as data analysis, there remains a strong requirement for human involvement in the process. AutoGPT [67], as a recently developed, advanced, and experimental open-source autonomous AI agent, has the capacity to string together LLM-generated “thoughts” to autonomously achieve user-defined objectives. Nevertheless, given the intricate and specialized nature of bioinformatics tasks, the direct application of AutoGPT in this field still presents significant challenges.

In this study, we introduce Automated Bioinformatics Analysis (AutoBA), a groundbreaking autonomous AI agent tailored for comprehensive and conventional bioinformatics analysis. AutoBA simplifies user interactions to just three

inputs: data path, data description, and the final objective. This powerful tool autonomously proposes analysis plans, generates code, executes codes, and conducts subsequent data analysis by using our well-designed prompts. We implemented AutoBA as an open-source software that offers five LLM backends, with options for both online and local usage, prioritizing data security and user privacy. To show the reliability of AutoBA, we tested it in a large number of real-world bioinformatics analysis scenarios. In summary, AutoBA represents a pioneering leap in the application of Large Language Models (LLMs) and automated AI agents within the domain of bioinformatics, highlighting their potential to accelerate future research in this field.

## 2 METHODS

### 2.1 The overall framework design of AutoBA

AutoBA is the first autonomous AI agent tailor-made for conventional bioinformatics analysis. As illustrated in Fig. 1, conventional bioinformatics typically entails the use of pipelines to analyze diverse data types such as WGS, WES, RNA-seq, single-cell RNA-seq, ChIP-seq, ATAC-seq, spatial transcriptomics, and more, all requiring the utilization of various software tools. Users are traditionally tasked with selecting the appropriate software tools based on their specific analysis needs. In practice, this process involves configuring the environment, installing software, writing code, and addressing code-related issues, which are time-consuming and labor-intensive.

With the advent of AutoBA, this labor-intensive process is revolutionized. Users are relieved from the burden of dealing with multiple software packages and need only provide three key inputs: the data path (e.g., */data/SRR1374921.fasta.gz*), data description (e.g., *single-end reads in condition A*), and the ultimate analysis goal (e.g., *identify differentially expressed genes*). AutoBA takes over by autonomously analyzing the data, generating comprehensive step-by-step plans, composing code for each step, executing the generated code, and conducting in-depth analysis. Depending on the complexity and difficulty of the tasks, users can expect AutoBA to complete the tasks within a

143 matter of minutes to a few hours, all without the need for  
 144 additional manual labor (Table 1 and Fig. 2).

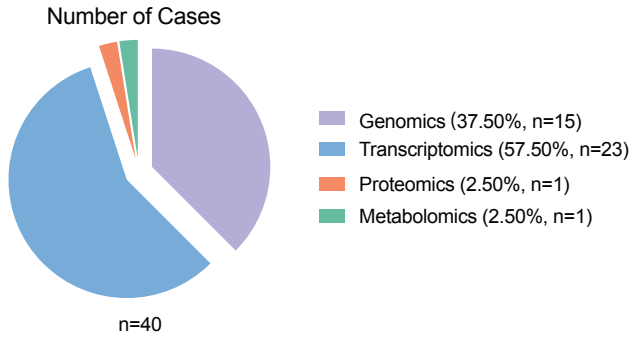


Fig. 2. Pie chart of all cases used for validating AutoBA.

## 145 2.2 Prompt engineering of AutoBA

146 To initiate AutoBA, users provide three essential inputs: the  
 147 data path, data description, and the previously mentioned  
 148 analysis objective. AutoBA comprises three distinct phases:  
 149 the planning phase, the code generation phase, and the  
 150 execution phase as shown in Step 2 of Fig. 1. During  
 151 the planning phase, AutoBA meticulously outlines a com-  
 152 prehensive step-by-step analysis plan. This plan includes  
 153 details such as the software name and version to be used at  
 154 each step, along with guided actions and specific sub-tasks  
 155 for each stage. Subsequently, in the code generation phase,  
 156 AutoBA systematically follows the plan and generates codes  
 157 for sub-tasks, which entails procedures like configuring the  
 158 environment, installing the necessary software, and writing  
 159 code. Then, in the execution phase, AutoBA executes the  
 160 generated code. In light of this workflow, AutoBA incor-  
 161 porates two distinct prompts: one tailored for the planning  
 162 phase and the other for the code generation phase. Intensive  
 163 experiments have shown that these two sets of prompts are  
 164 essential for the proper functioning of AutoBA in automated  
 165 bioinformatics analysis tasks.

166 The prompt for the planning phase is displayed as  
 167 follows:

```
168 prompt = {
169     "role": "Act as a bioinformatician, the
170         rules must be strictly followed!",
171     "rules": [
172         "When acting as a bioinformatician,
173         you strictly cannot stop acting
174         as a bioinformatician.",
175     ]
176 }
```

```

176     "All rules must be followed strictly
177     .",
178     "You should use information in input
179     to write a detailed plan to
180     finish your goal.",
181     f"You should include the software
182     name and should not use those
183     software: {self.blacklist}.",
184     "You should only respond in JSON
185     with the required format.",
186     "Your JSON should only enclosed in
187     double quotes."
188 ],
189     "input": [
190         "You have the following information
191         in a list with the format file
192         path: file description. I
193         provide those files to you, so
194         you don't need to prepare the
195         data.",
196         data_list
197     ],
198     "goal": self.current_goal,
199     "format": {
200         "plan": [
201             "Your detailed step-by-step sub-
202             tasks to finish your goal."
203         ]
204     }
205 }
206 
```

The prompt for the code generation phase is displayed  
 as follows:

```

208 prompt = {
209     "role": "Act as a bioinformatician, the
210         rules must be strictly followed!",
211     "rules": [
212         "When acting as a bioinformatician,
213         you strictly cannot stop acting
214         as a bioinformatician.",
215         "All rules must be followed strictly
216         .",
217         "You are provided a system with
218         specified constraints."
219     ],
220     "system": [
221         "The history of what you have done
222         is provided, you should take the
223         name changes of some files into
224         account, or use some output
225         from previous steps.",
226         "You should use all information you
227         have to write bash codes to
228         finish your current task.",
229         "All code requirements must be
230         followed strictly when you write
231         codes.",
232         "You should only respond in JSON
233         with the required format.",
234         "Your JSON should only enclosed in
235         double quotes."
236     ],
237     "input": [
238         "You have the following information
239         in a list with the format file
240         path: file description. I
241         provide those files to you, so
242         you don't need to prepare the
243         data.",
244         data_list
245     ],
246     "goal": self.current_goal,
247     "format": {
248         "plan": [
249             "Your detailed step-by-step sub-
250             tasks to finish your goal."
251         ]
252     }
253 }
254 
```

```

251     data_list
252 ],
253 "history": self.history_summary,
254 "current task": self.current_goal,
255 "code requirement": [
256     f"You should not use those software:
257       {self.blacklist}.",
258     'You should always source activate
259       the environment abc first, add
260       conda-forge and bioconda to the
261       list of channels',
262     'You should always install
263       dependencies with -y with conda
264       or pip.',
265     'You should pay attention to the
266       number of input files and do not
267       miss any.',
268     'You should process each file
269       independently and can not use
270       FOR loop.',
271     'You should use the path for all
272       files according to input and
273       history.',
274     'You should use the default values
275       for all parameters that are not
276       specified.',
277     'You should not repeat what you have
278       done in history.',
279     'You should only use software
280       directly you installed with
281       conda.',
282     'If you use Rscript -e, you should
283       make sure all variables exist in
284       your command, otherwise you
285       need check your history to
286       repeat previous steps and
287       generate those variables.'
288 ],
289 "format": {
290     "tool": "name of the tool you use",
291     "code": "bash code to finish the
292       current task"
293 }
294 }
295

```

In the two aforementioned prompt designs, the term *blacklist* pertains to the user's personalized list of prohibited software. The current default blacklist contains several tools frequently caused errors during our testing processes. Meanwhile, *data list* encompasses the inputs necessary for AutoBA, encompassing data paths and data descriptions. The term *current goal* serves as the final objective during the planning phase and as the sub-goal in the execution phase, while *history summary* encapsulates AutoBA's memory of previous actions and information.

## 2.3 Memory management of AutoBA

A memory mechanism is incorporated within AutoBA to enable it to generate code more effectively by drawing from past actions, thus avoiding unnecessary repetition of certain steps. AutoBA meticulously logs the outcome of each step

in a specific format, and all these historical records become part of the input for the subsequent prompt. In the planning phase, memories are structured as follows: "Firstly, you provided input in the format 'file path: file description' in a list: <data list>. You devised a detailed plan to accomplish your overarching objective. Your overarching goal is <global goal>. Your plan involves <tasks>." In the code generation phase, memories follow this format: "Then, you successfully completed the task: <task> with the corresponding code: <code>."

## 2.4 Evaluation of AutoBA

The results produced by AutoBA undergo thorough validation by bioinformatics experts. This validation process encompasses a comprehensive review of the proposed plans, generated codes, execution of the code, and confirmation of the results for accuracy and reliability. AutoBA's development and validation are built upon a specific environment and software stack, which includes Ubuntu version 18.04, Python 3.10.0, and openai version 0.27.6. These environment and software specifications form the robust foundation for AutoBA's functionality in the field of bioinformatics, ensuring its reliability and effectiveness.

## 2.5 Online and local LLM backends of AutoBA

AutoBA offers several versions of LLM backends, including online backends based on ChatGPT-3.5 and ChatGPT-4, and local LLMs, including CodeLlama-7B-Instruct, CodeLlama-13B-Instruct and CodeLlama-34B-Instruct [68].

# 3 RESULTS

## 3.1 AutoBA proposes detailed analysis plans for tasks

AutoBA offers a robust capability to generate a highly detailed and customized analysis plan, leveraging the user's input, which encompasses critical elements such as data paths, data descriptions, and objective descriptions.

As an example, in Fig. 3, the user supplied four RNA-Seq samples: two from the LoGlu group (SRR1374921.fastq.gz and SRR1374922.fastq.gz, mouse pancreatic islets cultured at low ambient glucose) and two from the HiGlu group

(SRR1374923.fastq.gz and SRR1374924.fastq.gz, mouse pancreatic islets cultured at high ambient glucose) from Benner et al.'s paper [69]. Additionally, the user also provided the mouse reference genome (mm39.fa) and genome annotation (mm39.ncbiRefSeq.gtf). The primary objective of this case was to identify differentially expressed genes between the two data groups. Using textual inputs only, AutoBA generated a detailed, step-by-step analysis plan during the planning phase, as outlined below:

```

1. First, I will use the "Trimmomatic"
   software to trim the adapters from the
   single-end mouse rna-seq reads. The
   input files for this step will be "./
   examples/case1.1/data/SRR1374921.fastq.
   gz", "./examples/case1.1/data/SRR1374922
   .fastq.gz", "./examples/case1.1/data/
   SRR1374923.fastq.gz", and the
   adapter file "./examples/case1.1/data/
   TruSeq3-SE.fa".

2. Next, I will use the "Hisat2" software to
   align the trimmed reads to the mouse
   mm39 genome. The input files for this
   step will be the output files from the
   previous step and the genome file "./
   examples/case1.1/data/mm39.fa".

3. Then, I will use the "Samtools" software
   to convert the alignment output files
   from Hisat2 (in SAM format) to BAM
   format.

4. After that, I will use the "HTSeq"
   software to count the number of reads
   mapped to each gene. The input files for
   this step will be the BAM files from
   the previous step and the annotation
   file "./examples/case1.1/data/mm39.
   ncbiRefSeq.gtf".

5. Finally, I will use the "DESeq2" software
   to identify the differentially
   expressed genes. The input files for
   this step will be the count files from
   the previous step.

```

Within each step of the plan outlined above, AutoBA provides precise instructions regarding the required software, including names like Trimmomatic, Hisat2, Samtools, HTSeq, and DESeq2, along with clear sub-tasks for each analytical stage. This level of tailored planning ensures that the analysis process aligns precisely with the user's objectives, promoting both efficiency and accuracy in data processing and results generation.

### 3.2 AutoBA generates precise codes for sub-tasks

During the code generation phase, AutoBA generates code in bash format for every sub-task of the plan established in the planning phase. These scripts encompass environment setup, software installation, and tailored code for software utilization. Parameters and data paths specific to the software are meticulously incorporated. As exemplified in Fig. 3, the preliminary phase of the differentially expressed genes (DEG) analysis constitutes the essential process of adapter trimming, an indispensable preprocessing step in the context of raw RNA-Seq data. Within this critical step, AutoBA automatically generated code, including activating the conda environment, installing software packages, and calling software to analyze data as shown below:

```

source activate abc
conda config --add channels conda-forge
conda config --add channels bioconda
conda install -y Trimmomatic
Trimmomatic SE -phred33 ./examples/case1.1/
data/SRR1374921.fastq.gz ./examples/
case1.1/output/SRR1374921_trimmed.fastq.
gz ILLUMINACLIP:./examples/case1.1/data/
TruSeq3-SE.fa:2:30:10
Trimmomatic SE -phred33 ./examples/case1.1/
data/SRR1374922.fastq.gz ./examples/
case1.1/output/SRR1374922_trimmed.fastq.
gz ILLUMINACLIP:./examples/case1.1/data/
TruSeq3-SE.fa:2:30:10
Trimmomatic SE -phred33 ./examples/case1.1/
data/SRR1374923.fastq.gz ./examples/
case1.1/output/SRR1374923_trimmed.fastq.
gz ILLUMINACLIP:./examples/case1.1/data/
TruSeq3-SE.fa:2:30:10
Trimmomatic SE -phred33 ./examples/case1.1/
data/SRR1374924.fastq.gz ./examples/
case1.1/output/SRR1374924_trimmed.fastq.
gz ILLUMINACLIP:./examples/case1.1/data/
TruSeq3-SE.fa:2:30:10

```

The generated code further underwent a meticulous and thorough validation process, which included a manual verification and execution performed by experienced and proficient bioinformaticians, as meticulously documented in Table 2 and Supplementary information. This critical validation step ensures the accuracy and reliability of the code, reaffirming the robustness of AutoBA.

### 3.3 AutoBA adeptly manages similar tasks with robustness

In practical bioinformatics applications, even when researchers are working with similar data types, such as RNA-Seq, it's noteworthy that analyses often manifest variations



**RNA-Seq Analysis****Case 1.1: Find differentially expressed genes****User's Inputs****Data and Description**

SRR1374921.fastq.gz: single-end mouse rna-seq reads, replicate 1 in LoGlu group  
 SRR1374922.fastq.gz: single-end mouse rna-seq reads, replicate 2 in LoGlu group  
 SRR1374923.fastq.gz: single-end mouse rna-seq reads, replicate 1 in HiGlu group  
 SRR1374924.fastq.gz: single-end mouse rna-seq reads, replicate 2 in HiGlu group  
 TruSeq3-SE.fa: trimming adapter  
 mm39.fa: mouse mm39 genome fasta  
 mm39.ncbiRefSeq.gtf: mouse mm39 genome annotation

**Goal**

Find the differentially expressed genes

**Generated Plans by AutoBA**

1. First, I will use the 'Trimmomatic' software to trim the adapters from the single-end mouse rna-seq reads. The input files for this step will be './examples/case1.1/data/SRR1374921.fastq.gz', './examples/case1.1/data/SRR1374922.fastq.gz', './examples/case1.1/data/SRR1374923.fastq.gz', and the adapter file './examples/case1.1/data/TruSeq3-SE.fa'.
2. Next, I will use the 'Hisat2' software to align the trimmed reads to the mouse mm39 genome. The input files for this step will be the output files from the previous step and the genome file './examples/case1.1/data/mm39.fa'.
3. Then, I will use the 'Samtools' software to convert the alignment output files from Hisat2 (in SAM format) to BAM format.
4. After that, I will use the 'HTSeq' software to count the number of reads mapped to each gene. The input files for this step will be the BAM files from the previous step and the annotation file './examples/case1.1/data/mm39.ncbiRefSeq.gtf'.
5. Finally, I will use the 'DESeq2' software to identify the differentially expressed genes. The input files for this step will be the count files from the previous step.

**Generated Codes by AutoBA**

```
# Step 1
source activate abc
conda config --add channels conda-forge
conda config --add channels bioconda
conda install -y trimmomatic
trimmomatic SE -phred33 ./examples/case1.1/data/SRR1374921.fastq.gz ./examples/case1.1/output/SRR1374921_trimmed.fastq.gz
ILLUMINACLIP:./examples/case1.1/data/TruSeq3-SE.fa:2:30:10
trimmomatic SE -phred33 ./examples/case1.1/data/SRR1374922.fastq.gz ./examples/case1.1/output/SRR1374922_trimmed.fastq.gz
ILLUMINACLIP:./examples/case1.1/data/TruSeq3-SE.fa:2:30:10
trimmomatic SE -phred33 ./examples/case1.1/data/SRR1374923.fastq.gz ./examples/case1.1/output/SRR1374923_trimmed.fastq.gz
ILLUMINACLIP:./examples/case1.1/data/TruSeq3-SE.fa:2:30:10
trimmomatic SE -phred33 ./examples/case1.1/data/SRR1374924.fastq.gz ./examples/case1.1/output/SRR1374924_trimmed.fastq.gz
ILLUMINACLIP:./examples/case1.1/data/TruSeq3-SE.fa:2:30:10

# Step 2
source activate abc
conda config --add channels conda-forge
conda config --add channels bioconda
conda install -y hisat2
hisat2 -x ./examples/case1.1/data/mm39.fa -U ./examples/case1.1/output/SRR1374921_trimmed.fastq.gz -S ./examples/case1.1/output/SRR1374921.sam
hisat2 -x ./examples/case1.1/data/mm39.fa -U ./examples/case1.1/output/SRR1374922_trimmed.fastq.gz -S ./examples/case1.1/output/SRR1374922.sam
hisat2 -x ./examples/case1.1/data/mm39.fa -U ./examples/case1.1/output/SRR1374923_trimmed.fastq.gz -S ./examples/case1.1/output/SRR1374923.sam
hisat2 -x ./examples/case1.1/data/mm39.fa -U ./examples/case1.1/output/SRR1374924_trimmed.fastq.gz -S ./examples/case1.1/output/SRR1374924.sam

# Step 3
source activate abc
conda config --add channels conda-forge
conda config --add channels bioconda
conda install -y samtools
samtools view -S -b ./examples/case1.1/output/SRR1374921.sam > ./examples/case1.1/output/SRR1374921.bam
samtools view -S -b ./examples/case1.1/output/SRR1374922.sam > ./examples/case1.1/output/SRR1374922.bam
samtools view -S -b ./examples/case1.1/output/SRR1374923.sam > ./examples/case1.1/output/SRR1374923.bam
samtools view -S -b ./examples/case1.1/output/SRR1374924.sam > ./examples/case1.1/output/SRR1374924.bam

# Step 4
source activate abc
conda config --add channels conda-forge
conda config --add channels bioconda
conda install -y htseq
htseq-count -f bam -s no ./examples/case1.1/output/SRR1374921.bam ./examples/case1.1/data/mm39.ncbiRefSeq.gtf
> ./examples/case1.1/output/SRR1374921.count
htseq-count -f bam -s no ./examples/case1.1/output/SRR1374922.bam ./examples/case1.1/data/mm39.ncbiRefSeq.gtf
> ./examples/case1.1/output/SRR1374922.count
htseq-count -f bam -s no ./examples/case1.1/output/SRR1374923.bam ./examples/case1.1/data/mm39.ncbiRefSeq.gtf
> ./examples/case1.1/output/SRR1374923.count
htseq-count -f bam -s no ./examples/case1.1/output/SRR1374924.bam ./examples/case1.1/data/mm39.ncbiRefSeq.gtf
> ./examples/case1.1/output/SRR1374924.count

# Step 5
source activate abc
conda install -y bioconductor-deseq2
Rscript -e "library(DESeq2); countData <- read.table('./examples/case1.1/output/SRR1374921.count', './examples/case1.1/output/SRR1374922.count',
'./examples/case1.1/output/SRR1374923.count', './examples/case1.1/output/SRR1374924.count', header=TRUE, row.names=1); colData <-
DataFrame(condition=factor(c('LoGlu', 'LoGlu', 'HiGlu', 'HiGlu'))); dds <- DESeqDataSetFromMatrix(countData=countData, colData=colData, design=~condition);
dds <- DESeq(dds); res <- results(dds); write.csv(as.data.frame(res), file='./examples/case1.1/output/differential_expression_results.csv')"
```

Fig. 3. **Example of applying AutoBA to find differentially expressed genes with RNA-Seq data.** In this case, the user supplied four RNA-seq datasets, comprising two from the LoGlu group and two from the HiGlu group. The primary objective of this analysis was to identify differentially expressed genes across the two datasets.



stemming from diverse sources. These variations are primarily attributed to disparities in the characteristics of input data and the distinct objectives pursued in the analytical process.

As exemplified in Case 10.1 (find differentially expressed genes), Case 10.2 (identify the top five down-regulated genes in HiGlu group), and Case 10.3 (predict fusion genes), when performing RNA-Seq analysis, users may have distinct final goals, necessitating adjustments in software and parameter selection during the actual execution. In comparison to case 10.1, AutoBA introduces a 6th step in case 10.2, tailored for screening the top five differentially expressed genes to fulfill the user's specific requirements as shown in the code below:

```
Rscript -e "library('pheatmap'); library('DESeq2'); res <- read.csv('./examples/case1.2/output/differential_expression_results.csv', row.names=1); res_ordered <- res[order(res$log2FoldChange),]; top5_downregulated <- head(res_ordered, 5);"
```

### 3.4 AutoBA adjusts analysis based on task and input data variations

Alignment is an essential step for bioinformatic analysis, for which multiple tools have been developed for distinct tasks. For instance, tools including STAR [70] and HISAT2 [71] designed for RNA-seq data analysis are splicing aware, which is efficient in identifying junction reads that map to two distal positions in the reference genome. Besides, long-read sequencing data from Pacific Bioscience (PacBio) and Oxford Nanopore Technology (ONT) also require specialized tools for the alignment, for which Minimap2 [72] is the most widely used method. Moreover, each read from single-cell sequencing data contains barcodes for UMI and cell labels, which needs to be integrated with the alignment. Cell Ranger is a popular software with this capacity. Therefore, bioinformatic analysis should use appropriate tools for the alignment based on the types of tasks. Interestingly, we found that AutoBA has learned this knowledge and can correctly employ the tool for the alignment (Fig. 4a).

For many bioinformatic analysis, multiple tools are available but require different conditions of inputs. For instance,

to identify structural variations from tumor WGS/WES data, the method "manta" [73] can handle the analysis against the matched normal. On the other hand, tools like "Pindel" [74] that relies on the detection of breakpoints with the reference genome, only conduct analysis on the tumor samples. We found that AutoBA can automatically select "manta" when the matched normal samples were provided and correctly utilized the parameters "--normalBam" and "--tumorBam". However, if only the tumor samples were provided in the input data, AutoBA will select "Pindel" for the analysis (Fig. 4b). These results suggest that AutoBA learned the requirements of different bioinformatic tools and is capable of selecting appropriate tools based on different conditions of the input data.

```
manta --normalBam ../output/case2.3/SRR23015874.recalibrated.bam --tumorBam ../output/case2.3/SRR23015876.recalibrated.bam --referenceFasta ./examples/case5.1/hg38.fa --runDir ../output/case2.3/manta_SRR23015874
```

### 3.5 Apply AutoBA to a variety of conventional bioinformatics analysis scenarios

To evaluate the robustness of AutoBA, we conducted assessments involving a total of 40 cases spanning four distinct types of omics data: genomics, transcriptomics, proteomics, and metabolomics as shown in Table 1 and Supplementary information.

All cases underwent an independent analysis process conducted by AutoBA and were subsequently subjected to validation by experienced bioinformatics experts. The collective results underscore the versatility and robustness of AutoBA across a spectrum of multi-omics analysis procedures in the field of bioinformatics as shown in Table 2. AutoBA demonstrates its capability to autonomously devise novel analysis processes based on varying input data, showcasing its adaptability to diverse input data and analysis objectives with a success rate of 90% (36 out of 40) for proposing plans, 82.5% (33 out of 40) for generating codes, and 65% (26 out of 40) for automated end-to-end analysis.

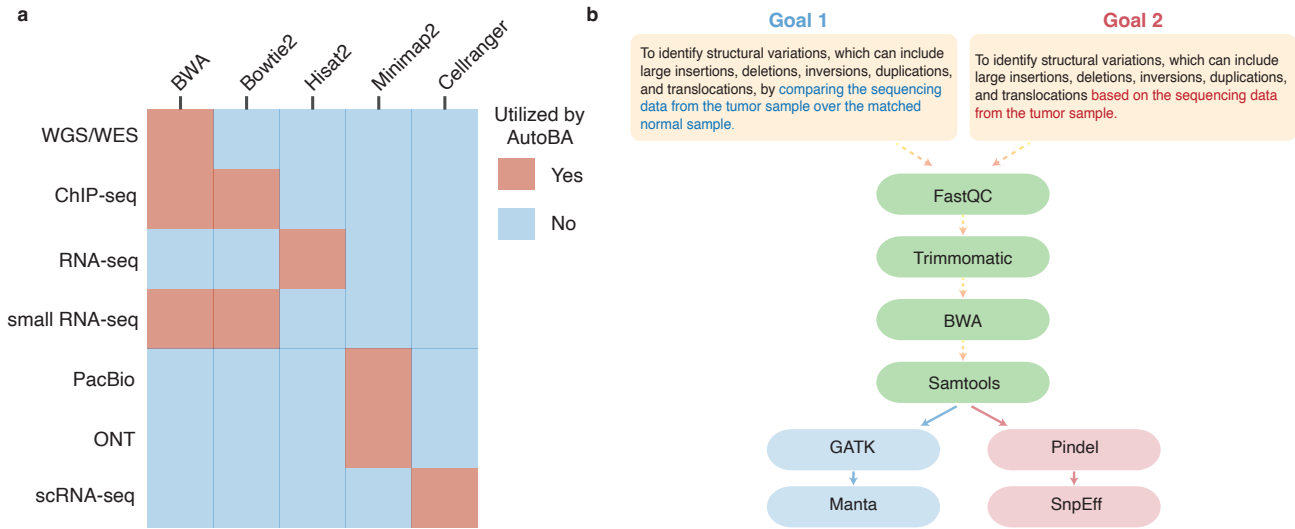


Fig. 4. **AutoBA adjusts analysis based on task and input data variations** **a** Heatmap illustrating options of utilizing different alignment tools for multiple tasks planned by AutoBA. **b** AutoBA utilizes the tools for identifying structure variations in tumor samples with or without the matched normal samples. The highlight shows the difference between Goal 1 and Goal 2.

4 DISCUSSION

To our knowledge, AutoBA is the first and a pioneering autonomous AI agent tailored explicitly for conventional bioinformatics analysis for omics data. AutoBA streamlines the analytical process, requiring minimal user input while providing detailed step-by-step plans for various bioinformatics tasks. The results of our investigation reveal that AutoBA excels in accurately handling a diverse array of omics analysis tasks, such as RNA-seq, scRNA-seq, ChIP-seq, spatial transcriptomics, and so on. One of the key strengths of AutoBA is its adaptability to variations in analysis objectives. As demonstrated in the cases presented, even with similar data types, such as RNA-Seq, users often have distinct goals, necessitating modifications in software and parameter selection during execution. AutoBA effectively accommodates these variations, allowing users to tailor their analyses to specific research needs without compromising accuracy. Furthermore, AutoBA’s versatility is highlighted by its ability to self-design new analysis processes based on differing input data. This autonomous adaptability makes AutoBA a valuable tool for bioinformaticians working on novel or unconventional research questions, as it can adjust its approach to the unique characteristics of the data.

Online bioinformatics analysis platforms are currently in vogue, but they often necessitate the uploading of either

raw data or pre-processed statistics by users, which could potentially give rise to privacy concerns and data leakage risks. In contrast, AutoBA offers a local solution that effectively addresses these privacy issues. Moreover, AutoBA showcases its adaptability in sync with emerging bioinformatics tools, with LLM seamlessly incorporating these latest tools into the database. Furthermore, AutoBA is inclined towards selecting the most popular analytical frameworks or widely applicable tools in the planning phase, underscoring its robustness. Another distinguishing feature is AutoBA’s transparent and interpretable execution process. This transparency allows professional bioinformaticians to easily modify and customize AutoBA’s outputs, leveraging AutoBA to expedite the data analysis process.

Given that classical bioinformatic analysis encompasses a far broader spectrum of tasks and challenges than the 40 cases studied in this work (Table 1 and 2), it is essential to conduct more real-world applications by our potential users to further comprehensively validate the robustness of AutoBA. We found that a large proportion (36%, 5 out of 14) of failed cases in executing code is due to the tools in conda being problematic, not in a regular form (end with .sh, .pl et al), or requiring an edited config file, suggesting a demand for more standard bioinformatics tools. Furthermore, taking into account the timeliness of the training data used for

TABLE 2

**Summary of AutoBA generated results evaluated by bioinformatics experts.** The table presents an assessment conducted by bioinformatics experts on the analysis plan proposed by AutoBA, along with the generated codes and the code execution. If the evaluation passes, it is displayed as success, while instances of failure are accompanied by detailed explanations of the specific reasons for the failure. Additionally, we provide a summary of the software tools automatically chosen by AutoBA for each case, as well as the total time taken to generate the corresponding code.

Case ID	Propose Plans	Generate Codes	Execute Codes	Tools Used	Time Cost (without Executing Codes) in Minutes
1.1	Success	Success	Success	FastQC, Trimmomatic, SPAdes, QUAST	3
2.1	Success	Success	Success	FastQC, Trimmomatic, BWA, Samtools, GATK	8
2.2	Success	Success	Success	FastQC, Trimmomatic, BWA, Samtools, GATK, ensemble-vep	8
2.3	Success	Success	Success	FastQC, Trimmomatic, BWA, Samtools, GATK, manta	18
2.4	Success	Success	Failed: pindel it requires configuration file	FastQC, Trimmomatic, BWA, Samtools, pindel, SnpEff	6
3.1	Success	Success	Success	FastQC, Trim Galore, Bowtie 2, Samtools, MACS2, BEDTools, IGV	6
3.2	Success	Success	Success	FastQC, Trim Galore, Bowtie2, MACS2, HOMER, MEME	4
3.3	Failed: DESeq2 is not suitable for peaks identified by MACS2	-	-	FastQC, BWA, MACS, BEDTools, DESeq2, g:Profiler, R	6
4.1	Success	Success	Success	Trim Galore, Bismark, IGV	9
5.1	Success	Success	Failed: (wrongly usedBEDTools)	Trim Galore, BWA, Samtools, MACS2, BEDTools	8
6.1	Success	Success	Success	FastQC, Cutadapt, BWA, MACS2, IGV, GREAT	5
7.1	Success	Success	Success	FastQC, BEDTools, Samtools, Bowtie 2, R	6
8.1	Success	Success	Failed: racon medaka wrongly used the parameters	canu, Minimap2, Racon, Flye, Medaka, Bandage	7
8.2	Failed: cannot find a correct pipeline	-	-	Minimap2, Samtools, trf	7
9.1	Success	Failed: install the wrong tool, pb-falcon rather than falcon	-	Canu, FALCON, Quiver, MUMmer,	7
10.1	Success	Success	Success	FASTQC, Trimmomatic, HISAT2, htseq, DESeq2	5
10.2	Success	Success	Success	FASTQC, Trimmomatic, HISAT2, htseq, DESeq2, gprofileR	5
10.3	Success	Success	Success	gunzip, HISAT2, fusioncatcher, gffcompare	6
10.4	Success	Success	Success	Trim Galore, HISAT2, Samtools, StringTie	5
10.5	Success	Success	Success	Trimmomatic, HISAT2, Samtools, StringTie, featureCounts, rMATs	6
10.6	Success	Failed: DaPars (not available in conda)	-	Trim Galore, HISAT2, StringTie, DaPars	7
10.7	Failed: cannot find a correct pipeline	-	-	FastQC, Trimmomatic, HISAT2, Samtools, StringTie, ballgown, GATK	7
10.8	Success	Failed: CIRI2 (not available in conda)	-	Trim Galore, HISAT2, CIRI2, CIRIQuant	5
11.1	Success	Success	Success	Fastqc, Cutadapt, Bowtie, Samtools, subread/featureCounts, DESeq2, edgeR	11
11.2	Success	Success	Failed: conda of miRDeep2 is problematic	Fastqc, Cutadapt, Bowtie, Samtools, featureCounts, miRDeep2, DESeq2, edgeR	11
12.1	Success	Success	Success	Fastqc, Trimmomatic, HISAT2, HTSeq/htseq-count, CAGEr	6
13.1	Failed: cannot find a correct pipeline	-	-	Trim Galore, HISAT2, StringTie, DaPars	5
14.1	Success	Success	Failed: prepDE.py no need to run with 'python prepDE.py'	Minimap2, Samtools, StringTie, DESeq2	9
15.1	Success	Success	Success	Minimap2, Samtools, StringTie, cufflinks	5
16.1	Success	Success	Failed: conda of Piranha is problematic	FastQC, Cutadapt, Bowtie2, Samtools, BEDTools, Piranha	6
16.2	Success	Success	Success	FastQC, Trim Galore, HISAT2, htseq, DESeq2	4
17.1	Success	Success	Failed: not regular conda of ribotaper	FastQC, Trim Galore, HISAT2, Samtools, StringTie, RiboTaper	7
18.1	Success	Success	Success	Cell Ranger, Seurat	5
18.2	Success	Success	Success	Scanpy	8
18.3	Success	Success	Success	Scanpy	6
18.4	Success	Success	Success	Scanpy	5
19.1	Success	Success	Success	Squidpy, AnnData	5
19.2	Success	Success	Success	AnnData, Scanpy, Tangram	3
20.1	Success	Success	Success	proteowizard, OpenMS	15
21.1	Success	Success	Success	pymzml, pandas, numpy, scipy	13

large language models, it's important to note that some of the most recently proposed methods in bioinformatics may still pose challenges in automatically generating code by AutoBA. Therefore, a future endeavor to train an up-to-date

large language model explicitly tailored for bioinformatics can significantly enhance AutoBA's ability to maintain up-to-date code generation capabilities. Nevertheless, AutoBA represents a significant advancement in the field of bioin-

formatics, offering a user-friendly, efficient, and adaptable solution for a wide range of omics analysis tasks. Its capacity to handle diverse data types and analysis goals, coupled with its robustness and adaptability, positions AutoBA as a valuable asset in the pursuit of accelerating bioinformatics research. We anticipate that AutoBA will find extensive utility in the scientific community, supporting researchers in their quest to extract meaningful insights from complex biological data.

## 5 DATA AVAILABILITY

The RNA-seq dataset could be downloaded from Sequence Read Archive (SRA) with IDs: SRR1374921, SRR1374922, SRR1374923, and SRR1374924. The dataset for case 1.3 could be downloaded from <https://github.com/STAR-Fusion/STAR-Fusion-Tutorial/wiki>. The scRNA-seq dataset could be downloaded from [http://cf.10xgenomics.com/samples/cell-exp/1.1.0/pbmc3k/pbmc3k\\_filtered\\_gene\\_bc\\_matrices.tar.gz](http://cf.10xgenomics.com/samples/cell-exp/1.1.0/pbmc3k/pbmc3k_filtered_gene_bc_matrices.tar.gz). The ChIP-seq dataset could be downloaded with IDs: SRR620204, SRR620205, SRR620206, and SRR620208. The Spatial Transcriptomics dataset could be downloaded from <https://doi.org/10.5281/zenodo.6334774>. The CAGE-seq dataset could be downloaded from SRA with IDs: SRR11351697, SRR11351698, SRR11351700, and SRR11351701. The 3'-end-seq dataset could be downloaded from SRA with IDs: SRR17422754, SRR17422755, SRR17422756, and SRR17422757. The CLIP-seq dataset could be downloaded from ENCODE (<https://www.encodeproject.org>) with IDs: ENCLB742AYH and ENCLB770EDJ. The Ribo-seq data could be downloaded from SRA with IDs: RR12354645 and RR12354646. The raw single-cell RNA sequencing data could be downloaded from 10X genomics. The PacBio long-read sequencing data could be downloaded from SRA with IDs: SRR19552218 and SRR19785215. The small RNA-seq data could be downloaded from the previous study [75].

## 6 CODE AVAILABILITY

The AutoBA software is publicly available at <https://github.com/JoshuaChou2018/AutoBA>.

## 7 CREDIT AUTHOR STATEMENT

Conceptualization: J.Z. and X.G. Design: J.Z., B.Z. and X.G. Code implementation: J.Z. Application: J.Z., B.Z., X.C., H.L. Drafting of the manuscript: J.Z. and B.Z. Critical revision of the manuscript for important intellectual content: J.Z., B.Z., X.X., S.C., X.G. Supervision: J.Z. and X.G. Funding acquisition: X.G.

## 8 ACKNOWLEDGEMENTS

Juexiao Zhou, Bin Zhang, Xiuying Chen, Haoyang Li, Xiaopeng Xu, Siyuan Chen and Xin Gao were supported in part by grants from the Office of Research Administration (ORA) at King Abdullah University of Science and Technology (KAUST) under award number FCC/1/1976-44-01, FCC/1/1976-45-01, REI/1/5202-01-01, REI/1/5234-01-01, REI/1/4940-01-01, RGC/3/4816-01-01, and REI/1/0018-01-01.

## 9 COMPETING INTERESTS

The authors have declared no competing interests.

## REFERENCES

- [1] N. M. Luscombe, D. Greenbaum, and M. Gerstein, "What is bioinformatics? a proposed definition and overview of the field," *Methods of information in medicine*, vol. 40, no. 04, pp. 346–358, 2001.
- [2] J. Gauthier, A. T. Vincent, S. J. Charette, and N. Derome, "A brief history of bioinformatics," *Briefings in bioinformatics*, vol. 20, no. 6, pp. 1981–1996, 2019.
- [3] A. D. Baxevanis, G. D. Bader, and D. S. Wishart, *Bioinformatics*. John Wiley & Sons, 2020.
- [4] P. Munk, C. Brinch, F. D. Møller, T. N. Petersen, R. S. Hendriksen, A. M. Seyfarth, J. S. Kjeldgaard, C. A. Svendsen, B. Van Bunnik, F. Berglund *et al.*, "Genomic analysis of sewage from 101 countries reveals global landscape of antimicrobial resistance," *Nature Communications*, vol. 13, no. 1, p. 7251, 2022.
- [5] F. Hemmerling and J. Piel, "Strategies to access biosynthetic novelty in bacterial genomes for drug discovery," *Nature Reviews Drug Discovery*, vol. 21, no. 5, pp. 359–378, 2022.

- 676 [6] E. H. Lips, T. Kumar, A. Megalios, L. L. Visser, M. Sheinman,  
677 A. Fortunato, V. Shah, M. Hoogstraat, E. Sei, D. Mallo *et al.*, "Ge-  
678 nomic analysis defines clonal relationships of ductal carcinoma in  
679 situ and recurrent invasive breast cancer," *Nature genetics*, vol. 54,  
680 no. 6, pp. 850–860, 2022.
- 681 [7] G. Orlando, D. Raimondi, R. Duran-Romaña, Y. Moreau,  
682 J. Schymkowitz, and F. Rousseau, "Pyuul provides an interface  
683 between biological structures and deep learning algorithms," *Nature communications*, vol. 13, no. 1, p. 961, 2022.
- 684 [8] D. T. Jones and J. M. Thornton, "The impact of alphafold2 one year  
685 on," *Nature methods*, vol. 19, no. 1, pp. 15–20, 2022.
- 686 [9] M. L. Hekkelman, I. de Vries, R. P. Joosten, and A. Perrakis, "Al-  
687 phafill: enriching alphafold models with ligands and cofactors,"  
688 *Nature Methods*, vol. 20, no. 2, pp. 205–213, 2023.
- 689 [10] N. Sapoval, A. Aghazadeh, M. G. Nute, D. A. Antunes, A. Balaji,  
690 R. Baraniuk, C. Barberan, R. Dannenfelser, C. Dun, M. Edrisi *et al.*,  
691 "Current progress and open challenges for applying deep learning  
692 across the biosciences," *Nature Communications*, vol. 13, no. 1, p.  
693 1728, 2022.
- 694 [11] T. Gupta, M. Zaki, and N. A. Krishnan, "Matscibert: A materials  
695 domain language model for text mining and information extrac-  
696 tion," *npj Computational Materials*, vol. 8, no. 1, p. 102, 2022.
- 697 [12] A. Santos, A. R. Colaço, A. B. Nielsen, L. Niu, M. Strauss, P. E.  
698 Geyer, F. Coscia, N. J. W. Albrechtsen, F. Mundt, L. J. Jensen *et al.*,  
699 "A knowledge graph to interpret clinical proteomics data," *Nature*  
700 *Biotechnology*, vol. 40, no. 5, pp. 692–702, 2022.
- 701 [13] Z. Zeng, Y. Yao, Z. Liu, and M. Sun, "A deep-learning system  
702 bridging molecule structure and biomedical text with comprehen-  
703 sion comparable to human professionals," *Nature communications*,  
704 vol. 13, no. 1, p. 862, 2022.
- 705 [14] S. W. Attwood, S. C. Hill, D. M. Aanensen, T. R. Connor, and  
706 O. G. Pybus, "Phylogenetic and phylodynamic approaches to  
707 understanding and combating the early sars-cov-2 pandemic,"  
708 *Nature Reviews Genetics*, vol. 23, no. 9, pp. 547–562, 2022.
- 709 [15] N. De Maio, P. Kalaghatgi, Y. Turakhia, R. Corbett-Detig, B. Q.  
710 Minh, and N. Goldman, "Maximum likelihood pandemic-scale  
711 phylogenetics," *Nature Genetics*, pp. 1–7, 2023.
- 712 [16] A. S. Chanderbali, L. Jin, Q. Xu, Y. Zhang, J. Zhang, S. Jian,  
713 E. Carroll, D. Sankoff, V. A. Albert, D. G. Howarth *et al.*, "Buxus  
714 and tetracentron genomes help resolve eudicot genome history,"  
715 *Nature communications*, vol. 13, no. 1, p. 643, 2022.
- 716 [17] J. Rhodes, A. Abdolrasouli, K. Dunne, T. R. Sewell, Y. Zhang,  
717 E. Ballard, A. P. Brackin, N. van Rhijn, H. Chown, A. Tsitsopoulou  
718 *et al.*, "Population genomics confirms acquisition of drug-resistant  
719 aspergillus fumigatus infection by humans from the environ-  
720 ment," *Nature microbiology*, vol. 7, no. 5, pp. 663–674, 2022.
- 721 [18] R. J. Rockett, J. Draper, M. Gall, E. M. Sim, A. Arnott, J. E. Agius,  
722 J. Johnson-Mackinnon, W. Fong, E. Martinez, A. P. Drew *et al.*,  
723 "Co-infection with sars-cov-2 omicron and delta variants revealed  
724 by genomic surveillance," *Nature communications*, vol. 13, no. 1, p.  
725 2745, 2022.
- 726 [19] A. Heinken, J. Hertel, G. Acharya, D. A. Ravcheev, M. Nyga,  
727 O. E. Okpala, M. Hogan, S. Magnúsdóttir, F. Martinelli, B. Nap  
728 *et al.*, "Genome-scale metabolic reconstruction of 7,302 human  
729 microorganisms for personalized medicine," *Nature Biotechnology*,  
730 pp. 1–12, 2023.
- 731 [20] J. Zhou, B. Zhang, H. Li, L. Zhou, Z. Li, Y. Long, W. Han, M. Wang,  
732 H. Cui, J. Li *et al.*, "Annotating tss in multiple cell types based  
733 on dna sequence and rna-seq data via deerec-tss," *Genomics*,  
734 *Proteomics & Bioinformatics*, vol. 20, no. 5, pp. 959–973, 2022.
- 735 [21] H. Li, H. Li, J. Zhou, and X. Gao, "Sd2: spatially resolved transcrip-  
736 tomics deconvolution through integration of dropout and spatial  
737 information," *Bioinformatics*, vol. 38, no. 21, pp. 4878–4884, 2022.
- 738 [22] T. Zhang, L. Li, H. Sun, D. Xu, and G. Wang, "Deepicsh: a complex  
739 deep learning framework for identifying cell-specific silencers and  
740 their strength from the human genome," *Briefings in Bioinformatics*,  
741 p. bbad316, 2023.
- 742 [23] Z. Li, E. Gao, J. Zhou, W. Han, X. Xu, and X. Gao, "Applications  
743 of deep learning in understanding gene regulation," *Cell Reports*  
744 *Methods*, 2023.
- 745 [24] Y. Long, B. Zhang, S. Tian, J. J. Chan, J. Zhou, Z. Li, Y. Li, Z. An,  
746 X. Liao, Y. Wang *et al.*, "Accurate transcriptome-wide identification  
747 and quantification of alternative polyadenylation from rna-seq  
748 data with apaiq," *Genome Research*, vol. 33, no. 4, pp. 644–657, 2023.
- 749 [25] A. F. Bardet, Q. He, J. Zeitlinger, and A. Stark, "A computational  
750 pipeline for comparative chip-seq analyses," *Nature protocols*,  
751 vol. 7, no. 1, pp. 45–61, 2012.
- 752 [26] B. Vieth, S. Parekh, C. Ziegenhain, W. Enard, and I. Hellmann,  
753 "A systematic evaluation of single cell rna-seq analysis pipelines,"  
754 *Nature communications*, vol. 10, no. 1, p. 4667, 2019.
- 755 [27] M. D. Luecken and F. J. Theis, "Current best practices in single-  
756 cell rna-seq analysis: a tutorial," *Molecular systems biology*, vol. 15,  
757 no. 6, p. e8746, 2019.
- 758 [28] F. C. Grandi, H. Modi, L. Kampman, and M. R. Corces, "Chromatin  
759 accessibility profiling by atac-seq," *Nature protocols*, vol. 17, no. 6,  
760 pp. 1518–1552, 2022.
- 761 [29] P. C. Ng and E. F. Kirkness, "Whole genome sequencing," *Genetic*  
762 *variation: Methods and protocols*, pp. 215–226, 2010.
- 763 [30] Z. Wang, M. Gerstein, and M. Snyder, "Rna-seq: a revolutionary  
764 tool for transcriptomics," *Nature reviews genetics*, vol. 10, no. 1, pp.  
765 57–63, 2009.
- 766 [31] A.-E. Saliba, A. J. Westermann, S. A. Gorski, and J. Vogel, "Single-  
767 cell rna-seq: advances and future challenges," *Nucleic acids research*,  
768 vol. 42, no. 14, pp. 8845–8860, 2014.
- 769 [32] J. D. Buenrostro, B. Wu, H. Y. Chang, and W. J. Greenleaf, "Atac-  
770 seq: a method for assaying chromatin accessibility genome-wide,"  
771 *Current protocols in molecular biology*, vol. 109, no. 1, pp. 21–29, 2015.
- 772 [33] P. J. Park, "Chip-seq: advantages and challenges of a maturing  
773 technology," *Nature reviews genetics*, vol. 10, no. 10, pp. 669–680,  
774 2009.
- 775 [34] D. J. Burgess, "Spatial transcriptomics coming of age," *Nature*  
776 *Reviews Genetics*, vol. 20, no. 6, pp. 317–317, 2019.
- 777 [35] A. Conesa, P. Madrigal, S. Tarazona, D. Gomez-Cabrero,  
778 A. Cervera, A. McPherson, M. W. Szczesniak, D. J. Gaffney, L. L.  
779

- Elo, X. Zhang *et al.*, "A survey of best practices for rna-seq data analysis," *Genome biology*, vol. 17, no. 1, pp. 1–19, 2016.
- [36] L. Wang, S. Wang, and W. Li, "Rseqc: quality control of rna-seq experiments," *Bioinformatics*, vol. 28, no. 16, pp. 2184–2185, 2012.
- [37] M. Martin, "Cutadapt removes adapter sequences from high-throughput sequencing reads," *EMBnet. journal*, vol. 17, no. 1, pp. 10–12, 2011.
- [38] A. Dobin and T. R. Gingeras, "Mapping rna-seq reads with star," *Current protocols in bioinformatics*, vol. 51, no. 1, pp. 11–14, 2015.
- [39] C. Trapnell, L. Pachter, and S. L. Salzberg, "Tophat: discovering splice junctions with rna-seq," *Bioinformatics*, vol. 25, no. 9, pp. 1105–1111, 2009.
- [40] Y. Liao, G. K. Smyth, and W. Shi, "featurecounts: an efficient general purpose program for assigning sequence reads to genomic features," *Bioinformatics*, vol. 30, no. 7, pp. 923–930, 2014.
- [41] F. Rapaport, R. Khanin, Y. Liang, M. Pirun, A. Krek, P. Zumbo, C. E. Mason, N. D. Socci, and D. Betel, "Comprehensive evaluation of differential gene expression analysis methods for rna-seq data," *Genome biology*, vol. 14, no. 9, pp. 1–13, 2013.
- [42] S. Shen, J. W. Park, Z.-x. Lu, L. Lin, M. D. Henry, Y. N. Wu, Q. Zhou, and Y. Xing, "rmats: robust and flexible detection of differential alternative splicing from replicate rna-seq data," *Proceedings of the National Academy of Sciences*, vol. 111, no. 51, pp. E5593–E5601, 2014.
- [43] Y. Katz, E. T. Wang, E. M. Airoidi, and C. B. Burge, "Analysis and design of rna sequencing experiments for identifying isoform regulation," *Nature methods*, vol. 7, no. 12, pp. 1009–1015, 2010.
- [44] X. Wang and M. J. Cairns, "Gene set enrichment analysis of rna-seq data: integrating differential expression and splicing," in *BMC bioinformatics*, vol. 14, no. 5. BioMed Central, 2013, pp. 1–10.
- [45] R. Thomas, S. Thomas, A. K. Holloway, and K. S. Pollard, "Features that define the best chip-seq peak calling algorithms," *Briefings in bioinformatics*, vol. 18, no. 3, pp. 441–450, 2017.
- [46] T. L. Bailey, "Dreme: motif discovery in transcription factor chip-seq data," *Bioinformatics*, vol. 27, no. 12, pp. 1653–1659, 2011.
- [47] G. Yu, L.-G. Wang, and Q.-Y. He, "Chipseeker: an r/bioconductor package for chip peak annotation, comparison and visualization," *Bioinformatics*, vol. 31, no. 14, pp. 2382–2383, 2015.
- [48] S. X. Ge, E. W. Son, and R. Yao, "idep: an integrated web application for differential expression and pathway analysis of rna-seq data," *BMC bioinformatics*, vol. 19, no. 1, pp. 1–24, 2018.
- [49] A. Jiang, K. Lehnert, L. You, and R. G. Snell, "Icarus, an interactive web server for single cell rna-seq analysis," *Nucleic acids research*, vol. 50, no. W1, pp. W427–W433, 2022.
- [50] X. Li, C. Xiao, J. Qi, W. Xue, X. Xu, Z. Mu, J. Zhang, C.-Y. Li, and W. Ding, "Stellaris: a web server for accurate spatial mapping of single cells based on spatial transcriptomics data," *Nucleic Acids Research*, p. gkad419, 2023.
- [51] J. Zhou, S. Chen, Y. Wu, H. Li, B. Zhang, L. Zhou, Y. Hu, Z. Xiang, Z. Li, N. Chen *et al.*, "Ppml-omics: a privacy-preserving federated machine learning method protects patients' privacy in omic data," *bioRxiv*, pp. 2022–03, 2022.
- [52] S. Roy, C. Coldren, A. Karunamurthy, N. S. Kip, E. W. Klee, S. E. Lincoln, A. Leon, M. Pullambhatla, R. L. Temple-Smolkin, K. V. Voelkerding *et al.*, "Standards and guidelines for validating next-generation sequencing bioinformatics pipelines: a joint recommendation of the association for molecular pathology and the college of american pathologists," *The Journal of Molecular Diagnostics*, vol. 20, no. 1, pp. 4–27, 2018.
- [53] P. A. Ewels, A. Peltzer, S. Fillinger, H. Patel, J. Alneberg, A. Wilm, M. U. Garcia, P. Di Tommaso, and S. Nahnsen, "The nf-core framework for community-curated bioinformatics pipelines," *Nature biotechnology*, vol. 38, no. 3, pp. 276–278, 2020.
- [54] L. Wratten, A. Wilm, and J. Göke, "Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers," *Nature methods*, vol. 18, no. 10, pp. 1161–1168, 2021.
- [55] E. B. Işık, M. D. Brazas, R. Schwartz, B. Gaeta, P. M. Palagi, C. W. van Gelder, P. Suravajhala, H. Singh, S. L. Morgan, H. Zahroh *et al.*, "Grand challenges in bioinformatics education and training," *Nature Biotechnology*, vol. 41, no. 8, pp. 1171–1174, 2023.
- [56] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler *et al.*, "Emergent abilities of large language models," *arXiv preprint arXiv:2206.07682*, 2022.
- [57] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, "Large language models in medicine," *Nature Medicine*, pp. 1–11, 2023.
- [58] A. Madani, B. Krause, E. R. Greene, S. Subramanian, B. P. Mohr, J. M. Holton, J. L. Olmos Jr, C. Xiong, Z. Z. Sun, R. Socher *et al.*, "Large language models generate functional protein sequences across diverse families," *Nature Biotechnology*, pp. 1–8, 2023.
- [59] B. Meskó and E. J. Topol, "The imperative for regulatory oversight of large language models (or generative ai) in healthcare," *npj Digital Medicine*, vol. 6, no. 1, p. 120, 2023.
- [60] S. Wang, Z. Zhao, X. Ouyang, Q. Wang, and D. Shen, "Chatcad: Interactive computer-aided diagnosis on medical image using large language models," *arXiv preprint arXiv:2302.07257*, 2023.
- [61] J. Zhou, X. He, L. Sun, J. Xu, X. Chen, Y. Chu, L. Zhou, X. Liao, B. Zhang, and X. Gao, "Skingpt-4: An interactive dermatology diagnostic system with visual large language model," *medRxiv*, pp. 2023–06, 2023.
- [62] J. Zhou, X. Chen, and X. Gao, "Path to medical agi: Unify domain-specific medical llms with the lowest cost," *arXiv preprint arXiv:2306.10765*, 2023.
- [63] T. Tu, S. Azizi, D. Driess, M. Schaeckermann, M. Amin, P.-C. Chang, A. Carroll, C. Lau, R. Tanno, I. Ktena *et al.*, "Towards generalist biomedical ai," *arXiv preprint arXiv:2307.14334*, 2023.
- [64] D. Flam-Shepherd, K. Zhu, and A. Aspuru-Guzik, "Language models can learn complex molecular distributions," *Nature Communications*, vol. 13, no. 1, p. 3293, 2022.
- [65] E. Shue, L. Liu, B. Li, Z. Feng, X. Li, and G. Hu, "Empowering beginners in bioinformatics with chatgpt," *bioRxiv*, pp. 2023–03, 2023.
- [66] S. R. Piccolo, P. Denny, A. Luxton-Reilly, S. Payne, and P. G. Ridge,

- 883 “Many bioinformatics programming tasks can be automated with  
884 chatgpt,” *arXiv preprint arXiv:2303.13528*, 2023.
- 885 [67] S. Gravitass, “Auto-gpt: An autonomous gpt-4 experiment,” 2023.
- 886 [68] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi,  
887 J. Liu, T. Remez, J. Rapin *et al.*, “Code llama: Open foundation  
888 models for code,” *arXiv preprint arXiv:2308.12950*, 2023.
- 889 [69] C. Benner, T. van der Meulen, E. Cacères, K. Tigyi, C. J. Don-  
890 aldson, and M. O. Huising, “The transcriptional landscape of  
891 mouse beta cells compared to human beta cells reveals notable  
892 species differences in long non-coding rna and protein-coding  
893 gene expression,” *BMC genomics*, vol. 15, no. 1, pp. 1–17, 2014.
- 894 [70] A. Dobin, C. A. Davis, F. Schlesinger, J. Drenkow, C. Zaleski, S. Jha,  
895 P. Batut, M. Chaisson, and T. R. Gingeras, “Star: ultrafast universal  
896 rna-seq aligner,” *Bioinformatics*, vol. 29, no. 1, pp. 15–21, 2013.
- 897 [71] D. Kim, J. M. Paggi, C. Park, C. Bennett, and S. L. Salzberg,  
898 “Graph-based genome alignment and genotyping with hisat2 and  
899 hisat-genotype,” *Nature biotechnology*, vol. 37, no. 8, pp. 907–915,  
900 2019.
- 901 [72] H. Li, “Minimap2: pairwise alignment for nucleotide sequences,”  
902 *Bioinformatics*, vol. 34, no. 18, pp. 3094–3100, 2018.
- 903 [73] X. Chen, O. Schulz-Trieglaff, R. Shaw, B. Barnes, F. Schlesinger,  
904 M. Källberg, A. J. Cox, S. Kruglyak, and C. T. Saunders, “Manta:  
905 rapid detection of structural variants and indels for germline and  
906 cancer sequencing applications,” *Bioinformatics*, vol. 32, no. 8, pp.  
907 1220–1222, 2016.
- 908 [74] K. Ye, M. H. Schulz, Q. Long, R. Apweiler, and Z. Ning, “Pindel:  
909 a pattern growth approach to detect break points of large dele-  
910 tions and medium sized insertions from paired-end short reads,”  
911 *Bioinformatics*, vol. 25, no. 21, pp. 2865–2871, 2009.
- 912 [75] Z. H. Kwok, B. Zhang, X. H. Chew, J. J. Chan, V. Teh, H. Yang,  
913 D. Kappei, and Y. Tay, “Systematic analysis of intronic mirnas  
914 reveals cooperativity within the multicomponent ftx locus to pro-  
915 mote colon cancer development,” *Cancer Research*, vol. 81, no. 5,  
916 pp. 1308–1320, 2021.