

BIOMARS: A MULTI-AGENT ROBOTIC SYSTEM FOR AUTONOMOUS BIOLOGICAL EXPERIMENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) and vision-language models (VLMs) have the potential to transform biological research by enabling autonomous experimentation. Yet, their application remains constrained by rigid protocol design, limited adaptability to dynamic lab conditions, inadequate error handling, and high operational complexity. Here we introduce BioMARS (Biological Multi-Agent Robotic System), an intelligent platform that integrates LLMs, VLMs, and modular robotics to autonomously design, plan, and execute biological experiments. BioMARS uses a hierarchical architecture: the Biologist Agent synthesizes protocols via retrieval-augmented generation; the Technician Agent translates them into executable robotic pseudo-code; and the Inspector Agent ensures procedural integrity through multimodal perception and anomaly detection. The system autonomously conducts cell passaging and culture tasks, matching or exceeding manual performance in viability, consistency, and morphological integrity. It also supports context-aware optimization, outperforming conventional strategies in differentiating retinal pigment epithelial cells. A web interface enables real-time human-AI collaboration, while a modular backend allows scalable integration with laboratory hardware. These results highlight the feasibility of generalizable, AI-driven laboratory automation and the transformative role of language-based reasoning in biological research.

1 INTRODUCTION

The convergence of robotic automation and artificial intelligence is reshaping experimental biology, promising greater reproducibility, throughput, and independence from human variability [Holland & Davies \(2020\)](#). However, the complexity of biological protocols—which demand adaptive decision-making, multi-stage coordination, and interpretation of nuanced environmental feedback—has hindered the realization of fully autonomous systems. Existing automation solutions, ranging from specialized liquid handling robots [Dettinger et al. \(2022\)](#); [Novak et al. \(2020\)](#); [Taguchi et al. \(2023\)](#), to modular single-arm platforms for cell culture automation [Hamm et al. \(2024\)](#); [Tristan et al. \(2021\)](#), and dual-arm platforms enabling automated cell production [Königer et al. \(2024\)](#); [Yachie & Natsume \(2017\)](#); [Ochiai et al. \(2021\)](#), often require extensive manual oversight and lack the flexibility to navigate unanticipated procedural deviations. Early systems focused on streamlining specific tasks, including biofoundries [Chao et al. \(2017\)](#), IoT-enabled experimental platforms [Miles & Lee \(2018\)](#), and clinical sample preparation [Müller et al. \(2020\)](#), but these non-robotic arm systems still faced hardware limitations that prompted the development of robotic arm solutions.

Concurrently, large language models (LLMs) and vision-language models (VLMs) are transforming scientific problem-solving by enabling machines to parse literature, synthesize knowledge, and execute multi-modal reasoning across diverse domains [Vaswani et al. \(2017\)](#); [Wang et al. \(2024\)](#); [Luu & Buehler \(2024\)](#); [Zhang et al. \(2025\)](#). Recent efforts leveraging LLMs in chemical experimentation [Boiko et al. \(2023\)](#); [Darvish et al. \(2025\)](#); [Cooper et al. \(2025\)](#) and biological protocol generation [O'Donoghue et al. \(2023\)](#); [Huang et al. \(2024\)](#) signal a paradigm shift toward AI-native experimentation. Yet, their integration with physical robotic systems for biological execution remains underexplored.

Here we introduce BioMARS (Biological Multi-Agent Robotic System), a dual-arm robotic platform orchestrated by LLMs and VLMs [Zhu et al. \(2023\)](#); [Zhang et al. \(2024\)](#) for fully autonomous

054 execution of biological experiments. BioMARS performs end-to-end protocol design, environmen-
055 tal coordination, and robotic manipulation through adaptive multimodal reasoning. By converting
056 research literature into actionable procedures and coupling them with error-aware execution strate-
057 gies, the system ensures both flexibility and robustness in complex biological tasks.

058 We demonstrate BioMARS across five experimental capabilities: (1) efficiently searching and ana-
059 lyzing online research documentation to design experimental protocols for diverse cell types under
060 varying conditions; (2) accurately translating and executing these protocols using a dual-arm biolog-
061 ical laboratory; (3) detecting experimental errors via keyframe analysis; (4) performing end-to-end
062 cell culturing; and (5) resolving optimization issues through the analysis of historical experimental
063 data.

064 2 RELATED WORK

065 2.1 AUTOMATION IN LIFE SCIENCE RESEARCH LABORATORY

066 In recent years, laboratory automation technology has rapidly advanced in the life sciences, with
067 its core objectives being to enhance experimental efficiency and reproducibility through standard-
068 ized operational workflows and reduced human intervention. Initial studies have explored various
069 systems without robotic arms: the biofoundry system proposed by Chao et al. [Chao et al. \(2017\)](#)
070 for synthetic biology automation, the IoT-enabled closed-loop experimental platform developed by
071 Miles and Lee [Miles & Lee \(2018\)](#) which integrates real-time data feedback, and Müller’s team
072 [Müller et al. \(2020\)](#) that demonstrated automated clinical sample preparation. Additionally, Det-
073 tinger’s team [Dettinger et al. \(2022\)](#) developed an open-source pipetting robot featuring customiz-
074 able liquid-handling protocols. These systems laid down essential workflows that boosted experi-
075 mental efficiency while also highlighting limitations in hardware adaptability.

076 To address these constraints, researchers shifted focus to robotic-arm systems. Initial efforts centered
077 on single-arm platforms: Li et al. [Li et al. \(2023\)](#) pioneered a robotic-arm-based 3D bioprinting
078 system for tissue engineering. Precision in these systems was further refined through innovations
079 like Zhang et al.’s [Zhang et al. \(2023\)](#) vision-assisted pipetting method, which dynamically adjusted
080 for container deformation. Hamm et al. [Hamm et al. \(2024\)](#) extended single-arm capabilities through
081 modular designs optimized for automated cell culture workflows.

082 Nevertheless, inherent limitations in single-arm flexibility for multi-step coordination drove explo-
083 ration of dual-arm systems. Koniger’s team [Königer et al. \(2024\)](#) achieved a breakthrough with their
084 dual-arm platform enabling uninterrupted 3D epithelial tissue production through synchronized me-
085 chanical manipulation and fluidic control. Ochiai et al. [Ochiai et al. \(2021\)](#) advanced this paradigm
086 by developing a variable-scheduling system capable of parallel maintenance for eight mammalian
087 cell lines through dynamic task prioritization. These developments collectively signify a paradigm
088 shift in laboratory automation, transitioning from discrete single-function devices to integrated sys-
089 tems capable of executing complex experimental protocols through coordinated multi-arm opera-
090 tions.

091 2.2 LARGE MODEL BASED MULTI-AGENT SYSTEMS

092 With breakthroughs in cross-modal reasoning and task planning capabilities of large language mod-
093 els (LLMs) [Achiam et al. \(2023\)](#), researchers have begun exploring their application paradigms in
094 multi-agent collaboration. Early foundational concepts of agent coordination and communication
095 were proposed by Wooldridge and Jennings [Wooldridge & Jennings \(1995\)](#), while the emergence of
096 powerful LLMs has injected new vitality into this field. Brown et al. demonstrated the potential of
097 LLMs in simulating human dialogues, which can be applied to inter-agent communication [Brown
098 et al. \(2020\)](#). In terms of architectural design, Patil et al. [Patil et al. \(2025\)](#) proposed connecting mul-
099 tiple functional APIs via LLMs for task decomposition, while Cai et al. [Cai et al. \(2023\)](#) validated
100 LLMs’ potential as “tool generators” capable of dynamically creating specialized agents. For multi-
101 robot collaboration scenarios, Chen et al. [Chen et al. \(2024\)](#) compared the performance boundaries
102 of centralized versus distributed multi-robot systems, and the RoCo framework developed by Mandi
103 et al. [Mandi et al. \(2024\)](#) enables semantic-level conversational planning among multiple robots
104 through LLMs. Zheng et al. [Zheng et al. \(2023\)](#) utilized multi-agent LLMs to optimize material
105 crystal structures, Hua et al. [Hua et al. \(2023\)](#) constructed a war simulation multi-agent system, and
106
107

Aher et al. [Aher et al. \(2023\)](#) explored LLMs’ applications in replicating psychological experiments. These interdisciplinary advancements highlight the universal advantages of LLM-based multi-agent systems in complex decision-making scenarios.

2.3 LLMs AND VLMS AS ROBOT PLANNERS

In the field of robotic task planning, LLMs demonstrate value through two complementary pathways: (1) Semantic reasoning-based high-level instruction decomposition, exemplified by Song et al. [Song et al. \(2023\)](#) with their few-shot planning framework, and Silver et al. [Silver et al. \(2024\)](#) who leverage GPT-4 to generate PDDL domain models. (2) Direct generation of executable control code, as seen in Liang et al. [Liang et al. \(2023\)](#)’s ”code-as-policies” paradigm, Huang et al. [Huang et al. \(2023\)](#)’s 3D value map synthesis method, and Wu et al. [Wu et al. \(2023\)](#)’s implementation of personalized household robots.

To address the physical perception limitations of pure language modalities, vision-language models (VLMS) have been integrated into planning pipelines: Hu et al. [Hu et al. \(2023\)](#) developed a visual foresight mechanism that significantly reduces physical interaction error rates, while Shirai et al. [Shirai et al. \(2024\)](#) created a vision-language interpreter for dynamic task sequence adjustment. Mei et al. [Mei et al. \(2024\)](#) further enhanced this through their ReplanVLM system, achieving real-time plan corrections via multimodal feedback.

These advancements collectively improve robotic adaptability in dynamic environments. Notably, Gao et al. [Gao et al. \(2024\)](#) proposed a physically-grounded model that simultaneously processes visual signals and manipulation constraints, while Wake et al. [Wake et al. \(2024\)](#) validated GPT-4V’s multimodal planning capabilities in imitation learning scenarios.

3 ARCHITECTURE OF BIOMARS SYSTEM

BioMARS (Biological Multi-Agent Robotic System) enables end-to-end autonomous execution of biological experiments through a network of specialized LLM- and VLM-based agents (Fig. 1a). Built on an enhanced Agentic Retrieval-Augmented Generation (RAG) framework with modular error correction [Singh et al. \(2025\)](#), BioMARS decomposes complex protocols, interprets unstructured literature, and dynamically synthesizes findings into executable procedures.

The Biologist Agent ingests diverse open-access research documents, generates executable protocol steps by leveraging biological domain knowledge to create structured, constraint-aware queries. By incorporating constraints such as container type (e.g., petri dishes, flasks) and platform capacity, it tailors each protocol to the laboratory’s operational environment. The Technician Agent transforms high-level plans into fine-grained control primitives for robotic execution. These primitives are allocated across dual robotic arms and coordinated with environmental modules such as incubator and centrifuge.

To ensure execution robustness, the Inspector Agent—powered by ViTs and VLMS—performs rapid anomaly detection. It identifies procedural deviations including geometric misalignments (e.g., unattached pipette tips, misaligned petri dishes) and mechanical failures, prompting replanning or user notification. This tri-agent system mirrors the modularity and task specialization seen in other autonomous platforms [Boiko et al. \(2023\)](#); [M. Bran et al. \(2024\)](#), enabling BioMARS to operate adaptively under changing experimental conditions.

The platform supports natural language prompts (e.g., “How to passage HeLa cells”) via a web interface. Users can initiate, monitor, and modify experiments interactively. Critically, BioMARS’s modular architecture allows seamless integration of new hardware and protocol domains through programmable function modules, facilitating extensibility across diverse biological workflows.

3.1 PROTOCOL SYNTHESIS UNDER ENVIRONMENTAL CONSTRAINTS

Reliable generation of biological protocols from literature poses challenges due to procedural complexity, heterogeneous experimental conditions, instrumentation constraints, and output formatting requirements. BioMARS addresses this through a multi-agent reasoning framework that integrates

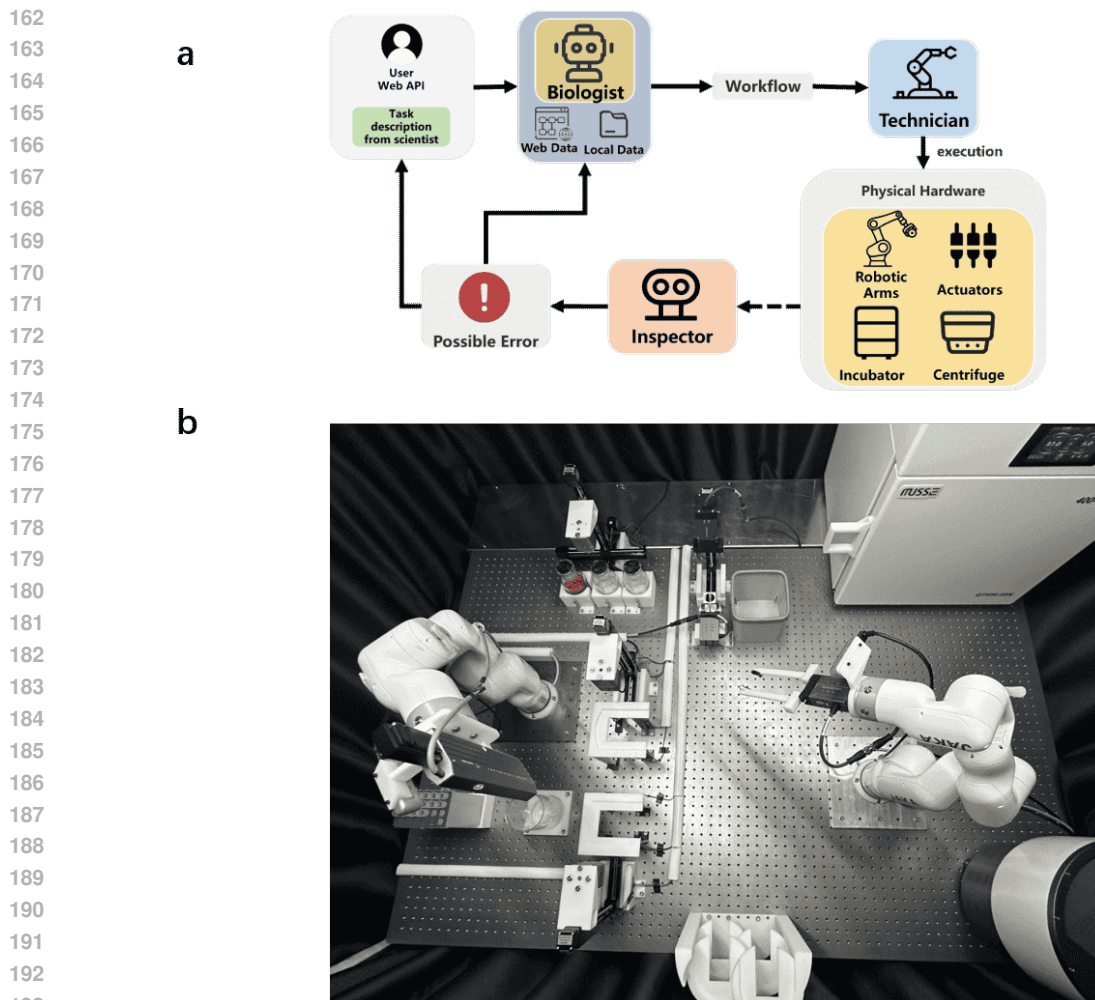


Figure 1: **System architecture and robotic setup.** **a**, Multi-agent workflow of BioMARS, comprising Biologist, Technician, and Inspector agents. **b**, Dual-arm robotic platform configured for autonomous biological experimentation.

LLM-based planning with vector-based retrieval and verification mechanisms to generate biologically accurate, context-aware procedures.

At the center of this system is the Biologist Agent, which operates within an enhanced Agentic Retrieval-Augmented Generation (RAG) architecture (Fig. 2a). The agent retrieves relevant knowledge using online query APIs (Google and Bing), extracting three PDFs and three high-relevance web snippets per query. Full paragraphs associated with each snippet are selected to preserve semantic context. These passages, along with the embedded user query (using OpenAI’s Ada model), undergo vector similarity ranking. The top five text chunks are used as context for downstream protocol generation.

Protocol construction is distributed across three sub-agents: the Knowledge Checker (KC), which filters domain-inconsistent content; the Workflow Generator (WG), which formulates stepwise procedures; and the Workflow Checker (WC), which iteratively refines outputs for logical coherence. The system accounts for laboratory constraints, such as limited stock of specific containers (e.g. 10 cm culture dishes), pipette tip volume (10 ml), and robotic station limits, ensuring that all outputs are executable on the BioMARS platform.

System performance was evaluated using a 70-query benchmark comprising 10 procedural categories across seven cell lines, ranging from routine tasks (e.g., cell passaging, thawing) to complex

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

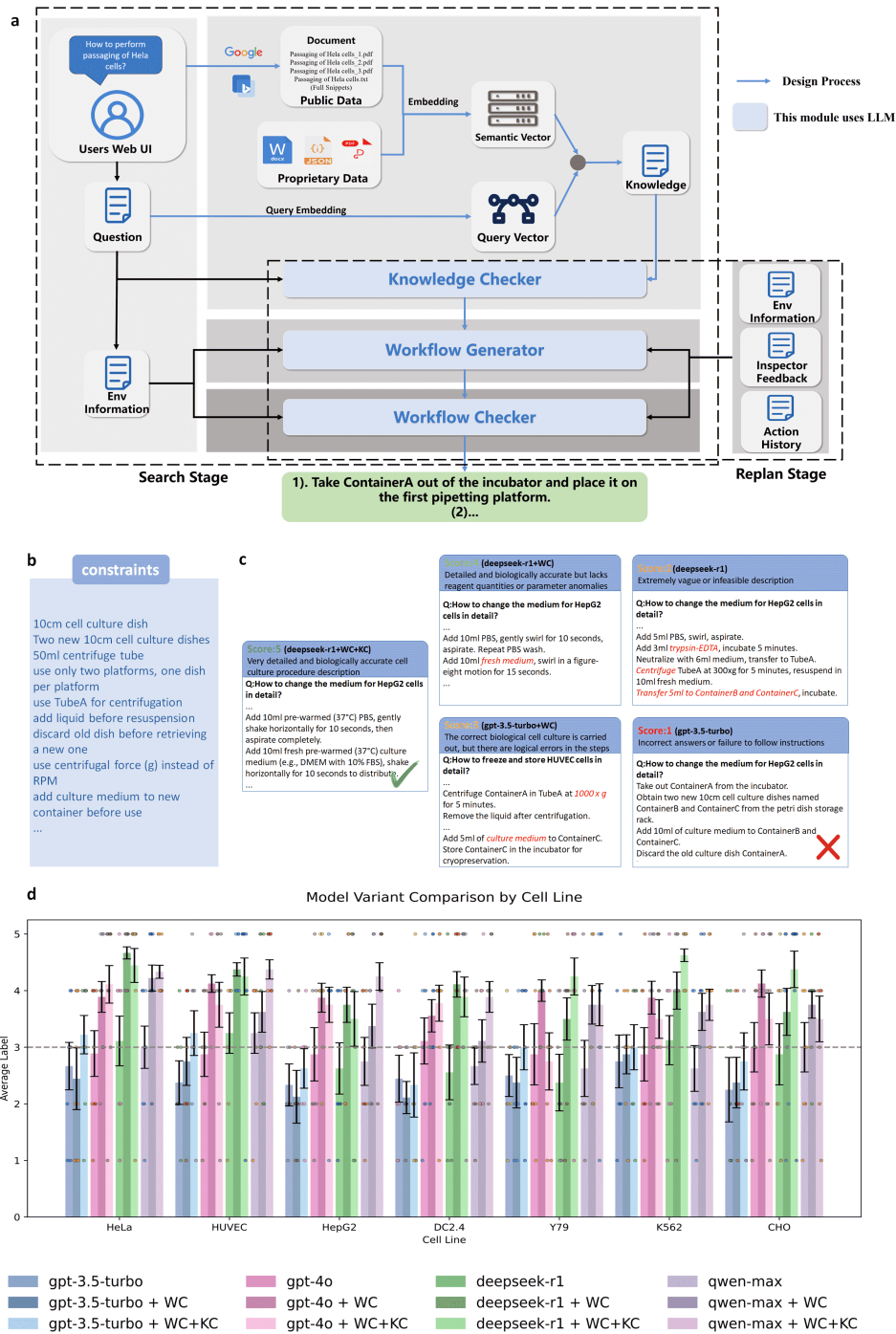


Figure 2: **Biologist Agent architecture and evaluation.** **a**, Biologist agent pipeline integrating document retrieval, semantic matching, and workflow refinement under constraints. **b**, Representative experimental constraints. **c**, Example protocol outputs with scores and errors. **d**, Performance comparison of four models (GPT-3.5-Turbo, GPT-4o, Deepseek-R1, Qwen-Max) and their variants on seven cell lines.

protocols (e.g., 3D culture, apoptosis analysis). Following Boiko et al. [Boiko et al. \(2023\)](#), model outputs were scored on a 5-point scale: 5 for fully detailed and accurate procedures; 4 for biologically sound steps with minor omissions; 3 for logically flawed but conceptually plausible outputs; 2 for vague or infeasible workflows; and 1 for incorrect or non-compliant procedures. Outputs below a score of 3 were considered task failures. Fig. 2c presents representative outputs with annotations.

Without WC or KC modules, base models—including GPT-4o, Qwen-Max, and DeepSeek-R1—did not exceed a mean score of 3. GPT-3.5 Turbo consistently underperformed; in one instance, it misinterpreted “How to change the HepG2 culture medium” by suggesting disposal of viable dishes and initiating culture from scratch (score: 1). DeepSeek-R1 proposed cell redistribution via trypsinization (score: 2), demonstrating procedural confusion.

Incorporation of the WC module significantly improved structural logic. For example, DeepSeek-R1+WC successfully outlined PBS rinsing and medium replacement steps but omitted critical conditions (temperature, CO₂ levels), yielding a score of 4. Further integration with the KC module provided domain-specific validations: in the cryopreservation task for HUVECs, KC-corrected protocols mitigated centrifugation errors and ensured cryostorage in liquid nitrogen.

The best-performing configuration—DeepSeek-R1+WC+KC—achieved consistent scores of 5. Its output for HepG2 medium replacement detailed exact reagent volumes, environmental settings (37°C, 5% CO₂), and handling protocols (PBS rinse with horizontal agitation), aligning closely with expert protocols. These results affirm the critical role of domain validation (KC) and procedural refinement (WC) in transforming LLM outputs into executable, high-fidelity biological protocols (Fig. 2c,d).

3.2 PROTOCOL-TO-CODE TRANSLATION FOR ROBOTICS

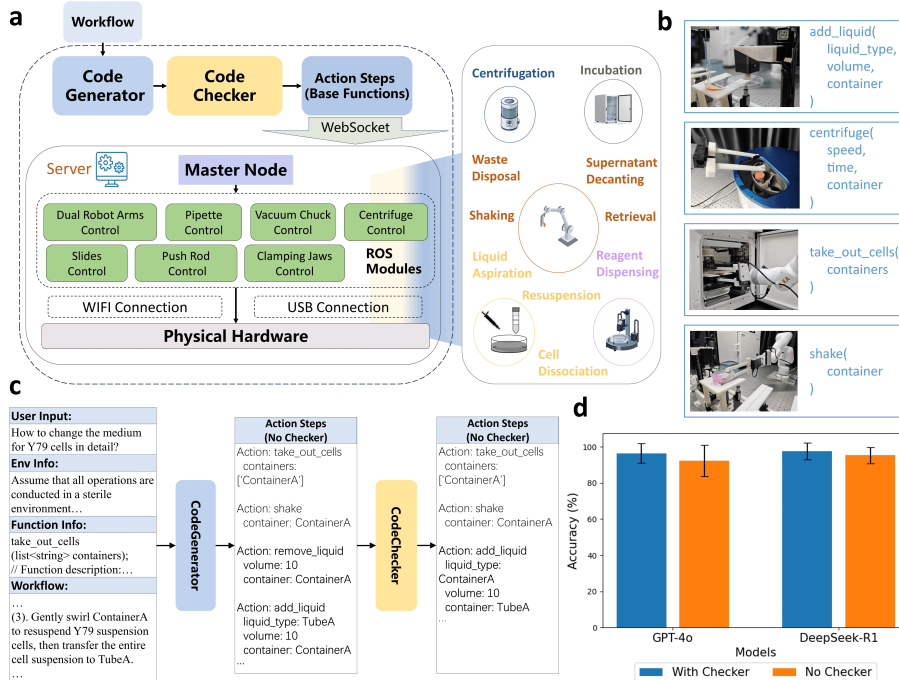


Figure 3: **Technician Agent architecture and performance of protocol translation and execution framework.** **a**, System workflow of Technician Agent, including CodeGenerator, CodeChecker, ROS node and the corresponding hardware module. **b**, Example pseudo-code instructions and corresponding robotic actions. **c**, The specific workflow of Technician Agent. **d**, Instruction accuracy comparison with and without CodeChecker for GPT-4o and DeepSeek-R1.

324 Translating free-text experimental protocols into executable robotic commands remains a central
325 bottleneck in laboratory automation. Existing systems typically rely on rigid, manually curated
326 command sequences [Kanda et al. \(2022\)](#); [Königer et al. \(2024\)](#), which limits their adaptability to
327 diverse and unstructured inputs. To address this constraint, we developed the Technician Agent—a
328 dual-module system that autonomously interprets natural language protocols and converts them into
329 validated robotic instructions.

330 The Technician Agent operates through a cooperative pipeline comprising a CodeGenerator and
331 a CodeChecker module (Fig. 3a). The CodeGenerator, powered by an LLM, maps protocol
332 descriptions into pseudo-code composed of primitive robotic operations such as `add_liquid`,
333 `centrifuge`, and `shake` (Fig. 3b). The CodeChecker subsequently performs rule-based vali-
334 dation, enforcing functional correctness and environmental compatibility based on the predefined
335 specification set.

336 The pipeline structure is illustrated in Fig. 3c. Given a protocol input, the CodeGenerator produces
337 candidate instructions tailored to the lab environment. These instructions are then parsed by the
338 CodeChecker, which applies logical and semantic checks including parameter validation, function
339 relevance, and argument structure. This ensures that all generated commands adhere to the opera-
340 tional and safety constraints of the BioMARS platform.

341 To assess performance, we benchmarked the Technician Agent across 300 experimental protocol
342 steps. As shown in Fig. 3d, the full pipeline (CodeGenerator + CodeChecker, GPT-4o) achieved a
343 96.4% instruction-matching accuracy, outperforming a single-module baseline (92.4%). The impact
344 of the CodeChecker module is particularly evident in complex procedural constructs. For example,
345 when parsing the instruction “resuspend the cell pellet in 10 mL fresh complete growth medium,” the
346 baseline failed to recognize the prerequisite transfer step. In contrast, the Technician Agent inserted
347 an implicit `add_liquid` operation before resuspension, preserving procedural logic.

348 Beyond resolving implicit steps, the CodeChecker module also corrects parameter mismatches, en-
349 forces range constraints, and eliminates superfluous instructions. For instance, it detects and corrects
350 overfilled volumes relative to container capacity and replaces invalid data types in function argu-
351 ments. This systematic refinement substantially improves the robustness of the robotic instruction
352 set.

353 By converting ambiguous natural language into explicit, verifiable pseudo-code, the Technician
354 Agent enhances experimental reproducibility, reduces human error, and simplifies execution on
355 robotic platforms. This capability shifts the experimental burden away from manual coding, en-
356 abling researchers to focus on scientific inquiry rather than operational encoding.

359 3.3 HIERARCHICAL VLM-BASED ERROR DETECTION

361 Biological experimentation demands strict precision, where minor procedural errors can compro-
362 mise outcomes. Conventional automation platforms typically rely on basic object detection without
363 semantic context awareness, limiting their robustness in dynamic laboratory environments [Jiang
364 et al. \(2022\)](#). To address this, we developed the Inspector Agent—a hierarchical visual monitor-
365 ing system integrating vision-language models (VLMs) and vision transformers (ViTs) [Han et al.
366 \(2022\)](#) for multi-stage perception and error detection (Fig. 4a).

367 The first stage performs visual segmentation of experimental scenes using few-shot prompting with
368 a VLM. Key objects—such as pipette tips, culture plates, and tubes—are segmented from raw RGB
369 inputs. To enhance spatial resolution and minimize background interference, the bounding boxes
370 generated by the VLM are manually refined. These cropped subregions are converted to grayscale,
371 preserving structural cues like pipette orientation and tube angles while reducing color-based noise.

372 In the second stage, a ViT-based keyframe detection module encodes 23 visually discriminative ac-
373 tions (selected from 11 control primitives) into a reference embedding library. This module enables
374 sub-second recognition of procedural steps. In benchmark tests, ViT achieved a mean inference
375 latency of 0.3066 s - 91.9% faster than GPT-4o (3.7960 s) - with lower temporal variability (co-
376 efficient of variation: 13.08% vs. 42.80%; Fig. 4d). In real-world experimental settings, the ViT
377 achieved an F1 score of 88.7% and a recall of 94.0%, demonstrating high temporal stability and
operational fidelity (Fig. 4c).

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

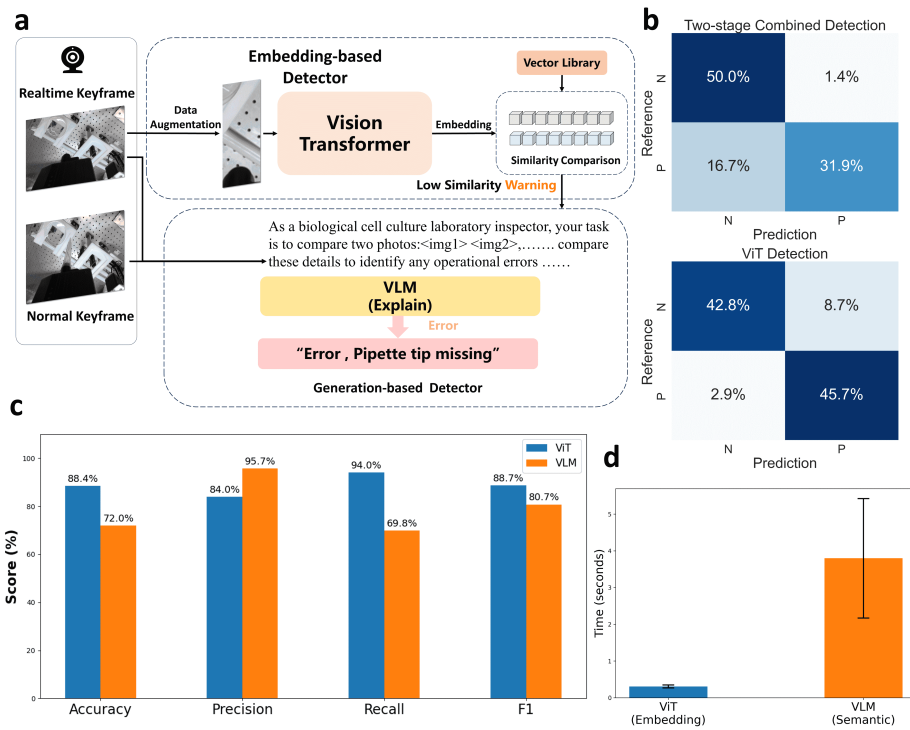


Figure 4: Inspector Agent Overview and Performance Metrics. **a**, Workflow diagram of the Technician Agent. **b**, Confusion matrix of two-stage combined detection and ViT detection. **c**, Performance of ViT and VLM on four evaluation metrics. **d**, The time performance of the two detection methods, ViT and VLM.

432 A final stage introduces zero-shot semantic validation using the VLM. When the ViT flags anoma-
433 lies, frames are semantically compared with idealized keyframes using language-guided prompts
434 (e.g., “attach pipette tip”). This semantic differential analysis enables detection of contextual errors
435 beyond geometry alone. In validation, this mechanism achieved 95.7% precision and 80.7% F1
436 score(Fig. 4c), reducing the false positives rate from 8.7% to 1.4%—an 83% improvement (Fig. 4b).
437 For example, detecting a detached pipette tip without a visible pipette is correctly flagged as an
438 action violation. Upon confirmation, robotic operations are automatically paused and visual alerts
439 issued.

440 By integrating geometric and semantic vision processing, the Inspector Agent ensures procedural
441 robustness, accelerates feedback response times, and significantly reduces downstream execution
442 failures.

444 4 BIOLOGICAL EXPERIMENT DESIGN

447 4.1 INTEGRATED BIOLOGICAL EXPERIMENT DESIGN

448 To evaluate the biological reliability and operational efficiency of BioMARS, we conducted a com-
449 parative study between automated and manual cell passaging protocols across three representative
450 cell types: HeLa (adherent), Y79 (suspension), and DC2.4 (semi-adherent/suspension). Experi-
451 mental evaluation included metabolic viability, survival consistency, morphological preservation,
452 and coefficient of variation (CV) analysis. All workflows adhered to established protocols, with
453 BioMARS dynamically adapting process parameters to each cell line.

454 Cells were cultured in standard media: HeLa in DMEM with 10% FBS and 1% peni-
455 cillin–streptomycin, Y79 in RPMI-1640 with 20% FBS and DC2.4 in RPMI-1640 with 10% FBS -
456 under 5% CO₂ at 37 ° C. Media changes were performed every 2–3 days. For passaging, adherent
457 cells were detached with 0.25% trypsin–EDTA. The BioMARS system adjusted enzymatic diges-
458 tion time and centrifugation based on cell type: 6 minutes for HeLa and 3 minutes for Y79, ensuring
459 optimal yield and viability.

460 Metabolic viability was assessed 48 hours post-passaging using the CCK-8 assay. Optical density
461 (OD) measurements showed no significant difference between BioMARS and manual protocols
462 across all three cell types (Fig. 5d), indicating that automated processing maintained normal cellular
463 proliferation. CV analysis revealed enhanced reproducibility in the BioMARS group: HeLa and
464 Y79 samples exhibited 12–18% lower variability compared to manual handling (Fig. 5e).

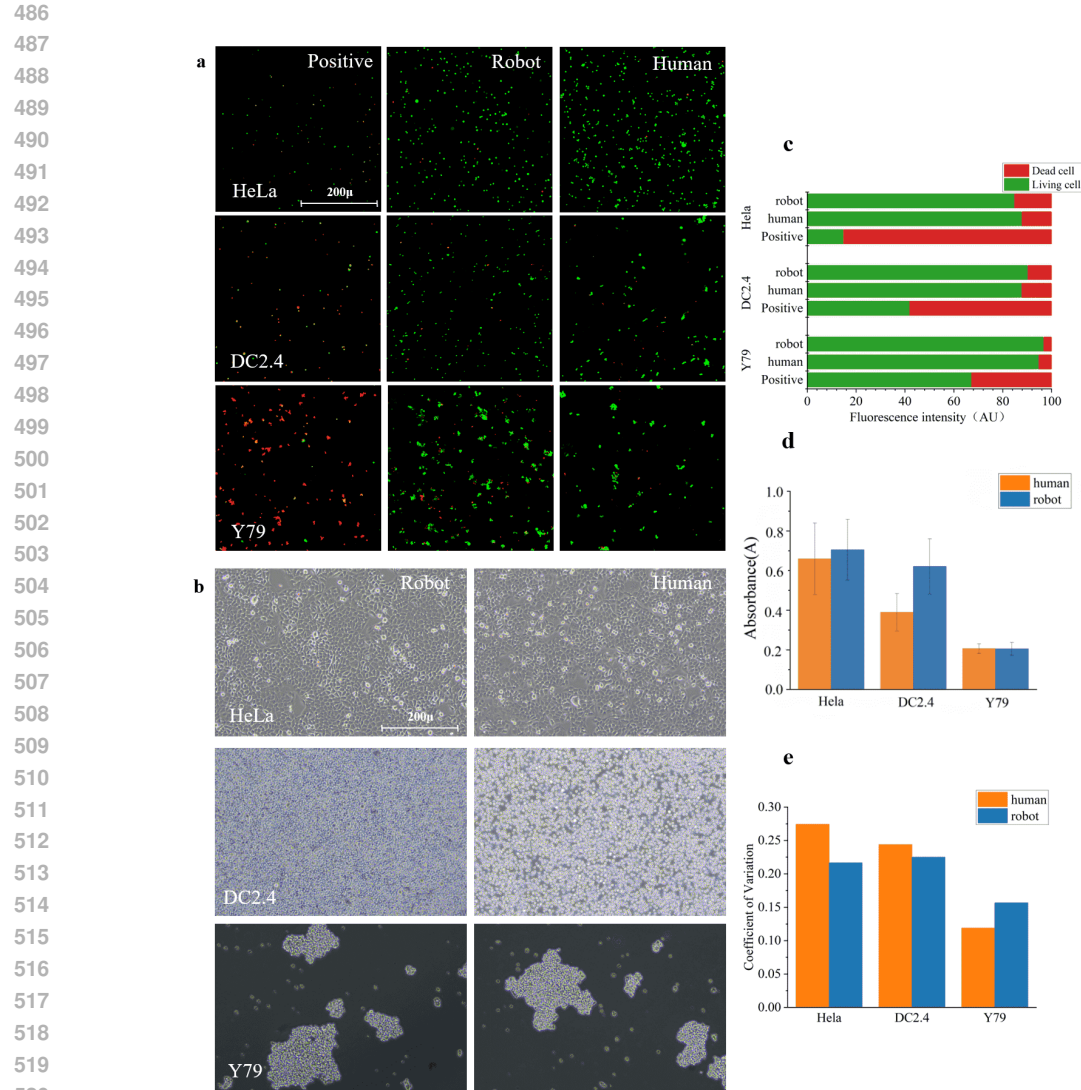
465 Live/dead staining confirmed high post-passaging viability, with over 92% concordance between
466 automated and manual groups (Fig. 5a,b). Green fluorescence indicated dominant live-cell popu-
467 lations, with clear contrast to the red-stained positive control. Morphological evaluation (Fig. 5c)
468 showed no detectable structural abnormalities, further confirming the BioMARS system’s ability to
469 preserve cell integrity.

470 In addition to biological fidelity, BioMARS markedly improved operational efficiency. Manual
471 passaging required approximately 60 minutes per cell line, whereas the BioMARS system reduced
472 hands-on time to 5–8 minutes—representing a 90% reduction. This time savings translates into
473 higher throughput and improved standardization, minimizing human error and procedural variability.

474 Collectively, these results establish that BioMARS performs comparably or superior to manual pro-
475 tocols in biological outcome metrics while offering significant gains in consistency, reproducibility,
476 and efficiency.

478 4.2 BIOLOGICAL OPTIMIZING CAPABILITY

480 Beyond static protocol generation, the Biologist Agent was evaluated for its capacity to perform
481 biological optimization—an advanced task requiring iterative reasoning, mechanistic understanding,
482 and strategic parameter adjustment. We assessed this capability using a publicly available dataset
483 for optimizing differentiation efficiency of induced pluripotent stem cell-derived retinal pigment
484 epithelial (iPSC-RPE) cells [Kanda et al. \(2022\)](#), which defines a high-dimensional experimental
485 space grounded in biological constraints.



522 **Figure 5: Comparison of automated vs. manual cell passaging outcomes.** **a**, Fluorescence images of live/dead-stained cells (automated vs. manual) at 48 h post-passaging. **b**, Bright-field images of cell morphology post-passaging. **c**, Live/dead cell ratio comparison after passaging. **d**, Cell viability comparison between methods. **e**, CV of CCK-8 viability across repeats (reproducibility).

523
524
525
526
527
528

529 The optimization target was the pigment score, a key phenotypic marker of iPSC-RPE maturation. Seven tunable parameters were considered across preconditioning, detachment, and differentiation stages: FGFR1 concentration (PC: 0–505 nM) and exposure duration (PP: 1–6 days); trypsin incubation time (DP: 5–23 min), pipetting strength (DS: 10–100 mm / s) and pipetting length (DL: short/long); KSR withdrawal schedule (KP: 1–19 days); and three-supplement exposure duration (3P: 3–19 days). This setup presents a biologically grounded, combinatorially complex optimization challenge.

530
531
532
533
534
535
536
537
538
539

To simulate realistic experimental conditions, optimization was constrained to 20 iterations, initialized from 10 randomly selected low-performing conditions (pigment score ≤ 0.6). Parameter selection used KDTree-based nearest-neighbor interpolation [Friedman et al. \(1977\)](#), with outputs formatted in structured JSON for reproducibility. We compared three strategies: DeepSeek-R1, GPT-4o, and Bayesian optimization under identical initialization settings.

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

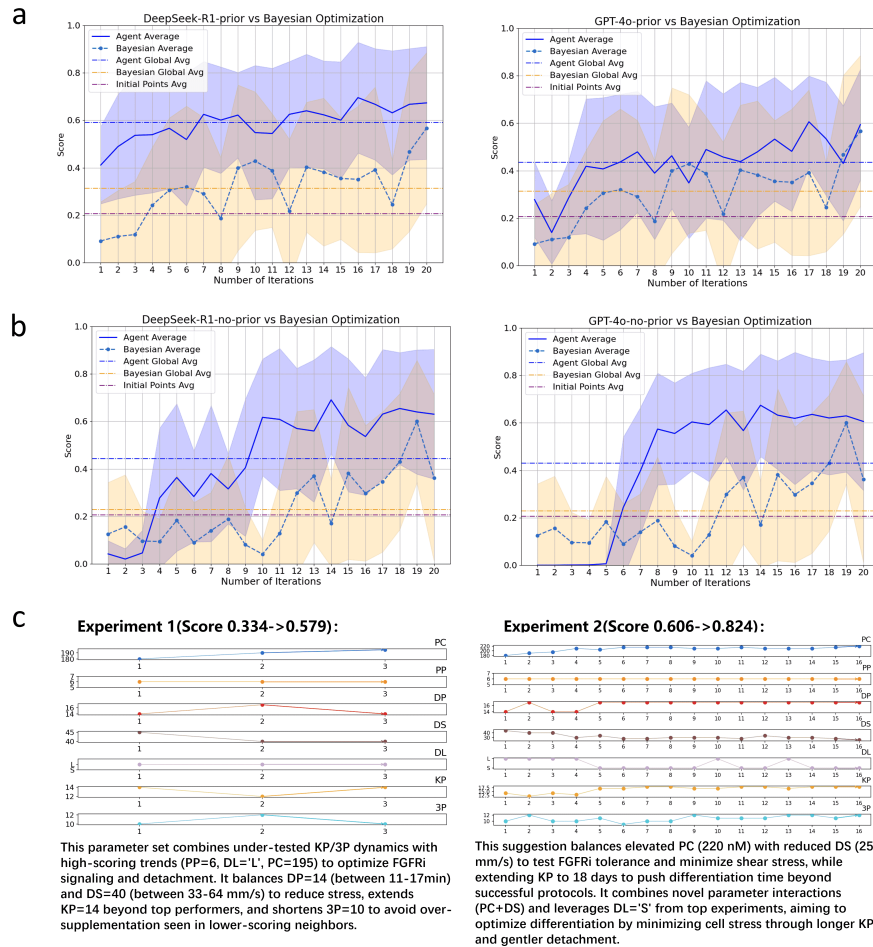


Figure 6: **Results of the iPSC-RPE optimization experiments.** **a**, Performance comparison between DeepSeek-R1 and GPT-4o models versus Bayesian Optimization using 10 prior experimental data points. **b**, Comparative analysis of DeepSeek-R1 and GPT-4o against Bayesian Optimization without leveraging prior experimental knowledge. **c**, Parameter recommendations from LLM-based optimizers across selected iteration rounds.

In the prior-informed setting (Fig. 6a), DeepSeek-R1 consistently outperformed baselines, reaching a final pigment score of 0.5913—surpassing GPT-4o (0.4344; +15.8%) and Bayesian optimization (0.3130; +28.5%). By iteration 7, it achieved 0.6252 and continued steady improvement. GPT-4o plateaued at 0.606, while Bayesian optimization peaked early at 0.5671. DeepSeek-R1’s advantage stems from its ability to encode mechanistic constraints; for instance, in one high-scoring trial (Fig. 6c), it selected PC = 220 nM (balancing efficacy and toxicity), DS = 25 mm/s (minimizing shear stress), and KP = 18 days (prolonging Wnt signaling), reflecting domain-consistent reasoning.

GPT-4o occasionally produced viable configurations but lacked consistent convergence, likely due to reliance on pretrained heuristics. Bayesian optimization, devoid of biological priors, frequently proposed implausible combinations (e.g., PC = 405.17 nM; KP = 2 days), resulting in limited progress.

In the no-prior setting (Fig. 6b), DeepSeek-R1 again demonstrated robust generalization, reaching performance comparable to the prior-informed case. GPT-4o improved after iteration 8, ultimately reaching a moderate score of 0.6303. Bayesian optimization showed minimal learning, with scores remaining near baseline. DeepSeek-R1 also exhibited superior balance between exploration and

594 exploitation, as evidenced by a lower standard deviation in output scores (0.2366 vs. 0.2447 for
595 GPT-4o and 0.2785 for Bayesian optimization), enabling more stable convergence.

596
597 These results validate the potential of knowledge-integrated LLMs to optimize complex biological
598 systems under data-sparse conditions. By combining contextual reasoning with structured decision-
599 making, such agents reduce dependency on manual tuning and offer scalable solutions for exper-
600 imental design. Future directions include reinforcement learning frameworks to further enhance
601 adaptive feedback integration in regenerative biology workflows.

602 603 5 DISCUSSION

604
605 This study introduces BioMARS, an intelligent agent system driven by LLMs and VLMs, capable of
606 autonomously designing, planning and executing biological experiments. By integrating language-
607 driven reasoning with multimodal perception and robotic control, BioMARS addresses the procedu-
608 ral complexity of biological workflows and generates reproducible, high-quality outcomes. These
609 capabilities are enabled by granting LLMs and VLMs access to essential research tools, including
610 scientific literature, programming environments and robotic execution platforms. The development
611 of such integrated AI systems holds substantial promise for accelerating discovery in the life sci-
612 ences.

613 While BioMARS demonstrates robust performance in standard cell culture tasks, several challenges
614 remain. Its operation under atypical or highly customized experimental conditions is limited, with
615 occasional human oversight required for critical steps such as pipetting volumes and centrifugation
616 parameters. Furthermore, although BioMARS integrates multimodal reasoning to interpret and
617 execute experimental protocols, its dependence on existing online procedures limits its capacity
618 for adaptive parameter tuning across diverse laboratory contexts. Its responsiveness to unexpected
619 experimental deviations is also limited, as real-time judgment remains an open challenge. Ongoing
620 efforts focus on enhancing the system’s adaptability and fault tolerance using advanced learning
621 algorithms, with preliminary improvements observed.

622 BioMARS represents a substantial step toward scalable, reproducible automation in biological re-
623 search. By addressing key barriers in protocol interpretation and execution, it lays the groundwork
624 for more reliable, flexible and scalable research practices. Its ability to ensure consistent procedural
625 replication enhances reproducibility and quality control, both of which are critical in applications
626 such as drug discovery and cell model production. In parallel, automation reduces operational bur-
627 den, with notable gains in time, material efficiency and labor reduction. Given that labor constitutes
628 a major cost component in biological production, deployment of BioMARS offers the potential to
629 significantly lower operational expenses. This economic benefit is expected to increase proportion-
630 ally with production scale.

631 REFERENCES

632
633 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
634 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
635 report. *arXiv preprint arXiv:2303.08774*, 2023.

636
637 Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate
638 multiple humans and replicate human subject studies. In *International Conference on Machine*
639 *Learning*, pp. 337–371. PMLR, 2023.

640
641 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research
642 with large language models. *Nature*, 624(7992):570–578, 2023.

643
644 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
645 Arvind Neelakantan, Pranav Shyam, Girish Sastry, and Amanda Askell. Language models are
646 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

647
648 Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as
649 tool makers. *arXiv preprint arXiv:2305.17126*, 2023.

- 648 Ran Chao, Shekhar Mishra, Tong Si, and Huimin Zhao. Engineering biological systems using
649 automated biofoundries. *Metabolic Engineering*, 42:98–108, 2017.
- 650
- 651 Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas Roy, and Chuchu Fan. Scalable multi-robot
652 collaboration with large language models: Centralized or decentralized systems? In *2024 IEEE*
653 *International Conference on Robotics and Automation (ICRA)*, pp. 4311–4317. IEEE, 2024.
- 654
- 655 Andrew Cooper, Zhengxue Zhou, Satheeshkumar Veeramani, Francisco Galeano, and Hatem
656 FakhruLdeen. Lira: Localization, inspection, and reasoning module for autonomous workflows
657 in self-driving labs. 2025.
- 658 Kourosh Darvish, Marta Skreta, Yuchi Zhao, Naruki Yoshikawa, Sagnik Som, Miroslav Bogdanovic,
659 Yang Cao, Han Hao, Haoping Xu, Alán Aspuru-Guzik, et al. Organa: a robotic assistant for
660 automated chemistry experimentation and characterization. *Matter*, 8(2), 2025.
- 661 Philip Dettinger, Tobias Kull, Geethika Arekatla, Nouraiz Ahmed, Yang Zhang, Florin Schneider,
662 Arne Wehling, Daniel Schirmacher, Shunsuke Kawamura, Dirk Loeffler, et al. Open-source per-
663 sonal pipetting robots with live-cell incubation and microscopy compatibility. *Nature Communi-*
664 *cations*, 13(1):2999, 2022.
- 665
- 666 Jerome H Friedman, Jon Louis Bentley, and Raphael Ari Finkel. An algorithm for finding best
667 matches in logarithmic expected time. *ACM Transactions on Mathematical Software (TOMS)*, 3
668 (3):209–226, 1977.
- 669 Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and
670 Dorsa Sadigh. Physically grounded vision-language models for robotic manipulation. In *2024*
671 *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 12462–12469. IEEE,
672 2024.
- 673
- 674 Jungmin Hamm, Seonghyeon Lim, Jiae Park, Jiwon Kang, Injun Lee, Yoongeun Lee, Jiseok Kang,
675 Youngjun Jo, Jaemin Lee, Seoyeong Lee, et al. A modular robotic platform for biological re-
676 search: Cell culture automation and remote experimentation. *Advanced Intelligent Systems*, 6(5):
677 2300566, 2024.
- 678 Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang,
679 An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on*
680 *pattern analysis and machine intelligence*, 45(1):87–110, 2022.
- 681
- 682 Ian Holland and Jamie A Davies. Automation in the life science research laboratory. *Frontiers in*
683 *bioengineering and biotechnology*, 8:571777, 2020.
- 684
- 685 Yingdong Hu, Fanqi Lin, Tong Zhang, Li Yi, and Yang Gao. Look before you leap: Unveiling the
686 power of gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- 687
- 688 Wenye Hua, Lizhou Fan, Lingyao Li, Kai Mei, Jianchao Ji, Yingqiang Ge, Libby Hemphill, and
689 Yongfeng Zhang. War and peace (waragent): Large language model-based multi-agent simulation
of world wars. *arXiv preprint arXiv:2311.17227*, 2023.
- 690
- 691 Kaixuan Huang, Yuanhao Qu, Henry Cousins, William A Johnson, Di Yin, Mihir Shah, Denny
692 Zhou, Russ Altman, Mengdi Wang, and Le Cong. Crispr-gpt: An llm agent for automated design
693 of gene-editing experiments. *arXiv preprint arXiv:2404.18021*, 2024.
- 694
- 695 Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer:
696 Composable 3d value maps for robotic manipulation with language models. *arXiv preprint*
arXiv:2307.05973, 2023.
- 697
- 698 Peiyuan Jiang, Daji Ergu, Fangyao Liu, Ying Cai, and Bo Ma. A review of yolo algorithm develop-
699 ments. *Procedia computer science*, 199:1066–1073, 2022.
- 700
- 701 Genki N Kanda, Taku Tsuzuki, Motoki Terada, Noriko Sakai, Naohiro Motozawa, Tomohiro Ma-
suda, Mitsuhiro Nishida, Chihaya T Watanabe, Tatsuki Higashi, Shuhei A Horiguchi, et al.
Robotic search for optimal cell culture in regenerative medicine. *Elife*, 11:e77007, 2022.

- 702 Lukas Königer, Christoph Malkmus, Dalia Mahdy, Thomas Däullary, Susanna Götz, Thomas
703 Schwarz, Marius Gensler, Niklas Pallmann, Danjouma Cheufou, Andreas Rosenwald, et al. Re-
704 bia—robotic enabled biological automation: 3d epithelial tissue production. *Advanced Science*,
705 11(45):2406608, 2024.
- 706 Kai Li, WenHui Huang, HaiTao Guo, YanYan Liu, Shuxian Chen, Heng Liu, and Qi Gu. Ad-
707 vancements in robotic arm-based 3d bioprinting for biomedical applications. *Life Medicine*, 2(6):
708 lnad046, 2023.
- 709 Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and
710 Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE*
711 *International Conference on Robotics and Automation (ICRA)*, pp. 9493–9500. IEEE, 2023.
- 712 Rachel K Luu and Markus J Buehler. Bioinspiredllm: Conversational large language model for the
713 mechanics of biological and bio-inspired materials. *Advanced Science*, 11(10):2306724, 2024.
- 714 Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe
715 Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelli-*
716 *gence*, 6(5):525–535, 2024.
- 717 Zhao Mandi, Shreeya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large
718 language models. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*,
719 pp. 286–299, 2024. doi: 10.1109/ICRA57147.2024.10610855.
- 720 Aoran Mei, Guo-Niu Zhu, Huaxiang Zhang, and Zhongxue Gan. Replanvlm: Replanning robotic
721 tasks with visual language models. *IEEE Robotics and Automation Letters*, 2024.
- 722 Ben Miles and Peter L Lee. Achieving reproducibility and closed-loop automation in biological
723 experimentation with an iot-enabled lab of the future. *SLAS TECHNOLOGY: Translating Life*
724 *Sciences Innovation*, 23(5):432–439, 2018.
- 725 Torsten Müller, Mathias Kalxdorf, Rémi Longuespée, Daniel N Kazdal, Albrecht Stenzinger, and
726 Jeroen Krijgsveld. Automated sample preparation with sp 3 for low-input clinical proteomics.
727 *Molecular systems biology*, 16(1):e91111, 2020.
- 728 Richard Novak, Miles Ingram, Susan Clauson, Debarun Das, Aaron Delahanty, Anna Herland,
729 Ben M Maoz, Sauveur SF Jeanty, Mahadevabharath R Somayaji, Morgan Burt, et al. A robotic
730 platform for fluidically-linked human body-on-chips experimentation. *Nature biomedical engi-*
731 *neering*, 4(4):407, 2020.
- 732 Koji Ochiai, Naohiro Motozawa, Motoki Terada, Takaaki Horinouchi, Tomohiro Masuda, Taku
733 Kudo, Motohisa Kamei, Akitaka Tsujikawa, Kenji Matsukuma, Tohru Natsume, et al. A variable
734 scheduling maintenance culture platform for mammalian cells. *SLAS TECHNOLOGY: Translat-*
735 *ing Life Sciences Innovation*, 26(2):209–217, 2021.
- 740 Odhran O’Donoghue, Aleksandar Shtedritski, John Ginger, Ralph Abboud, Ali Essa Ghareeb, Justin
741 Booth, and Samuel G Rodriques. Bioplanner: automatic evaluation of llms on protocol planning
742 in biology. *arXiv preprint arXiv:2310.10632*, 2023.
- 743 Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model
744 connected with massive apis. *Advances in Neural Information Processing Systems*, 37:126544–
745 126565, 2025.
- 746 Keisuke Shirai, Cristian C Beltran-Hernandez, Masashi Hamaya, Atsushi Hashimoto, Shohei
747 Tanaka, Kento Kawaharazuka, Kazutoshi Tanaka, Yoshitaka Ushiku, and Shinsuke Mori. Vision-
748 language interpreter for robot task planning. In *2024 IEEE International Conference on Robotics*
749 *and Automation (ICRA)*, pp. 2051–2058. IEEE, 2024.
- 750 Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B Tenenbaum, Leslie Kaelbling, and Michael
751 Katz. Generalized planning in pddl domains with pretrained large language models. In *Proceed-*
752 *ings of the AAAI conference on artificial intelligence*, volume 38, pp. 20256–20264, 2024.
- 753 Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. Agentic retrieval-augmented
754 generation: A survey on agentic rag. *arXiv preprint arXiv:2501.09136*, 2025.

- 756 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
757 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
758 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2998–3009, 2023.
759
- 760 Shodai Taguchi, Yasuyuki Suda, Kenji Irie, and Haruka Ozaki. Automation of yeast spot assays
761 using an affordable liquid handling robot. *SLAS technology*, 28(2):55–62, 2023.
- 762 Carlos A Tristan, Pinar Ormanoglu, Jaroslav Slamecka, Claire Malley, Pei-Hsuan Chu, Vukasin M
763 Jovanovic, Yeliz Gedik, Yogita Jethmalani, Charles Bonney, Elena Barnaeva, et al. Robotic high-
764 throughput biomanufacturing and functional differentiation of human pluripotent stem cells. *Stem*
765 *Cell Reports*, 16(12):3076–3092, 2021.
- 766 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
767 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
768 *tion processing systems*, 30, 2017.
- 769 Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. Gpt-4v
770 (ision) for robotics: Multimodal task planning from human demonstration. *IEEE Robotics and*
771 *Automation Letters*, 2024.
- 772 Chengrui Wang, Qingqing Long, Meng Xiao, Xunxin Cai, Chengjun Wu, Zhen Meng, Xuezhi Wang,
773 and Yuanchun Zhou. Biorag: A rag-llm framework for biological question reasoning. *arXiv*
774 *preprint arXiv:2408.01107*, 2024.
- 775 Michael Wooldridge and Nicholas Jennings. Intelligent agents: theory and practice. *The Knowledge*
776 *Engineering Review*, 10(2):115–152, 6 1995. doi: 10.1017/s0269888900008122.
- 777 Jimmy Wu, Rika Antonova, Adam Kan, Marion Lepert, Andy Zeng, Shuran Song, Jeannette Bohg,
778 Szymon Rusinkiewicz, and Thomas Funkhouser. Tidybot: Personalized robot assistance with
779 large language models. *Autonomous Robots*, 47(8):1087–1102, 2023.
- 780 Nozomu Yachie and Tohru Natsume. Robotic crowd biology with maholo labdroids. *Nature biotech-*
781 *nology*, 35(4):310–312, 2017.
- 782 Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks:
783 A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- 784 Junbo Zhang, Weiwei Wan, Nobuyuki Tanaka, Miki Fujita, Koichi Takahashi, and Kensuke Harada.
785 Integrating a pipette into a robot manipulator with uncalibrated vision and tcp for liquid handling.
786 *IEEE Transactions on Automation Science and Engineering*, 21(4):5503–5522, 2023.
- 787 Qiang Zhang, Keyan Ding, Tianwen Lv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao
788 Wang, Xiaotong Li, Zhuoyi Xiang, et al. Scientific large language models: A survey on biological
789 & chemical domains. *ACM Computing Surveys*, 57(6):1–38, 2025.
- 790 Zhiling Zheng, Oufan Zhang, Ha L Nguyen, Nakul Rampal, Ali H Alawadhi, Zichao Rong, Teresa
791 Head-Gordon, Christian Borgs, Jennifer T Chayes, and Omar M Yaghi. Chatgpt research group
792 for optimizing the crystallinity of mofs and cofs. *ACS Central Science*, 9(11):2161–2170, 2023.
- 793 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: En-
794 hancing vision-language understanding with advanced large language models. *arXiv preprint*
795 *arXiv:2304.10592*, 2023.
- 800

801
802
803
804
805
806
807
808
809