

# Autonomous Quantum Simulation through Large Language Model Agents

Weitang Li,<sup>1,\*</sup> Jiajun Ren,<sup>2</sup> Lixue Cheng,<sup>3</sup> and Cunxi Gong<sup>1</sup>

<sup>1</sup>*Guangdong Basic Research Center of Excellence for Aggregate  
Science, School of Science and Engineering, The Chinese University  
of Hong Kong, Shenzhen, Shenzhen, Guangdong, 518172, P.R. China*

<sup>2</sup>*MOE Key Laboratory of Theoretical and Computational Photochemistry,  
College of Chemistry, Beijing Normal University, Beijing, 100875, P. R. China*

<sup>3</sup>*Department of Chemistry, The Hong Kong University  
of Science and Technology, Hong Kong, P. R. China*

(Dated: January 16, 2026)

## Abstract

We demonstrate that large language model (LLM) agents can autonomously perform tensor network simulations of quantum many-body systems, achieving approximately 90% success rate across representative benchmark tasks. Tensor network methods are powerful tools for quantum simulation, but their effective use requires expertise typically acquired through years of graduate training. By combining in-context learning with curated documentation and multi-agent decomposition, we create autonomous AI agents that can be trained in specialized computational domains within minutes. We benchmark three configurations (baseline, single-agent with in-context learning, and multi-agent with in-context learning) on problems spanning quantum phase transitions, open quantum system dynamics, and photochemical reactions. Systematic evaluation using DeepSeek-V3.2, Gemini 2.5 Pro, and Claude Opus 4.5 demonstrates that both in-context learning and multi-agent architecture are essential. Analysis of failure modes reveals characteristic patterns across models, with the multi-agent configuration substantially reducing implementation errors and hallucinations compared to simpler architectures.

## I. INTRODUCTION

Large language models (LLMs) have emerged as transformative tools for scientific research [1–5]. Models such as GPT-4 [6], Claude, Gemini [7], and the open-weight DeepSeek model [8] demonstrate remarkable capabilities across diverse scientific tasks [9, 10]. When deployed as autonomous agents with external tools, memory, and planning abilities [11], these models become substantially more powerful. Tool augmentation through approaches like Toolformer [12] and retrieval-augmented generation [13] equips LLMs with external capabilities while reducing hallucinations. Building on decades of work toward self-driving laboratories [14–16], LLM-driven agents have achieved notable successes across scientific domains. Examples include ChemCrow, which outperforms GPT-4 alone on chemical research tasks [17]; Coscientist, which autonomously executes chemical syntheses [18, 19]; and El Agente Q, which achieves 88% success on quantum chemistry exercises [20]. Other agents have advanced quantum chemistry accessibility [21, 22], organic semiconductor optimization [23], single-cell transcriptomics [24], autonomous microscopy [25], and nanobody

---

\* liwt31@gmail.com

design [26].

Tensor network methods represent another powerful paradigm in computational science [27, 28]. By exploiting entanglement structure, tensor networks overcome the exponential barrier that prohibits direct classical simulation of quantum many-body systems. Over three decades, algorithmic development from White’s density matrix renormalization group (DMRG) [29, 30] and matrix product states (MPS) [31] to tree tensor networks [32–34] and time-evolution methods like the time-dependent variational principle (TDVP) [35–38] has produced a rich ecosystem of methods tailored to quantum chemistry [39–41], condensed matter physics [42, 43], and quantum computing [44, 45]. This diversity is reflected in more than 30 actively maintained software packages [46]. The proliferation of specialized tools underscores both the versatility and the complexity of tensor network methods, as effective use requires not only selecting an appropriate package but also mastering its particular conventions for network structures [47, 48], bond dimension control [49, 50], symmetry implementations [51–53], and observable computation.

Despite the individual strengths of LLMs and tensor networks, applying LLMs to automate tensor network simulations presents substantial challenges. Tensor network simulation requires simultaneous reasoning across three closely integrated layers: the underlying physics, code implementation, and numerical data. This physics-code-data coupling creates three difficulties. First, tensor network methods receive sparse coverage in LLM training data, causing severe hallucinations when models attempt simulations from incomplete knowledge. Second, simulations produce dense numerical outputs requiring precise quantitative analysis that LLMs cannot reliably perform through direct reasoning. Third, the tight coupling makes it difficult for a single agent to validate its own results, as failures may involve physics, code, or data artifacts simultaneously.

We address these challenges through in-context learning [54] with curated documentation and multi-agent decomposition for specialized validation. We systematically compare three configurations: a baseline single agent without documentation, a single agent with in-context learning, and a multi-agent architecture combining both approaches. Our multi-agent system interfaces with Renormalizer [55, 56], a general tensor network package with a focus on electron-phonon dynamics [57–60]. Benchmarking on quantum phase transitions in the two-dimensional Ising model, spin dynamics in the sub-Ohmic spin-boson model, and retinal photoisomerization using DeepSeek-V3.2, Gemini 2.5 Pro, and Claude Opus 4.5 reveals

that in-context learning is critical for success and that multi-agent decomposition further improves reliability. Different models exhibit characteristic failure modes.

## II. RESULTS

### A. Agent Architecture

Our approach combines two key strategies to address the challenges of autonomous tensor network simulation: in-context learning and multi-agent decomposition. The first strategy, in-context learning, addresses the sparse coverage of tensor network methods in LLM training data. Rather than relying on retrieval-augmented generation or expensive fine-tuning, we embed carefully curated Renormalizer documentation directly in the system prompt, totaling approximately 43k tokens. This documentation comprises Jupyter notebook tutorials (22k tokens), Python script examples (12k tokens), and refactored source code snippets (9k tokens). Embedding comprehensive documentation ensures that agents have instant access to all relevant information without relying on potentially imperfect retrieval queries. The complete Renormalizer source code would require approximately 0.8 million tokens, nearly 20 times the size of the curated documentation, and raw retrieval from source code often returns fragmented results or those lacking context that confuse the agent rather than help it. When debugging is required, agents can still access the full source code through file reading tools. While our implementation targets Renormalizer, the approach generalizes to other tensor network packages. The key requirement is comprehensive documentation that can fit within the context window. We expect similar agent systems can be developed for packages such as ITensor [61], TeNPy [62], and Block2 [63] with comparable effort.

The second strategy, multi-agent decomposition, addresses the tight coupling between physics, code, and data that makes validation difficult for a single agent. Scientific workflows naturally decompose into specialized subtasks, from high-level research planning to code implementation and result analysis. A single agent attempting to handle all aspects must maintain lengthy context and switch frequently between distinct reasoning modes. Our multi-agent architecture instead assigns each subtask to a dedicated agent with its own isolated context, a design pattern known as context quarantine. This isolation enables clear separation of concerns and more focused prompts, preventing irrelevant information from

one phase of the workflow from interfering with reasoning in another.

As shown in Fig. 1, our system comprises a central **Conductor** that coordinates seven specialized agents. The **Strategist** designs the overall research plan by decomposing complex scientific questions into manageable computational tasks. The **Guide** monitors progress through the simulation workflow and decides on the next step. The **Programmer** implements simulation code based on the research specification, with access to the embedded Renormalizer documentation. The **Executor** runs numerical simulations, monitors job progress, and handles minor debugging. The **Aggregator** collects and organizes numerical data from completed simulations, then performs data processing and analysis. The **Validator** critically reviews simulation results, checking for numerical artifacts, convergence issues, and physical inconsistencies. The **Visualizer** generates publication-quality figures from aggregated data. In the following sections, we include figures generated by the agent without manual modification or post-processing.

The **Conductor** orchestrates these agents iteratively, routing tasks to appropriate specialists and maintaining coherent progress toward the research objective. When an agent encounters difficulties, the **Conductor** routes the problem back for revision with relevant error information. For example, when the **Executor** encounters a bug beyond its capability, the **Conductor** reports it to the **Programmer** for fixing. Similarly, when the **Validator** raises concerns about numerical convergence, the **Conductor** directs the **Executor** to perform additional calculations with refined parameters.

The workflow begins with an interactive dialogue between the user and the agent to establish a complete research specification. Through this conversation, the user provides a high-level description of the research goal, and the agent asks clarifying questions about physical parameters, numerical settings, and computational constraints. The resulting specification captures all details needed for simulation and is then fed into the multi-agent system, which autonomously executes the simulation workflow. A detailed example of this process is presented in the following sections as well as in the Supporting Information.

We compare three agent configurations throughout this work. The baseline configuration uses a single agent without in-context learning, resembling general-purpose coding agents. When attempting simulations, the baseline agent must search and grep through the Renormalizer source code to learn usage patterns, a process prone to retrieving incomplete or misleading information. The single-agent configuration adds curated documentation to the

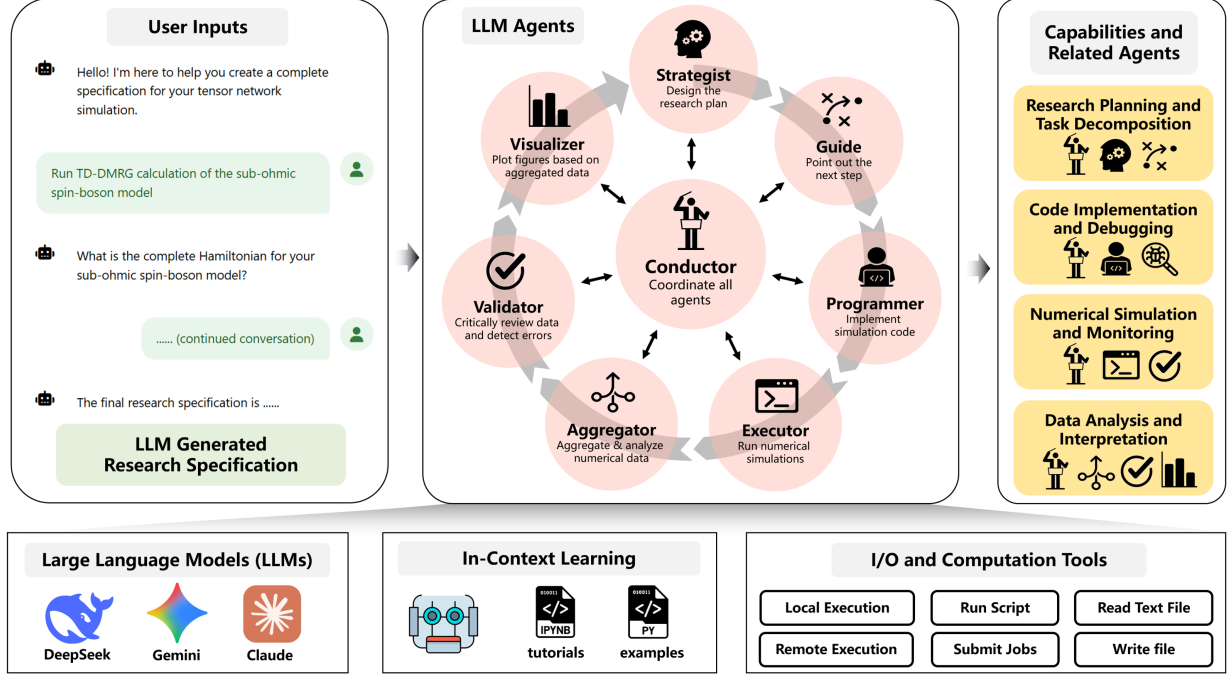


FIG. 1. **Multi-agent architecture for autonomous tensor network simulation.** The **Conductor** agent orchestrates seven specialized agents: **Strategist**, **Guide**, **Programmer**, **Executor**, **Aggregator**, **Validator**, and **Visualizer**. Agents communicate through structured messages and share access to tools such as file operations and code execution. The user interacts through natural language problem specifications and receives publication-ready outputs.

system prompt, providing comprehensive examples and API references. The multi-agent configuration combines in-context learning with the architecture shown in Fig. 1. In this work, we evaluate three state-of-the-art LLMs as backend engines, including DeepSeek-V3.2, Gemini 2.5 Pro, and Claude Opus 4.5. All three models achieve high success rates on the benchmark tasks presented below. The architecture is compatible with other LLMs, and we expect comparable frontier models to perform similarly.

## B. Benchmark Tasks

We evaluate our agent on three benchmark tasks spanning different areas of quantum many-body physics, with detailed physics background provided in the Methods section. The two-dimensional transverse-field Ising model represents a canonical ground-state problem with a well-established critical point ( $h_c/J \approx 3.04$  in the thermodynamic limit). We study

an  $8 \times 8$  lattice to ensure calculations complete within a reasonable time frame. This task tests whether agents can correctly implement DMRG calculations and interpret the resulting phase transition. The well-established critical value, however, creates a risk of hallucinations, as LLMs may output expected results even when simulations fail. The sub-Ohmic spin-boson model involves real-time quantum dynamics and phase classification. Unlike the Ising model, the coherent-incoherent phase boundary lacks a definitive literature value despite decades of numerical investigation, making validation more subtle. The fine phase boundary mapping presented in this work constitutes new scientific results rather than reproduction of existing literature. This task tests whether agents can correctly implement time evolution using TDVP, perform convergence analysis, and classify dynamical behavior from numerical trajectories. The retinal photoisomerization model demands implementation of specialized basis sets (exponential DVR) not available in Renormalizer. It also requires careful handling of unit and symbol conventions that differ between the literature and the software. This task tests the agent’s ability to translate complex physical specifications into working code, integrating information from multiple sources and understanding implicit conventions in scientific notation. For an expert already familiar with tensor network theory and proficient with a tensor network package, we estimate these tasks would require one day to one week to complete, depending on prior familiarity with the specific physical models.

Figure 2 presents model schematics alongside representative successful and failed runs for each task. The top row illustrates the three physical models. Figure 2(a) shows the 2D transverse-field Ising model, where spins on a square lattice interact ferromagnetically with coupling  $J$ , while a transverse field  $h$  induces quantum fluctuations. Figure 2(b) depicts the spin-boson model, where a two-level system with states  $|\uparrow\rangle$  and  $|\downarrow\rangle$  couples to a bath of harmonic oscillators, with the right panel showing the diabatic potential energy surfaces as a function of spin polarization  $\sigma_z$ . Figure 2(c) illustrates the retinal photoisomerization model, where the torsional coordinate  $\theta$  connects cis and trans configurations, and the 3D surface shows the coupled  $S_0$  and  $S_1$  potential energy surfaces that govern the ultrafast isomerization dynamics.

The middle row (d–f) shows successful runs produced by the multi-agent architecture. All figures are generated by the agent without manual modification or post-processing. Figure 2(d) shows absolute magnetization versus transverse field for the Ising model. The figure correctly captures that periodic boundary conditions stabilize the ferromagnetic phase

to higher field values due to additional bonds at the boundaries. The user’s initial prompt was deliberately brief: “Run DMRG simulation of 2d ising model and see how  $|M_z|$  changes with  $h$ . compare open and periodic boundary conditions.” From this minimal input, the agent initiated a clarification process to fill in the technical details required for a rigorous simulation. After receiving all necessary information from the user, the agent compiled a comprehensive research specification and passed it to the **Conductor** for the simulation workflow. Details of this dialogue are provided in the Supporting Information. Figure 2(e) shows a successful phase diagram for the spin-boson model. The critical coupling strength  $\alpha_c$  for the coherent-to-incoherent transition increases monotonically with the spectral exponent  $s$ , consistent with the physical expectation that sub-Ohmic baths with smaller  $s$  induce stronger dissipation. The agent autonomously performed parameter searches to determine appropriate bond dimensions and the number of phonon modes, calculated the dynamics across the parameter grid, and classified each point as coherent or incoherent based on the specified criteria. Figure 2(f) shows successful retinal photoisomerization dynamics. The  $S_0$  population and trans population evolution are consistent with reference calculations from the literature [64], demonstrating that the agent correctly implemented the exponential DVR basis and handled the coordinate conventions properly.

The bottom row (g-i) shows failed runs from the baseline configuration, illustrating characteristic failure modes. Figure 2(g) exhibits hallucinations: the agent prematurely concluded that the simulation was finished, resulting in missing periodic boundary condition data. The agent also computed  $\sqrt{\langle M^2 \rangle}$  in addition to  $|M|$ , which was not requested by the research specification. This completion bias reflects the tendency of LLMs to generate coherent narratives matching expected patterns regardless of whether the computation supports them. The overlapping curves obscure interpretation as visual inspection of the output figure was not performed. Figure 2(h) illustrates an implementation error. The agent failed to expand the bond dimension before TDVP time evolution, resulting in dynamics computed with bond dimension 1. Nearly the entire phase diagram is marked as “Unknown/Error,” yet the outputs appear complete and it is difficult to identify that the results are unreliable without close inspection of the underlying data. Figure 2(i) shows another implementation error specific to the retinal model. The highly oscillatory dynamics arise from two bugs: failing to expand the bond dimension as in Fig. 2(h), and incorrectly handling coordinate transformations between dimensionless and dimensional coordinates.



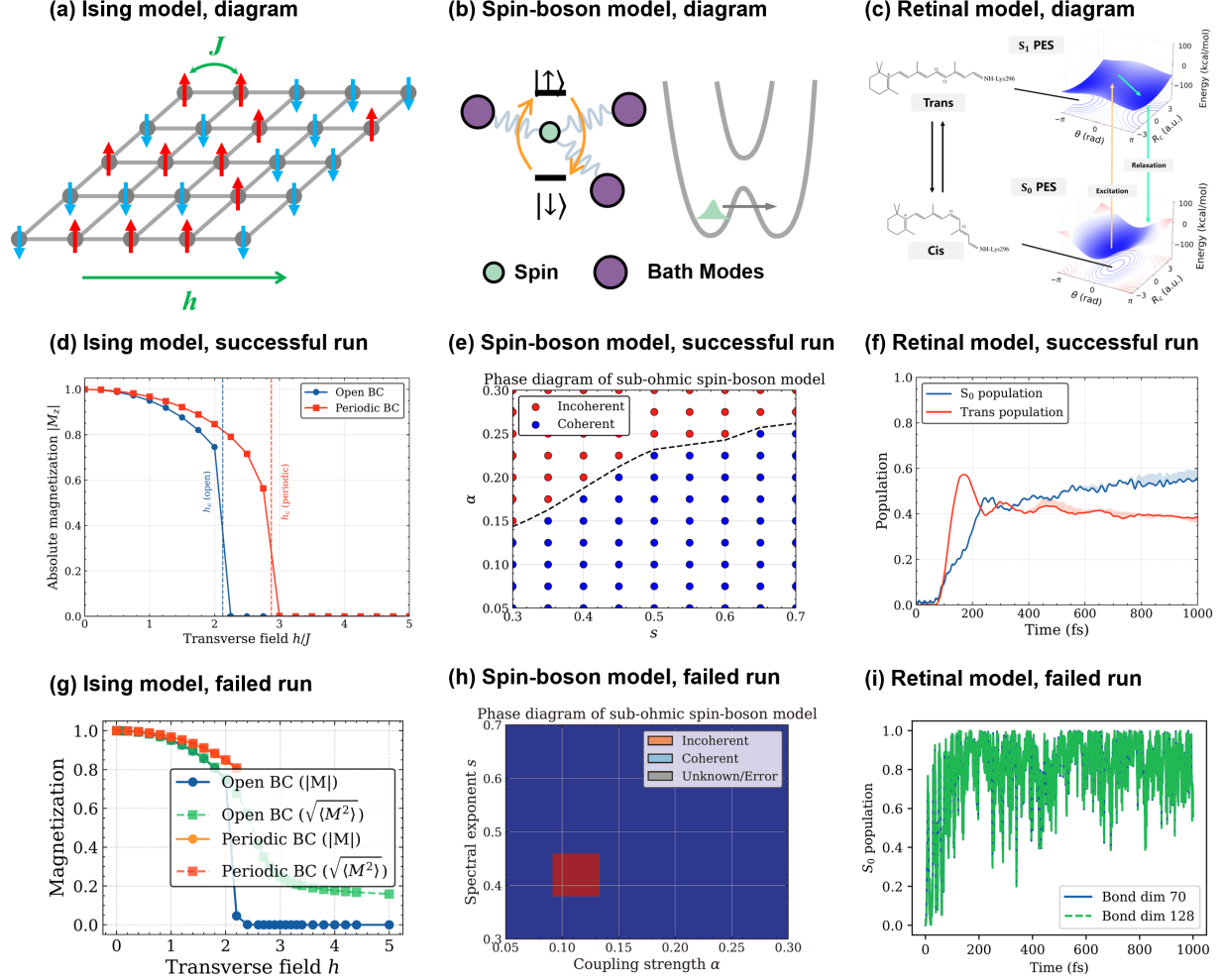


FIG. 2. **Model schematics and representative simulation results.** (a–c) Schematic diagrams of the three benchmark models: (a) 2D transverse-field Ising model with ferromagnetic coupling  $J$  and transverse field  $h$ ; (b) spin-boson model showing the two-level system coupled to bath modes and the diabatic potential energy surfaces; (c) retinal photoisomerization model with torsional coordinate connecting cis and trans configurations. (d–f) Successful runs from the multi-agent architecture showing correct phase transition behavior, accurate phase classification, and converged population dynamics. (g–i) Failed runs from the baseline configuration illustrating characteristic failure modes, which are missing data and overlapping curves, incomplete phase diagrams dominated by errors, and unphysical oscillatory dynamics.

Other commonly encountered bugs for the retinal model include implementing the  $\hat{p}^2$  operator as the second-order differencing operator while omitting the negative sign, and using energy values in eV directly from the literature without converting to atomic units.

The failure modes during the simulations are classified into four categories: implementation error, hallucinations, response error, and figure defects. Implementation errors occur when code runs without exceptions but produces incorrect results. These errors can arise during initial code generation, from incorrect physical parameters such as insufficient bond dimension and inadequate phonon basis truncation, or during post-processing when aggregating and analyzing results. Hallucinations encompass cases where the agent ignores user instructions, declares tasks complete prematurely, or makes judgments and decisions that contradict physical principles. Response errors arise from malformed outputs. DeepSeek-V3.2 occasionally produces responses that cannot be parsed as valid JSON, while Claude Opus 4.5 sometimes generates empty tool calls with missing arguments, both causing task termination. These response errors may stem from the complexity of maintaining structured output over extended multi-step reasoning. Finally, figure defects include visualization problems such as overlapping curves or missing labels that impair interpretation without affecting the underlying computation. We do not show results from runs terminated by response errors in Fig. 2, as these typically produce no usable output. We have confirmed through appropriate retry logic that such invalid responses do not arise from network errors, API rate limits, or transient failures in the agent architecture. While better prompt engineering or technical adjustments might reduce response errors, the fact that such errors occur more frequently in complex tasks like the retinal model than in simpler tasks like the Ising model suggests underlying robustness limitations in current LLMs.

### C. Typical Workflow

Figure 3 illustrates a complete workflow for the sub-Ohmic spin-boson model phase diagram task completed by DeepSeek-V3.2, demonstrating how multi-agent collaboration addresses the core challenges of tensor network simulation.

Tensor network simulations require careful navigation of multiple interdependent numerical parameters. For the spin-boson model, reliable results depend on adequate bond dimension  $M$  to capture entanglement, sufficient bath modes  $N_b$  to represent the continuous spectral density, and appropriate time evolution settings. These parameters interact in complex ways. Insufficient bond dimension produces qualitatively wrong dynamics, while excessive values waste computational resources. Determining appropriate parameters typi-

cally requires expert judgment informed by convergence analysis, a process that is tedious, error-prone, and rarely documented in published work.

The multi-agent architecture transforms this expert-driven process into a systematic, reproducible workflow. The **Strategist** decomposes the research goal into a seven-step plan that explicitly includes convergence testing before production runs. After an initial test run verifies the implementation, the **Executor** performs convergence sweeps across bond dimensions ( $M = 8, 12, 16$ ) and bath modes ( $N_b = 100, 200, 300, 400$ ) for representative parameter combinations. The **Validator** confirms convergence at each stage before the **Guide** proceeds to production runs, which scan the full parameter grid with 54 simulations. The workflow concludes with phase classification and visualization, producing a publication-ready phase diagram in 20 coordinated steps.

#### D. Performance Analysis

We evaluate each benchmark run on a 0–10 scale (see Methods Section IV E for detailed rubrics). A score of 10 indicates a perfect run producing publication-quality figures with correct scientific content, while a score of 0 represents complete failure due to response errors or severe hallucinations. The successful runs shown in Fig. 2(d–f) all receive scores of 10, demonstrating that the multi-agent architecture can produce publication-ready results. The failed runs in Fig. 2(g–i) receive scores of 6, 4, and 3, respectively. The Ising model run (g) loses points for hallucinations and figure defects. The spin-boson run (h) and the retinal run (i) receive lower scores due to implementation errors.

Figure 4 presents a comprehensive comparison of agent performance across all three benchmark tasks and agent architectures. The top row shows score distributions for each model, revealing that multi-agent architecture consistently outperforms both baseline and single-agent configurations across all LLM backends. The baseline configuration is plagued by implementation errors, hallucinations, and response errors, with each model exhibiting distinct failure patterns. DeepSeek-V3.2 performs relatively well on the Ising and spin-boson tasks in baseline mode, but all baseline runs on the retinal model suffer from implementation errors that produce incorrect results. Gemini 2.5 Pro exhibits severe hallucinations in baseline mode across all three tasks, as the model concludes that numerical results have been obtained and figures have been plotted, when in fact the simulation code failed and nothing

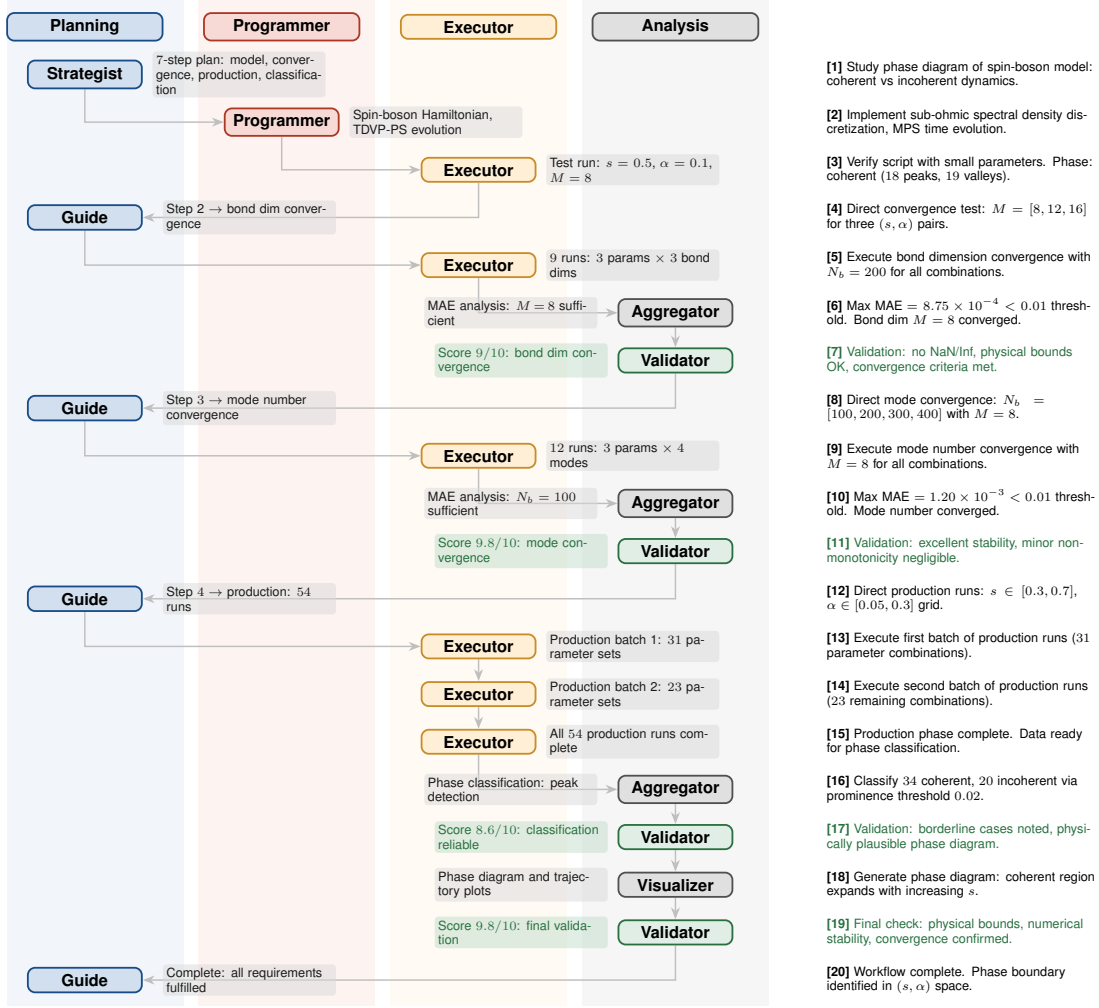


FIG. 3. **Complete workflow for spin-boson model phase diagram construction.** The workflow shows coordination between specialized agents across four tracks: planning (**Strategist**, **Guide**), programming (**Programmer**), execution (**Executor**), and analysis (**Aggregator**, **Validator**, **Visualizer**). The right column provides a step-by-step narrative of the 20 actions taken during the simulation. The **Validator** provides quality scores at each stage, ensuring numerical reliability before production runs and flagging borderline cases in the final phase classification.

was produced. Claude Opus 4.5 shows promise for the Ising and spin-boson tasks in baseline and single-agent configurations, but struggles with the retinal model due to response errors. Specifically, the model generates invalid tool calls with empty arguments when attempting to write Python code files, causing immediate task termination.

The multi-agent architecture with in-context learning improves performance across all

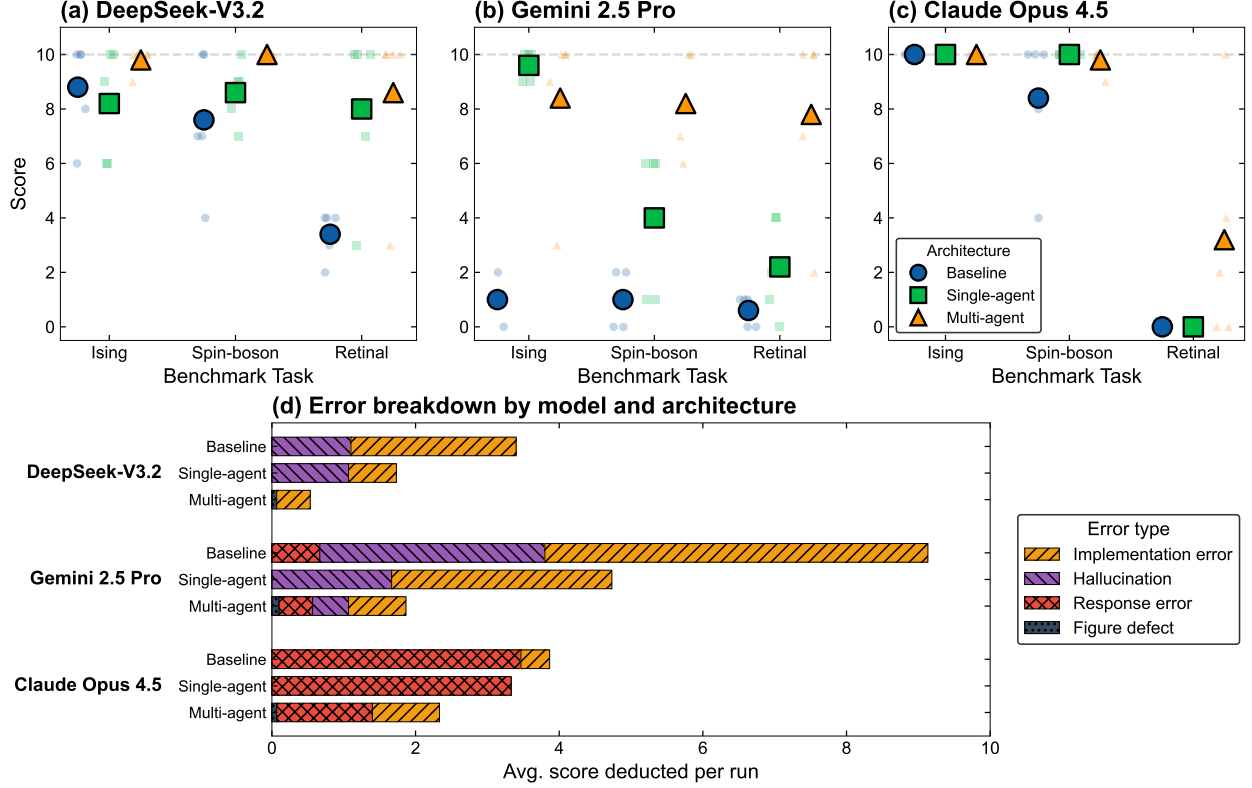


FIG. 4. **Performance comparison across agent architectures and LLM backends.** (a–c) Score distributions for each model across the three benchmark tasks (Ising, Spin-boson, Retinal). Small markers show individual runs (5 per task), while large markers indicate mean scores. (d) Average score deduction per run decomposed by error type for each model and architecture. Implementation errors and hallucination dominate failures for DeepSeek-V3.2 and Gemini 2.5 Pro, while response errors are the primary failure mode for Claude Opus 4.5.

models. DeepSeek-V3.2 with multi-agent architecture achieves the highest success rate, with 87% of runs receiving perfect scores and an average score of 9.47. Gemini 2.5 Pro shows the most dramatic improvement from multi-agent coordination. The baseline configuration achieves only 0.87 average score with no perfect runs, while the multi-agent configuration reaches an average score of 8.13 with 53% of runs receiving perfect scores.

Figure 4(d) decomposes the average score deduction per run by error type, revealing distinct failure patterns. For DeepSeek-V3.2 and Gemini 2.5 Pro, implementation errors and hallucinations constitute the dominant failure modes, with Gemini 2.5 Pro showing particularly severe issues in the baseline configuration. The multi-agent architecture substantially mitigates these errors, reducing average deductions to 0.53 for DeepSeek-V3.2 and

1.87 for Gemini 2.5 Pro. Claude Opus 4.5 exhibits a qualitatively different failure pattern, where response errors account for the majority of score deductions across all architectures, primarily arising from invalid outputs in the retinal photoisomerization task. These results demonstrate that multi-agent coordination with in-context learning provides consistent improvements over baseline approaches, and that different LLM backends exhibit characteristic failure modes.

We further compare token consumption and agent behavior across tasks in Figure 5. Figure 5(a) shows input token counts for each model-task combination. The number of input tokens is similar across the three tasks for each model, indicating that all three have comparable complexity in terms of the number of steps required for completion. This similarity results from our design choice that all three tasks should finish within a reasonable time. For complex tasks such as the retinal model, we did not ask the agent to perform extensive parameter convergence tests. Figure 5(b) shows output token counts, which are dramatically smaller than input tokens. Input tokens are on the order of 10 million, while output tokens are on the order of 0.1 million. This disparity arises because tool calls typically require few output tokens, whereas processing the Renormalizer documentation, script output, and full research context consumes substantial input tokens. Across the three models, DeepSeek-V3.2 consumes significantly more tokens than the others. For example, for the Ising model, DeepSeek-V3.2 consumes roughly 12 million input tokens, while Gemini 2.5 Pro and Claude Opus 4.5 consume roughly 3–4 million tokens. Figure 5(c) reveals the cause of the more token consumption by DeepSeek-V3.2, which makes approximately 300 tool calls per task, compared to 100–150 for the other models. This difference arises because DeepSeek-V3.2 tends to perform more granular actions, executing file operations such as listing directories, checking file existence, and reading files as separate calls, whereas Gemini 2.5 Pro and Claude Opus 4.5 batch operations more efficiently. For example, DeepSeek-V3.2 frequently reads simulation log files to monitor progress and check results, whereas Gemini 2.5 Pro and Claude Opus 4.5 tend to skip these intermediate checks. Since each tool call requires sending the full conversation context, more frequent calls lead to proportionally higher token consumption. Despite higher token consumption, DeepSeek-V3.2 remains cost-effective due to its lower API price. During the course of this study, the API prices for the three models were approximately \$0.3, \$1.25, and \$5 per million tokens, respectively. This makes DeepSeek-V3.2 the most cost-effective option among the three. Gemini 2.5 Pro has

a slightly higher cost, while Claude Opus 4.5 is typically 5 to 10 times more expensive than the other two models. The typical cost per task is estimated to be \$2–3 for DeepSeek-V3.2 and Gemini 2.5 Pro, and \$20–30 for Claude Opus 4.5. Cache reading may further reduce costs but is not considered here.

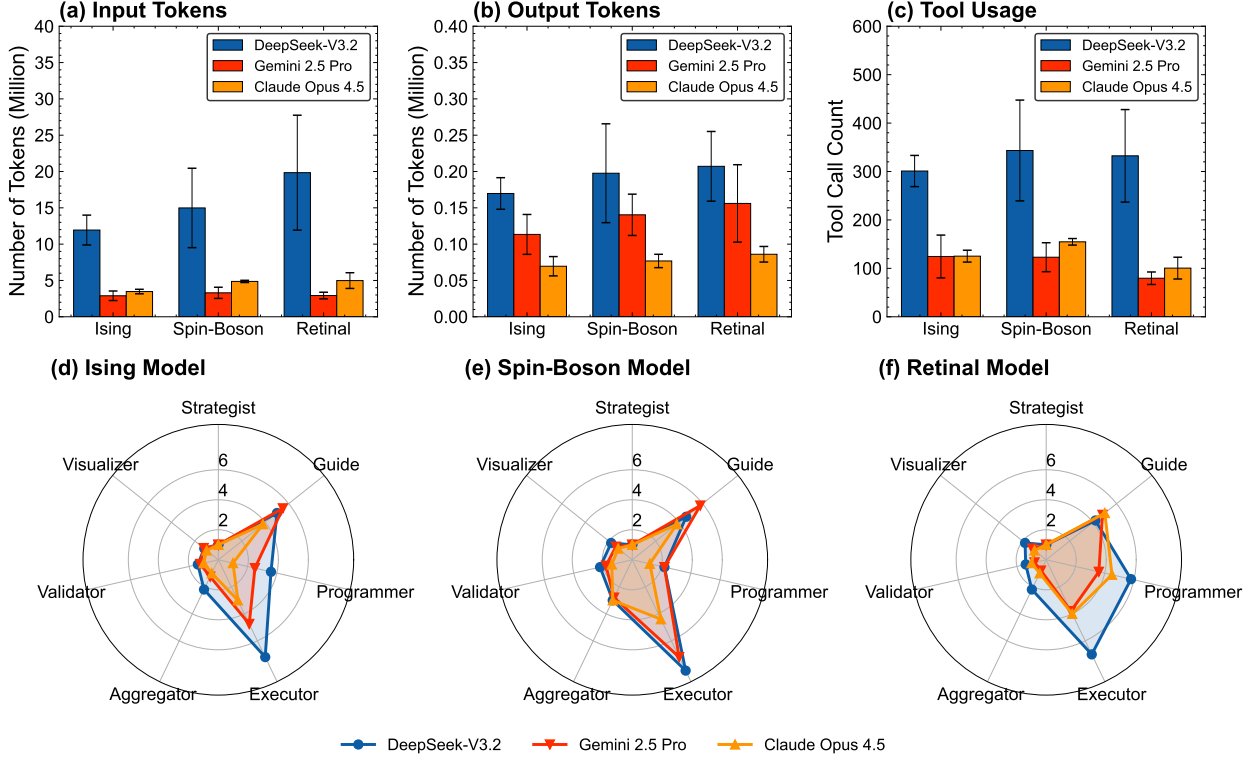


FIG. 5. **Token consumption and agent behavior analysis.** (a) Input and (b) output token counts for each model-task combination. (c) Total tool call counts. (d–f) Distribution of high-level agent calls across the seven specialized agents for the Ising model, spin-boson model, and retinal model, respectively.

The radar plots in Figure 5(d–f) show the distribution of high-level agent calls to the seven specialized agents for each model-task combination. The Ising and spin-boson model tasks require extensive numerical calculations, resulting in more calls to the **Executor** and **Aggregator**. In contrast, the major challenge for the retinal model is correct implementation of the simulation code, leading to relatively more calls to the **Programmer**. Among the LLM models, DeepSeek-V3.2 exhibits a higher number of agent calls, reflecting more iterative debugging cycles and data processing steps.

During preliminary development, we tested DeepSeek-V3.2 without reasoning enabled

and observed severe hallucination in the **Validator**. For example, the **Validator** frequently claimed that Ising model results were incorrect because “the energy should increase with  $h$ ,” and suggested the Hamiltonian implementation was wrong. When tested in a standalone chat environment without reasoning, DeepSeek-V3.2 almost consistently answered that energy increases with transverse field when asked to respond quickly without explanation. After receiving the warning from the **Validator**, the **Conductor** would then attempt to verify the result, and many steps were taken before realizing the Validator had made a mistake. This hallucination disappeared when reasoning was enabled, motivating our decision to use reasoning-enhanced models for all production runs reported in this work. This suggests that scientific agents must be designed to encourage deliberate reasoning rather than rapid response, particularly for validation tasks where physical intuition matters.

### III. DISCUSSION

We have demonstrated that LLM-driven agents can autonomously perform tensor network simulations of quantum many-body systems. Our approach addresses two key challenges: the sparse coverage of tensor network methods in LLM training data through in-context learning with curated documentation, and the tight coupling between physics, code, and data through multi-agent decomposition. Systematic comparison across three agent configurations reveals that both components are essential. The baseline configuration, which relies on the LLM to search through source code, fails frequently due to hallucinations and implementation errors. Adding in-context learning substantially improves performance, while multi-agent decomposition provides further gains by enabling focused expertise and iterative validation. The multi-agent architecture achieves approximately 90% success rate across three representative tasks spanning ground-state calculations, open quantum system dynamics, and photochemical reactions.

Analysis of both successful and failed runs reveals characteristic patterns across LLM backends. DeepSeek-V3.2 achieves the highest success rate but consumes more tokens due to more granular actions. Gemini 2.5 Pro shows the most dramatic improvement from multi-agent coordination, rising from near-complete failure in baseline mode to reasonable performance with the full architecture. Claude Opus 4.5 achieves consistent results but occasionally produces malformed outputs that terminate workflows prematurely.



Despite these encouraging results, the generated outputs require human verification due to the non-negligible error rate. In practice, failures are typically detectable, because the **Validator** flags convergence issues, response errors terminate workflows with clear error messages, and implementation errors often produce obviously unphysical results (as in Fig. 2(i)). Some failures, however, produce plausible-looking but incorrect results that require expert inspection to identify, representing a minority of cases. The failure modes we observed, including numerical parameter insufficiency, implementation errors despite correct reasoning, and subtle mistakes in translating mathematical specifications to code, mirror challenges faced by human researchers, suggesting that complete automation of scientific computation remains an open problem.

Looking forward, we anticipate rapid progress on several fronts. Continued improvements in LLM reasoning capabilities and context windows will likely reduce error rates and expand the complexity of problems that can be tackled autonomously. Advances in agent architectures, including improved multi-agent collaboration and self-verification protocols, may enhance output quality. Integration with literature retrieval could enable agents to autonomously identify research questions and design simulations. The emergence of multi-modal models will enable agents to directly parse equations, tables, and figures from published literature, and to critically examine their own generated plots for anomalies, iteratively refining them.

As these capabilities mature, LLM-driven agents may democratize access to tensor network methods. Currently, effective use of these powerful techniques requires years of specialized training, limiting their application to a small community of experts. By encoding domain expertise in documentation and agent workflows, the barrier to entry could be dramatically lowered, enabling researchers from adjacent fields to leverage tensor network methods for their own problems. Domain experts, freed from routine implementation and convergence test tasks, could focus on more creative and challenging aspects of research, accelerating scientific discovery in quantum many-body physics and beyond.

## IV. METHODS

### A. Matrix Product States and Renormalizer

The matrix product state is a representative tensor network ansatz that decomposes the exponentially large state vector into a product of local tensors. For a system of  $N$  sites, an MPS represents the state as

$$|\Psi\rangle = \sum_{\{a\},\{\sigma\}} A_{a_1}^{\sigma_1} A_{a_1 a_2}^{\sigma_2} \cdots A_{a_{N-1}}^{\sigma_N} |\sigma_1 \sigma_2 \cdots \sigma_N\rangle \quad (1)$$

where  $\sigma_i$  labels the local basis states with dimension  $d$ , and  $a_i$  are auxiliary bond indices with dimension up to  $M$ , the bond dimension. The bond dimension controls the amount of entanglement that can be captured. Larger  $M$  allows more accurate representation of highly entangled states but increases computational cost, which scales as  $O(NM^3d)$  for typical operations. For ground state calculations, the density matrix renormalization group algorithm optimizes the MPS tensors variationally. All results reported in this work use MPS as the simulation ansatz.

For time evolution, the time-dependent variational principle projects the Schrödinger equation onto the MPS manifold. The key idea is to find the optimal time derivative of the MPS that best approximates the exact evolution within the manifold of fixed bond dimension. This is achieved by applying the projector  $\hat{P}_{T_{|\Psi\rangle}\mathcal{M}}$  onto the tangent space of the MPS manifold at the current state:

$$\frac{\partial}{\partial t} |\Psi\rangle = -i \hat{P}_{T_{|\Psi\rangle}\mathcal{M}} \hat{H} |\Psi\rangle \quad (2)$$

The projector can be decomposed into local terms involving forward and backward sweeps through the MPS sites, leading to an efficient algorithm with computational cost scaling as  $O(NM^3d)$  per time step. TDVP preserves important physical properties including unitarity and energy conservation, making it particularly suitable for long-time dynamics simulations.

Renormalizer is an open-source Python package for tensor network simulations with a special focus on electron-phonon quantum dynamics [55, 56]. The package provides MPS and matrix product operator algorithms for ground state search and time evolution, tree tensor network support for more complex network topologies, and automatic operator construction from first and second-quantized Hamiltonians [65, 66]. Renormalizer also supports finite

temperature simulations through imaginary time evolution or thermofield transformation, GPU acceleration via CuPy, and quantum number conservation for exploiting symmetries. A comprehensive review of the underlying methods is available in Ref. [67].

## B. Agent Implementation and LLM Configuration

The agent system is implemented using the LangChain framework for LLM orchestration. Each agent is instantiated as a LangChain tool with a dedicated system prompt that encodes its specialized role and constraints. The **Conductor** operates with a recursion limit of 200 steps, allowing complex multi-step workflows while preventing infinite loops. All agents share access to common utility tools for file operations, including reading, writing, and listing directories, as well as code execution. The **Strategist** and **Guide** roles are implemented as a single LLM with shared conversation history, enabling coherent planning across multiple consultation rounds. A call to the **Strategist** generates a detailed step-by-step research plan, while subsequent calls to the **Guide** continue the same conversation thread, allowing the **Guide** to track progress against the original plan. The **Programmer** agent receives requirements from the **Conductor** along with the original user input for context and uses a reasoning-enhanced model at temperature 0 to ensure careful code generation. The **Executor** agent runs simulations by invoking the main script with different parameter sets. It supports parallel execution of up to 32 jobs per batch and automatically handles job scheduling for larger parameter sweeps. When encountering errors, the **Executor** attempts to debug by analyzing log files and applying fixes via diff patches. The **Aggregator** agent processes output files from multiple simulation runs, computing derived quantities and preparing data for visualization and human inspection. It answers specific scientific questions posed by the **Conductor** using quantitative analysis of the simulation outputs. The **Validator** agent critically examines numerical results for unphysical values, convergence failures, and systematic errors, providing severity assessments with recommendations to continue, investigate, or stop execution. The **Visualizer** agent generates publication-quality figures with a custom style template. It inspects data structures, makes plots using matplotlib, and saves output as PDF files.

We evaluate three LLM backends: DeepSeek-V3.2 with reasoning enabled, Gemini 2.5 Pro, and Claude Opus 4.5. All models use temperature 0 for the **Conductor** and **Pro-**

**grammer** agents to ensure consistent outputs, while other agents use default temperature settings. The system implements automatic retry logic with exponential backoff (initial delay of 120 seconds, factor of 2.0, maximum of 3 retries) to handle response errors. Invalid tool calls trigger automatic retries with the same backoff policy.

### C. Benchmark Tasks

The three benchmark tasks were selected according to the following criteria: (1) coverage of distinct algorithmic challenges (ground-state optimization, real-time dynamics, and custom basis implementation), (2) representation of different physical domains where tensor networks are applied (condensed matter, open quantum systems, and photochemistry), (3) varying degrees of information availability in LLM training data (well-documented critical point versus debated phase boundary versus specialized model), and (4) completion within a reasonable time frame (hours rather than days) to enable systematic evaluation across multiple runs and configurations.

#### 1. Two-Dimensional Transverse-Field Ising Model

The two-dimensional transverse-field Ising model on an  $N_x \times N_y$  lattice is described by the Hamiltonian

$$\hat{H} = -J \sum_{\langle i,j \rangle} \sigma_i^z \sigma_j^z - h \sum_i \sigma_i^x \quad (3)$$

where  $J > 0$  is the ferromagnetic coupling between nearest-neighbor spins,  $h$  is the transverse field strength, and  $\sigma_i^{x,z}$  are Pauli matrices at site  $i$ . The ground state of this model is governed by the competition between  $J$  and  $h$ . At zero temperature, the system undergoes a quantum phase transition as the ratio  $h/J$  is varied. For small  $h/J$ , the ground state exhibits spontaneous magnetization in the ferromagnetic phase, while for large  $h/J$ , quantum fluctuations dominate and the system becomes paramagnetic. The quantum transverse-field Ising model in two dimensions has no exact analytical solution, and extensive numerical studies using cluster Monte Carlo and DMRG have established the critical field at  $h_c/J \approx 3.04$  in the thermodynamic limit [68, 69]. The order parameter  $|M_z|$  versus  $h/J$  provides a direct probe of the phase transition.

We tasked our LLM-driven agent with studying the quantum phase transition in this model. We limited the computation to an  $8 \times 8$  lattice. Although finite-size effects are substantial at this system size, this choice ensures that calculations complete within a reasonable time frame. The methodology readily generalizes to larger lattices for production research.

## 2. Sub-Ohmic Spin-Boson Model

The spin-boson model describes a two-level system coupled to a bosonic bath and serves as a paradigmatic model for quantum dissipation [70]. The Hamiltonian is

$$\hat{H} = \frac{\Delta}{2}\sigma_x + \frac{\epsilon}{2}\sigma_z + \sum_k \omega_k b_k^\dagger b_k + \frac{\sigma_z}{2} \sum_k g_k (b_k^\dagger + b_k) \quad (4)$$

where  $\Delta$  is the tunneling amplitude,  $\epsilon$  is the bias,  $\omega_k$  and  $g_k$  are the frequency and coupling strength of the  $k$ -th bath mode, and  $b_k^{(\dagger)}$  are bosonic annihilation (creation) operators. The bath is characterized by the spectral density

$$J(\omega) = 2\pi\alpha\omega_c^{1-s}\omega^s e^{-\omega/\omega_c} \quad (5)$$

where  $\alpha$  is the dimensionless coupling strength,  $\omega_c$  is the cutoff frequency, and  $s$  determines the bath type: sub-Ohmic ( $s < 1$ ), Ohmic ( $s = 1$ ), or super-Ohmic ( $s > 1$ ). The sub-Ohmic spin-boson model exhibits a quantum phase transition at zero temperature. Below a critical coupling  $\alpha_c(s)$ , the spin dynamics shows coherent oscillations, while above  $\alpha_c$ , the dynamics becomes incoherent with monotonic decay toward a localized state. The boundary between coherent and incoherent regimes defines a dynamical phase diagram in the  $\alpha$ - $s$  parameter space. Determining this phase boundary has been the subject of extensive numerical investigation using a variety of advanced methods, including multilayer multiconfiguration time-dependent Hartree [71], the Davydov ansatz [72], hierarchical equations of motion [73], and quasi-adiabatic propagator path integral [74]. Remarkably, despite decades of effort, the exact location of the coherent-incoherent phase boundary remains under debate.

For numerical simulations, the continuous bath is discretized into  $N_b$  modes using an exponential discretization scheme. The mode frequencies are determined by

$$\omega_k = -\omega_c \ln \left[ 1 - \frac{k}{N_b + 1} \right], \quad k = 1, 2, \dots, N_b \quad (6)$$

with corresponding coupling strengths derived from the spectral density. The exponential discretization scheme in Eq. 6 implicitly determines the discrete frequencies  $\omega_k$  by inverting a cumulative distribution. Deriving this expression requires correctly evaluating the integral and obtaining the analytical form  $\omega_k = -\omega_c \ln[1 - k/(N_b + 1)]$ . All three LLM backends successfully performed this derivation.

Obtaining converged results requires careful selection of numerical parameters. The agent must perform convergence analysis over the bond dimension  $M$  and the number of bath modes  $N_b$ , and choose an appropriate basis truncation for each bosonic mode to ensure reliable dynamics.

### 3. Retinal Photoisomerization Model

The retinal photoisomerization in rhodopsin represents one of the fastest and most efficient photochemical reactions in nature. Upon absorption of a photon, the 11-cis retinal chromophore undergoes isomerization to the 11-trans configuration. This ultrafast reaction has attracted widespread interest due to its importance as the first step in vision [64, 75, 76]. We employ a well-established two-state 25-mode model to describe the retinal photoisomerization reaction. The model comprises two diabatic electronic states  $|0\rangle$ , the ground state  $S_0$ , and  $|1\rangle$ , the excited state  $S_1$ , one torsional degree of freedom  $\theta$  describing rotation about the  $C_{11}=C_{12}$  double bond, one stretching mode  $R_c$  that mediates electronic transitions near the conical intersection, and a bath of 23 harmonic oscillators  $\{R_j\}$  representing the protein environment. The full Hamiltonian takes the form:

$$\hat{H} = \hat{H}_0|0\rangle\langle 0| + \hat{H}_1|1\rangle\langle 1| + \hat{V}_{01}(|0\rangle\langle 1| + |1\rangle\langle 0|) \quad (7)$$

where the diabatic state Hamiltonians are:

$$\hat{H}_0 = \frac{\hat{p}_\theta^2}{2I} + \frac{W_0}{2}(1 - \cos \theta) + \hat{H}_{\text{vib}} \quad (8)$$

$$\hat{H}_1 = \frac{\hat{p}_\theta^2}{2I} + E_1 - \frac{W_1}{2}(1 - \cos \theta) + \hat{H}_{\text{vib}} + \kappa_c \hat{R}_c + \sum_{j=3}^{25} \kappa_j \hat{R}_j \quad (9)$$

Here  $I$  is the moment of inertia for torsional motion,  $W_0$  and  $W_1$  are the barrier heights for the ground and excited state potentials respectively,  $E_1$  is the vertical excitation energy, and the  $\kappa$  terms represent linear vibronic coupling. The vibrational Hamiltonian  $\hat{H}_{\text{vib}}$  describes

the harmonic bath modes:

$$\hat{H}_{\text{vib}} = \frac{\Omega_c}{2}(\hat{p}_c^2 + \hat{R}_c^2) + \sum_{j=3}^{25} \frac{\omega_j}{2}(\hat{p}_j^2 + \hat{R}_j^2) \quad (10)$$

The diabatic coupling is mediated by the coupling mode,  $\hat{V}_{01} = \lambda \hat{R}_c$ , where  $\lambda$  controls the strength of non-adiabatic transitions.

This benchmark poses a distinctive challenge of implementing a specialized basis set not available in Renormalizer. The torsional coordinate  $\theta$  is periodic with period  $2\pi$ , requiring basis functions that respect this periodicity. The appropriate choice is the exponential discrete variable representation (exponential DVR). Another challenge is the use of dimensionless coordinates  $\hat{R}$  and momenta  $\hat{p}$  in the Hamiltonian, a convention that differs from both the spin-boson model described above and the default convention in Renormalizer. The agent must correctly identify this difference and implement the Hamiltonian accordingly. This knowledge is typically acquired through graduate-level training in quantum dynamics.

#### D. Prompts of the Benchmark Tasks

For the Ising model task, the prompt contains minimal information, as described in the Results section. For the spin-boson model task, the prompt includes the necessary equations for deterministic implementation, namely the Hamiltonian and the exponential discretization of the spectral density. The prompt also includes an explicit definition of the phase classification criterion to ensure reproducible results across different runs. Coherent dynamics is defined as damped oscillatory behavior exhibiting at least one valley and one subsequent peak in  $\langle \sigma_z(t) \rangle$ , while incoherent dynamics encompasses either purely monotonic decay or a single valley followed by localization without oscillation. This level of methodological detail is necessary because the coherent-incoherent boundary is precisely where dynamics are most ambiguous. For the retinal model task, the prompt provided to the agent consists of a description of the model from the literature containing the complete model specification, including all Hamiltonian parameters and bath mode frequencies [64]. The description was generated by taking a screenshot of the published paper and converting it to LaTeX using LLMs. Similarly, the agent was provided with a description of exponential DVR from the literature [77]. The agent was instructed to run simulations with two different bond

dimensions, 70 and 128. Details of all prompts are provided in the Supporting Information.

### **E. Evaluation Rubrics**

Each benchmark run is evaluated on a 0–10 scale based on the quality of the simulation code, numerical results, and final interpretation and visualization. The rubrics are designed to capture the error types discussed in the Results section, including implementation errors, hallucinations, response errors, and figure defects.

A score of 10 represents a perfect run, which means that the agent produces publication-quality figures with correct scientific content and professional interpretation, fully addressing the research specification. A score of 8 indicates near-perfect execution where the results are scientifically correct but minor issues exist, such as suboptimal figure aesthetics or incomplete adherence to the research specification that does not affect scientific conclusions. This score typically results from figure defects or minor implementation errors during post-processing and analysis.

A score of 6 reflects runs where the simulation code executes correctly and produces some valid results, but significant errors occur during execution or analysis. This includes cases where key information required by the research specification is missing from the produced figure, or where some scientific content in the figure is incorrect. Such scores typically arise from major implementation errors or hallucinations during the analysis phase. A score of 4 indicates partial failure, which means that the code is nearly correct but contains incorrect parameters or 1–2 bugs that an expert could easily fix, with only partial results of limited reference value. This score is typically caused by major implementation errors in the main simulation script or strong hallucinations during execution and analysis.

A score of 2 represents mostly failed runs where code implementation is incomplete with significant bugs, producing no usable results. A score of 0 indicates complete failure: no working code is produced, the agent terminates due to response errors, or severe hallucinations result in entirely unreliable outputs. Scores of 0–2 are typically associated with severe hallucinations, fundamental implementation errors, or response errors that terminate the workflow prematurely.



## DATA AVAILABILITY

The Renormalizer package is available at <https://github.com/shuaigroup/Renormalizer>. The source code of the agent system and the log files for all benchmark tasks will be made available upon publication.

## ACKNOWLEDGEMENTS

The authors thank Shi-Xin Zhang, Xiang Sun, Zhenggang Lan, and Ningyi Lyu for insightful discussions. Weitang Li is supported by the Guangdong Basic Research Center of Excellence for Aggregate Science, the Shenzhen Science and Technology Program (No. KQTD20240729102028011), and the National Natural Science Foundation of China (Grant No. 22595405). Jiajun Ren is supported by the National Natural Science Foundation of China (Grant No. 22273005 and No. 22422301).

## COMPETING INTERESTS

The authors declare no competing interests.

- 
- [1] K. M. Jablonka, Q. Ai, A. Al-Feghali, S. Badhwar, J. D. Bocarsly, A. M. Bran, S. Bringuier, L. C. Brinson, K. Choudhary, D. Circi, *et al.*, 14 examples of how LLMs can transform materials science and chemistry: A reflection on a large language model hackathon, *Digit. Discov.* **2**, 1233 (2023).
  - [2] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao, Large language models: A survey, *arXiv preprint arXiv:2402.06196* (2024).
  - [3] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, X. W. Zhao, Z. Wei, and W. Ji-Rong, A survey on large language model based autonomous agents, *Front. Comput. Sci.* **18**, 186345 (2024).
  - [4] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, *et al.*, The rise and potential of large language model based agents: A survey, *Sci. China Inf. Sci.* **68**, 121101 (2025).

- [5] M. C. Ramos, C. J. Collison, and A. D. White, A review of large language models and autonomous agents in chemistry, *Chem. Sci.* **16**, 2514 (2025).
- [6] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, GPT-4 technical report, arXiv preprint arXiv:2303.08774 (2023).
- [7] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, *et al.*, Gemini: A family of highly capable multimodal models, arXiv preprint arXiv:2312.11805 (2023).
- [8] D. Guo, D. Yang, H. Zhang, J. Song, P. Wang, Q. Zhu, R. Xu, R. Zhang, S. Ma, X. Bi, *et al.*, Deepseek-R1 incentivizes reasoning in LLMs through reinforcement learning, *Nature* **645**, 633 (2025).
- [9] M. R. AI4Science and M. A. Quantum, The impact of large language models on scientific discovery: a preliminary study using GPT-4, arXiv preprint arXiv:2311.07361 (2023).
- [10] K. M. Jablonka, P. Schwaller, A. Ortega-Guerrero, and B. Smit, Leveraging large language models for predictive chemistry, *Nat. Mach. Intell.* **6**, 161 (2024).
- [11] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao, React: Synergizing reasoning and acting in language models, in *The eleventh international conference on learning representations* (2022).
- [12] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, and T. Scialom, Toolformer: Language models can teach themselves to use tools, *Advances in Neural Information Processing Systems* **36**, 68539 (2023).
- [13] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, H. Wang, and H. Wang, Retrieval-augmented generation for large language models: A survey, arXiv preprint arXiv:2312.10997 **2** (2023).
- [14] S. Steiner, J. Wolf, S. Glatzel, A. Andreou, J. M. Granda, G. Keenan, T. Hinkley, G. Aragon-Camarasa, P. J. Kitson, D. Angelone, *et al.*, Organic synthesis in a modular robotic system driven by a chemical programming language, *Science* **363**, eaav2211 (2019).
- [15] M. Abolhasani and E. Kumacheva, The rise of self-driving labs in chemical and materials sciences, *Nat. Synth.* **2**, 483 (2023).
- [16] T. Song, M. Luo, X. Zhang, L. Chen, Y. Huang, J. Cao, Q. Zhu, D. Liu, B. Zhang, G. Zou, *et al.*, A multiagent-driven robotic ai chemist enabling autonomous chemical research on de-

- mand, J. Am. Chem. Soc. **147**, 12534 (2025).
- [17] A. M. Bran, S. Cox, O. Schilter, C. Baldassari, A. D. White, and P. Schwaller, Augmenting large language models with chemistry tools, Nat. Mach. Intell. **6**, 525 (2024).
  - [18] D. A. Boiko, R. MacKnight, B. Kline, and G. Gomes, Autonomous chemical research with large language models, Nature **624**, 570 (2023).
  - [19] L. Huang, C. Zhang, Y. Fu, Y. Jiang, E. He, M.-Q. Qi, M.-H. Du, X.-J. Kong, J. Cheng, L. Cronin, *et al.*, Natural-language-interfaced robotic synthesis for AI-Copilot-assisted exploration of inorganic materials, (2025).
  - [20] Y. Zou, A. H. Cheng, A. Aldossary, J. Bai, S. X. Leong, J. A. Campos-Gonzalez-Angulo, C. Choi, C. T. Ser, G. Tom, A. Wang, *et al.*, El agente: An autonomous agent for quantum chemistry, Matter **8** (2025).
  - [21] J. Hu, H. Nawaz, Y. Rui, L. Chi, A. Ullah, and P. O. Dral, Aitomia: Your intelligent assistant for ai-driven atomistic and quantum chemical simulations, arXiv preprint arXiv:2505.08195 (2025).
  - [22] R. S. Gadde, S. Devaguptam, F. Ren, R. Mittal, L. Dong, Y. Wang, and F. Liu, Chatbot-assisted quantum chemistry for explicitly solvated molecules, Chem. Sci. **16**, 3852 (2025).
  - [23] Q. Zhang, Y. Hu, J. Yan, H. Zhang, X. Xie, J. Zhu, H. Li, X. Niu, L. Li, Y. Sun, *et al.*, Large-language-model-based ai agent for organic semiconductor device research, Adv. Mater. **36**, 2405163 (2024).
  - [24] E. Xie, L. Cheng, J. Shireman, Y. Cai, J. Liu, C. Mohanty, M. Dey, and C. Kendzierski, Cassia: A multi-agent large language model for automated and interpretable cell annotation, Nat. Commun. (2025).
  - [25] I. Mandal, J. Soni, M. Zaki, M. M. Smedskjaer, K. Wondraczek, L. Wondraczek, N. N. Gosvami, and N. A. Krishnan, Evaluating large language model agents for automation of atomic force microscopy, Nat. Commun. **16**, 9104 (2025).
  - [26] K. Swanson, W. Wu, N. L. Bulaong, J. E. Pak, and J. Zou, The virtual lab of ai agents designs new sars-cov-2 nanobodies, Nature **646**, 716 (2025).
  - [27] R. Orús, Tensor networks for complex quantum systems, Nat. Phys. Rev. **1**, 538 (2019).
  - [28] M. C. Bañuls, Tensor network algorithms: A route map, Annu. Rev. Condens. Matter Phys. **14**, 173 (2023).

- [29] S. R. White, Density matrix formulation for quantum renormalization groups, *Phys. Rev. Lett.* **69**, 2863 (1992).
- [30] S. R. White, Density-matrix algorithms for quantum renormalization groups, *Phys. Rev. B* **48**, 10345 (1993).
- [31] U. Schollwöck, The density-matrix renormalization group in the age of matrix product states, *Ann. Phys. (NY)* **326**, 96 (2011).
- [32] Y.-Y. Shi, L.-M. Duan, and G. Vidal, Classical simulation of quantum many-body systems with a tree tensor network, *Phys. Rev. A* **74**, 022320 (2006).
- [33] G. Vidal, Entanglement renormalization, *Phys. Rev. Lett.* **99**, 220405 (2007).
- [34] J. I. Cirac, D. Perez-Garcia, N. Schuch, and F. Verstraete, Matrix product states and projected entangled pair states: Concepts, symmetries, theorems, *Rev. Mod. Phys.* **93**, 045003 (2021).
- [35] J. Haegeman, J. I. Cirac, T. J. Osborne, I. Pižorn, H. Verschelde, and F. Verstraete, Time-dependent variational principle for quantum lattices, *Phys. Rev. Lett.* **107**, 070601 (2011).
- [36] J. Haegeman, C. Lubich, I. Oseledets, B. Vandereycken, and F. Verstraete, Unifying time evolution and optimization with matrix product states, *Phys. Rev. B* **94**, 165116 (2016).
- [37] C. Lubich, I. V. Oseledets, and B. Vandereycken, Time integration of tensor trains, **53**, 917 (2015).
- [38] S. Paeckel, T. Köhler, A. Swoboda, S. R. Manmana, U. Schollwöck, and C. Hubig, Time-evolution methods for matrix-product states, *Ann. Phys. (NY)* **411**, 167998 (2019).
- [39] Z. Li, S. Guo, Q. Sun, and G. K.-L. Chan, Electronic landscape of the p-cluster of nitrogenase as revealed through many-electron quantum wavefunction simulations, *Nat. Chem.* **11**, 1026 (2019).
- [40] P. Sharma, D. G. Truhlar, and L. Gagliardi, Magnetic coupling in a tris-hydroxo-bridged chromium dimer occurs through ligand mediated superexchange in conjunction with through-space coupling, **142**, 16644 (2020).
- [41] H. R. Larsson, H. Zhai, C. J. Umrigar, and G. K.-L. Chan, The chromium dimer: closing a chapter of quantum chemistry, **144**, 15932 (2022).
- [42] H.-J. Liao, Z.-Y. Xie, J. Chen, Z.-Y. Liu, H.-D. Xie, R.-Z. Huang, B. Normand, and T. Xiang, Gapless spin-liquid ground state in the  $s = 1/2$  kagome antiferromagnet, *Phys. Rev. Lett.* **118**, 137202 (2017).

- [43] P. Fowler-Wright, B. W. Lovett, and J. Keeling, Efficient many-body non-markovian dynamics of organic polaritons, *Phys. Rev. Lett.* **129**, 173001 (2022).
- [44] C. Huang, F. Zhang, M. Newman, X. Ni, D. Ding, J. Cai, X. Gao, T. Wang, F. Wu, G. Zhang, *et al.*, Efficient parallelization of tensor network contraction for simulating quantum computation, *Nat. Comput. Sci.* **1**, 578 (2021).
- [45] A. Berezutskii, M. Liu, A. Acharya, R. Ellerbrock, J. Gray, R. Haghshenas, Z. He, A. Khan, V. Kuzmin, D. Lyakh, *et al.*, Tensor networks for quantum computing, *Nat. Rev. Phys.* **7**, 581 (2025).
- [46] P. Sehlstedt, J. Brandejs, P. Bientinesi, and L. Karlsson, The software landscape for the density matrix renormalization group, *arXiv preprint arXiv:2506.12629* (2025).
- [47] N. Nakatani and G. K. Chan, Efficient tree tensor network states (TTNS) for quantum chemistry: Generalizations of the density matrix renormalization group algorithm, *J. Chem. Phys.* **138**, 134113 (2013).
- [48] M. M. Rams and M. Zwolak, Breaking the entanglement barrier: Tensor network simulation of quantum transport, *Phys. Rev. Lett.* **124**, 137701 (2020).
- [49] L. Tagliacozzo, G. Evenbly, and G. Vidal, Simulation of two-dimensional quantum systems using a tree tensor network that exploits the entropic area law, *Phys. Rev. B* **80**, 235127 (2009).
- [50] J. Eisert, M. Cramer, and M. B. Plenio, Colloquium: Area laws for the entanglement entropy, *Rev. Mod. Phys.* **82**, 277 (2010).
- [51] G. K.-L. Chan, An algorithm for large scale density matrix renormalization group calculations, *J. Comp. Phys.* **120**, 3172 (2004).
- [52] A. Weichselbaum, Non-abelian symmetries in tensor networks: A quantum symmetry space approach, *Ann. Phys.* **327**, 2972 (2012).
- [53] Q. Li, Y. Gao, Y.-Y. He, Y. Qi, B.-B. Chen, and W. Li, Tangent space approach for thermal tensor network simulations of the 2D Hubbard model, *Phys. Rev. Lett.* **130**, 226502 (2023).
- [54] Q. Dong, L. Li, D. Dai, C. Zheng, J. Ma, R. Li, H. Xia, J. Xu, Z. Wu, B. Chang, *et al.*, A survey on in-context learning, in *Proceedings of the 2024 conference on empirical methods in natural language processing* (2024) pp. 1107–1128.
- [55] J. Ren, W. Li, T. Jiang, Y. Wang, C. Gong, and Z. Shuai, The RENORMALIZER package. <https://github.com/shuaigroup/renormalizer> (2025).

- [56] J. Ren, Z. Shuai, and G. Kin-Lic Chan, Time-dependent density matrix renormalization group algorithms for nearly exact absorption and fluorescence spectra of molecular aggregates at both zero and finite temperature, **14**, 5027 (2018).
- [57] W. Li, J. Ren, and Z. Shuai, Finite-temperature TD-DMRG for the carrier mobility of organic semiconductors, **11**, 4930 (2020).
- [58] W. Li, J. Ren, and Z. Shuai, A general charge transport picture for organic semiconductors with nonlocal electron-phonon couplings, **12**, 4260 (2021).
- [59] Y. Wang, J. Ren, and Z. Shuai, Minimizing non-radiative decay in molecular aggregates through control of excitonic coupling, **14**, 5056 (2023).
- [60] Z. Sheng, T. Jiang, W. Li, and Z. Shuai, TD-DMRG study of exciton dynamics with both thermal and static disorders for fenna-matthews-olson complex, **20**, 6470 (2024).
- [61] M. Fishman, S. White, and E. M. Stoudenmire, The ITensor software library for tensor network calculations, SciPost Phys. Codebases , 004 (2022).
- [62] J. Hauschild and F. Pollmann, Efficient numerical simulations with tensor networks: Tensor network Python (TeNPy), SciPost Phys. Lect. Notes , 5 (2018).
- [63] H. Zhai, H. R. Larsson, S. Lee, Z.-H. Cui, T. Zhu, C. Sun, L. Peng, R. Peng, K. Liao, J. Tölle, *et al.*, Block2: A comprehensive open source framework to develop and apply state-of-the-art DMRG algorithms in electronic structure and beyond, J. Comp. Phys. **159**, 234801 (2023).
- [64] Z. Liu, N. Lyu, Z. Hu, H. Zeng, V. S. Batista, and X. Sun, Benchmarking various nonadiabatic semiclassical mapping dynamics methods with tensor-train thermo-field dynamics, J. Comp. Phys. **161** (2024).
- [65] J. Ren, W. Li, T. Jiang, and Z. Shuai, A general automatic method for optimal construction of matrix product operators using bipartite graph theory, J. Comp. Phys. **153**, 084118 (2020).
- [66] W. Li, J. Ren, H. Yang, H. Wang, and Z. Shuai, Optimal tree tensor network operators for tensor network simulations: Applications to open quantum systems, J. Chem. Phys. **161** (2024).
- [67] J. Ren, W. Li, T. Jiang, Y. Wang, and Z. Shuai, Time-dependent density matrix renormalization group method for quantum dynamics in complex systems, **12**, e1614 (2022).
- [68] M. d. C. de Jongh and J. Van Leeuwen, Critical behavior of the two-dimensional ising model in a transverse field: A density-matrix renormalization calculation, Phys. Rev. B **57**, 8494 (1998).

- [69] H. W. Blöte and Y. Deng, Cluster monte carlo simulation of the transverse ising model, *Phys. Rev. E* **66**, 066110 (2002).
- [70] A. J. Leggett, S. Chakravarty, A. T. Dorsey, M. P. Fisher, A. Garg, and W. Zwerger, Dynamics of the dissipative two-state system, *Rev. Mod. Phys.* **59**, 1 (1987).
- [71] H. Wang and M. Thoss, From coherent motion to localization: II. dynamics of the spin-boson model with sub-Ohmic spectral density at zero temperature, *Chem. Phys.* **370**, 78 (2010).
- [72] L. Wang, L. Chen, N. Zhou, and Y. Zhao, Variational dynamics of the sub-Ohmic spin-boson model on the basis of multiple Davydov D1 states, *J. Comp. Phys.* **144**, 024101 (2016).
- [73] C. Duan, Z. Tang, J. Cao, and J. Wu, Zero-temperature localization in a sub-ohmic spin-boson model investigated by an extended hierarchy equation of motion, *Phys. Rev. B* **95**, 214308 (2017).
- [74] F. Otterpohl, P. Nalbach, and M. Thorwart, Hidden phase of the spin-boson model, *Phys. Rev. Lett.* **129**, 120406 (2022).
- [75] S. Hahn and G. Stock, Femtosecond secondary emission arising from the nonadiabatic photoisomerization in rhodopsin, *Chem. Phys.* **259**, 297 (2000).
- [76] N. Lyu, M. B. Soley, and V. S. Batista, Tensor-train split-operator KSL (TT-SOKSL) method for quantum dynamics simulations, **18**, 3327 (2022).
- [77] M. H. Beck, A. Jäckle, G. A. Worth, and H.-D. Meyer, The multiconfiguration time-dependent Hartree (MCTDH) method: a highly efficient algorithm for propagating wavepackets, **324**, 1 (2000).