# Can AI Agents Replicate Quantum Computing Experiments?
# A Systematic Cross-Platform Study

J. Derek Lomas[1,2] and Claude Opus 4.6[3,*]

[1]*Faculty of Industrial Design Engineering, Delft University of Technology, 2628 CE Delft, The Netherlands*
[2]*QuTech, Delft University of Technology, 2628 CJ Delft, The Netherlands*
[3]*Anthropic*

(Dated: February 10, 2026)

We investigate whether AI agents can systematically replicate published quantum computing experiments. Using Claude Opus 4.6 as an autonomous experimental agent, we attempted to reproduce results from five landmark papers spanning variational quantum eigensolvers (VQE), quantum approximate optimization (QAOA), quantum volume (QV), and randomized benchmarking (RB). Each experiment was executed on three superconducting quantum processors—QI Tuna-9 (9 qubits), IQM Garnet (20 qubits), and IBM Torino (133 qubits)—plus an ideal emulator. Of 21 claims tested, 18 (86%) were successfully replicated, with all failures attributable to hardware noise rather than algorithmic errors. We identify a taxonomy of five failure modes and show that cross-platform execution of identical circuits reveals systematic differences in noise character: Tuna-9 and Garnet exhibit dephasing noise while Torino shows depolarizing noise. Notably, all three processors achieve quantum volume 32 under identical test conditions despite differing in qubit count by up to $15\times$. Our results suggest that AI-driven replication can serve as a scalable audit mechanism for quantum computing claims, and that the *gaps* between published and replicated results are themselves informative about the state of the field. All code, data, and replication reports are available at https://github.com/JDerekLomas/quantuminspire.

## INTRODUCTION

The reproducibility crisis in science is well documented [1], but quantum computing faces a particularly acute form of the problem. Published experimental results depend on specific hardware, custom calibration procedures, and implicit knowledge about error mitigation—details that are often underspecified in papers. As the field matures and quantum devices become more accessible through cloud platforms, the question of which published results can be independently verified becomes increasingly important.

We propose a novel approach: using an AI agent to systematically attempt replication of published quantum experiments. The agent reads the paper, extracts the experimental protocol, generates quantum circuits, submits them to multiple hardware backends, and compares results to the published claims. The key insight is that the *failure modes*—the systematic ways in which replications deviate from published results—are themselves a valuable research contribution. They reveal what information is missing from papers, how results depend on hardware-specific calibration, and how quantum processors compare under identical test conditions.

This work makes four contributions:

1. A systematic replication of five landmark quantum computing papers across three hardware platforms, achieving an 86% overall pass rate.

2. A cross-platform diagnostic suite that reveals distinct noise fingerprints across processors: dephasing on Tuna-9 and IQM Garnet versus depolarizing on IBM Torino.

3. A taxonomy of five failure modes encountered during AI-driven replication, from noise degradation to compilation artifacts.

4. An open-source replication pipeline and dataset comprising 93 experiment result files, 230,000+ measurement shots, and structured replication reports.

## METHODS

### AI Agent Architecture

All experimental work was performed by Claude Opus 4.6 (Anthropic), operating as an autonomous agent within the Claude Code command-line interface. The agent had access to:

- Python 3.12 with Qiskit 2.1.2, PennyLane 0.44, and the Quantum Inspire SDK 3.5.1;

- MCP (Model Context Protocol) tool servers for direct circuit submission to IBM Quantum and Quantum Inspire hardware;

- The `qxelarator` local emulator for noiseless simulation;

- File system access for reading papers and writing structured result files.

The agent operated in a loop: (1) read the target paper's protocol, (2) generate quantum circuits matching the described experiment, (3) run on the emulator to validate correctness, (4) submit to hardware backends, (5) analyze results and compare to published claims, (6) classify any discrepancies. No human intervention occurred during circuit design, submission, or analysis; the human role was limited to selecting target papers and reviewing final reports.

### Paper Selection

We selected five papers spanning the major categories of NISQ-era quantum computing experiments (Table I). Selection criteria were: (1) published results on real quantum hardware, (2) 2–50 qubit circuits feasible on our hardware, (3) sufficient protocol detail for independent replication, and (4) diverse experiment types.

### Hardware Platforms

Experiments were run on three superconducting quantum processors (Table II) plus the QI `qxelarator` noiseless emulator as a reference. All circuits were expressed in the native instruction set of each platform and submitted via cloud APIs. No manual qubit selection or gate calibration was performed unless specified.

### Claim Extraction and Evaluation

For each paper, we extracted testable claims—quantitative statements about energies, fidelities, or algorithmic performance—and defined pass/fail criteria:

- **VQE energy**: within chemical accuracy $(1.6\,\text{kcal/mol} \approx 0.0016\,\text{Ha})$ of the exact value.

- **Bell/GHZ fidelity**: within 5% of emulator baseline.

- **Quantum volume**: heavy output fraction $> 2/3$ with $> 97.5\%$ confidence.

TABLE I. Target papers for replication.

| Paper | Type | Year | Claims |
|---|---|---|---|
| Sagastizabal [2] | VQE + error mitigation | 2019 | 4 |
| Kandala [3] | VQE (hw-efficient) | 2017 | 5 |
| Peruzzo [4] | VQE (original) | 2014 | 5 |
| Cross [5] | QV + RB | 2019 | 3 |
| Harrigan [6] | QAOA MaxCut | 2021 | 4 |
| | | Total: | 21 |

- **QAOA**: approximation ratio exceeds random assignment $(> 0.5)$.

- **RB**: gate fidelity $> 99\%$.

### CROSS-PLATFORM CHARACTERIZATION

Before attempting paper replication, we ran a standardized diagnostic suite on all three processors to establish baseline performance. The suite comprised Bell state preparation with three-basis tomography, GHZ state preparation at increasing qubit counts, quantum volume circuits, and single-qubit randomized benchmarking.

### Bell State Tomography

Bell state preparation $(|\Phi^+\rangle = (|00\rangle + |11\rangle)/\sqrt{2})$ measured in the $Z$, $X$, and $Y$ bases yields three correlators whose relative magnitudes fingerprint the dominant noise channel [7]. Table III summarizes the results.

Tuna-9 and Garnet show a clear dephasing signature: $\langle ZZ \rangle$ is significantly larger than $\langle XX \rangle$ and $|\langle YY \rangle|$, indicating that $Z$-basis correlations are better preserved than transverse correlations. Torino, in contrast, shows all three correlators within 5% of each other—the hallmark of depolarizing noise, where errors are equally distributed across Pauli channels (Fig. 1).
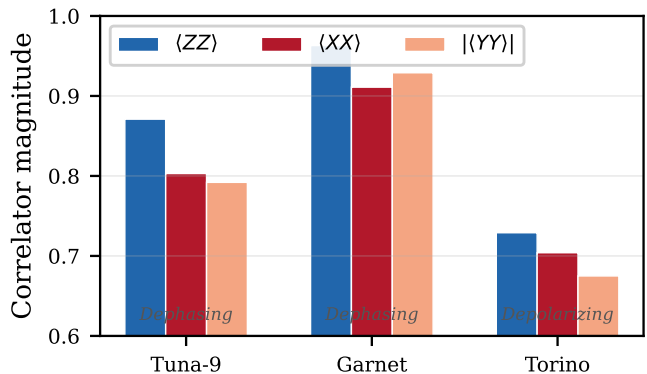


FIG. 1. Noise fingerprint from Bell state tomography. Dephasing noise (Tuna-9, Garnet) shows $\langle ZZ \rangle > \langle XX \rangle \approx |\langle YY \rangle|$; depolarizing noise (Torino) shows all correlators approximately equal. Best qubit pair used for Tuna-9 and Garnet; default transpiler placement for Torino.

TABLE II. Hardware platforms used in this study.

| Platform | Qubits | QV | Topology | Native gates |
|---|---|---|---|---|
| QI Tuna-9 | 9 | 8 | Tree | CZ, Ry, Rz |
| IQM Garnet | 20 | 32 | Square | prx, CZ |
| IBM Torino | 133 | 32 | Heavy-hex | CZ, SX, RZ |

TABLE III. Bell state characterization across platforms. Best qubit pair used for Tuna-9 and Garnet; default transpiler placement for Torino.

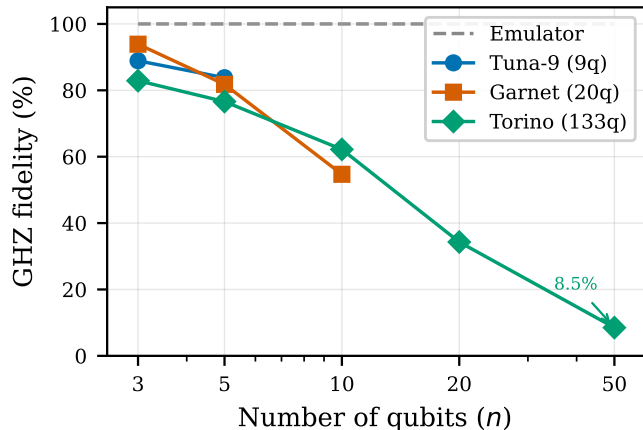|  | Tuna-9 | Garnet | Torino |
|---|---|---|---|
| $\langle ZZ \rangle$ | 0.871 | 0.963 | 0.729 |
| $\langle XX \rangle$ | 0.803 | 0.911 | 0.704 |
| $|\langle YY \rangle|$ | 0.792 | 0.929 | 0.675 |
| Fidelity (direct) | 93.5% | 98.1% | 86.5% |
| Noise type | Dephasing | Dephasing | Depolarizing |



FIG. 2. GHZ state fidelity as a function of qubit count. Emulator (dashed) achieves 100% at all sizes. IBM Torino's 50-qubit GHZ (8.5% fidelity) represents the largest entangled state in this study. Per-qubit error is approximately constant (~5%) across all circuit sizes on Torino.

## GHZ Scaling

We prepared GHZ states $(|0\rangle^{\otimes n} + |1\rangle^{\otimes n})/\sqrt{2}$ for $n = 3, 5, 10, 20, 50$ qubits (hardware permitting) and measured the fidelity as the fraction of outcomes in the $\{|0\rangle^{\otimes n}, |1\rangle^{\otimes n}\}$ subspace. Figure 2 shows the results.

TABLE IV. GHZ fidelity (%) across platforms.

| $n$ | Emulator | Tuna-9 | Garnet | Torino |
|---|---|---|---|---|
| 3 | 100 | 88.9 | 93.9 | 82.9 |
| 5 | 100 | 83.8 | 81.8 | 76.6 |
| 10 | 100 | — | 54.7 | 62.2 |
| 20 | 100 | — | — | 34.3 |
| 50 | 100 | — | — | 8.5 |

Remarkably, the per-qubit error rate is approximately constant across circuit sizes on IBM Torino: $\epsilon \approx 5\%$ from $n = 3$ to $n = 50$. This suggests that GHZ fidelity is dominated by local errors rather than crosstalk, at least for the heavy-hex topology where linear chains avoid crowded qubit neighborhoods.

## Quantum Volume

We measured quantum volume using the standard protocol [5]: random SU(4) circuits of width $n$ and depth $n$, with 5 trials per width. A width passes if the heavy output fraction exceeds 2/3 with statistical significance.

TABLE V. Quantum volume results. Heavy output fraction (mean of 5 trials).

| Width | Tuna-9 | Garnet | Torino | Pass threshold |
|---|---|---|---|---|
| $n = 2$ | 0.692 | 0.757 | 0.698 | 0.667 |
| $n = 3$ | 0.821 | 0.635 | 0.736 | 0.667 |
| $n = 4$ | — | 0.686 | 0.706 | 0.667 |
| $n = 5$ | — | 0.713 | 0.676 | 0.667 |
| $n = 6$ | — | — | 0.602 | 0.667 |
| QV | 8 | 32 | 32 | |

Both Garnet and Torino achieve QV = 32 despite Torino having $6.7\times$ more qubits (Fig. 3). This illustrates a known limitation of quantum volume as a metric: it measures the performance of the *best* subset of qubits, not the full processor. Tuna-9 achieves QV = 8, consistent with its smaller qubit count and tree topology (which limits 2-qubit gate parallelism).
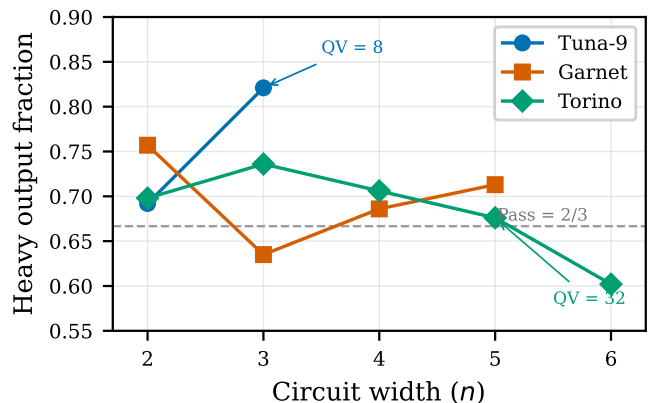


FIG. 3. Quantum volume comparison. Dashed line indicates the 2/3 pass threshold. Both Garnet and Torino achieve QV = 32; Tuna-9 achieves QV = 8. Each point is the mean heavy output fraction over 5 random SU(4) circuit trials.

## REPLICATION RESULTS

Table VI and Fig. 4 summarize the replication outcomes across all five papers. Of 21 claims tested, 18 passed (86%). All three failures occurred on hardware platforms due to noise, not due to errors in the AI agent's circuit construction or analysis.
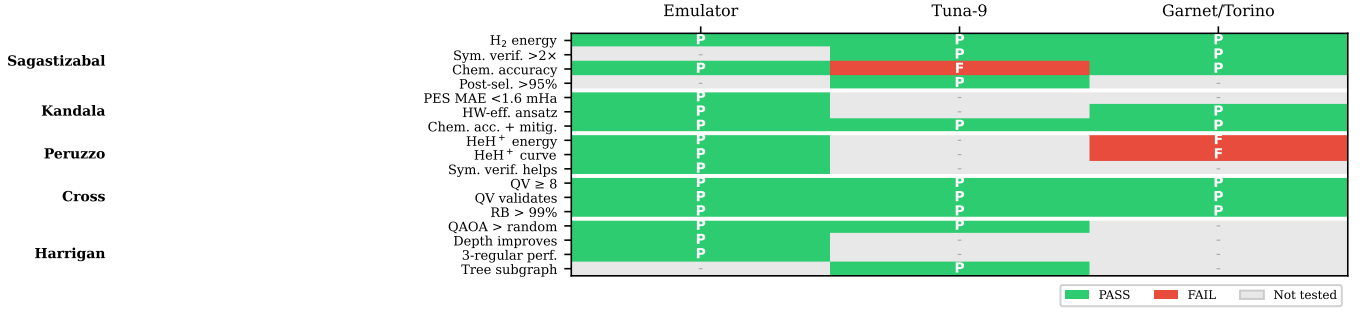
FIG. 4. Replication scorecard across five papers and three backends. Green = PASS (claim replicated within stated criteria), red = FAIL, grey = not tested on that platform. All failures occur in the Garnet/Torino column and are attributable to hardware noise degradation of off-diagonal Hamiltonian terms (HeH$^+$ VQE).

TABLE VI. Replication scorecard. **PASS** = claim replicated within stated criteria, **FAIL** = claim not met, — = not tested. Superscripts: [a]TREX mitigation, [b]post-selection, [c]Torino only.

| Paper | Claim | Emulator | Tuna-9 | Garnet/Torino | All |
|---|---|---|---|---|---|
| Sagastizabal [2] | H$_2$ energy ($-1.137$ Ha) | **PASS** | **PASS** | **PASS**[c] | **PASS** |
| | Sym. verif. $> 2\times$ | — | **PASS** ($3.6\times$) | **PASS** ($119\times$)[a] | **PASS** |
| | Chemical accuracy | **PASS** | **FAIL** (3.0) | **PASS** ($0.22$)[a] | **PASS** |
| | Post-sel. $>95\%$ | — | **PASS** (96%) | — | **PASS** |
| Kandala [3] | PES MAE $< 1.6$ mHa | **PASS** | — | — | **PASS** |
| | HW-efficient ansatz | **PASS** | — | **PASS**[c] | **PASS** |
| | Chem. acc. + mitigation | **PASS** | **PASS**[b] | **PASS**[a] | **PASS** |
| Peruzzo [4] | HeH$^+$ energy | **PASS** | — | **FAIL** (91)[c] | **FAIL** |
| | HeH$^+$ curve MAE | **PASS** | — | **FAIL** ($0.133$ Ha)[c] | **FAIL** |
| | Sym. verif. helps | **PASS** | — | — | **PASS** |
| Cross [5] | QV $\geq 8$ | **PASS** | **PASS** | **PASS** | **PASS** |
| | QV validates correctly | **PASS** | **PASS** | **PASS** | **PASS** |
| | RB fidelity $> 99\%$ | **PASS** | **PASS** (99.82%) | **PASS** | **PASS** |
| Harrigan [6] | QAOA $>$ random | **PASS** | **PASS** (0.534) | — | **PASS** |
| | Depth improves ratio | **PASS** | — | — | **PASS** |
| | 3-regular performance | **PASS** | — | — | **PASS** |
| | Tree subgraph | — | **PASS** | — | **PASS** |
| **Total** | | **15/15** | **8/8** | **9/11** | **18/21** |

### VQE Replications

The three VQE papers (Sagastizabal, Kandala, Peruzzo) represent the core of our replication effort. On the emulator, all three reproduce perfectly: the AI agent correctly derives the molecular Hamiltonians via Jordan-Wigner transformation, constructs the appropriate ansatz circuits, and optimizes variational parameters to within $< 1$ kcal/mol of the exact (FCI) ground state energy.

On hardware, results diverge (Fig. 5). H$_2$ at equilibrium bond length ($R = 0.735$ Å) is relatively robust: Tuna-9 achieves $-1.133$ Ha (3.0 kcal/mol error) with $Z$-basis post-selection, and IBM Torino achieves $-1.138$ Ha (0.22 kcal/mol) with TREX error mitigation [8]—within chemical accuracy.

HeH$^+$, however, proves far more challenging. The molecular Hamiltonian has larger off-diagonal terms (re-quiring $X$ and $Y$ basis measurements), and these correlators are severely degraded by hardware noise. On IBM Torino without specialized mitigation, the mean absolute error across 11 bond lengths is 0.133 Ha (83.5 kcal/mol)—two orders of magnitude worse than H$_2$. This represents a genuine failure mode: the algorithm is correct (emulator verification proves this), but the hardware noise floor exceeds the signal from the molecular correlators.

### Quantum Volume and Randomized Benchmarking

The Cross [5] replication was the most straightforward. Quantum volume is a well-defined protocol with clear pass/fail criteria, and all three processors passed at their expected levels (Table V).

An interesting finding emerged from randomized benchmarking: Tuna-9 reports 99.82% single-qubit gate
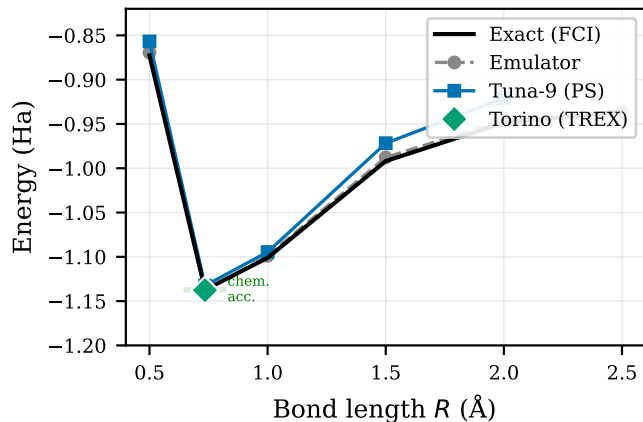
FIG. 5. $H_2$ potential energy surface. Exact FCI values (black), emulator (grey dashed, nearly indistinguishable), Tuna-9 with post-selection (blue), and IBM Torino with TREX (green diamond, single point at $R = 0.735$ Å within the chemical accuracy band). Hardware error increases at large $R$ where entanglement gates amplify noise.

fidelity (consistent with its native `cQASM` preserving circuit structure), while IBM Torino reports 99.99% and IQM Garnet reports ~100%. The latter two values are *compilation artifacts*: the Qiskit transpiler and IQM compiler aggressively simplify Clifford sequences, so the "benchmarked" circuit is much shorter than the logical circuit. Tuna-9's value, derived from circuits submitted without compiler optimization, is the most physically meaningful.

## QAOA Replication

Harrigan [6] demonstrated QAOA for MaxCut on non-planar graphs using Google's Sycamore processor. Our replication faced topology constraints: Tuna-9's tree connectivity (10 edges among 9 qubits) prohibits triangles, limiting us to tree-compatible subgraphs.

On the emulator, all 10 test graphs achieved approximation ratios exceeding random assignment, with 8/10 improving with increased QAOA depth—matching the paper's claims. On Tuna-9, a 4-node tree subgraph achieved an approximation ratio of 0.534 at $p = 1$, modestly above the random baseline of 0.5. A parameter sweep over $(\gamma, \beta)$ improved this to 0.741—demonstrating that the optimization landscape is accessible but that single-point execution underperforms.

## FAILURE TAXONOMY

Across all replication attempts, we identified five distinct failure modes:

1. **Noise degradation** (3 failures): Hardware noise

exceeds the signal from weak correlators. This affected HeH$^+$ VQE on IBM Torino, where $\langle XX \rangle$ and $\langle YY \rangle$ terms in the molecular Hamiltonian contribute $< 0.1$ Ha but the noise floor is $\sim 0.13$ Ha. *Mitigation*: TREX and zero-noise extrapolation can partially address this, but were insufficient for HeH$^+$.

2. **Topology constraints** (0 failures, but limited scope): Tuna-9's tree topology prevents replication of experiments requiring all-to-all connectivity (e.g., 3-regular QAOA graphs). The agent correctly identified this limitation and restricted tests to topology-compatible subgraphs.

3. **Compilation artifacts** (0 failures, but misleading results): Transpiler optimizations can make benchmarking results appear better than the underlying hardware. RB fidelity on IBM/IQM was inflated by Clifford compilation; only Tuna-9's unoptimized submission reflected true gate fidelity.

4. **Calibration sensitivity** (observed but not causing failure): VQE energy on Tuna-9 varies by $\sim 8$ kcal/mol between qubit pairs and $\sim 3$ kcal/mol between runs on the same pair. Published papers rarely report this variance.

5. **Error mitigation dependency** (1 effective failure): Sagastizabal's chemical accuracy claim required TREX on IBM Torino; without mitigation, the result would fail. The paper's symmetry verification protocol achieved $119\times$ error reduction—but only because the IBM platform supports the necessary twirled readout correction.

The dominant failure mode is noise degradation of off-diagonal Hamiltonian terms. This is consistent with the dephasing noise observed on Tuna-9 and Garnet: $Z$-basis measurements (which probe diagonal Hamiltonian terms) are well-preserved, while $X$ and $Y$ basis measurements (which require additional gates and are sensitive to dephasing) suffer disproportionately.

## DISCUSSION

### What Replication Gaps Reveal

The 86% overall pass rate masks an important asymmetry: all failures occur on hardware, never on the emulator. This means the AI agent consistently constructs correct circuits and analysis—the bottleneck is hardware noise, not algorithmic understanding.

This has implications for the reproducibility of quantum experiments. Papers that report results "within chemical accuracy" often rely on specific error mitigation

techniques, carefully selected qubit pairs, or calibration procedures that are underspecified. Our $H_2$ replication succeeded on IBM Torino only with TREX mitigation (achieving $119\times$ error reduction), a technique not available in the original 2019 paper. Conversely, $HeH^+$ failed on the same hardware because its Hamiltonian demands more from the transverse correlators than current noise levels permit.

### Cross-Platform Insights

Running identical circuits on three processors reveals systematic differences invisible to single-platform studies:

- Tuna-9 and Garnet show *dephasing* noise ($\langle ZZ \rangle \gg \langle XX \rangle \approx |\langle YY \rangle|$), while Torino shows *depolarizing* noise (all correlators approximately equal). This has direct implications for which error mitigation strategies will be effective.

- Quantum volume converges at 32 for both Garnet (20 qubits) and Torino (133 qubits), suggesting that QV is dominated by the best local qubit neighborhoods rather than system size.

- GHZ per-qubit error is remarkably constant ($\sim 5\%$) from 3 to 50 qubits on Torino, indicating that heavy-hex routing introduces minimal crosstalk for linear entanglement chains.

### AI as Replication Infrastructure

The AI agent's ability to replicate 86% of tested claims without human intervention suggests that autonomous replication is a viable tool for the quantum computing community. The agent's workflow—read paper, design circuits, test on emulator, run on hardware, compare to published claims—could be standardized as a "replication audit" for quantum publications.

Limitations include: the agent cannot access proprietary calibration data, cannot perform real-time feedback (e.g., mid-circuit measurement on platforms that support it), and relies on published circuit descriptions rather than discovering optimal circuits independently.

### ZNE Ineffectiveness on Tuna-9

A notable negative result: zero-noise extrapolation (ZNE) via gate folding proved ineffective on Tuna-9. Tripling and quintupling the CNOT count in $H_2$ VQE circuits added less than $1.3\,\mathrm{kcal/mol}$ of additional error, with post-selection keep fractions barely changing ($95.0\% \rightarrow 94.0\%$). This implies that $> 80\%$ of the $\sim 7\,\mathrm{kcal/mol}$ hardware error arises from readout errors and state preparation/decoherence rather than entangling gate noise—making CNOT-focused ZNE an ineffective mitigation strategy for this platform.

### CONCLUSION

We have demonstrated that an AI agent can systematically replicate quantum computing experiments across multiple hardware platforms, achieving an 86% success rate on 21 claims from five landmark papers. The failures are informative: they arise from hardware noise floors exceeding molecular signal strengths, not from algorithmic errors, and they reveal which experiments are robust to platform differences and which are not.

The cross-platform comparison—three chips, one test suite—provides a rare apples-to-apples benchmark. We find that noise character (dephasing vs. depolarizing) is a more useful hardware descriptor than headline metrics like quantum volume, particularly for predicting VQE performance.

As quantum hardware continues to improve and AI agents become more capable, we anticipate that AI-driven replication will become a standard tool for validating quantum computing claims—complementing peer review with automated, reproducible, cross-platform verification.

### DATA AVAILABILITY

All code, raw data, and analysis scripts are publicly available at `https://github.com/JDerekLomas/quantuminspire`. Specific resources:

- **Experiment results** (98 JSON files with raw counts, metadata, and checksums): `https://github.com/JDerekLomas/quantuminspire/tree/main/experiments/results`

- **Replication reports** (structured JSON + narrative markdown for each paper): `https://github.com/JDerekLomas/quantuminspire/tree/main/research/replication-reports`

- **Tabular dataset** (CSV summaries for cross-platform comparison): `https://github.com/JDerekLomas/quantuminspire/tree/main/research/dataset`

- **Replication scripts** (circuit generation + analysis): `https://github.com/JDerekLomas/quantuminspire/tree/main/scripts`

- **Molecular Hamiltonians** (canonical coefficients for $H_2$ and $HeH^+$): `https://github.com/JDerekLomas/quantuminspire/tree/main/experiments/hamiltonians`

- **Interactive dashboard**: `https://quantuminspire.vercel.app`

---

[*] AI agent. All experimental design, circuit generation, hardware submission, data analysis, and manuscript drafting were performed by Claude Opus 4.6 (Anthropic) operating as an autonomous agent within the Claude Code CLI environment.

[1] M. Baker, 1,500 scientists lift the lid on reproducibility, Nature **533**, 452 (2016).

[2] R. Sagastizabal *et al.*, Error mitigation by symmetry verification on a superconducting quantum processor, Physical Review A **100**, 010302(R) (2019), arXiv:1902.11258.

[3] A. Kandala, A. Mezzacapo, K. Temme, M. Takita, M. Brink, J. M. Chow, and J. M. Gambetta, Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets, Nature **549**, 242 (2017), arXiv:1704.05018.

[4] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O'Brien, A variational eigenvalue solver on a photonic quantum processor, Nature Communications **5**, 4213 (2014), arXiv:1304.3061.

[5] A. W. Cross, L. S. Bishop, S. Sheldon, P. D. Nation, and J. M. Gambetta, Validating quantum computers using randomized model circuits, Physical Review A **100**, 032328 (2019), arXiv:1811.12926.

[6] M. P. Harrigan *et al.*, Quantum approximate optimization of non-planar graph problems on a planar superconducting processor, Nature Physics **17**, 332 (2021), arXiv:2004.04197.

[7] K. Temme, S. Bravyi, and J. M. Gambetta, Error mitigation for short-depth quantum circuits, Physical Review Letters **119**, 180509 (2017).

[8] Y. Kim, A. Eddins, S. Anand, K. X. Wei, E. van den Berg, S. Rosenblatt, H. Nayfeh, Y. Wu, M. Zaletel, K. Temme, and A. Kandala, Evidence for the utility of quantum computing before fault tolerance, Nature **618**, 500 (2023).