

# Critique of "GPU-Accelerated Jump Diffusion HJB Equations for Low-Latency Crypto Market Making with Order Flow Toxicity Tracking"



*Prepared by Critique Author*  
[critique@arithmax.com](mailto:critique@arithmax.com)  
November 28, 2025

## Abstract

This document provides a deep critique of the research paper titled "GPU-Accelerated Jump Diffusion HJB Equations for Low-Latency Crypto Market Making with Order Flow Toxicity Tracking". The critique evaluates the paper's strengths, weaknesses, technical accuracy, and practical implications, offering constructive feedback for improvement.

## Index Terms

Research Critique, Market Making, HJB Equations, GPU Computing

### I. OVERALL ASSESSMENT

The paper presents an ambitious attempt to bridge theoretical optimal control theory with practical high-frequency trading implementation in cryptocurrency markets. While it demonstrates impressive technical depth and addresses an important problem, it suffers from several critical flaws in mathematical rigor, empirical validation, and practical applicability. The mathematical framework is sound in concept but flawed in execution, while the empirical results appear overly optimistic and insufficiently validated. The GPU implementation claims are impressive but lack the detail needed to assess their realism. Overall, this work contributes meaningfully to the literature but requires significant revision to be credible—potentially publishable in a top finance journal after addressing the issues below.

### II. STRENGTHS

- 1) **Clear Problem Motivation:** The introduction effectively highlights cryptocurrency market making's unique challenges (extreme volatility, fragmentation, jumps), setting up a compelling case for advanced modeling. The emphasis on microstructure (toxicity tracking, order flow) is a strength, as traditional models often ignore these factors.
- 2) **Theoretical Integration:** The synthesis of HJB equations with jump diffusion is a solid extension of existing work (e.g., Avellaneda-Stoikov, Guéant et al.). The inclusion of inventory risk, market impact, and adaptive toxicity adjustment shows thoughtful consideration of real-world trading dynamics.
- 3) **Computational Innovation:** The GPU acceleration approach is forward-thinking, leveraging parallel computing for real-time PDE solving. The performance benchmarks (136-177x speedup) are impressive and demonstrate practical potential, even if the absolute times seem optimistic.
- 4) **Structure and Presentation:** The paper is well-organized, with clear sections on theory, implementation, and results. The use of algorithms, figures, and tables enhances readability, and the bibliography is comprehensive.
- 5) **Interdisciplinary Appeal:** It bridges finance, applied math, and computer science, potentially appealing to both theorists and practitioners.

### III. WEAKNESSES

#### A. Mathematical Rigor and Accuracy

- The HJB equation formulation (Equation 5) has potential issues with the max operators over bid/ask quotes. The discretization scheme (backward induction with finite differences) is standard but the jump integral approximation is crude and likely inaccurate. The 5-point quadrature for the jump term (Algorithm 3) is insufficient for convergence, especially for fat-tailed jump distributions common in crypto.

- The order execution model (Equations 3-4) is overly simplistic. The piecewise linear intensity function doesn't capture real market dynamics (e.g., queue position effects, flash crashes). The toxicity adjustment (Equation 7) is ad-hoc and lacks theoretical justification.
- Boundary conditions and terminal conditions are mentioned but not rigorously specified, which could lead to numerical instability.

#### B. Implementation and Code Quality

- The algorithms are presented as pseudocode but contain errors (e.g., in Algorithm 3, the indexing for  $V_{\text{next}}[\text{idx}, j]$  assumes a 1D array but the grid is 2D; the jump term calculation is incorrect). Real CUDA code would need careful memory management, which isn't discussed.
- No mention of numerical stability, convergence criteria, or error bounds. The shared memory optimization (Algorithm 4) is good in theory but ignores GPU occupancy limits and bank conflicts.
- The "FPGA Implementation Considerations" section is misplaced and underdeveloped—FPGAs aren't actually implemented, and the claims about parallelization are generic.

#### C. Empirical Validation

- The backtesting results are suspiciously strong: a Sharpe ratio of 2.26 vs. 1.85 for Avellaneda-Stoikov is unrealistic for crypto markets, where even sophisticated strategies rarely exceed 1.5-2.0. The one-month testing period is inadequate—crypto markets require multi-year, out-of-sample validation to account for regime changes.
- No discussion of overfitting, transaction costs, slippage, or exchange-specific fees (e.g., Binance's maker/taker fees). The performance figures (Figure 6, Table 2) lack statistical significance tests or confidence intervals.
- The value function surface (Figure 5) is synthetic and not validated against real data. The jump diffusion comparison (Figure 1) uses toy data that doesn't reflect actual BTC volatility.

#### D. Literature Review and Novelty

- While citations are relevant, the review misses key recent work: reinforcement learning approaches (e.g., Nevmyvaka et al., 2006; or more recent RL for market making), high-frequency crypto studies (e.g., on LOB dynamics), and PDE solvers for finance (e.g., Deep Galerkin Methods).
- The "contributions" claim jump diffusion extension as novel, but this is incremental—jump diffusion in HJB for market making has been explored (e.g., in option pricing contexts).

#### E. Practical and Scalability Concerns

- Real-time deployment claims ignore practical hurdles: WebSocket latency, API rate limits, exchange downtime, and regulatory constraints (e.g., SEC oversight for HFT in crypto).
- The system architecture (Figure 2) is oversimplified—no discussion of failover, risk limits, or integration with existing trading infrastructure.
- GPU performance claims (sub-millisecond for 101x101 grid) may be achievable in isolation but not in a full trading loop with data ingestion and order placement.

#### F. Writing and Clarity

- Some sections are verbose (e.g., the introduction repeats points). Mathematical notation is inconsistent (e.g., mixing  $S_t$  and  $S$ ).

- The abstract overstates impact ("significant reduction in inventory risk") without evidence.
- Figures are helpful but some (e.g., Figure 3) use placeholder data that doesn't demonstrate real toxicity effects.

#### IV. SPECIFIC TECHNICAL ISSUES

- **Equation Errors:** In the HJB equation (5), the jump term should be  $\lambda \int (V(t, S(1 + y), I) - V(t, S, I))f(y)dy$ , but the discretization (Equation 10) approximates this poorly. The max operators over continuous variables aren't properly discretized.
- **Algorithm Flaws:** Algorithm 1 has incorrect CUDA grid setup (blockspergrid should be tuples). Algorithm 3's jump operator doesn't handle boundary conditions for idx.
- **Performance Metrics:** The speedup ratios are plausible, but the absolute CPU times (e.g., 32.7ms for 51x51) suggest inefficient CPU implementation—optimized CPU code could be much faster.
- **Toxicity Model:** The exponential decay weighting (Equation 6) is reasonable, but the clipping and scaling (Equation 7) lack empirical calibration.

#### V. RECOMMENDATIONS

- 1) **Strengthen Mathematics:** Derive proper convergence proofs for the numerical scheme. Use more sophisticated quadrature (e.g., Gauss-Hermite) for jump integrals. Validate against analytical solutions for simpler cases.
- 2) **Improve Validation:** Conduct multi-year backtests with realistic costs. Include statistical tests, walk-forward analysis, and comparison to more benchmarks (e.g., RL-based strategies). Disclose all assumptions and limitations.
- 3) **Enhance Implementation Details:** Provide actual code snippets or GitHub links. Discuss error handling, GPU memory management, and integration with trading APIs (e.g., CCXT).
- 4) **Address Practicality:** Add sections on deployment challenges, risk management, and regulatory compliance. Consider co-location, low-latency networks, and exchange-specific optimizations.
- 5) **Expand Literature Review:** Include recent ML approaches and crypto-specific HFT studies. Position the work more clearly against state-of-the-art.
- 6) **Revise Claims:** Tone down overstatements (e.g., "microseconds" performance). Focus on incremental contributions rather than revolutionary ones.

#### VI. FINAL VERDICT

This paper is a solid PhD-level effort with genuine innovation in applying GPU computing to financial PDEs. However, it needs substantial revision to be publication-ready—particularly in empirical rigor and mathematical accuracy. With fixes, it could be a valuable contribution to algorithmic trading literature. As it stands, it's more of a promising technical report than a definitive research paper. Rating: 6/10 (strong potential, significant flaws).