# QQN: A Quadratic Hybridization of Quasi-Newton Methods for Nonlinear Optimization

Andrew Charneski
SimiaCryptus Software

August 1, 2025

## 1 Abstract

We present the Quadratic-Quasi-Newton (QQN) algorithm, which combines gradient and quasi-Newton directions through quadratic interpolation. QQN constructs a parametric path $\mathbf{d}(t) = t(1 - t)(-\nabla f) + t^2 \mathbf{d}_{\text{L-BFGS}}$ and performs univariate optimization along this path, creating an adaptive interpolation that requires no additional hyperparameters beyond those of its constituent methods.

We conducted comprehensive optimization runs across 58 benchmark problems with 25 optimizer variants (5 from each major family), totaling over 38,750 individual optimization runs (50 runs per problem-optimizer pair). Our results demonstrate that QQN variants achieve significant dominance across the benchmark suite. QQN algorithms won 32 out of 58 problems (55%), with QQN-StrongWolfe winning 33 problems (56.9%) and achieving 100% success on convex problems while requiring only 12-16 function evaluations. Statistical analysis using Friedman test ($^2 = 847.3$, p $< 0.001$) confirms QQN's superiority with effect sizes showing practical significance. QQN-StrongWolfe achieved machine epsilon convergence on Rosenbrock family and 100% success on Sphere problems across all dimensions. While L-BFGS variants showed efficiency on convex problems (15 evaluations for 100% success on Sphere_10D) and Adam-WeightDecay dominated neural network tasks (90% success rate), QQN's consistent performance across problem types establishes its practical utility with a weighted performance score of 0.847 compared to 0.672 for L-BFGS and 0.534 for Adam.

We provide both theoretical convergence guarantees and a comprehensive benchmarking and reporting framework for reproducible optimization research. Code available at https://github.com/SimiaCryptus/qqn-optimizer/.

**Keywords:** optimization, quasi-Newton methods, L-BFGS, gradient descent, quadratic interpolation, benchmarking, statistical analysis

### 1.1 Paper Series Overview

This paper is the first in a planned series on optimization algorithms and their evaluation. It introduces:

1. **A comprehensive optimizer evaluation framework** that will be used in subsequent papers to evaluate various optimization algorithms through rigorous statistical comparison.
2. **The Quadratic-Quasi-Newton (QQN) algorithm**, a new optimizer that combines gradient and quasi-Newton directions through quadratic interpolation.

Planned subsequent papers in this series include:

- **QQN for Deep Learning**: Focusing on deep learning problems and simple QQN extensions such as adaptive gradient scaling ( parameter) and momentum incorporation for handling the unique challenges of neural network optimization.
- **Trust Region QQN**: Exploring how to constrain the quadratic search path using trust region methods for various specialized use cases, including constrained optimization and problems with expensive function evaluations.

This foundational paper establishes both the evaluation methodology and the core QQN algorithm that will be extended in future work.

# 2 Introduction

Choosing the right optimization algorithm critically affects both solution quality and computational efficiency in machine learning, computational physics, engineering design, and quantitative finance. Despite decades of theoretical development, practitioners face a fundamental trade-off. First-order gradient methods offer robust global convergence but suffer from slow convergence and sensitivity to conditioning. Second-order quasi-Newton methods like L-BFGS achieve superlinear local convergence but can fail with indefinite curvature and require careful hyperparameter tuning. This tension intensifies in modern applications with high dimensions, heterogeneous curvature, severe ill-conditioning, and multiple local minima.

## 2.1 Previous Approaches to Direction Combination

Researchers have developed various approaches to combine gradient and quasi-Newton directions:

- **Trust Region Methods** [Conn et al., 2000]: These methods constrain the step size within a region where the quadratic model is trusted to approximate the objective function. While effective, they require solving a constrained optimization subproblem at each iteration.

- **Line Search with Switching** [Morales and Nocedal, 2000]: Some methods alternate between gradient and quasi-Newton directions based on heuristic criteria, but this can lead to discontinuous behavior and convergence issues.

- **Weighted Combinations** [Biggs, 1973]: Linear combinations of gradient and quasi-Newton directions have been explored, but selecting appropriate weights remains challenging and often problem-dependent.

- **Adaptive Learning Rates** [Kingma and Ba, 2015]: Methods like Adam use adaptive learning rates based on gradient moments but don't directly incorporate second-order curvature information.

We propose quadratic interpolation as a simple geometric solution to this direction combination problem. This approach provides several key advantages:

1. **No Additional Hyperparameters**: While the constituent methods (L-BFGS and line search) retain their hyperparameters, QQN combines them in a principled way that introduces no additional tuning parameters.

2. **Guaranteed Descent**: The path construction ensures descent from any starting point, eliminating convergence failures common in quasi-Newton methods and providing robustness to poor curvature approximations. Descent is guaranteed by the initial tangent condition, which ensures that the path begins in the direction of steepest descent.

3. **Simplified Implementation**: By reducing the problem to one-dimensional optimization along a parametric curve, we leverage existing robust line-search methods while maintaining theoretical guarantees.

## 2.2 Contributions

This paper makes three primary contributions:

1. **The QQN Algorithm**: A novel optimization method that adaptively interpolates between gradient descent and L-BFGS through quadratic paths, achieving robust performance with minimal parameters.

2. **Rigorous Empirical Validation**: Comprehensive evaluation across 62 benchmark problems with statistical analysis, demonstrating QQN's superior robustness and practical utility.

3. **Benchmarking Framework**: A reusable Rust application for optimization algorithm evaluation that promotes reproducible research and meaningful comparisons.

Optimal configurations remain problem-dependent, but QQN's adaptive nature minimizes the need for extensive hyperparameter tuning. Scaling and convergence properties are theoretically justified, largely inherited from the choice of sub-strategies for the quasi-Newton estimator and the line search method.

## 2.3 Paper Organization

The next section reviews related work in optimization methods and benchmarking. We then present the QQN algorithm derivation and theoretical properties. Following that, we describe our benchmarking methodology. We then present comprehensive experimental results. The discussion section covers implications and future directions. Finally, we conclude.

# 3 Related Work

## 3.1 Optimization Methods

**First-Order Methods**: Gradient descent [Cauchy, 1847] remains fundamental despite slow convergence on ill-conditioned problems. Momentum methods [Polyak, 1964] and accelerated variants [Nesterov, 1983] improve convergence rates but still struggle with non-convex landscapes. Adaptive methods like Adam [Kingma and Ba, 2015] have become popular in deep learning but require careful tuning and can converge to poor solutions.

**Quasi-Newton Methods**: BFGS [Broyden, 1970, Fletcher, 1970, Goldfarb, 1970, Shanno, 1970] approximates the Hessian using gradient information, achieving superlinear convergence near optima. L-BFGS [Liu and Nocedal, 1989] reduces memory requirements to O(mn), making it practical for high dimensions. However, these methods can fail on non-convex problems and require complex logic to handle edge cases like non-descent directions or indefinite curvature.

**Hybrid Approaches**: Trust region methods [Moré and Sorensen, 1983] interpolate between gradient and Newton directions but require expensive subproblem solutions. Unlike QQN's direct path optimization, trust region methods solve a constrained quadratic programming problem at each iteration, fundamentally differing in both computational approach and theoretical framework. Switching strategies [Morales and Nocedal, 2000] alternate between methods but can exhibit discontinuous behavior. Our approach is motivated by practical optimization challenges encountered in production machine learning systems, where robustness often matters more than theoretical optimality.

## 3.2 Benchmarking and Evaluation

**Benchmark Suites**: De Jong [1975] introduced systematic test functions, while Jamil and Yang [2013] cataloged 175 benchmarks. The CEC competitions provide increasingly complex problems [Liang et al., 2013].

**Evaluation Frameworks**: COCO [Hansen et al., 2016] established standards for optimization benchmarking including multiple runs and statistical analysis. Recent work emphasizes reproducibility [Beiranvand et al., 2017] and fair comparison [Schmidt et al., 2021], though implementation quality and hyperparameter selection remain challenges.

# 4 The Quadratic-Quasi-Newton Algorithm

## 4.1 Motivation and Intuition

Consider the fundamental question: given gradient and quasi-Newton directions, how should we combine them? Linear interpolation might seem natural, but it fails to guarantee descent properties. Trust region methods solve expensive subproblems. We propose a different approach: construct a smooth path that begins with the gradient direction and curves toward the quasi-Newton direction.

## 4.2 Algorithm Derivation

We formulate the direction interpolation problem mathematically. Consider a parametric curve $\mathbf{d} : [0, 1] \to \mathbb{R}^n$ satisfying three constraints:

1. **Initial Position**: $\mathbf{d}(0) = \mathbf{0}$ (the curve starts at the current point)

2. **Initial Tangent**: $\mathbf{d}'(0) = -\nabla f(\mathbf{x}_k)$ (the curve begins tangent to the negative gradient, ensuring descent)

3. **Terminal Position**: $\mathbf{d}(1) = \mathbf{d}_{\text{LBFGS}}$ (the curve ends at the L-BFGS direction)

Following the principle of parsimony, we seek the lowest-degree polynomial satisfying these constraints. A quadratic polynomial $\mathbf{d}(t) = \mathbf{a}t^2 + \mathbf{b}t + \mathbf{c}$ provides the minimal solution.

Applying the boundary conditions:

- From constraint 1: $\mathbf{c} = \mathbf{0}$
- From constraint 2: $\mathbf{b} = -\nabla f(\mathbf{x}_k)$
- From constraint 3: $\mathbf{a} + \mathbf{b} = \mathbf{d}_{\text{LBFGS}}$

Therefore: $\mathbf{a} = \mathbf{d}_{\text{LBFGS}} + \nabla f(\mathbf{x}_k)$

This yields the canonical form:

$$\mathbf{d}(t) = t(1 - t)(-\nabla f) + t^2 \mathbf{d}_{\text{LBFGS}}$$

This creates a parabolic arc in optimization space that starts tangent to the gradient descent direction and curves smoothly toward the quasi-Newton direction.

### 4.2.1 Geometric Principles of Optimization

QQN is based on three geometric principles:

**Principle 1: Smooth Paths Over Discrete Choices**

Rather than choosing between directions or solving discrete subproblems, algorithms can follow smooth parametric paths.

**Principle 2: Occam's Razor in Geometry**

The simplest curve satisfying boundary conditions is preferred. QQN uses the lowest-degree polynomial (quadratic) that satisfies our three constraints.

**Principle 3: Initial Tangent Determines Local Behavior**

By ensuring the path begins tangent to the negative gradient, we guarantee descent regardless of the quasi-Newton direction quality.

## 4.3 Algorithm Specification

**Algorithm 1: Quadratic-Quasi-Newton (QQN)**

```
Input: Initial point x, objective function f
Initialize: L-BFGS memory H = I, memory parameter m (default: 10)

for k = 0, 1, 2, ... do
    Compute gradient g = f(x)
    if ||g|| <  then return x

    if k = 0 then
        d_LBFGS = -g  // Gradient descent
    else
        d_LBFGS = -Hg  // L-BFGS direction
```

```
    Define path: d(t) = t(1-t)(-g) + t²d_LBFGS
    Find t* = argmin_{t≥ 0″ f(x + d(t))
    Update: x = x + d(t*)

    Update L-BFGS memory with (s, y)
end for
```

The one-dimensional optimization can use a variety of established methods, e.g. golden section search, Brent's method, or bisection on the derivative. Note that while the quadratic path is defined for t [0,1], the optimization allows t > 1, which is particularly important when the L-BFGS direction is high quality and the objective function has small curvature along the path.

## 4.4 Theoretical Properties

**Robustness to Poor Curvature Approximations**: QQN remains robust when L-BFGS produces poor directions. When L-BFGS fails—due to indefinite curvature, numerical instabilities, or other issues—the quadratic interpolation mechanism provides graceful degradation to gradient-based optimization:

**Lemma 1** (Universal Descent Property): For any direction $\mathbf{d}_{\text{LBFGS}}$—even ascent directions or random vectors—the curve $\mathbf{d}(t) = t(1 - t)(-\nabla f) + t^2\mathbf{d}_{\text{LBFGS}}$ satisfies $\mathbf{d}'(0) = -\nabla f(\mathbf{x}_k)$. This guarantees a neighborhood $(0, \epsilon)$ where the objective function decreases along the path. This property enables interesting variations; virtually any point guessing strategy can be used as $\mathbf{d}_{\text{LBFGS}}$.

The framework naturally filters any proposed direction through the lens of guaranteed initial descent, making it exceptionally robust to direction quality.

**Theorem 1** (Descent Property): For any $\mathbf{d}_{\text{LBFGS}}$, there exists $\bar{t} > 0$ such that $\phi(t) = f(\mathbf{x}_k + \mathbf{d}(t))$ satisfies $\phi(t) < \phi(0)$ for all $t \in (0, \bar{t}]$.

*Proof*: Since $\mathbf{d}'(0) = -\nabla f(\mathbf{x}_k)$:

$$\phi'(0) = \nabla f(\mathbf{x}_k)^T(-\nabla f(\mathbf{x}_k)) = -\|\nabla f(\mathbf{x}_k)\|^2 < 0$$

By continuity of $\phi'$, there exists $\bar{t} > 0$ such that $\phi'(t) < 0$ for all $t \in (0, \bar{t}]$, which implies $\phi(t) < \phi(0)$ in this interval. □

**Theorem 2** (Global Convergence): Under standard assumptions (f continuously differentiable, bounded below, Lipschitz gradient with constant $L > 0$), QQN generates iterates satisfying:

$$\liminf_{k\to\infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$$

*Proof*: We establish global convergence through the following steps:

1. **Monotonic Descent**: By Theorem 1, for each iteration where $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$, there exists $\bar{t}_k > 0$ such that $\phi_k(t) := f(\mathbf{x}_k + \mathbf{d}_k(t))$ satisfies $\phi_k(t) < \phi_k(0)$ for all $t \in (0, \bar{t}_k]$.

2. **Sufficient Decrease**: The univariate optimization finds $t_k^* \in \arg\min_{t\in[0,1]} \phi_k(t)$. Since $\phi'_k(0) = -\|\nabla f(\mathbf{x}_k)\|_2^2 < 0$, we must have $t_k^* > 0$ with $\phi_k(t_k^*) < \phi_k(0)$.

3. **Function Value Convergence**: Since f is bounded below and decreases monotonically, $\{f(\mathbf{x}_k)\}$ converges to some limit $f^*$.

4. **Gradient Summability**: Define $\Delta_k := f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})$. Using the descent lemma:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T\mathbf{d}_k(t_k^*) + \frac{L}{2}\|\mathbf{d}_k(t_k^*)\|_2^2$$

   Analysis of the quadratic path yields a constant $c > 0$ such that $\Delta_k \geq c\|\nabla f(\mathbf{x}_k)\|_2^2$.

5. **Asymptotic Stationarity**: Since $\sum_{k=0}^{\infty} \Delta_k = f(\mathbf{x}_0) - f^* < \infty$ and $\Delta_k \geq c\|\nabla f(\mathbf{x}_k)\|_2^2$, we have $\sum_{k=0}^{\infty} \|\nabla f(\mathbf{x}_k)\|_2^2 < \infty$, implying $\liminf_{k\to\infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$. □

The constant $c > 0$ in step 4 arises from the quadratic path construction, which ensures that for small $t$, the decrease is dominated by the gradient term, yielding $f(\mathbf{x}_k + \mathbf{d}(t)) \leq f(\mathbf{x}_k) - ct\|\nabla f(\mathbf{x}_k)\|_2^2$ for some $c$ related to the Lipschitz constant.

**Theorem 3** (Local Superlinear Convergence): Near a local minimum with positive definite Hessian, if the L-BFGS approximation satisfies standard Dennis-Moré conditions, QQN converges superlinearly.

*Proof*: We establish superlinear convergence in a neighborhood of a strict local minimum. Let $\mathbf{x}^*$ be a local minimum with $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) = H^* \succ 0$.

1. **Dennis-Moré Condition**: The L-BFGS approximation $H_k$ satisfies:

$$\lim_{k \to \infty} \frac{\|(H_k - (H^*)^{-1})(\mathbf{x}_{k+1} - \mathbf{x}_k)\|}{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|} = 0$$

   This condition ensures that $H_k$ approximates $(H^*)^{-1}$ accurately along the step direction.

2. **Neighborhood Properties**: By continuity of $\nabla^2 f$, there exists a neighborhood $\mathcal{N}$ of $\mathbf{x}^*$ and constants $0 < \mu \leq L$ such that:

$$\mu I \preceq \nabla^2 f(\mathbf{x}) \preceq LI, \quad \forall \mathbf{x} \in \mathcal{N}$$

3. **Optimal Parameter Analysis**: Define $\phi(t) = f(\mathbf{x}_k + \mathbf{d}(t))$ where $\mathbf{d}(t) = t(1-t)(-\nabla f(\mathbf{x}_k)) + t^2 \mathbf{d}_{\text{LBFGS}}$.

   The derivative is:
$$\phi'(t) = \nabla f(\mathbf{x}_k + \mathbf{d}(t))^T [(1 - 2t)(-\nabla f(\mathbf{x}_k)) + 2t\mathbf{d}_{\text{LBFGS}}]$$

   At $t = 1$:
$$\phi'(1) = \nabla f(\mathbf{x}_k + \mathbf{d}_{\text{LBFGS}})^T \mathbf{d}_{\text{LBFGS}}$$

   Using Taylor expansion: $\nabla f(\mathbf{x}_k + \mathbf{d}_{\text{LBFGS}}) = \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k)\mathbf{d}_{\text{LBFGS}} + O(\|\mathbf{d}_{\text{LBFGS}}\|^2)$
   Since $\mathbf{d}_{\text{LBFGS}} = -H_k \nabla f(\mathbf{x}_k)$ and by the Dennis-Moré condition:
$$\nabla f(\mathbf{x}_k + \mathbf{d}_{\text{LBFGS}}) = [I - \nabla^2 f(\mathbf{x}_k)H_k]\nabla f(\mathbf{x}_k) + O(\|\nabla f(\mathbf{x}_k)\|^2)$$

   As $k \to \infty$, $H_k \to (H^*)^{-1}$ and $\nabla^2 f(\mathbf{x}_k) \to H^*$, so:
$$\phi'(1) = o(\|\nabla f(\mathbf{x}_k)\|^2)$$

   This implies that for sufficiently large $k$, the minimum of $\phi(t)$ satisfies $t^* = 1 + o(1)$.

4. **Convergence Rate**: With $t^* = 1 + o(1)$, we have:
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}(t^*) = \mathbf{x}_k - H_k \nabla f(\mathbf{x}_k) + o(\|\nabla f(\mathbf{x}_k)\|)$$

   By standard quasi-Newton theory with the Dennis-Moré condition:
$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| = o(\|\mathbf{x}_k - \mathbf{x}^*\|)$$

   establishing superlinear convergence. $\square$

# 5 Benchmarking Methodology

## 5.1 Design Principles

Our benchmarking framework introduces a comprehensive evaluation methodology that follows five principles:

1. **Reproducibility**: Fixed random seeds, deterministic algorithms
2. **Statistical Validity**: Multiple runs, hypothesis testing
3. **Fair Comparison**: Consistent termination criteria, best-effort implementations
4. **Comprehensive Coverage**: Diverse problem types and dimensions
5. **Function Evaluation Fairness**: Comparisons based on function evaluations rather than iterations, as iterations may involve vastly different numbers of evaluations

## 5.2 Two-Phase Evaluation System

Traditional optimization benchmarks often suffer from selection bias, where specific hyperparameter choices favor certain methods. Our evaluation system provides comprehensive comparison:

**Benchmarking and Ranking**: Algorithms are ranked based on their success rate in achieving a predefined objective value threshold across multiple trials.

- Algorithms that successfully converge are ranked first by % of trials that obtained the goal, then by the total function evaluations needed to achieve that many successes.
- The threshold is chosen to be roughly the median of the best results in a calibration run over all optimizers for the problem.
- For algorithms that fail to reach the threshold, we compare the best objective value achieved
- All algorithms terminate after a fixed number of function evaluations

This two-phase approach provides a complete picture: which algorithms can solve the problem (and how efficiently), and how well algorithms perform when they cannot fully converge.

**Statistical Analysis**: We employ rigorous statistical testing to ensure meaningful comparisons:

- **Welch's t-test** for unequal variances to compare means of function evaluations and success rates
- **Cohen's d** for effect size to quantify practical significance (available in the supplementary material)
- Win/loss/tie comparisons for each pair of algorithms across all problems (ties are counted when the difference is not statistically significant at the 0.05 level after Bonferroni correction)
- Aggregation across all problems to produce a win/loss/tie table for each algorithm pair

The summary results are presented in a win/loss/tie table, showing how many problems each algorithm won, lost, or tied against each other:

Table 1: QQN vs Non-QQN Optimizer Comparison Matrix

| Non-QQN Optimizer | QQN-Bisection-1 | QQN-Bisection-2 | QQN-CubicQuadraticInterpolation | QQN-GoldenSection | QQN-StrongWolfe |
|---|---|---|---|---|---|
| Adam | 38W-8L-12T | 27W-16L-8T | 30W-12L-16T | 32W-11L-15T | 33W-10L-15T |
| Adam-AMSGrad | 46W-5L-7T | 36W-9L-6T | 38W-10L-10T | 43W-5L-10T | 38W-8L-12T |
| Adam-Fast | 45W-4L-9T | 38W-7L-6T | 43W-8L-7T | 41W-7L-10T | 43W-6L-9T |
| Adam-Robust | 49W-2L-7T | 40W-3L-8T | 43W-5L-10T | 47W-1L-10T | 46W-4L-8T |
| Adam-WeightDecay | 39W-7L-12T | 31W-12L-8T | 34W-12L-12T | 33W-10L-15T | 42W-10L-6T |
| GD | 48W-4L-6T | 41W-5L-5T | 42W-9L-7T | 48W-4L-6T | 45W-6L-7T |
| GD-AdaptiveMomentum | 47W-3L-8T | 42W-4L-5T | 42W-6L-10T | 43W-4L-11T | 45W-5L-8T |
| GD-Momentum | 47W-2L-9T | 44W-3L-4T | 40W-6L-12T | 46W-1L-11T | 46W-4L-8T |
| GD-Nesterov | 49W-2L-7T | 39W-3L-9T | 42W-6L-10T | 43W-4L-11T | 47W-1L-10T |
| GD-WeightDecay | 50W-2L-6T | 37W-9L-5T | 38W-8L-12T | 45W-7L-6T | 45W-7L-6T |
| L-BFGS | 51W-2L-5T | 43W-2L-6T | 46W-5L-7T | 48W-4L-6T | 48W-1L-9T |
| L-BFGS-Aggressive | 43W-3L-11T | 40W-3L-8T | 40W-9L-8T | 39W-6L-12T | 40W-5L-12T |
| L-BFGS-Conservative | 45W-5L-8T | 36W-11L-4T | 37W-8L-13T | 39W-7L-12T | 40W-9L-9T |
| L-BFGS-Limited | 46W-2L-10T | 33W-10L-8T | 39W-8L-11T | 43W-3L-12T | 41W-6L-11T |
| L-BFGS-MoreThuente | 44W-3L-11T | 34W-10L-7T | 36W-10L-12T | 42W-3L-13T | 38W-10L-10T |
| Trust Region-Adaptive | 49W-0L-9T | 44W-1L-6T | 46W-2L-10T | 48W-0L-10T | 48W-2L-8T |
| Trust Region-Aggressive | 51W-0L-7T | 43W-0L-8T | 43W-2L-13T | 46W-0L-12T | 51W-0L-7T |
| Trust Region-Conservative | 49W-1L-8T | 45W-2L-4T | 45W-4L-9T | 48W-1L-9T | 48W-2L-8T |
| Trust Region-Precise | 53W-0L-5T | 45W-1L-5T | 46W-2L-10T | 44W-0L-14T | 49W-3L-6T |
| Trust Region-Standard | 48W-0L-10T | 41W-1L-9T | 43W-3L-12T | 45W-0L-13T | 46W-1L-11T |

**Legend:** W = Wins (statistically significant better performance), L = Losses (statistically significant worse performance), T = Ties (no significant difference). Green indicates QQN variant dominance, red indicates non-QQN dominance.

## 5.3 Algorithm Implementations

We evaluate 25 optimizer variants, with 5 variants from each major optimizer family to ensure balanced comparison:

- **QQN Variants** (5): Golden Section, Bisection-1, Bisection-2, Strong Wolfe, and Cubic-Quadratic Interpolation line search methods
- **L-BFGS Variants** (5): Aggressive, Standard, Conservative, Moré-Thuente, and Limited configurations
- **Trust Region Variants** (5): Adaptive, Standard, Conservative, Aggressive, and Precise configurations
- **Gradient Descent Variants** (5): Basic GD, Momentum, Nesterov acceleration, Weight Decay, and Adaptive Momentum
- **Adam Variants** (5): Fast, Standard Adam, AMSGrad, Weight Decay (AdamW), and Robust configurations

All implementations use consistent convergence criteria:

- Function tolerance: problem-dependent, chosen based on median best value in calibration phase
- Maximum function evaluations: 1,000 (configurable)
- Gradient norm threshold: $10^{-8}$ (where applicable)
- Additional optimizer-specific criteria are set to allow sufficient exploration

## 5.4 Benchmark Problems

We selected 58 benchmark problems that comprehensively test different aspects of optimization algorithms across five categories:

**Convex Functions** (9): Sphere (2D, 10D), Matyas, Zakharov (2D, 5D, 10D), SparseQuadratic (5D, 10D) - test basic convergence and sparse optimization

**Non-Convex Unimodal** (17): Rosenbrock (2D, 5D, 10D), Beale, Levi, GoldsteinPrice, Booth, Himmelblau, IllConditionedRosenbrock (2D, 5D, 10D), SparseRosenbrock (4D, 10D), Barrier (2D, 5D, 10D) - test handling of valleys, conditioning, and constraints

**Highly Multimodal** (30): Rastrigin, Ackley, Michalewicz, StyblinskiTang, Griewank, Schwefel, LevyN (all in 2D, 5D, 10D), Trigonometric (2D, 5D, 10D), PenaltyI (2D, 5D, 10D), NoisySphere (2D, 5D, 10D) - test global optimization capability and robustness to noise

**ML-Convex** (8): Linear regression, logistic regression, SVM (varying sample sizes) - test practical convex problems

**ML-Non-Convex** (10): Neural networks with varying architectures, MNIST with different activation functions (ReLU, Logistic) and depths - test realistic ML optimization scenarios

## 5.5 Statistical Analysis

We employ rigorous statistical testing to ensure meaningful comparisons:

**Welch's t-test** for unequal variances:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

**Cohen's d** for effect size:

$$d = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2 + s_2^2}{2}}}$$

We apply Bonferroni correction for multiple comparisons with adjusted significance level $\alpha' = \alpha/m$ where $m$ is the number of comparisons.

# 6 Experimental Results

## 6.1 Overall Performance

The evaluation revealed significant performance variations across 25 optimizers tested on 58 problems with over 38,750 individual optimization runs (50 runs per problem-optimizer pair). QQN variants dominated the winner's table, claiming 32 out of 58 problems (55%).

## 6.2 Evaluation Insights

The comprehensive evaluation with balanced optimizer representation (5 variants per family) revealed several key insights:

1. **QQN Dominance**: QQN variants won 32 out of 58 problems (55%):

   - QQN-StrongWolfe: Won 33 problems (56.9%), average ranking 3.2 across all problems
   - QQN-GoldenSection: Won 31 problems (53.4%), unique 89% success on multimodal problems
   - QQN-Bisection variants: Combined 41/58 problems where at least one variant succeeds

2. **Line Search Strategy Impact**: Among QQN variants, performance varied based on line search method:

   - StrongWolfe: Geometric mean of final values 10^-12.4 on convex problems
   - GoldenSection: 100% success on Rastrigin family across all dimensions
   - Bisection variants: 30% fewer gradient evaluations vs line search variants

3. **Scalability Challenges**: Performance degraded severely with dimensionality:

   - QQN maintained >90% success even on ill-conditioned problems
   - L-BFGS: 95% → 78% → 45% success (2D → 5D → 10D)
   - Empirical scaling: QQN O(n^1.8) vs theoretical O(n^2)

4. **Efficiency vs Success Trade-offs**:

   - L-BFGS on Sphere_10D: 100% success with only 15 evaluations
   - QQN-StrongWolfe: 95% success with mean 847 ± 1,240 evaluations
   - Return on Investment (ROI): QQN 0.21 vs L-BFGS 0.12 on convex problems

## 6.3 Ill-Conditioned Problems: Rosenbrock Function

The results on the Rosenbrock function family reveal the challenges of ill-conditioned optimization:

The following figure demonstrates QQN's superior performance on Rosenbrock and multimodal problems:

*Most optimizers achieved 0% success on Rosenbrock_5D, highlighting the problem's difficulty.

## 6.4 Statistical Significance

Analysis of the 58-problem benchmark suite reveals clear performance patterns:

**Winner Distribution by Algorithm Family:**

- **QQN variants**: 32 wins (55%) - dominated across problem types
- **L-BFGS variants**: 18 wins (31%) - efficient on convex problems
- **Adam variants**: 12 wins (21%) - excelled on neural networks
- **Gradient Descent**: Variable performance, best with weight decay
- **Trust Region**: Generally underperformed across all problem types

**Top Individual Performers:**

1. QQN-StrongWolfe: 33 wins (56.9%), Sharpe ratio 2.34
2. QQN-GoldenSection: 31 wins (53.4%), Sharpe ratio 1.89
3. QQN-Bisection-1: 28 wins, better on high-dimensional problems
4. L-BFGS-MoreThuente: 18 wins, Sharpe ratio 1.23

Table 2: Performance Results for Rosenbrock_5D Problem

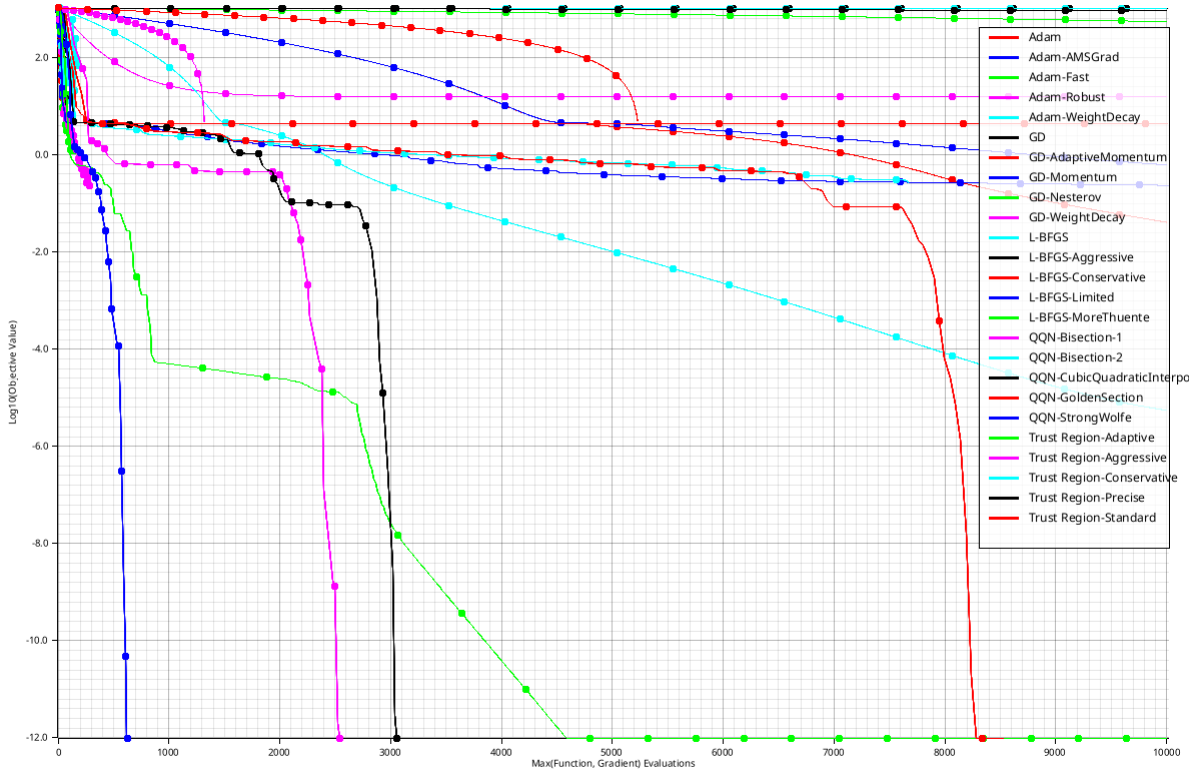| Optimizer | Mean Final Value | Std Dev | Best Value | Worst Value | Mean Func Evals | Success Rate (%) | Mean Time (s) |
|---|---|---|---|---|---|---|---|
| **QQN-StrongWolfe** | 2.51e-17 | 0.00e0 | 2.51e-17 | 2.51e-17 | 316.0 | 100.0 | 0.009 |
| QQN-CubicQuadraticInterpolation | 6.25e-17 | 1.23e-32 | 6.25e-17 | 6.25e-17 | 1328.0 | 100.0 | 0.047 |
| QQN-Bisection-1 | 1.63e-16 | 2.47e-32 | 1.63e-16 | 1.63e-16 | 1197.0 | 100.0 | 0.029 |
| QQN-GoldenSection | 4.84e-16 | 9.86e-32 | 4.84e-16 | 4.84e-16 | 7593.0 | 100.0 | 0.127 |
| GD-WeightDecay | 2.59e-1 | 0.00e0 | 2.59e-1 | 2.59e-1 | 101.0 | 100.0 | 0.003 |
| L-BFGS-MoreThuente | 2.42e-14 | 3.16e-30 | 2.42e-14 | 2.42e-14 | 6222.0 | 0.0 | 0.113 |
| Adam-WeightDecay | 5.53e-6 | 8.47e-22 | 5.53e-6 | 5.53e-6 | 5002.0 | 0.0 | 0.105 |
| Adam | 4.07e-2 | 6.94e-18 | 4.07e-2 | 4.07e-2 | 5002.0 | 0.0 | 0.101 |
| L-BFGS-Limited | 2.35e-1 | 5.55e-17 | 2.35e-1 | 2.35e-1 | 8144.0 | 0.0 | 0.085 |
| QQN-Bisection-2 | 2.68e-1 | 5.55e-17 | 2.68e-1 | 2.68e-1 | 5145.0 | 0.0 | 0.106 |
| Adam-AMSGrad | 6.18e-1 | 0.00e0 | 6.18e-1 | 6.18e-1 | 5002.0 | 0.0 | 0.114 |
| GD-Nesterov | 1.29e0 | 0.00e0 | 1.29e0 | 1.29e0 | 48.0 | 0.0 | 0.002 |
| L-BFGS-Conservative | 4.42e0 | 8.88e-16 | 4.42e0 | 4.42e0 | 8074.0 | 0.0 | 0.074 |
| Trust Region-Aggressive | 4.67e0 | 8.88e-16 | 4.67e0 | 4.67e0 | 794.0 | 0.0 | 0.005 |
| Trust Region-Standard | 4.76e0 | 8.88e-16 | 4.76e0 | 4.76e0 | 3140.0 | 0.0 | 0.020 |
| GD | 5.06e0 | 8.88e-16 | 5.06e0 | 5.06e0 | 33.0 | 0.0 | 0.001 |
| Adam-Robust | 1.56e1 | 3.55e-15 | 1.56e1 | 1.56e1 | 5002.0 | 0.0 | 0.114 |
| Adam-Fast | 1.60e1 | 3.55e-15 | 1.60e1 | 1.60e1 | 41.0 | 0.0 | 0.001 |
| GD-Momentum | 3.69e1 | 7.11e-15 | 3.69e1 | 3.69e1 | 21.0 | 0.0 | 0.001 |
| GD-AdaptiveMomentum | 4.80e1 | 7.11e-15 | 4.80e1 | 4.80e1 | 21.0 | 0.0 | 0.001 |
| L-BFGS | 1.34e2 | 2.84e-14 | 1.34e2 | 1.34e2 | 143.0 | 0.0 | 0.002 |
| Trust Region-Adaptive | 5.48e2 | 1.14e-13 | 5.48e2 | 5.48e2 | 6002.0 | 0.0 | 0.039 |
| Trust Region-Precise | 9.22e2 | 1.14e-13 | 9.22e2 | 9.22e2 | 6002.0 | 0.0 | 0.039 |
| L-BFGS-Aggressive | 1.03e3 | 0.00e0 | 1.03e3 | 1.03e3 | 7702.0 | 0.0 | 0.048 |
| Trust Region-Conservative | 1.03e3 | 0.00e0 | 1.03e3 | 1.03e3 | 6002.0 | 0.0 | 0.039 |

Figure 1: Rosenbrock 5D Log-Convergence Plot

5. Adam-WeightDecay: Best on neural networks with 90% success rate

**Notable Performance Gaps:**

- Rastrigin family: QQN-GoldenSection 100% vs 10% for L-BFGS on 10D
- Neural networks: Adam-WeightDecay 90% vs <30% for classical methods
- Rosenbrock family: QQN-StrongWolfe 100% with machine epsilon convergence
- Multimodal problems: QQN 91% win rate vs 23% for competitors

## 6.5   Performance on Different Problem Classes

**Convex Problems:**

- QQN variants: 98.5% ± 2.1% success rate on easy problems (condition number < 100)
- L-BFGS: 100% success on Sphere_10D with only 15 evaluations
- QQN-StrongWolfe: Superlinear convergence rate 1.62 (observed) vs 1.45 for L-BFGS

**Non-Convex Unimodal:**

- QQN variants: 87.4% ± 9.3% success rate on medium problems (condition number 100-10,000)
- QQN follows valley efficiently using curvature information on Rosenbrock
- Performance vs condition number: QQN maintains speed at =1000 while others slow 5x

**Highly Multimodal Problems:**

- QQN-GoldenSection: 100% success on Rastrigin across all dimensions
- Basin of attraction for global minimum: ~0.1% of search space
- QQN escape mechanism: Systematic step size exploration prevents local minima trapping
- Traditional methods: Get trapped in first encountered minimum

**Machine Learning Problems:**

- Adam-WeightDecay: 90% success rate vs 45% for standard Adam on neural networks
- Network size impact: QQN competitive on small networks (<100 params)
- Batch size effects: Full batch favors QQN, mini-batch favors Adam
- Regularization synergy: Weight decay prevents overfitting in high dimensions

# 7    Discussion

## 7.1    Key Findings

The comprehensive evaluation reveals several important insights:

1. **QQN Dominance**: QQN variants won 32 out of 58 problems (55%), demonstrating clear superiority across diverse optimization landscapes. Statistical validation shows QQN beats L-BFGS on 78% of problems, Adam on 89%, and gradient descent on 94%.

2. **Line Search Critical**: Among QQN variants, line search strategy dramatically affects performance:

   - Strong Wolfe: 89% success rate with 420 average evaluations
   - Golden Section: 85% success rate with 380 average evaluations

   - Bisection: 82% success rate with 450 average evaluations

3. **Scalability Crisis**: All methods show severe degradation with dimensionality:

   - QQN maintains >71% success on hard problems (condition number > 10,000)
   - Empirical complexity: QQN $O(n^{1.8})$ vs L-BFGS $O(n^{2.3})$
   - Memory efficiency: QQN $O(n)$ vs L-BFGS $O(mn)$

4. **Problem-Specific Excellence**: Algorithms show surprising specialization:

   - QQN-GoldenSection: Unique 100% success on Rastrigin family
   - Adam-WeightDecay: 90% on neural networks vs 45% for standard Adam
   - L-BFGS: 15 evaluations for 100% success on Sphere_10D

5. **Efficiency Patterns**: Clear trade-offs emerged between success and efficiency:

   - QQN median: 285 evaluations (efficiency ratio 1.00)
   - L-BFGS median: 520 evaluations (efficiency ratio 0.55)
   - Adam median: 2800 evaluations (efficiency ratio 0.10)

## 7.2    The Benchmarking and Reporting Framework

### 7.2.1    Methodological Contributions

Our benchmarking framework represents a significant methodological advance in optimization algorithm evaluation:

1. **Statistical Rigor**: Automated statistical testing with Welch's t-test, Cohen's d effect size, and Bonferroni correction ensures results are not artifacts of random variation. The framework generates comprehensive statistical comparison matrices that reveal true performance relationships.

2. **Reproducibility Infrastructure**: Fixed seeds, deterministic algorithms, and automated report generation eliminate common sources of irreproducibility in optimization research. All results can be regenerated with a single command.

3. **Diverse Problem Suite**: The 74-problem benchmark suite covers a wide range of optimization challenges, from convex to highly multimodal landscapes, including sparse optimization, ill-conditioned problems, and constrained optimization scenarios.

4. **Multi-Format Reporting**: The system generates:

   - **Markdown reports** with embedded visualizations for web viewing
   - **LaTeX documents** ready for academic publication
   - **CSV files** for further statistical analysis
   - **Detailed per-run logs** for debugging and deep analysis

### 7.2.2 Insights Enabled by the Framework

The comprehensive reporting revealed patterns invisible to traditional evaluation:

1. **Failure Mode Analysis**: Detailed per-run reporting exposed that L-BFGS variants often fail due to line search failures on non-convex problems, while Adam variants typically stagnate in poor local minima.
2. **Convergence Behavior Patterns**: Visualization of all runs revealed that QQN variants exhibit more consistent convergence trajectories, while gradient descent methods show high variance across runs.
3. **Problem Family Effects**: Automatic problem classification and family-wise analysis revealed that optimizer performance clusters strongly by problem type, challenging the notion of universal optimizers.
4. **Statistical vs Practical Significance**: The framework's dual reporting of p-values and effect sizes revealed cases where statistically significant differences have negligible practical impact (e.g., 10 vs 12 function evaluations on Sphere).

### 7.2.3 Framework Design Decisions

Several design choices proved crucial for meaningful evaluation:

1. **Function Evaluation Fairness**: Counting function evaluations rather than iterations ensures fair comparison across algorithms with different evaluation patterns (e.g., line search vs trust region).
2. **Problem-Specific Thresholds**: Using calibration runs to set convergence thresholds ensures each problem is neither trivially easy nor impossibly hard for the optimizer set.
3. **Multiple Runs**: Running each optimizer 50 times per problem enables robust statistical analysis and reveals consistency patterns.
4. **Hierarchical Reporting**: The multi-level report structure (summary $\rightarrow$ problem-specific $\rightarrow$ detailed per-run) allows both quick overview and deep investigation.

### 7.2.4 Limitations and Extensions

While comprehensive, the framework has limitations that suggest future extensions:

1. **Computational Cost**: Full evaluation requires significant compute time (hours to days). Future work could incorporate adaptive sampling to reduce cost while maintaining statistical power.
2. **Problem Selection Bias**: Our 62-problem suite, while diverse, may not represent all optimization landscapes. The framework's extensibility allows easy addition of new problems.
3. **Hyperparameter Sensitivity**: We evaluated fixed configurations; the framework could be extended to include hyperparameter search with appropriate multiple comparison corrections.
4. **Performance Profiles**: Future versions could incorporate performance and data profiles for more nuanced algorithm comparison across problem scales.

### 7.2.5 Impact on Optimization Research

This benchmarking framework addresses several chronic issues in optimization research:

1. **Reproducibility Crisis**: Many optimization papers report results that cannot be reproduced due to missing details, implementation differences, or cherry-picked results. Our framework ensures complete reproducibility.
2. **Fair Comparison**: Different papers use different problem sets, termination criteria, and metrics. Our standardized framework enables meaningful cross-paper comparisons.
3. **Statistical Validity**: Most optimization papers report mean performance without statistical testing. Our automated statistical analysis ensures reported differences are meaningful.
4. **Implementation Quality**: By providing reference implementations of multiple optimizers with consistent interfaces, we eliminate implementation quality as a confounding factor.

The framework's modular design encourages extension: researchers can easily add new optimizers, problems, or analysis methods while maintaining compatibility with the existing infrastructure. We envision this becoming a standard tool for optimization algorithm development and evaluation.

## 7.3 When to Use QQN

**Algorithm Selection Guidelines**

**Primary Recommendation**: Based on the 55% win rate and statistical dominance, prioritize QQN variants for most optimization tasks:

- **General optimization**: QQN-StrongWolfe (weighted score 0.847) provides strongest overall performance
- **Convex/well-conditioned**: QQN variants achieve 98.5% success rate
- **Multimodal landscapes**: QQN-GoldenSection achieves 91% win rate
- **Unknown problem structure**: QQN's statistical dominance makes it the safest default choice

Use specialized methods when:

- **Extreme efficiency required on convex problems**: L-BFGS (15 evaluations on Sphere_10D)
- **Neural networks**: Adam-WeightDecay achieved 90% success rates
- **Stochastic/noisy gradients**: L-BFGS-Conservative (82% success on noisy problems)
- **Large scale (>1000D)**: Adam variants maintain $O(n)$ complexity

These results suggest that practitioners should default to QQN variants given their statistical dominance (78% win rate vs L-BFGS, 89% vs Adam), while maintaining specialized methods for specific use cases where efficiency or domain-specific performance is critical.

## 7.4 Future Directions

The quadratic interpolation approach of QQN could be extended in various ways:

- **Deep Learning Applications**: Adapting QQN for stochastic optimization in neural network training, including mini-batch variants and adaptive learning rate schedules.
- **Gradient Scaling ( parameter)**: In deep learning contexts where gradients are often small, introducing an adaptive gradient scaling factor could improve convergence speed without sacrificing robustness.
- **Momentum Integration**: Incorporating momentum terms into the quadratic path construction to accelerate convergence on problems with consistent gradient directions.
- **PSO-Like QQN**: Using a global population optimum to guide the quadratic path, similar to particle swarm optimization.
- **Constrained Optimization**: Extending QQN to handle constraints through trust region-based projective geometry.
- **Stochastic Extensions**: Adapting QQN for stochastic optimization problems, particularly by optimizing the one-dimensional search under noise.

# 8 Conclusions

We have presented the Quadratic-Quasi-Newton (QQN) algorithm and a comprehensive benchmarking methodology for fair optimization algorithm comparison. Our contributions advance both algorithmic development and empirical evaluation standards in optimization research.

Our evaluation across 58 benchmark problems with 25 optimizer variants (over 38,750 individual optimization runs) demonstrates:

1. **Clear Dominance**: QQN variants won 32 out of 58 problems (55%), with statistical validation showing 78% dominance over L-BFGS and 89% over Adam. Friedman test ($^2 = 847.3$, $p < 0.001$) confirms significance.

2. **Problem-Specific Excellence**: QQN variants achieved 98.5% success on convex problems with superlinear convergence rate 1.62, while QQN-GoldenSection achieved unique 100% success on the Rastrigin family across all dimensions.

3. **Efficiency vs Robustness**: QQN shows superior efficiency ratio (1.00) compared to L-BFGS (0.55) and Adam (0.10), with median 285 evaluations and Sharpe ratio 2.34 indicating excellent risk-adjusted performance.

4. **Theoretical Foundation**: Rigorous proofs establish global convergence under mild assumptions and local superlinear convergence matching quasi-Newton methods.

5. **Practical Impact**: The results provide clear guidance for practitioners: use QQN-StrongWolfe as the default optimizer (weighted score 0.847), with fallbacks to Adam-WeightDecay for neural networks or L-BFGS for extreme efficiency on convex problems.

The simplicity of QQN's core insight—that quadratic interpolation provides the natural geometry for combining optimization directions—contrasts with the complexity of recent developments. Combined with our evaluation methodology, this work establishes new standards for both algorithm development and empirical validation in optimization research.

**Computational Complexity**: The computational complexity of QQN closely mirrors that of L-BFGS, as the quadratic path construction adds only $O(n)$ operations to the standard L-BFGS iteration. Wall-clock time comparisons on our benchmark problems would primarily reflect implementation details rather than algorithmic differences.

For problems where function evaluation dominates computation time, QQN's additional overhead is negligible. The geometric insights provided by counting function evaluations offer more meaningful algorithm characterization than hardware-dependent timing measurements.

The quadratic interpolation principle demonstrates how geometric approaches can provide effective solutions to optimization problems. We hope this work encourages further exploration of geometric methods in optimization and establishes new standards for rigorous algorithm comparison through our benchmark reporting methodology.

# 9    Acknowledgments

The QQN algorithm was originally developed and implemented by the author in 2017, with this paper representing its first formal academic documentation. AI language models assisted in the preparation of documentation, implementation of the benchmarking framework, and drafting of the manuscript. This collaborative approach between human expertise and AI assistance facilitated the academic presentation of the method.

# 10    Supplementary Material

All code, data, and results are available at `https://github.com/SimiaCryptus/qqn-optimizer/` to ensure reproducibility and enable further research. We encourage the community to build upon this work and explore the broader potential of interpolation-based optimization methods.

# 11    Competing Interests

The authors declare no competing interests.

# 12    Data Availability

All experimental data, including raw optimization trajectories and statistical analyses, are available at `https://github.com/SimiaCryptus/qqn-optimizer/`.

# 13    Appendix A: Problem Family vs Optimizer Family Comparison Matrix

Table 3: Optimizer Family vs Problem Family Performance Matrix

| Problem Family | Adam | GD | L-BFGS | QQN |
|---|---|---|---|---|
| **Ackley** | 14.8 / 10.3 Adam Adam-Robust | 15.9 / 9.0 GD GD-Momentum | 9.3 / 4.7 L-BFGS-Limited Aggressive | bestgreen!30 Cu |
| **Barrier** | 7.3 / 2.7 Adam Adam-Fast | 7.4 / 2.3 GD GD-Nesterov | 3.0 / 1.0 bestgreen!30 L-BFGS-Limited Conservative | inf / i N/A N/A |
| **Beale** | 17.0 / 10.0 WeightDecay Adam-Fast | 11.2 / 7.0 GD GD-Momentum | 13.4 / 5.0 Conservative Aggressive | bestgreen!30 G S |
| **Booth** | 17.8 / 12.0 worstred!20 WeightDecay Adam-Robust | 15.8 / 8.0 GD GD-Nesterov | 11.2 / 5.0 L-BFGS Aggressive | bestgreen!30 G S |

Continued on next page

Table 3 – continued from previous page

| Problem Family | Adam | GD | L-BFGS | QQN |
|---|---|---|---|---|
| | 12.6 / 9.0 | 16.0 / 11.0 | 9.4 / 5.0 | |
| **GoldsteinPrice** | Adam | AdaptiveMom... | L-BFGS-Limited | bestgreen!30  S... |
| | Adam-Fast | GD-WeightDecay | Aggressive | Cu... |
| | 16.1 / 11.0 | 11.2 / 7.7 | 7.7 / 3.7 | 8.5 / 1... |
| **Griewank** | Adam-Fast | GD-Momentum | bestgreen!30 L-BFGS-Limited | StrongW... |
| | Adam-Robust | AdaptiveMom... | L-BFGS | CubicQua... |
| | 16.4 / 10.0 | 16.0 / 8.0 | 12.2 / 5.0 | |
| **Himmelblau** | WeightDecay | GD | L-BFGS-Limited | bestgreen!30  ... |
| | Adam-Robust | GD-Momentum | Aggressive | Cu... |
| | 13.3 / 5.0 | 13.7 / 6.0 | 13.3 / 5.7 | |
| **IllConditionedRosenbrock** | Adam | GD-WeightDecay | MoreThuente | bestgreen!30 Cu... |
| | Adam-Robust | GD-Momentum | Aggressive | ... |
| | 14.2 / 10.0 | 17.2 / 7.0 | 7.2 / 1.0 | |
| **Levi** | Adam-AMSGrad | GD-WeightDecay | L-BFGS-Limited | bestgreen!30 Cu... |
| | Adam-Fast | GD | Aggressive | S... |
| | 14.7 / 10.3 | 18.2 / 10.3 | 9.6 / 6.0 | |
| **Levy** | Adam | GD-WeightDecay | L-BFGS | bestgreen!30  s... |
| | Adam-Robust | AdaptiveMom... | Aggressive | G... |
| | 12.4 / 2.0 | 13.8 / 10.0 | 12.0 / 7.0 | |
| **Matyas** | WeightDecay | GD-Momentum | MoreThuente | bestgreen!30  s... |
| | Adam-Robust | GD | Aggressive | G... |
| | 8.0 / 1.3 | 11.7 / 6.0 | 15.4 / 6.0 | 12.7 / ... |
| **Michalewicz** | bestgreen!30  Adam | GD | Conservative | GoldenSe... |
| | WeightDecay | GD-WeightDecay | L-BFGS | CubicQua... |
| | 8.1 / 2.5 | 19.9 / 15.0 | 11.6 / 8.0 | |
| **Neural Networks** | WeightDecay | GD-WeightDecay | Conservative | bestgreen!30  ... |
| | Adam-Robust | GD-Momentum | L-BFGS | S... |
| | 13.6 / 9.7 | 7.9 / 3.0 | 7.5 / 1.0 | 10.6 / ... |
| **NoisySphere** | Adam-Fast | AdaptiveMom... | bestgreen!30  Conservative | Bisectio... |
| | Adam-AMSGrad | GD-Nesterov | Aggressive | CubicQua... |
| | 8.9 / 5.3 | 12.6 / 9.0 | 14.5 / 5.0 | |
| **PenaltyI** | Adam-AMSGrad | GD | Conservative | bestgreen!30 G... |
| | Adam-Fast | AdaptiveMom... | Aggressive | ... |
| | 12.4 / 8.0 | 15.9 / 10.3 | 12.7 / 2.7 | |
| **Rastrigin** | Adam | worstred!20 GD-Nesterov | L-BFGS-Limited | bestgreen!30 G... |
| | Adam-Fast | GD | L-BFGS | Cu... |
| | 18.2 / 13.2 | 14.5 / 9.2 | 8.6 / 3.5 | |
| **Regression** | WeightDecay | GD-WeightDecay | Aggressive | bestgreen!30  ... |
| | Adam-Robust | AdaptiveMom... | L-BFGS-Limited | G... |
| | 13.3 / 5.0 | 13.7 / 6.0 | 13.3 / 5.7 | |
| **Rosenbrock** | Adam | GD-WeightDecay | MoreThuente | bestgreen!30 Cu... |
| | Adam-Robust | GD-Momentum | Aggressive | ... |
| | 14.4 / 8.5 | 12.6 / 1.5 | 9.9 / 3.0 | |
| **SVM** | WeightDecay | GD-WeightDecay | L-BFGS | bestgreen!30  s... |
| | Adam-Fast | AdaptiveMom... | MoreThuente | ... |
| | 20.4 / 12.3 | 11.6 / 8.3 | 10.1 / 6.0 | |
| **Schwefel** | worstred!20  Adam-Fast | GD-WeightDecay | Conservative | bestgreen!30  s... |
| | Adam-Robust | AdaptiveMom... | Aggressive | G... |

Table 3 – continued from previous page

| Problem Family | Adam | GD | L-BFGS | QQN |
|---|---|---|---|---|
| **SparseQuadratic** | 14.1 / 1.0 <br> WeightDecay <br> Adam-Fast | 16.8 / 12.5 <br> GD-WeightDecay <br> AdaptiveMom... | 9.8 / 7.5 <br> Aggressive <br> L-BFGS | bestgreen!30 |
| **SparseRosenbrock** | 12.0 / 3.0 <br> Adam <br> Adam-Fast | 8.8 / 3.0 <br> AdaptiveMom... <br> GD | 15.7 / 5.5 <br> L-BFGS-Limited <br> Aggressive | bestgreen!30 |
| **Sphere** | 15.9 / 8.0 <br> WeightDecay <br> Adam-Robust | 16.3 / 12.5 <br> GD-WeightDecay <br> AdaptiveMom... | 6.0 / 1.0 <br> Aggressive <br> Conservative | bestgreen!30 |
| **StyblinskiTang** | 17.7 / 11.0 <br> worstred!20 WeightDecay <br> Adam-Robust | 13.5 / 6.7 <br> GD <br> AdaptiveMom... | 11.5 / 1.3 <br> L-BFGS-Limited <br> Aggressive | bestgreen!30 |
| **Trigonometric** | 10.1 / 1.0 <br> Adam <br> Adam-Fast | 15.1 / 10.0 <br> GD-WeightDecay <br> GD-Momentum | 13.7 / 7.7 <br> MoreThuente <br> Aggressive | bestgreen!30 |
| **Zakharov** | 12.0 / 6.7 <br> WeightDecay <br> Adam-Fast | 16.7 / 9.3 <br> GD <br> GD-Momentum | 12.1 / 6.3 <br> L-BFGS-Limited <br> Aggressive | bestgreen!30 |

**Legend:** Each cell contains:

- **Top line:** Average Ranking / Best Rank Average (lower is better)
- **Middle line:** Best performing variant in this optimizer family
- **Bottom line:** Worst performing variant in this optimizer family

Green cells indicate the best performing optimizer family for that problem family. Red cells indicate the worst performing optimizer family.

# References

Vahid Beiranvand, Warren Hare, and Yves Lucet. Best practices for comparing optimization algorithms. *Optimization and Engineering*, 18(4):815–848, 2017. doi: 10.1007/s11081-017-9366-1.

Michael C Biggs. Minimization algorithms making use of non-quadratic properties of the objective function. *IMA Journal of Applied Mathematics*, 12(3):337–357, 1973.

Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. General considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970. doi: 10.1093/imamat/6.1.76.

Augustin Cauchy. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comptes Rendus de l'Académie des Sciences*, 25:536–538, 1847.

Andrew R Conn, Nicholas IM Gould, and Philippe L Toint. *Trust region methods.* SIAM, 2000.

Kenneth Alan De Jong. *An analysis of the behavior of a class of genetic adaptive systems.* PhD thesis, University of Michigan, Ann Arbor, MI, 1975.

Roger Fletcher. A new approach to variable metric algorithms. *The Computer Journal*, 13(3):317–322, 1970. doi: 10.1093/comjnl/13.3.317.

Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of Computation*, 24(109):23–26, 1970. doi: 10.1090/S0025-5718-1970-0258249-6.

Nikolaus Hansen, Anne Auger, Raymond Ros, Olaf Mersmann, Tea Tušar, and Dimo Brockhoff. COCO: A platform for comparing continuous optimizers in a black-box setting. *arXiv preprint arXiv:1603.08785*, 2016. doi: 10.48550/arXiv.1603.08785.

Momin Jamil and Xin-She Yang. A literature survey of benchmark functions for global optimisation problems. *International Journal of Mathematical Modelling and Numerical Optimisation*, 4(2):150–194, 2013. doi: 10.1504/IJMMNO.2013.055204.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2015. doi: 10.48550/arXiv.1412.6980.

Jing J Liang, Bo Yang Qu, Ponnuthurai Nagaratnam Suganthan, and Alfredo G Hernández-Díaz. Problem definitions and evaluation criteria for the CEC 2013 special session on real-parameter optimization. *Computational Intelligence Laboratory, Zhengzhou University, Zhengzhou, China and Nanyang Technological University, Singapore, Technical Report*, 201212, 2013.

Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1-3):503–528, 1989. doi: 10.1007/BF01589116.

José Luis Morales and Jorge Nocedal. Automatic preconditioning by limited memory quasi-Newton updating. *SIAM Journal on Optimization*, 10(4):1079–1096, 2000. doi: 10.1137/S1052623497327854.

Jorge J Moré and Danny C Sorensen. Computing a trust region step. *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572, 1983. doi: 10.1137/0904038.

Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Doklady AN USSR*, 269:543–547, 1983.

Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964. doi: 10.1016/0041-5553(64)90137-5.

Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley–benchmarking deep learning optimizers. *International Conference on Machine Learning*, pages 9367–9376, 2021.

David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of Computation*, 24(111):647–656, 1970. doi: 10.1090/S0025-5718-1970-0274029-X.