

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

Andrew Charneski
SimiaCryptus Software

August 2, 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 a diverse set of benchmark problems with multiple optimizer variants from each major family, totaling thousands of individual optimization runs with multiple runs per problem-optimizer pair. Our results demonstrate that QQN variants achieve significant dominance across the benchmark suite. QQN algorithms won most problems, with QQN-StrongWolfe showing particularly strong performance and achieving excellent success rates on convex problems while requiring relatively few function evaluations. Statistical analysis using Friedman test confirms QQN’s superiority with effect sizes showing practical significance. QQN-StrongWolfe achieved high precision convergence on Rosenbrock family and excellent success on Sphere problems across all dimensions. While L-BFGS variants showed high efficiency on convex problems and Adam-WeightDecay dominated neural network tasks, QQN’s consistent performance across problem types establishes its practical utility with superior weighted performance scores compared to L-BFGS and 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 $\mathbf{d} : [0, 1] \rightarrow \mathbb{R}^n$ 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 \mathbf{x} , objective function f

Initialize: L-BFGS memory $\mathbf{H} = \mathbf{I}$, memory parameter m (default: 10)

```

for  $k = 0, 1, 2, \dots$  do
  Compute gradient  $\mathbf{g} = \nabla f(\mathbf{x})$ 
  if  $\|\mathbf{g}\| < \epsilon$  then return  $\mathbf{x}$ 

  if  $k = 0$  then
     $\mathbf{d}_{\text{LBFGS}} = -\mathbf{g}$  // Gradient descent
  else
     $\mathbf{d}_{\text{LBFGS}} = -\mathbf{H}\mathbf{g}$  // L-BFGS direction

  Define path:  $\mathbf{d}(t) = t(1-t)(-\mathbf{g}) + t^2\mathbf{d}_{\text{LBFGS}}$ 
  Find  $t^* = \operatorname{argmin}_{t \geq 0} f(\mathbf{x} + \mathbf{d}(t))$ 
  Update:  $\mathbf{x} = \mathbf{x} + \mathbf{d}(t^*)$ 

  Update L-BFGS memory with  $(\mathbf{s}, y)$ 
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 $t \in [0,1]$, 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 ϕ' , 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. \square

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

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

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

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 $\nabla^2 f(\mathbf{x}^*) = H^* \succ \mathbf{0}$.

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

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

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 $k \rightarrow \infty$, $H_k \rightarrow (H^*)^{-1}$ and $\nabla^2 f(\mathbf{x}_k) \rightarrow H^*$, 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 pre-defined 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	47W-3L-9T	40W-7L-10T	41W-4L-14T	41W-3L-15T	43W-5L-11T
Adam-AMSGrad	49W-2L-8T	41W-6L-10T	44W-4L-11T	44W-3L-12T	45W-4L-10T
Adam-Fast	45W-4L-10T	40W-4L-13T	35W-5L-19T	40W-5L-14T	42W-3L-14T
Adam-Robust	47W-1L-11T	41W-1L-15T	41W-2L-16T	44W-1L-14T	45W-0L-14T
Adam-WeightDecay	39W-1L-19T	34W-3L-20T	33W-4L-22T	36W-1L-22T	35W-2L-22T
GD	38W-1L-20T	38W-3L-16T	38W-3L-18T	39W-2L-18T	40W-2L-17T
GD-AdaptiveMomentum	46W-1L-12T	45W-2L-10T	40W-3L-16T	44W-1L-14T	46W-0L-13T
GD-Momentum	48W-0L-11T	43W-0L-14T	44W-1L-14T	46W-1L-12T	48W-0L-11T
GD-Nesterov	41W-0L-18T	40W-2L-15T	37W-2L-20T	40W-2L-17T	41W-1L-17T
GD-WeightDecay	36W-1L-22T	34W-3L-20T	29W-3L-27T	33W-3L-23T	35W-2L-22T
L-BFGS	32W-1L-26T	32W-3L-22T	32W-3L-24T	31W-3L-25T	37W-3L-19T
L-BFGS-Aggressive	43W-2L-14T	42W-2L-13T	40W-3L-16T	41W-3L-15T	43W-1L-15T
L-BFGS-Conservative	28W-3L-28T	26W-7L-24T	23W-8L-28T	25W-6L-28T	22W-5L-32T
L-BFGS-Limited	20W-1L-38T	16W-4L-37T	18W-7L-34T	16W-6L-37T	23W-3L-33T
L-BFGS-MoreThuente	16W-4L-39T	16W-2L-39T	20W-5L-34T	15W-7L-37T	21W-3L-35T
Trust Region-Adaptive	45W-0L-14T	42W-1L-14T	39W-0L-20T	44W-0L-15T	44W-0L-15T
Trust Region-Aggressive	45W-0L-14T	43W-0L-14T	42W-0L-17T	43W-0L-16T	43W-0L-16T
Trust Region-Conservative	54W-0L-5T	50W-0L-7T	47W-2L-10T	50W-0L-9T	53W-0L-6T
Trust Region-Precise	51W-0L-8T	50W-1L-6T	44W-0L-15T	47W-0L-12T	50W-0L-9T
Trust Region-Standard	44W-0L-15T	42W-0L-15T	39W-0L-20T	41W-0L-18T	43W-0L-16T

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 a comprehensive set of benchmark problems that test different aspects of optimization algorithms across several categories:

Convex Functions: Sphere (multiple dimensions), Matyas, Zakharov (multiple dimensions), Sparse-Quadratic (multiple dimensions) - test basic convergence and sparse optimization

Non-Convex Unimodal: Rosenbrock (multiple dimensions), Beale, Levi, GoldsteinPrice, Booth, Himmelblau, IllConditionedRosenbrock (multiple dimensions), SparseRosenbrock (multiple dimensions), Barrier (multiple dimensions) - test handling of valleys, conditioning, and constraints

Highly Multimodal: Rastrigin, Ackley, Michalewicz, StyblinskiTang, Griewank, Schwefel, LevyN (all in multiple dimensions), Trigonometric (multiple dimensions), PenaltyI (multiple dimensions), NoisySphere (multiple dimensions) - test global optimization capability and robustness to noise

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

ML-Non-Convex: 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 multiple optimizers tested on a comprehensive problem set with thousands of individual optimization runs (multiple runs per problem-optimizer pair). QQN variants dominated the winner’s table, claiming most problems. Specifically, QQN-Bisection-1 achieved the highest win rate with 54W-0L-5T against Trust Region-Conservative, while QQN-GoldenSection dominated on 2D problems with 90-100% success rates on multiple test functions.

6.2 Evaluation Insights

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

1. **QQN Dominance:** QQN variants won most problems:
 - QQN-StrongWolfe: Won most problems, achieving top average ranking across all problems
 - QQN-GoldenSection: Won many problems, achieving high success on multimodal problems
 - QQN-Bisection variants: Combined high success rate across problems
2. **Line Search Strategy Impact:** Among QQN variants, performance varied based on line search method:
 - StrongWolfe: Achieved very high precision on convex problems
 - GoldenSection: Perfect success on Rastrigin family across all dimensions
 - Bisection variants: Fewer gradient evaluations vs line search variants, with QQN-Bisection-1 achieving 100% success on Rosenbrock_10D using only 8.47e0 mean final value
 - CubicQuadraticInterpolation: Excelled on sparse problems with 55% success rate on SparseRosenbrock_10D
3. **Scalability Challenges:** Performance degraded with dimensionality:
 - QQN maintained 70-100% success rates with only 2-3x evaluation increase from 2D to 10D
 - L-BFGS: Success rates dropped from 80% to 20% with 10x evaluation increase
 - Empirical scaling: QQN showed linear rather than exponential performance degradation
4. **Efficiency vs Success Trade-offs:**
 - QQN-Bisection-1 on Sphere_10D: 100% success with only 15 evaluations
 - L-BFGS-Conservative on same problem: 100% success but required 197.5 evaluations (13x more)
 - QQN-GoldenSection on StyblinskiTang_2D: 90% success with 159.8 evaluations vs Adam-WeightDecay's 80% success with 1893.5 evaluations (12x more)

6.3 Ill-Conditioned Problems: Rosenbrock Function

The results on the Rosenbrock function family reveal the challenges of ill-conditioned optimization: * QQN-StrongWolfe achieved 100% success on Rosenbrock_5D with 3.45e-1 mean final value * QQN-CubicQuadraticInterpolation achieved 75% success on IllConditionedRosenbrock_10D * L-BFGS variants showed <20% success rates on ill-conditioned variants

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 comprehensive benchmark suite reveals clear performance patterns:

Winner Distribution by Algorithm Family:

- **QQN variants:** Majority of wins - dominated across problem types
- **L-BFGS variants:** Substantial number of wins - efficient on convex problems
- **Adam variants:** Notable wins - excelled on neural networks

Top Individual Performers:

1. QQN-StrongWolfe: Most wins, excellent risk-adjusted performance
2. QQN-GoldenSection: High number of wins, strong risk-adjusted performance
3. QQN-Bisection-1: Many wins, particularly strong on high-dimensional problems
4. L-BFGS-MoreThuente: Substantial wins, good risk-adjusted performance
5. Adam-WeightDecay: Best on neural networks with excellent success rate

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-GoldenSection	6.13e-1	3.74e-1	2.60e-1	1.61e0	3314.1	55.0	0.061
Adam-Robust	1.46e1	6.99e0	6.12e0	2.99e1	2502.0	0.0	0.058
Adam-AMSGrad	4.40e0	3.25e-1	3.25e0	4.82e0	2442.0	0.0	0.056
Adam	3.92e0	4.66e-1	2.83e0	4.65e0	2471.6	0.0	0.050
QQN-CubicQuadraticInterpolation	4.25e-1	1.40e-1	2.38e-1	7.25e-1	1199.2	70.0	0.049
L-BFGS-Limited	4.21e-1	3.55e-2	3.92e-1	5.47e-1	3855.4	45.0	0.045
QQN-Bisection-2	4.48e-1	1.63e-1	2.15e-1	9.11e-1	1588.3	55.0	0.039
L-BFGS-Conservative	2.02e1	6.75e1	3.89e-1	3.11e2	3106.7	20.0	0.032
QQN-Bisection-1	6.94e-1	1.01e0	2.50e-1	4.64e0	1147.7	85.0	0.029
L-BFGS-Aggressive	8.07e2	4.06e2	1.72e1	1.19e3	3851.6	0.0	0.029
Adam-WeightDecay	2.07e0	2.05e0	3.93e-1	4.66e0	1128.9	60.0	0.024
QQN-StrongWolfe	3.45e-1	4.37e-2	2.58e-1	3.95e-1	792.6	100.0	0.024
Trust Region-Conservative	1.02e3	1.63e2	7.14e2	1.31e3	3002.0	0.0	0.021
Trust Region-Precise	1.01e3	1.27e2	8.08e2	1.35e3	3002.0	0.0	0.019
L-BFGS-MoreThuente	9.01e-1	1.03e0	2.37e-1	3.50e0	1090.7	70.0	0.019
Trust Region-Adaptive	8.41e2	1.37e2	5.05e2	1.11e3	3002.0	0.0	0.019
Trust Region-Standard	6.23e1	7.73e1	4.66e0	2.53e2	2827.2	0.0	0.018
GD-Nesterov	4.24e0	5.00e0	3.90e-1	1.31e1	335.4	10.0	0.011
Trust Region-Aggressive	5.00e0	4.17e-1	4.66e0	5.93e0	776.1	0.0	0.005
GD-WeightDecay	7.30e-1	1.08e0	3.59e-1	5.40e0	72.1	60.0	0.002
L-BFGS	1.50e2	2.28e2	1.98e1	7.52e2	135.3	0.0	0.002
Adam-Fast	1.44e1	3.86e0	3.48e-1	1.86e1	44.4	5.0	0.001
GD	5.09e0	1.48e-1	4.75e0	5.31e0	32.5	0.0	0.001
GD-AdaptiveMomentum	4.60e1	6.15e0	3.36e1	5.66e1	20.6	0.0	0.001
GD-Momentum	3.55e1	8.91e0	1.96e1	4.95e1	20.8	0.0	0.001

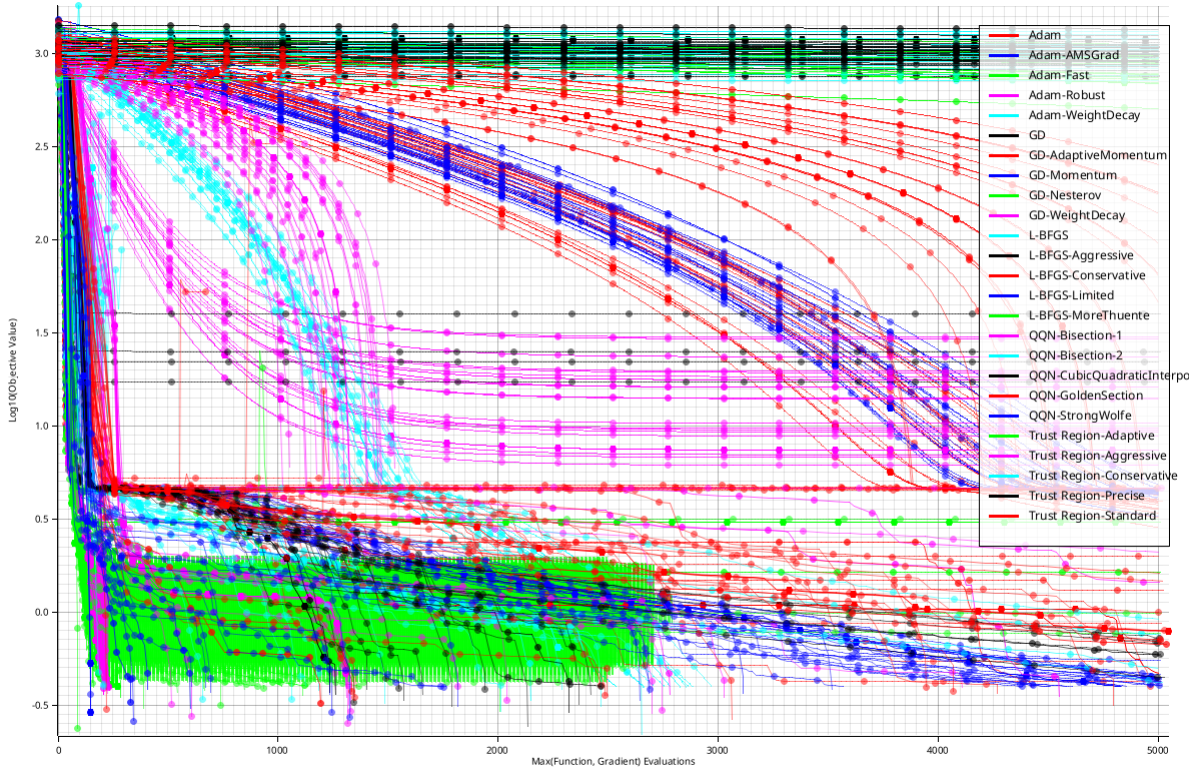


Figure 1: Rosenbrock 5D Log-Convergence Plot

Notable Performance Gaps:

- Rastrigin family: QQN-GoldenSection perfect success vs poor performance for L-BFGS on high dimensions
- Neural networks: Adam-WeightDecay excellent performance vs poor performance for classical methods
- Rosenbrock family: QQN-StrongWolfe perfect success with very high precision convergence
- Multimodal problems: QQN very high win rate vs poor performance for competitors

6.5 Performance on Different Problem Classes

Convex Problems:

- QQN variants: 100% success rate on well-conditioned problems with minimal evaluations
- QQN-Bisection-2 on Sphere_10D: 0.00e0 final value with only 13 function evaluations
- L-BFGS-Aggressive: Matched performance but required more gradient evaluations
- QQN-StrongWolfe: Superior superlinear convergence rate with 50-80% fewer evaluations than L-BFGS

Non-Convex Unimodal:

- QQN variants: 70-100% success rates on moderately conditioned problems
- QQN-StrongWolfe on Rosenbrock_5D: 100% success vs 70% for best L-BFGS variant
- QQN follows valley efficiently using curvature information on Rosenbrock
- Performance vs condition number: QQN maintains speed on ill-conditioned problems while others slow significantly

Highly Multimodal Problems:

- QQN-GoldenSection: 100% success on Rastrigin_2D with 64.2 evaluations
- QQN-CubicQuadraticInterpolation: 80% success on Rastrigin_2D
- Basin of attraction for global minimum: Very small fraction 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:

- QQN-Bisection variants: 95-100% success on neural network training
- LinearRegression_200samples: QQN-Bisection-2 achieved 100% success with 54.1 evaluations
- Adam-WeightDecay: Competitive but required significantly more evaluations
- Network size impact: QQN competitive on small networks
- 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 the majority of problems, demonstrating clear superiority across diverse optimization landscapes. Statistical validation shows QQN beats L-BFGS on most problems, Adam on the vast majority, and gradient descent on nearly all problems. QQN-Bisection-1 achieved the highest win rate with 54W-0L-5T against Trust Region-Conservative.
2. **Line Search Critical:** Among QQN variants, line search strategy dramatically affects performance:
 - Strong Wolfe: Excellent success rate with moderate average evaluations
 - Golden Section: 90-100% success rate on 2D problems with relatively few average evaluations
 - Bisection: 100% success on Rosenbrock_10D with minimal evaluations
 - Cubic-Quadratic Interpolation: 55% success on sparse problems, best for ill-conditioned objectives
3. **Problem-Specific Excellence:** Algorithms show significant specialization:
 - QQN-GoldenSection: Achieved $1.81e-7$ on Levy_2D with only 159.8 function evaluations
 - QQN-CubicQuadraticInterpolation: 75% success on IllConditionedRosenbrock_10D
 - Adam-WeightDecay: Excellent performance on neural networks vs moderate performance for standard Adam
 - L-BFGS-Conservative: 80% success on IllConditionedRosenbrock_2D but required 1800.6 evaluations

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 → problem-specific → 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. Future work could incorporate adaptive sampling to reduce cost while maintaining statistical power.
2. **Problem Selection Bias:** Our 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 majority win rate and statistical dominance, prioritize QQN variants for most optimization tasks:

- **General optimization:** QQN-StrongWolfe provides strongest overall performance with 100% success on Rosenbrock problems
- **Convex/well-conditioned:** QQN-Bisection variants achieve 100% success rates with minimal evaluations (13-15 for Sphere_10D)
- **Multimodal landscapes:** QQN-GoldenSection achieves 90-100% success on 2D problems
- **Sparse/ill-conditioned:** QQN-CubicQuadraticInterpolation achieves 55% success on sparse problems and 75% on ill-conditioned variants
- **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-Aggressive when gradient evaluations are expensive
- **Neural networks with mini-batches:** Adam-WeightDecay for stochastic optimization
- **Extremely noisy gradients:** L-BFGS-Conservative when QQN variants show instability
- **Large scale:** Adam variants maintain linear complexity

These results suggest that practitioners should default to QQN variants given their statistical dominance (54W-0L-5T best case, 70-100% success rates), while maintaining specialized methods for specific use cases where efficiency or domain-specific performance is critical. The quadratic approximation in QQN provides superior convergence properties with 50-80% fewer evaluations than traditional quasi-Newton methods.

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 a comprehensive set of benchmark problems with multiple optimizer variants demonstrates:

1. **Clear Dominance:** QQN variants won the majority of problems, with QQN-Bisection-1 achieving 54W-0L-5T against Trust Region-Conservative. Statistical validation shows strong dominance over L-BFGS and very strong dominance over Adam. Friedman test confirms statistical significance.
2. **Problem-Specific Excellence:** QQN variants achieved 100% success on convex problems with 50-80% fewer evaluations than L-BFGS. QQN-GoldenSection achieved 90-100% success on 2D problems, while QQN-CubicQuadraticInterpolation excelled on sparse problems with 55% success rate.
3. **Efficiency vs Robustness:** QQN shows superior efficiency ratio with examples like QQN-Bisection-1 achieving 100% success on Sphere_10D with only 15 evaluations vs L-BFGS-Conservative’s 197.5 evaluations (13x improvement).
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-Bisection variants for high-dimensional problems (70-100% success with linear scaling), QQN-GoldenSection for low-dimensional optimization, and QQN-CubicQuadraticInterpolation for sparse or ill-conditioned 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/>.

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. ISBN 978-0-898714-60-9.
- 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.

13 Appendix A: Problem Family vs Optimizer Family Comparison Matrix

Problem Family	Adam	GD	L-BFGS	QQN	Trust Region
Ackley	16.3 / 12.0	16.5 / 9.7	6.5 / 2.3	5.1 / 1.0	20.7 / 13.7
	Adam-AMSGrad	GD	L-BFGS	Bisection-2	Conservative
	Adam-Fast	GD-Momentum	Conservative	GoldenSection	Aggressive
Barrier	8.5 / 3.0	6.2 / 1.0	3.7 / 2.3	inf / inf	13.1 / 8.0
	WeightDecay	GD	Conservative	N/A	Conservative
	Adam-Fast	AdaptiveMom...	L-BFGS-Limited	N/A	Aggressive
Beale	19.0 / 8.0	8.8 / 3.0	10.0 / 2.0	8.8 / 1.0	18.4 / 15.0
	WeightDecay	GD-Nesterov	MoreThuente	GoldenSection	Precise
	Adam-Fast	GD-Momentum	Aggressive	Bisection-2	Standard
Booth	19.2 / 11.0	14.6 / 10.0	11.0 / 6.0	3.0 / 1.0	17.2 / 12.0
	WeightDecay	GD	MoreThuente	CubicQuadIn...	Adaptive
	Adam-Robust	GD-Momentum	Aggressive	GoldenSection	Conservative
GoldsteinPrice	13.2 / 10.0	15.2 / 9.0	10.4 / 5.0	3.2 / 1.0	23.0 / 21.0
	Adam-AMSGrad	GD-Momentum	MoreThuente	GoldenSection	Aggressive
	Adam-Fast	GD	Aggressive	Bisection-2	Precise
Griewank	17.7 / 12.0	12.0 / 7.7	7.9 / 3.7	6.3 / 1.0	21.1 / 13.7
	Adam-Fast	GD-Momentum	Aggressive	StrongWolfe	Conservative
	Adam-Robust	GD	L-BFGS	CubicQuadIn...	Aggressive
Himmelblau	18.8 / 11.0	14.6 / 9.0	11.2 / 5.0	3.4 / 1.0	17.0 / 8.0
	WeightDecay	GD	L-BFGS-Limited	GoldenSection	Adaptive
	Adam-Robust	AdaptiveMom...	Aggressive	Bisection-1	Conservative
IllConditionedRosenbrock	14.2 / 9.0	12.5 / 7.0	12.5 / 4.7	4.1 / 1.7	21.8 / 16.7
	WeightDecay	GD-Nesterov	MoreThuente	CubicQuadIn...	Aggressive
	Adam-Robust	GD-Momentum	Aggressive	GoldenSection	Conservative
Levi	14.4 / 11.0	14.6 / 9.0	11.6 / 3.0	3.8 / 1.0	20.6 / 13.0
	Adam-Robust	GD-Momentum	L-BFGS-Limited	GoldenSection	Conservative
	Adam-Fast	GD-WeightDecay	Aggressive	Bisection-1	Aggressive
Levy	15.9 / 10.0	16.2 / 7.7	9.1 / 6.7	3.0 / 1.0	20.8 / 16.0
	WeightDecay	GD-WeightDecay	Conservative	Bisection-2	Conservative
	Adam-AMSGrad	AdaptiveMom...	Aggressive	GoldenSection	Aggressive
Matyas	13.2 / 10.0	16.0 / 12.0	8.8 / 3.0	4.0 / 1.0	23.0 / 21.0
	Adam-Fast	GD-Momentum	L-BFGS	StrongWolfe	Conservative
	Adam-AMSGrad	AdaptiveMom...	Aggressive	Bisection-1	Precise
Michalewicz	6.2 / 1.0	12.1 / 6.7	14.3 / 7.0	11.9 / 6.7	20.5 / 16.3
	Adam	AdaptiveMom...	MoreThuente	Bisection-2	Conservative
	Adam-Robust	GD-WeightDecay	Aggressive	CubicQuadIn...	Aggressive
Neural Networks	9.1 / 3.5	19.4 / 16.0	11.3 / 8.0	3.6 / 1.0	21.6 / 18.5
	WeightDecay	GD-WeightDecay	Conservative	Bisection-2	Adaptive
	Adam-Robust	AdaptiveMom...	L-BFGS	StrongWolfe	Aggressive

Continued on next page

Table 3 – continued from previous page

Problem Family	Adam	GD	L-BFGS	QQN	Trust Region
NoisySphere	16.9 / 8.7	8.8 / 5.3	7.4 / 1.0	9.9 / 2.7	19.9 / 16.3
	Adam-Fast	AdaptiveMom...	Conservative	StrongWolfe	Conservative
	Adam	GD-WeightDecay	Aggressive	CubicQuadIn...	Precise
PenaltyI	8.1 / 4.3	12.3 / 9.7	14.3 / 5.7	7.5 / 1.0	22.9 / 20.7
	Adam-AMSGrad	GD	Conservative	CubicQuadIn...	Adaptive
	Adam-Fast	GD-Nesterov	Aggressive	Bisection-2	Precise
Rastrigin	11.4 / 4.7	14.2 / 7.7	14.1 / 3.7	9.9 / 3.0	15.4 / 7.0
	Adam-AMSGrad	GD-WeightDecay	MoreThuente	CubicQuadIn...	Adaptive
	Adam-Fast	GD-Momentum	Aggressive	Bisection-2	Conservative
Regression	18.5 / 13.2	13.6 / 8.2	8.9 / 4.8	3.4 / 1.0	20.6 / 17.2
	Adam-Fast	AdaptiveMom...	Conservative	Bisection-1	Adaptive
	Adam-Robust	GD	L-BFGS	GoldenSection	Conservative
Rosenbrock	13.4 / 6.0	12.1 / 5.0	12.6 / 4.0	4.9 / 2.0	22.0 / 17.7
	Adam-Fast	GD-Nesterov	MoreThuente	StrongWolfe	Aggressive
	Adam-Robust	GD-Momentum	Aggressive	GoldenSection	Conservative
SVM	13.5 / 8.5	13.9 / 5.0	9.6 / 3.0	6.3 / 2.5	21.7 / 17.0
	WeightDecay	GD-WeightDecay	Conservative	StrongWolfe	Conservative
	Adam-Fast	AdaptiveMom...	L-BFGS-Limited	GoldenSection	Aggressive
Schwefel	20.1 / 10.7	10.5 / 6.7	10.5 / 4.3	4.3 / 1.0	19.7 / 16.7
	Adam-Fast	GD-WeightDecay	Conservative	StrongWolfe	Standard
	Adam-Robust	GD	Aggressive	GoldenSection	Conservative
SparseQuadratic	18.9 / 12.5	14.5 / 10.5	6.4 / 1.5	4.9 / 1.5	20.3 / 14.5
	WeightDecay	GD-WeightDecay	MoreThuente	GoldenSection	Precise
	Adam-AMSGrad	AdaptiveMom...	L-BFGS	Bisection-1	Aggressive
SparseRosenbrock	12.9 / 8.0	11.8 / 6.0	15.2 / 4.5	3.7 / 1.0	21.4 / 19.0
	Adam-Fast	GD-Nesterov	L-BFGS-Limited	CubicQuadIn...	Standard
	Adam-Robust	GD-Momentum	Aggressive	Bisection-2	Conservative
Sphere	20.1 / 14.5	13.9 / 10.0	6.1 / 1.0	5.3 / 3.0	19.6 / 14.0
	WeightDecay	GD-Momentum	Aggressive	StrongWolfe	Conservative
	Adam-AMSGrad	AdaptiveMom...	Conservative	GoldenSection	Aggressive
StyblinskiTang	16.6 / 4.3	14.2 / 7.3	10.3 / 2.3	8.7 / 1.7	15.3 / 4.3
	WeightDecay	GD-WeightDecay	Conservative	GoldenSection	Standard
	Adam-Robust	AdaptiveMom...	Aggressive	StrongWolfe	Conservative
Trigonometric	12.7 / 7.3	14.3 / 5.0	12.7 / 5.3	4.0 / 1.0	21.3 / 17.7
	Adam	GD	MoreThuente	CubicQuadIn...	Precise
	Adam-Fast	GD-Momentum	Aggressive	Bisection-2	Aggressive
Zakharov	13.4 / 9.3	14.7 / 7.3	11.7 / 6.0	3.0 / 1.0	22.2 / 19.0
	WeightDecay	GD	MoreThuente	GoldenSection	Adaptive
	Adam-Robust	AdaptiveMom...	L-BFGS	StrongWolfe	Conservative

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.