# Stochastic Adaptive Optimization with Unreliable Inputs: A Unified Framework for High-Probability Complexity Analysis

Katya Scheinberg[*]     Miaolan Xie[†]

### Abstract

We consider an unconstrained continuous optimization problem where, in each iteration, gradient estimates may be arbitrarily corrupted with a probability greater than $\frac{1}{2}$. Additionally, function value estimates may exhibit heavy-tailed noise. This setting captures challenging scenarios where both gradient and function value estimates can be unreliable, making it applicable to many real-world problems, which can have outliers and data anomalies. We introduce an algorithmic and analytical framework that provides high-probability bounds on iteration complexity for this setting. The analysis offers a unified approach, encompassing methods such as line search and trust region.

## 1 Introduction

In this paper, we are interested in finding an $\epsilon$-stationary point of an unconstrained, differentiable, possibly non-convex function

$$\phi : \mathbb{R}^n \to \mathbb{R}.$$

We make the standard assumption that $\nabla \phi$ is $L$-Lipschitz and bounded below by some constant $\phi^*$, but knowledge of $L$ and $\phi^*$ is not assumed by the algorithm.

**Assumption 1.** *The gradient $\nabla \phi$ is $L$-Lipschitz continuous and $\phi$ is bounded below by some constant $\phi^*$.*

We consider a setting where neither the true function value $\phi(x)$ nor the true gradient $\nabla \phi(x)$ can be computed directly. Instead, the algorithm obtains all necessary function-related information by querying zeroth- and first-order stochastic oracles.

A key feature of our framework is its minimal assumptions on the oracles, which are only required to provide a reasonably accurate estimate with a certain probability. This flexibility allows for the modeling of highly unreliable or noisy inputs, a challenge in real-world optimization. Our setup is motivated by the prevalence of data anomalies and outliers in stochastic optimization, which can arise from diverse sources such as system errors, malfunctioning or byzantine machines, adversarial attacks and corruption, insufficient sampling, outdated data, or extreme, rare events. When such issues propagate into optimization, they can lead to poor performance and, consequently, suboptimal decision-making. This highlights the critical need for optimization algorithms that can reliably handle inputs containing bias, noise, and even adversarial corruption. Below, we give two motivating examples.

**Example 1: Expected Risk Minimization in Machine Learning.** Consider the stochastic optimization problem of expected risk minimization, where 5% of the dataset consists of outliers. In each iteration, a stochastic gradient estimate is obtained by averaging gradients over a randomly sampled mini-batch of data points. With a batch size of 32 (a common choice alongside 64, 128, and larger sizes), there is only about a 19% probability that the mini-batch contains no outliers. This means that gradient estimates are free from corruption with probability 19%, while with probability 81%, the mini-batch contains at least one outlier,

---

[*]School of Industrial and Systems Engineering, Georgia Tech, Atlanta, GA, USA; E-mail: katya.scheinberg@isye.gatech.edu
[†]Edwardson School of Industrial Engineering, Purdue University, West Lafayette, IN, USA; E-mail: miaolanx@purdue.edu

potentially corrupting the gradient estimate by an arbitrarily large amount. In contrast, function value estimates are often more robust to outliers. When the loss function is bounded (e.g., 0-1 loss), significantly corrupting the function value estimate with a mini-batch of 32 requires multiple outliers. For instance, the probability of sampling five or more outliers is only approximately 2%. Consequently, the probability of a large error in the function value estimate can be substantially lower than that for gradient estimates.

**Example 2: Derivative-Free Optimization.** In many real-world problems – such as those in chemical engineering and variational quantum algorithms (VQAs) – gradient information is unavailable, and only noisy or biased estimates of function values can be obtained. A common approach in these settings is to estimate gradients using finite-difference methods, which typically require $n + 1$ function value samples to construct a gradient estimate in dimension $n$. Even under relatively benign function value estimates, this approach can lead to frequent gradient corruption. Suppose each individual function value estimate is accurate with probability 99%, and is severely corrupted with probability 1%. In moderate dimensions, such as $n = 100$, the probability that all $n + 1$ function value estimates are accurate is only approximately $0.99^{101} \approx 36\%$. This leaves a $1 - 36\% = 64\%$ probability that the gradient estimate is significantly corrupted due to at least one inaccurate function value estimate.

**The Challenge.** In both settings described above, data anomalies can severely corrupt the inputs to optimization algorithms, highlighting a critical need for methods that are resilient to bias, noise, and adversarial corruption. However, most existing stochastic optimization algorithms, including popular methods like ADAM and AdaGrad, are not designed for such severely corrupted inputs, where gradient estimates may be arbitrarily corrupted with probability greater than 50% and function value estimates are subject to heavy-tailed noise. Most theoretical analyses typically rely on strong assumptions, such as bounded variance of gradient estimates, which fail to hold in the presence of heavy-tailed noise or adversarial attacks.

In this work, we introduce an adaptive algorithmic and analytical framework designed for such common yet challenging scenarios. The key idea is to leverage function value information to ensure reliability when gradient estimates are severely corrupted. Although function value estimates themselves can be subject to heavy-tailed noise or adversarial corruption, combining them with gradient estimates improves the overall reliability of the algorithm.

A key aspect of our framework is that our analysis is conducted with respect to the true, underlying objective function $\phi(x)$ (e.g., the true expected risk in ERM), which we assume is free from data anomalies. All sources of data corruption and noise are modeled as inaccuracies in the oracle outputs. We begin by formalizing the probabilistic oracles that model the noisy and corrupted inputs available to the algorithm.

## 1.1 Oracles

**Stochastic Zeroth-Order Oracle ($\mathbf{SZO}(\epsilon_f, q, \zeta_q)$)** Given a point $x$, the oracle computes $f(x, \Xi^0(x))$, a (random) estimate of the function value $\phi(x)$. $\Xi^0(x)$ is a random variable (whose distribution may depend on $x$). We assume the absolute value of the estimation error $E(x) = |f(x, \Xi^0(x)) - \phi(x)|$ (we omit the dependence on $\Xi^0(x)$ for brevity) to have bounded expectation and a bounded $q$-th centered moment, for some $q \geq 2$.

$$\mathbb{E}_{\Xi^0(x)}[E(x)] \leq \epsilon_f \text{ and } \mathbb{E}_{\Xi^0(x)}(E(x) - \mathbb{E}[E(x)])^q \leq \zeta_q. \tag{1}$$

The input to the oracle is $x$, the output is $f(x, \Xi^0(x))$ or $f(x, \Xi^0)$ for short, and the values $(\epsilon_f, q, \zeta_q)$ are intrinsic to the oracle.

**Remark**

- In the case where $q = 2$, the assumption on the zeroth-order oracle reduces to the noise having bounded expectation and bounded variance

$$\mathbb{E}[E(x)] \leq \epsilon_f \text{ and } \text{Var}(E(x)) \leq \zeta_2 = \sigma^2. \tag{2}$$

- In the case where the estimation error is a subexponential random variable, then the noise has exponentially decaying tail and bounded moments of all orders. In this case, for some parameters $(\nu, b)$,

$$\mathbb{E}\left[E(x)\right] \le \epsilon_f \text{ and } \mathbb{E}\left[\exp\{\lambda(E(x) - \mathbb{E}[E(x)])\}\right] \le \exp\left(\frac{\lambda^2 \nu^2}{2}\right), \quad \forall \lambda \in \left[0, \frac{1}{b}\right], \tag{3}$$

or equivalently, for some parameters $(\lambda, a)$,

$$\mathbb{E}\left[E(x)\right] \le \epsilon_f \text{ and } \mathbb{P}\left(E(x) \ge t\right) \le e^{\lambda(a-t)}, \text{ for any } t > 0 \tag{4}$$

or equivalently, for some constant $K \ge 1$,

$$\mathbb{E}\left[E(x)\right] \le \epsilon_f \text{ and } \mathbb{E}(E(x) - \mathbb{E}[E(x)])^q \le \zeta_q, \text{ where } \zeta_q = (Kq)^q, \text{ for all } q \ge 1. \tag{5}$$

This is the zeroth-order oracle assumption used in [JSX24, CBS24, BXZ25, MWX23, SX23]. For a proof of the equivalence of (3), (4), and (5), see Proposition 2.7.1 in [Ver18].

**Alternative assumption on function value differences.** Instead of assuming bounded moments on individual function value errors, we may consider a more general assumption on the difference of function values. For any two points $x$ and $x^+$, we could assume that the random variable

$$D(x, x^+) = \left|(f(x, \Xi^0) - f(x^+, \Xi^0)) - (\phi(x) - \phi(x^+))\right| \tag{6}$$

satisfies the bounded moments condition:

$$\mathbb{E}_{\Xi^0}\left[D(x, x^+)\right] \le \epsilon_f \quad \text{and} \quad \mathbb{E}_{\Xi^0}\left(D(x, x^+) - \mathbb{E}[D(x, x^+)]\right)^q \le \zeta_q. \tag{7}$$

This assumption is more general than the individual bounded moments assumption above. It is particularly beneficial when using common random numbers. Although all results extend under this broader assumption, we adopt the individual bounded-moments assumption for clarity of exposition.

In real-life, there may be settings that are adversarial rather than stochastic. To address such settings, we introduce the following corrupted zeroth-order oracle designed for adversarial settings.

**Corrupted Zeroth-Order Oracle (CZO($\epsilon_f, \epsilon_c, \delta_0$)).** Given a point $x$, the oracle computes $f(x, \Xi^0(x))$, where $\Xi^0(x)$ is a random variable, whose distribution may depend on $x$, $\epsilon_f, \epsilon_c$ and $\delta_0$, that satisfies

$$\mathbb{P}_{\Xi^0(x)}(|\phi(x) - f(x, \Xi^0(x))| \le \epsilon_f) \ge 1 - \delta_0, \quad \mathbb{P}_{\Xi^0(x)}(|\phi(x) - f(x, \Xi^0(x))| \le \epsilon_f + \epsilon_c) = 1. \tag{8}$$

where $\delta_0 \in (0, 1)$. We view $x$ as the input to the oracle, $f(x, \Xi^0(x))$ as the output and the values $(\epsilon_f, \epsilon_c, \delta_0)$ as values intrinsic to the oracle. By definition, $|f(x, \Xi^0(x)) - \phi(x)|$ is a bounded random variable. While the CZO can be viewed as a special case of the SZO—since a bounded random variable has all moments bounded—it is instructive to analyze it separately. The CZO framework is particularly well-suited for modeling noise from adversarial settings or data outliers, as it explicitly distinguishes between two scenarios:

- With probability $1 - \delta_0$, the error is small, resulting from typical stochastic sampling or minor bias.
- With probability $\delta_0$, the error is large, due to outliers or adversarial corruption.

This modeling differs from the SZO, and the distinction is crucial because the oracle parameters have distinct interpretations. For instance, $\epsilon_f$ in the CZO represents a probabilistic error bound, whereas in the SZO it represents a bound on the expected error.

**Remark** [Expected Risk Minimization with Corrupted Data] An example that motivates the CZO oracle is expected risk minimization (ERM) in machine learning with a corrupted dataset. Consider an objective $\phi(x) = \mathbb{E}_{d \sim \mathcal{D}}[l(x, d)]$, where $l(x, d)$ is the loss for a model with parameters $x$ on a data point $d$. In practice, $\phi(x)$ is estimated by an empirical average over a mini-batch of data, $f(x) = \frac{1}{B}\sum_{i=1}^{B} l(x, d_i)$.

Now, suppose a fraction of the underlying data distribution is corrupted. When we draw a mini-batch, the number of corrupted samples in it is random. We consider a setting where the loss function is bounded (e.g., 0-1 bounded loss or some other bounded loss), which ensures the error from corrupted samples is also bounded. We can categorize the outcome based on the level of corruption in the mini-batch:

- With probability $1 - \delta_0$, the number of corrupted data points in the mini-batch is small, such that their effect on the function value estimate is limited. In this case, the estimation error $|f(x) - \phi(x)|$ is bounded by a baseline error term $\epsilon_f$.

- With probability $\delta_0$, the mini-batch contains a significant number of corrupted points. This leads to a larger estimation error, which is bounded by $\epsilon_f + \epsilon_c$, where $\epsilon_c$ represents the maximum additional error due to the high level of corruption.

This scenario maps directly to the CZO$(\epsilon_f, \epsilon_c, \delta_0)$ definition, illustrating its utility in modeling practical machine learning problems. Note that under this model, the expected error is bounded by $\epsilon_f + \delta_0 \cdot \epsilon_c$.

**Stochastic First-Order Oracle (SFO$(r(\cdot), \delta_1)$).** Given a point $x$ and a parameter $\alpha > 0$ (e.g., a step size or trust-region radius), the oracle computes and outputs $g(x, \Xi^1(x))$ or $g(x, \Xi^1)$ for short, which is a (random) estimate of the gradient $\nabla \phi(x)$. The distribution of the random variable $\Xi^1(x)$ may depend on $x$. We assume that with probability at least $1 - \delta_1$, the estimation error is bounded by an accuracy tolerance function $r(\alpha)$:

$$\mathbb{P}_{\Xi^1(x)} \left( \|\nabla \phi(x) - g(x, \Xi^1(x))\| \leq r(\alpha) \right) \geq 1 - \delta_1. \tag{9}$$

The function $r(\alpha)$ defines the required relative accuracy, which may depend on algorithm parameters like $\alpha$ and intrinsic oracle parameters (related to bias and variance terms). With probability $\delta_1 \in [0, 1)$, the gradient estimate may be arbitrarily corrupted.

**Connections to Algorithmic Frameworks.** The general form in (9) unifies accuracy conditions for many optimization methods in the literature, including:

- **Trust-Region Methods:** In trust-region-like methods, the model accuracy is typically required to be proportional to the trust-region radius $\alpha$. For instance, one typical condition [CBS24] is:

$$\mathbb{P}_{\Xi^1(x)} \left( \|\nabla \phi(x) - g(x, \Xi^1(x))\| \leq \epsilon_g + \kappa \alpha \right) \geq 1 - \delta_1. \tag{10}$$

- **Line-Search Methods:** In line-search or step-search methods, the accuracy requirement might take the following form [JSX24]:

$$\mathbb{P}_{\Xi^1(x)} \left( \|g(x, \Xi^1(x)) - \nabla \phi(x)\| \leq \max\{\epsilon_g, \min\{\tau, \kappa \alpha\} \|g(x, \Xi^1(x))\|\} \right) \geq 1 - \delta_1, \tag{11}$$

or

$$\mathbb{P}_{\Xi^1(x)} \left( \|g(x, \Xi^1(x)) - \nabla \phi(x)\| \leq \max\{\epsilon_g, \tau \|\nabla \phi(x)\|\} \right) \geq 1 - \delta_1. \tag{12}$$

While the analyses in this paper apply equally well to both (11) and (12), we will primarily work with (11) due to its more common use in the literature.

In both examples, the presence of $\epsilon_g \geq 0$ in the accuracy requirement allows for a constant, irreducible amount of bias in the gradient estimation, while $\kappa \geq 0$ and $\tau \geq 0$ are constants controlling the relative accuracy. The values $(\epsilon_g, \tau, \kappa, \delta_1)$ are intrinsic to the oracle.

**Remark** By definition, we allow the first-order oracle to be corrupted by an arbitrarily large amount in each iteration with probability $\delta_1$. This model is well-suited for scenarios like Empirical Risk Minimization (ERM) with corrupted data. For instance, consider estimating the gradient from a mini-batch of data where a fraction of the underlying dataset is corrupted by outliers or adversarial examples. The probability of sampling at least one corrupted data point in a mini-batch depends on the corruption ratio and the batch size. If a mini-batch is contaminated, the resulting gradient estimate can be arbitrarily poor, pointing in any direction with any magnitude. This corresponds to the failure case of the SFO, which occurs with probability $\delta_1$. Conversely, if the mini-batch is clean, the gradient estimate is still stochastic, but its error can typically be bounded, satisfying the accuracy condition $\|\nabla \phi(x) - g(x, \Xi^1)\| \leq r(\alpha)$ with probability $1 - \delta_1$. The SFO framework thus naturally models the behavior of gradient estimation in the presence of data corruption.

The stochastic derivative-free optimization setting discussed in the introduction provides another natural application of our framework. In the scenario with $n = 100$, where each function value estimate is corrupted

or has a large noise with probability $\delta_0 = 0.01$, the resulting finite-difference gradient estimate can be arbitrarily corrupted with a probability as high as $\delta_1 = 0.64$. This scenario is precisely what our SZO and SFO definitions are designed to handle.

To summarize, the key oracle parameters that characterize the quality of the estimates are:

- For SZO: $\epsilon_f$ (expected error bound), $q \geq 2$ (moment order), $\zeta_q$ (moment bound)

- For CZO: $\epsilon_f$ (baseline error), $\epsilon_c$ (additional corruption), $\delta_0$ (corruption probability)

- For SFO: $r(\cdot)$ (accuracy tolerance function), $\delta_1$ (failure probability)

**Optimization in Highly Unreliable or Noisy Environments.** A central challenge addressed in this paper is optimization in what we term a *highly unreliable* or *highly noisy* environment. This setting is characterized by a low probability of obtaining reasonable gradient and function value estimates from the oracles, meaning in every iteration the algorithm must navigate using local models that are more often misleading than helpful. Specifically, we consider scenarios where the joint probability of receiving sufficiently accurate gradient and function value estimates in any given iteration can be less than $\frac{1}{2}$. Our analysis is specifically designed to ensure reliability in this challenging regime. The formal probabilistic framework and related notions for this setting will be made precise in the subsequent sections.

**Related literature.** Our work contributes to the growing literature on adaptive optimization methods with stochastic or probabilistic oracles. The foundation for these methods was established in [CS17] and [GRVZ18]. In [CS17], the authors analyzed the expected convergence of line-search and cubic regularized Newton methods using probabilistic models for gradient information. In [GRVZ18], the authors establish high-probability convergence for a trust-region method [BSV14] that permits arbitrarily inaccurate gradient estimates with probability greater than $\frac{1}{2}$. However, both of these works assume that all function value estimates are exact.

The line search method was extended by [BCS21] to allow for bounded noise in function values, and by [PS20] to a more general setting with unbounded noise with small variance, while under more complex step and accuracy adjustment process and zeroth-order oracle requirements. These studies initially concentrated on achieving convergence in expectation. More recently, [JSX24] introduced the first high-probability iteration complexity bounds for a stochastic step-search method, while [JSX25] developed a comprehensive framework for analyzing the expected and high-probability sample complexity of this class of adaptive algorithms. All of these contributions assume that the first- and higher-order oracle outputs are sufficiently accurate with probability greater than $\frac{1}{2}$ in every iteration.

The stochastic trust-region literature has been developed in parallel. [BCMS19] established convergence rate in expectation for a version of stochastic trust-region methods. Building on this and [JSX24], [CBS24] established first- and second-order high-probability complexity bounds for trust-region methods with noisy oracles, assuming the oracle outputs are sufficiently accurate with high probability in every iteration. Similar progress has been made for adaptive cubic regularization methods. For instance, [BGMT20, BGMT22] developed stochastic cubic regularization methods for finite-sum minimization and nonconvex optimization, and [SX23] provided a high-probability complexity bound for a stochastic adaptive regularization method with cubics. Recent work has also extended these adaptive stochastic methods to more specialized settings. [BXZ25] introduced a sequential quadratic programming (SQP) method with high-probability complexity bounds for solving nonlinear equality-constrained stochastic optimization, and [MWX23] developed stochastic quasi-Newton methods with high-probability iteration complexity bounds for settings where utilizing common random numbers in gradient estimation is not feasible.

Importantly, most existing works in the area assume that the first- and higher-order oracle outputs are sufficiently accurate with probability greater than $\frac{1}{2}$ in every iteration. A notable exception is [GRVZ18], which establishes a high-probability iteration complexity bound for a trust-region method that allows the gradient estimate in each iteration to have arbitrarily large error with probability greater than $\frac{1}{2}$. However, a key limitation of [GRVZ18] is the assumption that function value estimates are exact. In this paper, we tackle the challenging scenario where *both* the gradient estimate can have arbitrarily large error with probability

greater than $\frac{1}{2}$, *and* function value estimates may be corrupted by heavy-tailed noise. As a result, we extend the application and scope of previous works by allowing significantly more relaxed assumptions on the inputs of the algorithm. We develop a unified algorithmic framework that encompasses both trust-region and line-search methods, together with a unified analysis that applies to all algorithms within this framework.

**Our contribution.** The main contributions of this work are as follows:

- **A Unified Algorithmic Framework:** We introduce a unified algorithmic and analytical framework for adaptive optimization that is resilient to highly unreliable inputs. This framework is general enough to encompass both trust-region and line-search methods and is specifically designed to handle settings where gradient estimates can be arbitrarily corrupted with high probability and function value estimates are subject to heavy-tailed noise.

- **Analysis for Highly Unreliable or Noisy Environments:** We remove the restrictive assumption that in every iteration, a reliable model (i.e. a sufficiently accurate gradient and function value estimate) must be obtained with a probability greater than $\frac{1}{2}$. Our analysis is the first to formally handle scenarios where with high probability (could be close to one), the model can be arbitrarily corrupted in every iteration with the function value estimates subject to heavy-tailed noise, a critical feature for ensuring robustness against data anomalies and corruption.

- **High-Probability Complexity Bounds:** We provide the first iteration complexity analysis for this setting that establishes (overwhelmingly) high-probability bounds on the stopping time. We show that the tail probability of this stopping time decays either exponentially or polynomially, depending on the assumptions on the zeroth-order oracle (e.g., whether the noise is sub-exponential or heavy-tailed).

In the remainder of this paper, we will first present the unified algorithmic framework in Section 2, along with its instantiations for both trust region and line search methods. Following this, Section 3 will introduce the general stochastic process that underpins the framework, as well as the key definitions and notation used throughout our discussion. Finally, Section 4 will detail the main theoretical results, providing high-probability iteration complexity bounds and proofs for the general framework, as well as for its specific instantiations in trust region and line search methods.

## 2 Algorithm

In this section, we present our unified algorithmic framework for adaptive optimization with highly unreliable inputs. The framework handles both gradient and function value noise through the oracles defined in the introduction. It encompasses several algorithms in the literature, including trust region and line search methods, as special cases.

### 2.1 General Algorithmic Framework

The algorithmic framework employs an adaptive step-size strategy in which the step size parameter is adjusted dynamically in each iteration based on the estimated progress of the algorithm. In each iteration, we construct a local model of the objective function using gradient and function value estimates obtained from the oracles. A candidate point is then computed, typically by minimizing this model. The quality of this candidate point is evaluated using a *sufficient descent test*, which compares the estimated reduction in function value (obtained from the zeroth-order oracle) to a quantity related to the reduction in the model value. Since function value estimates are noisy, the sufficient descent test is typically relaxed compared to its deterministic counterpart to account for this noise.

Based on the outcome of the sufficient descent test, the algorithm takes one of three actions:

- If the candidate point passes the sufficient descent test and meets additional acceptance criteria, it is accepted and the step size parameter is increased. This encourages larger steps in subsequent iterations, allowing for potentially faster progress. Such steps are called *successful steps*.

- If sufficient reduction is achieved but additional acceptance criteria are not met, the candidate point is accepted but the step size parameter is decreased.

- If sufficient reduction is not achieved, the candidate point is rejected and the step size parameter is decreased.

Steps in these latter two categories are called *unsuccessful steps*. Decreasing the step size parameter encourages more conservative steps within a smaller neighborhood where the model is more likely to be accurate.

Importantly, this mechanism allows the algorithm to dynamically increase or decrease the step size based on the estimated progress of the algorithm from iteration to iteration, allowing it to adapt to the local geometry of the function and the level of noise from the oracles. The framework is presented in Algorithm 1.

---

**Algorithm 1** General Algorithmic Framework

---

1: **Initialization:** Choose parameters $\theta \in (0, 1)$ (sufficient reduction parameter), $\gamma_{\text{inc}} > 1$ (step size increase factor), $\gamma_{\text{dec}} \in (0, 1)$ (step size decrease factor), initial point $x_0 \in \mathbb{R}^n$, and initial step size $\alpha_0 > 0$. Set $k \leftarrow 0$.

2: **Main Loop:** For $k = 0, 1, 2, \ldots$

3: **Step 1: Model Construction and Step Computation**
    Query the first-order oracle to obtain gradient estimate $g_k = g(x_k, \Xi_{1,k})$, query the zeroth-order oracle to obtain function value estimate $f_k = f(x_k, \Xi_{0,k})$ at $x_k$, construct a local model $m_k$ of $\phi$ around $x_k$ using $g_k$ and $f_k$, and compute a step $s_k(\alpha_k)$ that achieves sufficient model reduction.

4: **Step 2: Compute Trial Point and Acceptance Test**
    Set trial point $x_k^+ \leftarrow x_k + s_k(\alpha_k)$, query the zeroth-order oracle to obtain $f_k^+ = f(x_k^+, \Xi_{0,k}^+)$ (function value estimate at $x_k^+$), compute reduction estimate $f_k - f_k^+$, and check sufficient reduction conditions (parameterized by $\theta$) with respect to the predicted model reduction $m_k(x_k) - m_k(x_k^+)$.

5: **Step 3: Update Rule**

6: **if** sufficient reduction is achieved **then**

7:     $x_{k+1} \leftarrow x_k^+$ (accept step)

8:     **if** additional acceptance criteria are met **then**

9:         $\alpha_{k+1} \leftarrow \gamma_{\text{inc}}\alpha_k$ (increase step size)

10:     **else**

11:         $\alpha_{k+1} \leftarrow \gamma_{\text{dec}}\alpha_k$ (decrease step size)

12: **else**

13:     $x_{k+1} \leftarrow x_k$ (reject step)

14:     $\alpha_{k+1} \leftarrow \gamma_{\text{dec}}\alpha_k$ (decrease step size)

15: **Step 4:** Set $k \leftarrow k + 1$ and continue.

---

The framework above is deliberately general; the specific forms of the model $m_k$, step $s_k(\alpha_k)$, and acceptance criteria depend on the specific instantiations of the framework. We next show how a trust region and line search method fit into this framework.

## 2.2 Trust Region Instantiation

In the trust region approach, the local model $m_k$ is typically as follows:

$$m_k(x_k + s) = f_k + g_k^T s + \frac{1}{2} s^T H_k s, \tag{13}$$

where $H_k$ can be any matrix as long as there exists some constant $\kappa_H$ such that $\|H_k\| \leq \kappa_H$ for all $k$.

The candidate step is computed by approximately minimizing the model in the ball of radius $\alpha_k$ around the current iterate:

$$s_k(\alpha_k) = \arg \min_{s \in B(x_k, \alpha_k)} m_k(x_k + s). \tag{14}$$

The accuracy requirement for the first-order oracle in the trust region method, following [CBS24], is given by Equation (10), where a smaller trust region radius translates to a higher accuracy requirement.

---

**Algorithm 2** Trust Region Based Algorithm

---

**Initialization:** Starting point $x_0$; initial radius $\alpha_0 > 0$; $\eta_1 > 0$, $\eta_2 > 0$, $\delta_1 \in [0, 1)$, $\kappa_{\text{fcd}} \in (0, 1]$, $\gamma_{\text{inc}} > 1$, $\gamma_{\text{dec}} \in (0, 1)$, $\epsilon_f > 0$.

**For** $k = 0, 1, 2, \ldots$ **do**

Compute vector $g_k$ using $\text{SFO}(r(\alpha_k), \delta_1)$ and matrix $H_k$ that has a norm bounded by some constant $\kappa_H$.

Compute $s_k$ by approximately minimizing $m_k$ in $B(x_k, \alpha_k)$ so that it satisfies

$$m_k(x_k) - m_k(x_k + s_k) \geq \frac{\kappa_{\text{fcd}}}{2} \|g_k\| \min\left\{\frac{\|g_k\|}{\|H_k\|}, \alpha_k\right\}.$$

Compute $f_k$ and $f_k^+$ using Zeroth-order oracle, with input $x_k$, and $x_k + s_k$, and then compute

$$\rho_k = \frac{f_k - f_k^+ + 2\epsilon_f}{m_k(x_k) - m_k(x_k + s_k)}.$$

**if** $\rho_k \geq \eta_1$ **then**

    Set $x_{k+1} = x_k + s_k$ and

$$\alpha_{k+1} = \begin{cases} \gamma_{\text{inc}}\alpha_k & \text{if } \|g_k\| \geq \eta_2\alpha_k \\ \gamma_{\text{dec}}\alpha_k & \text{if } \|g_k\| < \eta_2\alpha_k \end{cases}$$

**else**

    Set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_{\text{dec}}\alpha_k$.

---

This algorithm follows the standard trust-region framework, where the additional term $2\epsilon_f$ accounts for noise in function evaluations. It matches Algorithm 1 of [CBS24], with the sole difference that the trust-region radius is updated using $\gamma_{\text{inc}}$ and $\gamma_{\text{dec}}$ (rather than $\gamma$ and $\gamma^{-1}$). In this trust region algorithm, an iteration is defined as successful if the sufficient reduction condition $\rho_k \geq \eta_1$ and the gradient norm condition $\|g_k\| \geq \eta_2\alpha_k$ are both satisfied.

## 2.3 Line Search Instantiation

One can obtain a step search or line search method by setting the model to be

$$m_k(x_k + s) = f_k + g_k^T s + \frac{1}{2\alpha_k}s^T s$$

and $s_k(\alpha_k) = -\alpha_k g_k$.

The accuracy requirement for the first-order oracle in the line search method, following [JSX24], is given by Equation (11). A smaller step size parameter translates to a higher accuracy requirement.

Algorithm 3 is an instantiation of the general framework in Algorithm 1 for a line search method. It is identical to the SASS Algorithm from [JSX24] but with an additional criterion of:

$$\|g_k\| < \epsilon_{\text{rej}},$$

for some $\epsilon_{\text{rej}} > 0$. The condition $f_k^+ \leq f_k - \alpha_k\theta\|g_k\|^2 + 2\epsilon_f$ checks if sufficient descent is achieved. If it is, the step is accepted. The additional condition $\|g_k\| \geq \epsilon_{\text{rej}}$ then determines whether the step size should be increased. If the suffcient descent condition is not met, or the norm of the gradient estimate is too small, then the iteration is unsuccessful, and the step size is decreased. This mechanism prevents the algorithm

from increasing the step size when the gradient estimate is small, which might indicate unreliable estimates rather than actual progress towards a solution.

---

**Algorithm 3** Line Search Based Algorithm

---

**Initialization:** Choose $\theta \in (0,1)$, $\gamma_{\text{inc}} > 1$, $\gamma_{\text{dec}} \in (0,1)$, $\epsilon_f > 0$, $\epsilon_{\text{rej}} > 0$, $x_0 \in \mathbb{R}^n$, and $\alpha_0 > 0$. Set $k \leftarrow 0$.

**For** $k = 0, 1, 2, \ldots$ **do**

Compute gradient estimate $g_k$ using $\text{SFO}(r(\alpha_k), \delta_1)$.

Compute search direction $d_k = -g_k$ and trial point $x_k^+ = x_k + \alpha_k d_k$.

Compute $f_k$ and $f_k^+$ using Zeroth-order oracle, with input $x_k$ and $x_k^+$.

**if** $f_k^+ \leq f_k - \alpha_k \theta \|g_k\|^2 + 2\epsilon_f$ **then**

    Set $x_{k+1} = x_k^+$ and

$$\alpha_{k+1} = \begin{cases} \gamma_{\text{inc}}\alpha_k, & \text{if } \|g_k\| \geq \epsilon_{\text{rej}} \\ \gamma_{\text{dec}}\alpha_k & \text{if } \|g_k\| < \epsilon_{\text{rej}} \end{cases}$$

**else**

    Set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_{\text{dec}}\alpha_k$.

---

## 2.4 Key Adaptive Features

In the algorithm, both the **step size parameter** and the **gradient accuracy requirement** are adaptive. The step size parameter is adjusted based on the estimated progress of the algorithm from iteration to iteration. It can increase or decrease in each iteration, and it is not monotonic. This allows the algorithm to adapt to the local landscape of the function.

In each iteration, there are two reasons the algorithm might reject a candidate step. Either:

- The oracle estimates are so inaccurate so that the candidate step $s_k$ is not a descent direction, or

- The step size parameter is so large so that taking this step would overshoot.

When the algorithm rejects a step, it decreases the step size parameter and computes a new candidate step. This reduction addresses the second potential cause of failure. Additionally, because the accuracy requirement of the first-order oracle is governed by $r(\alpha_k)$, a smaller step size parameter demands higher accuracy from the oracle. Thus, by adaptively adjusting the step size, the algorithm simultaneously addresses both potential causes of failure.

In Section 4, we present a **unified analysis framework** that allows *any* algorithm within this algorithmic framework to be analyzed using the same theoretical approach, despite differences in the choice of model $m_k$, candidate step calculation, and sufficient descent condition.

# 3 Definitions and Notation

We introduce the key definitions and notation used throughout the paper. The algorithm generates a stochastic process, and we define the relevant random variables and concepts needed for the analysis.

## 3.1 Stochastic Process and Filtration

Let $M_k$ denote the collection of random variables $\{\Xi_{0,k}, \Xi_{0,k}^+, \Xi_{1,k}\}$, whose realizations are $\{\xi_{0,k}, \xi_{0,k}^+, \xi_{1,k}\}$, where:

- $\Xi_{0,k}$ and $\Xi_{0,k}^+$ are the random variables from the zeroth-order oracle at $X_k$ and $X_k^+$, respectively,

- $\Xi_{1,k}$ is the random variable from the first-order oracle at $X_k$.

At iteration $k$, we define the following random variables:

- $X_k$: the random iterate at step $k$, with realization $x_k$,

- $G_k$: the random gradient estimate from the first-order oracle, with realization $g_k$,

- $E_k = |f(X_k, \Xi_{0,k}) - \phi(X_k)|$: the random absolute error in the zeroth-order oracle at $X_k$, with realization $e_k$,

- $E_k^+ = |f(X_k^+, \Xi_{0,k}^+) - \phi(X_k^+)|$: the random absolute error in the zeroth-order oracle at $X_k^+$, with realization $e_k^+$,

- $\mathcal{A}_k$: the random step size parameter at step $k$, with realization $\alpha_k$.

The algorithm generates a stochastic process $\{(G_k, E_k, E_k^+, X_k, \mathcal{A}_k)\}$ with realizations $(g_k, e_k, e_k^+, x_k, \alpha_k)$ adapted to the filtration $\{\mathcal{F}_k : k \geq 0\}$, where $\mathcal{F}_k = \sigma(M_0, M_1, \ldots, M_k)$ is the $\sigma$-algebra generated by $M_0, \ldots, M_k$. All the random variables defined above are measurable with respect to $\mathcal{F}_k$. Note that $X_k$ and $\mathcal{A}_k$ are in addition measurable with respect to $\mathcal{F}_{k-1}$.

## 3.2 Key Definitions

**Definition 1** (True Iteration). *We say that iteration $k$ is **true** if both the gradient estimate and function value estimates are sufficiently accurate:*

$$\|\nabla\phi(X_k) - G_k\| \leq r(\mathcal{A}_k), \quad and \tag{15}$$

$$E_k + E_k^+ \leq 2\epsilon_f, \tag{16}$$

*where $r(\mathcal{A}_k)$ is the accuracy tolerance function from the first-order oracle definition, and $\epsilon_f$ is an upper bound on the expected error of the zeroth-order oracle.*

The probability of an iteration being true, denoted by $p$, is central to our framework. A lower bound on $p$ can be established from the oracle failure probabilities. The first-order oracle fails with probability $\delta_1$. For the corrupted zeroth-order oracle, since $\delta_0$ is the failure probability of a single query, and an iteration requires two CZO queries (at $x_k$ and $x_k^+$), the probability that at least one of these is corrupted can be bounded by $2\delta_0$ using a union bound. This provides a lower bound on the probability of a true iteration: $p \geq 1 - \delta_1 - 2\delta_0$.

A key focus of this paper is the challenging regime where this probability $p$ is low, which we formalize as follows.

**Definition 2** (Highly Unreliable Environment). *An optimization environment is considered **highly unreliable** if the probability $p$ of an iteration being true satisfies $p < 1/2$.*

When working with the corrupted zeroth-order oracle, this condition implies $\delta_1 + 2\delta_0 > 1/2$ in terms of the oracle failure probabilities. In such a regime, the local model provided to the algorithm is more likely to be misleading than helpful in any given iteration.

**Definition 3** (Successful Iteration). *We say that iteration $k$ is **successful** if the iterate is updated and the step size is increased, i.e., $X_{k+1} = X_k^+$ and $\mathcal{A}_{k+1} \geq \mathcal{A}_k$. Otherwise, the iteration is **unsuccessful**.*

We define the following indicator random variables:

- $I_k := \mathbb{1}\{\text{iteration } k \text{ is true}\}$

- $\Theta_k := \mathbb{1}\{\text{iteration } k \text{ is successful}\}$

Next, we define a random variable $U_k$ that measures whether the step size parameter is large or small. The threshold $\bar{\alpha}$ is a parameter that will depend on the algorithm being used; we will later give its concrete definition for the trust region and line search methods.

**Definition 4** (Large and Small Steps). *For a threshold $\bar{\alpha} > 0$, we define the random variable $U_k$ as:*

$$U_k = \begin{cases} 0, & \text{if } \max\{\mathcal{A}_k, \mathcal{A}_{k+1}\} \leq \bar{\alpha} \text{ (small step)} \\ 1, & \text{otherwise (large step)} \end{cases} \tag{17}$$

*Note that $U_k$ is measurable with respect to $\mathcal{F}_k$, since $\mathcal{A}_{k+1}$ is completely determined by $\mathcal{A}_k$ and $\Theta_k$.*

Since $\phi$ is non-convex, our goal is to find an approximately stationary point.

**Definition 5** (Stopping Time). *For some $\epsilon > 0$, let $T_\epsilon$ be the first time such that $\|\nabla\phi(X_{T_\epsilon})\| \leq \epsilon$. We refer to $T_\epsilon$ as the* stopping time *of the algorithm.*

Finally, the following measure of progress will be used as a potential function in our analysis.

**Definition 6** (Measure of Progress). *For each $k \geq 0$, define*

$$Z_k := \phi(X_k) - \phi^*,$$

*where $\phi^*$ is the infimum of the function $\phi$. Making progress means decreasing $Z_k$.*

# 4 Main Bound

We now present our main results. We first introduce some key assumptions and lemmas that will be used in the analysis. In Section 4.1, we present the high probability iteration complexity bound for the general algorithmic framework, assuming these assumptions hold. In Section 4.3, we prove that the assumptions hold for the trust region method, and specialize the general bound for the trust region method. Then, in Section 4.4, we prove that the assumptions hold for the line search method, and specialize the general bound for the line search method.

A key building block in our analysis is that iterations are true with some positive probability:

**Proposition 1** (Probabilistic Accuracy). *There exists a constant $p \in (0, 1]$ such that for all $k < T_\epsilon$:*

$$\mathbb{P}(I_k = 1 \mid \mathcal{F}_{k-1}) \geq p.$$

This condition ensures that, regardless of the algorithm's history, each iteration is true with probability at least $p > 0$. The value of $p$ depends on the specific oracle parameters. It is straightforward to see that:

- For SZO + SFO: $p \geq 1 - \delta_1 - \mathbb{P}(E_k + E_k^+ > 2\epsilon_f)$
- For CZO + SFO: $p \geq 1 - \delta_1 - 2\delta_0$

The following lemma will be useful in the high probability iteration result with the stochastic zeroth-order oracle.

**Lemma 1** (Fuk–Nagaev inequality). *Let $Y_1, \ldots, Y_t$ be independent random variables with $\mathbb{E}Y_i = 0$ and $\mathbb{E}|Y_i|^q < \infty$ for some $q \geq 2$. Define $S_t = \sum_{i=1}^{t} Y_i$. Then, for any $u > 0$,*

$$\mathbb{P}(S_t \geq u) \leq \exp\left(-\frac{u^2}{2V^2}\right) + \frac{C_q}{u^q},$$

11

*where*

$$V^2 := \frac{1}{4}(q+2)^2 e^q \sum_{i=1}^{t} \mathbb{E} Y_i^2, \qquad C_q := \left(1 + \frac{2}{q}\right)^q \sum_{i=1}^{t} \mathbb{E}|Y_i|^q.$$

*In particular, if $Y_i$ are i.i.d. with $\mathbb{E} Y_1^2 \leq \zeta_2$ and $\mathbb{E}|Y_1|^q \leq \zeta_q$, then*

$$\mathbb{P}(S_t \geq u) \ \leq \ \exp\left(-\frac{2\,u^2}{(q+2)^2 e^q \zeta_2\, t}\right) \ + \ \frac{\left(1 + \dfrac{2}{q}\right)^q \zeta_q\, t}{u^q}.$$

See, e.g., [FGL17, Corollary 2.5] for this one-sided form with explicit constants, and also [FN71, Pet95] for classical statements and proofs.

In our general algorithmic framework, we make two key assumptions. We will later show that these assumptions hold for both the trust-region and line search algorithms as instantiated in Section 2.

**Assumption 2.** *Any iteration $k < T_\epsilon$ that is small and true must be successful.*

**Assumption 3.** *There exists a function $h : \mathbb{R} \to \mathbb{R}$ such that for any large and successful iteration, the function value decreases by at least $h(\epsilon) - (2\epsilon_f + E_k + E_k^+)$.*

**Lemma 2.** *Under Assumption 1 and Assumption 2, for any $t \leq T_\epsilon$, Algorithm 1 satisfies*

$$\sum_{k=0}^{t-1} I_k \leq \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k \Theta_k + \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k\,(1 - \Theta_k) + \frac{1}{\lfloor m \rfloor + 1} t,$$

*where $m = -\frac{\ln \gamma_{\mathrm{inc}}}{\ln \gamma_{\mathrm{dec}}}$.*

*Proof.* Observe that

$$\begin{aligned}
\sum_{k=0}^{t-1} I_k &= \sum_{k=0}^{t-1} I_k U_k + \sum_{k=0}^{t-1} I_k(1 - U_k) \\
&\leq \sum_{k=0}^{t-1} U_k + \sum_{k=0}^{t-1} I_k(1 - U_k).
\end{aligned}$$

We first bound the second term $\sum_{k=0}^{t-1} I_k(1 - U_k)$, which is the total number of small and true iterations. By Assumption 2, if an iteration is small and true, it must be a successful iteration. This implies

$$\sum_{k=0}^{t-1} I_k(1 - U_k) \leq \sum_{k=0}^{t-1} \Theta_k(1 - U_k).$$

In other words, the number of steps that are small and true is upper bounded by the number of steps that are small and successful. As a result,

$$\sum_{k=0}^{t-1} I_k \leq \sum_{k=0}^{t-1} U_k + \sum_{k=0}^{t-1} I_k(1 - U_k) \leq \sum_{k=0}^{t-1} U_k + \sum_{k=0}^{t-1} \Theta_k(1 - U_k). \tag{18}$$

By the dynamics of how the the step size changes, we have the following inequality:

$$\lfloor m \rfloor \sum_{k=0}^{t-1} \Theta_k(1 - U_k) \leq \sum_{k=0}^{t-1} (1 - \Theta_k)(1 - U_k).$$

Using the inequality above, we have:

$$t - \sum_{k=0}^{t-1} U_k = \sum_{k=0}^{t-1}(1 - U_k) = \sum_{k=0}^{t-1} \Theta_k(1 - U_k) + \sum_{k=0}^{t-1}(1 - \Theta_k)(1 - U_k) \geq (\lfloor m \rfloor + 1) \sum_{k=0}^{t-1} \Theta_k(1 - U_k),$$

which implies

$$\sum_{k=0}^{t-1} \Theta_k(1 - U_k) \leq \frac{1}{\lfloor m \rfloor + 1} \left( t - \sum_{k=0}^{t-1} U_k \right).$$

Putting Equation (18) and the inequality above together, we have:

$$\sum_{k=0}^{t-1} I_k \leq \sum_{k=0}^{t-1} U_k + \sum_{k=0}^{t-1} \Theta_k(1 - U_k) \leq \sum_{k=0}^{t-1} U_k + \frac{1}{\lfloor m \rfloor + 1} \left( t - \sum_{k=0}^{t-1} U_k \right) = \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k + \frac{1}{\lfloor m \rfloor + 1} t.$$

We conclude that

$$\sum_{k=0}^{t-1} I_k \leq \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k + \frac{1}{\lfloor m \rfloor + 1} t = \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k \Theta_k + \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \sum_{k=0}^{t-1} U_k(1 - \Theta_k) + \frac{1}{\lfloor m \rfloor + 1} t.$$

$\square$

Without loss of generality, we assume that the algorithm starts with a step size parameter $\alpha_0 \geq \bar{\alpha}$. This is because if the initial step size parameter is less than $\bar{\alpha}$, we can simply redefine $\bar{\alpha}$ to be the initial step size parameter, which is still a constant.

**Lemma 3.** *Under Assumption 1, for any $t \leq T_\epsilon$, Algorithm 1 satisfies*

$$\sum_{k=0}^{t-1} U_k (1 - \Theta_k) \leq \lceil m \rceil \sum_{k=0}^{t-1} U_k \Theta_k + \left\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \right\rceil.$$

*Proof.* Recall that after any unsuccessful iteration, we decrease the step size by a factor of $\gamma_{\text{dec}}$. On the other hand, after any successful iteration, we increase the step size by a factor of $\gamma_{\text{inc}}$. As a result, if there are many large unsuccessful iterations, there must also be many large successful iterations. In particular, since we assume that the initial step size $\alpha_0 \geq \bar{\alpha}$, by the dynamics of the step size, we have that

$$\sum_{k=0}^{t-1} U_k (1 - \Theta_k) \leq \lceil m \rceil \cdot \sum_{k=0}^{t-1} U_k \Theta_k + \left\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \right\rceil.$$

In other words, the total number of large and unsuccessful iterations is within a multiplicative factor away from the total number of large and successful iterations, up to an additive constant. $\square$

We now prove an inequality that will be essential for our iteration complexity bound.

**Proposition 2.** *Under Assumptions 1 to 3, for any $t \leq T_\epsilon$, Algorithm 1 satisfies*

$$\sum_{k=0}^{t-1} I_k \leq \frac{\lfloor m \rfloor (\lceil m \rceil + 1)}{\lfloor m \rfloor + 1} \cdot \frac{Z_0 + \sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+)}{h(\epsilon)} + \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \left\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \right\rceil + \frac{1}{\lfloor m \rfloor + 1} t.$$

*In particular, if $m = -\frac{\ln \gamma_{\text{inc}}}{\ln \gamma_{\text{dec}}}$ is an integer, we have:*

$$\sum_{k=0}^{t-1} I_k \leq m \cdot \frac{Z_0 + \sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+)}{h(\epsilon)} + \frac{m}{m + 1} \left\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \right\rceil + \frac{1}{m + 1} t.$$

*Proof.* We first bound the total number of steps that are large and successful. Note that by the design of the algorithm, a step can increase the function value by at most $2\epsilon_f + E_k + E_k^+$ in each iteration due to the noise in the zeroth order oracle, and by Assumption 3 any large and successful step decreases the function value by at least $h(\epsilon) - (2\epsilon_f + E_k + E_k^+)$. As a result, there can be at most $\frac{Z_0 + \sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+)}{h(\epsilon)}$ iterations that are both large and successful.

Together with Lemma 2 and Lemma 3, we conclude that

$$\sum_{k=0}^{t-1} I_k \leq \frac{\lfloor m \rfloor (\lceil m \rceil + 1)}{\lfloor m \rfloor + 1} \cdot \frac{Z_0 + \sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+)}{h(\epsilon)} + \frac{\lfloor m \rfloor}{\lfloor m \rfloor + 1} \left\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \right\rceil + \frac{1}{\lfloor m \rfloor + 1} t.$$

$\square$

## 4.1 Iteration Complexity for General Framework

For the rest of the paper, we will assume for simplicity that $m$ is an integer. This is to avoid having to carry around $\lfloor m \rfloor$ and $\lceil m \rceil$. We can now obtain the high probability iteration complexity bound as follows for the general algorithmic framework.

**Theorem 1** (General high probability iteration complexity). *Suppose Assumptions 1 to 3 hold. Suppose $p > p_m$, where $p_m = \frac{1}{m+1} + \frac{2m\epsilon_f + m\mu}{h(\epsilon)}$, and $\mu := \sup_k \mathbb{E}[E_k + E_k^+]$ denote a uniform upper bound on the expected error across all iterations. Then for any $s \geq 0$, any $\hat{p} \in (p_m + \frac{s}{h(\epsilon)}, p)$, and any*

$$t > \frac{R}{\hat{p} - p_m - \frac{s}{h(\epsilon)}},$$

*the iteration complexity of Algorithm 1 satisfies*

$$\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p - \hat{p})^2}{2p^2} t\right) - \mathbb{P}(\bar{B}).$$

*Here, the event $B = \{\sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+) \leq t(2\epsilon_f + \mu + s)\}$, $\bar{B}$ is the complement of $B$, $R = \frac{mZ_0}{h(\epsilon)} + d$, and $d = \frac{m}{m+1}\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \rceil$.*

*Proof.* Define the event $A = \{\sum_{k=0}^{t-1} I_k \geq \hat{p}t\}$, and $B = \{\sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+) \leq t(2\epsilon_f + \mu + s)\}$. By Azuma–Hoeffding inequality applied to the submartingale $\sum_{k=0}^{t-1} I_k - pt$, we have for all $t \geq 1$, and any $\hat{p} \in [0, p)$,

$$\mathbb{P}(\bar{A}) = \mathbb{P}\left(\sum_{k=0}^{t-1} I_k < \hat{p}t\right) \leq \exp\left(-\frac{(p - \hat{p})^2}{2p^2} t\right). \tag{19}$$

Consider some positive integer $t > \frac{R}{\hat{p} - p_m - \frac{s}{h(\epsilon)}}$. We will show that $\mathbb{P}(T_\epsilon > t, A, B) = 0$ by contradiction. Suppose $T_\epsilon > t$. Then by Proposition 2, we have:

$$\sum_{k=0}^{t-1} I_k \leq m \cdot \frac{Z_0 + \sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+)}{h(\epsilon)} + d + \frac{1}{m+1} t.$$

Furthermore, under the events of $A$ and $B$, we have: $\sum_{k=0}^{t-1} I_k \geq \hat{p}t$ and $\sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+) \leq t(2\epsilon_f + \mu + s)$. Putting them together, we get $t \leq \frac{R}{\hat{p} - p_m - \frac{s}{h(\epsilon)}}$. This contradicts $t > \frac{R}{\hat{p} - p_m - \frac{s}{h(\epsilon)}}$. As a result, $\mathbb{P}(T_\epsilon > t, A, B) = 0$.

Using this fact, we obtain:

$$
\begin{aligned}
\mathbb{P}(T_\epsilon > t) &= \mathbb{P}(T_\epsilon > t, A) + \mathbb{P}(T_\epsilon > t, \bar{A}) \\
&\leq \mathbb{P}(T_\epsilon > t, A) + \mathbb{P}(\bar{A}) \\
&= \mathbb{P}(T_\epsilon > t, A, B) + \mathbb{P}(T_\epsilon > t, A, \bar{B}) + \mathbb{P}(\bar{A}) \\
&\leq \mathbb{P}(T_\epsilon > t, A, B) + \mathbb{P}(\bar{B}) + \mathbb{P}(\bar{A}) \\
&= 0 + \mathbb{P}(\bar{B}) + \mathbb{P}(\bar{A}).
\end{aligned}
$$

As a result, for any $t > \frac{R}{\bar{p} - p_m - \frac{s}{h(\epsilon)}}$, with probability at least $1 - \mathbb{P}(\bar{A}) - \mathbb{P}(\bar{B})$ we have $T_\epsilon \leq t$, i.e., the stopping time has been reached. Together with Equation (19), the result follows. □

**Discussion of Theorem 1** The general bound above applies to any algorithm that falls under the framework of Algorithm 1 as long as Assumptions 1 to 3 hold.

For the requirement $p > p_m$, let's first consider what this means in idealized setting with exact functions values, where $\epsilon_f = 0$ and $\mu = 0$. In this case, the condition simplifies to $p > \frac{1}{m+1}$. Recalling that $m = -\frac{\ln \gamma_{\text{inc}}}{\ln \gamma_{\text{dec}}}$, this is equivalent to requiring $p \ln \gamma_{\text{inc}} + (1-p) \ln \gamma_{\text{dec}} > 0$. This inequality provides a guideline for setting the step-size update parameters: the expected change in the step size parameter, $\alpha_k$, must ensure an "upward drift" and prevent it from shrinking to zero. A key challenge is that $p$ is usually unknown. In practice, one could simply set $\gamma_{\text{inc}}$ and $\gamma_{\text{dec}}$ so that the inequality is satisfied even for a conservative lower bound of $p$. That said, the tighter lower bound we have for $p$, the more efficent the algorithm will behave.

When function values are inexact, the requirement $p > p_m$ is related to a lower bound on the achievable target accuracy $\epsilon$; there exists a problem-dependent lower bound determined by the intrinsic oracle bias and noise levels (captured by $\epsilon_f$ and $\epsilon_g$). For trust region and line search algorithms, $\epsilon$ scales on the order of $\max\{\epsilon_g, \sqrt{\epsilon_f}\}$. Precise formulas are provided in the trust-region and line-search sections. We will see later that the condition $p > p_m$ implies the size of the convergence neighborhood of the algorithm is determined by the amount of noise in the oracles.

### 4.1.1 General Iteration Complexity with Stochastic Zeroth-Order Oracle

We now specialize the general bound to the case where the zeroth-order oracle is SZO. We first define the event $B_{\text{SZO}}(s, t)$ that is related to the accumulated function-value noise. For any $s > 0$ and horizon $t \geq 1$, let

$$
B_{\text{SZO}}(s, t) := \left\{ \sum_{k=0}^{t-1} (2\epsilon_f + E_k + E_k^+) \leq t(4\epsilon_f + s) \right\}.
$$

Since in the SZO setting $\mathbb{E}[E_k + E_k^+] \leq 2\epsilon_f$, we may use a concentration inequality to yield a tail bound $\delta_t(s)$ such that $\mathbb{P}(\bar{B}_{\text{SZO}}(s, t)) \leq \delta_t(s)$. The specific form of the tail bound is dependent on the distribution of the zeroth-order noise.

In general, if the zeroth-order noise is heavy-tailed with bounded $q$-th centered moment as in the SZO definition, then we can utilize Lemma 1 to obtain a tail bound. Define the centered variables

$$
Y_k := (E_k + E_k^+) - \mathbb{E}[E_k + E_k^+] \qquad (k = 0, \ldots, t-1),
$$

which satisfy $\mathbb{E}[Y_k] = 0$. Together with the SZO moment bound, we have

$$
\mathbb{E}|Y_k|^q \leq 2^{q-1} \left( \mathbb{E}|E_k - \mathbb{E}E_k|^q + \mathbb{E}|E_k^+ - \mathbb{E}E_k^+|^q \right) \leq 2^q \zeta_q, \quad \mathbb{E}Y_k^2 \leq 4\zeta_2.
$$

Let $S_t := \sum_{k=0}^{t-1} Y_k$. Since $\mathbb{E}[E_k + E_k^+] \leq 2\epsilon_f$, the event $\bar{B}_{\text{SZO}}(s, t)$ implies $S_t > ts$. Applying Lemma 1 yields

$$
\mathbb{P}(\bar{B}_{\text{SZO}}(s, t)) \leq \mathbb{P}(S_t \geq ts) \leq \exp\left( -\frac{s^2}{2(q+2)^2 e^q \zeta_2} t \right) + \frac{\left(1 + \frac{2}{q}\right)^q 2^q \zeta_q}{s^q \, t^{q-1}}.
$$

Applying Theorem 1 with event $B = B_{\text{SZO}}(s, t)$, we obtain the following theorem for the case where the zeroth-order oracle is SZO.

**Theorem 2** (High probability iteration complexity with SZO). *Suppose Assumptions 1 to 3 hold, and $p > p_m^{\text{SZO}}$, where $p_m^{\text{SZO}} = \frac{1}{m+1} + \frac{4m\epsilon_f}{h(\epsilon)}$. Then for any $s \geq 0$, any $\hat{p} \in (p_m^{\text{SZO}} + \frac{s}{h(\epsilon)}, p)$, and any*

$$t > \frac{R}{\hat{p} - p_m^{\text{SZO}} - \frac{s}{h(\epsilon)}},$$

*the iteration complexity of Algorithm 1 when using SZO satisfies*

$$\mathbb{P}(T_\epsilon \leq t) \;\geq\; 1 - \exp\left(-\frac{(p - \hat{p})^2}{2p^2}t\right) \;-\; \left[\exp\left(-\frac{s^2}{2(q+2)^2 e^q \zeta_2}t\right) + \frac{\left(1 + \frac{2}{q}\right)^q 2^q \zeta_q}{s^q t^{q-1}}\right].$$

Here, $R = \frac{mZ_0}{h(\epsilon)} + d$, and $d = \frac{m}{m+1}\lceil \frac{\ln \bar{\alpha} - \ln \alpha_0}{\ln \gamma_{\text{dec}}} \rceil$.

**Remark**  In the following special cases of the stochastic zeroth-order oracle, we can obtain tighter probability bounds:

- $q = 2$ (bounded variance; see definition as in Equation (2)). Using Chebyshev inequality in place of Lemma 1, one can show that

$$\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p - \hat{p})^2}{2p^2}t\right) - \frac{2\sigma^2}{s^2 t}.$$

- $q = \infty$ (sub-exponential tails; see definition as in Equation (3)). Using Bernstein inequality in place of Lemma 1, one can show that

$$\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p - \hat{p})^2}{2p^2}t\right) - \exp\left(-\min\left\{\frac{s^2}{8\nu^2}, \frac{s}{4b}\right\}t\right).$$

In other words, if the noise in the zeroth-order oracle follows a light-tailed distribution, then the failure probability decays exponentially in $t$.

### 4.1.2  General Iteration Complexity with Corrupted Zeroth-Order Oracle

We now specialize the general bound to the corrupted zeroth-order oracle (CZO) defined in Equation (8). Under CZO, we have $\mathbb{E}[E_k + E_k^+] \leq 2\epsilon_f + 2\delta_0\epsilon_c$. Define, for any $s > 0$ and horizon $t \geq 1$, the event

$$B_{\text{CZO}}(s, t) := \left\{\sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+) \leq t(4\epsilon_f + 2\delta_0\epsilon_c + s)\right\}.$$

Since any function value corruption lies in $[0, \epsilon_f + \epsilon_c]$ almost surely, Hoeffding's inequality yields the tail bound

$$\mathbb{P}(\bar{B}_{\text{CZO}}(s, t)) \;=\; \mathbb{P}\left(\sum_{k=0}^{t-1}(2\epsilon_f + E_k + E_k^+) > t(4\epsilon_f + 2\delta_0\epsilon_c + s)\right) \;\leq\; \exp\left(-\frac{s^2 t}{2(\epsilon_f + \epsilon_c)^2}\right).$$

Applying Theorem 1 with $B = B_{\text{CZO}}(s, t)$ gives the following result.

**Theorem 3** (High probability iteration complexity with CZO). *Suppose Assumptions 1 to 3 hold, and $p > p_m^{\text{CZO}}$, where $p_m^{\text{CZO}} := \frac{1}{m+1} + \frac{m\left(4\epsilon_f + 2\delta_0\epsilon_c\right)}{h(\epsilon)}$. Then for any $s \geq 0$, any $\hat{p} \in (p_m^{\text{CZO}} + \frac{s}{h(\epsilon)}, p)$, and any*

$$t > \frac{R}{\hat{p} - p_m^{\text{CZO}} - \frac{s}{h(\epsilon)}}$$

*the iteration complexity of Algorithm 1 when using CZO satisfies*

$$\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p - \hat{p})^2}{2p^2} t\right) - \exp\left(-\frac{s^2}{2(\epsilon_f + \epsilon_c)^2} t\right).$$

*Here, $R = \frac{mZ_0}{h(\epsilon)} + d$ and $d = \frac{m}{m+1}\left\lceil \frac{\ln\bar{\alpha} - \ln\alpha_0}{\ln\gamma_{\text{dec}}} \right\rceil$.*

## 4.2 Iteration Complexity in Expectation and Almost Sure Convergence

The high-probability bounds in Theorems 2 and 3 imply the following bounds in expectation.

**Corollary 1** (Iteration complexity in expectation). *Suppose Assumptions 1 to 3 hold, and let $p > p_m$ as in Theorem 1. Fix any $s \geq 0$, $\epsilon > 0$, $\hat{p} \in (p_m + \frac{s}{h(\epsilon)}, p)$, and define*

$$t_s = \left\lceil \frac{R}{\hat{p} - p_m - \frac{s}{h(\epsilon)}} \right\rceil, \qquad c = \frac{(p - \hat{p})^2}{2p^2},$$

*where $R$ is as in Theorem 1. If the zeroth-order oracle is SZO with bounded $q$-th centered moment with $q > 2$, or is CZO, then*

$$\mathbb{E}[T_\epsilon] = O(t_s).$$

*Proof.* By $\mathbb{E}[T_\epsilon] = \sum_{t=0}^{\infty} \mathbb{P}(T_\epsilon > t)$ and Theorem 1, we have

$$\mathbb{E}[T_\epsilon] = \sum_{t=0}^{\infty} \mathbb{P}(T_\epsilon > t) \leq \sum_{t=0}^{\infty} [e^{-ct} + \mathbb{P}(\bar{B}(s,t))].$$

The tail is summable (polynomial with exponent $q - 1 > 1$ for SZO; exponential for CZO). Hence, $\mathbb{E}[T_\epsilon] \leq t_s + O(1) = O(t_s)$. $\square$

For the general framework, using either SZO or CZO also ensures almost sure convergence.

**Theorem 4** (Almost sure convergence). *Suppose the zeroth-order oracle is SZO and the conditions of Theorem 2 hold, or the zeroth-order oracle is CZO and the conditions of Theorem 3 hold. The algorithm reaches an $\epsilon$-stationary iterate in finite time almost surely. That is,*

$$\mathbb{P}\left[\bigcap_{k=1}^{\infty} \{T_\epsilon > k\}\right] = 0.$$

*Proof.* We see that $\mathbb{P}\left[\cap_{k=1}^{\infty}(T_\epsilon > k)\right] = 0$ since the failure probability $\mathbb{P}(T_\epsilon > k)$ is going to 0 as $k \to \infty$, for both cases.

$\square$

## 4.3 Trust Region

In this section, we specialize our general analytical framework for the trust region method (Algorithm 2). To apply our main complexity theorems, we will define the concrete forms of the key analytical quantities from our general framework: the small-step threshold $\bar{\alpha}$, the progress measure $h(\epsilon)$, and the lower bound on $\epsilon$.

First, we define the threshold $\bar{\alpha}$ that distinguishes between "small" and "large" trust regions. This threshold is carefully chosen to be small enough to satisfy Assumption 2—ensuring that any "small and true" iteration is guaranteed to be a "successful" iteration. Its specific form is:

$$\bar{\alpha} := \min\left\{\frac{(1-\eta_1)\kappa_{\text{fcd}}(1-\eta) - 2\eta}{L + \kappa_{\text{H}} + 2\kappa + (1-\eta_1)\kappa_{\text{fcd}}\kappa}, \frac{1-\eta}{\kappa+\eta_2}\right\}\epsilon,$$

for some $\eta \in \left(0, \frac{(1-\eta_1)\kappa_{\text{fcd}}}{(1-\eta_1)\kappa_{\text{fcd}}+2}\right)$.

By the definition of a large step, if an iteration $k$ is large and successful, then its step size $\mathcal{A}_k$ must satisfy $\mathcal{A}_k \geq \frac{1}{\gamma_{\text{inc}}}\bar{\alpha}$. For the trust region method, we will show a large and successful iteration provides progress proportional to the square of the trust-region radius $\alpha_k^2$, up to the noise related term $2\epsilon_f + E_k + E_k^+$. This guarantees any large and successful iteration provides progress of at least $h(\epsilon) - \left(2\epsilon_f + E_k + E_k^+\right)$, where $h(\epsilon)$ is defined as:

$$h(\epsilon) = C_{\text{prog}}\left(\frac{\bar{\alpha}}{\gamma_{\text{inc}}}\right)^2 = C_{\text{prog}}\left(\frac{1}{\gamma_{\text{inc}}}\min\left\{\frac{(1-\eta_1)\kappa_{\text{fcd}}(1-\eta)-2\eta}{L+\kappa_{\text{H}}+2\kappa+(1-\eta_1)\kappa_{\text{fcd}}\kappa}, \frac{1-\eta}{\kappa+\eta_2}\right\}\right)^2\epsilon^2, \tag{20}$$

where $C_{\text{prog}} = \frac{1}{2}\eta_1\eta_2\kappa_{\text{fcd}}\min\left\{\frac{\eta_2}{\kappa_{\text{H}}},1\right\}$. Note that $h(\epsilon)$ is of the order $O(\epsilon^2)$.

Finally, due to the presence of bias in the oracle outputs, the achievable target accuracy $\epsilon$ is inherently limited. Specifically, $\epsilon$ is lower bounded by a certain threshold determined by the magnitude of the bias in zeroth- and first-order oracles. The specific lower bound on $\epsilon$ is given in Inequality 1.

**Inequality 1** (Lower bound on $\epsilon$ for Trust Region).

- **For SZO Oracle:**

$$\epsilon > \max\left\{\frac{\epsilon_g}{\eta}, \sqrt{\frac{4m\,\epsilon_f}{C_{prog}\left(p - \frac{1}{m+1}\right)}}\right\}.$$

- **For CZO Oracle:**

$$\epsilon > \max\left\{\frac{\epsilon_g}{\eta}, \sqrt{\frac{m(4\epsilon_f + 2\delta_0\epsilon_c)}{C_{prog}\left(p - \frac{1}{m+1}\right)}}\right\}.$$

for some $\eta \in (0, \frac{(1-\eta_1)\kappa_{\text{fcd}}}{2+(1-\eta_1)\kappa_{\text{fcd}}})$ and $p > \frac{1}{m+1}$.

With these specific definitions, we can now proceed to formally verify the two key assumptions of our general framework are satisfied by the trust region algorithm.

**Lemma 4** (Verification of Assumption 2 for Trust Region). *Consider the trust region algorithm (Algorithm 2). Under Assumption 1, for any iteration $k < T_\epsilon$ that is true and small, the iteration is also successful. That is, $\rho_k \geq \eta_1$ and $\|g_k\| \geq \eta_2\alpha_k$.*

*Proof.* First, we show that $\|g_k\| \geq \eta_2 \alpha_k$. Since iteration $k$ is true and $k < T_\epsilon$, we have $\|\nabla\phi(x_k)\| > \epsilon$ and $\|\nabla\phi(x_k) - g_k\| \leq \epsilon_g + \kappa\alpha_k$. By triangle inequality, we have

$$
\begin{aligned}
\|g_k\| &\geq \|\nabla\phi(x_k)\| - (\epsilon_g + \kappa\alpha_k) > \epsilon - \epsilon_g - \kappa\alpha_k \\
&\geq (1-\eta)\epsilon - \kappa\alpha_k \qquad\qquad\qquad\qquad\qquad\qquad \text{(by Inequality 1)} \\
&\geq \eta_2\alpha_k. \qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{(since } \alpha_k \leq \bar{\alpha})
\end{aligned}
$$

Next, we show $\rho_k \geq \eta_1$. Using Lemma 4.1 in [CBS24], we have that $\rho_k \geq \eta_1$ if the iteration is true and the following condition on $\alpha_k$ is met:

$$
\alpha_k \leq \frac{(1-\eta_1)\kappa_{\mathrm{fcd}}(\epsilon - \epsilon_g - \kappa\alpha_k) - 2\epsilon_g}{L + \kappa_{\mathrm{H}} + 2\kappa}.
$$

Using $\epsilon_g < \eta\epsilon$ from Inequality 1, this inequality is satisfied if

$$
\alpha_k \leq \frac{((1-\eta_1)\kappa_{\mathrm{fcd}}(1-\eta) - 2\eta)\,\epsilon}{L + \kappa_{\mathrm{H}} + 2\kappa + (1-\eta_1)\,\kappa_{\mathrm{fcd}}\kappa}.
$$

This inequality holds by the definition of $\bar{\alpha}$. Thus, both conditions for a successful iteration are met. $\qquad\square$

**Lemma 5** (Verification of Assumption 3 for Trust Region)**.** *Consider the trust region algorithm (Algorithm 2). Under Assumption 1, if an iteration $k < T_\epsilon$ is large and successful, then the function value decrease satisfies*

$$
\phi(x_k) - \phi(x_{k+1}) \geq h(\epsilon) - (2\epsilon_f + e_k + e_k^+).
$$

*Proof.* As established in Lemma 4.3 of [CBS24], any successful iteration provides a function value decrease of at least

$$
\phi(x_k) - \phi(x_{k+1}) \geq C_{\mathrm{prog}}\alpha_k^2 - (2\epsilon_f + e_k + e_k^+).
$$

Since the iteration is large, by definition we have $\alpha_k \geq \frac{1}{\gamma_{\mathrm{inc}}}\bar{\alpha}$. Therefore,

$$
C_{\mathrm{prog}}\alpha_k^2 \geq C_{\mathrm{prog}}\left(\frac{\bar{\alpha}}{\gamma_{\mathrm{inc}}}\right)^2.
$$

By definition, $h(\epsilon) = C_{\mathrm{prog}}\left(\frac{\bar{\alpha}}{\gamma_{\mathrm{inc}}}\right)^2$. Thus, the function value decrease is at least $h(\epsilon) - (2\epsilon_f + e_k + e_k^+)$, which verifies the assumption. $\qquad\square$

With both key assumptions verified, we can now present the main iteration complexity result for the trust-region method by specializing our general theorems. The high-probability iteration complexity is characterized by specializing Theorems 2 and 3 with the progress measure $h(\epsilon)$ defined for the trust-region method.

**Theorem 5** (High-Probability Iteration Complexity for Trust Region)**.** *Consider the trust region algorithm (Algorithm 2). Suppose Assumption 1 holds, the target accuracy $\epsilon$ satisfies Inequality 1, and $p > p_m$. Let $s > 0$ and $\hat{p} \in (p_m + \frac{s}{h(\epsilon)}, p)$. For any iteration count $t$ satisfying*

$$
t > \frac{R}{\hat{p} - p_m - s/h(\epsilon)},
$$

*the stopping time $T_\epsilon$ satisfies*

$$
\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p-\hat{p})^2}{2p^2}t\right) - \delta_t(s).
$$

*Here, $h(\epsilon)$ is as defined in Equation (20), $R = mZ_0/h(\epsilon) + d$ and $d = \frac{m}{m+1}\left\lceil\frac{\ln\bar{\alpha} - \ln\alpha_0}{\ln\gamma_{dec}}\right\rceil$. The threshold $p_m$ and the tail probability $\delta_t(s)$ depend on the zeroth-order oracle:*

- **For SZO Oracle:** $p_m = p_m^{\text{SZO}} = \frac{1}{m+1} + \frac{4m\epsilon_f}{h(\epsilon)}$, *and the tail probability is*

$$\delta_t(s) \;=\; \exp\!\left(-\frac{s^2}{2(q+2)^2 e^q \zeta_2}\,t\right) \;+\; \frac{\left(1+\dfrac{2}{q}\right)^q 2^q \zeta_q}{s^q\, t^{\,q-1}}$$

for bounded q-th centered moment with $q \geq 2$.

- **For CZO Oracle:** $p_m = p_m^{\text{CZO}} = \frac{1}{m+1} + \frac{m(4\epsilon_f + 2\delta_0\epsilon_c)}{h(\epsilon)}$, *and the tail probability is* $\delta_t(s) = \exp\!\left(-\frac{s^2}{2(\epsilon_f + \epsilon_c)^2}t\right)$.

*Since $h(\epsilon) = O(\epsilon^2)$, if the zeroth-order oracle is SZO with bounded q-th moment with $q > 2$, or is CZO, then the expected complexity is $\mathbb{E}[T_\epsilon] = O(\epsilon^{-2})$.*

Theorem 5 establishes a high-probability bound on the iteration complexity for the trust region method. This result not only provides insights into the expected complexity but also characterizes the tail behavior of the random variable of the stopping time $T_\epsilon$.

## 4.4  Line Search

We now turn our attention to the line search algorithm (Algorithm 3). The small-step threshold is defined as $\bar{\alpha} = \min\left\{\frac{1-\theta}{0.5L+\kappa}, \frac{2(1-2\eta-\theta(1-\eta))}{L(1-\eta)}\right\}$, for some $\eta \in \left(0, \frac{1-\theta}{2-\theta}\right)$. The progress measure $h(\epsilon)$ is given by:

$$h(\epsilon) = \frac{\theta\bar{\alpha}}{\gamma_{\text{inc}}\left(\max\left\{\frac{1}{\eta}, 1+\tau\right\}\right)^2}\epsilon^2. \tag{21}$$

In this algorithm, the achievable accuracy, denoted by $\epsilon$, is intrinsically linked to a rejection threshold, $\epsilon_{\text{rej}}$, which serves as the effective accuracy target. However, similar to the trust region setting, $\epsilon_{\text{rej}}$ cannot be set to an arbitrarily small value due to the inherent biases in the zeroth- and first-order oracles ($\epsilon_f$ and $\epsilon_g$). For the convergence analysis to hold, $\epsilon_{\text{rej}}$ must be sufficiently large with respect to the level of bias in the zeroth- and first-order oracles. The relationship between the oracle biases, $\epsilon_{\text{rej}}$, and the final stationarity guarantee $\epsilon$ is formalized below.

**Inequality 2** (Lower bound on $\epsilon_{\text{rej}}$ and resulting accuracy)**.**

- **For SZO Oracle:**
$$\epsilon_{\text{rej}} \geq \max\left\{\epsilon_g, \sqrt{\frac{4m\gamma_{inc}\epsilon_f}{\theta\bar{\alpha}(p-\frac{1}{m+1})}}\right\}.$$

- **For CZO Oracle:**
$$\epsilon_{\text{rej}} \geq \max\left\{\epsilon_g, \sqrt{\frac{2m\gamma_{inc}(2\epsilon_f + \delta_0\epsilon_c)}{\theta\bar{\alpha}\left(p-\frac{1}{m+1}\right)}}\right\}.$$

*Given such a $\epsilon_{\text{rej}}$, the algorithm is able to find an $\epsilon$-stationary point, where the achievable accuracy $\epsilon$ is defined as:*
$$\epsilon = \max\left\{\frac{1}{\eta}, 1+\tau\right\}\epsilon_{\text{rej}}.$$

**Lemma 6** (Verification of Assumption 2 for Line Search)**.** *Consider the line search algorithm (Algorithm 3). Under Assumption 1, for any iteration $k < T_\epsilon$ that is true and small, the iteration is also successful.*

*Proof.* An iteration is successful if it satisfies the sufficient decrease condition and is not rejected by the criterion $\|g_k\| \geq \epsilon_{\text{rej}}$. As established in [JSX24], any true and small step satisfies the sufficient decrease condition. We thus focus on showing that the criterion of $\|g_k\| \geq \epsilon_{\text{rej}}$ is also met.

Specifically, we need to show that for any true iteration $k < T_\epsilon$, we have $\|g_k\| \geq \epsilon_{\text{rej}}$. The derivation is as follows:

$$\|g_k\| \geq \min\left\{\frac{1}{1+\tau}, 1-\eta\right\} \|\nabla f(x_k)\| \quad \text{(by Proposition 1 (iii) in [JSX24], which uses } \epsilon \geq \frac{\epsilon_{\text{rej}}}{\eta} \geq \frac{\epsilon_g}{\eta})$$

$$\geq \min\left\{\frac{1}{1+\tau}, 1-\eta\right\} \epsilon \qquad\qquad \text{(since } k < T_\epsilon).$$

$$\geq \epsilon_{\text{rej}} \qquad\qquad \text{(by the lower bound on } \epsilon \text{ in Inequality 2)}$$

Here, the last inequality holds because $\eta < \frac{1-\theta}{2-\theta} < \frac{1}{2}$. This implies that $\frac{1}{\eta} > \frac{1}{1-\eta}$, which implies that $\epsilon = \max\left\{\frac{\epsilon_{\text{rej}}}{\eta}, (1+\tau)\epsilon_{\text{rej}}\right\} \geq \max\left\{\frac{\epsilon_{\text{rej}}}{1-\eta}, (1+\tau)\epsilon_{\text{rej}}\right\}$.

$\square$

**Lemma 7** (Verification of Assumption 3 for Line Search). *Consider the line search algorithm (Algorithm 3). Under Assumption 1, if an iteration $k < T_\epsilon$ is large and successful, then the function value decrease satisfies*

$$\phi(x_k) - \phi(x_{k+1}) \geq h(\epsilon) - (2\epsilon_f + e_k + e_k^+).$$

*Proof.* For a successful iteration, the sufficient decrease condition implies a lower bound on the objective function decrease of

$$\phi(x_k) - \phi(x_{k+1}) \geq \theta\alpha_k\|g_k\|^2 - (2\epsilon_f + e_k + e_k^+).$$

By definition, a successful iteration requires $\|g_k\| \geq \epsilon_{\text{rej}}$. For a large iteration, the step size is bounded below by $\alpha_k \geq \frac{\bar{\alpha}}{\gamma_{\text{inc}}}$. Substituting these bounds into the inequality gives:

$$\phi(x_k) - \phi(x_{k+1}) \geq \theta\left(\frac{\bar{\alpha}}{\gamma_{\text{inc}}}\right)\epsilon_{\text{rej}}^2 - (2\epsilon_f + e_k + e_k^+).$$

From the relationship between $\epsilon$ and $\epsilon_{\text{rej}}$ established in Inequality 2, we know that $\epsilon_{\text{rej}} = \epsilon/\max\{\frac{1}{\eta}, 1+\tau\}$. By substituting this into the expression for the progress, we recover the definition of $h(\epsilon)$:

$$\theta\left(\frac{\bar{\alpha}}{\gamma_{\text{inc}}}\right)\epsilon_{\text{rej}}^2 = \theta\frac{\bar{\alpha}}{\gamma_{\text{inc}}}\left(\frac{\epsilon}{\max\{\frac{1}{\eta}, 1+\tau\}}\right)^2 = h(\epsilon),$$

which completes the proof. $\square$

With these pieces in place, we can now state the main high-probability iteration complexity result for the line search method.

**Theorem 6** (High-Probability Iteration Complexity for Line Search). *Consider the line search algorithm (Algorithm 3) under Assumption 1, with the $\epsilon_{\text{rej}}$ parameter in Algorithm 3 satisfying Inequality 2. Let $p > p_m$, $s > 0$ and $\hat{p} \in (p_m + \frac{s}{h(\epsilon)}, p)$. For any iteration count $t$ satisfying*

$$t > \frac{R}{\hat{p} - p_m - s/h(\epsilon)},$$

*the algorithm finds an $\epsilon$-stationary point with probability*

$$\mathbb{P}(T_\epsilon \leq t) \geq 1 - \exp\left(-\frac{(p-\hat{p})^2}{2p^2}t\right) - \delta_t(s).$$

*Here, $h(\epsilon)$ is as defined in Equation (21), $R = mZ_0/h(\epsilon) + d$ and $d = \frac{m}{m+1}\left\lceil\frac{\ln\bar{\alpha} - \ln\alpha_0}{\ln\gamma_{dec}}\right\rceil$. The threshold $p_m$ and the tail probability $\delta_t(s)$ depend on the zeroth-order oracle:*

- **For SZO Oracle:** $p_m = p_m^{\text{SZO}} = \frac{1}{m+1} + \frac{4m\epsilon_f}{h(\epsilon)}$, and the tail probability is

$$\delta_t(s) \;=\; \exp\left(-\frac{s^2}{2(q+2)^2 e^q \zeta_2}\,t\right) \;+\; \frac{\left(1+\frac{2}{q}\right)^q 2^q \zeta_q}{s^q\,t^{\,q-1}}$$

  for bounded $q$-th centered moment with $q \geq 2$.

- **For CZO Oracle:** $p_m = p_m^{\text{CZO}} = \frac{1}{m+1} + \frac{m(4\epsilon_f + 2\delta_0 \epsilon_c)}{h(\epsilon)}$, and the tail probability is $\delta_t(s) = \exp\left(-\frac{s^2}{2(\epsilon_f + \epsilon_c)^2}\,t\right)$.

Since $h(\epsilon) = O(\epsilon^2)$, if the zeroth-order oracle is SZO with bounded $q$-th moment with $q > 2$, or is CZO, then the expected complexity is $\mathbb{E}[T_\epsilon] = O(\epsilon^{-2})$.

Theorem 6 establishes a high-probability bound on the iteration complexity for the line search method. This result not only provides insights into the expected complexity but also characterizes the tail behavior of the stopping time $T_\epsilon$. Notably, in the special case where $p > 1/2$ and the zeroth-order oracle noise is sub-exponential (a specific instance of our SZO oracle), our result recovers the high-probability complexity bound in [JSX24].

# 5  Concluding Remarks

A practical consideration when implementing the algorithms in our framework is the choice of the step-size update parameters, $\gamma_{\text{inc}}$ and $\gamma_{\text{dec}}$. Theoretically, these parameters must be set to ensure an "upward drift" of the stochastic process of the step size, which requires $p > \frac{\ln(\gamma_{\text{dec}})}{\ln(\gamma_{\text{dec}}/\gamma_{\text{inc}})}$, where $p$ is the probability of a true iteration. This can be achieved in practice by estimating a lower bound for $p$ and selecting the parameters accordingly.

It is also important to highlight a key implication of our results in the context of Expected Risk Minimization (ERM). In this setting, our objective function, $\phi(x)$, represents the true underlying risk, which is free from data anomalies. Our theoretical analysis is conducted with respect to this true function, while all sources of data corruption, noise, and error are modeled as inaccuracies in the oracle outputs. Therefore, by providing convergence guarantees to the true risk $\phi(x)$, our results also demonstrate generalization properties of the algorithm. Moreover, the convergence neighborhood it can achieve is cleanly dictated by the level of bias and noise in the oracles.

# References

[BCMS19]  Jose Blanchet, Coralia Cartis, Matt Menickelly, and Katya Scheinberg. Convergence rate analysis of a stochastic trust-region method via supermartingales. *INFORMS Journal on Optimization*, 1(2):92–119, 2019.

[BCS21]  Albert S. Berahas, Liyuan Cao, and Katya Scheinberg. Global convergence rate analysis of a generic line search algorithm with noise. *SIAM Journal on Optimization*, 31(2):1489–1518, 2021.

[BGMT20]  Stefania Bellavia, Gianmarco Gurioli, Benedetta Morini, and Philippe L Toint. A stochastic cubic regularisation method with inexact function evaluations and random derivatives for finite sum minimisation. In *Thirty-seventh International Conference on Machine Learning: ICML2020*, 2020.

[BGMT22]  Stefania Bellavia, Gianmarco Gurioli, Benedetta Morini, and Ph L Toint. Adaptive regularization for nonconvex optimization using inexact function values and randomly perturbed derivatives. *Journal of Complexity*, 68:101591, 2022.

[BSV14]     Afonso S Bandeira, Katya Scheinberg, and Luis Nunes Vicente. Convergence of trust-region methods based on probabilistic models. *SIAM Journal on Optimization*, 24(3):1238–1264, 2014.

[BXZ25]     Albert S Berahas, Miaolan Xie, and Baoyu Zhou. A sequential quadratic programming method with high-probability complexity bounds for nonlinear equality-constrained stochastic optimization. *SIAM Journal on Optimization*, 35(1):240–269, 2025.

[CBS24]     Liyuan Cao, Albert S Berahas, and Katya Scheinberg. First-and second-order high probability complexity bounds for trust-region methods with noisy oracles. *Mathematical Programming*, 207(1):55–106, 2024.

[CS17]      C. Cartis and K. Scheinberg. Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Mathematical Programming*, 169(2):337–375, 2017.

[FGL17]     Xiequan Fan, Ion Grama, and Quansheng Liu. Deviation inequalities for martingales with applications. *Journal of Mathematical Analysis and Applications*, 448(1):538–566, 2017.

[FN71]      D. H. Fuk and S. V. Nagaev. Probability inequalities for sums of independent random variables. *Theory of Probability & Its Applications*, 16(4):643–660, 1971.

[GRVZ18]    Serge Gratton, Clément W Royer, Luís N Vicente, and Zaikun Zhang. Complexity and global rates of trust-region methods based on probabilistic models. *IMA Journal of Numerical Analysis*, 38(3):1579–1597, 2018.

[JSX24]     Billy Jin, Katya Scheinberg, and Miaolan Xie. High probability complexity bounds for adaptive step search based on stochastic oracles. *SIAM Journal on Optimization*, 34(3):2411–2439, 2024.

[JSX25]     Billy Jin, Katya Scheinberg, and Miaolan Xie. Sample complexity analysis for adaptive optimization algorithms with stochastic oracles. *Mathematical Programming*, 209(1):651–679, 2025.

[MWX23]     Matt Menickelly, Stefan M Wild, and Miaolan Xie. A stochastic quasi-newton method in the absence of common random numbers. *arXiv preprint arXiv:2302.09128*, 2023.

[Pet95]     Valentin V. Petrov. *Limit Theorems of Probability Theory: Sequences of Independent Random Variables*. Oxford Studies in Probability. Oxford University Press, Oxford, 1995.

[PS20]      Courtney Paquette and Katya Scheinberg. A stochastic line search method with expected complexity analysis. *SIAM Journal on Optimization*, 30(1):349–376, 2020.

[SX23]      Katya Scheinberg and Miaolan Xie. Stochastic adaptive regularization method with cubics: A high probability complexity bound. In *2023 Winter Simulation Conference (WSC)*, pages 3520–3531. IEEE, 2023.

[Ver18]     Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.