

# Multifidelity preconditioning of the cross-entropy method for rare event simulation and failure probability estimation\*

Technical report, University of Wisconsin-Madison

Benjamin Peherstorfer<sup>†</sup>    Boris Kramer<sup>‡</sup>    Karen Willcox<sup>‡</sup>

April 11, 2017

Accurately estimating rare event probabilities with Monte Carlo can become costly if for each sample a computationally expensive high-fidelity model evaluation is necessary to approximate the system response. Variance reduction with importance sampling significantly reduces the number of required samples if a suitable biasing density is used. This work introduces a multifidelity approach that leverages a hierarchy of low-cost surrogate models to efficiently construct biasing densities for importance sampling. Our multifidelity approach is based on the cross-entropy method that derives a biasing density via an optimization problem. We approximate the solution of the optimization problem at each level of the surrogate-model hierarchy, reusing the densities found on the previous levels to precondition the optimization problem on the subsequent levels. With the preconditioning, an accurate approximation of the solution of the optimization problem at each level can be obtained from a few model evaluations only. In particular, at the highest level, only few evaluations of the computationally expensive high-fidelity model are necessary. Our numerical results demonstrate that our multifidelity approach achieves speedups of several orders of magnitude in a thermal and a reacting-flow example compared to the single-fidelity cross-entropy method that uses a single model alone.

**Keywords:** multifidelity, Monte Carlo, rare event simulation, failure probability estimation, surrogate models, reduced models, importance sampling, multilevel, cross-entropy method, variance reduction

---

\*This work was supported by the DARPA EQUiPS Program, Award UTA15-001067, Program Manager F. Fahroo, and by the AFOSR MURI on multi-information sources of multi-physics systems, Award Number FA9550-15-1-0038, Program Manager J.-L. Cambier. Several examples were computed on the computer clusters of the Munich Centre of Advanced Computing.

<sup>†</sup>Department of Mechanical Engineering and Wisconsin Institute for Discovery, University of Wisconsin-Madison (peherstorfer@wisc.edu).

<sup>‡</sup>Department of Aeronautics and Astronautics, Massachusetts Institute of Technology (bokramer@mit.edu, kwillcox@mit.edu).

# 1 Introduction

Rare event simulation with standard Monte Carlo typically requires a large number of samples to derive accurate estimates of rare event probabilities, which can become computationally infeasible if for each sample a computationally expensive high-fidelity model evaluation is necessary to simulate the system response. Importance sampling is a variance reduction strategy for Monte Carlo estimation that samples from a problem-dependent biasing distribution. The biasing distribution is chosen such that fewer samples are necessary to obtain an acceptable estimate of the rare event probability than with standard Monte Carlo. The bias introduced by the sampling from the biasing distribution is corrected by reweighing the samples in the importance sampling estimator [14, 28].

Traditionally, importance sampling consists of two steps. In the first step, the biasing distribution is constructed, and in the second step, samples are drawn from the biasing distribution and the estimate is derived [4, 33]. The challenge of rare event probability estimation with importance sampling is the construction of a suitable biasing distribution that leads to variance reduction. In principle, the optimal biasing distribution that leads to an estimator with zero variance is known, but evaluating the density of this zero-variance distribution requires the probability of the rare event, i.e., the quantity that is to be estimated. The cross-entropy (CE) method [31, 30, 32, 8] provides a practical way to approximate the zero-variance density. The CE method optimizes for a density that minimizes the Kullback-Leibler divergence from the zero-variance density in a set of feasible densities. Even though solving for a biasing density with the CE method typically requires fewer high-fidelity model evaluations than estimating the rare event probability with a standard Monte Carlo approach, the costs of the optimization problem in the CE method can still be significant if the high-fidelity model is expensive to evaluate.

In this paper, we introduce a multifidelity method that leverages a hierarchy of low-cost surrogate models to reduce the costs of constructing a CE-optimal biasing density. Examples of surrogate models include projection-based reduced models [29, 2], data-fit interpolation and regression models [15], machine-learning-based models such as support vector machines [7], and other simplified models [21]. At each level of the hierarchy, a CE-optimal density is derived with respect to the surrogate model corresponding to that current level. The optimization is initialized with the CE-optimal density of the previous level, which leads to preconditioned optimization problems that can be solved accurately with few model evaluations only. Thus, at higher levels, where the models are expensive to evaluate, only few model evaluations are necessary to obtain an accurate approximation of the solution of the optimization problem, which can lead to significant runtime savings while obtaining biasing densities that lead to a similar variance reduction as the biasing densities derived with the single-fidelity CE method that uses the high-fidelity model alone.

Multifidelity methods have been extensively used to speedup rare event probability estimation. We distinguish here between three categories of multifidelity methods for rare event simulation. First, there are two-fidelity methods that use a single surrogate model and combine it with the high-fidelity model. The work [19, 18, 20] introduces a two-fidelity approach that switches between a single surrogate model and the high-fidelity model depending on the error of the surrogate model, which can lead to unbiased estimators if the error of the surrogate model is known. In [6], an error estimator of a reduced-basis model is used to decide whether to evaluate the reduced or the high-fidelity model. In [10], the zero-variance biasing density is approximated with a Kriging model. Unbiasedness of the estimator is guaranteed by using the Kriging model as a proxy in the biasing density only. Similarly, in [24], unbiased estimators of

rare event probabilities are obtained by using a surrogate model to construct biasing densities and the high-fidelity model to derive the actual estimates.

Second, there are methods that use a multilevel hierarchy of models of a single type to speedup the estimation. Typically, these methods are developed for high-fidelity models that stem from partial differential equations (PDEs). The model hierarchy then often corresponds to different discretizations of the underlying PDE. There are extensions [11, 12, 13] of the multilevel Monte Carlo method [17, 16] for rare event probability estimation, which are based on variance reduction with control variates, instead of importance sampling. The subset method [1, 37] is another approach that has been extended to exploit a hierarchy of coarse-grid approximations in [35]. The subset method has also been combined with classification methods of machine learning such as support vector machines and neural networks in, e.g., [3, 23].

Third, multifidelity methods have been proposed that use multiple surrogate models of any type and combine them with the high-fidelity model. The method introduced in [22, 27] uses a control variate framework based on multiple surrogate models to accelerate the Monte Carlo estimation of statistics of the outputs of the high-fidelity model; however, the approach does not target rare event probability estimation. The multifidelity approach presented in [25] uses multiple surrogate models for speeding up the construction of biasing densities in importance sampling and guarantees unbiased estimators of the rare event probabilities by using the high-fidelity model to derive the estimate; however, the approach proposed in [25] has not been demonstrated on small probabilities below  $10^{-6}$ . The new multifidelity approach proposed in this paper also falls in this third category of multifidelity methods because we aim to exploit a hierarchy of surrogate models of any type. In contrast to [22, 27, 25], our approach explicitly targets rare event probabilities and we show that we successfully estimate probabilities as low as  $\approx 10^{-9}$ .

Section 2 of this paper provides preliminaries and the problem setup. Section 3 introduces our multifidelity preconditioner for the CE method, provides an error analysis, and summarizes our multifidelity approach in Algorithm 1. Section 4 demonstrates that our multifidelity approach achieves up to two orders of magnitude speedup compared to the single-fidelity CE method in a thermal and a reacting-flow example. Section 5 gives concluding remarks.

## 2 Importance sampling with the cross-entropy method

We first introduce the problem setup and then discuss importance sampling with the classical CE method that uses a single model alone.

### 2.1 Notation and problem setup

Let the value of the function  $f : \mathcal{D} \rightarrow \mathbb{R}$  denote the system response to an input  $\mathbf{z} \in \mathcal{D}$  with the input domain  $\mathcal{D} \subset \mathbb{R}^d$  in  $d \in \mathbb{N}$  dimensions. For example, if the system of interest is a cantilever beam, then the input could define material properties and the system response could be the displacement of the tip of the beam. Let  $Z : \Omega \rightarrow \mathcal{D}$  be a random variable with sample space  $\Omega$  and with probability density function  $p$ . We denote a realization of  $Z$  as  $\mathbf{z} \in \mathcal{D}$ .

Let  $t \in \mathbb{R}$  with  $t > 0$  be a rare event threshold and define the rare event probability as

$$P_t = \mathbb{P}_p[f \leq t] = \int_{\Omega} I_t(\mathbf{z})p(\mathbf{z})d\mathbf{z},$$

with the indicator function  $I_t : \mathcal{D} \rightarrow \{0, 1\}$  defined as

$$I_t(\mathbf{z}) = \begin{cases} 1, & f(\mathbf{z}) \leq t, \\ 0, & f(\mathbf{z}) > t \end{cases}.$$

Note that  $P_t = \mathbb{E}_p[I_t]$ , where  $\mathbb{E}_p$  denotes the expected value with respect to  $p$ . Let  $\text{Var}_p[I_t]$  be the variance of  $I_t$  with respect to  $p$  and assume  $\text{Var}_p[I_t] \in \mathbb{R}$  such that  $\text{Var}_p[I_t] = P_t(1 - P_t)$ . Let  $\rho \in (0, 1)$  and define the  $\rho$ -quantile of  $Z$  as  $\gamma \in \mathcal{D}$ , i.e.,

$$\mathbb{P}_p[Z \leq \gamma] = \rho.$$

Note that  $t$  is the  $P_t$ -quantile of  $f(Z)$ .

Consider now models  $f^{(\ell)} : \mathcal{D} \rightarrow \mathbb{R}$  of the system of interest, where  $\ell \in \mathbb{N}$  is a level parameter. For example, the models  $f^{(\ell)}$  can be derived via finite-element discretization from the governing equations of the system of interest; in this case, the level parameter  $\ell$  determines the mesh width. Note, however, that we will also consider models where the level parameter refers to more general concepts than mesh widths, e.g., the number of reduced basis vectors in reduced models and the number of data points in support vector regression machines. The costs of evaluating a model  $f^{(\ell)}$  are denoted as  $0 < w_\ell \in \mathbb{R}$ . For each model  $f^{(\ell)}$ , we define the indicator function  $I_t^{(\ell)} : \mathcal{D} \rightarrow \{0, 1\}$  as

$$I_t^{(\ell)}(\mathbf{z}) = \begin{cases} 1, & f^{(\ell)}(\mathbf{z}) \leq t, \\ 0, & f^{(\ell)}(\mathbf{z}) > t, \end{cases}$$

with the rare event threshold  $t$ . The rare event probability with respect to a model  $f^{(\ell)}$  is  $P_t^{(\ell)} = \mathbb{P}_p[f^{(\ell)} \leq t]$ . In the following, we choose a maximal level  $L \in \mathbb{N}$  such that the indicator function  $I_t^{(L)}$  leads to a rare event probability  $P_t^{(L)}$  that is a sufficiently accurate approximation of the rare event probability  $P_t$  of the system of interest for the current application at hand.

Let  $\hat{P}_t^{(L)}$  be an unbiased estimator of the rare event probability  $P_t^{(L)}$ . We assess the quality of an estimator with respect to its error and costs. We measure the error of an estimator  $\hat{P}_t^{(L)}$  with the squared coefficient of variation

$$e(\hat{P}_t^{(L)}) = \frac{\text{Var}_p[\hat{P}_t^{(L)}]}{\left(\mathbb{E}_p[\hat{P}_t^{(L)}]\right)^2}. \quad (1)$$

The costs  $c(\hat{P}_t^{(L)})$  are quantified with the costs of the model evaluations required by the estimator.

## 2.2 Standard Monte Carlo estimators

Let  $\mathbf{z}_1, \dots, \mathbf{z}_m \in \mathcal{D}$  be  $m \in \mathbb{N}$  realizations of the random variable  $Z$  and let

$$\hat{P}_t^{\text{MC}} = \frac{1}{m} \sum_{i=1}^m I_t^{(L)}(\mathbf{z}_i) \quad (2)$$

be the standard Monte Carlo estimator of  $P_t^{(L)}$ . The squared coefficient of variation  $e(\hat{P}_t^{\text{MC}})$  of  $\hat{P}_t^{\text{MC}}$  is

$$e(\hat{P}_t^{\text{MC}}) = \frac{1 - P_t^{(L)}}{mP_t^{(L)}}.$$

To achieve  $e(\hat{P}_t^{(L)}) < \epsilon$  for a given tolerance  $0 < \epsilon \in \mathbb{R}$ , the standard Monte Carlo estimator requires

$$m > \frac{1 - P_t^{(L)}}{\epsilon P_t^{(L)}}$$

evaluations of  $I_t^{(L)}$ , and thus  $m$  evaluations of  $f^{(L)}$ . Since  $m$  depends inverse proportionally on the rare event probability  $P_t^{(L)}$ , the number of evaluations  $m$  can become large for small rare event probabilities, which means that standard Monte Carlo estimators become computationally infeasible if the costs  $w_L$  of evaluating  $f^{(L)}$  are high.

## 2.3 Importance sampling with the cross-entropy method

Importance sampling estimators draw samples from a problem-dependent biasing distribution with the aim of reducing the variance compared to standard Monte Carlo estimators. This section discusses the CE method that iteratively constructs biasing distributions.

### 2.3.1 Importance sampling

Let  $\text{supp}(p) = \{z \in \mathcal{D} \mid p(z) > 0\}$  be the support of the density  $p$ . For a biasing density  $q$  with  $\text{supp}(p) \subseteq \text{supp}(q)$ , the importance sampling estimator  $\hat{P}_t^{\text{IS}}$  of  $P_t^{(L)}$  is

$$\hat{P}_t^{\text{IS}} = \frac{1}{m} \sum_{i=1}^m I_t^{(L)}(z'_i) \frac{p(z'_i)}{q(z'_i)},$$

with  $m$  realizations  $z'_1, \dots, z'_m$  of the random variable  $Z_q$  with the biasing density  $q$ . Because  $\text{supp}(p) \subseteq \text{supp}(q)$ , and if the variance of the importance sampling estimator

$$\text{Var}_q[\hat{P}_t^{\text{IS}}] = \frac{1}{m} \text{Var}_q \left[ I_t^{(L)} \frac{p}{q} \right]$$

is finite, then the importance sampling estimator  $\hat{P}_t^{\text{IS}}$  is an unbiased estimator of  $P_t^{(L)}$ . The biasing density  $q^*$  that minimizes the variance  $\text{Var}_q[\hat{P}_t^{\text{IS}}]$  is

$$q^*(z) = \frac{I_t^{(L)}(z)p(z)}{P_t^{(L)}},$$

which leads to an importance sampling estimator with variance 0. The density  $q^*$ , however, depends on  $P_t^{(L)}$ , which is the quantity we want to estimate.

### 2.3.2 CE-optimal biasing density

The CE method [31, 30, 32, 8] provides a practical way of approximating the zero-variance density  $q^*$ . Consider a set of parametrized densities  $\mathcal{Q} = \{q_{\mathbf{v}} \mid \mathbf{v} \in \mathcal{P}\}$ , where  $\mathbf{v} \in \mathcal{P}$  is a parameter in the set  $\mathcal{P}$ . For example,  $\mathcal{Q}$  could be the set of normal distributions with the parameter  $\mathbf{v}$  corresponding to the mean and covariance matrix. To ease the presentation, we assume in the following without loss of generality that the nominal density  $p$  of the random variable  $Z$  is in the set  $\mathcal{Q}$ . The CE method optimizes for a parameter  $\mathbf{v}_* \in \mathcal{P}$  such that the

corresponding density  $q_{\mathbf{v}_*} \in \mathcal{Q}$  minimizes the Kullback-Leibler divergence (also called the cross entropy) from the zero-variance density  $q^*$ . Transformations show that a solution of the problem

$$\mathbf{v}_* = \arg \max_{\mathbf{v} \in \mathcal{P}} \mathbb{E}_p[I_t^{(L)} \log(q_{\mathbf{v}})] \quad (3)$$

is a parameter  $\mathbf{v}_*$  that corresponds to a CE-optimal density  $q_{\mathbf{v}_*}$ , see, e.g., [8]. Solving the stochastic counterpart of (3)

$$\max_{\hat{\mathbf{v}} \in \mathcal{Q}} \frac{1}{m} \sum_{i=1}^m I_t^{(L)}(\mathbf{z}_i) \log(q_{\hat{\mathbf{v}}}(\mathbf{z}_i)), \quad (4)$$

with realizations  $\mathbf{z}_1, \dots, \mathbf{z}_m$  of  $Z$ , typically fails, because the stochastic counterpart (4) is affected by the rareness of  $I_t^{(L)}(Z)$ , just as the standard Monte Carlo estimator (2).

In [31, 30, 32], the CE method is proposed. The CE method iteratively derives an estimate  $\hat{\mathbf{v}}_*$  of the solution of the optimization problem (3). Our description of the CE method follows [8]. Consider the first iteration  $k = 1$ . In the first iteration, the CE method is initialized with the nominal random variable  $Z$  and the nominal density  $p \in \mathcal{Q}$ . Define the rare event threshold for the first iteration  $t_1 \in \mathbb{R}$  to be the  $\rho$ -quantile of the distribution of  $f^{(L)}(Z)$ , where  $\rho \in (0, 1)$  is a parameter that is typically in the range  $[10^{-2}, 10^{-1}]$ . Note that typically  $t_1 > t$ . Then, a solution  $\hat{\mathbf{v}}_1 \in \mathcal{Q}$  of the optimization problem

$$\max_{\hat{\mathbf{v}} \in \mathcal{Q}} \frac{1}{m} \sum_{i=1}^m I_{t_1}^{(L)}(\mathbf{z}_i) \log(q_{\hat{\mathbf{v}}}(\mathbf{z}_i)), \quad \mathbf{z}_1, \dots, \mathbf{z}_m \sim Z, \quad (5)$$

is obtained, where  $\mathbf{z}_1, \dots, \mathbf{z}_m \sim Z$  denotes that  $\mathbf{z}_1, \dots, \mathbf{z}_m$  are realizations of  $Z$ . The optimization problem (5) uses the indicator function  $I_{t_1}^{(L)}$  with threshold  $t_1$  instead of  $t$ , and therefore (5) avoids the rare event induced by the original threshold  $t$ . Note that the gradient of  $\log(q_{\hat{\mathbf{v}}})$  with respect to the parameter  $\hat{\mathbf{v}}$  is known analytically for certain sets of distributions  $\mathcal{Q}$ , see Section 3.5. In the second iteration  $k = 2$ , the threshold  $t_2$  is selected with respect to the distribution of  $f^{(L)}(Z_1)$ , where  $Z_1$  is the random variable with density  $q_{\hat{\mathbf{v}}_1}$  derived in the first iteration. To guarantee termination of the CE method, the threshold  $t_2$  is set to the minimum of the  $\rho$ -quantile of  $f^{(L)}(Z_1)$  and  $t_1 - \delta$ , where  $0 < \delta \in \mathbb{R}$  is a small constant [9, 8]. Then, the parameter  $\hat{\mathbf{v}}_2$  is derived from the optimization problem

$$\max_{\hat{\mathbf{v}} \in \mathcal{Q}} \frac{1}{m} \sum_{i=1}^m I_{t_2}^{(L)}(\mathbf{z}_i) \frac{p(\mathbf{z}_i)}{q_{\hat{\mathbf{v}}_1}(\mathbf{z}_i)} \log(q_{\hat{\mathbf{v}}}(\mathbf{z}_i)), \quad \mathbf{z}_1, \dots, \mathbf{z}_m \sim Z_1,$$

which is formulated with respect to the indicator function  $I_{t_2}^{(L)}$  that depends on the threshold  $t_2$ . This process is continued until step  $K \in \mathbb{N}$  where  $t_K \leq t$ , and where an estimate  $\hat{\mathbf{v}}_*$  of the CE-optimal parameter  $\mathbf{v}_*$  is obtained.

The CE method depends on two parameters: the quantile parameter  $\rho$  that determines the  $\rho$ -quantile for selecting the thresholds  $t_1, t_2, t_3, \dots, t_K$ , and the minimal-step-size parameter  $\delta$  that defines the minimal reduction of the threshold in each iteration. With the parameter  $\delta$ , the CE method terminates after at most

$$K = \frac{t_1 - t}{\delta} \quad (6)$$

iterations with an estimate  $\hat{\mathbf{v}}_*$  of  $\mathbf{v}_*$ . Thus,  $K$  is an upper bound on the number of CE iterations. Note that we sometimes use  $K$  but implicitly mean  $\lceil K \rceil$  to get an integer number. Details on the CE method, including a convergence analysis, are given in [8, 9].

In each iteration  $k = 1, \dots, K$  of the CE method, the model  $f^{(L)}$  is evaluated at  $m$  realizations. Therefore, a bound on the costs of deriving an importance-sampling estimate  $\hat{P}_t^{\text{IS}}$  of  $P_t^{(L)}$  with the CE method from  $m$  samples is

$$c(\hat{P}_t^{\text{IS}}) \leq K m w_L.$$

The squared coefficient of variation of the importance sampling estimator  $\hat{P}_t^{\text{IS}}$  depends on the variance reduction achieved by the biasing density

$$e(\hat{P}_t^{\text{IS}}) = \frac{\text{Var}_{\hat{v}_*} \left[ I_t^{(L)} \frac{p}{q_{\hat{v}_*}} \right]}{\left( \mathbb{E}_p \left[ \hat{P}_t^{\text{IS}} \right] \right)^2 m}. \quad (7)$$

Note that we abbreviate  $\text{Var}_{q_{\hat{v}_*}}$  with  $\text{Var}_{\hat{v}_*}$  in (7) and in the following.

### 3 A multifidelity preconditioner for the cross-entropy method

We propose a multifidelity-preconditioned CE (MFCE) method that exploits a hierarchy of models  $f^{(1)}, \dots, f^{(L)}$  to reduce the runtime of constructing a biasing density compared to the classical, single-fidelity CE method that uses  $f^{(L)}$  only. Section 3.1 introduces our MFCE approach. Section 3.2 and Section 3.3 formalize our MFCE method and present an analysis of the savings obtained with our MFCE method compared to the classical, single-fidelity CE method in terms of the bounds on the number of CE iterations. Section 3.4 summarizes the MFCE method in Algorithm 1, and Section 3.5 provides practical considerations.

#### 3.1 The MFCE method

Let  $p$  be the nominal density and let  $q \in \mathcal{Q}$  be a density in  $\mathcal{Q}$ . Consider the classical, single-fidelity CE method that uses model  $f^{(L)}$  alone, as discussed in Section 2.3.2. Let the CE method be initialized with density  $p$  and let  $t_p$  be the  $\rho$ -quantile of  $f^{(L)}(Z)$ . Then, the bound  $K_p = (t_p - t)/\delta$  on the number of CE iterations is obtained from (6). Similarly, the bound on the number of CE iterations is  $K_q = (t_q - t)/\delta$  if the CE method is initialized with  $q$ , where  $t_q$  is the  $\rho$ -quantile of  $f^{(L)}(Z_q)$  and where  $Z_q$  is a random variable with density  $q$ . This shows that the bound on the number of iterations of the CE method depends on the density with which the CE method is initialized. If  $t_q \leq t_p$ , then the bound on the number of iterations of the CE method initialized with  $q$  is lower or equal than the bound on the number of iterations of the CE method initialized with  $p$ .

We propose to exploit that the bound on the number of CE iterations can be reduced by a suitable choice of the density with which the CE method is initialized. Our MFCE method iterates through the levels  $\ell = 1, \dots, L$ . At level  $\ell = 1$ , our MFCE method constructs a biasing density  $q_{\hat{v}_*^{(1)}}$  with parameter  $\hat{v}_*^{(1)} \in \mathcal{P}$  with the classical CE method, initialized with the nominal density  $p$  and using model  $f^{(1)}$ . At level  $\ell = 2$ , our MFCE method uses the CE method to derive a density  $q_{\hat{v}_*^{(2)}}$  with model  $f^{(2)}$ ; however, the CE method on level  $\ell = 2$  is initialized with the density  $q_{\hat{v}_*^{(1)}}$  of the previous level, instead of the nominal density  $p$  as in the classical CE method. This hierarchical process is continued until level  $\ell = L$ , where density  $q_{\hat{v}_*^{(L-1)}}$  and model  $f^{(L)}$  are used to derive density  $q_{\hat{v}_*^{(L)}}$ .

### 3.2 Effect of the MFCE preconditioning

Consider now our MFCE method on level  $\ell$ , where we have obtained an estimate  $\hat{\mathbf{v}}_*^{(\ell)}$  and the corresponding biasing density  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$ . Using  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$  on level  $\ell + 1$  to obtain an estimate of the solution of the stochastic counterpart of

$$\max_{\mathbf{v} \in \mathcal{Q}} \mathbb{E}_{\hat{\mathbf{v}}_*^{(\ell)}} \left[ I_t^{(\ell+1)} \frac{p}{q_{\hat{\mathbf{v}}_*^{(\ell)}}} \log(q_{\mathbf{v}}) \right]$$

means that in the first CE iteration on level  $\ell + 1$  the threshold parameter  $t_1^{(\ell+1)}$  is set to the  $\rho$ -quantile of  $f^{(\ell+1)}(Z_\ell)$ , where  $Z_\ell$  is a random variable with density  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$ . In contrast, the classical CE method uses the  $\rho$ -quantile of  $f^{(\ell+1)}(Z)$  instead, where  $Z$  corresponds to the nominal density. If the  $\rho$ -quantile  $t_1^{(\ell+1)}$  of  $f^{(\ell+1)}(Z_\ell)$  is smaller than the  $\rho$ -quantile of  $f^{(\ell+1)}(Z)$ , then the bound on the iterations of the CE method on level  $\ell + 1$  is smaller if the CE method is initialized with  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$  than if the CE method is initialized with the nominal density  $p$ . The following proposition formalizes this notion.

**Proposition 1.** *Let  $\ell \in \mathbb{N}$  and let  $q_{\hat{\mathbf{v}}_*^{(\ell)}} \in \mathcal{Q}$  be the biasing density obtained with the CE method on level  $\ell$  of our MFCE approach. Let further  $t_1^{(\ell+1)}$  be the  $\rho$ -quantile of  $f^{(\ell+1)}(Z_\ell)$  with respect to the random variable  $Z_\ell$  with density  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$ . If*

$$\mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}} \left[ f^{(\ell+1)} \leq t_1^{(\ell+1)} \right] \geq \mathbb{P}_p \left[ f^{(\ell+1)} \leq t_1^{(\ell+1)} \right], \quad (8)$$

*then the bound on the number of iterations of the CE method initialized with  $q_{\hat{\mathbf{v}}_*^{(\ell)}}$  on level  $\ell + 1$  of our MFCE approach is less or equal to the bound of the classical CE method initialized with the nominal density  $p$ .*

*Proof.* Let  $t_p$  be the  $\rho$ -quantile of  $f^{(\ell+1)}(Z)$ , then we obtain with the monotonicity of the cumulative distribution function of  $f^{(\ell+1)}(Z)$  and (8) that  $t_p \geq t_1^{(\ell+1)}$ . The proposition follows with (6).  $\square$

### 3.3 Error analysis of multifidelity approach

Proposition 1 states under which condition the bound on the number of iterations of our MFCE approach is lower than the bound on the number of iterations of the classical CE method. In this section, we analyze which properties of the models  $f^{(1)}, \dots, f^{(L)}$  are required such that the condition (8) of Proposition 1 is met. The following analysis is based on the framework introduced in [11, 12]. We first make similar assumptions on the models as in [11, 12].

**Assumption 1.** *Let  $0 < \alpha < 1$  and let  $t$  be a threshold parameter. The models  $f^{(\ell)}$  satisfy*

$$|f^{(\ell)}(\mathbf{z}) - f^{(\ell+1)}(\mathbf{z})| \leq \alpha^\ell \text{ or } |f^{(\ell)}(\mathbf{z}) - f^{(\ell+1)}(\mathbf{z})| \leq |f^{(\ell)}(\mathbf{z}) - t|, \quad \mathbf{z} \in \mathcal{D},$$

*for  $\ell = 1, \dots, L - 1$ .*

**Assumption 2.** *Consider a density  $q \in \mathcal{Q}$  and the corresponding random variable  $Z_q$ . Let further  $F_q^{(\ell)}$  be the cumulative distribution function of  $f^{(\ell)}(Z_q)$  for  $\ell = 1, \dots, L$ . The cumulative distribution function  $F_q^{(\ell)}$  is Lipschitz continuous with Lipschitz constant  $C_q^{(\ell)}$ . Furthermore,*

there exists a constant  $C \in \mathbb{R}$  that bounds  $C_q^{(\ell)} \leq C$  for all  $\ell = 1, \dots, L$  and  $q \in \mathcal{Q}$ . We therefore have

$$|F_q^{(\ell)}(t) - F_q^{(\ell)}(t')| \leq C|t - t'|,$$

where  $t, t' \in \mathbb{R}$ .

Under Assumption 1 and Assumption 2 we obtain the following proposition.

**Proposition 2.** *Let  $\ell \in \{1, \dots, L-1\}$  and let  $\hat{\mathbf{v}}_*^{(\ell)}$  be the parameter of the biasing density estimated on level  $\ell$ . Let further  $t_1^{(\ell+1)}$  be the  $\rho$ -quantile with  $\mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell+1)} \leq t_1^{(\ell+1)}] = \rho$  and let  $\mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq t_1^{(\ell+1)}] \geq \mathbb{P}_p[f^{(\ell)} \leq t_1^{(\ell+1)}]$ , then*

$$\mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell+1)} \leq t_1^{(\ell+1)}] \geq \mathbb{P}_p[f^{(\ell+1)} \leq t_1^{(\ell+1)}] - 8C\alpha^\ell, \quad (9)$$

where  $\alpha$  and  $C$  are the constants of Assumptions 1–2.

Before we prove Proposition 2, we first show Lemma 1 and Lemma 2.

**Lemma 1.** *Define  $\gamma = t_1^{(\ell+1)}$  and the set  $\mathcal{B} = \{\mathbf{z} \in \mathcal{D} : |f^{(\ell)}(\mathbf{z}) - \gamma| \leq \alpha^\ell\}$ . With Assumption 1 follows that  $I_\gamma^{(\ell)}(\mathbf{z}) = I_\gamma^{(\ell+1)}(\mathbf{z})$  for  $\mathbf{z} \in \mathcal{D} \setminus \mathcal{B}$ .*

*Proof.* This proof follows similar arguments as the proof of Lemma 3.3 in [12]. We show that  $f^{(\ell)}(\mathbf{z}) \leq \gamma \iff f^{(\ell+1)}(\mathbf{z}) \leq \gamma$  holds for  $\mathbf{z} \in \mathcal{D} \setminus \mathcal{B}$ , which is equivalent to  $I_\gamma^{(\ell)}(\mathbf{z}) = I_\gamma^{(\ell+1)}(\mathbf{z})$  for  $\mathbf{z} \in \mathcal{D} \setminus \mathcal{B}$ . Consider first  $f^{(\ell)}(\mathbf{z}) \leq \gamma \Rightarrow f^{(\ell+1)}(\mathbf{z}) \leq \gamma$ . We obtain

$$0 \leq \gamma - f^{(\ell)}(\mathbf{z}) \leq \gamma - f^{(\ell+1)}(\mathbf{z}) + |f^{(\ell+1)}(\mathbf{z}) - f^{(\ell)}(\mathbf{z})| \leq \gamma - f^{(\ell+1)}(\mathbf{z}) + |f^{(\ell)}(\mathbf{z}) - \gamma|,$$

because for all  $\mathbf{z} \in \mathcal{D} \setminus \mathcal{B}$  we have  $|f^{(\ell)}(\mathbf{z}) - \gamma| > \alpha^\ell$  by definition of  $\mathcal{B}$  and therefore  $|f^{(\ell+1)}(\mathbf{z}) - f^{(\ell)}(\mathbf{z})| \leq |f^{(\ell)}(\mathbf{z}) - \gamma|$  because of Assumption 1. This shows

$$\begin{aligned} 0 \leq \gamma - f^{(\ell)}(\mathbf{z}) &= |\gamma - f^{(\ell)}(\mathbf{z})| \leq \gamma - f^{(\ell+1)}(\mathbf{z}) + |f^{(\ell)}(\mathbf{z}) - \gamma| \\ 0 &\leq \gamma - f^{(\ell+1)}(\mathbf{z}) \end{aligned}$$

and therefore  $f^{(\ell)}(\mathbf{z}) \leq \gamma \Rightarrow f^{(\ell+1)}(\mathbf{z}) \leq \gamma$ . For  $f^{(\ell+1)}(\mathbf{z}) \leq \gamma \Rightarrow f^{(\ell)}(\mathbf{z}) \leq \gamma$ , we show  $f^{(\ell)}(\mathbf{z}) \geq \gamma \Rightarrow f^{(\ell+1)}(\mathbf{z}) \geq \gamma$  with similar arguments. We obtain

$$0 \leq f^{(\ell)}(\mathbf{z}) - \gamma \leq |f^{(\ell)}(\mathbf{z}) - f^{(\ell+1)}(\mathbf{z})| + f^{(\ell+1)}(\mathbf{z}) - \gamma \leq |f^{(\ell)}(\mathbf{z}) - \gamma| + f^{(\ell+1)}(\mathbf{z}) - \gamma,$$

and therefore  $0 \leq f^{(\ell+1)}(\mathbf{z}) - \gamma$ .  $\square$

**Lemma 2.** *With Assumptions 1–2, we obtain*

$$P_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq \gamma] \leq P_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell+1)} \leq \gamma] + 4C\alpha^\ell, \quad (10)$$

where  $\gamma = t_1^{(\ell+1)}$  as in Lemma 1.

*Proof.* Let  $\mathcal{B}$  be the set defined in Lemma 1. For  $\mathbf{z} \in \mathcal{B}$ , we obtain with Assumption 1 that  $|f^{(\ell+1)}(\mathbf{z}) - f^{(\ell)}(\mathbf{z})| \leq \alpha^\ell$  and  $|f^{(\ell+1)}(\mathbf{z}) - \gamma| \leq 2\alpha^\ell$ . Consider now  $\mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq \gamma]$ , which we write as

$$\begin{aligned} \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq \gamma] &= \int_{\mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} + \int_{\mathcal{D} \setminus \mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} + \int_{\mathcal{D} \setminus \mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} \end{aligned} \quad (11)$$

$$\leq \int_{\mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} + \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell+1)} \leq \gamma], \quad (12)$$

where we obtain equality in (11) because  $I_\gamma^{(\ell)}(\mathbf{z}) = I_\gamma^{(\ell+1)}(\mathbf{z})$  for  $\mathbf{z} \in \mathcal{D} \setminus \mathcal{B}$ , see Lemma 1, and  $\leq$  in (12) because  $I_\gamma^{(\ell+1)}$  is non-negative. Consider now the first term in (12), for which we obtain

$$\begin{aligned} \int_{\mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} &\leq \int_{\mathcal{B}} q_{\hat{\mathbf{v}}_*^{(\ell)}}(\mathbf{z}) d\mathbf{z} \\ &\leq \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}} \left[ |f^{(\ell+1)} - \gamma| \leq 2\alpha^\ell \right] \\ &= F_{\hat{\mathbf{v}}_*^{(\ell)}}^{(\ell+1)}(\gamma - 2\alpha^\ell) - F_{\hat{\mathbf{v}}_*^{(\ell)}}^{(\ell+1)}(\gamma + 2\alpha^\ell) \\ &\leq 4C\alpha^\ell, \end{aligned} \tag{13}$$

where we used Assumption 2 in (13). Combining the bound in (13) with (12) leads to (10).  $\square$

We now state the proof of Proposition 2.

*Proof of Proposition 2.* Let  $\gamma$  and  $\mathcal{B}$  be defined as in Lemma 1. Consider  $\mathbb{P}_p[f^{(\ell+1)} \leq \gamma]$ , which we write as

$$\begin{aligned} \mathbb{P}_p[f^{(\ell+1)} \leq \gamma] &= \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} + \int_{\mathcal{D} \setminus \mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} + \int_{\mathcal{D} \setminus \mathcal{B}} I_\gamma^{(\ell)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \end{aligned} \tag{14}$$

$$\begin{aligned} &\leq \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} + \mathbb{P}_p[f^{(\ell)} \leq \gamma] \\ &\leq \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} + \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq \gamma] \end{aligned} \tag{15}$$

$$\leq \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} + 4C\alpha^\ell + \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell+1)} \leq \gamma], \tag{16}$$

where we obtain equality in (14) because of Lemma 1, and  $\leq$  in (15) because  $\mathbb{P}_p[f^{(\ell)} \leq \gamma] \leq \mathbb{P}_{\hat{\mathbf{v}}_*^{(\ell)}}[f^{(\ell)} \leq \gamma]$  as assumed in the statement of Proposition 2. The inequality  $\leq$  in (16) is obtained because of Lemma 2. Consider now the first term in (16), for which we obtain

$$\begin{aligned} \int_{\mathcal{B}} I_\gamma^{(\ell+1)}(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} &\leq \int_{\mathcal{B}} p(\mathbf{z}) d\mathbf{z} \\ &\leq \mathbb{P}_p \left[ |f^{(\ell+1)} - \gamma| \leq 2\alpha^\ell \right] \end{aligned} \tag{17}$$

$$\begin{aligned} &= F_p^{(\ell+1)}(\gamma - 2\alpha^\ell) - F_p^{(\ell+1)}(\gamma + 2\alpha^\ell) \\ &\leq 4C\alpha^\ell, \end{aligned} \tag{18}$$

where we used  $|f^{(\ell+1)} - \gamma| \leq 2\alpha^\ell$  for  $\mathbf{z} \in \mathcal{B}$  in (17) as in the proof of Lemma 2 and Assumption 2 in (18). Combining the bound in (18) with (16) leads to (9).  $\square$

Proposition 2 shows that with Assumptions 1–2 we obtain condition (8) of Proposition 1 up to the factor  $8C\alpha^\ell$ . Note that the factor  $8C\alpha^\ell$  decays with the level  $\ell$ , because we have  $|\alpha| < 1$ . With the parameter  $\hat{\mathbf{v}}_*^{(L)}$  derived at level  $L$ , we define our MFCE estimator as

$$\hat{P}_t^{\text{MFCE}} = \frac{1}{m} \sum_{i=1}^m I_t^{(L)}(\mathbf{z}_i) \frac{p(\mathbf{z}_i)}{q_{\hat{\mathbf{v}}_*^{(L)}}(\mathbf{z}_i)},$$

---

**Algorithm 1** Cross-entropy method with multifidelity preconditioning
 

---

```

1: procedure MFCE( $f^{(1)}, \dots, f^{(L)}, t, p, m, \rho, \delta$ )
2:   Initialize  $\hat{\mathbf{v}}_*^{(0)} \in \mathcal{P}$  such that  $q_{\hat{\mathbf{v}}_*^{(0)}} = p$ ; redefine  $\mathcal{Q} = \mathcal{Q} \cup \{p\}$  if necessary
3:   for  $\ell = 1, \dots, L$  do
4:     Initialize  $\hat{\mathbf{v}}_0^{(\ell)} = \hat{\mathbf{v}}_*^{(\ell-1)}$  and set  $k = 0$ 
5:     while 1 do
6:       Draw realizations  $\mathbf{z}_1, \dots, \mathbf{z}_m$  from random variable with density  $q_{\hat{\mathbf{v}}_k^{(\ell)}}$ 
7:       Compute model outputs  $\mathcal{G}_k^{(\ell)} = \{f^{(\ell)}(\mathbf{z}_1), \dots, f^{(\ell)}(\mathbf{z}_m)\}$ 
8:       Estimate  $\rho$ -quantile  $\hat{\gamma}_k^{(\ell)}$  from  $\mathcal{G}_k^{(\ell)}$ 
9:       if  $k == 0$  then  $t_0^{(\ell)} = \hat{\gamma}_k^{(\ell)} + \delta$ 
10:      end if
11:      Set  $t_{k+1}^{(\ell)} = \max\{t, \min\{t_{k-1}^{(\ell)} - \delta, \hat{\gamma}_k^{(\ell)}\}\}$ 
12:      Estimate parameter  $\hat{\mathbf{v}}_{k+1}^{(\ell)} \in \mathcal{P}$  by solving
          
$$\max_{\hat{\mathbf{v}} \in \mathcal{P}} \frac{1}{m} \sum_{i=1}^m I_{t_{k+1}^{(\ell)}}^{(\ell)}(\mathbf{z}_i) \frac{p(\mathbf{z}_i)}{q_{\hat{\mathbf{v}}_k^{(\ell)}}(\mathbf{z}_i)} \log(q_{\hat{\mathbf{v}}}(\mathbf{z}_i)) \quad (19)$$

13:      if  $t_{k+1}^{(\ell)} == t$  then break
14:      end if
15:      Set  $k = k + 1$ 
16:    end while
17:    Set  $\hat{\mathbf{v}}_*^{(\ell)} = \hat{\mathbf{v}}_k^{(\ell)}$  and  $\mathcal{G}_*^{(\ell)} = \mathcal{G}_k^{(\ell)}$ 
18:  end for
19:  Return estimate  $\hat{P}_t^{\text{MFCE}}$  with  $\mathcal{G}_*^{(L)} = \{f^{(L)}(\mathbf{z}_1), \dots, f^{(L)}(\mathbf{z}_m)\}$  and  $\hat{\mathbf{v}}_*^{(L)}$  as

```

$$\hat{P}_t^{\text{MFCE}} = \frac{1}{m} \sum_{i=1}^m I_t^{(L)}(\mathbf{z}_i) \frac{p(\mathbf{z}_i)}{q_{\hat{\mathbf{v}}_*^{(L)}}(\mathbf{z}_i)}$$

20: **end procedure**

---

where  $\mathbf{z}_1, \dots, \mathbf{z}_m$  are realizations of the random variable with density  $q_{\hat{\mathbf{v}}_*^{(L)}}$ . The MFCE estimator is unbiased with respect to the rare event probability  $P_t^{(L)}$  if the biasing density defined by the parameter  $\hat{\mathbf{v}}_*^{(L)}$  has a support that is a superset of the support of the nominal density  $p$  and the variance of the MFCE estimator is finite. The squared coefficient of variation of the MFCE estimator is

$$e(\hat{P}_t^{\text{MFCE}}) = \frac{\text{Var}_{\hat{\mathbf{v}}_*^{(L)}} \left[ I_t^{(L)} \frac{p}{q_{\hat{\mathbf{v}}_*^{(L)}}} \right]}{\left( \mathbb{E}[\hat{P}_t^{\text{MFCE}}] \right)^2 m}$$

and depends on the parameter  $\hat{\mathbf{v}}_*^{(L)}$ .

### 3.4 Computational procedure

Algorithm 1 summarizes our MFCE method. Inputs are the models  $f^{(1)}, \dots, f^{(L)}$ , the threshold  $t$ , the nominal density  $p$ , the number of samples  $m \in \mathbb{N}$ , the quantile parameter  $\rho$ , and the

minimal step size  $\delta$ . Note that the parameters  $\rho$  and  $\delta$  are the same parameters as in the classical, single-fidelity CE method, see Section 2.3.2. The `for` loop in line 3 iterates through the levels  $\ell = 1, \dots, L$ . At each level  $\ell$ , the density with parameter  $\hat{\mathbf{v}}_*^{(\ell)}$  with respect to model  $f^{(\ell)}$  is derived. At iteration  $k = 0, 1, 2, \dots$  of the `while` loop in line 5, realizations  $\mathbf{z}_1, \dots, \mathbf{z}_m$  are drawn from the random variable with the density  $q_{\hat{\mathbf{v}}_k^{(\ell)}}$  of the current iteration, the model  $f^{(\ell)}$  is evaluated at the realizations, and the  $\rho$ -quantile is estimated from the model outputs. In line 9, the threshold  $t_0^{(\ell)}$  is set to the  $\rho$ -quantile estimate  $\hat{\gamma}_k^{(\ell)} + \delta$  in the first iteration of the loop. The threshold  $t_{k+1}^{(\ell)}$  is selected in line 11, where the  $\max\{\}$  operation guarantees that the threshold  $t_{k+1}^{(\ell)}$  is less or equal to  $t$  and the  $\min\{\}$  operation guarantees that the threshold parameter is reduced by at least  $\delta$  in each iteration of the `while` loop, except in the first and the last iteration. In line 12, the parameter  $\hat{\mathbf{v}}_{k+1}^{(\ell)}$  is estimated. Line 13 exits the `while` loop if  $t_{k+1}^{(\ell)}$  is equal to the threshold  $t$ . In line 17, the estimate  $\hat{\mathbf{v}}_k^{(\ell)}$  and the model outputs  $\mathcal{G}_k^{(\ell)}$  of the last iteration of the `while` loop are stored, and the `for` loop starts a new iteration for the next level  $\ell + 1$ . After the `for` loop iterated through all levels  $\ell = 1, \dots, L$ , the MFCE estimate  $\hat{P}_t^{\text{MFCE}}$  is returned using the density  $p_{\hat{\mathbf{v}}_*^{(L)}}$  and the model outputs in  $\mathcal{G}_*^{(L)}$ .

Typically, the computationally most expensive step of Algorithm 1 is the computation of the model outputs on line 7, which scales linearly with the number of realizations  $m$ . Solving the optimization problem (19) on line 12 typically incurs small costs if the gradients of the objective can be computed analytically. If the gradients of the objective have to be approximated numerically, then solving the optimization on line 12 can become expensive.

### 3.5 Practical considerations

In the following, we consider the set  $\mathcal{Q}$  of Gaussian distributions of dimension  $d$ , for which we derive the gradient of the objective of the stochastic counterparts (19) analytically. Let the parameter  $\mathbf{v} \in \mathcal{P}$  describe the mean  $\boldsymbol{\mu} \in \mathbb{R}^d$  and the covariance matrix  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  of a  $d$ -dimensional Gaussian distribution. The corresponding density is

$$q_{\mathbf{v}}(\mathbf{z}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})\right),$$

where  $|2\pi\boldsymbol{\Sigma}|$  denotes the determinant of the matrix  $2\pi\boldsymbol{\Sigma}$ . We obtain the gradient of  $\log(q_{\mathbf{v}}(\mathbf{z}))$  with respect to  $\boldsymbol{\mu}$  as

$$\nabla_{\boldsymbol{\mu}} \log(q_{\mathbf{v}}(\mathbf{z})) = \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu}),$$

and the gradient with respect to  $\boldsymbol{\Sigma}$  as

$$\nabla_{\boldsymbol{\Sigma}} \log(q_{\mathbf{v}}(\mathbf{z})) = -\frac{1}{2} \left( \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})^T(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \right).$$

We use the gradients of  $\log(q_{\mathbf{v}}(\mathbf{z}))$  and plug them into the gradient of the objective of the stochastic counterpart (19). Then, setting the gradient of the objective of (19) to zero leads to a system of equations that is linear in  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Therefore, in the case where  $\mathcal{Q}$  is the set of Gaussian distributions of dimension  $d$ , solving the optimization problem (19) means solving a system of linear equations.

## 4 Numerical results

We demonstrate the efficiency of our MFCE method on a heat transfer and a reacting flow example. In all of the following experiments, the quantile parameter is set to  $\rho = 0.1$  and

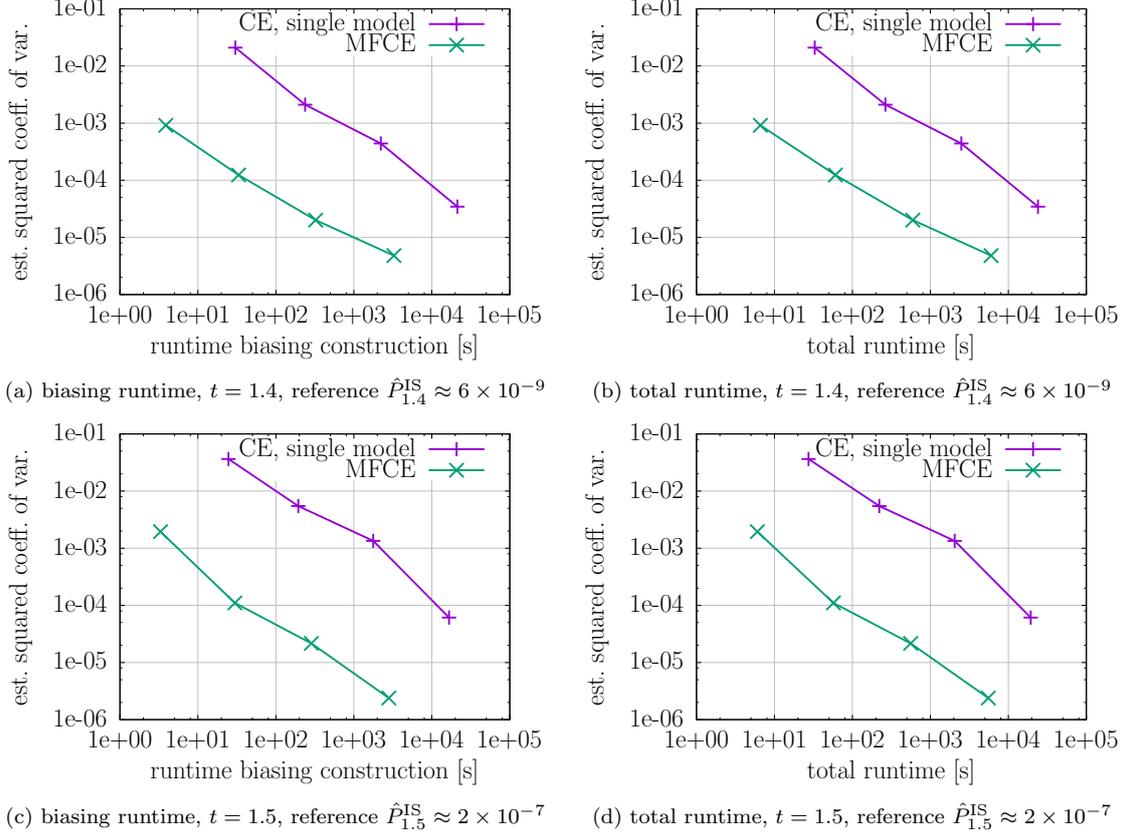


Figure 1: Heat transfer: Our MFCE approach achieves up to two orders of magnitude speedup compared to using the single-fidelity CE method with the model  $f^{(L)}$  for  $t \in \{1.4, 1.5\}$ .

the minimal step size is  $\delta = 10^{-2}$ , which are similar to the parameters chosen in, e.g., [8]. Furthermore, we set  $\mathcal{Q}$  to the set of Gaussian distributions of the respective dimensions. We constrain the optimization problem (19) in Algorithm 1 to covariance matrices  $\Sigma$  with a minimal absolute value of  $10^{-3}$ , which avoids convergence of the biasing distributions to outliers and single points, see [34] for a similar technique. All runtime measurements were performed on compute nodes with Intel Xeon E5-1620 and 32GB RAM on a single core using a MATLAB implementation.

## 4.1 Heat transfer

We consider rare event probability estimation with a one-dimensional heat problem with two inputs.

### 4.1.1 Problem setup

Let  $\mathcal{X} = (0, 1) \in \mathbb{R}$  be a domain with boundary  $\partial\mathcal{X} = \{0, 1\}$ . Consider the linear elliptic PDE with random coefficients

$$-\nabla \cdot (a(\omega, x) \nabla u(\omega, x)) = 1, \quad x \in \mathcal{X}, \quad (20)$$

$$u(\omega, 0) = 0, \quad (21)$$

$$\partial_n u(\omega, 1) = 0, \quad (22)$$

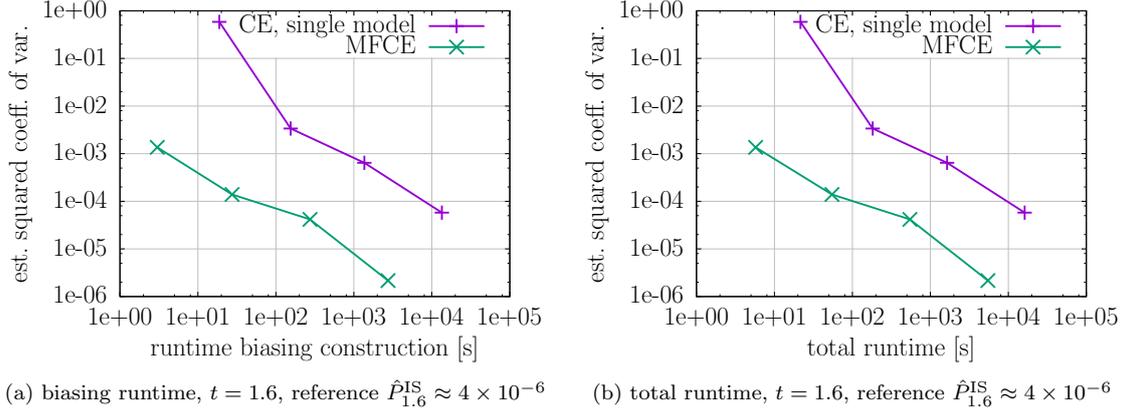


Figure 2: Heat transfer: Our MFCE approach achieves a speedup of about one order of magnitude for rare event probabilities of  $\approx 10^{-6}$  with  $t = 1.6$  in this example.

where  $u : \Omega \times \bar{\mathcal{X}} \rightarrow \mathbb{R}$  is the solution function defined on the set of outcomes  $\Omega$  and where  $\bar{\mathcal{X}}$  is the closure of  $\mathcal{X}$ . We impose homogeneous Dirichlet boundary condition on the left boundary  $x = 0$  of the domain  $\mathcal{X}$  and homogeneous Neumann boundary conditions on the right boundary  $x = 1$ . The coefficient  $a$  is given as

$$a(\omega, x) = \sum_{i=1}^n \exp(z_i(\omega)) \exp\left(-0.5 \frac{|x - v_i|}{0.0225}\right),$$

where  $n = 2$  and where  $Z = [z_1, z_2]^T$  is a random vector with components that are normally distributed with mean  $\boldsymbol{\mu} = [1, 1]^T$  and covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix} \in \mathbb{R}^{2 \times 2}.$$

The vector  $\boldsymbol{v} = [v_1, v_2]^T \in \mathbb{R}$  is  $\boldsymbol{v} = [0.5, 0.8]^T$ . The quantity of interest is the value of the solution function at the right boundary and given by the output of the function  $f : \mathcal{D} \rightarrow \mathbb{R}$  defined as

$$f(Z(\omega)) = u(\omega, 1).$$

We discretize (20)–(22) with linear finite elements on an equidistant grid with mesh width  $h^{(\ell)} = 2^{-\ell}$  on level  $\ell \in \mathbb{N}$ . The solution of the discretized problem on level  $\ell$  leads to models  $f^{(\ell)}$ . We set the maximal level to  $L = 8$ .

#### 4.1.2 Rare event probability estimation

Our goal is to estimate the rare event probabilities  $P_t^{(L)}$  for  $t \in \{1.4, 1.5, 1.6\}$ . We derive reference rare event probabilities with the classical, single-fidelity CE method with  $10^7$  realizations. To obtain these reference probabilities, we run Algorithm 1 with the model  $f^{(L)}$  only. We average over 30 runs and obtain the reference probability  $\hat{P}_{1.4}^{\text{IS}} \approx 6 \times 10^{-9}$  for  $t = 1.4$ , the reference probability  $\hat{P}_{1.5}^{\text{IS}} \approx 2 \times 10^{-7}$  for  $t = 1.5$ , and the reference probability  $\hat{P}_{1.6}^{\text{IS}} \approx 4 \times 10^{-6}$  for  $t = 1.6$ . For our MFCE method, we consider the levels  $\ell = 4, \dots, 8$  and run Algorithm 1 with  $m \in \{10^3, 10^4, 10^5, 10^6\}$  realizations. We repeat the estimation with MFCE 30 times and estimate the squared coefficient of variation (1) with respect to the reference probabilities.

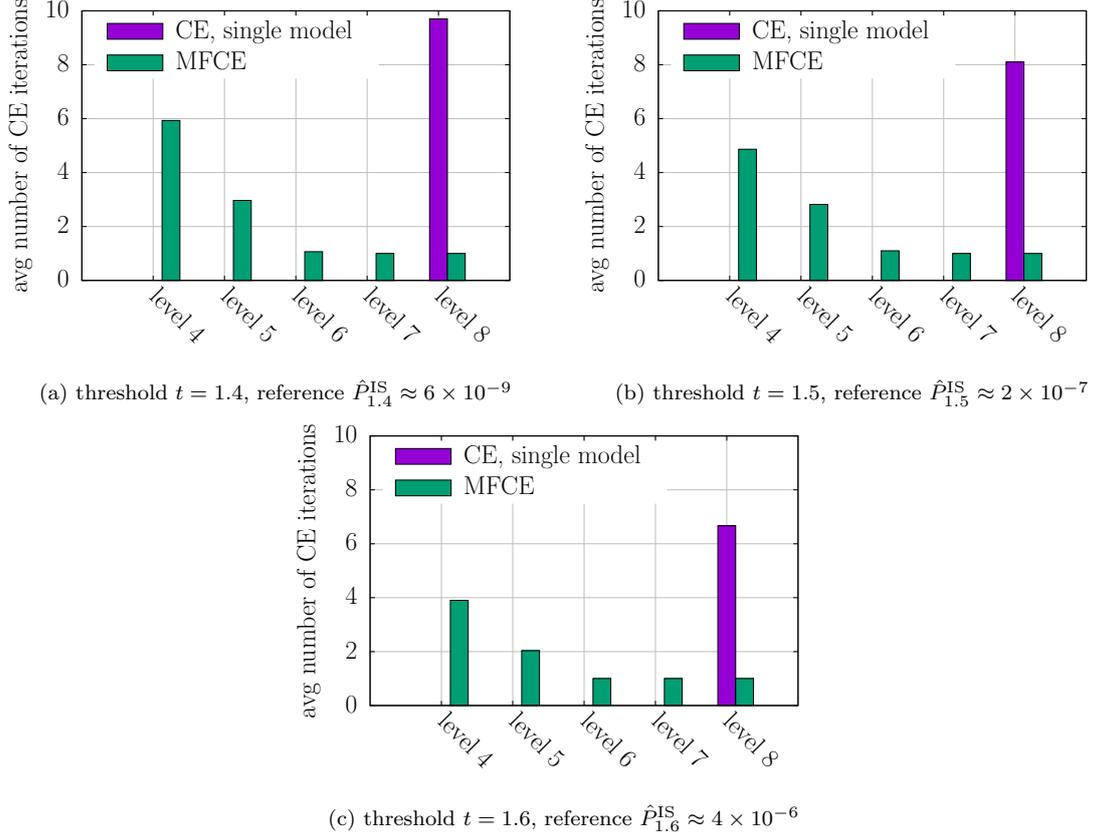


Figure 3: Heat transfer: Our MFCE approach achieves runtime speedups by shifting most of the iterations onto the models with coarse grids.

Figure 1a and Figure 1c compare the runtime of constructing the biasing density with our MFCE method to the runtime of the single-fidelity CE method that uses model  $f^{(L)}$  alone. Our MFCE approach achieves a speedup of up to two orders of magnitude. Figure 1b and Figure 1d show similar speedups for the total runtime, which includes the runtime of constructing the biasing density and the final estimation step. In this example, the total runtime is dominated by the runtime of constructing the biasing densities.

Figure 2 shows the speedup of our MFCE method for the threshold  $t = 1.6$ , which is smaller than the speedup obtained with the thresholds  $t = 1.4$  and  $t = 1.5$  shown Figure 1. The threshold  $t = 1.6$  corresponds to a reference probability of  $\approx 10^{-6}$ , which is significantly higher than the reference probabilities  $\approx 10^{-7}$  and  $\approx 10^{-9}$  corresponding to  $t = 1.5$  and  $t = 1.4$ , respectively. Typically fewer CE iterations are sufficient to construct a biasing density to estimate a rare event probability of  $\approx 10^{-6}$  than of  $\approx 10^{-9}$ . Thus, the results in Figure 2 confirm that our MFCE approach is particularly beneficial in cases where the CE method requires many iterations to obtain a biasing density, which typically is the case for small rare event probabilities. Overall, the results reported in Figure 1 and Figure 2 show that our MFCE method successfully leverages the hierarchy of models to estimate rare event probabilities that vary by about three orders of magnitude (i.e., from  $\approx 10^{-6}$  to  $\approx 10^{-9}$ ).

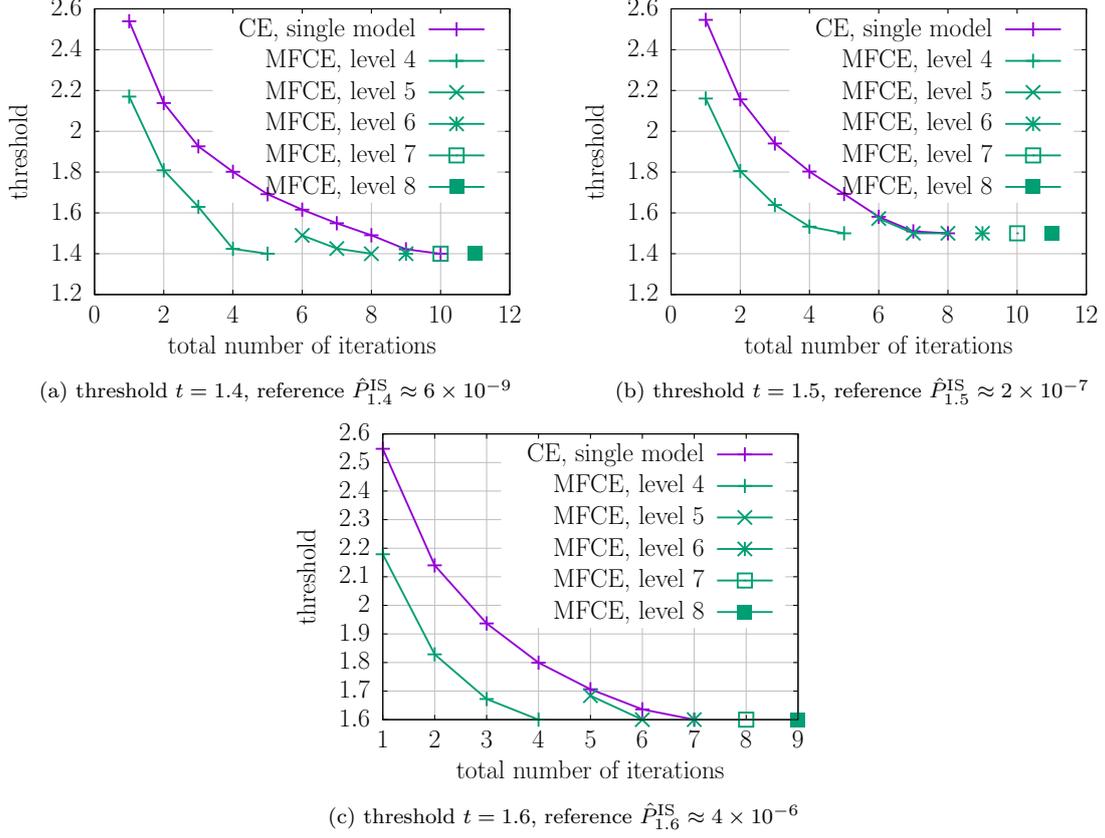


Figure 4: Heat transfer: The plot shows the intermediate thresholds selected by Algorithm 1. Our multifidelity approach achieves speedups because only a single iteration is required with the computationally expensive models on the fine grids (high level) to correct estimates obtained with the cheap models on the coarse grids (low levels).

#### 4.1.3 Number of model evaluations

Figure 3 compares the number of iterations spent at each level  $\ell = 4, \dots, 8$  of our MFCE method to the number of iterations of the single-fidelity CE method on level  $L = 8$ . The reported numbers of iterations are averaged over 30 runs. First, note that our MFCE approach and the single-fidelity CE method require most iterations for  $t = 1.4$ , which corresponds to the smallest rare event probability  $\approx 10^{-9}$  of the three cases  $t \in \{1.4, 1.5, 1.6\}$ . Second, the results confirm that our multifidelity approach spends most of the iterations with models on the coarse grids, where model evaluations are cheap compared to the model  $f^{(L)}$  on the finest grid. Consider now Figure 4, which reports the intermediate thresholds that are selected by Algorithm 1. Figure 4a shows that the single-fidelity CE method requires about 10 iterations on level  $L = 8$  to obtain an intermediate threshold that is equal to or below  $t = 1.4$ . Our MFCE approach requires five iterations on the lowest level  $\ell = 4$ ; however, the parameter  $\hat{v}_*^{(4)}$  estimated on level  $\ell = 4$  is then further corrected on level  $\ell = 5$  in three iterations. On levels  $\ell = 6, 7, 8$  only a single iteration is necessary to slightly correct the estimated parameter. Similar results are obtained for  $t \in \{1.5, 1.6\}$  shown in Figure 4b and Figure 4c.

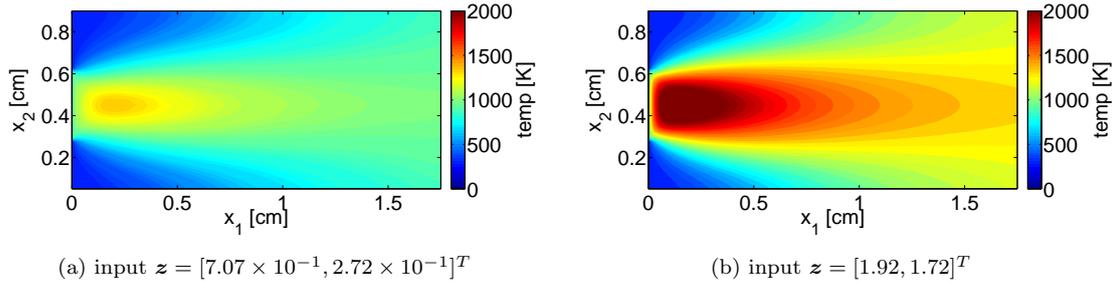


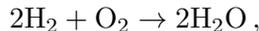
Figure 5: Reacting flow: The plots in (a) and (b) show the temperature of the reaction for two different inputs.

## 4.2 Reacting flow problem

This section demonstrates the MFCE approach on a reacting-flow problem.

### 4.2.1 Problem setup

We consider the simplified combustor model described in [5], which is based on the one-step reaction



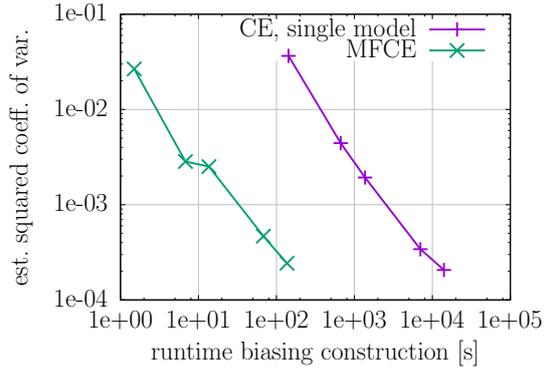
with the fuel  $\text{H}_2$ , the oxidizer  $\text{O}_2$ , and the product  $\text{H}_2\text{O}$ . The governing equations are nonlinear advection-diffusion-reaction equations that are discretized with finite differences on a mesh with equidistant grid points. The problem has two inputs  $\mathbf{z} = [z_1, z_2]^T$  that define properties of the reaction. The inputs are realizations of a random variable with normal distribution with mean  $\boldsymbol{\mu} = [1, 1]^T$  and covariance

$$\boldsymbol{\Sigma} = \begin{bmatrix} 0.0060 & 0 \\ 0 & 0.0037 \end{bmatrix}.$$

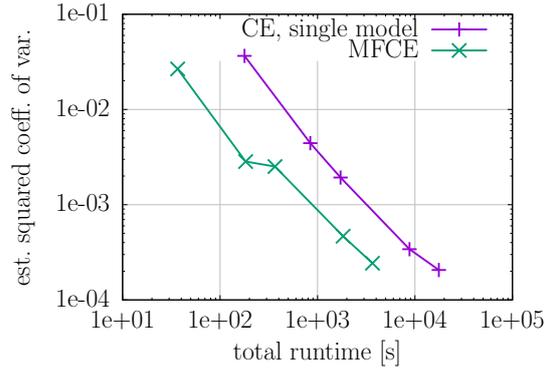
The output of the model is the maximum temperature in the combustion chamber, see Figure 5. We refer to [5, 26] for details.

The high-fidelity model in this experiment is given by the finite-difference model on a mesh with  $54 \times 27$  equidistant grid points. Furthermore, we derive a reduced model with proper orthogonal decomposition and the discrete empirical interpolation method as described in [36]. To construct the reduced model, we derive 100 snapshots with the high-fidelity model that correspond to inputs on an equidistant  $10 \times 10$  grid in the domain  $[0.7, 1.92] \times [0.27, 1.72]$  and derive proper orthogonal decomposition and empirical interpolation bases with 4 and 8 basis vectors, respectively. Additionally, we derive a piecewise-linear interpolant of the input-output map given by the high-fidelity model from four data points drawn from the distribution of the inputs. Thus, we have an interpolant  $f^{(1)}$ , a reduced model  $f^{(2)}$ , and a high-fidelity finite-difference model  $f^{(3)} = f^{(L)}$ .

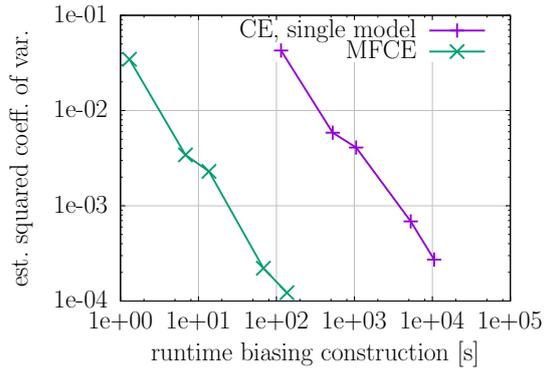
Our goal in this experiment is to estimate the probability that the temperature is below a threshold value, which can indicate a poor mixing in the reaction. We estimate the rare event probabilities for the thresholds  $t \in \{2021.3, 2043\}$ . We first run Algorithm 1 with the high-fidelity model  $f^{(L)}$  alone to obtain the reference probabilities  $\hat{P}_{2021.3}^{\text{IS}} \approx 2 \times 10^{-6}$  and  $\hat{P}_{2043}^{\text{IS}} \approx 2 \times 10^{-5}$ , respectively. The reference probabilities are estimated from  $10^4$  realizations and are averaged over 30 runs. We then run Algorithm 1 with the models  $f^{(1)}, f^{(2)}, f^{(3)}$  30 times with  $m \in \{10^2, 5 \times 10^2, 10^3, 5 \times 10^3, 10^4\}$  and estimate the squared coefficient of variation with respect to the reference probabilities.



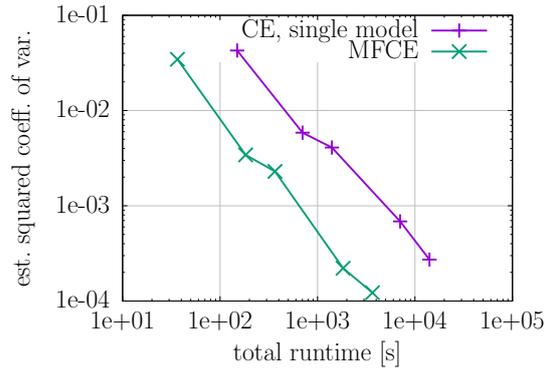
(a) biasing runtime,  $t = 2021.3$ ,  $\hat{P}_{2021.3}^{IS} \approx 2 \times 10^{-6}$



(b) total runtime,  $t = 2021.3$ , reference  $\hat{P}_{2021.3}^{IS} \approx 2 \times 10^{-6}$



(c) biasing runtime,  $t = 2043$ , reference  $\hat{P}_{2043}^{IS} \approx 2 \times 10^{-5}$



(d) total runtime,  $t = 2043$ , reference  $\hat{P}_{2043}^{IS} \approx 2 \times 10^{-5}$

Figure 6: Reacting flow: The total runtime is dominated by the estimation step with the high-fidelity model because a speedup of multiple orders of magnitude is achieved for constructing the biasing density and a speedup of at most one order of magnitude is obtained in the total runtime.

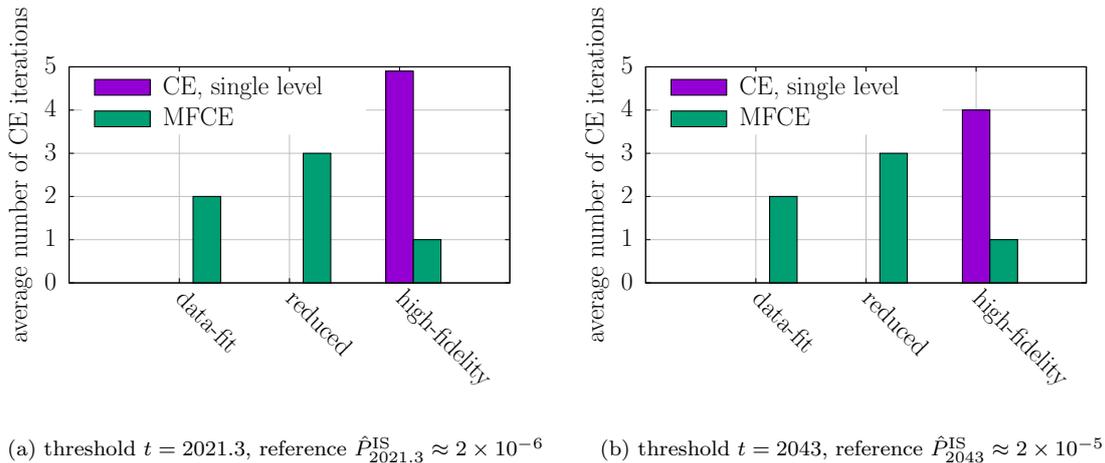


Figure 7: Reacting flow: Our MFCE approach shifts most of the iteration onto the reduced and data-fit interpolation models.

#### 4.2.2 Comparison of multi- and single-fidelity approaches

Figure 6 compares the runtime of our multifidelity approach with the runtime of the single-fidelity CE method that uses  $f^{(L)}$  alone. Our multifidelity approach achieves speedups of more than two orders of magnitude in the construction of the biasing densities. The large speedups are obtained because the data-fit  $f^{(1)}$  and the reduced model  $f^{(2)}$  are five and two orders of magnitude cheaper to evaluate than the high-fidelity model  $f^{(3)}$ , respectively. Consider now the total runtime, which includes the construction of the biasing densities and the final estimate step. The total runtime in this example is dominated by the costs of the final estimation step, which means that the speedup of our multifidelity approach obtained in constructing the biasing densities is smaller when considered with respect to the total runtime. Our method achieves up to an order of magnitude speedup in the total runtime.

Figure 7 reports the number of iterations per level. In our MFCE method, a single iteration with the high-fidelity model is sufficient, whereas the single-fidelity CE method requires up to almost 5 iterations. The intermediate thresholds selected by Algorithm 1 are shown in Figure 8. The results confirm that model  $f^{(1)}$  is a poor approximation of the high-fidelity model because the intermediate threshold selected on level  $\ell = 1$  needs to be corrected with three iterations on level  $\ell = 2$ .

## 5 Conclusions

We presented a multifidelity preconditioner for the CE method to accelerate the estimation of rare event probabilities. Our multifidelity approach leverages a hierarchy of surrogate models to reduce the costs of constructing biasing densities compared to the single-fidelity CE method that uses the high-fidelity model alone. Our approach can exploit multiple surrogate models that include general surrogate models such as projection-based reduced models and data-fit models, which goes beyond the classical setting of multilevel techniques that are often restricted to hierarchies of coarse-grid approximations. In our numerical examples, our approach achieved speedups of up to two orders of magnitude compared to the single-fidelity CE method in estimating probabilities on the order of  $10^{-5}$  to  $10^{-9}$ .

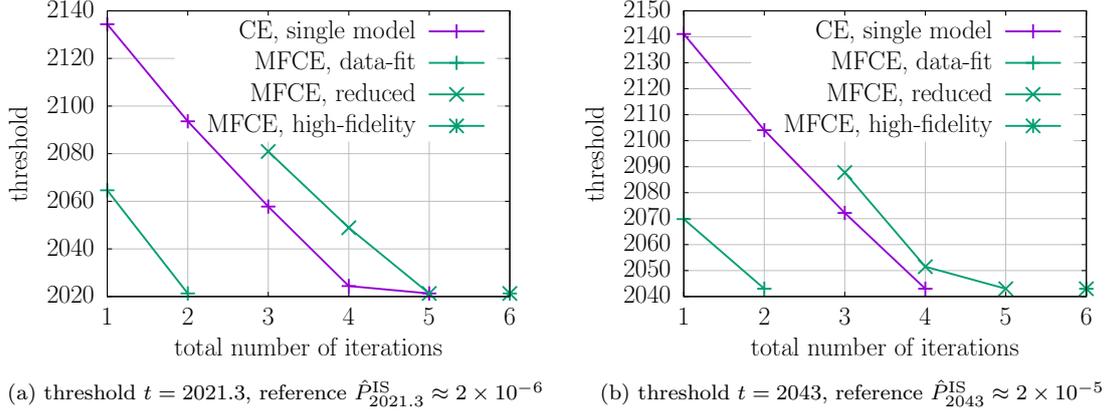


Figure 8: Reacting flow: The data-fit model  $f^{(1)}$  is a poor approximation of the high-fidelity model and therefore three iterations with the more accurate reduced model  $f^{(2)}$  are necessary to correct the intermediate thresholds. Overall, our multifidelity method leverages the data-fit and the reduced model to reduce the number of iterations required on the highest level with the high-fidelity model.

## References

- [1] S.-K. Au and J. L. Beck. Estimation of small failure probabilities in high dimensions by subset simulation. *Probabilistic Engineering Mechanics*, 16(4):263 – 277, 2001.
- [2] P. Benner, S. Gugercin, and K. Willcox. A survey of projection-based model reduction methods for parametric dynamical systems. *SIAM Review*, 57(4):483–531, 2015.
- [3] J.-M. Bourinet, F. Deheeger, and M. Lemaire. Assessing small failure probabilities by combined subset simulation and support vector machines. *Structural Safety*, 33(6):343 – 353, 2011.
- [4] J. Bucklew. *Introduction to Rare Event Simulation*. Springer, 2003.
- [5] M. Buffoni and K. Willcox. Projection-based model reduction for reacting flows. In *40th Fluid Dynamics Conference and Exhibit, Fluid Dynamics and Co-located Conferences*, pages 1–14. American Institute of Aeronautics and Astronautics, 2010.
- [6] P. Chen and A. Quarteroni. Accurate and efficient evaluation of failure probability for partial differential equations with random input data. *Computer Methods in Applied Mechanics and Engineering*, 267:233–260, 2013.
- [7] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [8] P.-T. de Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein. A tutorial on the cross-entropy method. *Annals of Operations Research*, 134(1):19–67, 2005.
- [9] T. H. de Mello and R. Y. Rubinstein. Estimation of rare event probabilities using cross-entropy. In *Winter Simulation Conference*, volume 01, pages 310–319, 2002.
- [10] V. Dubourg, B. Sudret, and F. Deheeger. Metamodel-based importance sampling for structural reliability analysis. *Probabilistic Engineering Mechanics*, 33:47 – 57, 2013.

- [11] D. Elfverson, D. J. Estep, F. Hellman, and A. Målqvist. Uncertainty quantification for approximate p-quantiles for physical models with stochastic inputs. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):826–850, 2014.
- [12] D. Elfverson, F. Hellman, and A. Målqvist. A multilevel Monte Carlo method for computing failure probabilities. *SIAM/ASA Journal on Uncertainty Quantification*, 4(1):312–330, 2016.
- [13] F. Fagerlund, F. Hellman, A. Målqvist, and A. Niemi. Multilevel monte carlo methods for computing failure probability of porous media flow systems. *Advances in Water Resources*, 94:498 – 509, 2016.
- [14] G. Fishman. *Monte Carlo*. Springer, 1996.
- [15] A. Forrester, A. Sóbester, and A. Keane. *Engineering design via surrogate modelling: a practical guide*. Wiley, 2008.
- [16] M. B. Giles. Multilevel monte carlo path simulation. *Operations Research*, 56(3):607–617, 2008.
- [17] S. Heinrich. Multilevel Monte Carlo Methods. In S. Margenov, J. Waśniewski, and P. Yalamov, editors, *Large-Scale Scientific Computing*, number 2179 in Lecture Notes in Computer Science, pages 58–67. Springer, 2001.
- [18] J. Li, J. Li, and D. Xiu. An efficient surrogate-based method for computing rare failure probability. *Journal of Computational Physics*, 230(24):8683–8697, 2011.
- [19] J. Li and D. Xiu. Evaluation of failure probability via surrogate models. *Journal of Computational Physics*, 229(23):8966–8980, 2010.
- [20] J. Li and D. Xiu. Surrogate based method for evaluation of failure probability under multiple constraints. *SIAM Journal on Scientific Computing*, 36(2):A828–A845, 2014.
- [21] A. J. Majda and B. Gershgorin. Quantifying uncertainty in climate change science through empirical information theory. *Proceedings of the National Academy of Sciences of the United States of America*, 107(34):14958–14963, Aug. 2010.
- [22] L. Ng and K. Willcox. Multifidelity approaches for optimization under uncertainty. *International Journal for Numerical Methods in Engineering*, 100(10):746–772, 2014.
- [23] V. Papadopoulos, D. G. Giovanis, N. D. Lagaros, and M. Papadrakakis. Accelerated subset simulation with neural networks for reliability analysis. *Computer Methods in Applied Mechanics and Engineering*, 223–224:70 – 80, 2012.
- [24] B. Peherstorfer, T. Cui, Y. Marzouk, and K. Willcox. Multifidelity importance sampling. *Computer Methods in Applied Mechanics and Engineering*, 300:490 – 509, 2016.
- [25] B. Peherstorfer, B. Kramer, and K. Willcox. Combining multiple surrogate models to accelerate failure probability estimation with expensive high-fidelity models. Technical report, Aerospace Computational Design Laboratory, Massachusetts Institute of Technology, 2016.
- [26] B. Peherstorfer and K. Willcox. Online adaptive model reduction for nonlinear systems via low-rank updates. *SIAM Journal on Scientific Computing*, 37(4):A2123–A2150, 2015.

- [27] B. Peherstorfer, K. Willcox, and M. Gunzburger. Optimal model management for multifidelity Monte Carlo estimation. *SIAM Journal on Scientific Computing*, 38(5):A3163–A3194, 2016.
- [28] C. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer, 2004.
- [29] G. Rozza, D. Huynh, and A. Patera. Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations. *Archives of Computational Methods in Engineering*, 15(3):1–47, 2007.
- [30] R. Rubinstein. The cross-entropy method for combinatorial and continuous optimization. *Methodology And Computing In Applied Probability*, 1(2):127–190, 1999.
- [31] R. Y. Rubinstein. Optimization of computer simulation models with rare events. *European Journal of Operational Research*, 99(1):89 – 112, 1997.
- [32] R. Y. Rubinstein. Combinatorial optimization, cross-entropy, ants and rare events. In S. Uryasev and P. M. Pardalos, editors, *Stochastic Optimization: Algorithms and Applications*, pages 303–363, Boston, MA, 2001. Springer US.
- [33] R. Srinivasan. *Importance Sampling*. Springer, 2002.
- [34] I. Szita and A. Lőrincz. Learning tetris using the noisy cross-entropy method. *Neural Computation*, 18:2936–2941, 2006.
- [35] E. Ullmann and I. Papaioannou. Multilevel estimation of rare events. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):922–953, 2015.
- [36] Y. B. Zhou. *Model reduction for nonlinear dynamical systems with parametric uncertainties*. Thesis (S.M.), Massachusetts Institute of Technology, Dept. of Aeronautics and Astronautics, 2012.
- [37] K. M. Zuev, J. L. Beck, S.-K. Au, and L. S. Katafygiotis. Bayesian post-processor and other enhancements of subset simulation for estimating failure probabilities in high dimensions. *Computers & Structures*, 92–93:283 – 296, 2012.