

1 **AN ENSEMBLE ALGORITHM FOR NUMERICAL SOLUTIONS TO**
2 **DETERMINISTIC AND RANDOM PARABOLIC PDES**

3 YAN LUO * AND ZHU WANG[†]

4 **Abstract.** In this paper, we develop an ensemble-based time-stepping algorithm to efficiently
5 find numerical solutions to a group of linear, second-order parabolic partial differential equations
6 (PDEs). Particularly, the PDE models in the group could be subject to different diffusion coefficients,
7 initial conditions, boundary conditions, and body forces. The proposed algorithm leads to a single
8 discrete system for the group with multiple right-hand-side vectors by introducing an ensemble
9 average of the diffusion coefficient functions and using a new semi-implicit time integration method.
10 The system could be solved more efficiently than multiple linear systems with a single right-hand-
11 side vector. We first apply the algorithm to deterministic parabolic PDEs and derive a rigorous
12 error estimate that shows the scheme is first-order accurate in time and is optimally accurate in
13 space. We then extend it to find stochastic solutions of parabolic PDEs with random coefficients
14 and put forth an ensemble-based Monte Carlo method. The effectiveness of the new approach is
15 demonstrated through theoretical analysis. Several numerical experiments are presented to illustrate
16 our theoretical results.

17 **Key words.** ensemble-based method, parabolic PDEs, random parabolic PDEs, Monte Carlo
18 method

19 **AMS subject classifications.** 65C05, 65C20, 65M60

20 **1. Introduction.** In many application problems involving numerical simulations
21 of PDEs, one is not interested in a single simulation, but a number of simulations
22 with different computational settings such as distinct initial condition, boundary con-
23 ditions, body force, and physical parameters. For instance, data assimilation methods
24 used in meteorology, such as the ensemble Kalman filter [14], run a numerical weather
25 model forward for many times by perturbing initial conditions and uncertain param-
26 eters. The ensemble of outputs is then used to update not only the forecast, but
27 also the covariance matrix. Similarly, when a random sampling method is used in
28 uncertainty quantification, a model is run forward for a large number of times at
29 the selected parameter sample values in order to collect outputs and determine the
30 underlying statistical information [7]. Such a group of simulations is computationally
31 expensive, especially, when the forward simulation is of a large scale.

32 Aiming at developing a fast algorithm for such applications, an ensemble time-
33 stepping algorithm was proposed in [12], where a group of size J Navier-Stokes equa-
34 tions (NSE) with different initial conditions and forcing terms is simulated. All so-
35 lutions are found by solving a single linear system with one shared coefficient matrix
36 and J right-hand-side (RHS) vectors at each time step. Thus, it reduces the storage
37 requirements and computational cost for the group of simulations. The algorithm
38 is first-order accurate in time, which was extended to higher-order accurate schemes
39 in [9, 10]. For high Reynolds number incompressible flows, ensemble regularization
40 methods were developed in [9, 13, 21], and a turbulence model based on ensemble

*School of Mathematics, Sichuan University, No.24 South Section 1, Yihuan Road, Chengdu, Sichuan 610064, P. R. China and School of Mathematical Sciences, University of Electronic Science and Technology of China, No.2006, Xiyuan Ave, West Hi-Tech Zone, Chengdu, Sichuan 611731, P. R. China. Research supported by the Young Scientists Fund of the National Natural Science Foundation of China grant 11501088.

[†]Department of Mathematics, University of South Carolina, 1523 Greene Street, Columbia, SC 29208, USA (wangzhu@math.sc.edu). Research supported by the U.S. National Science Foundation grant DMS-1522672 and the U.S. Department of Energy grant DE-SC0016540.

41 averaging was developed in [11]. The ensemble algorithm has also been extended
 42 to simulate MHD flows in [17] and to the reduced-order modeling setting in [4, 5].
 43 For parametrized flows, the ensemble algorithms were developed in [6] for multiple
 44 numerical simulations subject to not only different initial condition, boundary condi-
 45 tions and body force, but also physical parameters. It is worth mentioning that the
 46 ensemble method is only applied to problems with constant physical parameters so
 47 far.

48 In this paper, we consider a group of numerical solutions to second-order parabolic
 49 PDEs. We first develop an ensemble algorithm for deterministic problems, in which
 50 the physical parameter—diffusion coefficient—is a function of space and time; and
 51 then extend the method to stochastic problems. To our knowledge, this is the *first*
 52 time that an ensemble scheme is derived for problems with non-constant parameters
 53 and is further applied to PDEs with random coefficients.

54 The initial boundary value problem we consider is as follows.

$$55 \quad (1) \quad \begin{cases} u_t - \nabla \cdot [a(\mathbf{x}, t)\nabla u] = f(\mathbf{x}, t), & \text{in } D \times (0, T), \\ u(\mathbf{x}, t) = g(\mathbf{x}, t), & \text{on } \partial D \times (0, T), \\ u(\mathbf{x}, 0) = u^0(\mathbf{x}), & \text{in } D, \end{cases}$$

56 where D is a bounded Lipschitz domain in \mathbb{R}^d for $d = 1, 2, 3$, the diffusion coefficient
 57 $a(\mathbf{x}, t) \in L^2(W^{1,\infty}(D); 0, T)$, body force $f(\mathbf{x}, t) \in L^2(H^{-1}(D); 0, T)$, initial condition
 58 $u^0(\mathbf{x}) \in H_0^1(D) \cap H^{l+1}(D)$ with $l \geq 1$.

59 We consider the setting in which one needs to run a group of simulations. Each of
 60 these simulations is subject to independent initial and boundary conditions, body force
 61 and diffusion coefficient. Suppose the ensemble of simulations includes J independent
 62 members, in which the j -th member satisfies:

$$63 \quad (2) \quad \begin{cases} u_{j,t} - \nabla \cdot [a_j(\mathbf{x}, t)\nabla u_j] = f_j(\mathbf{x}, t), & \text{in } D \times (0, T), \\ u_j(\mathbf{x}, t) = g_j(\mathbf{x}, t), & \text{on } \partial D \times (0, T), \\ u_j(\mathbf{x}, 0) = u_j^0(\mathbf{x}), & \text{in } D, \end{cases}$$

64 for $j = 1, 2, \dots, J$. In general, when an implicit time stepping method is applied to the
 65 above system, the related discrete system would change as j varies. As a result, one
 66 needs to solve J linear systems at each time step. In order to improve its efficiency,
 67 we develop an ensemble-based time-stepping scheme. For illustrating the main idea,
 68 we first present the semi-discrete (in time) system.

69 For simplicity, we consider a uniform time partition on $[0, T]$ with the step size
 70 Δt . Denote by u_j^n , a_j^n and f_j^n the functions u_j , a_j and f_j evaluated at the time instant
 71 $t_n = n\Delta t$. Define the ensemble mean of the diffusion coefficient functions at time t_n
 72 by

$$73 \quad (3) \quad \bar{a}^n := \frac{1}{J} \sum_{j=1}^J a_j(\mathbf{x}, t_n).$$

74 The new ensemble-based time stepping scheme reads, for $j = 1, \dots, J$:

$$75 \quad (4) \quad \frac{u_j^{n+1} - u_j^n}{\Delta t} - \nabla \cdot (\bar{a}^{n+1}\nabla u_j^{n+1}) - \nabla \cdot [(a_j^{n+1} - \bar{a}^{n+1})\nabla u_j^n] = f_j^{n+1}$$

76 with the same boundary and initial conditions as those in (2). Rearranging the
 77 equation gives

$$78 \quad (5) \quad \frac{1}{\Delta t} u_j^{n+1} - \nabla \cdot (\bar{a}^{n+1}\nabla u_j^{n+1}) = f_j^{n+1} + \frac{1}{\Delta t} u_j^n + \nabla \cdot [(a_j^{n+1} - \bar{a}^{n+1})\nabla u_j^n].$$

79 It is easy to see, after a spatial discretization, the coefficient matrix of the resulting
80 linear system will be independent of j . This represents the key feature of the ensemble
81 method, that is, the discrete systems associated with an ensemble of simulations
82 share a unique coefficient matrix and only the RHS vectors vary among the ensemble
83 members. Thus, when the considered problem has a small scale, one only needs to
84 do LU factorization of the coefficient matrix once and use it to find the solutions of
85 the group; when the problem is of a large scale, the proposed scheme together with
86 the block Krylov subspace iterative methods will improve the efficiency of a number
87 of numerical simulations (see, e.g. [20] and references therein).

88 The rest of this paper is structured as follows. In Section 2, we introduce some
89 notations and mathematical preliminaries. In Section 3, we analyze the ensemble
90 scheme (4) under a full discretization in both space and time, and prove its stabil-
91 ity and convergence. In Section 4, we discuss the random parabolic equations and
92 provide stability analysis and error estimations. Several numerical experiments are
93 presented in Section 5, which illustrates the effectiveness of our proposed scheme on
94 both deterministic and random parabolic problems. A few concluding remarks are
95 given in Section 6.

96 **2. Notations and preliminaries.** Denote the $L^2(D)$ norm and inner prod-
97 uct by $\|\cdot\|$ and (\cdot, \cdot) , respectively. Let $L^\infty(D)$ be the set of bounded measurable
98 functions equipped with the norm $\|v\|_\infty := \text{ess sup}_{\mathbf{x} \in D} |v|$. Denoted by $W^{s,q}(D)$ the
99 Sobolev space of functions having generalized derivatives up to the order s in the
100 space $L^q(D)$, where s is a nonnegative integer and $1 \leq q \leq +\infty$. The equipped
101 Sobolev norm of $v \in W^{s,q}(D)$ is denoted by $\|v\|_{W^{s,q}(D)}$. When $q = 2$, we use the
102 notation $H^s(D)$ instead of $W^{s,2}(D)$. As usual, the function space $H_0^1(D)$ is the sub-
103 space of $H^1(D)$ consisting of functions that vanish on the boundary of D in the sense
104 of trace, equipped with the norm $\|v\|_{H_0^1(D)} = (\int_D |\nabla v|^2 d\mathbf{x})^{1/2}$. When $s = 0$, we
105 shall keep the notation with $L^q(D)$ instead of $W^{0,q}(D)$. The space $H^{-s}(D)$ is the
106 dual space of bounded linear functions on $H_0^s(D)$. A norm for $H^{-1}(D)$ is defined by
107 $\|f\|_{-1} = \sup_{0 \neq v \in H_0^1(D)} (f, v) / \|\nabla v\|$.

Stochastic functions have different structures. Let (Ω, \mathcal{F}, P) be a complete prob-
ability space, where Ω is the set of outcomes, $\mathcal{F} \subset 2^\Omega$ is the σ -algebra of events, and
 $P : \mathcal{F} \rightarrow [0, 1]$ is a probability measure. If Y is a random variable in the space and
belongs to $L_P^1(\Omega)$, its expected value is defined by

$$E[Y] = \int_{\Omega} Y(\omega) dP(\omega).$$

108 With the multi-index notation, $\alpha = (\alpha_1, \dots, \alpha_d)$ is a d -tuple of nonnegative inte-
109 gers with the length $|\alpha| = \sum_{i=1}^d \alpha_i$. The stochastic Sobolev spaces $\widetilde{W}^{s,q}(D) =$
110 $L_P^q(\Omega, W^{s,q}(D))$ containing stochastic function, $v : \Omega \times D \rightarrow R$, that are measur-
111 able with respect to the product σ -algebra $\mathcal{F} \otimes B(D)$ and equipped with the aver-
112 aged norms $\|v\|_{\widetilde{W}^{s,q}(D)} = \left(E \left[\|v\|_{W^{s,q}(D)}^q \right] \right)^{1/q} = \left(E \left[\sum_{|\alpha| \leq s} \int_D |\partial^\alpha v|^q d\mathbf{x} \right] \right)^{1/q}$, $1 \leq$
113 $q < +\infty$. Observe that if $v \in \widetilde{W}^{s,q}(D)$, then $v(\omega, \cdot) \in W^{s,q}(D)$ almost surely
114 (a.s.) and $\partial^\alpha v(\cdot, \mathbf{x}) \in L_P^q(\Omega)$ almost everywhere (a.e.) on the D for $|\alpha| \leq s$.
115 Whenever $q = 2$, the above space is a Hilbert space, i.e., $\widetilde{W}^{s,2}(D) = \widetilde{H}^s(D) \simeq$
116 $L_P^2(\Omega) \otimes H^s(D)$. In this paper, we consider the tensor product Hilbert space $H =$
117 $\widetilde{L}^2(H_0^1(D); 0, T) \simeq L_P^2(\Omega; H_0^1(D); 0, T)$ endowed with the inner product $(v, u)_H \equiv$
118 $E \left[\int_0^T \int_D \nabla v \cdot \nabla u d\mathbf{x} dt \right]$.

120 **3. The ensemble scheme of deterministic parabolic equations.** We first
 121 consider the deterministic parabolic equations and analyze the full discretization of
 122 the proposed ensemble scheme (4). For simplicity of presentation, we consider the
 123 problem with a homogeneous boundary condition, while the nonhomogeneous cases
 124 can be similarly analyzed incorporating the method of shifting (see Section 5.4 in
 125 [2]). Furthermore, we include numerical test cases with nonhomogeneous boundary
 126 conditions in Section 5.

127 Suppose the following two conditions are valid:

128 (i) There exists a positive constant θ such that, for any $t \in [0, T]$,

129 (6)
$$\min_{\mathbf{x} \in \bar{D}} a(\mathbf{x}, t) \geq \theta;$$

130 (ii) There exist positive constants θ_- and θ_+ such that, for any $t \in [0, T]$,

131 (7)
$$\theta_- \leq |a_j(\mathbf{x}, t) - \bar{a}(\mathbf{x}, t)|_\infty \leq \theta_+.$$

132 Obviously, Condition (i) guarantees the uniform coercivity of the problem; Condition
 133 (ii) states that the distance from coefficient $a_j(\mathbf{x}, t)$ to the ensemble average $\bar{a}(\mathbf{x}, t)$ is
 134 uniformly bounded.

Let \mathcal{T}_h be a quasi-uniform triangulation of the domain D , made of elements K ,
 such that $\bar{D} = \bigcup_{K \in \mathcal{T}_h} \bar{K}$. Define the mesh size $h := \max_{K \in \mathcal{T}_h} h_K$, where h_K is the
 diameter of the element K . Denoted by V_h the finite element space

$$V_h := \{v \in H_0^1(D) \cap H^{l+1}(D); v|_K \text{ is a polynomial of degree } l, \forall K \in \mathcal{T}_h\}.$$

135 With the assumed uniform time partition on $[0, T]$ and set $N = T/\Delta t$, the fully
 136 discrete approximation of (4) is as follows: Find $u_{j,h}^{n+1} \in V_h$ satisfying, for $n =$
 137 $0, \dots, N-1$ and $j = 1, \dots, J$:

138 (8)
$$\begin{aligned} & \left(\frac{u_{j,h}^{n+1} - u_{j,h}^n}{\Delta t}, v_h \right) + \left(\bar{a}^{n+1} \nabla u_{j,h}^{n+1}, \nabla v_h \right) + \left((a_j^{n+1} - \bar{a}^{n+1}) \nabla u_{j,h}^n, \nabla v_h \right) \\ & = (f_j^{n+1}, v_h), \quad \forall v_h \in V_h \end{aligned}$$

139 with the initial condition $u_{j,h}^0 \in V_h$ satisfying $(u_{j,h}^0, v_h) = (u_j^0, v_h), \forall v_h \in V_h$.

140 **3.1. Stability and convergence.** We first discuss the stability of the ensemble
 141 algorithm (8).

142 **THEOREM 1.** *Suppose $f_j \in L^2(H^{-1}(D); 0, T)$ and conditions (i) and (ii) are sat-*
 143 *isfied, the ensemble scheme (8) is stable provided that*

144 (9)
$$\theta > \theta_+.$$

145 *Furthermore, the numerical solution to (8) satisfies*

146 (10)
$$\begin{aligned} \|u_{j,h}^N\|^2 + \theta_- \Delta t \|\nabla u_{j,h}^N\|^2 + (\theta - \theta_+) \Delta t \sum_{n=1}^N \|\nabla u_{j,h}^n\|^2 & \leq C \Delta t \sum_{n=1}^N \|f_j^n\|_{-1}^2 \\ & + C \Delta t \|\nabla u_{j,h}^0\|^2 + \|u_{j,h}^0\|^2, \end{aligned}$$

147 *where C is a generic positive constant independent of J, h and Δt .*

148 *Proof.* Taking $v_h = u_{j,h}^{n+1}$ in (8), we have

$$149 \quad \frac{1}{\Delta t} (u_{j,h}^{n+1} - u_{j,h}^n, u_{j,h}^{n+1}) + (\bar{a}^{n+1} \nabla u_{j,h}^{n+1}, \nabla u_{j,h}^{n+1}) + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla u_{j,h}^n, \nabla u_{j,h}^{n+1}) \\ = (f_j^{n+1}, u_{j,h}^{n+1}).$$

150 Multiplying both sides by Δt , and using the polarization identity and coercivity of
151 \bar{a}^{n+1} , we get

$$152 \quad (11) \quad \frac{1}{2} \|u_{j,h}^{n+1}\|^2 - \frac{1}{2} \|u_{j,h}^n\|^2 + \frac{1}{2} \|u_{j,h}^{n+1} - u_{j,h}^n\|^2 + \Delta t \theta \|\nabla u_{j,h}^{n+1}\|^2 \\ \leq -\Delta t ((a_j^{n+1} - \bar{a}^{n+1}) \nabla u_{j,h}^n, \nabla u_{j,h}^{n+1}) + \Delta t (f_j^{n+1}, u_{j,h}^{n+1}).$$

153 By the Cauchy-Schwarz and Young's inequalities, we have, for some $\mu, \alpha > 0$,

$$154 \quad \Delta t \left| ((a_j^{n+1} - \bar{a}^{n+1}) \nabla u_{j,h}^n, \nabla u_{j,h}^{n+1}) \right| \leq \Delta t |a_j^{n+1} - \bar{a}^{n+1}|_\infty \|\nabla u_{j,h}^n\| \|\nabla u_{j,h}^{n+1}\| \\ 155 \quad (12) \quad \leq \Delta t |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\frac{1}{2\mu} \|\nabla u_{j,h}^n\|^2 + \frac{\mu}{2} \|\nabla u_{j,h}^{n+1}\|^2 \right)$$

156 and

$$157 \quad (13) \quad \Delta t |(f_j^{n+1}, u_{j,h}^{n+1})| \leq \Delta t \|f_j^{n+1}\|_{-1} \|\nabla u_{j,h}^{n+1}\| \leq \frac{\Delta t}{4\alpha} \|f_j^{n+1}\|_{-1}^2 + \alpha \Delta t \|\nabla u_{j,h}^{n+1}\|^2.$$

158 Substituting (12) and (13) into (11) and dropping the non-negative term $\frac{1}{2} \|u_{j,h}^{n+1} -$
159 $u_{j,h}^n\|^2$, we get

$$160 \quad \frac{1}{2} \left(\|u_{j,h}^{n+1}\|^2 - \|u_{j,h}^n\|^2 \right) + \Delta t \left[\theta - \alpha - \left(\frac{\mu}{2} + \frac{1}{2\mu} \right) |a_j^{n+1} - \bar{a}^{n+1}|_\infty \right] \|\nabla u_{j,h}^{n+1}\|^2 \\ 161 \quad + \frac{\Delta t}{2\mu} |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\|\nabla u_{j,h}^{n+1}\|^2 - \|\nabla u_{j,h}^n\|^2 \right) \leq \frac{\Delta t}{4\alpha} \|f_j^{n+1}\|_{-1}^2.$$

162 Multiplying both sides by 2, summing over n from 0 to $N-1$ and taking $\mu = 1$ yields

$$163 \quad \|u_{j,h}^N\|^2 - \|u_{j,h}^0\|^2 + 2\Delta t \sum_{n=0}^{N-1} \left(\theta - \alpha - |a_j^{n+1} - \bar{a}^{n+1}|_\infty \right) \|\nabla u_{j,h}^{n+1}\|^2 \\ 164 \quad (14) \quad + \Delta t \sum_{n=0}^{N-1} |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\|\nabla u_{j,h}^{n+1}\|^2 - \|\nabla u_{j,h}^n\|^2 \right) \leq \frac{\Delta t}{2\alpha} \sum_{n=0}^{N-1} \|f_j^{n+1}\|_{-1}^2.$$

165 We then choose $\alpha = (\theta - |a_j^{n+1} - \bar{a}^{n+1}|_\infty)/2$ and use the conditions (7) and (9), and
166 obtain

$$167 \quad \|u_{j,h}^N\|^2 + \theta_- \Delta t \|\nabla u_{j,h}^N\|^2 + (\theta - \theta_+) \Delta t \sum_{n=0}^{N-1} \|\nabla u_{j,h}^{n+1}\|^2 \leq \frac{\Delta t}{\theta - \theta_+} \sum_{n=0}^{N-1} \|f_j^{n+1}\|_{-1}^2 \\ + \theta_- \Delta t \|\nabla u_{j,h}^0\|^2 + \|u_{j,h}^0\|^2.$$

168 This completes the proof. \square

169 **REMARK 2.** *The stability condition (9) requires, for $\{a_j\}_{j=1}^J$, the deviation of a_j*
170 *from the ensemble average \bar{a} to be less than the coercivity constant θ . If this is not the*
171 *case, one might divide the ensemble into smaller groups so that the stability condition*
172 *holds in each of them and the algorithm is stable and applicable.*

173 Next, we estimate the approximation error of the ensemble algorithm (8). We
 174 assume the exact solution of the PDEs is smooth enough, in particular,

$$175 \quad u_j \in L^2(H_0^1(D) \cap H^{l+1}(D); 0, T) \cap H^1(H^{l+1}(D); 0, T) \cap H^2(L^2(D); 0, T).$$

176

177 **THEOREM 3.** *Let u_j^n and $u_{j,h}^n$ be the solutions of equations (2) and (8) at time t_n ,*
 178 *respectively. Assume $f_j \in L^2(H^{-1}(D); 0, T)$ and conditions (i) and (ii) hold. Then*
 179 *there exists a generic constant $C > 0$ independent of J , h and Δt such that*

$$180 \quad \|u_j^N - u_{j,h}^N\|^2 + \theta_- \Delta t \|\nabla(u_j^N - u_{j,h}^N)\|^2 + (\theta - \theta_+) \Delta t \sum_{n=1}^N \|\nabla(u_j^n - u_{j,h}^n)\|^2$$

$$181 \quad (15) \quad \leq C(\Delta t^2 + h^{2l}),$$

182 *provided that the stability condition (9) holds, that is, $\theta > \theta_+$.*

183 *Proof.* In order to estimate the approximation error of (8), we first find the error
 184 equation. Evaluating the equation (2) at $t = t_{n+1}$, tested by $v_h \in V_h$, yields

$$185 \quad (16) \quad \frac{1}{\Delta t} (u_j^{n+1} - u_j^n, v_h) + (a_j^{n+1} \nabla u_j^{n+1}, \nabla v_h) = (f_j^{n+1}, v_h) - (r_j^{n+1}, v_h),$$

186 where $r_j^{n+1} = u_{j,t}^{n+1} - \frac{u_j^{n+1} - u_j^n}{\Delta t}$. Denoted by $e_j^n := u_j^n - u_{j,h}^n$ the approximation error
 187 of the j -th simulation at time t_n . Subtracting (8) from (16), we have

$$188 \quad \frac{1}{\Delta t} (e_j^{n+1} - e_j^n, v_h) + (\bar{a}^{n+1} \nabla e_j^{n+1}, \nabla v_h) + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla e_j^n, \nabla v_h)$$

$$189 \quad (17) \quad + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla (u_j^{n+1} - u_j^n), \nabla v_h) + (r_j^{n+1}, v_h) = 0.$$

Now we decompose the error as

$$e_j^n = (u_j^n - P_h(u_j^n)) - (u_{j,h}^n - P_h(u_j^n)) = \rho_j^n - \phi_{j,h}^n,$$

190 where $\rho_j^n = u_j^n - P_h(u_j^n)$ and $\phi_{j,h}^n = u_{j,h}^n - P_h(u_j^n)$ with $P_h(u_j^n)$ the L^2 projection of
 191 u_j^n onto V_h , that is, $(\rho_j^n, v_h) = 0$ for any $v_h \in V_h$. We then use the decomposition in
 192 (17) and obtain

$$193 \quad \frac{1}{\Delta t} (\phi_{j,h}^{n+1} - \phi_{j,h}^n, v_h) + (\bar{a}^{n+1} \nabla \phi_{j,h}^{n+1}, \nabla v_h) + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla \phi_{j,h}^n, \nabla v_h)$$

$$194 \quad = (\bar{a}^{n+1} \nabla \rho_j^{n+1}, \nabla v_h) + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla \rho_j^n, \nabla v_h)$$

$$195 \quad (18) \quad + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla (u_j^{n+1} - u_j^n), \nabla v_h) + (r_j^{n+1}, v_h).$$

196 Letting $v_h = \phi_{j,h}^{n+1}$ and using the polarization identity and coercivity (6), we have

$$197 \quad \frac{1}{2\Delta t} (\|\phi_{j,h}^{n+1}\|^2 - \|\phi_{j,h}^n\|^2 + \|\phi_{j,h}^{n+1} - \phi_{j,h}^n\|^2) + \theta \|\nabla \phi_{j,h}^{n+1}\|^2$$

$$198 \quad \leq |((a_j^{n+1} - \bar{a}^{n+1}) \nabla \phi_{j,h}^n, \nabla \phi_{j,h}^{n+1})|$$

$$199 \quad + |(\bar{a}^{n+1} \nabla \rho_j^{n+1}, \nabla \phi_{j,h}^{n+1})| + |((a_j^{n+1} - \bar{a}^{n+1}) \nabla \rho_j^n, \nabla \phi_{j,h}^{n+1})|$$

$$200 \quad (19) \quad + |((a_j^{n+1} - \bar{a}^{n+1}) \nabla (u_j^{n+1} - u_j^n), \nabla \phi_{j,h}^{n+1})| + |(r_j^{n+1}, \phi_{j,h}^{n+1})|.$$

201 The Cauchy-Schwarz and Young's inequalities applied to the RHS terms of (19) give,
 202 for any positive constants β_i ($i = 0, \dots, 3$), that

$$203 \quad |((a_j^{n+1} - \bar{a}^{n+1})\nabla\phi_{j,h}^n, \nabla\phi_{j,h}^{n+1})| \leq |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\frac{\|\nabla\phi_{j,h}^n\|^2}{2} + \frac{\|\nabla\phi_{j,h}^{n+1}\|^2}{2} \right),$$

$$204 \quad |(\bar{a}^{n+1}\nabla\rho_j^{n+1}, \nabla\phi_{j,h}^{n+1})| \leq |\bar{a}^{n+1}|_\infty \left(\frac{\|\nabla\rho_j^{n+1}\|^2}{2\beta_0} + \frac{\beta_0\|\nabla\phi_{j,h}^{n+1}\|^2}{2} \right),$$

$$207 \quad |((a_j^{n+1} - \bar{a}^{n+1})\nabla\rho_j^n, \nabla\phi_{j,h}^{n+1})| \leq |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\frac{\|\nabla\rho_j^n\|^2}{2\beta_1} + \frac{\beta_1\|\nabla\phi_{j,h}^{n+1}\|^2}{2} \right),$$

$$208 \quad |((a_j^{n+1} - \bar{a}^{n+1})\nabla(u_j^{n+1} - u_j^n), \nabla\phi_{j,h}^{n+1})|$$

$$209 \quad \leq |a_j^{n+1} - \bar{a}^{n+1}|_\infty \left(\frac{\|\nabla u_j^{n+1} - \nabla u_j^n\|^2}{2\beta_2} + \frac{\beta_2\|\nabla\phi_{j,h}^{n+1}\|^2}{2} \right),$$

$$210 \quad |(\rho_j^{n+1}, \phi_{j,h}^{n+1})| \leq \left(\frac{\|r_j^{n+1}\|_{-1}^2}{2\beta_3} + \frac{\beta_3\|\nabla\phi_{j,h}^{n+1}\|^2}{2} \right).$$

212 Using the above inequalities in (19) and dropping the non-negative term, $\frac{1}{2\Delta t}\|\phi_{j,h}^{n+1} -$
 213 $\phi_{j,h}^n\|^2$ on the left hand side (LHS), we get

$$214 \quad \frac{1}{2\Delta t}(\|\phi_{j,h}^{n+1}\|^2 - \|\phi_{j,h}^n\|^2) + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2}(\|\nabla\phi_{j,h}^{n+1}\|^2 - \|\nabla\phi_{j,h}^n\|^2)$$

$$215 \quad + \left(\theta - \frac{\beta_0}{2}|\bar{a}^{n+1}|_\infty - \frac{\beta_1 + \beta_2 + 2}{2}|a_j^{n+1} - \bar{a}^{n+1}|_\infty - \frac{\beta_3}{2} \right) \|\nabla\phi_{j,h}^{n+1}\|^2$$

$$216 \quad \leq \left(\frac{|\bar{a}^{n+1}|_\infty}{2\beta_0}\|\nabla\rho_j^{n+1}\|^2 + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2\beta_1}\|\nabla\rho_j^n\|^2 + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2\beta_2}\|\nabla u_j^{n+1} - \nabla u_j^n\|^2 \right.$$

$$217 \quad \left. + \frac{\|r_j^{n+1}\|_{-1}^2}{2\beta_3} \right).$$

218 Selecting $\beta_0 = \frac{\delta|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2|\bar{a}^{n+1}|_\infty}$, $\beta_1 = \beta_2 = \frac{\delta}{2}$, and $\beta_3 = \frac{\delta|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2}$ for some positive
 219 δ , yields

$$220 \quad \frac{1}{2\Delta t}(\|\phi_{j,h}^{n+1}\|^2 - \|\phi_{j,h}^n\|^2) + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{2}(\|\nabla\phi_{j,h}^{n+1}\|^2 - \|\nabla\phi_{j,h}^n\|^2)$$

$$221 \quad + \left[\theta - (1 + \delta)|a_j^{n+1} - \bar{a}^{n+1}|_\infty \right] \|\nabla\phi_{j,h}^{n+1}\|^2$$

$$222 \quad \leq \frac{|\bar{a}^{n+1}|_\infty^2}{\delta|a_j^{n+1} - \bar{a}^{n+1}|_\infty}\|\nabla\rho_j^{n+1}\|^2 + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{\delta}\|\nabla\rho_j^n\|^2$$

$$223 \quad (20) \quad + \frac{|a_j^{n+1} - \bar{a}^{n+1}|_\infty}{\delta}\|\nabla u_j^{n+1} - \nabla u_j^n\|^2 + \frac{\|r_j^{n+1}\|_{-1}^2}{\delta|a_j^{n+1} - \bar{a}^{n+1}|_\infty}.$$

224 Taking $\delta = \frac{\theta - \theta_+}{2\theta_+}$, we have $\theta - (1 + \delta) |a_j^{n+1} - \bar{a}^{n+1}|_\infty > \frac{\theta - \theta_+}{2} > 0$ based on the
 225 stability condition (9) and the upper bound in condition (7). In the last two terms
 226 on the RHS of (20), note that

$$227 \quad \|\nabla u_j^{n+1} - \nabla u_j^n\|^2 = \int_D |\nabla u_j^{n+1} - \nabla u_j^n|^2 dx = \int_D \left| \int_{t_n}^{t_{n+1}} (\nabla u_j)_t dt \right|^2 dx$$

$$228 \quad \leq \Delta t \int_{t_n}^{t_{n+1}} \int_D |(\nabla u_j)_t|^2 dx dt = \Delta t \|\nabla u_{j,t}\|_{L^2(L^2(D); t_n, t_{n+1})}^2.$$

229 By the integral form of Taylor's theorem

$$230 \quad u_j^n = u_j^{n+1} - \Delta t u_{j,t}^{n+1} - \int_{t_n}^{t_{n+1}} u_{j,tt}(\cdot, s)(t_n - s) ds,$$

231 we have

$$232 \quad \|r_j^{n+1}\| = \frac{1}{\Delta t} \left\| \int_{t_n}^{t_{n+1}} u_{j,tt}(\cdot, s)(s - t_n) ds \right\| \leq \int_{t_n}^{t_{n+1}} \|u_{j,tt}(\cdot, s)\| \cdot 1 ds$$

$$233 \quad \leq \left[\int_{t_n}^{t_{n+1}} \|u_{j,tt}(\cdot, s)\|^2 ds \right]^{1/2} \left(\int_{t_n}^{t_{n+1}} 1^2 ds \right)^{1/2}$$

$$234 \quad \leq \sqrt{\Delta t} \|u_{j,tt}\|_{L^2(L^2(D); t_n, t_{n+1})}$$

235 and

$$236 \quad \|r_j^{n+1}\|_{-1}^2 \leq C \|r_j^{n+1}\|^2 \leq C \Delta t \|u_{j,tt}\|_{L^2(L^2(D); t_n, t_{n+1})}^2.$$

237 Substituting these inequalities in (20), considering the uniform bounds of $|a_j - \bar{a}|_\infty$
 238 given in (7), multiplying both sides of (20) by $2\Delta t$, and summing over n , we get

$$239 \quad \|\phi_{j,h}^N\|^2 + (\theta - \theta_+) \Delta t \sum_{n=1}^N \|\nabla \phi_{j,h}^n\|^2 + \theta_- \Delta t \|\nabla \phi_{j,h}^N\|^2$$

$$240 \quad \leq \frac{4\Delta t \theta_+}{\theta - \theta_+} \sum_{n=0}^{N-1} \left(\frac{1}{\theta_-} |\bar{a}^{n+1}|_\infty^2 \|\nabla \rho_j^{n+1}\|^2 + \theta_+ \|\nabla \rho_j^n\|^2 + \theta_+ \Delta t \|\nabla u_{j,t}\|_{L^2(L^2(D); t_n, t_{n+1})}^2 \right.$$

$$241 \quad \left. + \frac{C}{\theta_-} \Delta t \|u_{j,tt}\|_{L^2(L^2(D); t_n, t_{n+1})}^2 \right),$$

242 where we used the assumption that $u_{j,h}^0 = P_h(u_j^0)$, thus $\|\phi_{j,h}^0\| = \|\nabla \phi_{j,h}^0\| = 0$. By the
 243 regularity assumption and standard finite element estimates of the L^2 projection error
 244 in the H^1 norm (see, e.g., Section 4.4 in [2]), namely, for any $u_j^n \in H^{l+1}(D) \cap H_0^1(D)$,

$$245 \quad (21) \quad \|\nabla \rho_j^n\| = \|\nabla(P_h(u_j^n) - u_j^n)\| \leq Ch^{2l} \|u_j^n\|_{l+1}^2,$$

246 we have

$$247 \quad (22) \quad \|\phi_{j,h}^N\|^2 + (\theta - \theta_+) \Delta t \sum_{n=0}^{N-1} \|\nabla \phi_{j,h}^{n+1}\|^2 + \theta_- \Delta t \|\nabla \phi_{j,h}^N\|^2 \leq C(\Delta t^2 + h^{2l}),$$

248 where C is a generic constant independent of the time step Δt and mesh size h . By
 249 the triangle inequality, we have the error estimate (15). \square

250 The proposed ensemble scheme can be easily extended to more general parabolic
 251 equations such as those with random coefficients. Next we discuss the ensemble
 252 solution to unsteady random diffusion equations.

253 **4. The ensemble scheme of parabolic equations with random coeffi-**
 254 **icients.** We consider numerical simulations of an unsteady heat equation for a spa-
 255 tially and temporally varying medium, in the absence of convection. That is, to find
 256 a random function, $u : \Omega \times \overline{D} \times [0, T] \rightarrow \mathbb{R}$ satisfying a.s.,

$$257 \quad (23) \quad \begin{cases} u_t - \nabla \cdot [a(\omega, \mathbf{x}, t)\nabla u] = f(\omega, \mathbf{x}, t), & \text{in } \Omega \times D \times [0, T], \\ u(\omega, \mathbf{x}, t) = g(\omega, \mathbf{x}, t), & \text{on } \Omega \times \partial D \times [0, T], \\ u(\omega, \mathbf{x}, 0) = u^0(\omega, \mathbf{x}), & \text{in } \Omega \times D, \end{cases}$$

258 where D is a bounded Lipschitz domain in \mathbb{R}^d , Ω is the set of outcomes in the complete
 259 probability space, diffusion coefficient a , source term f and boundary condition g :
 260 $\Omega \times D \times [0, T] \rightarrow \mathbb{R}$, and initial condition $u^0: \Omega \times D \rightarrow \mathbb{R}$ are random fields with
 261 continuous and bounded covariance functions.

262 As a first step of the investigations on the ensemble method to random PDEs, we
 263 choose the Monte Carlo method for random sampling because it is nonintrusive, easy
 264 to implement, and its convergence is independent of the dimension of the uncertain
 265 model parameters. However, other methods like Quasi Monte Carlo, Latin Hypercube
 266 Sampling, Stochastic Collocation (see, e.g., [18, 8, 1, 22, 16, 3, 7, 23] and references
 267 therein) are applicable, too. When the Monte Carlo method is applied, a large num-
 268 ber of samples are randomly selected first, then a group of independent simulations
 269 needs to be implemented in order to quantify the underlying stochastic information
 270 of the problem. To improve its computational efficiency, we propose an *ensemble-*
 271 *based Monte-Carlo (EMC) method* for the purpose of uncertainty quantification. The
 272 method consists of the following steps:

- 273 1. Choose a set of random samples for the random medium coefficient, source
 274 term, boundary and initial conditions: $a_j \equiv a(\omega_j, \cdot, \cdot)$, $f_j \equiv f(\omega_j, \cdot, \cdot)$, $g_j \equiv$
 275 $g(\omega_j, \cdot, \cdot)$, and $u_j^0 \equiv u^0(\omega_j, \cdot)$ for $j = 1, \dots, J$. Note that the corresponding
 276 solutions $u(\omega_j, \cdot, \cdot)$ are independent, identically distributed (i.i.d.).
- 277 2. Use the uniform time partition on $[0, T]$ with the step size $\Delta t = T/N$. Define
 278 $u_j^n = u(\omega_j, \mathbf{x}, t_n)$ and $\bar{a}^n = \frac{1}{J} \sum_{j=1}^J a(\omega_j, \mathbf{x}, t_n)$. For $j = 1, \dots, J$ and $n =$
 279 $0, \dots, N - 1$, one finds u_j^{n+1} satisfying the ensemble scheme (4). In practice,
 280 an appropriate finite element space on the mesh \mathcal{T}_h could be chosen, on which
 281 one finds the finite dimensional approximation $u_h(\omega_j, \cdot, \cdot)$.
- 282 3. Approximate $E[u]$ by the EMC sample average $\frac{1}{J} \sum_{j=1}^J u_h(\omega_j, \cdot, \cdot)$. If a quan-
 283 tity of interest $Q(u)$ is given, one analyzes the outputs from the ensemble
 284 simulations, $Q(u_h(\omega_1, \cdot, \cdot)), \dots, Q(u_h(\omega_J, \cdot, \cdot))$, to extract the stochastic in-
 285 formation.

286 It is seen that the EMC method naturally synthesizes the ensemble-based time-
 287 stepping algorithm (4) with the Monte-Carlo random sampling approach. It keeps
 288 the same advantage of the ensemble algorithm when applying to the deterministic
 289 PDEs: all the simulations on the selected samples would share a single coefficient
 290 matrix at each time step, thus one only needs to solve a linear system with multiple
 291 RHS vectors, which leads to the reduction of computational cost. Next, we derive
 292 some numerical analysis for the proposed method.

293 **4.1. Stability and convergence.** Similar to the deterministic case, we con-
 294 sider problems with homogeneous boundary conditions in the following analysis,
 295 which could be extended to the inhomogeneous cases by means of the method of
 296 shifting. Choose the same finite element space V_h as defined in Section 3. Denote
 297 $u_{j,h}^n = u_h(\omega_j, \mathbf{x}, t_n)$. For the j -th ensemble member and for $n = 0, \dots, N - 1$, find an

298 approximation solution $u_{j,h}^{n+1} \in V_h$ such that

$$\begin{aligned}
299 \quad & \left(\frac{u_{j,h}^{n+1} - u_{j,h}^n}{\Delta t}, v_h \right) + (\bar{a}^{n+1} \nabla u_{j,h}^{n+1}, \nabla v_h) + ((a_j^{n+1} - \bar{a}^{n+1}) \nabla u_{j,h}^n, \nabla v_h) \\
300 \quad (24) \quad & = (f_j^{n+1}, v_h), \quad \forall v_h \in V_h
\end{aligned}$$

301 with the initial condition $u_{j,h}^0 \in V_h$ satisfying $(u_{j,h}^0, v_h) = (u_j^0, v_h)$, $\forall v_h \in V_h$. Al-
302 though (24) has the same form as (8), we still present it here because $u_{j,h}^n$ in (24)
303 changes from a real-valued function to a random variable.

304 Suppose the following two conditions are valid:

305 (iii) There exists a positive constant θ such that, for any $t \in [0, T]$,

$$306 \quad (25) \quad P\{\omega \in \Omega; \min_{x \in \bar{D}} a(\omega, \mathbf{x}, t) > \theta\} = 1.$$

307 (iv) There exist positive constants θ_- and θ_+ such that, for any $t \in [0, T]$,

$$308 \quad (26) \quad P\{\omega_j \in \Omega; \theta_- \leq |a(\omega_j, \mathbf{x}, t) - \bar{a}|_\infty \leq \theta_+\} = 1.$$

309 Here, condition (iii) guarantees the uniform coercivity a.s.; condition (iv) gives the
310 uniform bounds of the distance from coefficient $a(\omega_j, \mathbf{x}, t)$ to the ensemble average
311 $\bar{a} = \frac{1}{J} \sum_{j=1}^J a(\omega_j, \mathbf{x}, t)$ a.s.

312 Theorem 1 together with the property of expectation lead to the following stability
313 analysis for the finite element solution $u_{j,h}^n$:

314 THEOREM 4. Suppose $f_j \in \tilde{L}^2(H^{-1}(D); 0, T)$ and conditions (iii) and (iv) are
315 satisfied, the finite element solution $u_{j,h}^n$ to (24) is stable provided

$$316 \quad (27) \quad \theta > \theta_+.$$

317 Especially, for any $\Delta t > 0$, the solution satisfies

$$\begin{aligned}
318 \quad (28) \quad & E [\|u_{j,h}^N\|^2] + \theta_- \Delta t E [\|\nabla u_{j,h}^N\|^2] + (\theta - \theta_+) \Delta t \sum_{n=1}^N E [\|\nabla u_{j,h}^n\|^2] \\
& \leq C \Delta t \sum_{n=1}^N E [\|f_j^n\|_{-1}^2] + C \Delta t E [\|\nabla u_{j,h}^0\|^2] + E [\|u_{j,h}^0\|^2],
\end{aligned}$$

319 where C is a generic constant independent of J , h and Δt .

320 The stability condition (27) restricts the deviation of random diffusion coefficients
321 from the ensemble average. Similar to the deterministic case (see Remark 2), if it does
322 not hold, one might separate the entire ensemble into smaller groups to ensure that
323 (27) is true for each of the small groups, then the EMC method will be applicable to
324 all the groups.

325 The full-discrete EMC approximation is defined to be $\Psi_h^n \equiv \frac{1}{J} \sum_{j=1}^J u_{j,h}^n$. Next,
326 we will derive an estimate for $E[u^n] - \Psi_h^n$ in certain averaged norms. Note that
327 $E[u^n] - \Psi_h^n$ can be naturally split into two parts:

$$\begin{aligned}
328 \quad & E[u^n] - \Psi_h^n = (E[u_j^n] - E[u_{j,h}^n]) + (E[u_{j,h}^n] - \Psi_h^n) \\
329 \quad & = \mathcal{E}_h^n + \mathcal{E}_S^n,
\end{aligned}$$

330 where we use $E[u^n] = E[u_j^n]$ in the first equality. The first part, $\mathcal{E}_h^n = E[u_j^n - u_{j,h}^n]$, is
 331 related to the finite element discretization error controlled by the size of the spatial
 332 triangulation and time step; while the second part, $\mathcal{E}_S^n = E[u_{j,h}^n] - \Psi_h^n$, is the statistical
 333 error controlled by the number of realizations. In the following, we will analyze \mathcal{E}_h^n
 334 in Theorem 5, bound \mathcal{E}_S^n in Theorem 6, and obtain an error estimate of the EMC
 335 approximation in Theorem 7.

336 For \mathcal{E}_h^n , we have the following estimate:

337 **THEOREM 5.** *Let u_j^n be the solution to equation (23) when $\omega = \omega_j$ and $t =$
 338 t_n , and $u_{j,h}^n$ be the solution to (24). Suppose $u_j^0 \in \tilde{L}^2(H_0^1(D) \cap H^{l+1}(D))$, $f_j \in$
 339 $\tilde{L}^2(H^{-1}(D); 0, T)$. Under conditions (iii) and (iv), there exists a generic constant
 340 $C > 0$ independent of J, h and Δt such that*

$$341 \quad E [\|u_j^N - u_{j,h}^N\|^2] + (\theta - \theta_+) \Delta t \sum_{n=1}^N E [\|\nabla(u_j^n - u_{j,h}^n)\|^2]$$

$$342 \quad (29) \quad + \theta_- \Delta t E [\|\nabla(u_j^N - u_{j,h}^N)\|^2] \leq C(\Delta t^2 + h^{2l}),$$

343 provided that the stability condition (27) holds.

344 *Proof.* The conclusion follows Theorem 3 after applying the expectation on (15). \square

345 With the standard error estimate of the Monte Carlo method (e.g., see [15]), the
 346 statistical error \mathcal{E}_S^n can be bounded as follows:

347 **THEOREM 6.** *Suppose conditions (iii) and (iv), and the stability condition (27)*
 348 *hold, $f_j \in \tilde{L}^2(H^{-1}(D); 0, T)$ and $u_j^0 \in \tilde{L}^2(H_0^1(D) \cap H^{l+1}(D))$, then there is a generic*
 349 *constant $C > 0$ independent of J, h and Δt such that*

$$350 \quad (30) \quad E [\|\mathcal{E}_S^N\|^2] + \theta_- \Delta t E [\|\nabla \mathcal{E}_S^N\|^2] + (\theta - \theta_+) \Delta t \sum_{n=1}^N E [\|\nabla \mathcal{E}_S^n\|^2]$$

$$\leq \frac{C}{J} \left(\Delta t \sum_{n=1}^N E [\|f_j^n\|_{-1}^2] + \Delta t E [\|\nabla u_{j,h}^0\|^2] + E [\|u_{j,h}^0\|^2] \right).$$

351 *Proof.* We first estimate $E[\|\nabla \mathcal{E}_S^n\|]$, define $\langle u_h^n, u_h^n \rangle := (\nabla u_h^n, \nabla u_h^n)$, then we have

$$352 \quad E [\|\nabla \mathcal{E}_S^n\|^2] = E \left[\left\langle \frac{1}{J} \sum_{i=1}^J (E[u_h^n] - u_{i,h}^n), \frac{1}{J} \sum_{j=1}^J (E[u_h^n] - u_{j,h}^n) \right\rangle \right]$$

$$353 \quad = \frac{1}{J^2} \sum_{i=1}^J \sum_{j=1}^J E [\langle E[u_h^n] - u_{i,h}^n, E[u_h^n] - u_{j,h}^n \rangle]$$

$$354 \quad = \frac{1}{J^2} \sum_{j=1}^J E [\langle E[u_h^n] - u_{j,h}^n, E[u_h^n] - u_{j,h}^n \rangle].$$

The last equality is due to the fact that $u_h^n(\omega_1, \cdot), \dots, u_h^n(\omega_J, \cdot)$ are i.i.d., and thus the
 expected value of $\langle E[u_h^n] - u_{i,h}^n, E[u_h^n] - u_{j,h}^n \rangle$ is a zero for $i \neq j$. We now expand
 the quantity $\langle E[u_h^n] - u_{j,h}^n, E[u_h^n] - u_{j,h}^n \rangle$ and use the fact that $E[u_h^n] = E[u_{j,h}^n]$ and
 $E[(u_h^n)^2] = E[(u_{j,h}^n)^2]$ to obtain

$$E [\|\nabla \mathcal{E}_S^n\|^2] = -\frac{1}{J} \|\nabla E[u_{j,h}^n]\|^2 + \frac{1}{J} E [\|\nabla u_{j,h}^n\|^2].$$

Therefore, we have

$$E [\|\nabla \mathcal{E}_S^n\|^2] \leq \frac{1}{J} E [\|\nabla u_{j,h}^n\|^2].$$

355 By Theorem 4, we get

$$\begin{aligned} 356 \quad (\theta - \theta_+) \Delta t \sum_{n=1}^N E [\|\nabla u_{j,h}^n\|^2] &\leq C \Delta t \sum_{n=1}^N E [\|f_j^n\|_{-1}^2] \\ 357 \quad &+ C \Delta t E [\|\nabla u_{j,h}^0\|^2] + E [\|u_{j,h}^0\|^2]. \end{aligned}$$

358 The other terms on the LHS of (30) involving $E[\|\mathcal{E}_S^N\|^2]$ and $E[\|\nabla \mathcal{E}_S^N\|^2]$ can be
359 treated in the same manner. This completes the proof. \square

360 The combination of error contributions from the finite element approximation
361 and Monte Carlo sampling yields a bound for the EMC approximation error in the
362 following sense:

363 **THEOREM 7.** *For the given source function $f_j \in \tilde{L}^2(H^{-1}(D); 0, T)$ and $u_j^0 \in$
364 $\tilde{L}^2(H_0^1(D) \cap H^{l+1}(D))$. Under conditions (iii) and (iv), and suppose the stability
365 condition (27) is satisfied, that is, $\theta > \theta_+$, then there holds*

$$\begin{aligned} 366 \quad E [\|E[u^N] - \Psi_h^N\|^2] + \theta_- \Delta t E [\|\nabla(E[u^N] - \Psi_h^N)\|^2] \\ 367 \quad + (\theta - \theta_+) \Delta t \sum_{n=1}^N E [\|\nabla(E[u^n] - \Psi_h^n)\|^2] \\ 368 \quad \leq \frac{1}{J} \left(C \Delta t \sum_{n=1}^N E [\|f_j^n\|_{-1}^2] + C \Delta t E [\|\nabla u_{j,h}^0\|^2] + E [\|u_{j,h}^0\|^2] \right) \\ 369 \quad (31) \quad + C(\Delta t^2 + h^{2l}), \end{aligned}$$

370 where $C > 0$ is a constant independent of J , h and Δt .

371 *Proof.* Consider the first term on the LHS of (31). By the triangle and Young's
372 inequality, we have

$$373 \quad E [\|E[u^N] - \Psi_h^N\|^2] \leq 2 (E [\|E[u^N] - E[u_h^N]\|^2] + E [\|E[u_h^N] - \Psi_h^N\|^2]).$$

374 Applying Jensen's inequality to the first term on the RHS of the above inequality, we
375 have

$$376 \quad E [\|E[u^N] - E[u_h^N]\|^2] \leq E [E[\|u^N - u_h^N\|^2]] = E [\|u^N - u_h^N\|^2].$$

377 Then the conclusion follows from Theorems 5-6. The other terms on the LHS of (31)
378 can be estimated in a similar manner. \square

379 **5. Numerical experiments.** We present two numerical tests on the ensemble
380 schemes for second-order parabolic PDEs in this section: the first problem is deter-
381 ministic heat transfer with an *a priori* known exact solution, which aims to illustrate
382 Theorem 3; the second problem is random heat transfer without a known exact solu-
383 tion, which is used to illustrate Theorem 7 and shows the effectiveness of the ensemble
384 method by comparing the results with those of independent, individual simulations.

385 **5.1. Deterministic heat transfer.** We first test the numerical performance of
 386 the ensemble algorithm on the deterministic second-order parabolic equation (2). A
 387 group of simulations is considered, which contains $J = 3$ members. The diffusion
 388 coefficient and the exact solution of j -th simulation are selected as follows.

$$389 \quad \begin{aligned} a_j(\mathbf{x}, t) &= 1 + (1 + \epsilon_j) \sin(t) \sin(xy), \\ u_j(\mathbf{x}, t) &= (1 + \epsilon_j) [\sin(2\pi x) \sin(2\pi y) + \sin(4\pi t)], \end{aligned}$$

390 where ϵ_j is a perturbation randomly selected from $[0, 1]$, $t \in [0, 1]$ and $(x, y) \in [0, 1]^2$.
 391 The initial condition, Dirichlet boundary condition and source term are chosen to
 392 match the exact solution.

The group is simulated by using the ensemble scheme (8). In the test, the ensemble contains three members with $\epsilon_1 = 0.6207$, $\epsilon_2 = 0.1841$, and $\epsilon_3 = 0.2691$. In order to check the convergence order in time, we use quadratic finite elements, a uniform time partition, and uniformly refine the mesh size h and time step size Δt from the initial mesh size $\sqrt{2}/4$ and initial time step size $1/10$. Let the maximum numerical approximation errors in L^2 norm be

$$\mathcal{E}_{L^2}^j = \max_{n \in \{1, \dots, N\}} \|u_j^n - u_{j,h}^n\|$$

and errors in the time average H^1 semi-norm be

$$\mathcal{E}_{H^1}^j = \sqrt{\Delta t \sum_{n=1}^N \|\nabla u_j^n - \nabla u_{j,h}^n\|^2}$$

393 for $j = 1, 2, 3$, respectively. The approximation errors of the ensemble method (de-
 394 noted by $\mathcal{E}_{L^2}^{E,j}$ and $\mathcal{E}_{H^1}^{E,j}$) are listed in Table 1. It is seen that the rate of convergence
 395 is nearly linear, which matches our theoretical analysis in Theorem 3.

Table 1: Numerical errors of the ensemble simulations.

$\sqrt{2}/h$	$\mathcal{E}_{L^2}^{E,1}$	rate	$\mathcal{E}_{L^2}^{E,2}$	rate	$\mathcal{E}_{L^2}^{E,3}$	rate
4	2.2271×10^{-1}	-	2.2168×10^{-1}	-	2.2177×10^{-1}	-
8	1.1477×10^{-1}	0.96	1.1623×10^{-1}	0.93	1.1594×10^{-1}	0.94
16	5.9080×10^{-2}	0.96	5.9921×10^{-2}	0.96	5.9756×10^{-2}	0.96
32	3.0007×10^{-2}	0.98	3.0445×10^{-2}	0.98	3.0359×10^{-2}	0.98
$\sqrt{2}/h$	$\mathcal{E}_{H^1}^{E,1}$	rate	$\mathcal{E}_{H^1}^{E,2}$	rate	$\mathcal{E}_{H^1}^{E,3}$	rate
4	1.3678×10^0	-	1.0922×10^0	-	1.1437×10^0	-
8	4.7311×10^{-1}	1.53	4.2423×10^{-1}	1.36	4.3280×10^{-1}	1.40
16	1.9969×10^{-1}	1.24	1.9560×10^{-1}	1.12	1.9618×10^{-1}	1.14
32	9.5767×10^{-2}	1.06	9.6972×10^{-2}	1.01	9.6692×10^{-2}	1.02

396 To compare with the individual simulations, we list in Table 2 the numerical errors
 397 of independent simulations (denoted by $\mathcal{E}_{L^2}^{I,j}$ and $\mathcal{E}_{H^1}^{I,j}$) in the same computational
 398 setting. It is observed that the ensemble simulation results in Table 1 are close to
 399 those obtained from individual simulations in Table 2. Indeed, the errors are at the
 400 same order of magnitude and the convergence rates are almost same.

401 **5.2. Random heat transfer.** Next we consider the second-order parabolic
 402 equation with a random diffusion coefficient (23) on the unit square domain. The
 403 test problem is associated with the zero forcing term f , zero initial conditions, and

Table 2: Numerical errors of the independent simulations.

$\sqrt{2}/h$	$\mathcal{E}_{L^2}^{I,1}$	rate	$\mathcal{E}_{L^2}^{I,2}$	rate	$\mathcal{E}_{L^2}^{I,3}$	rate
4	2.2206×10^{-1}	-	2.2215×10^{-1}	-	2.2200×10^{-1}	-
8	1.1469×10^{-1}	0.95	1.1629×10^{-1}	0.93	1.1597×10^{-1}	0.94
16	5.9072×10^{-2}	0.96	5.9928×10^{-2}	0.96	5.9759×10^{-2}	0.96
32	3.0007×10^{-2}	0.98	3.0446×10^{-2}	0.98	3.0359×10^{-2}	0.98
$\sqrt{2}/h$	$\mathcal{E}_{H^1}^{I,1}$	rate	$\mathcal{E}_{H^1}^{I,2}$	rate	$\mathcal{E}_{H^1}^{I,3}$	rate
4	1.3641×10^0	-	1.0955×10^0	-	1.1453×10^0	-
8	4.7186×10^{-1}	1.53	4.2529×10^{-1}	1.37	4.3331×10^{-1}	1.40
16	1.9933×10^{-1}	1.24	1.9588×10^{-1}	1.12	1.9632×10^{-1}	1.14
32	9.5677×10^{-2}	1.06	9.7041×10^{-2}	1.01	9.6726×10^{-2}	1.02

404 homogeneous Dirichlet boundary conditions on the top, bottom and right edges of
 405 the domain but nonhomogeneous Dirichlet boundary condition, $u = y(1 - y)$, on the
 406 left edge. The random coefficient varies in the vertical direction and has the following
 407 form

408 (32) $a(\omega, \mathbf{x}) = a_0 + \sigma \sqrt{\lambda_0} Y_0(\omega) + \sum_{i=1}^{n_f} \sigma \sqrt{\lambda_i} [Y_i(\omega) \cos(i\pi y) + Y_{n_f+i}(\omega) \sin(i\pi y)]$

409 with $\lambda_0 = \frac{\sqrt{\pi L_c}}{2}$, $\lambda_i = \sqrt{\pi L_c} e^{-\frac{(i\pi L_c)^2}{4}}$ for $i = 1, \dots, n_f$ and Y_0, \dots, Y_{2n_f} are uncorre-
 410 lated random variables with zero mean and unit variance. In the following numerical
 411 test, we take $a_0 = 1$, $L_c = 0.25$, $\sigma = 0.15$, $n_f = 3$ and assume the random variables
 412 Y_0, \dots, Y_{2n_f} are independent and uniformly distributed in the interval $[-\sqrt{3}, \sqrt{3}]$.
 413 We use linear finite elements for spatial discretization and simulate the system over
 414 the time interval $[0, 0.5]$. This choice of final time guarantees a steady-state can be
 415 achieved at the end of simulations. A similar computational setting is used in [19]
 416 to compare several numerical methods for parabolic equations with random coeffi-
 417 cients. In the following tests, the uniform triangulation with the maximum mesh size
 418 $h = \sqrt{2}/32$ and the uniform time partition with the time step size $\Delta t = 2.5 \times 10^{-3}$
 419 are used.

For the implementation of the EMC method as discussed in Section 4, we first select a set of J random samples by the MC sampling, then run our deterministic code for simulating the ensemble of the deterministic PDEs associated with the J realizations. Since the numerical accuracy with respect to the mesh size and time step size for the deterministic case has been verified in the first example, here we only check the rate of convergence in the EMC approximation error with respect to the number of samples, J . As the exact solution is unknown, we choose the EMC solution using $J_0 = 5000$ samples as our benchmark, vary the values of J in the EMC simulations, and then evaluate the approximation errors based on the benchmark. Furthermore, we repeat such error analysis for $M = 10$ independent replicas and compute the average of the output errors. Denote the EMC solution at time t_n in the m -th independent replica by $\Psi_{J,h}^{n,m} = \frac{1}{J} \sum_{j=1}^J u_{j,h}^{n,m}$, where $u_{j,h}^{n,m}$ is the output of the ensemble scheme (8) in the m -th experiment. Define

$$\mathcal{E}_{L^2} = \max_{n \in \{1, \dots, N\}} \sqrt{\frac{1}{M} \sum_{m=1}^M \|\Psi_{J_0,h}^{n,m} - \Psi_{J,h}^{n,m}\|^2},$$

$$\mathcal{E}_{H^1} = \sqrt{\frac{\Delta t}{M} \sum_{m=1}^M \sum_{n=1}^N \|\nabla \Psi_{J_0,h}^{n,m} - \nabla \Psi_{J,h}^{n,m}\|^2}.$$

420 The numerical results at $J = 10, 20, 40, 80, 160$ are listed in Table 3. We further apply
 421 linear regression analysis on these numerical results, which shows $\mathcal{E}_{L^2} \approx 0.0032 J^{-0.5133}$
 422 and $\mathcal{E}_{H^1} \approx 0.0104 J^{-0.4877}$. The values of \mathcal{E}_{L^2} and \mathcal{E}_{H^1} together with their linear re-
 423 gression models are plotted in Figure 1, respectively. It is seen that the rate of convergence
 424 with respect to J is close to -0.5 , which coincides with our theoretical
 425 results in Theorem 7.

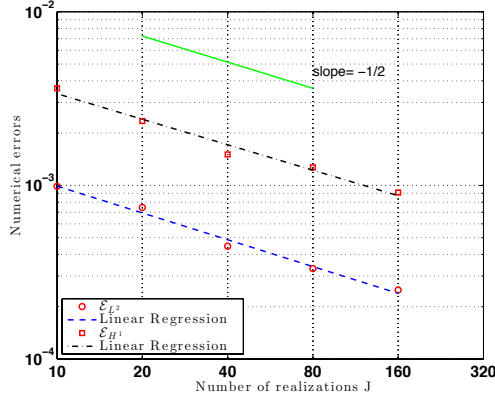


Fig. 1: Ensemble simulations: errors converge on the order of $\mathcal{O}(1/\sqrt{J})$.

Table 3: Numerical errors of ensemble simulations.

J	10	20	40	80	160
\mathcal{E}_{L^2}	9.8672e-04	7.4471e-04	4.4640e-04	3.3151e-04	2.4966e-04
\mathcal{E}_{H^1}	3.6240e-03	2.3475e-03	1.5080e-03	1.2696e-03	9.0903e-04

Next, we analyze some stochastic information of the system including the expectation of u at the final time and a quantity of interest. In particular, we are mainly interested in comparing the ensemble simulation outputs with those from the individual simulations when the same set of samples is used. More precisely, we approximate the expected value $E[u(\omega, \mathbf{x}, T)]$ by the EMC approximation

$$\Psi_h^E(\mathbf{x}) := \frac{1}{J} \sum_{j=1}^J u_h^E(\omega_j, \mathbf{x}, T),$$

where $u_h^E(\omega_j, \mathbf{x}, T)$ is the j -th member solution in the ensemble-based simulation at time T . Taking the number of sample points to be $J = 5000$, we compute the mean and standard deviation of the solutions at the final time, which are plotted in Figure 2 (left and middle). To quantify the performance of the EMC method, we compare the result with that of individual finite element Monte Carlo (FEMC) simulations using

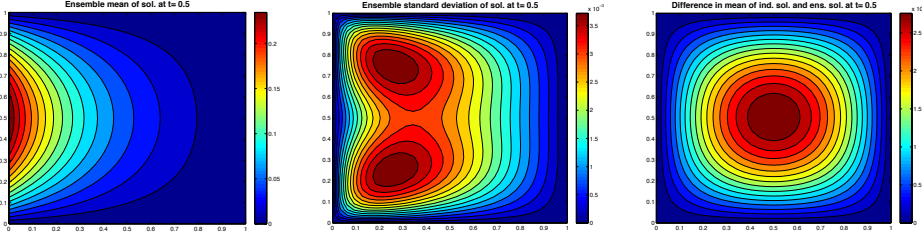


Fig. 2: Ensemble simulations: mean (left), standard deviation (middle) of solution at $t = 0.5$, difference of the mean from that of the individual FEMC simulations (right).

the same set of sample values. Denote the FEMC approximation solution by

$$\Psi_h^I(\mathbf{x}) := \frac{1}{J} \sum_{j=1}^J u_h^I(\omega_j, \mathbf{x}, T),$$

426 where $u_h^I(\omega_j, \mathbf{x}, T)$ is the j -th independent solution at time T . The difference between
 427 Ψ_h^E and Ψ_h^I is shown in Figure 2 (right). It is observed that the difference is on the
 428 order of 10^{-7} , which indicates the EMC method is able to provide the same accurate
 429 approximation as the individual FEMC simulations. However, the CPU time for the
 430 ensemble simulation is 1.4494×10^4 seconds, while that of the individual simulations
 431 is 4.9197×10^4 seconds. The former improves the computational efficiency by about
 432 70%. Here, because the size of discrete system is small, we use MATLAB and its LU
 433 matrix factorization for solving the linear system.

Following [19], we also calculate the quantity of interest

$$Q(\omega) = \int_D u(\omega, \mathbf{x}, T) dD.$$

434 When the sample size is $J = 5000$, the histogram of the quantity of interest obtained
 435 from the ensemble simulations, Q_h^E , is shown in Figure 3 (left) with a fitted gaussian
 436 distribution. We then do a comparison with the quantity of interest, Q_h^I , achieved
 437 from individual simulations at the same sampling set. The histogram of the differences
 438 in absolute value, $|Q_h^E - Q_h^I|$, is plotted in Figure 3 (right), which also illustrates that
 439 the EMC method outputs a close quantity of interest to the standard FEMC method.

440 **6. Conclusion.** We propose an ensemble-based algorithm in this paper to im-
 441 prove the computational efficiency for a group of numerical solutions to parabolic
 442 problems. The fundamental idea is to turn the linear systems associated to the group
 443 into a linear system with multiple right-hand-side vectors, which would reduce the
 444 computational time. We first analyze the ensemble scheme for deterministic equa-
 445 tions, then develop the ensemble-based Monte Carlo method for stochastic equations.
 446 The effectiveness of both cases is demonstrated through rigorous error estimates and
 447 illustrated with numerical experiments. The approach can be easily extended to more
 448 general, nonlinear parabolic equations, which is one of the research directions we are
 449 pursuing.

450 **Acknowledgments.** The first author would like to thank the support of China
 451 Scholarship Council for visiting the Interdisciplinary Mathematics Institute at Uni-
 452 versity of South Carolina during the year 2016–2017. We also thank the anonymous
 453 referees for their comments and suggestions, which significantly improved this paper.

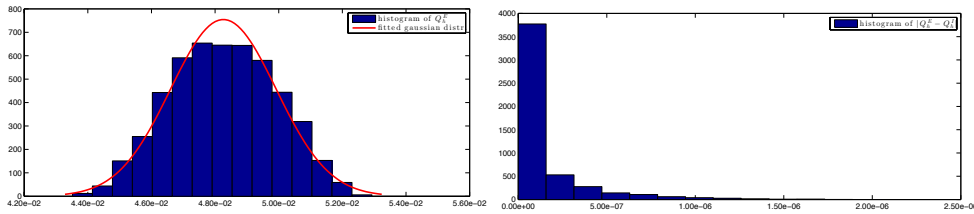


Fig. 3: Ensemble simulations: histogram of the quantity of interest Q_h^E with a fitted gaussian distribution (left); histogram of the difference of the quantity in absolute value, $|Q_h^E - Q_h^I|$, between ensemble simulations and individual simulations (right).

454

REFERENCES

- 455 [1] I. BABUŠKA, R. TEMPONE, AND G. E. ZOURARIS, *Solving elliptic boundary value problems with*
456 *uncertain coefficients by the finite element method: the stochastic formulation*, Computer
457 methods in applied mechanics and engineering, 194 (2005), pp. 1251–1294.
- 458 [2] S. BRENNER AND R. SCOTT, *The mathematical theory of finite element methods*, vol. 15,
459 Springer Science & Business Media, New York, USA, 2007.
- 460 [3] B. GANAPATHYSUBRAMANIAN AND N. ZABARAS, *Sparse grid collocation schemes for stochastic*
461 *natural convection problems*, Journal of Computational Physics, 225 (2007), pp. 652–685.
- 462 [4] M. GUNZBURGER, N. JIANG, AND M. SCHNEIER, *An ensemble-proper orthogonal decomposi-*
463 *tion method for the nonstationary navier–stokes equations*, SIAM Journal on Numerical
464 Analysis, 55 (2017), pp. 286–304.
- 465 [5] M. GUNZBURGER, N. JIANG, AND M. SCHNEIER, *A higher-order ensemble/proper orthogonal*
466 *decomposition method for the nonstationary navier–stokes equations*, (in press).
- 467 [6] M. GUNZBURGER, N. JIANG, AND Z. WANG, *An efficient algorithm for simulating ensembles of*
468 *parameterized flow problems*, (submitted).
- 469 [7] M. D. GUNZBURGER, C. G. WEBSTER, AND G. ZHANG, *Stochastic finite element methods for*
470 *partial differential equations with random input data*, Acta Numerica, 23 (2014), pp. 521–
471 650.
- 472 [8] J. C. HELTON AND F. J. DAVIS, *Latin hypercube sampling and the propagation of uncertainty in*
473 *analyses of complex systems*, Reliability Engineering & System Safety, 81 (2003), pp. 23–69.
- 474 [9] N. JIANG, *A higher order ensemble simulation algorithm for fluid flows*, Journal of Scientific
475 Computing, 64 (2015), pp. 264–288.
- 476 [10] N. JIANG, *A second-order ensemble method based on a blended backward differentiation formula*
477 *timestepping scheme for time-dependent navier–stokes equations*, Numerical Methods for
478 Partial Differential Equations, 33 (2017), pp. 34–61.
- 479 [11] N. JIANG, S. KAYA, AND W. LAYTON, *Analysis of model variance for ensemble based turbulence*
480 *modeling*, Computational Methods in Applied Mathematics, 15 (2015), pp. 173–188.
- 481 [12] N. JIANG AND W. LAYTON, *An algorithm for fast calculation of flow ensembles*, International
482 Journal for Uncertainty Quantification, 4 (2014).
- 483 [13] N. JIANG AND W. LAYTON, *Numerical analysis of two ensemble eddy viscosity numerical reg-*
484 *ularizations of fluid motion*, Numerical Methods for Partial Differential Equations, 31
485 (2015), pp. 630–651.
- 486 [14] E. KALNAY, *Atmospheric modeling, data assimilation and predictability*, Cambridge University
487 Press, New York, USA, 2003.
- 488 [15] K. LIU AND B. M. RIVIÈRE, *Discontinuous galerkin methods for elliptic partial differential*
489 *equations with random coefficients*, International Journal of Computer Mathematics, 90
490 (2013), pp. 2477–2490.
- 491 [16] L. MATHELIN, M. Y. HUSSAINI, AND T. A. ZANG, *Stochastic approaches to uncertainty quan-*
492 *tification in CFD simulations*, Numerical Algorithms, 38 (2005), pp. 209–236.
- 493 [17] M. MOHEBUJJAMAN AND L. G. REBHOLZ, *An efficient algorithm for computation of mhd flow*
494 *ensembles*, Computational Methods in Applied Mathematics, 17 (2017), pp. 121–137.
- 495 [18] H. NIEDERREITER, *Random number generation and quasi-Monte Carlo methods*, SIAM
496 Philadelphia, PA, USA, 1992.
- 497 [19] F. NOBILE AND R. TEMPONE, *Analysis and implementation issues for the numerical approxi-*

- 498 *mation of parabolic equations with random coefficients*, International journal for numerical
499 methods in engineering, 80 (2009), pp. 979–1006.
- 500 [20] M. L. PARKS, K. M. SOODHALTER, AND D. B. SZYLD, *A block recycled gmres method with*
501 *investigations into aspects of solver performance*, arXiv preprint arXiv:1604.01713, (2016).
- 502 [21] A. TAKHIROV, M. NEDA, AND J. WATERS, *Time relaxation algorithm for flow ensembles*, Nu-
503 merical Methods for Partial Differential Equations, 32 (2016), pp. 757–777.
- 504 [22] D. XIU AND J. S. HESTHAVEN, *High-order collocation methods for differential equations with*
505 *random inputs*, SIAM Journal on Scientific Computing, 27 (2005), pp. 1118–1139.
- 506 [23] X. ZHU, E. M. LINEBARGER, AND D. XIU, *Multi-fidelity stochastic collocation method for com-*
507 *putation of statistical moments*, Journal of Computational Physics, 341 (2017), pp. 386–
508 396.