1  We thank the reviewers for their constructive suggestions and insightful comments. We have (1) added simulations for
2  Bayesian Causal Forests and (2) have substantially expanded the discussion in the final version to address the various
3  reviewer comments. A summary of the added discussion is provided point-by-point below.

4  **Reviewer 1: $\pi$ is not assigned a prior.** We use an empirical Bayes approach, which is, as the reviewer points out,
5  computationally much faster than a fully hierarchical Bayesian approach of placing a prior on $\pi$. For Gaussian
6  processes (GPs), since our approach reduces to modifying the prior covariance, the posterior can be computed using
7  standard computational tools for GPs. Conversely, one can only sample from the hierarchical Bayes posterior using a
8  Metropolis-Hastings-within-Gibbs-sampling algorithm, which is far slower in practice.

9  $\pi$ **requires a prior.** Please note that we consider *independent* priors on $(m, \pi, F)$ (for computational reasons), so that
10 the posterior also factorizes and hence the $\pi$ term is conditionally independent of $\psi$ (see lines 94-107 and Appendix A).

11 **Comment on Table 1.** The reason for our high coverage is that our posterior bias, due to our explicit ATE bias
12 correction (4), is (much) smaller than the posterior variance by design. As the reviewer correctly points out, asymptotic
13 theory predicts the frequentist coverage should converge to 0.95 as the sample size increases due to the semiparametric
14 Bernstein-von Mises theorem (cf Ray & van der Vaart (2018)). However, it is a subtle question as to when the asymptotic
15 regime applies and our examples seem insufficiently data rich for this to be the case (e.g. $n = 1000$ observations but
16 $d = 100$ input features).

17 **Missing organizational section/reference for bias/language precision.** We have incorporated these suggestions.

18 **Reviewer 2: Variable selection for causal inference.** We consider here a GP with squared exponential kernel with
19 automatic relevance determination (ARD), i.e. whose data-driven lengthscale $\ell_i$ represents the relevance of the $ith$
20 feature to the response surface. ARD has been used successfully for removing irrelevant inputs by several authors (see
21 e.g. Chapter 5.1 of Rasmussen & Williams (2006)) and can thus be viewed as a form of automatic (causal) feature
22 selection. Diagnosing missing significant covariates (confounders) is an important and difficult problem which requires
23 further investigation in the future.

24 **Reviewer 3: Comparison with Hahn et al.** We would firstly like to clarify that our goal is to improve estimation of
25 the *average treatment effect* (ATE) in the *presence* of heterogeneous treatment effects. It is known that naively using
26 product priors in casual inference/missing data problems can yield biased inference [this goes back to at least Robins
27 and Ritov (1997). The 'regularization-induced confounding' of Hahn et al. is a very nice illustration of similar ideas
28 for the concrete and important cases of linear models and BART priors]. One solution is to reparametrize to force the
29 missing information into the likelihood (e.g. Ritov et al. (2014), Hahn et al. (2018)), while another is to use propensity
30 score (PS) information (Rosenbaum & Rubin (1983)).

31 In a nice paper Hahn et al. (2017) successfully show that this latter idea also helps Bayesian estimation using BART.
32 Their approach is designed to improve nonparametric estimation of the *whole* response surface, which will also lead to
33 some improvement when estimating the ATE. However, it is known that even when the prior is perfectly calibrated (i.e.
34 all tuning parameters are set optimally) and recovers the entire response surface at the optimal rate, the posterior can
35 still induce a bias in the *marginal posterior* for the ATE $\psi$ that prevents efficient estimation and destroys uncertainty
36 quantification (see e.g. Ray & van der Vaart (2018)).

37 The specific form in which we include the PS in our prior (4) is very deliberate - it arises as the 'least favorable direction'
38 of the ATE in semiparametric statistical theory and is specifically designed for estimating the ATE. When either the PS
39 or response surface are especially difficult to estimate, we expect that incorporating the PS as a covariate as in Hahn et
40 al. (2017) will still induce a bias for the ATE (the theory in Ray & van der Vaart (2018) predicts this). In fairness, we
41 wish to emphasize that Hahn et al.'s goal is to estimate the *entire* response surface, for which they provide excellent
42 results, which is a different problem to estimating the ATE we consider here.

43 **Bayesian Causal Forest (BCF) simulations** have been added to the paper. In summary, for estimation BCF performs
44 well on the synthetic data (but moderately worse than our method in both the homogeneous and heterogeneous cases)
45 and excellently on the semi-synthetic data (moderately better than our method). For uncertainty quantification, BCF
46 typically had the shortest credible intervals with suboptimal coverage (80-85%) on the various synthetic datasets and
47 excellent coverage ($\sim$95%) on the semi-snythetic data.

48 **Estimation of $F$.** We use the widely used 'Bayesian bootstrap' (BB) since (1) it is computationally much faster (you
49 need only one costly Cholesky decomposition to generate posterior samples of the ATE whereas with the full Dirichlet
50 process (DP) posterior we require one per posterior draw) and (2) for moderate/large sample sizes it is very close to the
51 true DP posterior. We do not assume that 'one has observed all possible covariates', rather that our covariate samples
52 are representative of the population. If the observed covariates greatly differ from the underlying population distribution
53 then indeed this will not generalize well, but then neither will any prior not involving detailed outside expert information
54 for that particular application.