

# Two steps at a time — taking GAN training in stride with Tseng’s method\*

Axel Böhm<sup>†,‡</sup>, Michael Sedlmayer<sup>†,‡</sup>, Ernő Robert Csetnek<sup>†</sup>, and Radu Ioan Boț<sup>†</sup>

**Abstract.** Motivated by the training of Generative Adversarial Networks (GANs), we study methods for solving minimax problems with additional nonsmooth regularizers. We do so by employing *monotone operator* theory, in particular the *Forward-Backward-Forward (FBF)* method, which avoids the known issue of limit cycling by correcting each update by a second gradient evaluation and does so requiring less projection steps compared to the Extragradient method in the presence of constraints. Furthermore, we propose a seemingly new scheme which recycles old gradients to mitigate the additional computational cost. In doing so we rediscover a known method, related to *Optimistic Gradient Descent Ascent (OGDA)*. For both schemes we prove novel convergence rates for convex-concave minimax problems via a unifying approach. The derived error bounds are in terms of the gap function for the ergodic iterates. For the deterministic and the stochastic problem we show a convergence rate of  $\mathcal{O}(1/k)$  and  $\mathcal{O}(1/\sqrt{k})$ , respectively. We complement our theoretical results with empirical improvements in the training of Wasserstein GANs on the CIFAR10 dataset.

**Key words.** min-max, convex-concave, stochastic gradient, GAN

**AMS subject classifications.** 65K15, 90C15, 90C47

**1. Introduction.** *Generative Adversarial Networks (GANs)* [19] have proven to be a powerful class of generative models, producing for example unseen realistic images. Two neural networks, called generator and discriminator, compete against each other in a game. In the special case of a zero sum game this task can be formulated as a minimax (aka saddle point) problem.

Conventionally, GANs are trained using variants of (stochastic) *Gradient Descent Ascent (GDA)* which are known to exhibit oscillatory behavior [39] and thus fail to converge even for simple bilinear saddle point problems, see [18]. We therefore propose the use of methods with provable convergence guarantees for (stochastic) convex-concave minimax problems, even though GANs are well known to not warrant these properties. Along similar considerations an adaptation of the *Extragradient method (EG)* [29] for the training of GANs was suggested in [15], whereas [11, 12, 31] studied *Optimistic Gradient Descent Ascent (OGDA)* based on *optimistic mirror descent* [48, 49]. We however investigate the *Forward-Backward-Forward (FBF)* method [55] from monotone operator theory, which uses two gradient evaluations per update, similar to EG, in order to circumvent the aforementioned issues but requires less projection/proximal steps per iteration.

Instead of trying to improve GAN performance with new architectures, loss functions, etc., we want to improve their training by contributing to the study of minimax problems. While the landscape of GAN training is far from matching the rigorous setting of monotonicity, the nonconvex-nonconcave setting remains either intractable in its full generality [13] or other simplifying assumption or other simplifying assumptions have to be made, for example the existence of Minty solutions [37, 32], which

---

<sup>‡</sup>Research Platform Data Science @ Uni Vienna

\*Submitted to the editors March 22, 2022.

**Funding:** This project has received funding from the doctoral programme *Vienna Graduate School on Computational Optimization (VGSCO)*, FWF (Austrian Science Fund), project W 1260, as well as project P 29809-N32.

<sup>†</sup>Faculty of Mathematics, University of Vienna, Austria ({axel.boehm, michael.sedlmayer, robert.csetnek, radu.bot}@univie.ac.at)

36 are again related to monotonicity.

37 We also want to point out that while extrapolation / optimistic steps are able to combat some of the  
 38 oscillatory behaviour of minimax problems, another set of problems arises from the stochastic noise  
 39 of the gradient evaluations which can be dealt with variance reduction techniques [8, 26] or reducing  
 40 the stepsize (as we do). Other considerations have proposed to use the same sample in the for the two  
 41 gradient evaluations leading to improved empirical performance [40]. Another approach was explored  
 42 in [25] to use very different stepsizes for the descent and ascent steps respectively (something that has  
 43 been observed to be very beneficial in practice as well). Known methods were outfitted with *negative*  
 44 *momentum* to improve stability in [16]. An additional technique proposed to deal with the oscillatory  
 45 behaviour of minimax problems called *crossing-the-curl* was suggested in [14], whose authors used  
 46 second order information to take a step perpendicular to the direction of rotation.

47 While our convergence results are stated in terms of the averaged iterates, a technique which can  
 48 prove beneficial in practice [9, 15], having guarantees on the last iterate would be more in the spirit  
 49 of nonconvex methods. Such results have been obtained in [11] and [8], but are for bilinear and  
 50 strongly-convex-strongly-concave problems, respectively, which are known to allow for the derivation  
 51 of better convergence rates. On the other hand the works [17, 33, 20] have been able to guarantee  
 52 rates of convergence for the last iterate for convex-concave problems, but all of them only in the  
 53 unconstrained setting. Furthermore [17] uses additional smoothness assumptions, [33] employs a  
 54 second-order method and [20] uses the norm of the gradient as the measure of optimality.

55 **Contribution.** Establishing the connection between GAN training and *monotone inclusions* moti-  
 56 vates to use the FBF method, originally designed to solve this type of problems. This approach allows  
 57 to naturally extend the constrained setting to a regularized one making use of the proximal operator.

58 We also propose a variant of FBF reusing previous gradients to reduce the computational cost per  
 59 iteration, which turns out to be a known method, related to OGDA. By developing a unifying scheme  
 60 that captures FBF and a generalization of OGDA, we reveal a hitherto unknown connection. Using  
 61 this approach we prove novel nonasymptotic convergence statements in terms of the minimax gap  
 62 for both methods in the context of saddle point problems. In the deterministic and stochastic setting  
 63 we obtain rates of  $\mathcal{O}(1/k)$  and  $\mathcal{O}(1/\sqrt{k})$ , respectively. Concluding, we highlight the relevance of our  
 64 proposed method as well as the role of regularizers by showing empirical improvements in the training  
 65 of Wasserstein GANs on the CIFAR10 dataset.

66 **Organization.** This paper is structured as follows. In [section 2](#) we highlight the connection of  
 67 GAN training and monotone inclusions and give an extensive review of methods with convergence  
 68 guarantees for the latter. The main results as well as a precise definition of the measure of optimality  
 69 are discussed in [section 3](#). Concluding, [section 4](#) illustrates the empirical performance in the training  
 70 of GANs as well as solving bilinear problems.

71 **2. GAN training as monotone inclusion.** The GAN objective was originally cast as a two-  
 72 player zero-sum game between the discriminator  $D_y$  and the generator  $G_x$  [19] given by

$$73 \min_x \max_y \mathbb{E}_{\rho \sim q} [\log(D_y(\rho))] + \mathbb{E}_{\zeta \sim p} [\log(1 - D_y(G_x(\zeta)))],$$

74 exhibiting the aforementioned minimax structure. Due to problems with vanishing gradients in the  
 75 training of such models, a successful alternative formulation called *Wasserstein GAN (WGAN)* [1] has  
 76 been proposed. In this case the minimization tries to reduce the Wasserstein distance between the true

77 distribution  $q$  and the one learned by the generator. Reformulating this distance via the Kantorovich  
 78 Rubinstein duality leads to an inner maximization over 1-Lipschitz functions which are approximated  
 79 via neural networks, yielding the saddle point problem

$$80 \quad \min_x \max_{y: \|D_y\|_{\text{Lip}} \leq 1} \mathbb{E}_{\rho \sim q}[D_y(\rho)] - \mathbb{E}_{\zeta \sim p}[D_y(G_x(\zeta))].$$

81 **2.1. Convex-concave minimax problems.** Due to the observations made in the previous  
 82 paragraph we study the following abstract minimax problem

$$83 \quad (2.1) \quad \min_{x \in \mathbb{R}^d} \max_{y \in \mathbb{R}^n} \Psi(x, y) := f(x) + \mathbb{E}_{\xi \sim Q} [\Phi(x, y; \xi)] - h(y),$$

84 where the convex-concave coupling function  $\Phi(x, y) := \mathbb{E}_{\xi \sim Q} [\Phi(x, y; \xi)]$ , which hides the stochas-  
 85 ticity for ease of notation, is differentiable with  $L$ -Lipschitz continuous gradient. The proper, convex  
 86 and lower semicontinuous functions  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  and  $h: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  act as regulariz-  
 87 ers. A solution of (2.1) is given by a so-called *saddle point*  $(x^*, y^*)$  fulfilling for all  $x$  and  $y$

$$88 \quad (2.2) \quad \Psi(x^*, y) \leq \Psi(x^*, y^*) \leq \Psi(x, y^*).$$

89 In the context of two-player games this corresponds to a pair of strategies, where no player can be  
 90 better off by changing just their own strategy.

91 For the purpose of this motivating section, we will restrict ourselves for now to the special case of  
 92 the *deterministic constrained* version of (2.1), given by

$$93 \quad (2.3) \quad \min_{x \in X} \max_{y \in Y} \Phi(x, y),$$

94 where  $f$  and  $h$  are given by indicator functions of closed convex sets  $X$  and  $Y$ , respectively. The  
 95 indicator function  $\delta_C$  of a set  $C$  is defined as  $\delta_C(z) = 0$  for  $z \in C$  and  $\delta_C(z) = +\infty$  otherwise.

96 **2.2. Minimax problems as monotone inclusions.** If the coupling function  $\Phi$  is convex-  
 97 concave and differentiable then solving (2.1) is equivalent to solving the first order optimality condi-  
 98 tions which can be written as a so-called *monotone inclusion* with  $w = (x, y) \in \mathbb{R}^m$  and  $m = d + n$ ,  
 99 given by

$$100 \quad (2.4) \quad 0 \in F(w) + N_\Omega(w).$$

101 The entities involved are

$$102 \quad (2.5) \quad F(x, y) := (\nabla_x \Phi(x, y), -\nabla_y \Phi(x, y)),$$

103 and the *normal cone*  $N_\Omega$  of the convex set  $\Omega := X \times Y$ . The normal cone mapping is given by

$$104 \quad (2.6) \quad N_\Omega(w) = \{v \in \mathbb{R}^m : \langle v, w' - w \rangle \leq 0 \quad \forall w' \in \Omega\},$$

105 for  $w \in \Omega$  and  $N_\Omega(w) = \emptyset$  for  $w \notin \Omega$ . Here, the operators  $F$  and  $N_\Omega$  satisfy well known properties  
 106 from convex analysis [4], in particular the first one is monotone (and Lipschitz if  $\nabla \Phi$  is so) whereas

107 the latter one is maximal monotone. We call a, possibly *set-valued*, operator  $A$  from  $\mathbb{R}^m$  to itself  
 108 monotone if

$$109 \quad (2.7) \quad \langle u - u', z - z' \rangle \geq 0 \quad \forall u \in A(z), u' \in A(z').$$

110 We say  $A$  is maximal monotone, if there exists no monotone operator  $A'$  such that the graph of  $A$  is  
 111 properly contained in the graph of  $A'$ .

112 Problems of type (2.4) have been studied thoroughly in convex optimization, with the most estab-  
 113 lished solution methods being *Extragradient* [29] and *Forward-Backward-Forward* [55]. Both meth-  
 114 ods are known to generate sequences of iterates converging to a solution of (2.4). Note that in the  
 115 unconstrained setting (i.e. if  $\Omega$  is the entire space) both of these algorithms even produce the same  
 116 iterates.

117 **2.3. Solving monotone inclusions.** The connection between monotone inclusions and sad-  
 118 dle point problems is of course not new. The application of Extragradient (EG) to minimax problems  
 119 has been studied in the seminal paper [43] under the name of *Mirror Prox* and a convergence rate of  
 120  $\mathcal{O}(1/k)$  in terms of the function values has been proven. Even a stochastic version of the Mirror Prox  
 121 algorithm has been studied in [27] with a convergence rate of  $\mathcal{O}(1/\sqrt{k})$ . Applied to problem (2.4), with  
 122  $P_\Omega$  being the projection onto  $\Omega$ , it iterates

$$123 \quad (2.8) \quad \text{EG: } \begin{cases} w_k = P_\Omega[z_k - \alpha_k F(z_k)] \\ z_{k+1} = P_\Omega[z_k - \alpha_k F(w_k)]. \end{cases}$$

124 The Forward-Backward-Forward (FBF) method, introduced and convergence of the iterates estab-  
 125 lished in [55], has not been studied rigorously for minimax problems in terms of function values yet,  
 126 despite promising applications in [7] and *its advantage of it only requiring one projection*, whereas  
 127 EG needs two. It is given by

$$128 \quad (2.9) \quad \text{FBF: } \begin{cases} w_k = P_\Omega[z_k - \alpha_k F(z_k)] \\ z_{k+1} = w_k + \alpha_k (F(z_k) - F(w_k)). \end{cases}$$

129 Both, EG and FBF, have the “disadvantage” of needing two gradient evaluations per iteration. A  
 130 possible remedy — suggested in [15] for EG under the name of *extrapolation from the past* — is to  
 131 reuse previous gradients. In a similar fashion we consider

$$132 \quad (2.10) \quad \text{FBFp: } \begin{cases} w_k = P_\Omega[z_k - \alpha_k F(w_{k-1})] \\ z_{k+1} = w_k + \alpha_k (F(w_{k-1}) - F(w_k)), \end{cases}$$

133 where we replaced  $F(z_k)$  by  $F(w_{k-1})$  twice in (2.9). As a matter of fact, the above method can be  
 134 written exclusively in terms of the first variable  $w_k$  by incrementing the index  $k$  in the first update and  
 135 then substituting in the second line. This results in

$$136 \quad (2.11) \quad w_{k+1} = P_\Omega \left[ w_k - \alpha_{k+1} F(w_k) + \alpha_k (F(w_{k-1}) - F(w_k)) \right].$$

137 This way we rediscover a known method which was studied in [36] and convergence of the iterates  
 138 established for general monotone inclusions under the name of *forward-reflected-backward*. It reduces  
 139 to *optimistic mirror descent* [48, 49] in the unconstrained case with constant step size  $\alpha_k = \alpha$ , giving

$$140 \quad (2.12) \quad w_{k+1} = w_k - \alpha (2F(w_k) - F(w_{k-1}))$$

141 which has been proposed for the training of GANs under the name of *Optimistic Gradient Descent*  
 142 *Ascent (OGDA)*, see [11, 12, 31]. Only recently a method very related to (2.11) was proposed in [10]  
 143 and is characterized by applying the correction term after the projection

$$144 \quad (2.13) \quad w_{k+1} = P_{\Omega} \left[ w_k - \alpha F(w_k) \right] - \alpha (F(w_k) - F(w_{k-1})).$$

145 Evidently, in the unconstrained case and for constant stepsize the methods (2.10), (2.11), (2.12)  
 146 and (2.13) are all equivalent.

147 All of the above methods and extensions rely solely on the monotone operator formulation of the  
 148 saddle point problem where the two components  $x$  and  $y$  play a symmetric role. Taking the special  
 149 minimax structure into consideration, [22] showed convergence of a method that uses an optimistic  
 150 step (2.12) in one component and a regular gradient step in the other, thus requiring less storing of past  
 151 gradients in comparison to (2.11).

152 On the downside, however, by reducing the number of required gradient evaluations per iteration,  
 153 the largest possible step size is reduced from  $1/L$  (see [29] or section 3) to  $1/2L$  (see [15, 36, 35] or  
 154 section 3). To summarize, the number of required gradient evaluations is halved, but so is the step  
 155 size, resulting in no clear net gain.

156 **2.4. Regularizers.** The role of regularizers is well studied in many fields such as statistics [54],  
 157 signal processing [45] or inverse problems [52]. They serve different purposes such as inducing spar-  
 158 sity in the solution or conditioning of the problem. In the context of deep learning this has been  
 159 explored from different perspectives, e.g. in incremental convex neural networks where neurons with  
 160 zero weights are removed from the network and new ones are inserted according to different poli-  
 161 cies, see [2, 5, 51, 46]. Other examples include the box-constraints for WGANs with weight clipping  
 162 (see [1]) or spectral normalization (see [41]) which has so far rather been considered as part of the  
 163 architecture, but could at the same time be seen as a regularization step or as a projection onto the set  
 164 of matrices with spectral norm less than 1 (again not rigorously).

165 In the framework of monotone operator theory the optimality condition of the regularized minimax  
 166 problem (2.1) can be written as

$$167 \quad (2.14) \quad 0 \in F(w) + \partial r(w),$$

168 where  $r$  is given by  $(x, y) \mapsto f(x) + h(y)$ . The possibly set-valued operator  $\partial r$  denotes the subdiffer-  
 169 ential of  $r$  and is given by

$$170 \quad (2.15) \quad \partial r(w) := \{v \in \mathbb{R}^m : \langle v, w' - w \rangle + r(w) \leq r(w') \quad \forall w' \in \mathbb{R}^m\}.$$

171 The monotone inclusion (2.14) generalizes (2.4) in a natural way, since  $N_{\Omega} = \partial \delta_{\Omega}$ . Similarly, the  
 172 projection constitutes a special case of the so-called *proximal mapping* which for the function  $r$  and  
 173  $\lambda > 0$  is given by

$$174 \quad (2.16) \quad \text{prox}_{\lambda r}(w) := \arg \min_{w' \in \mathbb{R}^m} \left\{ r(w') + \frac{1}{2\lambda} \|w' - w\|^2 \right\}.$$

175 In particular, the proximal mapping of the indicator  $\delta_{\Omega}$  yields the projection onto the set  $\Omega$ , i.e.  
 176  $\text{prox}_{\lambda \delta_{\Omega}} = P_{\Omega}$ .

177 **3. Main results.** Motivated by the considerations above we study the inclusion problem

$$178 \quad (3.1) \quad 0 \in F(w) + \partial r(w),$$

179 where  $F : \mathbb{R}^m \rightarrow \mathbb{R}^m$  is a monotone and Lipschitz operator and  $r : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  is a proper  
180 convex lower semicontinuous function.

181 **3.1. Measure of optimality.** There are two common quantities measuring the quality of a point  
182 with respect to the monotone inclusion (2.14). The most natural one is the distance to the solution set  
183 for which typically only asymptotic convergence can be proved. If  $F$  arises from a saddle point  
184 problem (2.1) meaning that  $F$  has the form (2.5), we want to use a more problem specific measure,  
185 the *minimax gap*, which for a point  $w = (u, v) \in \mathbb{R}^d \times \mathbb{R}^n$  is given by

$$186 \quad (3.2) \quad \sup_{y \in \mathbb{R}^n} \Psi(u, y) - \inf_{x \in \mathbb{R}^d} \Psi(x, v) \left( = \sup_{x \in \mathbb{R}^d, y \in \mathbb{R}^n} \Psi(u, y) - \Psi(x, v) \right).$$

187 This minimax gap can be interpreted from a game theoretic standpoint as the sum of the maximal  
188 payoffs achievable by the two players by playing their respective best responses, given the current  
189 strategy of the opponent. In the more general monotone inclusion setting where no function values are  
190 available, an appropriate generalization of (3.2) is given for any  $w \in \mathbb{R}^m$  by

$$191 \quad (3.3) \quad \sup_{z \in \mathbb{R}^m} \langle F(z), w - z \rangle + r(w) - r(z).$$

192 If  $r$  is the indicator  $\delta_\Omega$  of the compact and convex set  $\Omega$  it is clear that the supremum is only taken over  
193  $z \in \Omega$  and will thus be finite.

194 *The restricted gap.* Since the problem (3.1) is in general unconstrained and the supremum can be  
195 infinite we consider instead, as done for example in [44], the *restricted gap* where the above supremum  
196 is taken over an auxiliary compact set  $B \subset \mathbb{R}^m$  instead of the entire space. Note that the restricted gap  
197 is in general only a reasonable measure of optimality for elements of  $B$ . It is nonnegative on  $B$  and  
198 zero for points of  $B$  which solve (3.1). Additionally we want to be able to conclude that if a point  $w^*$   
199 has zero gap it solves (3.1). This is for example the case if  $w^*$  is in the interior of  $B$ , which can always  
200 be ensured if  $B$  is chosen large enough.

201 In order to capture both at the same time we define the following unifying gap

$$202 \quad (3.4) \quad G_B(w) := \begin{cases} \sup_{(x,y) \in B} \Psi(u, y) - \Psi(x, v) & \text{if } F \text{ and } r \text{ come from (2.1)} \\ \sup_{z \in B} \langle F(z), w - z \rangle + r(w) - r(z) & \text{otherwise.} \end{cases}$$

203 **3.2. Methods.** We now present a novel unifying scheme for solving problem (3.1), which gen-  
204 eralizes FBF (2.9) and in addition recovers the method motivated in (2.10) as FBFp. Let us point out  
205 again that the latter algorithm was already introduced in [36] and corresponds to OGDA [48, 11, 12]  
206 if  $F$  stems from the minimax setting (2.5).

207 **Algorithm 3.1 (generalized FBF).** For a starting point  $z_0 \in \mathbb{R}^m$  and step sizes  $\alpha_k > 0$  we  
208 consider for all  $k \geq 0$

$$209 \quad (3.5) \quad \begin{cases} w_k = \text{prox}_{\alpha_k r}(z_k - \alpha_k F(\diamond_k)) \\ z_{k+1} = w_k + \alpha_k (F(\diamond_k) - F(w_k)). \end{cases}$$

210 For  $\diamond_k = z_k$  this reduces to the well known FBF method, whereas  $\diamond_k = w_{k-1}$ , with the additional  
211 initial condition  $w_{-1} = z_0$ , recycles previous gradients (FBFp).



212 Consider the scenario where  $F$  is given as an expectation  $\mathbb{E}_\xi[F(\cdot; \xi)]$ , e.g. coming from (2.1), and only  
 213 a stochastic estimator  $F(\cdot; \xi)$  is accessible instead of  $F$  itself. In this case we adapt Algorithm 3.1 in  
 214 the following way.

215 **Algorithm 3.2 (generalized stochastic FBF).** For a starting point  $z_0 \in \mathbb{R}^m$  and step sizes  
 216  $\alpha_k > 0$  we consider for all  $k \geq 0$

$$217 \quad \begin{cases} \xi_k \sim Q & (\text{optionally } \eta_k \sim Q) \\ w_k = \text{prox}_{\alpha_k r}(z_k - \alpha_k F(\diamond_k; \triangle_k)) \\ z_{k+1} = w_k + \alpha_k (F(\diamond_k; \triangle_k) - F(w_k; \xi_k)). \end{cases}$$

218 For  $\diamond_k = z_k$  and  $\triangle_k = \eta_k$  this results in a stochastic version of FBF, whereas  $\diamond_k = w_{k-1}$  and  $\triangle_k =$   
 219  $\xi_{k-1}$  recycles previous gradients (stochastic FBFp) with the additional initial condition  $w_{-1} = z_0$  and  
 220  $\xi_{-1} = \eta_0$ .

221 Even though both methods encompassed by the unifying scheme Algorithm 3.1 have been studied  
 222 in the deterministic setting before, the stated convergence results are new. Note that while the rate  
 223 for FBF is completely new our result for FBFp provides only a generalization of the known rate for  
 224 OGD, see [42]. Similarly, the stochastic version of FBF has been considered before in [6] and rates  
 225 have been obtained, but only in terms of the fixed point residual and not the function values. However,  
 226 we want to point out that the stochastic version of FBFp has not been considered prior to this work.

227 **3.3. Convergence.** Let in the following  $B \subset \mathbb{R}^m$  be the compact set of the restricted (unify-  
 228 ing) gap function (3.4) with  $D := \sup_{w, z \in B} \|z - w\|$  denoting its diameter. For convenience in the  
 229 estimation we assume that the starting point  $z_0$  of the discussed methods is in  $B$ . Recall that  $L \geq 0$   
 230 denotes the Lipschitz constant of the operator  $F$ .

231 **Theorem 3.3 (deterministic).** Let  $(w_k)_{k \geq 0}$  be the sequence generated by Algorithm 3.1. If

232 (i) FBF, i.e.  $\diamond_k = z_k$ , with step size  $\alpha_k = \alpha \leq 1/L$ , or

233 (ii) FBFp, i.e.  $\diamond_k = w_{k-1}$ , with step size  $\alpha_k = \alpha \leq 1/2L$

234 is chosen, then for all  $K \geq 1$  the averaged iterates  $\bar{w}_K := \frac{1}{K} \sum_{k=0}^{K-1} w_k$  fulfill

$$235 \quad (3.6) \quad G_B(\bar{w}_K) \leq \frac{D^2}{2\alpha K},$$

236 where  $G_B$  is the restricted gap defined in (3.4).

237 The proof of the theorem relies on standard techniques making use of the monotonicity of the operator  
 238  $F$  where we are able to show that the distance to the solution set decreases in every iteration up to an  
 239 additive error of

$$240 \quad (3.7) \quad -\|z_k - w_k\|^2 + \alpha_k^2 L^2 \|\diamond_k - w_k\|.$$

241 Depending on which iterate we use for  $\diamond_k$  we choose the stepsize  $\alpha_k^2$  such that the two terms cancel  
 242 and obtain the desired result.

243 In order to derive similar convergence statements for the stochastic algorithm we need to assume  
 244 (standard) properties of the gradient estimator  $F(\cdot; \xi)$ .

245 **Assumption 3.4. Unbiasedness:**  $\mathbb{E}_\xi[F(w; \xi)] = F(w) \forall w \in \mathbb{R}^m$ .

246 **Assumption 3.5.** *Bounded variance:*  $\mathbb{E}_\xi[\|F(w; \xi) - F(w)\|^2] \leq \sigma^2 \forall w \in \mathbb{R}^m$ .

247 We want to point out that the latter assumption is weaker than the commonly used bounded  
248 (sub)gradient hypothesis ( $\mathbb{E}_\xi[\|F(w; \xi)\|^2] \leq \sigma^2$ ), which is known to conflict with properties like  
249 strong monotonicity, see [34].

250 In particular we actually only need the above assumption to hold for all iterates  $w_k$ . Such an  
251 hypothesis is in practice difficult to check, but could be exploited in special cases where additional  
252 properties of the variance and boundedness of the iterates are known a priori.

253 **Assumption 3.6.** *The samples  $\xi_k$  are independent of the iterates  $w_k$ , for all  $k \geq 0$ .*

254 Equipped with these assumptions we are now able to prove the statement.

255 **Theorem 3.7 (stochastic).** *Let Assumptions 3.4 to 3.6 hold and let  $(w_k)_{k \geq 0}$  be the sequence  
256 generated by Algorithm 3.2. If*

257(i) *stochastic FBF, i.e.  $\diamond_k = z_k$  and  $\triangle_k = \eta_k$ , with step size  $\alpha_k \leq \alpha \leq 1/\sqrt{2}L$ , or*

258(ii) *stochastic FBFp, i.e.  $\diamond_k = w_{k-1}$  and  $\triangle_k = \xi_{k-1}$ , with step size  $\alpha_k \leq \alpha \leq 1/3L$*

259 *is chosen, then for all  $K \geq 1$  the averaged iterates  $\bar{w}_K := \frac{\sum_{k=0}^{K-1} \alpha_k w_k}{\sum_{k=0}^{K-1} \alpha_k}$  fulfill*

$$260 \quad (3.8) \quad \mathbb{E}[G_B(\bar{w}_K)] \leq \frac{D^2 + 24\sigma^2 \sum_{k=0}^{K-1} \alpha_k^2}{\sum_{k=0}^{K-1} \alpha_k},$$

261 where  $G_B$  is the restricted gap defined in (3.4).

262 The above theorem exhibits a classical step size dependence [50], yielding convergence for sequences  
263  $(\alpha_k)_{k \geq 0}$  that are square summable  $\sum_{k=0}^{\infty} \alpha_k^2 < +\infty$  but not summable  $\sum_{k=0}^{\infty} \alpha_k = +\infty$ . Addition-  
264 ally, if in the setting of Theorem 3.7 the step size is chosen to be  $\alpha_k = \alpha/\sqrt{k+1}$ , a convergence rate  
265 can be obtained and is given by

$$266 \quad (3.9) \quad \mathbb{E}[G_B(\bar{w}_K)] = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

267 If the step size does not go to zero, the gap can usually not be expected to vanish either. However,  
268 we can still show decrease in the gap up to a residual stemming from the variance. In particular, for a  
269 constant step size  $\alpha_k = \alpha$  we have

$$270 \quad (3.10) \quad \mathbb{E}[G_B(\bar{w}_K)] \leq \frac{D^2}{\alpha K} + 24\sigma^2 \alpha.$$

271 Additionally, if the number of iterations  $K$  is fixed beforehand, a conclusion similar to (3.9) can be  
272 obtained by choosing  $\alpha = 1/\sqrt{K}$  in (3.10).

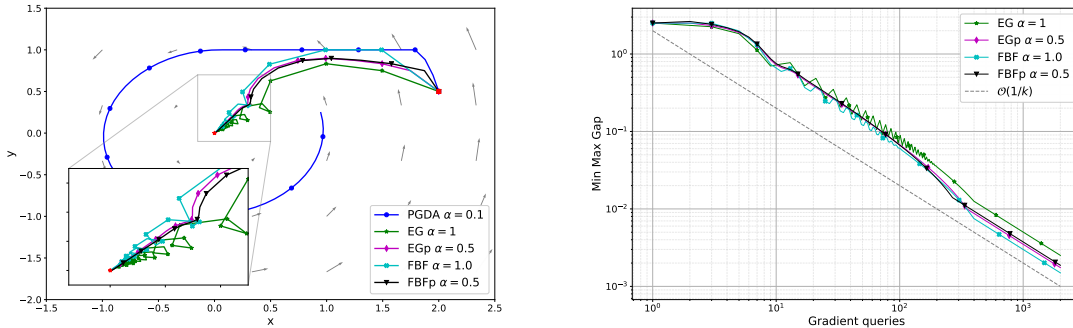
273 **4. Experiments.** The aim of this section is to show how the use of methods with convergence  
274 guarantees, albeit only in the monotone setting, can yield better training performance for different  
275 architectures and objectives. In particular, we demonstrate that FBF can perform at least as good as  
276 EG although requiring less evaluations of the regularizers.



277 **4.1. 2D toy example.** Following [18, 38] and others we consider the example  $\min_x \max_y xy$ ,  
 278 illustrating the cycling behavior of (even bilinear) minimax problems. We augment this approach by  
 279 adding a nonsmooth L1-regularizer for one player, with  $\kappa > 0$ , and constraints for the other, resulting  
 280 in

$$281 \quad (4.1) \quad \min_{x \in \mathbb{R}} \max_{y \in [-1,1]} \kappa|x| + xy.$$

282 The aforementioned issue of GDA (and its proximal extension PGDA) cycling around the solution  
 283 is highlighted in Figure 1. The other methods, for which we display the averaged iterates, however do  
 284 converge to a solution and show a decrease in the restricted gap according to Theorem 3.3. Interest-  
 285 ingly, the two methods using extrapolation from the past seem to exhibit a more monotone decrease  
 286 than FBF or EG.



(a) Trajectories converging to solution.

(b) Restricted gap function.

Figure 1: A comparison of the methods presented in subsection 2.3 applied to problem (4.1) with  $\kappa = 0.01$ . *PGDA* denotes (alternating) gradient descent ascent with proximal steps. As mentioned in the introduction it fails to converge (without averaging of the iterates). *EGp* denotes the method presented in [15] as extrapolation from the past. For the restricted gap we use  $B_1 = B_2 = [-1, 1]$ .

287 **4.2. WGAN trained on CIFAR10.** We now apply the above proposed techniques from mono-  
 288 tone inclusions to the training of Wasserstein GANs employing DCGAN [47] and ResNet [23] archi-  
 289 tectures. All models are trained on the CIFAR10 dataset [30] which consists of 60,000 images in 10  
 290 different classes (with 50,000 training images and 10,000 test images) using an NVIDIA RTX 2080Ti  
 291 GPU.

292 For the DCGAN experiments we work with the original WGAN formulation including weight  
 293 clipping, since it includes regularizers innately (the indicator of a box for the weights of the discrim-  
 294 inator). In addition we propose a modification of the WGAN formulation which replaces the box  
 295 constraint on the discriminator’s weights with an L1-regularization, under the name of *WGAN-L1*.  
 296 This results in a *soft thresholding* operation instead of the “harsh” clipping.

297 For the experiments on ResNet we use the WGAN-GP formulation [21] which penalizes the norm

298 of the gradient of the discriminator to enforce the Lipschitz constraint, together with spectral normal-  
 299 ization of the weight matrices [41] which can be seen as a projection as argued in [subsection 2.4](#).

Table 1: The best Inception Score (IS)<sup>1</sup> and Fréchet Inception Distance (FID). The column denoted by *WGAN*, *WGAN-L1* and *WGAN-GP* refers to the standard formulation with weight clipping, our regularized implementation using the 1-norm and the formulation with gradient penalty and spectral normalization, respectively.

| Method     | WGAN            |                  | WGAN-L1         |                  | WGAN-GP         |                   |
|------------|-----------------|------------------|-----------------|------------------|-----------------|-------------------|
|            | IS              | FID              | IS              | FID              | IS              | FID               |
| AltAdam1   | 4.12±.06        | 56.44±.62        | 4.43±.03        | 50.86±2.17       | 6.01±.31        | 28.11±3.65        |
| Extra Adam | 4.07±.05        | 56.67±.61        | 4.67±.11        | 47.24±1.21       | <b>6.58±.08</b> | 21.40±.58         |
| FBF Adam   | <b>4.54±.04</b> | <b>45.85±.35</b> | <b>4.68±.16</b> | <b>46.60±.76</b> | 6.57±.10        | <b>21.22±1.29</b> |
| Opt. Adam  | 4.35±.06        | 50.41±.46        | 4.63±.13        | 47.98±1.49       | 6.42±.10        | 23.01±.95         |

300 Given the ubiquity and dominance of Adam [28] as an optimizer for many deep learning related  
 301 training tasks, instead of using vanilla SGD we opt for Adam updates. This results in a method we call  
 302 *FBF Adam*. Analogous approaches have been applied in [15] and [11] resulting in *Extra Adam* and  
 303 *Optimistic Adam*, respectively. We compare the aforementioned methods with the status-quo in GAN  
 304 training, namely alternating one Adam step for each network: *AltAdam1*.

305 Our hyperparameter search was limited to the step sizes when using the WGAN-L1 and WGAN-  
 306 GP formulation, while all other parameters were kept the same as in [15, 7]. Note that we still report  
 307 different values for the IS because we used the updated implementation [3]. It seems noteworthy that in  
 308 the case of soft thresholding bigger step sizes performed better with the only exception of AltAdam1.

309 The two evaluation metrics used are the Inception Score (IS, higher is better) [53] and the Fréchet  
 310 inception distance (FID, lower is better) [24], both computed on 50,000 samples.

311 In [Table 1](#) the best IS<sup>1</sup> and FID for each method are reported. FBF Adam performs at least as good  
 312 as all considered competitors with respect to both evaluation metrics. One can also see that WGAN-L1  
 313 using the proximal operator improves the performance of all considered methods. [Figure 2](#) shows the  
 314 training progress regarding IS for each method and both problem formulations. The graphs suggest  
 315 that making use of WGAN-L1 objective has a stabilizing effect during training, leading to a smoother  
 316 and more consistent learning curve — a property that only FBF Adam seems to exhibit for weight  
 317 clipping. [Figure 3](#) as well as [Table 1](#) show that for the WGAN-GP formulation FBF Adam maintains  
 318 the improved performance of EG compared to GDA, while only requiring half the amount of spectral  
 319 normalizations, resulting in time savings of up to 10% as reported in [41]. An even greater improve-  
 320 ment by using FBF can be seen in [Table 2](#) where we additionally consider the averaged iterates using  
 321 an exponential moving average (EMA). As observed by others [15], this can have a very beneficial  
 322 effect.

323 **5. Conclusion.** By highlighting the connection between GAN objectives and monotone inclu-  
 324 sions, we are able to tackle their training via the Forward-Backward-Forward method which is known

<sup>1</sup>In the case of the IS we use the updated and corrected implementation from [3], leading to lower reported values.

Table 2: Fréchet Inception Distance (FID) of regular iterates and averaged iterates using an exponential moving average (EMA), from training a ResNet with WGAN-GP formulation.

| Method     | WGAN-GP           |                   |              |
|------------|-------------------|-------------------|--------------|
|            | non avg.          | EMA               | best         |
| AltAdam1   | 28.11±3.65        | —                 | —            |
| Extra Adam | 21.40±.58         | 18.29±0.26        | 17.86        |
| FBF Adam   | <b>21.22±1.29</b> | <b>17.58±0.73</b> | <b>16.65</b> |
| Opt. Adam  | 23.01±.95         | 21.00±1.33        | 19.28        |

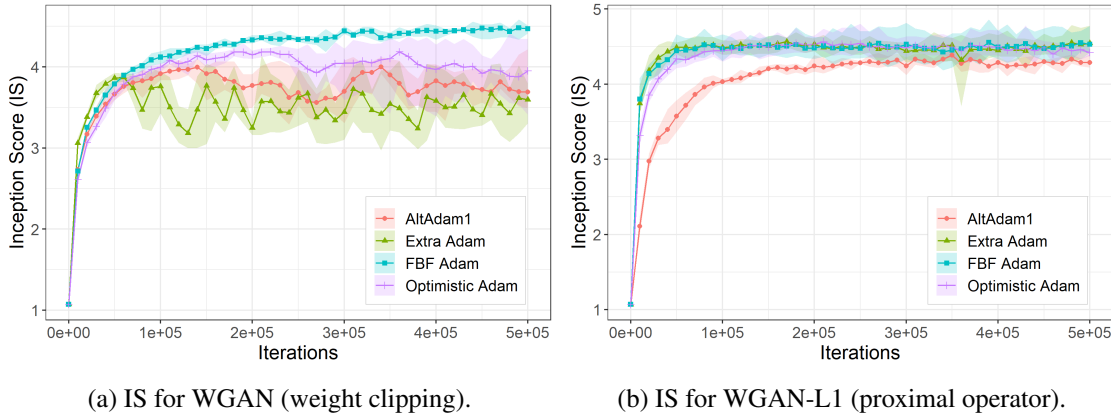


Figure 2: **DCGAN performance.** (a): Average and best/worst IS<sup>1</sup> on the WGAN objective with weight clipping. (b): Average and best/worst IS on the WGAN-L1 objective using the proximal operator (soft thresholding); The WGAN-L1 objective improves the IS in comparison to weight clipping and stabilizes the behavior of all considered methods during the training procedure.

325 to converge to a solution for convex-concave minimax problems. We extend this theoretical under-  
 326 standing by proving novel convergence rates in terms of the function values and highlighting the  
 327 connection to other known methods like OGD. We complement these rigorous considerations by  
 328 promising practical results, indicating that application of FBF can lead to improved performance and  
 329 saved computation time (compared to EG).

### 330 Appendix A. Definitions.

331 In subsection 2.4 we require the regularizers to be proper, convex and lower semicontinuous which  
 332 are common properties in convex analysis. We call a function  $r : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  *proper* if it is not  
 333 constant  $+\infty$ , which means that it takes a finite value for at least a single point. In addition, we say  
 334 that  $r$  is *lower semicontinuous* if for all  $z_0 \in \mathbb{R}^m$

$$335 \text{ (A.1)} \quad \liminf_{z \rightarrow z_0} r(z) \geq r(z_0).$$

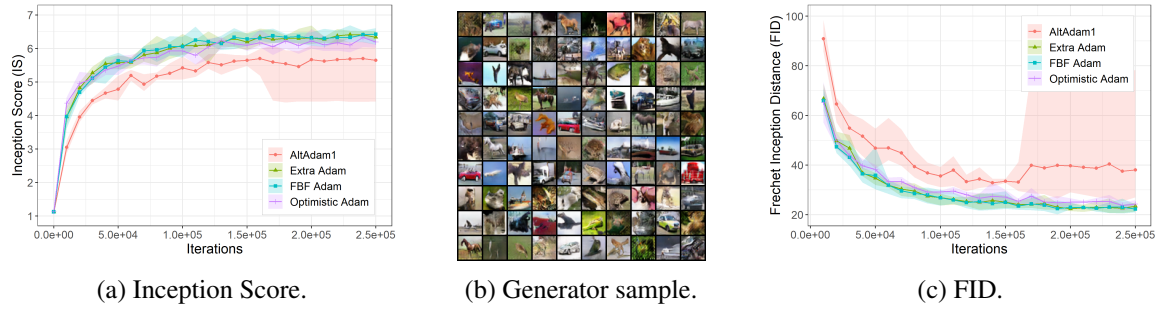


Figure 3: **ResNet performance.** Average and best/worst results regarding  $IS^1$  and FID, see (a) and (b) respectively, using ResNet architecture on the WGAN-GP objective including spectral normalization. (c): Samples from the generator trained with FBF Adam.

336 It is easy to see that if  $C \subset \mathbb{R}^m$  is nonempty, closed and convex, then the indicator  $\delta_C$  of this set, given  
 337 by

$$338 \quad (A.2) \quad \delta_C(z) = \begin{cases} 0 & \text{if } z \in C \\ +\infty & \text{otherwise} \end{cases}$$

339 fulfills the assumptions of being proper, convex and lower semicontinuous.

#### 340 **Appendix B. About the gap function.**

341 Typically in monotone inclusions, the distance to the set of solutions is used as a measure of  
 342 quality of a given point due to the lack of more specific structure in general. Asymptotic convergence  
 343 of the iterates has been established for FBF and FBFp in [4, Proposition 27.13] and [36], respectively.  
 344 Furthermore, no convergence rates can be expected without stronger monotonicity assumptions. We  
 345 want to take into account the special structure of the monotone inclusion coming from the minimax  
 346 problem (2.1). For this reason we use the following (restricted) *minimax gap*, common for saddle point  
 347 problems, which for a point  $(u, v)$  is given by

$$348 \quad (B.1) \quad G_B(u, v) = \sup_{(x, y) \in B} \Psi(u, y) - \Psi(x, v).$$

349 For the general case, i.e.  $F$  being an arbitrary monotone and Lipschitz operator this is connected to  
 350 the other measure of optimality we use in (3.4), for  $w \in \mathbb{R}^m$  given by

$$351 \quad (B.2) \quad G_B(w) = \sup_{z \in B} \langle F(z), w - z \rangle + r(w) - r(z),$$

352 where we interpret the possible occurrence of  $\infty - \infty$  as  $+\infty$ . It stems from the field of Variational  
 353 Inequalities where such a function is also known as *merit function* [44]. The relevance of the above  
 354 two quantities will be made clear by the following statements.

355 **Theorem B.1.** *Let  $\Phi : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable and  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ ,  
 356  $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be proper, convex and lower semicontinuous and  $B \subset \mathbb{R}^d \times \mathbb{R}^n$ . A point  $(x^*, y^*)$   
 357 in the interior of  $B$  solves the saddle point problem (2.1) if and only if its minimax gap (B.1) is zero,  
 358  $G_B(x^*, y^*) = 0$ . For all other elements of  $B$  the gap is nonnegative.*

359 *Proof.* A saddle point  $(x^*, y^*)$  clearly fulfills that  $\sup_{(x,y) \in \mathbb{R}^d \times \mathbb{R}^n} \Psi(x^*, y) - \Psi(x, y^*) = 0$ . On  
 360 the other hand let  $G_B(x^*, y^*) = 0$ . For an arbitrary point  $(x, y)$  we can choose  $\alpha \in (0, 1)$  large  
 361 enough such that  $(u, v) := \alpha(x^*, y^*) + (1 - \alpha)(x, y)$  is in the interior of  $B$ . Therefore,

$$362 \quad (\text{B.3}) \quad \Psi(x^*, v) - \Psi(u, y^*) = \Psi(x^*, \alpha y^* + (1 - \alpha)y) - \Psi(\alpha x^* + (1 - \alpha)x, y^*) \leq 0.$$

363 Using the convex-concave structure of  $\Psi$  we deduce that

$$364 \quad (\text{B.4}) \quad \alpha \Psi(x^*, y^*) + (1 - \alpha) \Psi(x^*, y) - \alpha \Psi(x^*, y^*) - (1 - \alpha) \Psi(x, y^*) \leq 0,$$

365 which implies that  $\Psi(x^*, y) \leq \Psi(x, y^*)$ . Since  $(x, y)$  was chosen arbitrary  $(x^*, y^*)$  is a saddle point. ■

366 Similarly, an analogous statement can be shown for (B.2). The proof, however is split up into  
 367 multiple lemmas to highlight the connection to Variational Inequalities.

368 **Theorem B.2.** Let  $F : \mathbb{R}^m \rightarrow \mathbb{R}^m$  be monotone and continuous,  $r : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  proper,  
 369 convex and lower semicontinuous and  $B \subset \mathbb{R}^m$ . A point  $w^*$  in the interior of  $B$  solves the monotone  
 370 inclusion

$$371 \quad (\text{B.5}) \quad 0 \in F(w) + \partial r(w)$$

372 if and only if its restricted gap (B.2) is zero,  $G_B(w^*) = 0$ . For all other elements of  $B$  the gap is  
 373 nonnegative.

374 Let the assumptions of Theorem B.2 hold true for the following lemmas as we break up the proof into  
 375 separate statements. We do so by making use of the associated Variational inequality (VI)

$$376 \quad (\text{B.6}) \quad \text{find } w \text{ such that } \langle F(w), z - w \rangle + r(z) - r(w) \geq 0 \quad \forall z \in \mathbb{R}^m.$$

377

378 **Lemma B.3.** The monotone inclusion (B.5) is equivalent to the VI (B.6).

379 *Proof.* The equivalence of (B.5) and (B.6) follows immediately from the definition of the subdif-  
 380 ferential of  $r$ . ■

381 The formulation (B.6) is typically referred to as the *strong* form of the VI, whereas

$$382 \quad (\text{B.7}) \quad \text{find } w \text{ such that } \langle F(z), z - w \rangle + r(z) - r(w) \geq 0 \quad \forall z \in \mathbb{R}^m,$$

383 is known as the *weak* formulation.

384 **Lemma B.4.** Under the given assumptions the notion of weak and strong VI are equivalent.

385 *Proof.* For the monotone operator  $F$  it is clear that if  $w^*$  is a solution to the strong formula-  
 386 tion (B.6), it is also a solution to the weak formulation (B.7). In fact, if  $F$  is continuous the reverse im-  
 387 plication also holds true. To see this, let  $w^*$  be a solution to the weak VI (B.7) and  $z = \alpha w^* + (1 - \alpha)u$   
 388 for an arbitrary  $u \in \mathbb{R}^m$  and  $\alpha \in (0, 1)$ , then

$$389 \quad (\text{B.8}) \quad \langle F(\alpha w^* + (1 - \alpha)u), (1 - \alpha)(u - w^*) \rangle + r(\alpha w^* + (1 - \alpha)u) - r(w^*) \geq 0.$$

390 This implies by the convexity of  $r$  that

$$391 \quad (\text{B.9}) \quad (1 - \alpha) \langle F(\alpha w^* + (1 - \alpha)u), (u - w^*) \rangle + (1 - \alpha)(r(u) - r(w^*)) \geq 0.$$

392 By dividing by  $(1 - \alpha)$  and then taking the limit  $\alpha \rightarrow 1$  we obtain that  $w^*$  is a solution of the strong  
 393 form (B.6). ■

394 With the notion of VIs in mind, the above defined gap (B.2) becomes natural as it measures how much  
395 the statement of (B.7) is violated.

396 **Lemma B.5.**  $G_B$  is nonnegative on  $B$  and zero for solutions of the weak VI.

397 *Proof.* It is clear that  $G_B(w) \geq 0$  for  $w \in B$  as  $z = w$  can be chosen in the supremum. On the  
398 other hand if  $w^* \in B$  is a solution to the weak VI (B.7) then  $G_B(w^*) = 0$ . This follows from the fact  
399 that for a solution of (B.7) for all  $z \in B$

$$400 \text{ (B.10)} \quad \langle F(z), w^* - z \rangle + r(w^*) - r(z) \leq 0.$$

401 Therefore the supremum over the above expression in  $z$  is also less than zero, but clearly zero is  
402 obtained for  $z = w^*$ . ■

403 For the reverse implication to hold true, we may not use points on the boundary of  $B$ .

404 **Lemma B.6.** If a point  $w^*$  in the interior of  $B$  exhibits zero gap  $G_B(w^*) = 0$ , then it is a solution  
405 to the weak VI (B.7).

406 *Proof.* Since  $w^*$  is in the interior of  $B$  we can, for an arbitrary  $w \in \mathbb{R}^m$ , choose  $\alpha \in (0, 1)$  large  
407 enough such that  $z := \alpha w^* + (1 - \alpha)w \in B$ . Using this  $z$  in the supremum of the gap we deduce that

$$408 \text{ (B.11)} \quad \langle F(\alpha w^* + (1 - \alpha)w), w^* - \alpha w^* - (1 - \alpha)w \rangle + r(w^*) - r(\alpha w^* + (1 - \alpha)w) \leq 0.$$

409 This implies that

$$410 \text{ (B.12)} \quad (1 - \alpha) \langle F(\alpha w^* + (1 - \alpha)w), w - w^* \rangle + (1 - \alpha)(r(w) - r(w^*)) \geq 0.$$

411 By dividing by  $(1 - \alpha)$  and then taking the limit  $\alpha \rightarrow 1$  we deduce that  $w^*$  solves the strong form of  
412 the VI (B.6). ■

413 Now, we can turn to proving the theorem.

414 *Proof of Theorem B.2.* Combine Lemmas B.3 to B.6. ■

415 **Appendix C. Refined theorems.** Recall that restricted (unifying) gap function  $G_B$  defined  
416 in (3.4) is computed with respect to a set  $B \subset \mathbb{R}^m$  where  $D := \sup_{w, z \in B} \|z - w\|$  denotes its diameter  
417 and it is assumed that  $z_0 \in B$ . Furthermore, the averaged iterates  $\bar{w}_K$  for  $K \geq 1$  are given by

$$418 \text{ (C.1)} \quad \bar{w}_K := \frac{\sum_{k=0}^{K-1} \alpha_k w_k}{\sum_{k=0}^{K-1} \alpha_k}.$$

419 **C.1. Deterministic statements.** The convergence statement of Theorem 3.3 actually holds  
420 true not just for a constant step size as presented in section 3, but for variable step sizes as well.

421 **Theorem C.1.** Let  $(w_k)_{k \geq 0}$  be the sequence generated by Algorithm 3.1. If

422 (i) FBF, i.e.  $\diamond_k = z_k$ , with step size  $0 < \alpha_k \leq \alpha \leq 1/L$ , or

423 (ii) FBFp, i.e.  $\diamond_k = w_{k-1}$ , with step size  $0 < \alpha_k \leq \alpha \leq 1/2L$

424 is chosen, then for all  $K \geq 1$

$$425 \text{ (C.2)} \quad G_B(\bar{w}_K) \leq \frac{D^2}{2 \sum_{k=0}^{K-1} \alpha_k}.$$



426 **C.2. Stochastic statements.** We actually prove a slightly more general version of [Theorem 3.7](#).  
 427 In particular the step size can be chosen larger than initially claimed, however, at the cost of  
 428 a worse constant.

429 **Theorem C.2.** *Let [Assumptions 3.4 to 3.6](#) hold and let  $(w_k)_{k \geq 0}$  be the sequence generated by  
 430 FBF, i.e. [Algorithm 3.2](#) with  $\diamond_k = z_k$  and  $\triangle_k = \eta_k$ . Let the step size  $\alpha_k \leq \alpha < \frac{1}{L}$ , then*

$$431 \quad (C.3) \quad \mathbb{E}[G_B(\bar{w}_K)] \leq \frac{D^2 + 4(1 - \alpha^2 L^2)^{-1} \sigma^2 \sum_{k=0}^{K-1} \alpha_k^2}{2 \sum_{k=0}^{K-1} \alpha_k},$$

432 for all  $K \geq 1$ .

433 [Theorem 3.7](#) (i) can be deduced from the above statement by using  $\alpha = 1/\sqrt{2}L$  which yields that  
 434  $(1 - \alpha^2 L^2)^{-1} = 2$ .

435 **Theorem C.3.** *Let [Assumptions 3.4 to 3.6](#) hold and let  $(w_k)_{k \geq 0}$  be the sequence generated by  
 436 FBFp, i.e. [Algorithm 3.2](#) with  $\diamond_k = w_{k-1}$  and  $\triangle_k = \xi_{k-1}$ . Let the step size  $\alpha_k \leq \alpha < \frac{1}{2\sqrt{2}L}$ , then*

$$437 \quad (C.4) \quad \mathbb{E}[G_B(\bar{w}_K)] \leq \frac{D^2 + 6 \left(1 + \frac{4\alpha^2 L^2}{1 - 8\alpha^2 L^2}\right) \sigma^2 \sum_{k=0}^{K-1} \alpha_k^2}{\sum_{k=0}^{K-1} \alpha_k},$$

438 for all  $K \geq 1$ .

439 [Theorem 3.7](#) (ii) is obtained from the above theorem by using the particular step size bound of  
 440  $\alpha = 1/3L$ , which yields that

$$441 \quad (C.5) \quad \frac{4\alpha^2 L^2}{1 - 8\alpha^2 L^2} = 4.$$

442 Although, the step size in the refined statements [Theorems C.2 and C.3](#) can be chosen arbitrarily  
 443 close to  $1/L$  and  $1/(2\sqrt{2}L)$  for stochastic FBF and stochastic FBFp, respectively. This does not mean it  
 444 should be — since the constant in the convergence rate deteriorates when the step size is close to its  
 445 allowed upper bound.

## 446 **Appendix D. Proofs.**

447 **D.1. Preparations.** We introduce the notation connected to the strong formulation of the VI [\(B.6\)](#)  
 448 associated to the monotone inclusion [\(3.1\)](#), given by

$$449 \quad (D.1) \quad g(w, z) := \langle F(w), w - z \rangle + r(w) - r(z),$$

450 for  $g : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ . Next we will establish the fact that this function can be used to  
 451 bound the (restricted) unifying gap function, which we remind, is defined as

$$452 \quad (D.2) \quad G_B(w) = \begin{cases} \sup_{(x,y) \in B} \Psi(u, y) - \Psi(x, v) & \text{if } F \text{ is } (2.5) \\ \sup_{z \in B} \langle F(z), w - z \rangle + r(w) - r(z) & \text{otherwise,} \end{cases}$$

453 where in the first case  $(u, v) \in \mathbb{R}^d \times \mathbb{R}^n$  is identified with  $w \in \mathbb{R}^m$ . In particular the dimensions fulfill  
 454  $d + n = m$ , and  $r(w)$  is given by  $f(u) + h(v)$ .

455 **Lemma D.1.** *It holds that for all  $K \geq 1$*

$$456 \quad (\text{D.3}) \quad \sup_{z \in B} \left\{ \frac{1}{\sum_{k=0}^{K-1} \alpha_k} \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \geq G_B(\bar{w}_K).$$

457 *Proof.* First we will prove the case if  $F$  is derived from a saddle point problem. Note that from  
458 the convex-concave structure of  $\Phi$  we get that

$$459 \quad (\text{D.4}) \quad \Phi(u, y) \leq \Phi(u, v) + \langle \nabla_y \Phi(u, v), y - v \rangle$$

460 and

$$461 \quad (\text{D.5}) \quad \Phi(u, v) + \langle \nabla_x \Phi(u, v), x - u \rangle \leq \Phi(x, v).$$

462 By summing the two up we obtain

$$463 \quad (\text{D.6}) \quad \Phi(u, y) - \Phi(x, v) \leq \left\langle \begin{array}{c} -\nabla_x \Phi(u, v), \quad x - u \\ \nabla_y \Phi(u, v), \quad y - v \end{array} \right\rangle.$$

464 We can reformulate the above inequality in terms of  $g$  to see that for  $z = (x, y) \in \mathbb{R}^d \times \mathbb{R}^n$

$$465 \quad (\text{D.7}) \quad \langle F(w), w - z \rangle \geq \Phi(u, y) - \Phi(x, v).$$

466 The statement of the first case is obtained by adding  $r(w) - r(z)$  on both sides and using the fact that  
467  $\Psi$  is convex-concave.

468 If  $F$  is a general monotone operator, then we use its monotonicity to deduce that

$$469 \quad (\text{D.8}) \quad \langle F(w), w - z \rangle \geq \langle F(z), w - z \rangle.$$

470 The desired result follows from using the linearity of the inner product. ■

471 **Notation.** We denote the error of the stochastic estimator via

$$472 \quad (\text{D.9}) \quad Z_k := F(\diamond_k; \triangle_k) - F(\diamond_k) \quad \text{and} \quad W_k := F(w_k; \xi_k) - F(w_k).$$

473 Furthermore, we will denote via  $\mathbb{E}[\cdot | U]$ , the conditional expectation with respect to the random variable  $U$ .

474 We will also need the following lemma.

475 **Lemma D.2.** *Let  $(p_k)_{k \geq 0} \in \mathbb{R}^d$  be a given sequence and  $(v_k)_{k \geq 0}$  recursively defined for all  
476  $k \geq 0$  by  $v_{k+1} := v_k - p_k$  for some  $v_0 \in \mathbb{R}^d$ , then*

$$478 \quad (\text{D.10}) \quad \sum_{k=0}^{K-1} \langle p_k, v_k - u \rangle \leq \frac{1}{2} \|v_0 - u\|^2 + \frac{1}{2} \sum_{k=0}^{K-1} \|p_k\|^2.$$

479 *Proof.* From the three point identity it follows immediately that

$$480 \quad (\text{D.11}) \quad \langle p_k, v_k - u \rangle = \langle v_k - v_{k+1}, v_k - u \rangle = \frac{1}{2} (\|v_k - u\|^2 - \|v_{k+1} - u\|^2 + \|v_{k+1} - v_k\|^2)$$

481 from which the statement of the lemma follows. ■

482 **D.2. A unified decrease result.** We will start with a unifying proposition which covers the  
 483 common parts of all convergence proofs.

484 **Proposition D.3.** For a  $\gamma > 0$  we have that for all  $k \geq 0$  and  $z \in \mathbb{R}^m$

(D.12)

$$485 \quad \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} \|z_k - w_k\|^2 + \frac{1}{2} (1 + \gamma) \alpha_k^2 L^2 \|\diamond_k - w_k\|^2 \\ + \alpha_k \langle W_k, z - w_k \rangle + (1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2).$$

486 *Proof.* Let  $k \geq 0$  and  $z \in \mathbb{R}^m$  be arbitrary. Using the decomposition (D.9) it follows that

$$487 \quad (D.13) \quad \langle \alpha_k F(w_k; \xi_k), w_k - z \rangle = \alpha_k \langle W_k, w_k - z \rangle + \alpha_k \langle F(w_k), w_k - z \rangle.$$

488 Since  $\text{prox}_{\alpha_k r} = (\text{Id} + \alpha_k \partial r)^{-1}$  we deduce that

$$489 \quad (D.14) \quad \langle z - w_k, w_k - z_k + \alpha_k F(\diamond_k; \triangle_k) \rangle \geq \alpha_k (r(w_k) - r(z)).$$

490 Adding (D.13) and (D.14) gives that

$$491 \quad (D.15) \quad \langle \alpha_k (F(w_k; \xi_k) - F(\diamond_k; \triangle_k)) + z_k - w_k, w_k - z \rangle \geq \alpha_k \langle W_k, w_k - z \rangle + \alpha_k g(w_k, z),$$

492 which, using the definition of  $z_{k+1}$ , is equivalent to

$$493 \quad (D.16) \quad \langle z - w_k, z_{k+1} - z_k \rangle \geq \alpha_k \langle W_k, w_k - z \rangle + \alpha_k g(w_k, z).$$

494 We estimate the inner product on the left side of the inequality by inserting and subtracting  $z_k$  and  
 495 using the three point identity twice to deduce

$$496 \quad (D.17) \quad \langle z - w_k, z_{k+1} - z_k \rangle = \langle z - z_k + z_k - w_k, z_{k+1} - z_k \rangle \\ = \frac{1}{2} (\|z - z_k\|^2 - \|z_{k+1} - z\|^2 + \|z_{k+1} - w_k\|^2 - \|z_k - w_k\|^2).$$

497 The first two summands are fine as they will telescope, so we are left with estimating  $\|z_{k+1} - w_k\|^2$ .

498 By the definition of  $z_{k+1}$  we have that

$$499 \quad (D.18) \quad \|z_{k+1} - w_k\|^2 = \alpha_k^2 \|F(\diamond_k; \triangle_k) - F(w_k; \xi_k)\|^2 \\ = \alpha_k^2 \|F(\diamond_k) - F(w_k) + Z_k - W_k\|^2 \\ \leq (1 + \gamma) \alpha_k^2 \|F(\diamond_k) - F(w_k)\|^2 + (1 + \gamma^{-1}) \alpha_k^2 \|Z_k - W_k\|^2 \\ \leq (1 + \gamma) \alpha_k^2 L^2 \|\diamond_k - w_k\|^2 + 2(1 + \gamma^{-1}) \alpha_k^2 (\|Z_k\|^2 + \|W_k\|^2),$$

500 where we inserted and subtracted  $F(\diamond_k)$  and  $F(w_k)$  and applied Young's inequality to deduce the  
 501 result. Adding (D.18), (D.17) and (D.16) we conclude that

(D.19)

$$502 \quad \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} \|z_k - w_k\|^2 + \frac{1}{2} (1 + \gamma) \alpha_k^2 L^2 \|\diamond_k - w_k\|^2 \\ + \alpha_k \langle W_k, z - w_k \rangle + (1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2). \quad \blacksquare$$

### 503 D.3. Forward-Backward-Forward.

504 *Proof for deterministic FBF, Theorem C.1 (i).* We start off by plugging  $\diamond_k = z_k$  into (D.12).  
 505 Since  $W_k = Z_k = 0$  we can use  $\gamma \rightarrow 0$  to deduce that for all  $k \geq 0$

$$506 \quad (D.20) \quad \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} (1 - \alpha_k^2 L^2) \|z_k - w_k\|^2.$$

507 From this it is clear that the step size is constrained by  $\alpha \leq 1/L$  as stated in the theorem. By summing  
 508 up from  $k = 0$  to  $K - 1$  and dividing by  $\sum_{k=0}^{K-1} \alpha_k$  we obtain

$$509 \quad (D.21) \quad \frac{1}{\sum_{k=0}^{K-1} \alpha_k} \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \leq \frac{\|z_0 - z\|^2}{2 \sum_{k=0}^{K-1} \alpha_k}.$$

510 The claimed statement is then derived by taking the supremum in  $z$  over  $B$  and applying Lemma D.1. ■

511 *Proof for stochastic FBF, Theorem C.2.* Plugging  $\diamond_k = z_k$  and  $\triangle_k = \eta_k$  into (D.12) gives for all  
 512  $k \geq 0$

$$513 \quad (D.22) \quad \begin{aligned} & \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \\ & \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} (1 - (1 + \gamma) \alpha_k^2 L^2) \|z_k - w_k\|^2 + \alpha_k \langle W_k, z - v_k \rangle \\ & \quad + \alpha_k \langle W_k, v_k - w_k \rangle + (1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2). \end{aligned}$$

514 By summing this inequality up and applying Lemma D.2 with  $v_0 = z_0$ ,  $p_k = -\alpha_k W_k$  and  $v_{k+1} :=$   
 515  $v_k - p_k$  we deduce that

$$516 \quad (D.23) \quad \sum_{k=0}^{K-1} \langle -\alpha_k W_k, v_k - z \rangle \leq \frac{1}{2} \|z_0 - z\|^2 + \frac{1}{2} \sum_{k=0}^{K-1} \alpha_k^2 \|W_k\|^2,$$

517 and therefore

$$518 \quad (D.24) \quad \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \leq \|z_0 - z\|^2 + \sum_{k=0}^{K-1} \alpha_k \langle W_k, v_k - w_k \rangle + 2(1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2).$$

519 By choosing  $\gamma$  such that  $\alpha = (\sqrt{1 + \gamma} L)^{-1}$  we deduce that  $1 + \gamma^{-1} = 1/(1 - \alpha^2 L^2)$ . Next, we take  
 520 the supremum over  $z \in B$  and the expectation to obtain

$$521 \quad (D.25) \quad \mathbb{E} \left[ \sup_{z \in B} \left\{ \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \right] \leq D^2 + 4(1 - \alpha^2 L^2)^{-1} \sigma^2 \sum_{k=0}^{K-1} \alpha_k^2,$$

522 where we used that

$$523 \quad (D.26) \quad \begin{aligned} \mathbb{E}[\langle W_k, v_k - w_k \rangle] &= \mathbb{E} \left[ \mathbb{E}[\langle W_k, v_k - w_k \rangle \mid w_{[k]}, \xi_{[k-1]}] \right] \\ &= \mathbb{E} \left[ \langle \mathbb{E}[W_k \mid w_{[k]}, \xi_{[k-1]}], v_k - w_k \rangle \right] = 0, \end{aligned}$$

524 with  $\xi_{[k-1]} = (\xi_0, \dots, \xi_{k-1})$  and  $w_{[k]} = (w_0, \dots, w_k)$ . The final statement follows by dividing by  
 525  $\sum_{k=0}^{K-1} \alpha_k$  and applying Lemma D.1. ■

#### 526 D.4. Forward-Backward-Forward-past.

527 *Proof for deterministic FBFp, Theorem C.1 (ii).* We start off by plugging  $\diamond_k = z_k$  into (D.12).

528 Since  $W_k = Z_k = 0$  we can use  $\gamma \rightarrow 0$  to conclude that for all  $k \geq 0$

$$529 \quad (\text{D.27}) \quad \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} \|z_k - w_k\|^2 + \frac{1}{2} \alpha_k^2 L^2 \|w_{k-1} - w_k\|^2.$$

530 Now we need to bound the term  $\|w_{k-1} - w_k\|^2$  by  $\|z_k - w_k\|^2$ . Since

$$531 \quad (\text{D.28}) \quad 2\|z_k - w_k\|^2 + 2\|z_k - w_{k-1}\|^2 \geq \|w_k - w_{k-1}\|^2$$

532 we have for all  $k \geq 1$

$$533 \quad (\text{D.29}) \quad \begin{aligned} \|z_k - w_k\|^2 &\geq -\|z_k - w_{k-1}\|^2 + \frac{1}{2} \|w_{k-1} - w_k\|^2 \\ &\geq -\alpha_{k-1}^2 L^2 \|w_{k-1} - w_{k-2}\|^2 + \frac{1}{2} \|w_{k-1} - w_k\|^2 \end{aligned}$$

534 whereas for  $k = 0$ , since  $w_{-1} = z_0$ , we have that

$$535 \quad (\text{D.30}) \quad \|z_0 - w_0\|^2 = \|w_{-1} - w_0\|^2.$$

536 Plugging (D.30) into (D.27) for  $k = 0$  we get that

$$537 \quad (\text{D.31}) \quad \alpha_0 g(w_0, z) + \frac{1}{2} \|z_1 - z\|^2 + \frac{1}{2} (1 - \alpha_0^2 L^2) \|w_0 - w_{-1}\|^2 \leq \frac{1}{2} \|z_0 - z\|^2.$$

538 Plugging (D.29) into (D.27) we get that for all  $k \geq 1$

$$539 \quad (\text{D.32}) \quad \begin{aligned} \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 + \frac{1}{2} \left( \frac{1}{2} - \alpha_k^2 L^2 \right) \|w_k - w_{k-1}\|^2 \\ \leq \frac{1}{2} \|z_k - z\|^2 + \frac{1}{2} \alpha_{k-1}^2 L^2 \|w_{k-1} - w_{k-2}\|^2. \end{aligned}$$

540 In order to be able to telescope we need to ensure that for all  $k \geq 0$

$$541 \quad (\text{D.33}) \quad \left( \frac{1}{2} - \alpha_k^2 L^2 \right) \geq \alpha_k^2 L^2.$$

542 This is equivalent to the condition  $\alpha_k \leq 1/2L$  which was required in the statement of the theorem. Now

543 we sum up (D.32) from  $k = 1$  to  $K - 1$  which yields

$$544 \quad (\text{D.34}) \quad \begin{aligned} \sum_{k=1}^{K-1} \alpha_k g(w_k, z) + \frac{1}{2} \|z_K - z\|^2 + \frac{1}{2} \left( \frac{1}{2} - \alpha_{K-1}^2 L^2 \right) \|w_{K-1} - w_{K-2}\|^2 \\ \leq \frac{1}{2} \|z_1 - z\|^2 + \frac{1}{2} \alpha_0^2 L^2 \|w_0 - w_{-1}\|^2. \end{aligned}$$

545 Adding (D.34) and (D.31) and dividing by  $\sum_{k=0}^{K-1} \alpha_k$  to deduce

$$546 \quad (\text{D.35}) \quad \frac{1}{\sum_{k=0}^{K-1} \alpha_k} \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \leq \frac{\|z_0 - z\|^2}{2 \sum_{k=0}^{K-1} \alpha_k},$$

547 where we used that  $1 - \alpha_0^2 L^2 \geq \alpha_0^2 L^2$  to get rid of  $\|w_0 - w_{-1}\|^2$ . The final statement follows by  
548 taking the supremum in  $z$  over  $B$  and applying Lemma D.1. ■

549 *Proof for stochastic FBFp, Theorem C.3.* By using  $\diamond_k = w_{k-1}$  we deduce from (D.12) for all  
 550  $k \geq 0$  that

$$(D.36) \quad \alpha_k g(w_k, z) + \frac{1}{2} \|z_{k+1} - z\|^2 \leq \frac{1}{2} \|z_k - z\|^2 - \frac{1}{2} \|z_k - w_k\|^2 + \frac{1}{2} (1 + \gamma) \alpha_k^2 L^2 \|w_{k-1} - w_k\|^2 \\ + \alpha_k \langle W_k, z - w_k \rangle + 2(1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2).$$

552 As in (D.23) we can split  $\langle \alpha_k W_k, z - w_k \rangle$  into  $\langle \alpha_k W_k, z - v_k \rangle + \langle \alpha_k W_k, v_k - w_k \rangle$  and use Lemma D.2  
 553 to deduce

$$(D.37) \quad \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \leq \|z_0 - z\|^2 - \sum_{k=0}^{K-1} \left( \frac{1}{2} \|z_k - w_k\|^2 + \frac{1}{2} (1 + \gamma) \alpha_k^2 L^2 \|w_{k-1} - w_k\|^2 \right. \\ \left. + \langle \alpha_k W_k, v_k - w_k \rangle + 3(1 + \gamma^{-1}) \alpha_k^2 (\|W_k\|^2 + \|Z_k\|^2) \right).$$

555 Taking now the supremum over  $z \in B$  and then the expectation we conclude that the inequality

$$(D.38) \quad \mathbb{E} \left[ \sup_{z \in B} \left\{ \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \right] \leq D^2 - \frac{1}{2} \sum_{k=0}^{K-1} (\|z_k - w_k\|^2 - (1 + \gamma) \alpha_k^2 L^2 \|w_{k-1} - w_k\|^2) \\ + 3(1 + \gamma^{-1}) \sigma^2 \sum_{k=0}^{K-1} \alpha_k^2$$

557 holds. Let from now on  $k \geq 1$  as we will treat the case  $k = 0$  separately. Using (D.28) we deduce that

$$(D.39) \quad \|z_k - w_k\|^2 \geq -\|z_k - w_{k-1}\|^2 + \frac{1}{2} \|w_{k-1} - w_k\|^2 \\ \geq -\alpha_{k-1}^2 \|F(w_{k-1}; \xi_{k-1}) - F(w_{k-2}; \xi_{k-2})\|^2 + \frac{1}{2} \|w_{k-1} - w_k\|^2.$$

559 Now we bound the difference of the two estimators by inserting  $\pm F(w_{k-1})$ ,  $\pm F(w_{k-2})$  and applying  
 560 the inequality  $\|a + b + c\|^2 \leq 3(\|a\|^2 + \|b\|^2 + \|c\|^2)$  which yields

$$(D.40) \quad \|F(w_{k-1}; \xi_{k-1}) - F(w_{k-2}; \xi_{k-2})\|^2 \leq 3\|W_{k-1}\|^2 + 3\|W_{k-2}\|^2 + 3\|F(w_{k-2}) - F(w_{k-1})\|^2.$$

562 We conclude that

$$(D.41) \quad \mathbb{E} [\|F(w_{k-1}; \xi_{k-1}) - F(w_{k-2}; \xi_{k-2})\|^2] \leq 6\sigma^2 + 3L^2 \mathbb{E} \|w_{k-1} - w_{k-2}\|^2.$$

564 Using (D.41) in (D.39) we deduce that

$$(D.42) \quad \mathbb{E} \|z_k - w_k\|^2 \geq -\alpha_{k-1}^2 (6\sigma^2 + 3L^2 \mathbb{E} \|w_{k-1} - w_{k-2}\|^2) + \frac{1}{2} \mathbb{E} \|w_{k-1} - w_k\|^2,$$



566 whereas for  $k = 0$  we have (D.30). Now we plug (D.42) into (D.38) to conclude that  
 (D.43)

$$\begin{aligned}
 & \mathbb{E} \left[ \sup_{z \in B} \left\{ \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \right] \\
 567 \quad & \leq D^2 - \frac{1}{2} \sum_{k=1}^{K-1} \left( -3\alpha_{k-1}^2 L^2 \mathbb{E} \|w_{k-1} - w_{k-2}\|^2 + \left( \frac{1}{2} - (1 + \gamma)\alpha_k^2 L^2 \right) \|w_{k-1} - w_k\|^2 \right) \\
 & \quad + \frac{1}{2} ((1 + \gamma)\alpha_0^2 L^2 - 1) \|w_{-1} - w_0\|^2 + 6(1 + \gamma^{-1})\sigma^2 \sum_{k=0}^{K-1} \alpha_k^2
 \end{aligned}$$

568 From this we conclude that in order to be able to telescope we need to enforce

$$569 \quad (D.44) \quad \left( \frac{1}{2} - (1 + \gamma)\alpha_k^2 L^2 \right) \geq 3\alpha_k^2 L^2$$

570 which is equivalent to

$$571 \quad (D.45) \quad \frac{1}{2(4 + \gamma)} \geq \alpha_k^2 L^2.$$

572 Since  $\alpha_k \leq \alpha$ , we can ensure this by choosing  $\gamma$  such that

$$573 \quad (D.46) \quad \frac{1}{2(4 + \gamma)} = \alpha^2 L^2.$$

574 With (D.46) in place conclude from (D.43) that the inequality

$$\begin{aligned}
 & \mathbb{E} \left[ \sup_{z \in B} \left\{ \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \right] \\
 575 \quad (D.47) \quad & \leq D^2 + \frac{1}{2} ((4 + \gamma)\alpha_0^2 L^2 - 1) \|w_{-1} - w_0\|^2 + 6(1 + \gamma^{-1})\sigma^2 \sum_{k=0}^{K-1} \alpha_k^2
 \end{aligned}$$

576 Using the fact that  $3\alpha_0^2 L^2 \leq 1 - (1 + \gamma)\alpha_0^2 L^2$  from (D.46) to discard the  $\|w_0 - w_{-1}\|^2$  term, yields

$$577 \quad (D.48) \quad \mathbb{E} \left[ \sup_{z \in B} \left\{ \sum_{k=0}^{K-1} \alpha_k g(w_k, z) \right\} \right] \leq D^2 + 6(1 + \gamma^{-1})\sigma^2 \sum_{k=0}^{K-1} \alpha_k^2$$

578 Through (D.46), we can estimate

$$579 \quad (D.49) \quad \frac{1}{\gamma} = \frac{2\alpha^2 L^2}{1 - 8\alpha^2 L^2}.$$

580 Plugging (D.49) into (D.48), dividing by  $\sum_{k=0}^{K-1} \alpha_k$  and applying Lemma D.1, deduces the final state-  
 581 ment. ■

## REFERENCES

- 582
- 583 [1] M. ARJOVSKY, S. CHINTALA, AND L. BOTTOU, *Wasserstein GAN*, arXiv:1701.07875, (2017).
- 584 [2] F. BACH, *Breaking the curse of dimensionality with convex neural networks*, The Journal of Machine Learning Re-
- 585 search, 18 (2017), pp. 629–681.
- 586 [3] S. BARRATT AND R. SHARMA, *A note on the inception score*, arXiv:1801.01973, (2018).
- 587 [4] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*,
- 588 vol. 408, Springer, 2011.
- 589 [5] Y. BENGIO, N. L. ROUX, P. VINCENT, O. DELALLEAU, AND P. MARCOTTE, *Convex neural networks*, in Advances
- 590 in Neural Information Processing Systems, 2006, pp. 123–130.
- 591 [6] R. I. BOŢ, P. MERTIKOPOULOS, M. STAUDIGL, AND P. T. VUONG, *Minibatch forward-backward-forward methods*
- 592 *for solving stochastic variational inequalities*, Stochastic Systems, 11 (2021), pp. 112–139.
- 593 [7] R. I. BOŢ, M. SEDLMAYER, AND P. T. VUONG, *A relaxed inertial forward-backward-forward algorithm for solving*
- 594 *monotone inclusions with application to GANs*, arXiv:2003.07886, (2020).
- 595 [8] T. CHAVDAROVA, G. GIDEL, F. FLEURET, AND S. LACOSTE-JULIEN, *Reducing noise in GAN training with vari-*
- 596 *ance reduced extragradient*, in Advances in Neural Information Processing Systems, 2019, pp. 391–401.
- 597 [9] T. CHAVDAROVA, M. PAGLIARDINI, S. U. STICH, F. FLEURET, AND M. JAGGI, *Taming gans with lookahead-*
- 598 *minmax*, arXiv preprint arXiv:2006.14567, (2020).
- 599 [10] E. R. CSETNEK, Y. MALITSKY, AND M. K. TAM, *Shadow douglas-rachford splitting for monotone inclusions*,
- 600 Applied Mathematics & Optimization, 80 (2019), pp. 665–678.
- 601 [11] C. DASKALAKIS, A. ILYAS, V. SYRGKANIS, AND H. ZENG, *Training GANs with optimism*, in International Con-
- 602 ference on Learning Representations, 2018, <https://openreview.net/forum?id=SJySbbAZ>.
- 603 [12] C. DASKALAKIS AND I. PANAGEAS, *The limit points of (optimistic) gradient descent in min-max optimization*, in
- 604 Advances in Neural Information Processing Systems, 2018, pp. 9236–9246.
- 605 [13] C. DASKALAKIS, S. SKOULAKIS, AND M. ZAMPETAKIS, *The complexity of constrained min-max optimization*, in
- 606 Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing, 2021, pp. 1466–1478.
- 607 [14] I. GEMP AND S. MAHADEVAN, *Global convergence to the equilibrium of gans using variational inequalities*, arXiv
- 608 preprint arXiv:1808.01531, (2018).
- 609 [15] G. GIDEL, H. BERARD, G. VIGNOUD, P. VINCENT, AND S. LACOSTE-JULIEN, *A variational inequality per-*
- 610 *spective on generative adversarial networks*, in International Conference on Learning Representations, 2019,
- 611 <https://openreview.net/forum?id=r1laEnA5Ym>.
- 612 [16] G. GIDEL, R. A. HEMMAT, M. PEZESHKI, R. LE PRIOL, G. HUANG, S. LACOSTE-JULIEN, AND I. MITLIAGKAS,
- 613 *Negative momentum for improved game dynamics*, in The 22nd International Conference on Artificial Intelligence
- 614 and Statistics, 2019, pp. 1802–1811.
- 615 [17] N. GOLOWICH, S. PATTATHIL, C. DASKALAKIS, AND A. OZDAGLAR, *Last iterate is slower than averaged iterate*
- 616 *in smooth convex-concave saddle point problems*, in Conference on Learning Theory, PMLR, 2020, pp. 1758–
- 617 1784.
- 618 [18] I. GOODFELLOW, *Nips 2016 tutorial: Generative adversarial networks*, arXiv:1701.00160, (2016).
- 619 [19] I. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. COURVILLE,
- 620 AND Y. BENGIO, *Generative adversarial nets*, in Advances in Neural Information Processing Systems, 2014,
- 621 pp. 2672–2680.
- 622 [20] E. GORBUNOV, N. LOIZOU, AND G. GIDEL, *Extragradient method:  $o(1/k)$  last-iterate convergence for monotone*
- 623 *variational inequalities and connections with cocoercivity*, arXiv preprint arXiv:2110.04261, (2021).
- 624 [21] I. GULRAJANI, F. AHMED, M. ARJOVSKY, V. DUMOULIN, AND A. C. COURVILLE, *Improved training of Wasser-*
- 625 *stein GANs*, in Advances in Neural Information Processing Systems, 2017, pp. 5767–5777.
- 626 [22] E. Y. HAMEDANI AND N. S. AYBAT, *A primal-dual algorithm for general convex-concave saddle point problems*,
- 627 arXiv:1803.01401, (2018).
- 628 [23] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in Proceedings of the IEEE
- 629 conference on computer vision and pattern recognition, 2016, pp. 770–778.
- 630 [24] M. HEUSEL, H. RAMSAUER, T. UNTERTHINER, B. NESSLER, AND S. HOCHREITER, *GANs trained by a two time-*
- 631 *scale update rule converge to a local nash equilibrium*, in Advances in Neural Information Processing Systems,
- 632 2017, pp. 6626–6637.
- 633 [25] Y.-G. HSIEH, F. IUTZELER, J. MALICK, AND P. MERTIKOPOULOS, *Explore aggressively, update conservatively:*
- 634 *Stochastic extragradient methods with variable stepsize scaling*, arXiv preprint arXiv:2003.10162, (2020).

- 635 [26] A. N. IUSEM, A. JOFRÉ, R. I. OLIVEIRA, AND P. THOMPSON, *Extragradient method with variance reduction for*  
636 *stochastic variational inequalities*, SIAM Journal on Optimization, 27 (2017), pp. 686–724.
- 637 [27] A. JUDITSKY, A. NEMIROVSKI, AND C. TAUVEL, *Solving variational inequalities with stochastic mirror-prox al-*  
638 *gorithm*, Stochastic Systems, 1 (2011), pp. 17–58.
- 639 [28] D. P. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, arXiv:1412.6980, (2014).
- 640 [29] G. KORPELEVICH, *The extragradient method for finding saddle points and other problems*, Matecon, 12 (1976),  
641 pp. 747–756.
- 642 [30] A. KRIZHEVSKY, G. HINTON, ET AL., *Learning multiple layers of features from tiny images*, (2009).
- 643 [31] T. LIANG AND J. STOKES, *Interaction matters: A note on non-asymptotic local convergence of generative adver-*  
644 *sarial networks*, in The 22nd International Conference on Artificial Intelligence and Statistics, K. Chaudhuri and  
645 M. Sugiyama, eds., vol. 89 of Proceedings of Machine Learning Research, PMLR, 2019, pp. 907–915.
- 646 [32] Q. LIN, M. LIU, H. RAFIQUE, AND T. YANG, *Solving weakly-convex-weakly-concave saddle-point problems as*  
647 *successive strongly monotone variational inequalities*, arXiv:1810.10207, (2018).
- 648 [33] N. LOIZOU, H. BERARD, A. JOLICOEUR-MARTINEAU, P. VINCENT, S. LACOSTE-JULIEN, AND I. MITLIAGKAS,  
649 *Stochastic hamiltonian gradient methods for smooth games*, in International Conference on Machine Learning,  
650 PMLR, 2020, pp. 6370–6381.
- 651 [34] N. LOIZOU, S. VASWANI, I. H. LARADJI, AND S. LACOSTE-JULIEN, *Stochastic polyak step-size for sgd: An*  
652 *adaptive learning rate for fast convergence*, in International Conference on Artificial Intelligence and Statistics,  
653 PMLR, 2021, pp. 1306–1314.
- 654 [35] Y. MALITSKY, *Projected reflected gradient methods for monotone variational inequalities*, SIAM Journal on Opti-  
655 mization, 25 (2015), pp. 502–520.
- 656 [36] Y. MALITSKY AND M. K. TAM, *A forward-backward splitting method for monotone inclusions without cocoercivity*,  
657 SIAM Journal on Optimization, 30 (2020), pp. 1451–1472.
- 658 [37] P. MERTIKOPOULOS, B. LECOAT, H. ZENATI, C.-S. FOO, V. CHANDRASEKHAR, AND G. PILIOURAS, *Opti-*  
659 *mistic mirror descent in saddle-point problems: Going the extra(-gradient) mile*, in International Conference on  
660 Learning Representations, 2019, <https://openreview.net/forum?id=Bkg8jjC9KQ>.
- 661 [38] L. MESCHEDER, A. GEIGER, AND S. NOWOZIN, *Which training methods for GANs do actually converge?*, in  
662 International Conference on Machine Learning, 2018.
- 663 [39] L. MESCHEDER, S. NOWOZIN, AND A. GEIGER, *The numerics of gans*, arXiv preprint arXiv:1705.10461, (2017).
- 664 [40] K. MISHCHENKO, D. KOVALEV, E. SHULGIN, P. RICHTÁRIK, AND Y. MALITSKY, *Revisiting stochastic extragrad-*  
665 *ient*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2020, pp. 4573–4582.
- 666 [41] T. MIYATO, T. KATAOKA, M. KOYAMA, AND Y. YOSHIDA, *Spectral normalization for generative adversar-*  
667 *ial networks*, in International Conference on Learning Representations, 2018, [https://openreview.net/forum?id=](https://openreview.net/forum?id=B1QRgziT)  
668 [B1QRgziT](https://openreview.net/forum?id=B1QRgziT).
- 669 [42] A. MOKHTARI, A. E. OZDAGLAR, AND S. PATTATHIL, *Convergence rate of  $O(1/k)$  for optimistic gradient and*  
670 *extra-gradient methods in smooth convex-concave saddle point problems*, SIAM Journal on Optimization, 30  
671 (2020), pp. 3230–3251.
- 672 [43] A. NEMIROVSKI, *Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz contin-*  
673 *uous monotone operators and smooth convex-concave saddle point problems*, SIAM Journal on Optimization, 15  
674 (2004), pp. 229–251.
- 675 [44] Y. NESTEROV, *Dual extrapolation and its applications to solving variational inequalities and related problems*,  
676 Mathematical Programming, 109 (2007), pp. 319–344.
- 677 [45] D. P. PALOMAR AND Y. C. ELДАР, *Convex optimization in signal processing and communications*, Cambridge  
678 university press, 2010.
- 679 [46] K. PIEPER AND A. PETROSYAN, *Nonconvex penalization for sparse neural networks*, arXiv:2004.11515, (2020).
- 680 [47] A. RADFORD, L. METZ, AND S. CHINTALA, *Unsupervised representation learning with deep convolutional gener-*  
681 *ative adversarial networks*, arXiv:1511.06434, (2015).
- 682 [48] A. RAKHLIN AND K. SRIDHARAN, *Online learning with predictable sequences*, in Proceedings of the 26th Annual  
683 Conference on Learning Theory, 2013, pp. 993–1019.
- 684 [49] S. RAKHLIN AND K. SRIDHARAN, *Optimization, learning, and games with predictable sequences*, in Advances in  
685 Neural Information Processing Systems, 2013, pp. 3066–3074.
- 686 [50] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, The annals of mathematical statistics, (1951),  
687 pp. 400–407.
- 688 [51] S. ROSSET, G. SWIRSZCZ, N. SREBRO, AND J. ZHU,  *$\ell_1$  regularization in infinite dimensional feature spaces*, in

- 689 International Conference on Computational Learning Theory, Springer, 2007, pp. 544–558.
- 690 [52] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, *Physica D:*  
691 *nonlinear phenomena*, 60 (1992), pp. 259–268.
- 692 [53] T. SALIMANS, I. GOODFELLOW, W. ZAREMBA, V. CHEUNG, A. RADFORD, AND X. CHEN, *Improved techniques*  
693 *for training GANs*, in *Advances in Neural Information Processing Systems*, 2016, pp. 2234–2242.
- 694 [54] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, *Journal of the Royal Statistical Society: Series B*  
695 *(Methodological)*, 58 (1996), pp. 267–288.
- 696 [55] P. TSENG, *Applications of a splitting algorithm to decomposition in convex programming and variational inequalities*,  
697 *SIAM Journal on Control and Optimization*, 29 (1991), pp. 119–138.