

## 521 A Appendix

### 522 A.1 Proof of Theorem 4.1

523 *Proof.* To begin with, let us consider Lv.1 GP:

$$\begin{aligned}\boldsymbol{\theta}_t^{(1)} &= \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(0)}) = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t) \\ \boldsymbol{\phi}_t^{(1)} &= \boldsymbol{\phi}_t - \eta \nabla_{\boldsymbol{\phi}} g(\boldsymbol{\theta}_t^{(0)}, \boldsymbol{\phi}_t) = \boldsymbol{\phi}_t - \eta \nabla_{\boldsymbol{\phi}} g(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t)\end{aligned}$$

524 The differences between  $(\boldsymbol{\theta}_t^{(1)}, \boldsymbol{\phi}_t^{(1)})$  and  $(\boldsymbol{\theta}_t^{(0)}, \boldsymbol{\phi}_t^{(0)})$  are:

$$\begin{aligned}\|\boldsymbol{\theta}_t^{(1)} - \boldsymbol{\theta}_t^{(0)}\| &= \eta \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t)\|, \\ \|\boldsymbol{\phi}_t^{(1)} - \boldsymbol{\phi}_t^{(0)}\| &= \eta \|\nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t)\|.\end{aligned}$$

525 Recall our definition of  $\Delta_{\max}$  we have:

$$\|\boldsymbol{\theta}_t^{(1)} - \boldsymbol{\theta}_t^{(0)}\| + \|\boldsymbol{\phi}_t^{(1)} - \boldsymbol{\phi}_t^{(0)}\| \leq \eta \Delta_{\max} \quad (9)$$

526 Then, with  $\boldsymbol{\omega}_t = [\boldsymbol{\theta}_t, \boldsymbol{\phi}_t]^T$  and  $\boldsymbol{\omega}_t^{(k)} = [\boldsymbol{\theta}_t^{(k)}, \boldsymbol{\phi}_t^{(k)}]^T$ , the differences between Lv.2 agents and Lv.1  
527 agents are:

$$\begin{aligned}\|\boldsymbol{\theta}_t^{(2)} - \boldsymbol{\theta}_t^{(1)}\| &= \eta \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(1)}) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(0)})\| \\ &\leq \eta L_{\theta\phi} \|\boldsymbol{\phi}_t^{(1)} - \boldsymbol{\phi}_t^{(0)}\|, \\ \|\boldsymbol{\phi}_t^{(2)} - \boldsymbol{\phi}_t^{(1)}\| &= \eta \|\nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(1)}, \boldsymbol{\phi}_t) - \nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(0)}, \boldsymbol{\phi}_t)\| \\ &\leq \eta L_{\phi\theta} \|\boldsymbol{\theta}_t^{(1)} - \boldsymbol{\theta}_t^{(0)}\|.\end{aligned}$$

528 Recall that  $L := \max\{L_{\theta\theta}, L_{\theta\phi}, L_{\phi\theta}, L_{\phi\phi}\}$ , using Equation (9) we have:

$$\|\boldsymbol{\theta}_t^{(2)} - \boldsymbol{\theta}_t^{(1)}\| + \|\boldsymbol{\phi}_t^{(2)} - \boldsymbol{\phi}_t^{(1)}\| \leq \eta^2 L \Delta_{\max} \quad (10)$$

529 Similarly, we can derive the differences between Lv.3 and Lv.2 agents:

$$\begin{aligned}\|\boldsymbol{\theta}_t^{(3)} - \boldsymbol{\theta}_t^{(2)}\| &= \eta \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(2)}) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(1)})\| \\ &\leq \eta L_{\theta\phi} \|\boldsymbol{\phi}_t^{(2)} - \boldsymbol{\phi}_t^{(1)}\| \\ \|\boldsymbol{\phi}_t^{(3)} - \boldsymbol{\phi}_t^{(2)}\| &= \eta \|\nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(2)}, \boldsymbol{\phi}_t) - \nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(1)}, \boldsymbol{\phi}_t)\| \\ &\leq \eta L_{\phi\theta} \|\boldsymbol{\theta}_t^{(2)} - \boldsymbol{\theta}_t^{(1)}\| \\ \|\boldsymbol{\theta}_t^{(3)} - \boldsymbol{\theta}_t^{(2)}\| + \|\boldsymbol{\phi}_t^{(3)} - \boldsymbol{\phi}_t^{(2)}\| &\leq \eta^3 L^2 \Delta_{\max}\end{aligned} \quad (11)$$

530 Consequently, the difference between any two consecutive states  $k$  and  $k-1$  are upper bounded by:

$$\begin{aligned}\|\boldsymbol{\theta}_t^{(k)} - \boldsymbol{\theta}_t^{(k-1)}\| &= \eta \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(k-1)}) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_t^{(k-2)})\| \\ &\leq \eta L_{\theta\phi} \|\boldsymbol{\phi}_t^{(k-1)} - \boldsymbol{\phi}_t^{(k-2)}\| \\ \|\boldsymbol{\phi}_t^{(k)} - \boldsymbol{\phi}_t^{(k-1)}\| &= \eta \|\nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(k-1)}, \boldsymbol{\phi}_t) - \nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_t^{(k-2)}, \boldsymbol{\phi}_t)\| \\ &\leq \eta L_{\phi\theta} \|\boldsymbol{\theta}_t^{(k-1)} - \boldsymbol{\theta}_t^{(k-2)}\| \\ \|\boldsymbol{\theta}_t^{(k)} - \boldsymbol{\theta}_t^{(k-1)}\| + \|\boldsymbol{\phi}_t^{(k)} - \boldsymbol{\phi}_t^{(k-1)}\| &\leq \eta \cdot (\eta L)^{(k-1)} \Delta_{\max}\end{aligned} \quad (12)$$

531 Since  $\|\boldsymbol{\omega}_t^{(k)} - \boldsymbol{\omega}_t^{(k-1)}\| \leq \|\boldsymbol{\theta}_t^{(k)} - \boldsymbol{\theta}_t^{(k-1)}\| + \|\boldsymbol{\phi}_t^{(k)} - \boldsymbol{\phi}_t^{(k-1)}\|$  we have:

$$\|\boldsymbol{\omega}_t^{(k)} - \boldsymbol{\omega}_t^{(k-1)}\| \leq \eta \cdot (\eta L)^{(k-1)} \Delta_{\max} \quad (13)$$

Suppose  $\eta < (2L)^{-1}$ , such that the difference between any two consecutive states is a contraction, then we consider the difference,  $\|\omega_t^{(a)} - \omega_t^{(b)}\|$ , where  $a > b > 0$ . We can rewrite it as:

$$\begin{aligned}
\|\omega_t^{(a)} - \omega_t^{(b)}\| &= \left\| \sum_{i=b+1}^a \omega_t^{(i)} - \omega_t^{(i-1)} \right\| \\
&\leq \sum_{i=b+1}^a \|\omega_t^{(i)} - \omega_t^{(i-1)}\| \\
&\leq \sum_{i=b+1}^a \eta \cdot (\eta L)^{(i-1)} \Delta_{\max} \\
&\leq \eta \Delta_{\max} \cdot [(\eta L)^{(b)} + \dots + (\eta L)^{(a-1)}] \\
&\leq \eta \Delta_{\max} \cdot (\eta L)^{(b-1)} \\
&\leq \eta^b L^{(b-1)} \Delta_{\max} = \mathcal{O}(\eta^b).
\end{aligned} \tag{14}$$

Since  $\eta < (2L)^{-1}$ , we have that  $\eta L < 1$  and for any  $\epsilon > 0$ , we can solve for  $b$  such that  $\eta^b L^{(b-1)} \Delta_{\max} < \epsilon$ . Therefore the sequence  $\{\omega_t^{(k)}\}_{k=0}^{\infty}$  is a Cauchy sequence. Moreover, in a complete space, every Cauchy sequence has a limit:  $\lim_{k \rightarrow \infty} \omega_t^{(k)} = \omega_t^*$   $\square$

## A.2 Proof of Theorem 5.1

**Theorem A.1.** Consider the (Minimax) problem under Assumption 3.1 and Lv.k GP. Let  $(\theta^*, \phi^*)$  be a stationary point. Suppose  $\theta_t - \theta^*$  not in kernel of  $\nabla_{\phi\theta} f(\theta^*, \phi^*)$ ,  $\phi_t - \phi^*$  not in kernel of  $\nabla_{\theta\phi} f(\theta^*, \phi^*)$  and  $\eta < (L)^{-1}$ . There exists a neighborhood  $\mathcal{U}$  of  $(\theta^*, \phi^*)$  such that if SPPM started at  $(\theta_0, \phi_0) \in \mathcal{U}$ , the iterates  $\{\theta_t, \phi_t\}_{t \geq 0}$  generated by SPPM satisfy:

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 \leq \frac{\rho^2(\mathbf{I} - \eta \nabla_{\theta\theta} f^*) \|\theta_t - \theta^*\|^2 + \rho^2(\mathbf{I} + \eta \nabla_{\phi\phi} f^*) \|\phi_t - \phi^*\|^2}{1 + \eta^2 \lambda_{\min}(\nabla_{\theta\phi} f^* \nabla_{\phi\theta} f^*)}$$

where  $f^* = f(\theta^*, \phi^*)$ . Moreover, for any  $\eta$  satisfying:

$$\frac{\max(\rho^2(\mathbf{I} - \eta \nabla_{\theta\theta} f^*), \rho^2(\mathbf{I} + \eta \nabla_{\phi\phi} f^*))}{1 + \eta^2 \lambda_{\min}(\nabla_{\theta\phi} f^* \nabla_{\phi\theta} f^*)} < 1, \tag{15}$$

SPPM converges asymptotically to  $(\theta^*, \phi^*)$ .

*Proof.* Consider the learning dynamics:

$$\begin{aligned}
\theta_{t+1} &= \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_{t+1}) \\
\phi_{t+1} &= \phi_t + \eta \nabla_{\phi} f(\theta_{t+1}, \phi_t)
\end{aligned}$$

Let us define

$$\begin{aligned}
\hat{\theta}_t &= \theta_t - \theta^* \\
\hat{\phi}_t &= \phi_t - \phi^*
\end{aligned}$$

It follows immediately by linearizing the system about the stationary point  $(\theta^*, \phi^*)$  that

$$\begin{bmatrix} \hat{\theta}_{t+1} \\ \hat{\phi}_{t+1} \end{bmatrix} \simeq \begin{bmatrix} \mathbf{I} - \eta \nabla_{\theta\theta}^2 f(\theta^*, \phi^*) & \mathbf{0} \\ \mathbf{0} & \mathbf{I} + \eta \nabla_{\phi\phi}^2 f(\theta^*, \phi^*) \end{bmatrix} \begin{bmatrix} \hat{\theta}_t \\ \hat{\phi}_t \end{bmatrix} + \begin{bmatrix} \mathbf{0} & -\eta \nabla_{\theta\phi}^2 f(\theta^*, \phi^*) \\ \eta \nabla_{\phi\theta}^2 f(\theta^*, \phi^*) & \mathbf{0} \end{bmatrix} \begin{bmatrix} \hat{\theta}_{t+1} \\ \hat{\phi}_{t+1} \end{bmatrix}$$

Let us denote the Jacobian by

$$\begin{bmatrix} -\nabla_{\theta\theta}^2 f(\theta^*, \phi^*) & -\nabla_{\theta\phi}^2 f(\theta^*, \phi^*) \\ \nabla_{\phi\theta}^2 f(\theta^*, \phi^*) & \nabla_{\phi\phi}^2 f(\theta^*, \phi^*) \end{bmatrix} = \begin{bmatrix} -\mathbf{A} & -\mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \tag{16}$$

548 Then we can rewrite the dynamics around the stationary point as

$$\begin{aligned}
\hat{\theta}_{t+1} &= \hat{\theta}_t - \eta A \hat{\theta}_t - \eta B \hat{\phi}_{t+1} \\
\hat{\theta}_{t+1} &= \hat{\theta}_t - \eta A \hat{\theta}_t - \eta B (\hat{\phi}_t + \eta B^T \hat{\theta}_{t+1} + \eta C \hat{\phi}_t) \\
(I + \eta^2 B B^T) \hat{\theta}_{t+1} &= (I - \eta A) \hat{\theta}_t - \eta B (I + \eta C) \phi_t \\
\hat{\theta}_{t+1} &= (I + \eta^2 B B^T)^{-1} \left[ (I - \eta A) \hat{\theta}_t - \eta B (I + \eta C) \phi_t \right] \tag{17}
\end{aligned}$$

549 Similarly, for the other player we have

$$\begin{aligned}
\hat{\phi}_{t+1} &= \hat{\phi}_t + \eta B^T \hat{\theta}_{t+1} + \eta C \hat{\phi}_t \\
\hat{\phi}_{t+1} &= \hat{\phi}_t + \eta B^T (\hat{\theta}_t - \eta A \hat{\theta}_t - \eta B \hat{\phi}_{t+1}) + \eta C \hat{\phi}_t \\
(I + \eta^2 B^T B) \hat{\phi}_{t+1} &= \eta B^T (I - \eta A) \hat{\theta}_t + (I + \eta C) \phi_t \\
\hat{\phi}_{t+1} &= (I + \eta^2 B^T B)^{-1} \left[ \eta B^T (I - \eta A) \hat{\theta}_t + (I + \eta C) \phi_t \right] \tag{18}
\end{aligned}$$

550 Let us define the symmetric matrices  $\mathbf{Q}_\theta = (I + \eta^2 B B^T)^{-1}$ ,  $\mathbf{Q}_\phi = (I + \eta^2 B^T B)^{-1}$  and  
551  $\mathbf{P}_\theta = (I - \eta A)$ ,  $\mathbf{P}_\phi = (I + \eta C)$ . Further we define  $r_t = \|\hat{\theta}_{t+1}\|^2 + \|\hat{\phi}_{t+1}\|^2$ . Based on these  
552 definitions, and the expressions in (53) and (54) we have

$$\begin{aligned}
\|\hat{\theta}_{t+1}\|^2 + \|\hat{\phi}_{t+1}\|^2 &= \|\mathbf{Q}_\theta \mathbf{P}_\theta \hat{\theta}_t\|^2 + \eta^2 \|\mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t\|^2 + \|\mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t\|^2 + \|\mathbf{Q}_\phi \mathbf{P}_\phi \hat{\phi}_t\|^2 \\
&\quad - 2\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^T \mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t + 2\eta \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{Q}_\phi^T \mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t \tag{19}
\end{aligned}$$

553 To simplify the expression in (55) we use the following lemma:

554 **Lemma A.1.** *The matrices  $\mathbf{Q}_\theta = (I + \eta^2 B B^T)^{-1}$ ,  $\mathbf{Q}_\phi = (I + \eta^2 B^T B)^{-1}$  satisfy the following*  
555 *properties:*

$$\mathbf{Q}_\theta B = B \mathbf{Q}_\phi \tag{20}$$

$$\mathbf{Q}_\phi B^T = B^T \mathbf{Q}_\theta \tag{21}$$

556 Using this lemma, we can show that

$$\hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^T \mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t = \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^T B \mathbf{Q}_\phi \mathbf{P}_\phi \hat{\phi}_t = \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{Q}_\phi^T B^T \mathbf{Q}_\theta \mathbf{P}_\theta \hat{\theta}_t = \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{Q}_\phi^T \mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t$$

557 where the intermediate equality holds as  $\mathbf{a}^T \mathbf{b} = \mathbf{b}^T \mathbf{a}$ . Hence, the expression in (55) can be simplified  
558 as

$$\|\hat{\theta}_{t+1}\|^2 + \|\hat{\phi}_{t+1}\|^2 = \|\mathbf{Q}_\theta \mathbf{P}_\theta \hat{\theta}_t\|^2 + \eta^2 \|\mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t\|^2 + \|\mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t\|^2 + \|\mathbf{Q}_\phi \mathbf{P}_\phi \hat{\phi}_t\|^2 \tag{22}$$

559 We simplify equation (58) as follows. Consider the term involving  $\hat{\theta}_t$ . We have

$$\begin{aligned}
\|\mathbf{Q}_\theta \mathbf{P}_\theta \hat{\theta}_t\|^2 + \eta^2 \|\mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t\|^2 &= \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^2 \mathbf{P}_\theta \hat{\theta}_t + \eta^2 \hat{\theta}_t^T \mathbf{P}_\theta^T B \mathbf{Q}_\phi^2 B^T \mathbf{P}_\theta \hat{\theta}_t \\
&= \hat{\theta}_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 B \mathbf{Q}_\phi^2 B^T) \mathbf{P}_\theta \hat{\theta}_t \\
&= \hat{\theta}_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 B \mathbf{Q}_\phi B^T \mathbf{Q}_\theta) \mathbf{P}_\theta \hat{\theta}_t \\
&= \hat{\theta}_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 B B^T \mathbf{Q}_\theta \mathbf{Q}_\theta) \mathbf{P}_\theta \hat{\theta}_t \\
&= \hat{\theta}_t^T \mathbf{P}_\theta^T (I + \eta^2 B B^T) \mathbf{Q}_\theta^2 \mathbf{P}_\theta \hat{\theta}_t \\
&= \hat{\theta}_t^T \mathbf{P}_\theta^T (I + \eta^2 B B^T)^{-1} \mathbf{P}_\theta \hat{\theta}_t \tag{23}
\end{aligned}$$

560 where the last equality follows by replacing  $\mathbf{Q}_\theta$  by its definition. The same procedure follows for the  
561 term involving  $\hat{\phi}_t$  which leads to the expression

$$\|\mathbf{Q}_\phi \mathbf{P}_\phi \hat{\phi}_t\|^2 + \eta^2 \|\mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t\|^2 = \hat{\phi}_t^T \mathbf{P}_\phi^T (I + \eta^2 B^T B)^{-1} \mathbf{P}_\phi \hat{\phi}_t. \tag{24}$$

562 Substitute  $\|\mathbf{Q}_\theta \mathbf{P}_\theta \hat{\theta}_t\|^2 + \eta^2 \|\mathbf{Q}_\phi B^T \mathbf{P}_\theta \hat{\theta}_t\|^2$  and  $\|\mathbf{Q}_\phi \mathbf{P}_\phi \hat{\phi}_t\|^2 + \eta^2 \|\mathbf{Q}_\theta B \mathbf{P}_\phi \hat{\phi}_t\|^2$  in (58) with the  
563 expressions in (23) and (24), respectively, to obtain

$$\|\hat{\theta}_{t+1}\|^2 + \|\hat{\phi}_{t+1}\|^2 = \hat{\theta}_t^T \mathbf{P}_\theta^T (I + \eta^2 B B^T)^{-1} \mathbf{P}_\theta \hat{\theta}_t + \hat{\phi}_t^T \mathbf{P}_\phi^T (I + \eta^2 B^T B)^{-1} \mathbf{P}_\phi \hat{\phi}_t. \tag{25}$$

564 Note that, we assume that the trajectory  $\{\hat{\theta}_t, \hat{\phi}_t\}_{t \geq 0}$  is not in the kernel of  $BB^T$  and  $B^TB$ , thus  
 565  $BB^T \hat{\theta}_t \neq 0$  and  $B^TB \hat{\phi}_t \neq 0$ . Now using the expression in (61) and the fact that  $P_\theta = P_\theta^T$ ,  
 566  $P_\phi = P_\phi^T$  and  $BB^T$  and  $B^TB$  have the same set of non-zero eigenvalues, if we denote the minimum  
 567 non-zero eigenvalues by  $\lambda_{\min}(BB^T)$  and  $\lambda_{\min}(B^TB)$ , we can write

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 \leq \frac{\rho^2(I - \eta A)\|\theta_t - \theta^*\|^2 + \rho^2(1 + \eta C)\|\phi_t - \phi^*\|^2}{I + \eta^2 \lambda_{\min}(B^TB)}.$$

568 Replacing  $\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2$  and  $\|\theta_t - \theta^*\|^2 + \|\phi_t - \phi^*\|^2$  with  $r_{t+1}$  and  $r_t$  we have:

$$r_{t+1} \leq \frac{\max(\rho^2(I - \eta A), \rho^2(I + \eta C))}{1 + \eta^2 \lambda_{\min}(B^TB)} r_t.$$

569 Recall that  $A = \nabla_{\theta\theta} f(\theta^*, \phi^*)$ ,  $B = \nabla_{\theta\phi} f(\theta^*, \phi^*)$  and  $C = \nabla_{\phi\phi} f(\theta^*, \phi^*)$ , therefore for any  $\eta$   
 570 satisfying that:

$$\frac{\max(\rho^2(I - \eta \nabla_{\theta\theta} f(\theta^*, \phi^*)), \rho^2(I + \eta \nabla_{\phi\phi} f(\theta^*, \phi^*)))}{1 + \eta^2 \lambda_{\min}(\nabla_{\theta\phi} f(\theta^*, \phi^*) \nabla_{\phi\theta} f(\theta^*, \phi^*))} < 1, \quad (26)$$

571 we have  $r_{t+1} < r_t$ . Since we linearize the system about the stationary point  $(\theta^*, \phi^*)$ , there exists  
 572 a neighborhood  $\mathcal{U}$  around the stationary point, such that, SPPM started at  $(\theta_0, \phi_0) \in \mathcal{U}$  converges  
 573 asymptotically to  $(\theta^*, \phi^*)$ .  $\square$

### 574 A.3 Proof of Remark 5.1

575 *Proof.* To prove the local convergence of Lv.k GP in non-convex non-concave games, we first  
 576 consider the update rule of Lv.k GP:

$$\text{Reasoning: } \begin{cases} \theta_t^{(k)} = \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t^{(k-1)}) \\ \phi_t^{(k)} = \phi_t - \eta \nabla_{\phi} g(\theta_t^{(k-1)}, \phi_t) \end{cases} \quad \text{Update: } \begin{cases} \theta_{t+1} = \theta_t^{(k)} \\ \phi_{t+1} = \phi_t^{(k)} \end{cases}$$

577 Similar to Section A.2 let us denote

$$\begin{bmatrix} -\nabla_{\theta\theta}^2 f(\theta^*, \phi^*) & -\nabla_{\theta\phi}^2 f(\theta^*, \phi^*) \\ \nabla_{\phi\theta}^2 f(\theta^*, \phi^*) & \nabla_{\phi\phi}^2 f(\theta^*, \phi^*) \end{bmatrix} = \begin{bmatrix} -A & -B \\ B^T & C \end{bmatrix}$$

578 and we define the difference between states and stationary points as

$$\begin{aligned} \hat{\theta}_t^{(k)} &= \theta_t^{(k)} - \theta^* \text{ and } \hat{\theta}_t = \theta_t - \theta^* \\ \hat{\phi}_t^{(k)} &= \phi_t^{(k)} - \phi^* \text{ and } \hat{\phi}_t = \phi_t - \phi^* \end{aligned}$$

579 Linearizing the dynamical system induced by Lv.k GP about the stationary point  $(\theta^*, \phi^*)$  we get:

$$\begin{cases} \hat{\theta}_{t+1} = \hat{\theta}_t^{(k)} = (I - \eta A)\hat{\theta}_t - \eta B \hat{\phi}_t^{(k-1)} \\ \hat{\phi}_{t+1} = \hat{\phi}_t^{(k)} = \eta B^T \hat{\theta}_t^{(k-1)} + (I + \eta C)\hat{\phi}_t \end{cases}$$

580 Note, in Lv.k GP, we define  $\theta_t^{(0)} = \theta_t$  and  $\phi_t^{(0)} = \phi_t$ , thus for Lv.1 GP, we have:

$$\begin{cases} \hat{\theta}_t^{(1)} = (I - \eta A)\hat{\theta}_t - \eta B \hat{\phi}_t \\ \hat{\phi}_t^{(1)} = \eta B^T \hat{\theta}_t + (I + \eta C)\hat{\phi}_t \end{cases}$$

581 For Lv.2 GP, we have:

$$\begin{cases} \hat{\theta}_t^{(2)} = (I - \eta A)\hat{\theta}_t - \eta B \hat{\phi}_t^{(1)} \\ \hat{\phi}_t^{(2)} = \eta B^T \hat{\theta}_t^{(1)} + (I + \eta C)\hat{\phi}_t \end{cases}$$

582 Substituting  $\hat{\theta}_t^{(1)}$  and  $\hat{\phi}_t^{(1)}$  into the update rule above we get:

$$\begin{cases} \hat{\theta}_t^{(2)} = (I - \eta A)\hat{\theta}_t - \eta B(I + \eta C)\hat{\phi}_t - \eta^2 BB^T \hat{\theta}_t \\ \hat{\phi}_t^{(2)} = \eta B^T(I - \eta A)\hat{\theta}_t + (I + \eta C)\hat{\phi}_t - \eta^2 B^T B \hat{\phi}_t \end{cases}$$

583 Similarly, for Lv.3 and Lv.4 GP we have:

$$\begin{cases} \hat{\theta}_t^{(3)} = (\mathbf{I} - \eta^2 \mathbf{B} \mathbf{B}^T)(\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t - \eta \mathbf{B}(\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t + \eta^3 \mathbf{B} \mathbf{B}^T \mathbf{B} \hat{\phi}_t \\ \hat{\phi}_t^{(3)} = \eta \mathbf{B}^T(\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t + (\mathbf{I} - \eta^2 \mathbf{B}^T \mathbf{B})(\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t - \eta^3 \mathbf{B} \mathbf{B}^T \mathbf{B} \hat{\theta}_t \end{cases}$$

584 and

$$\begin{cases} \hat{\theta}_t^{(4)} = (\mathbf{I} - \eta^2 \mathbf{B} \mathbf{B}^T) \left[ (\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t - \eta \mathbf{B}(\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t \right] + \eta^4 \mathbf{B} \mathbf{B}^T \mathbf{B} \mathbf{B}^T \hat{\theta}_t \\ \hat{\phi}_t^{(4)} = (\mathbf{I} - \eta^2 \mathbf{B}^T \mathbf{B}) \left[ \eta \mathbf{B}^T(\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t + (\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t \right] + \eta^4 \mathbf{B}^T \mathbf{B} \mathbf{B}^T \mathbf{B} \hat{\phi}_t \end{cases}$$

585 Summarizing the equations above we have that for Lv.2k GP, its update can be written as:

$$\begin{cases} \hat{\theta}_t^{(2k)} = (\sum_{i=0}^{k-1} (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \left[ (\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t - \eta \mathbf{B}(\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t \right] + (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \\ \hat{\phi}_t^{(2k)} = (\sum_{i=0}^{k-1} (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \left[ \eta \mathbf{B}^T(\mathbf{I} - \eta \mathbf{A})\hat{\theta}_t + (\mathbf{I} + \eta \mathbf{C})\hat{\phi}_t \right] + (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \end{cases}$$

586 Similar to Appendix A.2, let us define  $\mathbf{Q}_\theta = (\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1}$ ,  $\mathbf{Q}_\phi = (\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1}$  and  
 587  $\mathbf{P}_\theta = (\mathbf{I} - \eta \mathbf{A})$ ,  $\mathbf{P}_\phi = (\mathbf{I} + \eta \mathbf{C})$ . Further, we define  $\mathbf{R}_\theta^{(k)} = (\sum_{i=0}^{k-1} (-\eta^2 \mathbf{B} \mathbf{B}^T)^k)$ ,  $\mathbf{R}_\phi^{(k)} =$   
 588  $(\sum_{i=0}^{k-1} (-\eta^2 \mathbf{B}^T \mathbf{B})^k)$  and  $\mathbf{E}_\theta^{(k)} = \mathbf{R}_\theta^{(k)} - \mathbf{Q}_\theta$ ,  $\mathbf{E}_\phi^{(k)} = \mathbf{R}_\phi^{(k)} - \mathbf{Q}_\phi$ .

589 Since  $\eta < L^{-1}$ , we have that:

$$\begin{aligned} \mathbf{Q}_\theta &= \sum_{i=0}^{\infty} (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \\ \mathbf{Q}_\phi &= \sum_{i=0}^{\infty} (-\eta^2 \mathbf{B}^T \mathbf{B})^k \end{aligned}$$

590 and

$$\mathbf{E}_\theta^{(k)} = \mathbf{R}_\theta^{(k)} - \mathbf{Q}_\theta = - \sum_{i=k}^{\infty} (-\eta^2 \mathbf{B} \mathbf{B}^T)^k = -(\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1} \cdot (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \quad (27)$$

$$\mathbf{E}_\phi^{(k)} = \mathbf{R}_\phi^{(k)} - \mathbf{Q}_\phi = - \sum_{i=k}^{\infty} (-\eta^2 \mathbf{B}^T \mathbf{B})^k = -(\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1} \cdot (-\eta^2 \mathbf{B}^T \mathbf{B})^k \quad (28)$$

591 Also, from Lemma A.1 and the definition of the error terms, it can be verified that

$$\mathbf{E}_\theta^{(k)} \mathbf{B} = \mathbf{B} \mathbf{E}_\phi^{(k)} \quad (29)$$

$$\mathbf{E}_\phi^{(k)} \mathbf{B}^T = \mathbf{B}^T \mathbf{E}_\theta^{(k)} \quad (30)$$

592 Then we can rewrite the update rule of Lv.2k GP:

$$\begin{cases} \hat{\theta}_t^{(2k)} = (\mathbf{Q}_\theta + \mathbf{E}_\theta^{(k)}) \left[ \mathbf{P}_\theta \hat{\theta}_t - \eta \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \right] + (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \\ \hat{\phi}_t^{(2k)} = (\mathbf{Q}_\phi + \mathbf{E}_\phi^{(k)}) \left[ \eta \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t + \mathbf{P}_\phi \hat{\phi}_t \right] + (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \end{cases}$$

593 Let us consider the following sum:

$$\begin{aligned} &\|\hat{\theta}_t^{(2k)} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 + \|\hat{\phi}_t^{(2k)} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2 \\ &= \left\| (\mathbf{Q}_\theta + \mathbf{E}_\theta^{(k)}) \left[ \mathbf{P}_\theta \hat{\theta}_t - \eta \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \right] \right\|^2 + \left\| (\mathbf{Q}_\phi + \mathbf{E}_\phi^{(k)}) \left[ \eta \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t + \mathbf{P}_\phi \hat{\phi}_t \right] \right\|^2 \end{aligned} \quad (31)$$

594 The R.H.S. of Eq.(31) can be written as:

$$\begin{aligned} &\left\| (\mathbf{Q}_\theta + \mathbf{E}_\theta^{(k)}) \left[ \mathbf{P}_\theta \hat{\theta}_t - \eta \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \right] \right\|^2 + \left\| (\mathbf{Q}_\phi + \mathbf{E}_\phi^{(k)}) \left[ \eta \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t + \mathbf{P}_\phi \hat{\phi}_t \right] \right\|^2 \\ &= \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^2 \mathbf{P}_\theta \hat{\theta}_t - 2\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^2 \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t + \eta^2 \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{B}^T \mathbf{Q}_\theta^2 \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \\ &\quad + 2\hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta \mathbf{E}_\theta^{(k)} \mathbf{P}_\theta \hat{\theta}_t - 4\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta \mathbf{E}_\theta^{(k)} \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t + 2\eta^2 \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{B}^T \mathbf{Q}_\theta \mathbf{E}_\theta^{(k)} \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \\ &\quad + \hat{\theta}_t^T \mathbf{P}_\theta^T [\mathbf{E}_\theta^{(k)}]^2 \mathbf{P}_\theta \hat{\theta}_t - 2\eta \hat{\theta}_t^T \mathbf{P}_\theta^T [\mathbf{E}_\theta^{(k)}] \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t + \eta^2 \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{B}^T [\mathbf{E}_\theta^{(k)}] \mathbf{B} \mathbf{P}_\phi \hat{\phi}_t \\ &\quad + \hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{Q}_\phi^2 \mathbf{P}_\phi \hat{\phi}_t + 2\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} \mathbf{Q}_\phi^2 \mathbf{P}_\phi \hat{\phi}_t + \eta^2 \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} \mathbf{Q}_\phi^2 \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t \\ &\quad + 2\hat{\phi}_t^T \mathbf{P}_\phi^T \mathbf{Q}_\phi \mathbf{E}_\phi^{(k)} \mathbf{P}_\phi \hat{\phi}_t + 4\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} \mathbf{Q}_\phi \mathbf{E}_\phi^{(k)} \mathbf{P}_\phi \hat{\phi}_t + 2\eta^2 \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} \mathbf{Q}_\phi \mathbf{E}_\phi^{(k)} \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t \\ &\quad + \hat{\phi}_t^T \mathbf{P}_\phi^T [\mathbf{E}_\phi^{(k)}]^2 \mathbf{P}_\phi \hat{\phi}_t + 2\eta \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} [\mathbf{E}_\phi^{(k)}] \mathbf{P}_\phi \hat{\phi}_t + \eta^2 \hat{\theta}_t^T \mathbf{P}_\theta^T \mathbf{B} [\mathbf{E}_\phi^{(k)}] \mathbf{B}^T \mathbf{P}_\theta \hat{\theta}_t \end{aligned} \quad (32)$$

595 Now, before adding all terms in Eq.(32), note that all of the cross terms in Eq.(32) cancel out.

596 For instance, using Lemma A.1 and Eq.(29), Eq.(30) we can show that

$$\begin{aligned} 4\eta\hat{\theta}_t^T P_\theta^T BQ_\phi E_\phi^{(k)} P_\phi \hat{\phi}_t - 4\eta\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t &= 4\eta\hat{\theta}_t^T P_\theta^T Q_\theta B E_\phi^{(k)} P_\phi \hat{\phi}_t - 4\eta\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t \\ &= 4\eta\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t - 4\eta\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t \\ &= 0 \end{aligned}$$

597 By using similar arguments it can be shown that terms in Eq.(32) leads to:

$$\begin{aligned} &\left\| (Q_\theta + E_\theta^{(k)}) \left[ P_\theta \hat{\theta}_t - \eta B P_\phi \hat{\phi}_t \right] \right\|^2 + \left\| (Q_\phi + E_\phi^{(k)}) \left[ \eta B^T P_\theta \hat{\theta}_t + P_\phi \hat{\phi}_t \right] \right\|^2 \\ &= \hat{\theta}_t^T P_\theta^T Q_\theta^2 P_\theta \hat{\theta}_t + \eta^2 \hat{\phi}_t^T P_\phi^T B^T Q_\theta^2 B P_\phi \hat{\phi}_t \\ &\quad + 2\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} P_\theta \hat{\theta}_t + 2\eta^2 \hat{\phi}_t^T P_\phi^T B^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t \\ &\quad + \hat{\theta}_t^T P_\theta^T [E_\theta^{(k)}]^2 P_\theta \hat{\theta}_t + \eta^2 \hat{\phi}_t^T P_\phi^T B^T [E_\theta^{(k)}]^2 B P_\phi \hat{\phi}_t \\ &\quad + \hat{\phi}_t^T P_\phi^T Q_\phi^2 P_\phi \hat{\phi}_t + \eta^2 \hat{\theta}_t^T P_\theta^T B Q_\phi^2 B^T P_\theta \hat{\theta}_t \\ &\quad + 2\hat{\phi}_t^T P_\phi^T Q_\phi E_\phi^{(k)} P_\phi \hat{\phi}_t + 2\eta^2 \hat{\theta}_t^T P_\theta^T B Q_\phi E_\phi^{(k)} B^T P_\theta \hat{\theta}_t \\ &\quad + \hat{\phi}_t^T P_\phi^T [E_\phi^{(k)}]^2 P_\phi \hat{\phi}_t + \eta^2 \hat{\theta}_t^T P_\theta^T B [E_\phi^{(k)}]^2 B^T P_\theta \hat{\theta}_t \end{aligned} \quad (33)$$

598 Similar to Eq.(23) we have the following simplification:

$$\begin{aligned} \hat{\theta}_t^T P_\theta^T Q_\theta^2 P_\theta \hat{\theta}_t + \eta^2 \hat{\theta}_t^T P_\theta^T B Q_\phi^2 B^T P_\theta \hat{\theta}_t &= \hat{\theta}_t^T P_\theta^T Q_\theta P_\theta \hat{\theta}_t \\ \hat{\phi}_t^T P_\phi^T Q_\phi^2 P_\phi \hat{\phi}_t + \eta^2 \hat{\phi}_t^T P_\phi^T B^T Q_\theta^2 B P_\phi \hat{\phi}_t &= \hat{\phi}_t^T P_\phi^T Q_\phi P_\phi \hat{\phi}_t \\ 2\hat{\theta}_t^T P_\theta^T Q_\theta E_\theta^{(k)} P_\theta \hat{\theta}_t + 2\eta^2 \hat{\theta}_t^T P_\theta^T B Q_\phi E_\phi^{(k)} B^T P_\theta \hat{\theta}_t &= 2\hat{\theta}_t^T P_\theta^T E_\theta^{(k)} P_\theta \hat{\theta}_t \\ 2\hat{\phi}_t^T P_\phi^T Q_\phi E_\phi^{(k)} P_\phi \hat{\phi}_t + 2\eta^2 \hat{\phi}_t^T P_\phi^T B^T Q_\theta E_\theta^{(k)} B P_\phi \hat{\phi}_t &= 2\hat{\phi}_t^T P_\phi^T E_\phi^{(k)} P_\phi \hat{\phi}_t \end{aligned}$$

599 Now we can further simplify Eq.(32) as:

$$\begin{aligned} &\left\| (Q_\theta + E_\theta^{(k)}) \left[ P_\theta \hat{\theta}_t - \eta B P_\phi \hat{\phi}_t \right] \right\|^2 + \left\| (Q_\phi + E_\phi^{(k)}) \left[ \eta B^T P_\theta \hat{\theta}_t + P_\phi \hat{\phi}_t \right] \right\|^2 \\ &= (P_\theta \hat{\theta}_t)^T \left[ Q_\theta + 2E_\theta^{(k)} + [E_\theta^{(k)}]^2 + \eta^2 B [E_\phi^{(k)}]^2 B^T \right] (P_\theta \hat{\theta}_t) \\ &\quad + (P_\phi \hat{\phi}_t)^T \left[ Q_\phi + 2E_\phi^{(k)} + [E_\phi^{(k)}]^2 + \eta^2 B [E_\theta^{(k)}]^2 B^T \right] (P_\phi \hat{\phi}_t) \end{aligned}$$

600 Using Eq.(27) and Eq.(28) and definition of  $Q_\theta$  and  $Q_\phi$  we have:

$$\begin{aligned} &(P_\theta \hat{\theta}_t)^T \left[ Q_\theta + 2E_\theta^{(k)} + [E_\theta^{(k)}]^2 + \eta^2 B [E_\phi^{(k)}]^2 B^T \right] (P_\theta \hat{\theta}_t) \\ &= (P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T)^{-1} (P_\theta \hat{\theta}_t) - 2(P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T)^{-1} (-\eta^2 B B^T)^k (P_\theta \hat{\theta}_t) \\ &\quad + (P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T) (I + \eta^2 B B^T)^{-2} (-\eta^2 B B^T)^{2k} (P_\theta \hat{\theta}_t) \\ &= (P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T)^{-1} (I - 2(-\eta^2 B B^T)^k + (-\eta^2 B B^T)^{2k}) (P_\theta \hat{\theta}_t) \\ &= ((I - (-\eta^2 B B^T)^k) P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T)^{-1} ((I - (-\eta^2 B B^T)^k) P_\theta \hat{\theta}_t) \end{aligned}$$

601 Similarly, we have that

$$\begin{aligned} &(P_\phi \hat{\phi}_t)^T \left[ Q_\phi + 2E_\phi^{(k)} + [E_\phi^{(k)}]^2 + \eta^2 B^T [E_\theta^{(k)}]^2 B \right] (P_\phi \hat{\phi}_t) \\ &= ((I - (-\eta^2 B^T B)^k) P_\phi \hat{\phi}_t)^T (I + \eta^2 B^T B)^{-1} ((I - (-\eta^2 B^T B)^k) P_\phi \hat{\phi}_t) \end{aligned}$$

602 Thus we simplify the R.H.S. of Eq.(31) as

$$\begin{aligned} &\left\| (Q_\theta + E_\theta^{(k)}) \left[ P_\theta \hat{\theta}_t - \eta B P_\phi \hat{\phi}_t \right] \right\|^2 + \left\| (Q_\phi + E_\phi^{(k)}) \left[ \eta B^T P_\theta \hat{\theta}_t + P_\phi \hat{\phi}_t \right] \right\|^2 \\ &= ((I - (-\eta^2 B B^T)^k) P_\theta \hat{\theta}_t)^T (I + \eta^2 B B^T)^{-1} ((I - (-\eta^2 B B^T)^k) P_\theta \hat{\theta}_t) \\ &\quad + ((I - (-\eta^2 B^T B)^k) P_\phi \hat{\phi}_t)^T (I + \eta^2 B^T B)^{-1} ((I - (-\eta^2 B^T B)^k) P_\phi \hat{\phi}_t) \end{aligned} \quad (34)$$

603 Let us consider the L.H.S. of Eq.(31)

$$\begin{aligned}
& \|\hat{\theta}_t^{(2k)} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 + \|\hat{\phi}_t^{(2k)} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2 \\
&= \|\hat{\theta}_t^{(2k)}\|^2 - 2\langle \hat{\theta}_t^{(2k)}, (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \rangle + \|(-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 \\
&+ \|\hat{\phi}_t^{(2k)}\|^2 - 2\langle \hat{\phi}_t^{(2k)}, (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \rangle + \|(-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2
\end{aligned} \tag{35}$$

604 Substituting Eq.(35) and Eq.(34) into L.H.S. and R.H.S. of Eq.(31) respectively we get:

$$\begin{aligned}
& \|\hat{\theta}_t^{(2k)}\|^2 - 2\langle \hat{\theta}_t^{(2k)}, (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \rangle + \|(-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 \\
&+ \|\hat{\phi}_t^{(2k)}\|^2 - 2\langle \hat{\phi}_t^{(2k)}, (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \rangle + \|(-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2 \\
&= ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t)^T (\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t) \\
&+ ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t)^T (\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t)
\end{aligned}$$

605 Now we have the following equation:

$$\begin{aligned}
& \|\hat{\theta}_t^{(2k)}\|^2 + \|\hat{\phi}_t^{(2k)}\|^2 \\
&= ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t)^T (\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t) \\
&+ 2\langle \hat{\theta}_t^{(2k)}, (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \rangle - \|(-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 \\
&+ ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t)^T (\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t) \\
&+ 2\langle \hat{\phi}_t^{(2k)}, (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \rangle - \|(-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2
\end{aligned}$$

606 Note that

$$2\langle \hat{\theta}_t^{(2k)}, (-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \rangle = 2\langle (-\eta^2 \mathbf{B} \mathbf{B}^T)^{\frac{k}{2}} \hat{\theta}_t^{(2k)}, (\eta^2 \mathbf{B} \mathbf{B}^T)^{\frac{k}{2}} \hat{\theta}_t \rangle \tag{36}$$

$$\leq \eta^{2k} \hat{\theta}_t^{(2k)T} (\mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t^{(2k)} + \eta^{2k} \hat{\theta}_t^T (\mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t \tag{37}$$

607 Similarly

$$2\langle \hat{\phi}_t^{(2k)}, (-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \rangle = 2\langle (-\eta^2 \mathbf{B}^T \mathbf{B})^{\frac{k}{2}} \hat{\phi}_t^{(2k)}, (\eta^2 \mathbf{B}^T \mathbf{B})^{\frac{k}{2}} \hat{\phi}_t \rangle \tag{38}$$

$$\leq \eta^{2k} \hat{\phi}_t^{(2k)T} (\mathbf{B}^T \mathbf{B})^k \hat{\phi}_t^{(2k)} + \eta^{2k} \hat{\phi}_t^T (\mathbf{B}^T \mathbf{B})^k \hat{\phi}_t \tag{39}$$

608 Summing everything together we have:

$$\begin{aligned}
& (\hat{\theta}_t^{(2k)})^T (\mathbf{I} - (\eta^2 \mathbf{B} \mathbf{B}^T)^k) (\hat{\theta}_t^{(2k)}) + (\hat{\phi}_t^{(2k)})^T (\mathbf{I} - (\eta^2 \mathbf{B}^T \mathbf{B})^k) (\hat{\phi}_t^{(2k)}) \\
&\leq ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t)^T (\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B} \mathbf{B}^T)^k) \mathbf{P}_\theta \hat{\theta}_t) \\
&+ (\hat{\theta}_t)^T (\eta^2 \mathbf{B} \mathbf{B}^T)^k (\hat{\theta}_t) - \|(-\eta^2 \mathbf{B} \mathbf{B}^T)^k \hat{\theta}_t\|^2 \\
&+ ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t)^T (\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1} ((\mathbf{I} - (-\eta^2 \mathbf{B}^T \mathbf{B})^k) \mathbf{P}_\phi \hat{\phi}_t) \\
&+ (\hat{\phi}_t)^T (\eta^2 \mathbf{B}^T \mathbf{B})^k (\hat{\phi}_t) - \|(-\eta^2 \mathbf{B}^T \mathbf{B})^k \hat{\phi}_t\|^2
\end{aligned}$$

609 Note that, we assume that the trajectory  $\{\hat{\theta}_t, \hat{\phi}_t\}_{t \geq 0}$  is not in the kernel of  $\mathbf{B} \mathbf{B}^T$  and  $\mathbf{B}^T \mathbf{B}$ , thus  
610  $\mathbf{B} \mathbf{B}^T \hat{\theta}_t \neq 0$  and  $\mathbf{B}^T \mathbf{B} \hat{\phi}_t \neq 0$ . Now using the expression in (61) and the fact that  $\mathbf{P}_\theta = \mathbf{P}_\theta^T$ ,  
611  $\mathbf{P}_\phi = \mathbf{P}_\phi^T$  and  $\mathbf{B} \mathbf{B}^T$  and  $\mathbf{B}^T \mathbf{B}$  have the same set of non-zero eigenvalues, if we denote the minimum

612 non-zero eigenvalues by  $\lambda_{\min}(\mathbf{B}\mathbf{B}^T)$  and  $\lambda_{\min}(\mathbf{B}^T\mathbf{B})$ , we can write

$$\begin{aligned}
& (1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k) \|\hat{\boldsymbol{\theta}}_t^{(2k)}\|^2 + (1 - (\eta^2 \lambda_{\max}(\mathbf{B}^T\mathbf{B}))^k) \|\hat{\boldsymbol{\phi}}_t^{(2k)}\|^2 \\
& \leq \frac{(1 - (-\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k)^2 \rho^2 (1 - \eta \mathbf{A}) \|\hat{\boldsymbol{\theta}}_t\|^2}{1 + \eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T)} + (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k \|\hat{\boldsymbol{\theta}}_t\|^2 - (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^{2k} \|\hat{\boldsymbol{\theta}}_t\|^2 \\
& + \frac{(1 - (-\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k)^2 \rho^2 (1 + \eta \mathbf{C}) \|\hat{\boldsymbol{\phi}}_t\|^2}{1 + \eta^2 \lambda_{\min}(\mathbf{B}^T\mathbf{B})} + (\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k \|\hat{\boldsymbol{\phi}}_t\|^2 - (\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^{2k} \|\hat{\boldsymbol{\phi}}_t\|^2 \\
& \leq \frac{(1 - (-\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k)^2 \rho^2 (1 - \eta \mathbf{A}) \|\hat{\boldsymbol{\theta}}_t\|^2}{1 + \eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T)} + (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k (1 - (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k) \|\hat{\boldsymbol{\theta}}_t\|^2 \\
& + \frac{(1 - (-\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k)^2 \rho^2 (1 + \eta \mathbf{C}) \|\hat{\boldsymbol{\phi}}_t\|^2}{1 + \eta^2 \lambda_{\min}(\mathbf{B}^T\mathbf{B})} + (\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k \|\hat{\boldsymbol{\phi}}_t\|^2 - (\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^{2k} \|\hat{\boldsymbol{\phi}}_t\|^2 \\
& \leq \frac{(1 - (-\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k)^2 \rho^2 (1 - \eta \mathbf{A}) \|\hat{\boldsymbol{\theta}}_t\|^2}{(1 + \eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T))(1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k)} + \frac{(\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k (1 - (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k) \|\hat{\boldsymbol{\theta}}_t\|^2}{(1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k)} \\
& + \frac{(1 - (-\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k)^2 \rho^2 (1 + \eta \mathbf{C}) \|\hat{\boldsymbol{\phi}}_t\|^2}{(1 + \eta^2 \lambda_{\min}(\mathbf{B}^T\mathbf{B}))(1 - (\eta^2 \lambda_{\max}(\mathbf{B}^T\mathbf{B}))^k)} + \frac{(\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k (1 - (\eta^2 \lambda(\mathbf{B}^T\mathbf{B}))^k) \|\hat{\boldsymbol{\phi}}_t\|^2}{(1 - (\eta^2 \lambda_{\max}(\mathbf{B}^T\mathbf{B}))^k)}
\end{aligned} \tag{40}$$

613 Let us define the distance as:

$$\|r_t^{(k)}\|^2 = \|\hat{\boldsymbol{\theta}}_t^{(k)}\|^2 + \|\hat{\boldsymbol{\phi}}_t^{(k)}\|^2 \tag{41}$$

$$\|r_t\|^2 = \|\hat{\boldsymbol{\theta}}_t\|^2 + \|\hat{\boldsymbol{\phi}}_t\|^2 \tag{42}$$

614 Then we have

$$\begin{aligned}
\|r_t^{(2k)}\| & \leq \frac{(1 - (-\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k) (\rho^2 (1 - \eta \mathbf{A}) \|\hat{\boldsymbol{\theta}}_t\|^2 + \rho^2 (1 + \eta \mathbf{C}) \|\hat{\boldsymbol{\phi}}_t\|^2)}{(1 + \eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T))(1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k)} + \frac{(\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k (1 - (\eta^2 \lambda(\mathbf{B}\mathbf{B}^T))^k) \|r_t\|^2}{(1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k)} \\
& = a \left( \frac{\rho^2 (1 - \eta \mathbf{A}) \|\hat{\boldsymbol{\theta}}_t\|^2 + \rho^2 (1 + \eta \mathbf{C}) \|\hat{\boldsymbol{\phi}}_t\|^2}{1 + \eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T)} \right) + b
\end{aligned} \tag{43}$$

615 where

$$a = \begin{cases} \frac{1 + (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k}{1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k} & \text{for odd } k, \\ \frac{1 - (\eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T))^k}{1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k} & \text{for even } k \end{cases}$$

616 and

$$b = \frac{(\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k (1 - (\eta^2 \lambda_{\min}(\mathbf{B}\mathbf{B}^T))^k)}{1 - (\eta^2 \lambda_{\max}(\mathbf{B}\mathbf{B}^T))^k} \tag{44}$$

617  $\square$

### 618 A.3.1 Proof of Lemma A.1

619 Let  $\mathbf{B} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$  be the singular value decomposition of  $\mathbf{B}$ . Here  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal matrices  
620 and  $\boldsymbol{\Sigma}$  is a rectangular diagonal matrix. Then we have:

$$\begin{aligned}
\mathbf{Q}_\theta \mathbf{B} &= (\mathbf{I} + \eta^2 \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \mathbf{V}\boldsymbol{\Sigma}^T \mathbf{U}^T)^{-1} \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \\
&= (\mathbf{U}(\eta^2 \boldsymbol{\Sigma}\boldsymbol{\Sigma}^T + \mathbf{I})\mathbf{U}^T)^{-1} \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \\
&= \mathbf{U}(\eta^2 \boldsymbol{\Sigma}\boldsymbol{\Sigma}^T + \mathbf{I})^{-1} \mathbf{U}^T \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \\
&= \mathbf{U}(\eta^2 \boldsymbol{\Sigma}\boldsymbol{\Sigma}^T + \mathbf{I})^{-1} \boldsymbol{\Sigma}\mathbf{V}^T
\end{aligned}$$

621 Here we used the fact that  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal matrices. Now, we simplify the other side to  
622 get:

$$\begin{aligned}
\mathbf{B}\mathbf{Q}_\phi &= \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T (\mathbf{I} + \eta^2 \mathbf{V}\boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T)^{-1} \\
&= \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T (\mathbf{V}(\eta^2 \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \mathbf{I})\mathbf{V}^T)^{-1} \\
&= \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \mathbf{V}(\eta^2 \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \mathbf{I})^{-1} \mathbf{V}^T \\
&= \mathbf{U}\boldsymbol{\Sigma}(\eta^2 \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \mathbf{I})^{-1} \mathbf{V}^T
\end{aligned}$$



Now we consider the following equation:

$$\eta^2 \Sigma \Sigma^T \Sigma + \Sigma = \Sigma(\eta^2 \Sigma^T \Sigma + \mathbf{I}) = (\eta^2 \Sigma \Sigma^T + \mathbf{I}) \Sigma \quad (45)$$

which indicates that  $\Sigma(\eta^2 \Sigma^T \Sigma + \mathbf{I}) = (\eta^2 \Sigma \Sigma^T + \mathbf{I}) \Sigma$ . Multiplying both sides of this equation by  $(\eta^2 \Sigma \Sigma^T + \mathbf{I})^{-1}$  and  $(\eta^2 \Sigma^T \Sigma + \mathbf{I})^{-1}$  we have:

$$\begin{aligned} (\eta^2 \Sigma \Sigma^T + \mathbf{I})^{-1} \Sigma (\eta^2 \Sigma^T \Sigma + \mathbf{I}) (\eta^2 \Sigma^T \Sigma + \mathbf{I})^{-1} &= (\eta^2 \Sigma \Sigma^T + \mathbf{I})^{-1} (\eta^2 \Sigma \Sigma^T + \mathbf{I}) \Sigma (\eta^2 \Sigma^T \Sigma + \mathbf{I})^{-1} \\ (\eta^2 \Sigma \Sigma^T + \mathbf{I})^{-1} \Sigma &= \Sigma (\eta^2 \Sigma^T \Sigma + \mathbf{I})^{-1} \end{aligned}$$

Therefore, we have  $\mathbf{Q}_\theta \mathbf{B} = \mathbf{B} \mathbf{Q}_\phi$ . Using a similar argument, we can also prove the equality in Equation (21).

#### A.4 Theorem 5.1 without kernel assumption

**Theorem A.2.** Consider the (Minimax) problem under Assumption 3.1 and Lvk GP. Let  $(\theta^*, \phi^*)$  be a stationary point. Suppose  $\eta < (L)^{-1}$ . There exists a neighborhood  $\mathcal{U}$  of  $(\theta^*, \phi^*)$  such that if SPPM started at  $(\theta_0, \phi_0) \in \mathcal{U}$ , the iterates  $\{\theta_t, \phi_t\}_{t \geq 0}$  generated by SPPM satisfy:

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 \leq \frac{\rho^2(\mathbf{I} - \eta \mathbf{A}) \|\theta_t - \theta^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B} \mathbf{B}^T)} + \frac{\rho^2(1 + \eta \mathbf{C}) \|\phi_t - \phi^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B}^T \mathbf{B})}.$$

where  $f^* = f(\theta^*, \phi^*)$ . Moreover, for any  $\eta$  satisfying:

$$\frac{\rho^2(\mathbf{I} - \eta \mathbf{A}) \|\theta_t - \theta^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B} \mathbf{B}^T)} + \frac{\rho^2(1 + \eta \mathbf{C}) \|\phi_t - \phi^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B}^T \mathbf{B})} < \|\theta_t - \theta^*\|^2 + \|\phi_t - \phi^*\|^2 \quad (46)$$

SPPM converges asymptotically to  $(\theta^*, \phi^*)$ .

*Proof.* Following the same setting and procedure as in Appendix A.2, we have that

$$\|\hat{\theta}_{t+1}\|^2 + \|\hat{\phi}_{t+1}\|^2 = \hat{\theta}_t^T \mathbf{P}_\theta^T (\mathbf{I} + \eta^2 \mathbf{B} \mathbf{B}^T)^{-1} \mathbf{P}_\theta \hat{\theta}_t + \hat{\phi}_t^T \mathbf{P}_\phi^T (\mathbf{I} + \eta^2 \mathbf{B}^T \mathbf{B})^{-1} \mathbf{P}_\phi \hat{\phi}_t \quad (47)$$

Now using the fact that  $\mathbf{P}_\theta = \mathbf{P}_\theta^T$ ,  $\mathbf{P}_\phi = \mathbf{P}_\phi^T$ , we can write

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 \leq \frac{\rho^2(\mathbf{I} - \eta \mathbf{A}) \|\theta_t - \theta^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B} \mathbf{B}^T)} + \frac{\rho^2(1 + \eta \mathbf{C}) \|\phi_t - \phi^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B}^T \mathbf{B})}.$$

For any  $\eta$  that satisfying

$$\frac{\rho^2(\mathbf{I} - \eta \mathbf{A}) \|\theta_t - \theta^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B} \mathbf{B}^T)} + \frac{\rho^2(1 + \eta \mathbf{C}) \|\phi_t - \phi^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{B}^T \mathbf{B})} < \|\theta_t - \theta^*\|^2 + \|\phi_t - \phi^*\|^2 \quad (48)$$

we have that

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 < \|\theta_t - \theta^*\|^2 + \|\phi_t - \phi^*\|^2 \quad (49)$$

i.e., SPPM converges asymptotically towards  $(\theta^*, \phi^*)$ .  $\square$

#### A.5 Proof of Theorem 5.2

*Proof.* In order to proof the convergence of SPPM in bilinear games, we first show that the SPPM update rule is equivalent to that of the following Proximal Point Method:

$$\begin{cases} \theta_{t+1} = \theta_t - \eta \nabla_\theta f(\theta_{t+1}, \phi_{t+1}) \\ \phi_{t+1} = \phi_t + \eta \nabla_\phi f(\theta_{t+1}, \phi_{t+1}) \end{cases} \quad (50)$$

In the Bilinear game, the SPPM update is:

$$\begin{aligned} &\begin{cases} \theta_{t+1} = \theta_t - \eta \nabla_\theta f(\theta_t, \phi_{t+1}) \\ \phi_{t+1} = \phi_t + \eta \nabla_\phi f(\theta_t, \phi_{t+1}) \end{cases} \\ &= \begin{cases} \theta_{t+1} = \theta_t - \eta \mathbf{M} \phi_{t+1} \\ \phi_{t+1} = \phi_t + \eta \mathbf{M}^T \theta_{t+1} \end{cases} \end{aligned}$$

643 and the PPM update is:

$$\begin{aligned} & \begin{cases} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_{t+1}, \boldsymbol{\phi}_{t+1}) \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta \nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_{t+1}, \boldsymbol{\phi}_{t+1}) \end{cases} \\ & = \begin{cases} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{M} \boldsymbol{\phi}_{t+1} \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta \mathbf{M}^T \boldsymbol{\theta}_{t+1} \end{cases} \end{aligned}$$

644 Thus SPPM and PPM are equivalent in the Bilinear game. The convergence result of PPM in bilinear  
645 games has been proved in Theorem 2 of [49]:

646 **Theorem A.3.** *Consider the Bilinear game and the PPM method. Further, we define  $r_t = \|\boldsymbol{\theta}_t -$   
647  $\boldsymbol{\theta}^*\|^2 + \|\boldsymbol{\phi}_t - \boldsymbol{\phi}^*\|^2$ . Then, for any  $\eta > 0$ , the iterates  $\{\boldsymbol{\theta}_t, \boldsymbol{\phi}_t\}_{t \geq 0}$  generated by SPPM satisfy*

$$r_{t+1} \leq \frac{1}{1 + \eta^2 \lambda_{\min}(\mathbf{M}^T \mathbf{M})} r_t. \quad (51)$$

648 Therefore, SPPM and PPM have the same convergence property in bilinear games.  $\square$

### 649 A.6 Proof of Theorem 5.3

650 *Proof.* Consider the learning dynamics:

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_t, \boldsymbol{\phi}_{t+1}) \\ \boldsymbol{\phi}_{t+1} &= \boldsymbol{\phi}_t + \eta \nabla_{\boldsymbol{\phi}} f(\boldsymbol{\theta}_{t+1}, \boldsymbol{\phi}_t) \end{aligned}$$

651 In the Quadratic game, the SPPM update rule can be written as:

$$\begin{cases} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{A} \boldsymbol{\theta}_t - \mathbf{C} \boldsymbol{\phi}_{t+1} \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta \mathbf{B} \boldsymbol{\phi}_t + \mathbf{C}^T \boldsymbol{\theta}_{t+1} \end{cases} \quad (52)$$

652 Then we can rewrite the learning dynamics:

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \mathbf{A} \boldsymbol{\theta}_t - \eta \mathbf{C} \boldsymbol{\phi}_{t+1} \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \mathbf{A} \boldsymbol{\theta}_t - \eta \mathbf{C} (\boldsymbol{\phi}_t + \eta \mathbf{C}^T \boldsymbol{\theta}_{t+1} + \eta \mathbf{B} \boldsymbol{\phi}_t) \\ (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T) \boldsymbol{\theta}_{t+1} &= (\mathbf{I} - \eta \mathbf{A}) \boldsymbol{\theta}_t - \eta \mathbf{C} (\mathbf{I} + \eta \mathbf{B}) \boldsymbol{\phi}_t \\ \boldsymbol{\theta}_{t+1} &= (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T)^{-1} [(\mathbf{I} - \eta \mathbf{A}) \boldsymbol{\theta}_t - \eta \mathbf{C} (\mathbf{I} + \eta \mathbf{B}) \boldsymbol{\phi}_t] \end{aligned} \quad (53)$$

653 Similarly, for the other player we have

$$\begin{aligned} \boldsymbol{\phi}_{t+1} &= \boldsymbol{\phi}_t + \eta \mathbf{C}^T \boldsymbol{\theta}_{t+1} + \eta \mathbf{B} \boldsymbol{\phi}_t \\ \boldsymbol{\phi}_{t+1} &= \boldsymbol{\phi}_t + \eta \mathbf{C}^T (\boldsymbol{\theta}_t - \eta \mathbf{A} \boldsymbol{\theta}_t - \eta \mathbf{C} \boldsymbol{\phi}_{t+1}) + \eta \mathbf{B} \boldsymbol{\phi}_t \\ (\mathbf{I} + \eta^2 \mathbf{C}^T \mathbf{C}) \boldsymbol{\phi}_{t+1} &= \eta \mathbf{C}^T (\mathbf{I} - \eta \mathbf{A}) \boldsymbol{\theta}_t + (\mathbf{I} + \eta \mathbf{B}) \boldsymbol{\phi}_t \\ \boldsymbol{\phi}_{t+1} &= (\mathbf{I} + \eta^2 \mathbf{C}^T \mathbf{C})^{-1} [\eta \mathbf{C}^T (\mathbf{I} - \eta \mathbf{A}) \boldsymbol{\theta}_t + (\mathbf{I} + \eta \mathbf{B}) \boldsymbol{\phi}_t] \end{aligned} \quad (54)$$

654 Let us define the symmetric matrices  $\mathbf{Q}_{\boldsymbol{\theta}} = (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T)^{-1}$ ,  $\mathbf{Q}_{\boldsymbol{\phi}} = (\mathbf{I} + \eta^2 \mathbf{C}^T \mathbf{C})^{-1}$  and  
655  $\mathbf{P}_{\boldsymbol{\theta}} = (\mathbf{I} - \eta \mathbf{A})$ ,  $\mathbf{P}_{\boldsymbol{\phi}} = (\mathbf{I} + \eta \mathbf{B})$ . Further we define  $r_t = \|\boldsymbol{\theta}_{t+1}\|^2 + \|\boldsymbol{\phi}_{t+1}\|^2$ . Based on these  
656 definitions, and the expressions in (53) and (54) we have

$$\begin{aligned} \|\boldsymbol{\theta}_{t+1}\|^2 + \|\boldsymbol{\phi}_{t+1}\|^2 &= \|\mathbf{Q}_{\boldsymbol{\theta}} \mathbf{P}_{\boldsymbol{\theta}} \boldsymbol{\theta}_t\|^2 + \eta^2 \|\mathbf{Q}_{\boldsymbol{\theta}} \mathbf{C} \mathbf{P}_{\boldsymbol{\phi}} \boldsymbol{\phi}_t\|^2 + \|\mathbf{Q}_{\boldsymbol{\phi}} \mathbf{C}^T \mathbf{P}_{\boldsymbol{\theta}} \boldsymbol{\theta}_t\|^2 + \|\mathbf{Q}_{\boldsymbol{\phi}} \mathbf{P}_{\boldsymbol{\phi}} \boldsymbol{\phi}_t\|^2 \\ &\quad - 2\eta \boldsymbol{\theta}_t^T \mathbf{P}_{\boldsymbol{\theta}}^T \mathbf{Q}_{\boldsymbol{\theta}}^T \mathbf{Q}_{\boldsymbol{\theta}} \mathbf{C} \mathbf{P}_{\boldsymbol{\phi}} \boldsymbol{\phi}_t + 2\eta \boldsymbol{\phi}_t^T \mathbf{P}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}} \mathbf{C}^T \mathbf{P}_{\boldsymbol{\theta}} \boldsymbol{\theta}_t \end{aligned} \quad (55)$$

657 To simplify the expression in (55) we use Lemma A.1 to obtain the following equations:

$$\mathbf{Q}_{\boldsymbol{\theta}} \mathbf{C} = \mathbf{C} \mathbf{Q}_{\boldsymbol{\phi}} \quad (56)$$

$$\mathbf{Q}_{\boldsymbol{\phi}} \mathbf{C}^T = \mathbf{C}^T \mathbf{Q}_{\boldsymbol{\theta}} \quad (57)$$

658 Using this lemma, we can show that

$$\boldsymbol{\theta}_t^T \mathbf{P}_{\boldsymbol{\theta}}^T \mathbf{Q}_{\boldsymbol{\theta}}^T \mathbf{Q}_{\boldsymbol{\theta}} \mathbf{C} \mathbf{P}_{\boldsymbol{\phi}} \boldsymbol{\phi}_t = \boldsymbol{\theta}_t^T \mathbf{P}_{\boldsymbol{\theta}}^T \mathbf{Q}_{\boldsymbol{\theta}}^T \mathbf{C} \mathbf{Q}_{\boldsymbol{\phi}} \mathbf{P}_{\boldsymbol{\phi}} \boldsymbol{\phi}_t = \boldsymbol{\phi}_t^T \mathbf{P}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}} \mathbf{C}^T \mathbf{Q}_{\boldsymbol{\theta}} \mathbf{P}_{\boldsymbol{\theta}} \boldsymbol{\theta}_t = \boldsymbol{\phi}_t^T \mathbf{P}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}}^T \mathbf{Q}_{\boldsymbol{\phi}} \mathbf{C}^T \mathbf{P}_{\boldsymbol{\theta}} \boldsymbol{\theta}_t$$

where the intermediate equality holds as  $a^T C = C^T a$ . Hence, the expression in (55) can be simplified as

$$\|\theta_{t+1}\|^2 + \|\phi_{t+1}\|^2 = \|\mathbf{Q}_\theta \mathbf{P}_\theta \theta_t\|^2 + \eta^2 \|\mathbf{Q}_\theta \mathbf{C} \mathbf{P}_\phi \phi_t\|^2 + \|\mathbf{Q}_\phi \mathbf{C}^T \mathbf{P}_\theta \theta_t\|^2 + \|\mathbf{Q}_\phi \mathbf{P}_\phi \phi_t\|^2 \quad (58)$$

We simplify equation (58) as follows. Consider the term involving  $\theta_t$ . We have

$$\begin{aligned} \|\mathbf{Q}_\theta \mathbf{P}_\theta \theta_t\|^2 + \eta^2 \|\mathbf{Q}_\phi \mathbf{C}^T \mathbf{P}_\theta \theta_t\|^2 &= \theta_t^T \mathbf{P}_\theta^T \mathbf{Q}_\theta^2 \mathbf{P}_\theta \theta_t + \eta^2 \theta_t^T \mathbf{P}_\theta^T \mathbf{C} \mathbf{Q}_\phi^2 \mathbf{C}^T \mathbf{P}_\theta \theta_t \\ &= \theta_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 \mathbf{C} \mathbf{Q}_\phi^2 \mathbf{C}^T) \mathbf{P}_\theta \theta_t \\ &= \theta_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 \mathbf{C} \mathbf{Q}_\phi \mathbf{C}^T \mathbf{Q}_\theta) \mathbf{P}_\theta \theta_t \\ &= \theta_t^T \mathbf{P}_\theta^T (\mathbf{Q}_\theta^2 + \eta^2 \mathbf{C} \mathbf{C}^T \mathbf{Q}_\theta \mathbf{Q}_\theta) \mathbf{P}_\theta \theta_t \\ &= \theta_t^T \mathbf{P}_\theta^T (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T) \mathbf{Q}_\theta^2 \mathbf{P}_\theta \theta_t \\ &= \theta_t^T \mathbf{P}_\theta^T (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T)^{-1} \mathbf{P}_\theta \theta_t \end{aligned} \quad (59)$$

where the last equality follows by replacing  $\mathbf{Q}_\theta$  by its definition. The same procedure follows for the term involving  $\phi_t$  which leads to the expression

$$\|\mathbf{Q}_\phi \mathbf{P}_\phi \phi_t\|^2 + \eta^2 \|\mathbf{Q}_\theta \mathbf{C} \mathbf{P}_\phi \phi_t\|^2 = \phi_t^T \mathbf{P}_\phi^T (\mathbf{I} + \eta^2 \mathbf{C}^T \mathbf{C})^{-1} \mathbf{P}_\phi \phi_t. \quad (60)$$

Substitute  $\|\mathbf{Q}_\theta \mathbf{P}_\theta \theta_t\|^2 + \eta^2 \|\mathbf{Q}_\phi \mathbf{C}^T \mathbf{P}_\theta \theta_t\|^2$  and  $\|\mathbf{Q}_\phi \mathbf{P}_\phi \phi_t\|^2 + \eta^2 \|\mathbf{Q}_\theta \mathbf{C} \mathbf{P}_\phi \phi_t\|^2$  in (58) with the expressions in (59) and (60), respectively, to obtain

$$\|\theta_{t+1}\|^2 + \|\phi_{t+1}\|^2 = \theta_t^T \mathbf{P}_\theta^T (\mathbf{I} + \eta^2 \mathbf{C} \mathbf{C}^T)^{-1} \mathbf{P}_\theta \theta_t + \phi_t^T \mathbf{P}_\phi^T (\mathbf{I} + \eta^2 \mathbf{C}^T \mathbf{C})^{-1} \mathbf{P}_\phi \phi_t. \quad (61)$$

Now using the expression in (61) and the fact that  $\mathbf{P}_\theta = \mathbf{P}_\theta^T$ ,  $\mathbf{P}_\phi = \mathbf{P}_\phi^T$  and  $\lambda_{\min}(\mathbf{C}^T \mathbf{C}) = \lambda_{\min}(\mathbf{C} \mathbf{C}^T)$ , we can write

$$\|\theta_{t+1} - \theta^*\|^2 + \|\phi_{t+1} - \phi^*\|^2 \leq \frac{\rho^2 (\mathbf{I} - \eta \mathbf{A}) \|\theta_t - \theta^*\|^2 + \rho^2 (1 + \eta \mathbf{B}) \|\phi_t - \phi^*\|^2}{\mathbf{I} + \eta^2 \lambda_{\min}(\mathbf{C}^T \mathbf{C})}.$$

□

## A.7 Competitive Gradient Descent as an Approximation of SPPM

In this section, we justify our results in Section 4.1 that Competitive Gradient Descent is a first order Taylor approximation of SPPM. Firstly, we consider the standard definition of CGD:

$$\begin{cases} \theta_{t+1} = \theta_t - \eta (\mathbf{I} + \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t))^{-1} (\nabla_{\theta} f(\theta_t, \phi_t) + \eta \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t)) \\ \phi_{t+1} = \phi_t + \eta (\mathbf{I} + \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t))^{-1} (\nabla_{\phi} f(\theta_t, \phi_t) - \eta \nabla_{\phi\theta} f(\theta_t, \phi_t) \nabla_{\theta} f(\theta_t, \phi_t)) \end{cases}$$

Rewriting the update rules we can get:

$$\begin{aligned} (\mathbf{I} + \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t)) (\theta_{t+1} - \theta_t) &= -\eta (\nabla_{\theta} f(\theta_t, \phi_t) + \eta \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t)) \\ \theta_{t+1} &= \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t) - \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t) - \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t) (\theta_{t+1} - \theta_t) \end{aligned}$$

Similarly, we have:

$$\phi_{t+1} = \phi_t + \eta \nabla_{\phi} f(\theta_t, \phi_t) - \eta^2 \nabla_{\phi\theta} f(\theta_t, \phi_t) \nabla_{\theta} f(\theta_t, \phi_t) - \eta^2 \nabla_{\phi\theta} f(\theta_t, \phi_t) \nabla_{\theta\phi} f(\theta_t, \phi_t) (\phi_{t+1} - \phi_t)$$

Therefore, CGD is a first order approximation of SPPM. Then we prove that the standard definition of CGD is equivalent to the update rule in Table 1. Consider the update rule in Table 1 and its footnote, we have:

$$\begin{cases} \theta_{t+1} = \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t) - \eta \nabla_{\theta\phi} f(\theta_t, \phi_t) (\phi_{t+1} - \phi_t) \\ \phi_{t+1} = \phi_t + \eta \nabla_{\phi} f(\theta_t, \phi_t) + \eta \nabla_{\phi\theta} f(\theta_t, \phi_t) (\theta_{t+1} - \theta_t) \end{cases} \quad (62)$$

Substituting  $(\phi_{t+1} - \phi_t)$  into the first equation of (62) we get:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t) - \eta \nabla_{\theta\phi} f(\theta_t, \phi_t) (\eta \nabla_{\phi} f(\theta_t, \phi_t) + \eta \nabla_{\phi\theta} f(\theta_t, \phi_t) (\theta_{t+1} - \theta_t)) \\ \theta_{t+1} &= \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t) - \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t) - \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t) (\theta_{t+1} - \theta_t) \\ (\mathbf{I} + \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t)) (\theta_{t+1} - \theta_t) &= -\eta \nabla_{\theta} f(\theta_t, \phi_t) - \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t) \\ \theta_{t+1} &= \theta_t - \eta (\mathbf{I} + \eta^2 \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi\theta} f(\theta_t, \phi_t))^{-1} (\nabla_{\theta} f(\theta_t, \phi_t) + \eta \nabla_{\theta\phi} f(\theta_t, \phi_t) \nabla_{\phi} f(\theta_t, \phi_t)). \end{aligned}$$

Substituting  $(\theta_{t+1} - \theta_t)$  into the second equation of (62) and applying similar arguments we get:

$$\phi_{t+1} = \phi_t + \eta (\mathbf{I} + \eta^2 \nabla_{\phi\theta} f(\theta_t, \phi_t) \nabla_{\theta\phi} f(\theta_t, \phi_t))^{-1} (\nabla_{\phi} f(\theta_t, \phi_t) - \eta \nabla_{\phi\theta} f(\theta_t, \phi_t) \nabla_{\theta} f(\theta_t, \phi_t))$$

Thus the update rule in Table 1 is equivalent to the standard definition of CGD and it is equivalent to the first order Taylor approximation of SPPM.

## 681 A.8 Experiments on Bilinear and Quadratic Games

682 **Bilinear Game** Consider the following bilinear game:

$$\min_{\theta \in \mathbb{R}} \max_{\phi \in \mathbb{R}} a\theta\phi \quad (63)$$

683 the example presented in Figure 1 is an example of using different algorithms to solve the bilinear  
 684 game above with coefficient  $a = 10$ . For sake of completeness, we also provide a grid of experiment  
 685 results for different algorithms with different coefficients  $a$  and learning rates  $\eta$ , starting from the  
 same point  $(\theta_0, \phi_0) = (-12, 10)$ . The result is presented in Figure 6. The experiment demonstrates

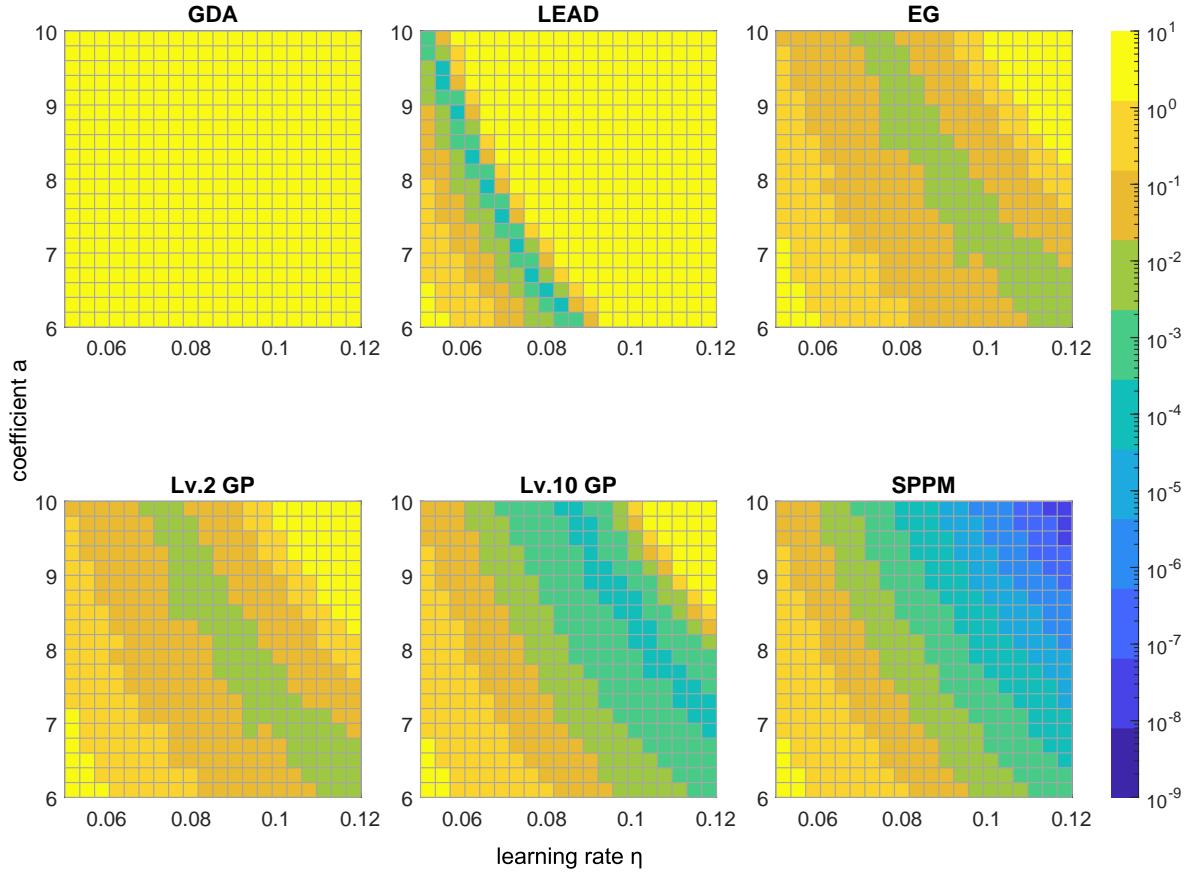


Figure 6: A grid of experiments on the bilinear game for different algorithms with different values of coefficient  $a$  and learning rates  $\eta$ . The color in each cell indicates the distance to the equilibrium after 50 iterations.

686 that, in a bilinear game, Lv.2 GP is equivalent to the extra-gradient method, and higher level Lv. $k$  GP  
 687 performs better with increased coefficient  $a$  and learning rate  $\eta$  as long as it remains a contraction  
 688 (i.e.,  $\eta < a^{-1}$ ).  
 689

690 **Quadratic Game** For the quadratic game presented in Figure 3, we randomly initialize the matrices  
691  $A$  and  $B$ :

$$A = \begin{bmatrix} 1.8398 & 0.5195 & 1.2537 & 1.7470 & 1.2769 \\ 0.5195 & 0.6586 & 0.4476 & 0.8898 & 1.1309 \\ 1.2537 & 0.4476 & 1.4440 & 1.3923 & 0.8877 \\ 1.7470 & 0.8898 & 1.3923 & 2.1249 & 1.7664 \\ 1.2769 & 1.1309 & 0.8877 & 1.7664 & 2.1553 \end{bmatrix} \quad B = - \begin{bmatrix} 1.0821 & 1.2427 & 1.0093 & 1.3335 & 0.6761 \\ 1.2427 & 2.2031 & 1.3236 & 1.8566 & 0.9394 \\ 1.0093 & 1.3236 & 1.2393 & 1.3675 & 0.9065 \\ 1.3335 & 1.8566 & 1.3675 & 1.9081 & 0.9693 \\ 0.6761 & 0.9394 & 0.9065 & 0.9693 & 0.7141 \end{bmatrix}$$

692 where  $A$  is symmetric and positive definite and  $B$  is symmetric and negative definite. The interaction  
693 matrix is defined as:

$$C = \begin{bmatrix} c & 0 & 0 & 0 & 0 \\ 0 & c & 0 & 0 & 0 \\ 0 & 0 & c & 0 & 0 \\ 0 & 0 & 0 & c & 0 \\ 0 & 0 & 0 & 0 & c \end{bmatrix}$$

694 where  $c$  represents the strength of the interaction between the two players. The starting point  $\theta_0$  and  
695  $\phi_0$  are  $[0.1270, 0.9667, 0.2605, 0.8972, 0.3767]^T$  and  $[0.3362, 0.4514, 0.8403, 0.1231, 0.5430]^T$  respectively.

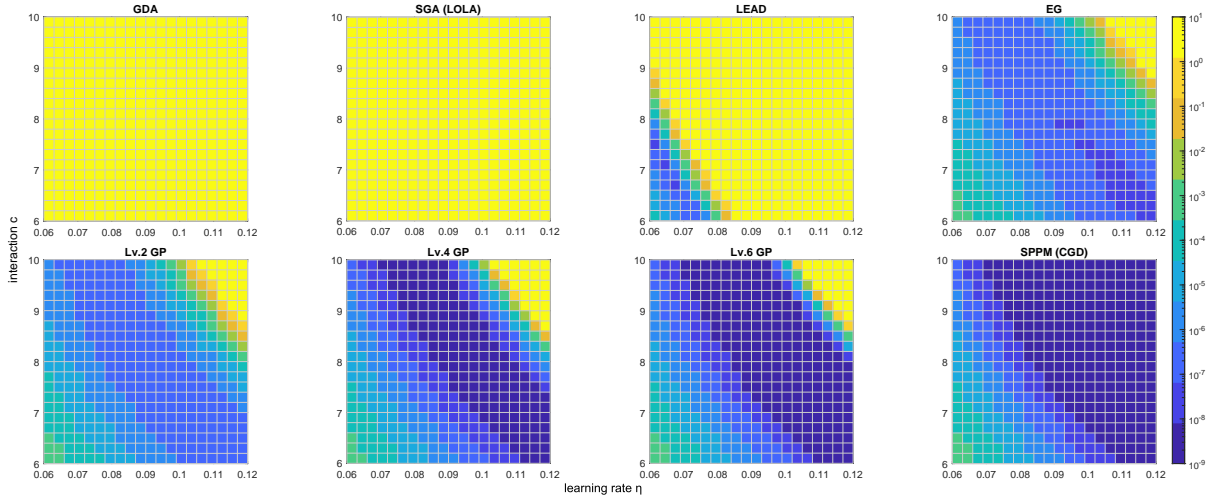


Figure 7: A grid of experiments on the quadratic game for different algorithms with different values of coefficient  $a$  and learning rates  $\eta$ . The color in each cell indicates the distance to the equilibrium after 50 iterations.

696

## 697 A.9 Experiments on 8-Gaussians

698 **Dataset** The target distribution is a mixture of 8-Gaussians with standard deviation equal to 0.05  
699 and modes uniformly distributed around a unit circle.

700 **Architecture** The architecture for the generator and the discriminator, each consists of four hidden  
701 layers followed by ReLU activation function. The weight initialization uses PyTorch’s default  
702 initialization scheme. The architecture is presented in Table 8:

703 **Experiment** For our experiments, we used the PyTorch framework. Furthermore, the batch size we  
704 used is 128.

Table 5: Architecture used for the Mixture of 8-Gaussians

Generator	Discriminator
<i>Input: <math>z \in \mathbb{R}^{64} \sim \mathcal{N}(0, I)</math></i>	<i>Input: <math>x \in \mathbb{R}^2</math></i>
Linear(64 $\rightarrow$ 2000)	Linear(2 $\rightarrow$ 2000)
ReLU	ReLU
Linear(2000 $\rightarrow$ 2000)	Linear(2000 $\rightarrow$ 2000)
ReLU	ReLU
Linear(2000 $\rightarrow$ 2000)	Linear(2000 $\rightarrow$ 2000)
ReLU	ReLU
Linear(2000 $\rightarrow$ 1)	Linear(2000 $\rightarrow$ 1)

## 705 A.10 Experiments on CIFAR-10 and STL-10

706 For our experiments, we used the PyTorch<sup>2</sup> framework. For experiments on CIFAR-10 and STL-  
 707 10, we compute the FID and IS metrics using the provided implementations in Tensorflow<sup>3</sup> for  
 708 consistency with related works.

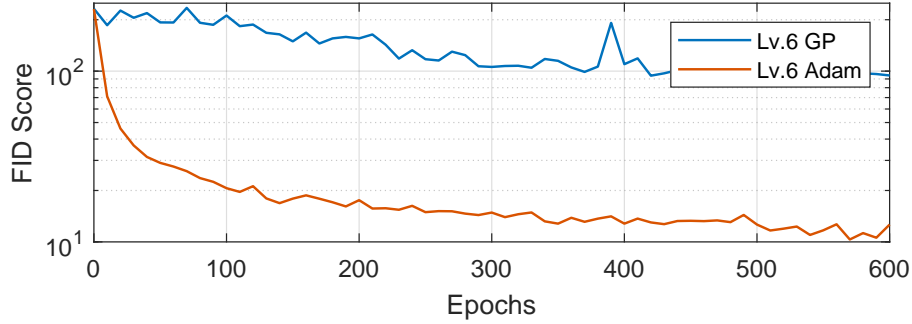


Figure 8: Comparison between  $Lv.k$  GP and  $Lv.k$  Adam on generating CIFAR-10 images. We can see significant improvements in FID when using the  $Lv.k$  Adam algorithm we proposed.

709  **$Lv.k$  GP vs  $Lv.k$  Adam** In experiments, we compare the performance of  $Lv.k$  GP and  $Lv.k$  Adam  
 710 on the task of CIFAR-10 image generation. The experiment results is presented in Figure 8. The  
 711 experiments on  $Lv.k$  GP and  $Lv.k$  Adam use the same initialization and hyperparameters. According  
 712 to our experiments,  $Lv.k$  Adam converges much faster than  $Lv.k$  GP for the same choice of  $k$  and  
 713 learning rates.

714 **Adam vs  $Lv.k$  Adam** We also present a comparison between the performance of Adam and  $Lv.k$   
 715 Adam optimizers on the task of STL-10 image generation. The experiment results is presented in  
 716 Figure 9. Under the same choice of hyperparameters and identical model parameter initialization,  
 717  $Lv.k$  Adam consistently outperforms the Adam optimizer in terms of FID score.

718 **Accelerated  $Lv.k$  Adam** In this section, we propose an accelerated version of  $Lv.k$  Adam. The  
 719 intuition is that we update the min player  $\theta$  and the max player  $\phi$  in an alternating order. The  
 720 corresponding  $Lv.k$  GP algorithm can be written as:

$$\text{Reasoning: } \begin{cases} \theta_t^{(k)} = \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t^{(k-1)}) \\ \phi_t^{(k)} = \phi_t - \eta \nabla_{\phi} g(\theta_t^{(k)}, \phi_t) \end{cases} \text{Update: } \begin{cases} \theta_{t+1} = \theta_t - \eta \nabla_{\theta} f(\theta_t, \phi_t^{(k)}) \\ \phi_{t+1} = \phi_t - \eta \nabla_{\phi} g(\theta_t^{(k)}, \phi_t) \end{cases} \text{ (Alt-} Lv.k \text{ GP)}$$

<sup>2</sup><https://pytorch.org/>

<sup>3</sup><https://tensorflow.org/>

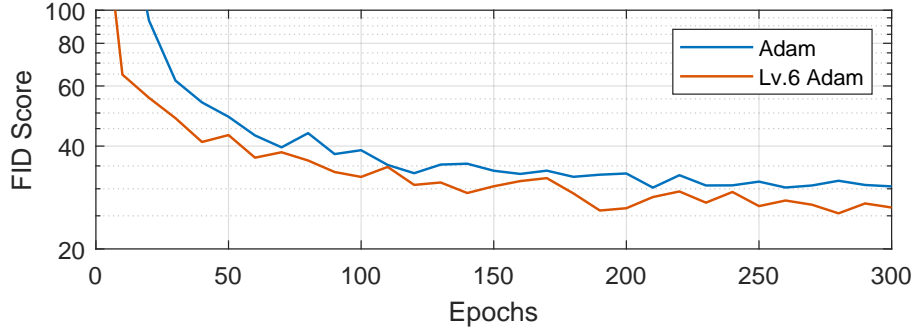


Figure 9: Comparison between Adam and Lv. $k$  Adam on generating STL-10 images. We can see that Lv. $k$  Adam consistently outperforms the Adam optimizer in terms of FID score.

721 Instead of responding to  $\theta_t^{(k-1)}$ , in Alt-Lv. $k$  GP, the max player  $\phi_t^{(k)}$  acts in response to the min  
 722 player's current action,  $\theta_t^{(k)}$ . A Lv. $k$  min player in Alt-Lv. $k$  GP is equivalent to a Lv. $2k - 1$  player in  
 723 the Lv. $k$  GP, and a Lv. $k$  max player in Alt-Lv. $k$  GP is equivalent to a Lv. $2k$  player in the Lv. $k$  GP,  
 724 respectively. Therefore, it is easy to verify that Alt-Lv. $k$  GP converges two times faster than Lv. $k$  GP  
 725 and the corresponding Alt-Lv. $k$  Adam algorithm is provided in Algorithm 2.

726 **Architecture** In this section, we describe the model we used to evaluate the performance of Lv. $k$   
 727 Adam for generating CIFAR-10<sup>4</sup> and STL-10 datasets. With 'conv' we denote a convolutional layer  
 728 and 'transposed conv' a transposed convolution layer. The models use Batch Normalization and  
 729 Spectral Normalization. The model's parameters are initialized with Xavier initialization.

730 **Images generated on CIFAR-10 and STL-10** In this section, we present sample images generated  
 731 by the best performing trained generators on CIFAR-10 and STL-10.

<sup>4</sup>CIFAR10 is released under the MIT license.

---

**Algorithm 2:** Accelerated Level  $k$  Adam: proposed Adam with recursive reasoning steps

---

**Input:** Stopping time  $T$ , reasoning steps  $k$ , learning rate  $\eta_\theta, \eta_\phi$ , decay rates for momentum estimates  $\beta_1, \beta_2$ , initial weight  $(\theta_0, \phi_0)$ ,  $P_x$  and  $P_z$  real and noise-data distributions, losses  $\mathcal{L}_G(\theta, \phi, x, z)$  and  $\mathcal{L}_D(\theta, \phi, x, z)$ ,  $\epsilon = 1e - 8$ .

**Parameters :** Initial parameters:  $\theta_0, \phi_0$

Initialize first moments:  $m_{\theta,0} \leftarrow 0, m_{\phi,0} \leftarrow 0$

Initialize second moments:  $v_{\theta,0} \leftarrow 0, v_{\phi,0} \leftarrow 0$

**for**  $t=0, \dots, T-1$  **do**

**Sample** new mini-batch:  $x, z \sim P_x, P_z$ ,

$\theta_t^{(0)} \leftarrow \theta_t, \phi_t^{(0)} \leftarrow \phi_t$ ,

**for**  $n=1, \dots, k$  **do**

        Compute stochastic gradient:  $g_{\theta,t}^{(n)} = \nabla_{\theta} \mathcal{L}_G(\theta_t, \phi_t^{(n-1)}, x, z)$ ;

        Update estimate of first moment:  $m_{\theta,t}^{(n)} = \beta_1 m_{\theta,t-1} + (1 - \beta_1) g_{\theta,t}^{(n)}$ ;

        Update estimate of second moment:  $v_{\theta,t}^{(n)} = \beta_2 v_{\theta,t-1} + (1 - \beta_2) (g_{\theta,t}^{(n)})^2$ ;

        Correct the bias for the moments:  $\hat{m}_{\theta,t}^{(n)} = \frac{m_{\theta,t}^{(n)}}{(1 - \beta_1^n)}, \hat{v}_{\theta,t}^{(n)} = \frac{v_{\theta,t}^{(n)}}{(1 - \beta_2^n)}$ ;

        Perform Adam update:  $\theta_t^{(n)} = \theta_t - \eta_\theta \frac{\hat{m}_{\theta,t}^{(n)}}{\sqrt{\hat{v}_{\theta,t}^{(n)} + \epsilon}}$ ;

        Compute stochastic gradient:  $g_{\phi,t}^{(n)} = \nabla_{\phi} \mathcal{L}_D(\theta_t^{(n)}, \phi_t, x, z)$ ;

        Update estimate of first moment:  $m_{\phi,t}^{(n)} = \beta_1 m_{\phi,t-1} + (1 - \beta_1) g_{\phi,t}^{(n)}$ ;

        Update estimate of second moment:  $v_{\phi,t}^{(n)} = \beta_2 v_{\phi,t-1} + (1 - \beta_2) (g_{\phi,t}^{(n)})^2$ ;

        Correct the bias for the moments:  $\hat{m}_{\phi,t}^{(n)} = \frac{m_{\phi,t}^{(n)}}{(1 - \beta_1^n)}, \hat{v}_{\phi,t}^{(n)} = \frac{v_{\phi,t}^{(n)}}{(1 - \beta_2^n)}$ ;

        Perform Adam update:  $\phi_t^{(n)} = \phi_t - \eta_\phi \frac{\hat{m}_{\phi,t}^{(n)}}{\sqrt{\hat{v}_{\phi,t}^{(n)} + \epsilon}}$ ;

$\theta_{t+1} \leftarrow \theta_t^{(k)}, \phi_{t+1} \leftarrow \phi_t^{(k)}$ ;

$m_{\theta,t} \leftarrow m_{\theta,t}^{(k)}, m_{\phi,t} \leftarrow m_{\phi,t}^{(k)}$ ,

$v_{\theta,t} \leftarrow v_{\theta,t}^{(k)}, v_{\phi,t} \leftarrow v_{\phi,t}^{(k)}$

---

Table 6: ResNet blocks used for the SN-GAN architectures on CIFAR-10 image generation, for the generator (left) and the discriminator (right).

G-ResBlock		D-ResBlock ( $l$ -th block)
<i>Shortcut:</i>		<i>Shortcut:</i>
Upsample( $\times 2$ )		[AvgPool (ker: $2 \times 2$ )], if $l = 1$
<i>Residual:</i>		conv (ker: $1 \times 1, 3_{l=1}/128_{l \neq 1} \rightarrow 128$ ; stride: 1)
Batch Normalization		Spectral Normalization
ReLU		[AvgPool (ker: $2 \times 2$ , stride: 2)], if $l = 1$
Upsample( $\times 2$ )		<i>Residual:</i>
conv (ker: $3 \times 3, 256 \rightarrow 256$ ; stride: 1; pad: 1)		[ReLU], if $l \neq 1$
Batch Normalization		conv (ker: $3 \times 3, 3_{l=1}/128_{l \neq 1} \rightarrow 128$ ; stride: 1; pad: 1)
ReLU		Spectral Normalization
conv (ker: $3 \times 3, 256 \rightarrow 256$ ; stride: 1; pad: 1)		ReLU
		conv (ker: $1 \times 1, 128 \rightarrow 128$ ; stride: 1)
		Spectral Normalization
		AvgPool (ker: $2 \times 2$ )

---



Table 7: SN-GAN architectures for experiments on CIFAR-10

Generator	Discriminator
<i>Input: <math>z \in \mathbb{R}^{128} \sim \mathcal{N}(0, \mathbf{I})</math></i>	<i>Input: <math>x \in \mathbb{R}^{3 \times 32 \times 32}</math></i>
Linear(128 $\rightarrow$ 4096)	D-ResBlock
G-ResBlock	D-ResBlock
G-ResBlock	D-ResBlock
G-ResBlock	D-ResBlock
Batch Normalization	ReLU
ReLU	AvgPool(ker: $8 \times 8$ )
conv (ker: $3 \times 3$ , 256 $\rightarrow$ 3; stride: 1; pad: 1)	Linear(128 $\rightarrow$ 1)
Tanh	Spectral Normalization

Table 8: SN-GAN architectures for experiments on STL-10

Generator	Discriminator
<i>Input: <math>z \in \mathbb{R}^{128} \sim \mathcal{N}(0, \mathbf{I})</math></i>	<i>Input: <math>x \in \mathbb{R}^{3 \times 48 \times 48}</math></i>
Linear(128 $\rightarrow$ $6 \times 6 \times 512$ )	D-ResBlock down 64 $\rightarrow$ 128
G-ResBlock up 512 $\rightarrow$ 256	D-ResBlock down 3 $\rightarrow$ 128
G-ResBlock up 256 $\rightarrow$ 128	D-ResBlock down 128 $\rightarrow$ 256
G-ResBlock up 128 $\rightarrow$ 64	D-ResBlock down 256 $\rightarrow$ 512
Batch Normalization	D-ResBlock 512 $\rightarrow$ 1024
ReLU	ReLU, AvgPool (ker: $8 \times 8$ )
conv (ker: $3 \times 3$ , 64 $\rightarrow$ 3; stride: 1; pad: 1)	Linear(128 $\rightarrow$ 1)
Tanh	Spectral Normalization



Figure 10: The presented samples are generated by the best performing trained generator on CIFAR-10, using Lv.6 Adam. This gives a FID score of 10.12.

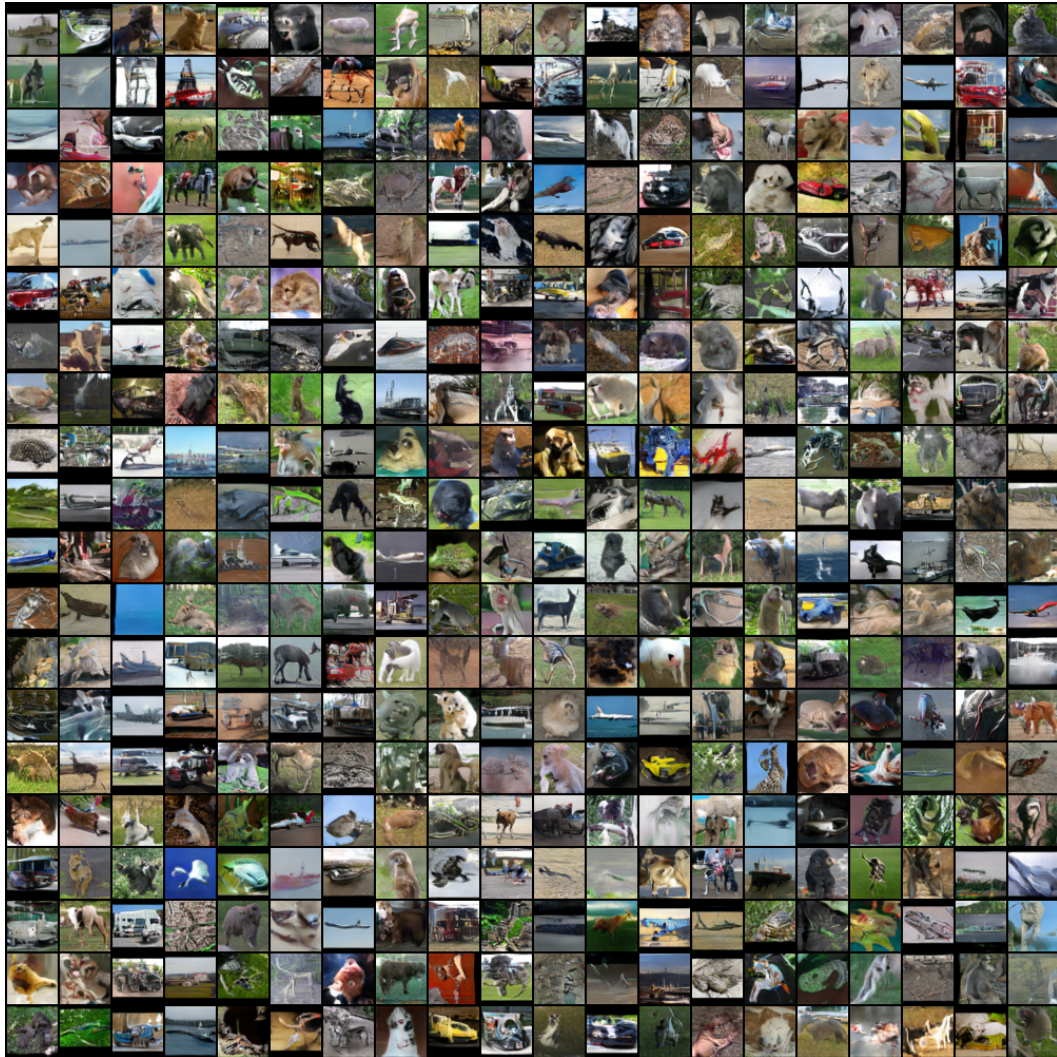


Figure 11: The presented samples are generated by the best performing trained generator on STL-10, using Lv.6 Adam. This gives a FID score of 25.43.