Supplementary Materials for "Fast and Scalable Algorithm for Detection of Structural Breaks in Big VAR Models"

Abolfazl Safikhani

Yue Bai

and

George Michailidis

Department of Statistics and Informatics Institute, University of Florida

June 26, 2021

In Appendix A, we collect technical lemmas needed to prove the main results, which are presented in Appendix B. Details of the algorithm for solving the optimization problem (4) are presented in Appendix C, while the selection of the tuning parameters is discussed in detail in Appendix D. Further, additional simulation results for the comparison with two competing methods are reported in Appendix E. Finally, note that all constants $c_1, c_2, \ldots, k_1, k_2, \ldots$ in the proofs of Lemmas and Theorems are generic constants, and their repetitions on different places in the proofs do not imply they are identical.

Appendix A: Technical Lemmas

Lemma A1. Under Assumption A1, there exist constants $c_i > 0$ such that with probability at least $1 - c_1 \exp(-c_2(\log(q) + 2\log(p)))$, for any sequence c_n ,

$$\sup_{1 \le j \le m_0 + 1, s \ge t_{j-1} + q, |t_j - s| > c_n} \left\| \left| (t_j - s)^{-1} \sum_{l=s}^{t_j - 1} Y_{l-1} Y_{l-1}' - \Gamma_j^q(0) \right| \right\|_{\infty} \le c_3 \sqrt{\frac{\log(q) + 2\log(p)}{c_n}},$$
(A.1)

where $\Gamma_j^q(0) = \mathbb{E}(Y_{t-1}Y'_{t-1})$ for $t_{j-1} + q \leq t < t_j$. Moreover,

$$\sup_{1 \le j \le m_0 + 1, s \ge t_{j-1} + q, |t_j - s| > c_n} \left\| \left| (t_j - s)^{-1} \sum_{l=s}^{t_j - 1} Y_{l-1} \varepsilon_l' \right| \right\|_{\infty} \le c_3 \sqrt{\frac{\log(q) + 2\log(p)}{c_n}}.$$
 (A.2)

Proof. The proof of this lemma is similar to that of Proposition 2.4 in Basu and Michailidis (2015). Next, we briefly outline the main steps of the proof, while omitting unnecessary

details. For (A.1), note that using an argument similar to Proposition 2.4(a) in Basu and Michailidis (2015), there exist $k_1, k_2 > 0$ such that for each fixed $k, l = 1, \dots, pq$,

$$\mathbb{P}\left(\left|e_{k}'\left(\frac{\sum_{l=s}^{t_{j}-1}Y_{l-1}Y_{l-1}'}{t_{j}-s}-\Gamma_{j}^{q}(0)\right)e_{l}\right| > k_{1}\eta\right) \le 6\exp(-k_{2}c_{n}\min(\eta,\eta^{2})).$$
(A.3)

Setting $\eta = k_3 \sqrt{\frac{\log(qp^2)}{c_n}}$, and taking union over all possible values of k, l, we obtain (A.1). Note that, there exist $k_1, k_2 > 0$ such that for each fixed k = 1, ..., pq, l = 1, ..., p,

$$\mathbb{P}\left(\left|e_k'\frac{\sum_{l=s}^{t_j-1}Y_{l-1}\varepsilon_l'}{t_j-s}e_l\right| > k_1\eta\right) \le 6\exp(-k_2c_n\min(\eta,\eta^2)).$$
(A.4)

Setting $\eta = k_3 \sqrt{\frac{\log(qp^2)}{c_n}}$, and taking union over all possible values of k, l, we get:

$$\left\| (t_j - s)^{-1} \sum_{l=s}^{t_j - 1} Y_{l-1} \varepsilon_l' \right\|_{\infty} \le c_3 \sqrt{\frac{\log(q) + 2\log(p)}{c_n}},\tag{A.5}$$

with high probability converging to 1 for any $j = 1, 2, \dots, m_0 + 1$, as long as $|t_j - s| > c_n$ and $s \ge t_{j-1} + q$. Note that the constants c_1, c_2 and c_3 can be chosen large enough such that the upper bounds above would be independent of the break point t_i . Therefore, we have the desired upper bounds verified with probability at least $1 - c_1 \exp(-c_2(\log(q) + 2\log(p)))$. \Box

Appendix B: Proof of Main Results

Proof of Theorem 1. The proof is similar to Theorem 2 in Safikhani and Shojaie (2020). For a matrix $A \in \mathbb{R}^{p \times pq}$, let $||A||_{1,\mathcal{I}_j} = \sum_{(k,h) \in \mathcal{I}_j} |a_{kh}|$. First, we focus on the second part. Suppose there exists a true break point t_{j_0} where

First, we focus on the second part. Suppose there exists a true break point t_{j_0} where $j_0 \in \{1, \dots, m_0\}$, which is isolated from all estimated points, i.e., $\min_{1 \le j \le \widehat{m}} |\widehat{t}_j - t_{j_0}| > n\gamma_n$. In other words, there exists an estimated break point \widehat{t}_j such that, $t_{j_0} - t_{j_0-1} \lor \widehat{t}_j \ge n\gamma_n$ and $t_{j_0+1} \land \widehat{t}_{j+1} - t_{j_0} \ge n\gamma_n$. The idea of the proof is to show the estimated AR parameter estimated in the interval $[t_{j_0-1} \lor \widehat{t}_j, t_{j_0+1} \land \widehat{t}_{j+1}]$ converges in L_2 to both $\Phi^{(.,j_0)}$ and $\Phi^{(.,j_0+1)}$ which contradicts Assumption A3. This is due to the fact that the length of the interval is large enough to verify restricted eigenvalue and deviation bound inequalities needed to show parameter estimation consistency.

Based on the definition of $\widehat{\Theta}$ in (4), the value of the function defined in (4) is minimized exactly at $\widehat{\Theta}$. This means that any other choice of parameters yields a higher value in (4). Denote the closest r_i to the right side of t_{j_0-1} plus q by s_{j_0-1} , i.e., $s_{j_0-1} = r_i + q$, and similarly, denote the closest r_i to the left side of t_{j_0} by s_{j_0} . First, we focus on the interval $[s_{j_0-1} \lor \hat{t}_j, s_{j_0}]$. Define a new parameter sequence ψ_k 's, $k = 1, ..., k_n$ with $\psi_k = \hat{\theta}_k$ except for two time points $k = \hat{i}_j \lor i_{j_0-1}$ and $k = i_{j_0}$. For these two points, if $\hat{t}_j > s_{j_0-1}$, set $\psi_{\hat{i}_j} = \Phi^{(..j_0)} - \hat{\Phi}_j$ and $\psi_{i_{j_0}} = \hat{\Phi}_{j+1} - \Phi^{(..j_0)}$; if $\hat{t}_j \le s_{j_0-1}$, set set $\psi_{i_{j_0-1}} = \Phi^{(..j_0)} - \hat{\Phi}_{j+1}$ and $\psi_{i_{j_0}} = \hat{\Phi}_{j+1} - \Phi^{(..j_0)}$, where $\hat{\Phi}_j = \sum_{k=1}^{\hat{i}_j-1} \hat{\theta}_k$ and $\hat{\Phi}_{j+1} = \sum_{k=1}^{\hat{i}_j} \hat{\theta}_k$, i.e. $\hat{\theta}_{\hat{i}_j} = \hat{\Phi}_{j+1} - \hat{\Phi}_j$, where i_{j_0-1} , i_{j_0} and \hat{i}_j are the corresponding indices of candidate points s_{j_0-1} , s_{j_0} and \hat{t}_j . Denoting $\Psi = \operatorname{vector}(\psi_1, ..., \psi_{k_n}) \in \mathbb{R}^{\pi_b \times 1}$, we have

$$\frac{1}{n} \|\mathbf{Y} - \mathbf{Z}\widehat{\mathbf{\Theta}}\|_{2}^{2} + \lambda_{1,n} \|\widehat{\mathbf{\Theta}}\|_{1} + \lambda_{2,n} \sum_{i=1}^{k_{n}} \left\| \sum_{j=1}^{i} \widehat{\theta}_{j} \right\|_{1} \leq \frac{1}{n} \|\mathbf{Y} - \mathbf{Z}\Psi\|_{2}^{2} + \lambda_{1,n} \|\Psi\|_{1} \qquad (B.1)$$

$$+ \lambda_{2,n} \sum_{i=1}^{k_{n}} \left\| \sum_{j=1}^{i} \psi_{j} \right\|_{1}.$$

When $\hat{t}_j > s_{j_0-1}$, some rearrangement of equation (B.1) leads to

$$\begin{array}{lll}
0 &\leq c \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{F}^{2} \\
&\leq \frac{1}{s_{j_{0}} - s_{j_{0}-1} \vee \widehat{t}_{j}} \sum_{l=s_{j_{0}-1} \vee \widehat{t}_{j}}^{s_{j_{0}} - 1} \left(\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\right)' Y_{l-1} Y_{l-1}' \left(\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\right) \\
&\leq \frac{2}{s_{j_{0}} - s_{j_{0}-1} \vee \widehat{t}_{j}} \sum_{l=s_{j_{0}-1} \vee \widehat{t}_{j}}^{s_{j_{0}} - 1} Y_{l-1}' \left(\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\right) \varepsilon_{l} \\
&+ \frac{n\lambda_{1,n}}{s_{j_{0}} - s_{j_{0}-1} \vee \widehat{t}_{j}} \left(\|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{1} + \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j}\|_{1} - \|\widehat{\Phi}_{j+1} - \widehat{\Phi}_{j}\|_{1}\right) \\
&+ \frac{n\lambda_{2,n}}{b_{n}} \left(\|\Phi^{(.,j_{0})}\|_{1} - \|\widehat{\Phi}_{j+1}\|_{1}\right) \\
&\leq \left(\frac{2n\lambda_{1,n}}{s_{j_{0}} - s_{j_{0}-1} \vee \widehat{t}_{j}} + C\sqrt{\frac{\log p}{n\gamma_{n}}}\right) \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{1} + \frac{n\lambda_{2,n}}{b_{n}} \left(\|\Phi^{(.,j_{0})}\|_{1} - \|\widehat{\Phi}_{j+1}\|_{1}\right) \\
&\leq \frac{1}{2} \frac{n\lambda_{2,n}}{b_{n}} \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{1} + \frac{n\lambda_{2,n}}{b_{n}} \left(\|\Phi^{(.,j_{0})}\|_{1} - \|\widehat{\Phi}_{j+1}\|_{1}\right) \\
&\leq \frac{3}{2} \frac{n\lambda_{2,n}}{b_{n}} \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{1,\mathcal{I}_{j_{0}}} - \frac{1}{2} \frac{n\lambda_{2,n}}{b_{n}} \|\Phi^{(.,j_{0})} - \widehat{\Phi}_{j+1}\|_{1,\mathcal{I}_{j_{0}}}. \tag{B.2}
\end{array}$$

When $\hat{t}_j \leq s_{j_0-1}$, a similar result can be obtained using an analogous rearrangement of equation (B.1). In equation (B.2), the second inequality holds with high probability converging to 1 due to first part of Lemma A1 and the fact that $s_{j_0} - s_{j_0-1} \vee \hat{t}_j \geq \frac{1}{2}n\gamma_n$ and $b_n \leq \frac{1}{4}n\gamma_n$ by Assumption A3. The fourth inequality holds with high probability converging to 1 due to second part of Lemma A1 and triangular inequality. The fifth inequality is based on Assumption A3 and the selection for $\lambda_{2,n}$ in the statement of the Theorem. The last inequality holds by the sparsity assumption. This implies that

$$\left\| \Phi^{(.,j_0)} - \widehat{\Phi}_{j+1} \right\|_F = O_p\left(\sqrt{\frac{d_n^* \log p}{n\gamma_n}}\right),\tag{B.3}$$

which means that $\|\Phi^{(.,j_0)} - \widehat{\Phi}_{j+1}\|_F$ converges to zero in probability based on Assumption A3. Similarly, the same procedure can be applied to the interval $[s_{j_0}, s_{j_0+1} \wedge \widehat{t}_{j+1}]$, where s_{j_0+1} is the closest r_i to the left side of t_{j_0+1} , which leads to $\|\Phi^{(.,j_0+1)} - \widehat{\Phi}_{j+1}\|_F$ converges to zero in probability as well. This yields a contradiction to Assumption A3, and therefore, the proof is complete.

The proof of the first part is similar to the second part. Hence, a brief sketch is provided. Assume $|\hat{\mathcal{A}}_n| < m_0$. This means there exist an isolated true break point, say t_{j_0} . More specifically, there exists an estimated break point \hat{t}_j such that, $t_{j_0} - t_{j_0-1} \vee \hat{t}_j \ge n\gamma_n/3$ and $t_{j_0+1} \wedge \hat{t}_{j+1} - t_{j_0} \ge n\gamma_n/3$. Now, similar arguments as explained in details in the second part can be applied to both intervals $[s_{j_0-1} \vee \hat{t}_j, s_{j_0}]$ and $[s_{j_0}, s_{j_0+1} \wedge \hat{t}_{j+1}]$ which leads to $\|\Phi^{(.,j_0+1)} - \Phi^{(.,j_0)}\|_F$ converges to zero and therefore contradicts Assumption A3. This completes the proof.

Proof of Theorem 2. Theorem 1 implies that for each t_i , $i = 1, \ldots, m_0$, there exist points $\hat{t}_j \in \hat{\mathcal{A}}_n$ such that $|\hat{t}_j - t_i| \leq n\gamma_n$. First, we show that

$$\left\{\widehat{t}_j \in \widehat{\mathcal{A}}_n \text{ s.t. } |\widehat{t}_j - t_i| \le n\gamma_n, \text{ for some } i = 1, \dots, m_0\right\} \subseteq \widetilde{\mathcal{A}}_n$$

with high probability converging to one as n tends to infinity. This shows that (a) $\mathbb{P}(\widetilde{m} < m_0) \to 0$; (b) $\mathbb{P}\left(d_H\left(\widetilde{\mathcal{A}}_n, \mathcal{A}_n\right) \leq n\gamma_n\right) \to 1$. We show this by contradiction. Suppose for some j, \hat{t}_j is not included in the optimal set $\widetilde{\mathcal{A}}_n$ which is a minimizer of (10). We show that adding \hat{t}_j costs less than removing it based on the LIC price.

Consider the case when \hat{t}_j is dropped in the screening step, while there exists a true change point t_i such that $|\hat{t}_j - t_i| \leq n\gamma_n$. Similar arguments as in Proposition 4.1 of Basu and Michailidis (2015) show that the tuning parameter selected based on Assumption A5 (i.e. $\eta_{\hat{t}_j} = c_1 \sqrt{\frac{d_n^* n\gamma_n}{a_n}} + c_2 d_n^*$) yield: $\left\| \Phi^{(.,i)} - \hat{\psi}_{\hat{t}_j} \right\|_1 \leq 4\sqrt{d_n^*} \left\| \Phi^{(.,i)} - \hat{\psi}_{\hat{t}_j} \right\|_F$ and $\left\| \Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_j} \right\|_1 \leq 4\sqrt{d_n^*} \left\| \Phi^{(.,i)} - \hat{\psi}_{\hat{t}_j} \right\|_F$. Note that in this case, the restricted eigenvalue condition does not hold. Therefore, the convergence of the $\hat{\psi}_{\hat{t}_j}$ cannot be verified. For this scenario, based on similar calculations as in case (c) of proof of Lemma 4 in Safikhani and Shojaie (2020) (see equations (22)-(26) in Safikhani and Shojaie (2020)'s supplement for details), after replacing s_{i-1} and s_i with $\hat{t}_j - a_n$ and $\hat{t}_j + a_n - 1$, we have the following lower bound:

$$\sum_{t=\hat{t}_j-a_n}^{\hat{t}_j+a_n-1} \|y_t - \hat{\psi}_{\hat{t}_j} Y_{t-1}\|_2^2 \ge \sum_{t=\hat{t}_j-a_n}^{\hat{t}_j+a_n-1} \|\varepsilon_t\|_2^2 + c_1 a_n - c_2 d_n^\star \log p, \tag{B.4}$$

for some positive c_1 and c_2 with high probability converging to one.

Denote by $B = \widetilde{\mathcal{A}}_n \cup \{\widehat{t}_j\}$. We only consider one segment $[\widehat{t}_j - a_n, \widehat{t}_j + a_n)$ for the case when \widehat{t}_j is kept as a selected break point. Similar arguments as in Lemma 4 of Safikhani and Shojaie (2020) yield:

$$\left\|\widehat{\psi}_{\widehat{t}_{j,1}} - \Phi^{(.,i)}\right\|_{1} \le 4\sqrt{d_{n}^{\star}} \left\|\widehat{\psi}_{\widehat{t}_{j,1}} - \Phi^{(.,i)}\right\|_{F}, \quad \left\|\widehat{\psi}_{\widehat{t}_{j,1}} - \Phi^{(.,i)}\right\|_{F} = O_{p}\left(d_{n}^{\star}\sqrt{\frac{n\gamma_{n}}{a_{n}}}\right), \quad (B.5)$$

and

$$\left\|\widehat{\psi}_{\widehat{t}_{j,2}} - \Phi^{(.,i+1)}\right\|_{1} \leq 4\sqrt{d_{n}^{\star}} \left\|\widehat{\psi}_{\widehat{t}_{j,2}} - \Phi^{(.,i+1)}\right\|_{F}, \quad \left\|\widehat{\psi}_{\widehat{t}_{j,2}} - \Phi^{(.,i+1)}\right\|_{F} = O_{p}\left(d_{n}^{\star}\sqrt{\frac{n\gamma_{n}}{a_{n}}}\right). \tag{B.6}$$

To see this, suppose $t_i + q < \hat{t}_j$ (the other case is similar). Observe that $\hat{\psi}_{\hat{t}_j,1}$ in (6) minimizes the least squares plus the ℓ_1 norm loss function. Therefore, the value of this objective function for $\hat{\psi}_{\hat{t}_j,1}$ will be smaller than any other choice of parameters, including $\Phi^{(.,i)}$. Hence,

$$\frac{1}{a_n} \sum_{t=\hat{t}_j-a_n}^{\hat{t}_j-1} \left\| y_t - \hat{\psi}_{\hat{t}_j,1} Y_{t-1} \right\|_2^2 + \eta_{\hat{t}_j,1} \left\| \hat{\psi}_{\hat{t}_j,1} \right\|_1 \leq \frac{1}{a_n} \sum_{t=\hat{t}_j-a_n}^{\hat{t}_j-1} \left\| y_t - \Phi^{(.,i)} Y_{t-1} \right\|_2^2 + \eta_{\hat{t}_j,1} \left\| \Phi^{(.,i)} \right\|_1.$$
(B.7)

Some rearrangements together with the use of proposition 4.2 of Basu and Michailidis (2015) lead to:

$$\begin{split} & c_{1} \left\| \Phi^{(.,i+1)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{F}^{2} - c_{2} \frac{\log p}{a_{n}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} &\leq \max \left(0, c_{1} \left\| \Phi^{(.,i+1)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{F}^{2} - c_{2} \frac{\log p}{a_{n}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} \right) \\ &\leq \frac{1}{a_{n}} \sum_{t=\widehat{t}_{j}-a_{n}}^{\widehat{t}_{j}-1} Y_{t-1}' \left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right)' \left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right) Y_{t-1} \\ &\leq \frac{2}{a_{n}} \sum_{t=\widehat{t}_{j}-a_{n}}^{\widehat{t}_{j}-1} Y_{t-1}' \left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right)' \left(y_{t} - \Phi^{(.,i)} Y_{t-1} \right) + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_{1}^{2} - \left\| \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} \right) \\ &= \frac{2}{a_{n}} \sum_{t=\widehat{t}_{j}-a_{n}}^{\widehat{t}_{j}-1} Y_{t-1}' \left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right)' \varepsilon_{t} + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_{1}^{2} - \left\| \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} \right) \\ &+ \frac{2}{a_{n}} \sum_{t=t_{i}}^{\widehat{t}_{j}-1} Y_{t-1}' \left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right)' \left(\Phi^{(.,i+1)} - \Phi^{(.,i)} \right) Y_{t-1} \\ &\leq c \sqrt{\frac{\log p}{a_{n}}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_{1}^{2} - \left\| \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1}^{2} \right) \\ &+ 2 \frac{\widehat{t}_{j} - t_{i}}{a_{n}} \operatorname{tr} \left(\left(\Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right)' \left(\frac{1}{\widehat{t}_{j} - t_{i}} \sum_{t=t_{i}}^{\widehat{t}_{j}-1} Y_{t-1} - \Gamma_{i+1}'(0) + \Gamma_{i+1}'(0) \right) \left(\Phi^{(.,i+1)} - \Phi^{(.,i)} \right) \right) \\ &= \mathrm{I}, \end{split}$$

where the second inequality is due to the restricted eigenvalue property (proposition 4.2 of Basu and Michailidis (2015)), the fourth inequality uses the circular invariance property of the trace function. Now, observe that

$$\begin{split} \mathbf{I} &\leq c\sqrt{\frac{\log p}{a_n}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_1 - \left\| \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 \right) \\ &+ 2c' \frac{\widehat{t}_j - t_i}{a_n} \left(\max\left(\sqrt{\frac{\log p}{\widehat{t}_j - t_i}}, \frac{\log p}{\widehat{t}_j - t_i} \right) + \lambda_{\max} \left(\Gamma_{i+1}^q(0) \right) \right) \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 \left\| \Phi^{(.,i+1)} - \Phi^{(.,i)} \right\|_1 \\ &\leq \left(c\sqrt{\frac{\log p}{a_n}} + 4c' M_{\Phi} d_n^* \frac{n\gamma_n}{a_n} \right) \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_1 - \left\| \widehat{\theta} \right\|_1 \right) \\ &\leq \max\left(c, 4c' M_{\Phi} \right) \sqrt{\frac{d_n^* n\gamma_n}{a_n}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_1 - \left\| \widehat{\theta} \right\|_1 \right) \\ &\leq \frac{\eta_{\widehat{t}_{j,1}}}{2} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 + \eta_{\widehat{t}_{j,1}} \left(\left\| \Phi^{(.,i)} \right\|_1 - \left\| \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1 \right) \\ &\leq \frac{3\eta_{\widehat{t}_{j,1}}}{2} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1,\mathcal{I}_i} - \frac{\eta_{\widehat{t}_{j,1}}}{2} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1,\mathcal{I}_i^c} \\ &\leq 2\eta_{\widehat{t}_{j,1}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_1, \end{split}$$
(B.8)

where the first inequality is based on the first part of Lemma A1 with $c_n = \hat{t}_j - t_i$, the third inequality is due to Assumptions A3 and A4 under which $\log p \leq n\gamma_n \leq a_n$ and $a_n > cd_n^{\star 3}n\gamma_n$, which lead to $\sqrt{\frac{\log p}{a_n}} \leq \sqrt{\frac{n\gamma_n}{a_n}}$ and $d_n^{\star}\frac{n\gamma_n}{a_n} \leq \sqrt{\frac{d_n^{\star}n\gamma_n}{a_n}}$; the fourth inequality is due the selection of the tuning parameter $\eta_{\hat{t}_j,1} = c\sqrt{\frac{d_n^{\star}n\gamma_n}{a_n}}$; the fifth inequality is due the triangular inequality.

This ensures that $\left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1,\mathcal{I}_{i}^{c}} \leq 3 \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1,\mathcal{I}_{i}}$, and hence $\left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1} \leq 4 \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{1,\mathcal{I}_{i}} \leq 4 \sqrt{d_{n}^{\star}} \left\| \Phi^{(.,i)} - \widehat{\psi}_{\widehat{t}_{j,1}} \right\|_{F}$. This comparison between L_{1} and L_{2} norms of the error term together with the bound in Equation (B.8) will get the desired consistency rates in (B.5). Verifying (B.6) follows exactly from proposition 4.1 of Basu and Michailidis (2015) since there are no break points in the interval $[\widehat{t}_{j}, \widehat{t}_{j} + a_{n}]$. Note that the convergence rate here is of order $\sqrt{\frac{d_{n}^{\star}n\gamma_{n}}{a_{n}}}$ instead of $\sqrt{d_{n}^{\star}}\sqrt{\frac{\log p}{a_{n}}}$ due to the different selection of tuning parameter in our setup compared to Basu and Michailidis (2015). In fact, one could select $\eta_{\widehat{t}_{j,2}}$ to be of order $\sqrt{\frac{\log p}{a_{n}}}$ in order to match the rate here with Basu and Michailidis (2015). However, since we don't know whether the true break point t_{i} will be on the left hand side or the right hand side of \widehat{t}_{j} , it's better for the tuning parameters $\eta_{\widehat{t}_{j,1}}$ and $\eta_{\widehat{t}_{j,2}}$ to have the same order from a practical viewpoint. This justifies setting them to be of the same order here.

Using the results of (B.5) and (B.6), one can find upper bounds for the sum of squared errors in the segment $(\hat{t}_j - a_n, \hat{t}_j + a_n)$ for the case when \hat{t}_j is kept as a selected break point. Specifically, we have:

$$\sum_{t=\hat{t}_j-a_n}^{\hat{t}_j-1} \|y_t - \widehat{\psi}_{\hat{t}_j,1} Y_{t-1}\|_2^2 \le \sum_{t=\hat{t}_j-a_n}^{\hat{t}_j-1} \|\varepsilon_t\|_2^2 + O_p\left(d_n^{\star 3} n \gamma_n\right), \tag{B.9}$$

and

$$\sum_{t=\hat{t}_{j}}^{\hat{t}_{j}+a_{n}-1} \|y_{t}-\hat{\psi}_{\hat{t}_{j},2}Y_{t-1}\|_{2}^{2} \leq \sum_{t=\hat{t}_{j}}^{\hat{t}_{j}+a_{n}-1} \|\varepsilon_{t}\|_{2}^{2} + O_{p}\left(d_{n}^{\star\,3}n\gamma_{n}\right).$$
(B.10)

To see (B.9), we follow the steps of the proof of Theorem 3 in Safikhani and Shojaie (2020). Observe that:

$$\sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1} \|y_{t} - \widehat{\psi}_{\hat{t}_{j},1}Y_{t-1}\|_{2}^{2} \leq \sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1} \|\varepsilon_{t}\|_{2}^{2} + c\sqrt{(t_{i} - \hat{t}_{j} + a_{n})\log p} \left\|\Phi^{(.,i)} - \widehat{\psi}_{\hat{t}_{j},1}\right\|_{1}^{1} \\
+ (t_{i} - \hat{t}_{j} + a_{n})\operatorname{tr}\left(\left(\Phi^{(.,i)} - \widehat{\psi}_{\hat{t}_{j},1}\right)'\left(\frac{1}{t_{i} - \hat{t}_{j} + a_{n}}\sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1}Y_{t-1} - \Gamma_{i}^{q}(0) + \Gamma_{i}^{q}(0)\right)\left(\Phi^{(.,i)} - \widehat{\psi}_{\hat{t}_{j},1}\right)\right) \\
\leq \sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1} \|\varepsilon_{t}\|_{2}^{2} + cd_{n}^{\star 2}\sqrt{n\gamma_{n}\log p} + (t_{i} - \hat{t}_{j} + a_{n})\left(\sqrt{\frac{\log p}{t_{i} - \hat{t}_{j} + a_{n}} + \lambda_{\max}\left(\Gamma_{i}^{q}(0)\right)\right)\left\|\Phi^{(.,i)} - \widehat{\psi}_{\hat{t}_{j},1}\right\|_{1}^{2} \\
\leq \sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1} \|\varepsilon_{t}\|_{2}^{2} + cd_{n}^{\star 2}\sqrt{n\gamma_{n}\log p} + c'd_{n}^{\star 3}n\gamma_{n} \\
= \sum_{t=\hat{t}_{j}-a_{n}}^{t_{i}-1} \|\varepsilon_{t}\|_{2}^{2} + O_{p}\left(d_{n}^{\star 3}n\gamma_{n}\right), \tag{B.11}$$

where the second inequality is based on the fact that $\left\|\widehat{\psi}_{\widehat{t}_{j},1} - \Phi^{(.,i)}\right\|_{1} \leq c d_{n}^{\star 3/2} \sqrt{\frac{n\gamma_{n}}{a_{n}}}$ and $t_{i} - \widehat{t}_{j} + a_{n} \leq a_{n}$. Also, for the segment $[t_{i}, \widehat{t}_{j}]$, note that:

$$\begin{split} \sum_{t=t_{i}}^{\hat{t}_{j}-1} \|y_{t} - \hat{\psi}_{\hat{t}_{j},1}Y_{t-1}\|_{2}^{2} &\leq \sum_{t=t_{i}}^{\hat{t}_{j}-1} \|\varepsilon_{t}\|_{2}^{2} + c(\hat{t}_{j} - t_{i}) \max\left(\frac{\log p}{\hat{t}_{j} - t_{i}}, \sqrt{\frac{\log p}{\hat{t}_{j} - t_{i}}}\right) \left\|\Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_{j},1}\right\|_{1} \\ &+ (\hat{t}_{j} - t_{i}) \operatorname{tr}\left(\left(\Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_{j},1}\right)'\left(\frac{1}{\hat{t}_{j} - t_{i}}\sum_{t=t_{i}}^{\hat{t}_{j}-1}Y_{t-1}Y_{t-1}' - \Gamma_{i+1}^{q}(0) + \Gamma_{i+1}^{q}(0)\right)\left(\Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_{j},1}\right)\right) \\ &\leq \sum_{t=t_{i}}^{\hat{t}_{j}-1} \|\varepsilon_{t}\|_{2}^{2} + cn\gamma_{n}\left(2M_{\Phi}d_{n}^{\star} + \left\|\Phi^{(.,i)} - \hat{\psi}_{\hat{t}_{j},1}\right\|_{1}\right) \\ &+ (\hat{t}_{j} - t_{i})\left(\max\left(\frac{\log p}{\hat{t}_{j} - t_{i}}, \sqrt{\frac{\log p}{\hat{t}_{j} - t_{i}}}\right) + \lambda_{\max}\left(\Gamma_{i+1}^{q}(0)\right)\right)\left\|\Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_{j},1}\right\|_{1}^{2} \\ &\leq \sum_{t=t_{i}}^{\hat{t}_{j}-1} \|\varepsilon_{t}\|_{2}^{2} + cd_{n}^{\star}n\gamma_{n} + c'n\gamma_{n}\left(2M_{\Phi}d_{n}^{\star} + \left\|\Phi^{(.,i)} - \hat{\psi}_{\hat{t}_{j},1}\right\|_{1}\right)^{2} \\ &= \sum_{t=t_{i}}^{\hat{t}_{j}-1} \|\varepsilon_{t}\|_{2}^{2} + O_{p}\left(d_{n}^{\star^{2}}n\gamma_{n}\right). \end{split}$$
(B.12)

where the second inequality is based on $\hat{t}_j - t_i \leq n\gamma_n$, $\log p \leq n\gamma_n$ and triangular inequality $\left\| \Phi^{(.,i+1)} - \hat{\psi}_{\hat{t}_j,1} \right\|_1 \leq \left\| \Phi^{(.,i+1)} - \Phi^{(.,i)} \right\|_1 + \left\| \Phi^{(.,i)} - \hat{\psi}_{\hat{t}_j,1} \right\|_1 \leq 2M_{\Phi}d_n^{\star} + \left\| \Phi^{(.,i)} - \hat{\psi}_{\hat{t}_j,1} \right\|_1$; the second term in the third inequality is based on Assumption A4 of $a_n > cd_n^{\star 3}n\gamma_n$ so that $d_n^{\star} \geq d_n^{\star 3/2} \sqrt{\frac{n\gamma_n}{a_n}}$.

Property (B.9) is now verified by combining (B.11) and (B.12). Property (B.10) is similar and therefore its proof is omitted.

Recall that $B = \{\hat{t}_1, \dots, \hat{t}_{m_0}\} = \widetilde{\mathcal{A}}_n \cup \{\hat{t}_j\}$. Applying the results (B.4), (B.9) and (B.10), we have

$$\begin{aligned} \operatorname{LIC}(\hat{t}_{1},...,\hat{t}_{\widetilde{m}}) - \operatorname{LIC}(B) \\ &= L_{n}(\tilde{t}_{1},...,\tilde{t}_{\widetilde{m}};\eta_{n}) + \widetilde{m}\omega_{n} - \left(L_{n}(\hat{t}_{1},...,\hat{t}_{m_{0}};\eta_{n}) + (\widetilde{m}+1)\omega_{n}\right) \\ &= \sum_{t=\hat{t}_{j}-a_{n}}^{\hat{t}_{j}+a_{n}-1} \|y_{t} - \widehat{\psi}_{j}Y_{t-1}\|_{2}^{2} - \sum_{t=\hat{t}_{j}-a_{n}}^{\hat{t}_{j}-1} \|y_{t} - \widehat{\psi}_{\hat{t}_{j},1}Y_{t-1}\|_{2}^{2} - \sum_{t=\hat{t}_{j}}^{\hat{t}_{j}+a_{n}-1} \|y_{t} - \widehat{\psi}_{\hat{t}_{j},2}Y_{t-1}\|_{2}^{2} - \omega_{n} \\ &\geq c_{1}a_{n} - c_{2}d_{n}^{\star}\log p - c_{3}n\gamma_{n}d_{n}^{\star^{3}} - \omega_{n} \\ &\geq c_{1}a_{n} - 2c_{2}n\gamma_{n}d_{n}^{\star^{3}} - \omega_{n} \\ &> 0. \end{aligned}$$

The second inequality hold due to condition $\log p \leq n\gamma_n$ in Assumption A3. The last inequality holds due to conditions $\omega_n = n\gamma_n d_n^{\star^3}$ and $\lim_{n\to\infty} \omega_n/a_n = 0$ in Assumption A4. This proves that $\mathbb{P}\left(\left\{\widehat{t}_1, \ldots, \widehat{t}_{m_0}\right\} \subseteq \widetilde{\mathcal{A}}_n\right) \to 1$.

It remains to establish that $\mathbb{P}\left(d_H\left(\mathcal{A}_n, \mathcal{\widetilde{A}}_n\right) \leq a_n\right) \to 1$. It suffices to show that there exists a true break point in the interval $(\tilde{t}_j - a_n, \tilde{t}_j + a_n)$ for any $\tilde{t}_j \in \mathcal{\widetilde{A}}_n$. We prove this statement by contradiction. Suppose that there is a time point $s \in \mathcal{\widetilde{A}}_n$, such that there exists no true break point in the interval $(s - a_n, s + a_n)$. Next, define a new set $C = \mathcal{\widetilde{A}}_n \setminus \{s\}$.

Since there is no true break point in the segment, $(s - a_n, s + a_n)$, it belongs to some stationary segment, say the *i*-th one. Select tuning parameter $\eta_s = \eta_{s,1} = \eta_{s,2} = c \sqrt{\frac{\log p}{a_n}}$. Denote the estimated AR parameters in $[s - a_n, s)$ by $\hat{\psi}_{s,1}$ and the estimated AR parameters in $[s, s + a_n)$ by $\hat{\psi}_{s,2}$. Keeping the point *s* yields to :

$$\left\|\widehat{\psi}_{s,1} - \Phi^{(.,i)}\right\|_{1} \le 4\sqrt{d_{n}^{\star}} \left\|\widehat{\psi}_{s,1} - \Phi^{(.,i)}\right\|_{F}, \ \left\|\widehat{\psi}_{s,1} - \Phi^{(.,i)}\right\|_{F} = O_{p}\left(\sqrt{\frac{d_{n}^{\star}\log p}{a_{n}}}\right), \quad (B.13)$$

and

$$\left\|\widehat{\psi}_{s,2} - \Phi^{(.,i)}\right\|_{1} \le 4\sqrt{d_{n}^{\star}} \left\|\widehat{\psi}_{s,2} - \Phi^{(.,i)}\right\|_{F}, \quad \left\|\widehat{\psi}_{s,2} - \Phi^{(.,i)}\right\|_{F} = O_{p}\left(\sqrt{\frac{d_{n}^{\star}\log p}{a_{n}}}\right). \quad (B.14)$$

Therefore, we have the following lower bound:

$$\sum_{t=s-a_n}^{s-1} \|y_t - \widehat{\psi}_{s,1} Y_{t-1}\|_2^2 + \sum_{t=s}^{s+a_n-1} \|y_t - \widehat{\psi}_{s,2} Y_{t-1}\|_2^2 \ge \sum_{t=s-a_n}^{s+a_n-1} \|\varepsilon_t\|_2^2 - c_1 d_n^* n \gamma_n, \qquad (B.15)$$

for some positive c_1 with high probability converging to one.

To see (B.15), we follow the case (b) in the proof of Lemma 4 in Safikhani and Shojaie (2020). For interval $[s - a_n, s)$,

$$\begin{split} \sum_{t=s-a_n}^{s-1} \|y_t - \widehat{\psi}_{s,1} Y_{t-1}\|_2^2 &\geq \sum_{t=s-a_n}^{s-1} \|\varepsilon_t\|_2^2 + ca_n \left\|\Phi^{(.,i)} - \widehat{\psi}_{s,1}\right\|_2^2 - c'\sqrt{a_n \log p} \left\|\Phi^{(.,i)} - \widehat{\psi}_{s,1}\right\|_1 \\ &\geq \sum_{t=s-a_n}^{s-1} \|\varepsilon_t\|_2^2 + ca_n \left\|\Phi^{(.,i)} - \widehat{\psi}_{s,1}\right\|_2 \left(\left\|\Phi^{(.,i)} - \widehat{\psi}_{s,1}\right\|_2 - \frac{4c'}{c}\sqrt{\frac{d_n^* \log p}{a_n}}\right) \\ &\geq \sum_{t=s-a_n}^{s-1} \|\varepsilon_t\|_2^2 - c'' d_n^* \log p. \end{split}$$

Similarly, for interval $[s, s + a_n)$, we have $\sum_{t=s}^{s+a_n-1} \|y_t - \widehat{\psi}_2 Y_{t-1}\|_2^2 \ge \sum_{t=s}^{s+a_n-1} \|\varepsilon_t\|_2^2 - c'' d_n^{\star} \log p.$

Next, we remove the break point s and denote the estimated parameters in $[s - a_n, s + a_n)$ by $\hat{\psi}_s$. We get

$$\sum_{t=s-a_n}^{s+a_n-1} \|y_t - \widehat{\psi}_s Y_{t-1}\|_2^2 \le \sum_{t=s-a_n}^{s+a_n-1} \|\varepsilon_t\|_2^2 + c_2 d_n^{\star 2} \log p \tag{B.16}$$

To see (B.16), we again follow the steps of the proof of Theorem 3 in Safikhani and Shojaie (2020):

$$\begin{split} \sum_{t=s-a_{n}}^{s+a_{n}-1} \|y_{t} - \widehat{\psi}_{s}Y_{t-1}\|_{2}^{2} &\leq \sum_{t=s-a_{n}}^{s+a_{n}-1} \|\varepsilon_{t}\|_{2}^{2} + c\sqrt{(2a_{n})\log p} \left\|\Phi^{(.,i)} - \widehat{\psi}_{s}\right\|_{1} \\ &+ (2a_{n})\mathrm{tr} \left(\left(\Phi^{(.,i)} - \widehat{\psi}_{s}\right)' \left(\frac{1}{2a_{n}} \sum_{t=s-a_{n}}^{s+a_{n}-1} Y_{t-1}Y_{t-1}' - \Gamma_{i}^{q}(0) + \Gamma_{i}^{q}(0)\right) \left(\Phi^{(.,i)} - \widehat{\psi}_{s}\right)\right) \\ &\leq \sum_{t=s-a_{n}}^{s+a_{n}-1} \|\varepsilon_{t}\|_{2}^{2} + cd_{n}^{\star}\log p + (2a_{n}) \left(\sqrt{\frac{\log p}{2a_{n}}} + \lambda_{\max}\left(\Gamma_{i}^{q}(0)\right)\right) \left\|\Phi^{(.,i)} - \widehat{\psi}_{s}\right\|_{1}^{2} \\ &\leq \sum_{t=s-a_{n}}^{s+a_{n}-1} \|\varepsilon_{t}\|_{2}^{2} + cd_{n}^{\star}\log p + c'd_{n}^{\star^{2}}\log p \\ &= \sum_{t=s-a_{n}}^{s+a_{n}-1} \|\varepsilon_{t}\|_{2}^{2} + O_{p}\left(d_{n}^{\star^{2}}\log p\right). \end{split}$$

where $\left\|\widehat{\psi}_{s} - \Phi^{(.,i)}\right\|_{1} \leq 4\sqrt{d_{n}^{\star}} \left\|\widehat{\psi}_{s} - \Phi^{(.,i)}\right\|_{2}, \left\|\widehat{\psi}_{s} - \Phi^{(.,i)}\right\|_{F} = O_{p}\left(\sqrt{\frac{d_{n}^{\star}\log p}{a_{n}}}\right),$ Combining (B.15) and (B.16) yields

$$D \geq \operatorname{LIC}(\widetilde{t}_{1},...,\widetilde{t}_{\widetilde{m}}) - \operatorname{LIC}(C)$$

$$\geq -c_{1}d_{n}^{\star}\log p - c_{2}d_{n}^{\star^{2}}\log p + \omega_{n}$$

$$\geq -2c_{2}d_{n}^{\star^{2}}\log p + \omega_{n} \qquad (B.17)$$

However, (B.17) contradicts Assumption A4 that $d_n^{\star 2} \log p/\omega_n \to 0$. This completes the proof of the second part.

It remains to show that $\mathbb{P}\left(\left|\operatorname{cluster}\left(\widetilde{\mathcal{A}}_{n}, 2a_{n}\right)\right| = m_{0}\right) \to 1$. For this, note that based on the previous calculations, all points in $\widetilde{\mathcal{A}}_{n}$ are either very close to a true break point,

i.e. in the $Kn\gamma_n$ -neighborhood of a true break point for some positive and finite K, or they are at least in the ca_n -neighborhood of a true break point for a constant 0 < c < 1. Next, fix one of the true break points, say t_j . Denote by B_j , all the selected break points in $\widetilde{\mathcal{A}}_n$ which are in the $Kn\gamma_n$ -neighborhood of t_j for some positive and finite K. Now, if there are any points in the a_n -neighborhood of t_j , add them to B_j as well. Note that since there is at least one estimated break point in the $n\gamma_n$ -neighborhood of t_j , B_j is not empty. Also, note that the diameter of each B_j is at most $2a_n$. Finally, based on the choice of a_n in Assumption A4, all B_j 's are disjoint, $j = 1, \ldots, m_0$. Therefore, the collection of all B_j 's form cluster $(\widetilde{\mathcal{A}}_n, 2a_n)$ which has cardinality equal to m_0 (note that due to Assumption A4, the collection of sets B_j 's is a minimal partitioning of set $\widetilde{\mathcal{A}}_n$ while the diameter of all B_j 's are at most $2a_n$). This completes the proof of the last part.

Proof of Theorem 3. Suppose for any constant K > 0, there exists some change point $t_j \in \mathcal{A}_n$ such that $\left| \tilde{t}_j^f - t_j \right| > K d_n^* \log p$. From Theorem 2, we know that for any subset $B_j \in$ cluster $\left(\tilde{\mathcal{A}}_n, 2a_n \right)$, there exists a true change point t_j lies in $(\min(B_j) - a_n, \max(B_j) + a_n)$. Therefore, there exists a true change point t_j lying within interval $[\tilde{t}_j^f - a_n, \tilde{t}_j^f + a_n)$. Assume, without loss of generality, that $\tilde{t}_j^f > t_j$.

Based on the local refinement algorithm, we define the loss function $L_n\left(\tilde{t}_j^f\right)$ as follows:

$$\sum_{k=\min(B_j)-a_n}^{\widetilde{t}_j^J-1} \left\| y_t - \widetilde{\psi}_{j,1} Y_{t-1} \right\|_2^2 + \sum_{t=\widetilde{t}_j^f}^{\max(B_j)+a_n-1} \left\| y_t - \widetilde{\psi}_{j,2} Y_{t-1} \right\|_2^2 \stackrel{\text{def}}{=} I_1 + I_2.$$

and the loss function $L_n(t_j)$ for the true change point t_j as:

$$\sum_{t=\min(B_j)-a_n}^{t_j-1} \left\| y_t - \widetilde{\psi}_{j,1} Y_{t-1} \right\|_2^2 + \sum_{t=t_j}^{\max(B_j)+a_n-1} \left\| y_t - \widetilde{\psi}_{j,2} Y_{t-1} \right\|_2^2 \stackrel{\text{def}}{=} I_3 + I_4.$$

where $\tilde{t}_j^f \in (l_j, u_j)$, $l_j = \min(B_j) - a_n$ and $u_j = \max(B_j) + a_n$, for $j = 1, \ldots, m_0 + 1$. $\tilde{\psi}_{j,1}$ and $\tilde{\psi}_{j,2}$ are the local coefficient parameter estimates given that a time point $s \in B_j$ is the break point.

Suppose $s > t_j$; then, similar arguments as in Lemma 4 of Safikhani and Shojaie (2020) yield:

$$\left\|\widetilde{\psi}_{j,1} - \Phi^{(.,j)}\right\|_{1} \le 4\sqrt{d_{n}^{\star}} \left\|\widetilde{\psi}_{j,1} - \Phi^{(.,j)}\right\|_{F}, \quad \left\|\widetilde{\psi}_{j,1} - \Phi^{(.,j)}\right\|_{F} = O_{p}\left(\sqrt{\frac{d_{n}^{\star}\log p}{\widetilde{R}_{n}}}\right), \quad (B.18)$$

and

t

$$\left\|\widetilde{\psi}_{j,2} - \Phi^{(.,j+1)}\right\|_{1} \le 4\sqrt{d_{n}^{\star}} \left\|\widetilde{\psi}_{j,2} - \Phi^{(.,j+1)}\right\|_{F}, \quad \left\|\widetilde{\psi}_{j,2} - \Phi^{(.,j+1)}\right\|_{F} = O_{p}\left(\sqrt{\frac{d_{n}^{\star}\log p}{\widetilde{R}_{n}}}\right). \tag{B.19}$$

To see this, observe that $\tilde{\psi}_{j,1}$ in (12) minimizes the least squares plus the ℓ_1 norm loss function. Therefore, the value of this objective function for $\tilde{\psi}_{j,1}$ will be smaller than any other choice of parameters, including $\Phi^{(.,j)}$. Hence,

$$\frac{1}{\widetilde{R}_{n}} \sum_{t=\widetilde{t}_{j}^{f}-\widetilde{R}_{n}}^{\widetilde{t}_{j}^{f}-1} \left\| y_{t}-\widetilde{\psi}_{j,1}Y_{t-1} \right\|_{2}^{2} + \widetilde{\eta}_{j,1} \left\| \widetilde{\psi}_{j,1} \right\|_{1} \leq \frac{1}{\widetilde{R}_{n}} \sum_{t=\widetilde{t}_{j}^{f}-\widetilde{R}_{n}}^{\widetilde{t}_{j}^{f}-1} \left\| y_{t}-\Phi^{(.,j)}Y_{t-1} \right\|_{2}^{2} + \widetilde{\eta}_{j,1} \left\| \Phi^{(.,j)} \right\|_{1}.$$
(B.20)

Note that the segment $[t_j, \tilde{t}_j^f)$ has at most length a_n , so the length of $[\tilde{t}_j^f - \tilde{R}_n, t_j)$ is at least $c\tilde{R}_n$, for some positive constant c > 0. In that case, the misspecification part is negligible. Some rearrangements together with the use of proposition 4.2 of Basu and Michailidis (2015) lead to:

$$\begin{split} &c_{1} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{F}^{2} - c_{2} \frac{\log p}{\tilde{R}_{n}} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1}^{2} &\leq \max \left(0, c_{1} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{F}^{2} - c_{2} \frac{\log p}{\tilde{R}_{n}} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1}^{2} \right) \\ &\leq \frac{1}{\tilde{R}_{n}} \sum_{t=\tilde{t}_{j}^{'} - \tilde{R}_{n}}^{\tilde{t}_{j}^{'} - 1} Y_{t-1}^{\prime} \left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right)^{\prime} \left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right) Y_{t-1} \\ &\leq \frac{2}{\tilde{R}_{n}} \sum_{t=\tilde{t}_{j}^{'} - \tilde{R}_{n}}^{\tilde{t}_{j}^{'} - 1} Y_{t-1}^{\prime} \left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right)^{\prime} \left(y_{t} - \Phi^{(.,j)} Y_{t-1} \right) + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1}^{-} - \left\| \tilde{\psi}_{j,1} \right\|_{1}^{2} \right) \\ &= \frac{2}{\tilde{R}_{n}} \sum_{t=\tilde{t}_{j}^{'} - \tilde{R}_{n}}^{\tilde{t}_{j}^{'} - 1} Y_{t-1}^{\prime} \left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right)^{\prime} \varepsilon_{t} + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1}^{-} - \left\| \tilde{\psi}_{j,1} \right\|_{1}^{2} \right) \\ &+ \frac{2}{\tilde{R}_{n}} \sum_{t=t_{j}}^{\tilde{t}_{j}^{'} - 1} Y_{t-1}^{\prime} \left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right)^{\prime} \left(\Phi^{(.,j+1)} - \Phi^{(.,j)} \right) Y_{t-1} \\ &\leq c \sqrt{\frac{\log p}{\tilde{R}_{n}}} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1}^{+} \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1}^{-} - \left\| \tilde{\psi}_{j,1} \right\|_{1}^{2} \right) \\ &+ 2\frac{\tilde{t}_{j}^{'} - t_{j}}{\tilde{R}_{n}} \operatorname{tr} \left(\left(\Phi^{(.,j)} - \tilde{\psi}_{j,1} \right)^{\prime} \left(\frac{1}{\tilde{t}_{j}^{'} - t_{j}} \sum_{t=t_{j}}^{\tilde{t}_{j}^{'} - 1} Y_{t-1} - \Gamma_{j+1}^{q} (0) + \Gamma_{j+1}^{q} (0) \right) \left(\Phi^{(.,j+1)} - \Phi^{(.,j)} \right) \right) \\ &= \mathrm{I}, \end{split}$$

where the second inequality is due to the restricted eigenvalue property (proposition 4.2 of Basu and Michailidis (2015)) and the fourth inequality uses the circular invariance property of the trace function. Next, observe that

$$\begin{split} \mathbf{I} &\leq c\sqrt{\frac{\log p}{\tilde{R}_{n}}} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1} + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1} - \left\| \tilde{\psi}_{j,1} \right\|_{1} \right) \\ &+ 2c' \frac{\tilde{t}_{j}^{f} - t_{j}}{\tilde{R}_{n}} \left(\max\left(\sqrt{\frac{\log p}{\tilde{t}_{j}^{f} - t_{j}}, \frac{\log p}{\tilde{t}_{j}^{f} - t_{j}} \right) + \lambda_{\max} \left(\Gamma_{j+1}^{q}(0) \right) \right) \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1} \left\| \Phi^{(.,j+1)} - \Phi^{(.,j)} \right\|_{1} \\ &\leq \left(c\sqrt{\frac{\log p}{\tilde{R}_{n}}} + 4c' M_{\Phi} d_{n}^{*} \frac{a_{n}}{\tilde{R}_{n}} \right) \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1} + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,i)} \right\|_{1} - \left\| \tilde{\psi}_{j,1} \right\|_{1} \right) \\ &\leq \max\left(c, 4c' M_{\Phi} \right) \sqrt{\frac{\log p}{\tilde{R}_{n}}} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1} + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1} - \left\| \tilde{\psi}_{j,1} \right\|_{1} \right) \\ &\leq \frac{\tilde{\eta}_{j,1}}{2} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1} + \tilde{\eta}_{j,1} \left(\left\| \Phi^{(.,j)} \right\|_{1} - \left\| \tilde{\psi}_{j,1} \right\|_{1} \right) \\ &\leq \frac{3\tilde{\eta}_{j,1}}{2} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1,\mathcal{I}_{j}} - \frac{\tilde{\eta}_{j,1}}{2} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1,\mathcal{I}_{j}^{c}} \\ &\leq 2\tilde{\eta}_{j,1} \left\| \Phi^{(.,j)} - \tilde{\psi}_{j,1} \right\|_{1}, \end{split}$$
(B.21)

where the first inequality is based on the first part of Lemma A1 with $c_n = \tilde{t}_j^f - t_j$; the third inequality is due to Assumption A3 and Assumption A4 under which $\log p \leq n\gamma_n \leq a_n$ and $\frac{\Delta_n}{4} \geq \tilde{R}_n = \frac{d_n^{\star 2} a_n^2}{\log p}$ such that $\sqrt{\frac{\log p}{\tilde{R}_n}} \geq d_n^{\star} \frac{a_n}{\tilde{R}_n}$; the fourth inequality is due the selection of the tuning parameter $\tilde{\eta}_{j,1} = \sqrt{\frac{\log p}{\tilde{R}_n}}$; and finally the fifth inequality is due the triangular inequality.

This ensures that $\left\| \Phi^{(.,j)} - \widetilde{\psi}_{j,1} \right\|_{1,\mathcal{I}_{j}^{c}} \leq 3 \left\| \Phi^{(.,j)} - \widetilde{\psi}_{j,1} \right\|_{1,\mathcal{I}_{j}}$, and hence $\left\| \Phi^{(.,j)} - \widetilde{\psi}_{j,1} \right\|_{1} \leq 4 \left\| \Phi^{(.,j)} - \widetilde{\psi}_{j,1} \right\|_{1,\mathcal{I}_{j}} \leq 4\sqrt{d_{n}^{\star}} \left\| \Phi^{(.,j)} - \widetilde{\psi}_{j,1} \right\|_{F}$. This comparison between L_{1} and L_{2} norms of the error term together with the bound in Equation (B.21) will get the desired consistency rates in (B.18). Verifying (B.19) follows exactly from proposition 4.1 of Basu and Michailidis (2015) since there are no break points in the interval $[\widetilde{t}_{j}^{f}, \widetilde{t}_{j}^{f} + \widetilde{R}_{n})$.

Using the results from (B.18) and (B.19), we obtain

$$\begin{split} I_{1} &= \sum_{i=l_{j}^{l}-a_{n}}^{l_{j}-1} \left\| y_{t} - \tilde{\psi}_{j,1}Y_{t-1} \right\|_{2}^{2} + \sum_{t-1}^{\tilde{r}_{j}^{l}-1} \left\| y_{t} - \tilde{\psi}_{j,1}Y_{t-1} \right\|_{2}^{2} \\ &\geq I_{3} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\tilde{\psi}_{j,1} - \Phi^{(,j+1)} \right) Y_{t-1} \right\|_{2}^{2} - 2 \left| \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} Y_{t-1}^{\prime} \left(\tilde{\psi}_{j,1} - \Phi^{(,j+1)} \right) e_{t} \right| \\ &= I_{3} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\tilde{\psi}_{j,1} - \Phi^{(,j)} + \Phi^{(,j)} - \Phi^{(,j+1)} \right) Y_{t-1} \right\|_{2}^{2} \\ &- 2 \left| \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\tilde{\psi}_{j,1} - \Phi^{(,j)} + \Phi^{(,j)} - \Phi^{(,j+1)} \right) Y_{t-1} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\tilde{\psi}_{j,1} - \Phi^{(,j)} \right) Y_{t-1} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\Phi^{(,j)} - \Phi^{(,j+1)} \right) Y_{t-1} \right\|_{2}^{2} \\ &- 2 \left| \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| \left(\Phi^{(,j)} - \Phi^{(,j+1)} \right) Y_{t-1} \right\|_{2}^{2} + \sum_{t-1,j}^{\tilde{r}_{j}^{l}-1} \left\| e_{t} \right\|_{2}^{2} + \int \left| e_{t} \right|_{2}^{2} + \int \left| e_{t} \right|_{$$

The second term inequality in (i) holds by Hölder's inequality; the third term inequality

in (i) holds by the deviation bound condition; inequality (ii) holds by $\|\Phi^{(.,j+1)} - \Phi^{(.,j)}\|_1 \leq \sqrt{d_n^{\star}} \|\Phi^{(.,j+1)} - \Phi^{(.,j)}\|_F$; inequality (iii) holds by the fact that $\|\Phi^{(.,j+1)} - \Phi^{(.,j)}\|_F^2 \geq v^2$, where v is a positive constant; and inequality (iv) holds by $\left|\tilde{t}_j^f - t_j\right| > K d_n^{\star} \log p$ and choosing large enough constant K > 0 such that $\sqrt{\frac{1}{K}} < \frac{1}{4}$.

Similarly, we have

$$\begin{split} I_{4} &= \sum_{t=t_{j}}^{\max(B_{j})+a_{n}-1} \left\| y_{t} - \widetilde{\psi}_{j,2}Y_{t-1} \right\|_{2}^{2} \\ &= \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \left\| y_{t} - \widetilde{\psi}_{j,2}Y_{t-1} \right\|_{2}^{2} + \sum_{t=\tilde{t}_{j}^{f}}^{\max(B_{j})+a_{n}-1} \left\| y_{t} - \widetilde{\psi}_{j,2}Y_{t-1} \right\|_{2}^{2} \\ &\leq I_{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \left\| \varepsilon_{t} \right\|_{2}^{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \left\| \left(\widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right) Y_{t-1} \right\|_{2}^{2} + 2 \left\| \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} Y_{t-1}^{\prime} \left(\widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right) \varepsilon_{t} \right\| \\ &\stackrel{(i)}{\leq} I_{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{r}-1} \left\| \varepsilon_{t} \right\|_{2}^{2} + c^{\prime}a_{n} \left\| \widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right\|_{F}^{2} + c^{\prime\prime}\sqrt{a_{n}(\log p)} \left\| \widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right\|_{1}^{1} \\ &\leq I_{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{r}-1} \left\| \varepsilon_{t} \right\|_{2}^{2} + c^{\prime}a_{n} \left\| \widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right\|_{F} \left(\left\| \widetilde{\psi}_{j,2} - \Phi^{(..j+1)} \right\|_{F} + \frac{c^{\prime\prime}}{c} \sqrt{\frac{d_{n}^{\star}(\log p)}{a_{n}}} \right) \\ &\stackrel{(i)}{\leq} I_{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \left\| \varepsilon_{t} \right\|_{2}^{2} + K_{2}d_{n}^{\star}(\log p), \end{split}$$
(B.23)

where the second term in inequality (i) holds by the upper-RE condition and the fact that $\left\| \Phi^{(.,j+1)} - \widetilde{\psi}_{j,2} \right\|_{1}^{2} \leq 16d_{n}^{\star} \left\| \Phi^{(.,j+1)} - \widetilde{\psi}_{j,2} \right\|_{F}^{2}$; the third term in inequality (i) holds by deviation bound condition; inequality (ii) holds by the result (B.19). Note that based on Assumption A4 that $\frac{\Delta_{n}}{4} \geq \widetilde{R}_{n} = \frac{d_{n}^{\star 2}a_{n}^{2}}{\log p}$, we have $\left\| \Phi^{(.,j+1)} - \widetilde{\psi}_{j,2} \right\|_{F} \leq c\sqrt{\frac{d_{n}\log p}{a_{n}}}$.

Next, based on the definition of (11), we have $L_n\left(\tilde{t}_j^f\right) \leq L_n(t_j)$ and then

$$I_{3} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \|\boldsymbol{\varepsilon}_{t}\|_{2}^{2} + K_{1} \left| \tilde{t}_{j}^{f} - t_{j} \right| + I_{2} \leq L_{n}(\tilde{t}_{j}^{f}) \leq L_{n}(t_{j}) \leq I_{3} + I_{2} + \sum_{t=t_{j}}^{\tilde{t}_{j}^{f}-1} \|\boldsymbol{\varepsilon}_{t}\|_{2}^{2} + K_{2}d_{n}^{\star}(\log p),$$

which leads to

$$\left|\tilde{t}_{j}^{f} - t_{j}\right| \le K_{j} d_{n}^{\star} \log p. \tag{B.24}$$

This contradicts the setting. Under Assumptions A1 and A2, we set $K^* = \max_{1 \le j \le m_0} K_j$ and complete the proof.

Proof of Theorem 4. The proof of this theorem follows along similar lines to Theorem 4 of Safikhani and Shojaie (2020), where it is also shown that, if m_0 is known, then it is enough to set $R_n = n\gamma_n$.

We need to firstly verify two important conditions: (1) the restricted eigenvalue (RE) condition for $\hat{\Gamma} = I_p \otimes (\mathcal{X}'_{\mathbf{r}} \mathcal{X}_{\mathbf{r}}/N)$; and (2) the deviation bound condition for $\|\hat{\gamma} - \hat{\Gamma}\Phi\|_{\infty}$ where $\hat{\gamma} = (I_p \otimes \mathcal{X}'_{\mathbf{r}})\mathbf{Y}_{\mathbf{r}}/N$. Once these two are verified, the rest of the proof is applying deterministic arguments used in Proposition 4.1 in Basu and Michailidis (2015).

Condition (1) implies that there exist $\alpha, \tau > 0$ such that for any $\theta \in \mathbb{R}^{\tilde{\pi}}$, we have

$$\theta' \hat{\Gamma} \theta \ge \alpha \|\theta\|_2^2 - \tau \|\theta\|_1^2$$

with probability at least $1 - c_1 \exp(-c_2 N)$ for large enough constants $c_1, c_2 > 0$. Based on Lemma B.1 in Basu and Michailidis (2015), it is enough to show the RE for $S = \mathcal{X}'_i \mathcal{X}_i / N$, where \mathcal{X}_i is the *i*th block component of \mathcal{X}_r . Applying Proposition 2.4 of Basu and Michailidis (2015), we have for any \mathbb{R}^{pq} with $||v||_2 \leq 1$, and any $\eta > 0$:

$$\mathbb{P}\Big(\Big|v'\big(S - \frac{N_i}{N}\Gamma_i(0)\big)v\Big| > c\eta\Big) \le 2\exp(-c_3N\min(\eta^2,\eta)).$$

Next, to make the above probability hold uniformly over all vectors v, we apply the discretization Lemma F2 in Basu and Michailidis (2015) and also Lemma 12 in the supplement of Loh and Wainwright (2012) to get:

$$\left|v'\left(S - \frac{N_i}{N}\Gamma_i(0)\right)v\right| \le \alpha \|v\|_2^2 + \alpha/k\|v\|_1^2,$$

with probability at least $1 - c_1 \exp(c_2 N)$, for all $v \in \mathbb{R}^{pq}$, some $\alpha > 0$ and with an integer $k = \lfloor c_4 N / \log(pq) \rfloor$ with some $c_4 > 0$. This implies that

$$v'Sv \ge v'\frac{N_i}{N}\Gamma_i(0)v - \alpha \|v\|_2^2 - \alpha/k\|v\|_1^2 \ge \alpha \|v\|_2^2 - \alpha/k\|v\|_1^2,$$

since $N_i \ge \Delta_n - 4R_n$, $N = n + q - 1 - 2m_0R_n$, and assuming $\Delta_n \ge \epsilon n$ implies that $N_i/N \ge \epsilon \ge 2\alpha$.

Condition (2) means that there exists a large enough constant C' > 0 such that

$$\|\hat{\gamma} - \hat{\Gamma}\Phi\|_{\infty} \le C' \sqrt{\frac{\log \tilde{\pi}}{N}}$$

with probability at least $1 - c_1 \exp(-c_2 \log \tilde{\pi})$. To verify this condition here, observe that $\hat{\gamma} - \hat{\Gamma} \Phi = \operatorname{vec}(\mathcal{X}'_{\mathbf{r}} E_r)/N$. Therefore, denoting the *h*-th column block of $\mathcal{X}_{\mathbf{r}}$ by $\mathcal{X}_{\mathbf{r},(h)}$, for $h = 1, \ldots, (m_0 + 1)q$, we have:

$$\|\hat{\gamma} - \hat{\Gamma}\Phi\|_{\infty} = \max_{1 \le k, l \le p; 1 \le h \le (m_0 + 1)q} \left| e'_k \mathcal{X}'_{\mathbf{r},(h)} E_r e_l \right|.$$
(B.25)

Now, for a fixed k, l, h, applying Proposition 2.4(b) in Basu and Michailidis (2015) gives:

$$\mathbb{P}(\left|e_k'\mathcal{X}_{\mathbf{r},(h)}'E_r e_l\right| > k_1\eta) \le 6\exp(-k_2N\min(\eta^2,\eta)),$$
(B.26)

for large enough $k_1, k_2 > 0$, and any $\eta > 0$. Now, setting $\eta = C' \sqrt{\frac{\log \pi}{N}}$, and taking the union over all the $\tilde{\pi}$ cases for k, l, h yield the desired result. This completes the proof of this theorem.

Proof of Corollary 1. The proof of this Corollary is similar to that of Proposition 4.1 in Basu and Michailidis (2015). The key steps in this proof are similar to the previous proof of Theorem 4. Hence, we only highlight the main differences compared to Theorem 4; see more details in the proof of Theorem 4.

Instead of estimating the $m_0 + 1$ transition matrices simultaneously, we separately estimate the model parameters in each segment. For the *i*-the segment, $i = 1, \dots, m_0 + 1$, let $\hat{\Gamma}_i = (I_p \otimes \mathcal{X}'_i \mathcal{X}_i / N_i)$ and $\hat{\gamma}_i = (I_p \otimes \mathcal{X}'_i) \mathcal{Y}_i / N_i$. We need to verify the restricted eigenvalue (RE) condition for $\hat{\Gamma}_i$ and the deviation bound condition for $\|\hat{\gamma}_i - \hat{\Gamma}_i \Phi^{(.,i)}\|_{\infty}$. Once these two are verified, the rest of the proof is applying deterministic arguments used in Proposition 4.1 in Basu and Michailidis (2015).

Similar to Theorem 4, it suffices to show the RE for $S = \mathcal{X}'_i \mathcal{X}_i / N_i$. Then, we have

$$v'Sv \ge v'\Gamma_i(0)v - \alpha \|v\|_2^2 - \alpha/k \|v\|_1^2 \ge \alpha \|v\|_2^2 - \alpha/k \|v\|_1^2,$$

with probability at least $1 - c_1 \exp(c_2 N_i)$, for all $v \in \mathbb{R}^{pq}$, some $\alpha > 0$ and with an integer $k = \lfloor c_4 N_i / \log(pq) \rfloor$ with some $c_4 > 0$

The DB condition implies that there exists a large enough constant C' > 0 such that

$$\|\hat{\gamma}_i - \hat{\Gamma}_i \Phi^{(.,i)}\|_{\infty} \le C' \sqrt{\frac{\log p^2 q}{N_i}},$$

with probability at least $1 - c_1 \exp(-c_2 \log(p^2 q))$. To verify this condition here, observe that $\hat{\gamma}_i - \hat{\Gamma}_i \Phi^{(.,i)} = \operatorname{vec}(\mathcal{X}'_i E_i)/N_i$. We replace $\mathcal{X}_{\mathbf{r}}$, E_r and N with \mathcal{X}_i , E_i and N_i in (B.25) and (B.26). Next, setting $\eta = C' \sqrt{\frac{\log p^2 q}{N_i}}$, and taking the union over all $p^2 q$ cases for k, l, h yields the desired result. This completes the proof.

Appendix C: Details on the Detection Algorithm

Next, we provide details of the algorithm that solves optimization problem (4).

Let $S(.; \lambda)$ be the element-wise soft-thresholding operator which maps its input x to $x - \lambda$ when $x > \lambda$, $x + \lambda$ when $x < -\lambda$, and 0 when $|x| \le \lambda$. Recall that throughout the paper, for a $m \times n$ matrix A, $||A||_{\infty} = \max_{1 \le i \le m, 1 \le j \le n} |a_{ij}|$. recall that $Y'_l = (y'_l \dots y'_{l-q+1})_{1 \times pq}$. For two time points s < t, define $Y_{(s,t)} = (Y_s, Y_{s+1}, \dots, Y_t) \in \mathbb{R}^{pq \times (t-s+1)}$ and $y_{(s,t)} = (y_s, y_{s+1}, \dots, y_t) \in \mathbb{R}^{p \times (t-s+1)}$.

The main steps of the algorithm are as follows:

- (i) Set the initial values for all parameters to zero; i.e. $\theta_i^{(0)} = 0$, for $i = 1, \ldots, k_n$.
- (ii) For each $i = 1, ..., k_n$, calculate the (h+1)-th iteration of the parameters $\theta_i^{(h+1)}$ using the KKT condition as follows:

$$\theta_{i}^{\prime\,(h+1)} = \left(\sum_{l=i}^{k_{n}} Y_{(r_{i-1},r_{i})} Y_{(r_{i-1},r_{i})}^{\prime}\right)^{-1} S\left(\sum_{l=i}^{k_{n}} Y_{(r_{i-1},r_{i})} y_{(r_{i},r_{i+1})}^{\prime} - \sum_{j \neq i} \left(\sum_{l=\max(i,j)}^{k_{n}} Y_{(r_{i-1},r_{i})} Y_{(r_{i-1},r_{i})}^{\prime}\right) \theta_{j}^{\prime\,(h)}; \lambda_{1,n}\right)$$

$$(C.1)$$

- (iii) (a) If $\max_{1 \le i \le k_n} \|\theta_i^{(h+1)} \theta_i^{(h)}\|_{\infty} < \delta$, where δ is a tolerance threshold set to 10^{-3} in the implementation, stop and denote the final estimate by $\Theta^{(intermediate)}$.
 - (b) If $\max_{1 \le i \le k_n} \|\theta_i^{(h+1)} \theta_i^{(h)}\|_{\infty} \ge \delta$, set h = h + 1. Go to step (ii).
- (iv) Apply soft-thresholding to $\Theta^{(intermediate)}$ to find the optimizer in equation (4). In other words, $\widehat{\Theta} = S(\Theta^{(intermediate)}; \lambda_{2,n}).$

Note that in this algorithm, the whole block of θ_i with p^2q elements is updated once, which reduces the computation time dramatically.

Appendix D: Tuning Parameters Selection

There are a number of tuning parameters in BSS: $\lambda_{1,n}$, $\lambda_{2,n}$, a_n , η_n , ω_n , ρ_n and R_n . Although their asymptotic values have been presented and discussed in the previous section, their selection in finite samples merits further discussion. Data-driven methods are developed to select them¹:

- $\lambda_{1,n}, \lambda_{2,n}$: Both $\lambda_{1,n}$ and $\lambda_{2,n}$ can be selected through cross-validation. For splitting the training set and the validation set, we randomly select 20% of the blocks equally spaced with a random initial point. Denote the last time point in these selected blocks by \mathcal{T} . The data without observations in \mathcal{T} can then be used in the first step of our procedure to estimate Θ for a range of values for $\lambda_{1,n}$ and $\lambda_{2,n}$. The parameters estimated in the first step are then used to predict the series at time points in \mathcal{T} . The value of $\lambda_{1,n}$ and $\lambda_{2,n}$ which minimize the mean squared prediction error over \mathcal{T} is the cross-validated choice of $\lambda_{1,n}$ and $\lambda_{2,n}$. The sequence for $\lambda_{1,n}, \lambda_{2,n}$ are selected as follows. Similar to Friedman et al. (2010), we construct a sequence of K_1 values for $\lambda_{1,n}$, decreasing from $\lambda_{1,\max}$ to $\lambda_{1,\min}$ on the log scale, where the maximum value $\lambda_{1,\max}$ is the smallest value for which the entire estimated parameter $\hat{\theta}_i = 0$, for all $i = 1, \ldots, k_n$ while the minimum value is set to be $\lambda_{1,\min} = \epsilon \lambda_{1,\max}$. We choose $\lambda_{2,n} = c \sqrt{\frac{\log p}{T}}$, where c is a decreasing sequence of K_2 values. In the simulation study, we choose $K_1 = 10$, $K_2 = 3, \epsilon = 10^{-3}$ when the blocks size b_n is smaller than 2p and $\epsilon = 10^{-4}$ otherwise.
 - a_n : This tuning parameter can be selected through a grid search. We apply an exhaustive search procedure on a grid of a_n 's ranked from the minimum to the maximum and record the number of selected break points for each a_n , and then stop this process when the number of break points selected does not change any more. In detail, we select $a_n^{(1)}, a_n^{(2)}, \ldots, a_n^{(\ell)}$ as an equally spaced sequence from the interval $[a_n^{(lb)}, a_n^{(ub)}]$ in a increasing order where $a_n^{(lb)} = \max(\lfloor \overline{b_n} \rfloor, \lfloor \log n \log p \rfloor), a_n^{(ub)} = \min(10a_n^{(lb)}, (\widehat{t}_1 - q - 1), (T - q - \widehat{t_m} - 1))$, and $\overline{b_n}$ is the mean of block sizes. Denote the number of selected break points using BSS with $a_n^{(i)}$ as the neighborhood size by n_i , for $i = 1, 2, \ldots, \ell$. The optimal neighborhood size can be defined as the first time n_i

¹The R/Rcpp codes to perform the BSS algorithm are available at the author's GitHub page: https://github.com/abolfazlsafikhani/BSS-ChangePoint-VAR

remains unchanged. In other words, $a_n^{(\ell^*)}$ is the optimal neighborhood size when $\ell^* = \min\{1 \le i \le \ell : n_i = n_{i+1} = n_{i+2}\}$. In all simulation scenarios, we set $\ell = 5$.

- η_n : Despite the fact that based on Assumption A5, different treatments are needed for selecting tuning parameters $\eta_{\hat{t}_i}$, $\eta_{\hat{t}_{i,1}}$, $\eta_{\hat{t}_{i,2}}$, all $\eta_{\hat{t}_i}$, $\eta_{\hat{t}_{i,1}}$ and $\eta_{\hat{t}_{i,2}}$ are suggested to be set to $(\log(2a_n)\log p)/(2a_n)$, $i = 1, \ldots, \hat{m}$. This choice was used in all of the numerical analyses undertaken and provides very good results.
- ω_n : The selection of ω_n is the most difficult, since it depends on how large changes in the VAR parameters must be in order to consider them as break points in finite sample applications. Here, we consider using a data-driven method for selecting ω_n . The idea is to cluster the changes in the objective function L_n into two subgroups, small and large. The proposed algorithm is summarized as follow:
 - Denote the candidate break points selected in the first step as $\hat{t}_1, \ldots, \hat{t}_{\widehat{m}}$. For each $k = 1, 2, \ldots, \hat{m}$, compute $v_k = L_n(\widehat{\mathcal{A}}_n \setminus \{\widehat{t}_k\}; \eta_n) L_n(\widehat{\mathcal{A}}_n; \eta_n)$.
 - Consider two boundary points $t_1^r = a_n + q$ and $t_2^r = T a_n$ as reference points. Compute $v_i^r = L_n(\widehat{\mathcal{A}}_n; \eta_n) - L_n(\widehat{\mathcal{A}}_n \cup \{t_i^r\}; \eta_n)$ for i = 1, 2, and set the reference value as $v^r = \max(v_1^r, v_2^r)$.
 - Combine the jumps v_k for each candidate break points and the reference value v^r (with 2 replicates) into one vector $V = (v_1, v_2, \ldots, v_{\hat{m}}, v^r, v^r)$. Apply kmeans clustering algorithm (Hartigan and Wong, 1979) to the vector V with two centers. Denote the sub-vector with smaller center as the small subgroup, V_S , and the other sub-vector as the large subgroup, V_L .
 - If (between-group SS/total SS) in k-means clustering is high and the reference value v^r is not in V_L , set $\omega_n = \min V_L$; otherwise, set $\omega_n = \max V$.

One could also combine the k-means clustering method (Hartigan and Wong, 1979) with the BIC criterion (Schwarz et al., 1978) to cluster the changes in the parameter matrix into two subgroups.

 ρ_n : We select ρ_n as the minimizer of the combined Bayesian Information Criterion (BIC) for all the segments. Following Lütkepohl (2005) and Zou et al. (2007), for $j = 0, \ldots, \tilde{m}$ we define the BIC on the interval $I_{j+1} = [r_{j2}, r_{(j+1)1}]$ as

$$\operatorname{BIC}(j,\rho_n) = \log(\operatorname{det}\widehat{\Sigma}_{\varepsilon,j}) + \frac{\log(r_{(j+1)1} - r_{j2})}{(r_{(j+1)1} - r_{j2})} \left\|\widehat{\beta}_{j+1}\right\|_0,$$

where $\widehat{\Sigma}_{\varepsilon,j}$ is the residual sample covariance matrix with $\widehat{\mathbf{B}}$ estimated in (16), and $\|\widehat{\beta}_{j+1}\|_0$ is the number of nonzero elements in $\widehat{\beta}_{j+1}$; then ρ_n is selected as

$$\widehat{\rho}_n = \operatorname{argmin}_{\rho_n} \sum_{j=0}^m \operatorname{BIC}(j, \rho_n).$$
(D.1)

 R_n : Recall from Section 3 that we need to remove the selected break points together with their R_n -radius neighborhood before estimating the parameters using (16). In practice, the radius R_n needs to be estimated. However, a closer look into Theorem 2 together with Assumption A4 suggest that a_n can be chosen as the radius R_n . Therefore, in all simulation scenarios and data applications, we set $R_n = a_n$.

Appendix E: Comparison with SBS and DCBS

We record the detection accuracy of the proposed BSS strategy and compare it to two binary segmentation-based methods, the SBS (Cho and Fryzlewicz, 2015) and DCBS (Cho, 2016) methods (both are implemented in the R package "hdbinseg"), geared towards detection of multiple break points in multivariate time series data. Here, we consider the BSS method with three different block size settings: the large block size $b_n = \lfloor n^{\frac{1}{2}} \rfloor$, the medium block size $b_n = \lfloor n^{\frac{1}{2}} \rfloor$ and the small block size $b_n = \lfloor n^{\frac{1}{2}} \rfloor$, where n = T - q. Six additional simulation scenarios (G.1 - G.6) are considered for this comparison with model parameter values summarized in Table 1. The true coefficient matrices are similar to simulation B, as depicted in Figure 2 (top panel) with repeated entries -0.6, 0.6, and -0.6 off the main diagonal. Details of the simulation settings are as follows:

Setting G (detection comparison). Similar to scenario E, there are many true break points in the data generating process. In scenario G, T = 4,000, p = 10, q = 1 with break points being equally spaced: $\lfloor \frac{T}{m_0+1} \rfloor$, $\lfloor \frac{2T}{m_0+1} \rfloor$, \dots , $\lfloor \frac{m_0T}{m_0+1} \rfloor$. In scenarios G.1 though G.6, the true number of break points are $m_0 = 2, 4, 6, 8, 10, 12$, respectively.

| Sim | T | p | AR order q | block size b_n | m_0 | AR structure |
|-----|-------|----|--------------|-----------------------------------------------------------------------------------------------------------|-------|--------------|
| G.1 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple |
| G.2 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 4 | simple |
| G.3 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 6 | simple |
| G.4 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 8 | simple |
| G.5 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 10 | simple |
| G.6 | 4,000 | 10 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 12 | simple |

Table 1: Details of model parameters for simulation settings G.

The median of selection rates among m_0 true break points for all scenarios in setting G are summarized in the left panel in Figure 1. This plot clearly shows that in all settings considered, BSS has an advantage over SBS and DCBS in detection performance while this advantage gets more significant by increasing m_0 . The selection rate for the BSS method with small and medium block sizes remains above 98%, while the selection rate for the SBS and DCBS methods are ~ 80% for $m_0 = 4, 6$ and drop to ~ 10% for $m_0 = 10, 12$. Note that the selection rate for BSS with large block size also remains above ~ 95% with the exception of ~ 40% for the case of $m_0 = 12$ (still higher than the selection rate of SBS and DCBS in this case). The Hausdorff distance $d_H\left(\widetilde{\mathcal{A}}_n^f, \mathcal{A}_n\right)$ between the set of estimated break points and the set of true break points is a reasonable measure for estimation accuracy. The right panel in Figure 1 illustrates the performance of all three methods in terms of this estimation accuracy (averaged over 100 replicates) in which the BSS outperforms the SBS

and DCBS methods in almost all settings, while the advantage of BSS is more significant for larger m_0 values. For example, in the case of 6 true break points, the average estimation accuracy for SBS is around 677.29, while the same quantity for BSS (large, medium, small block sizes) is around 3.50, 0.84 and 0.10, respectively. Therefore, even the BSS method with large block sizes reduces the estimation error in locating the break points by around 99.5% compared to SBS.



Figure 1: Simulation G results: (left) median selection rate for the BSS, SBS and DCBS methods; (right) median Hausdorff distance $d_H\left(\widetilde{\mathcal{A}}_n^f, \mathcal{A}_n\right)$ for the BSS, SBS and DCBS methods.

Next, we compare the computation time for three methods (BSS, SBS and DCBS). Six additional simulation scenarios (H.1-H.6) are considered for comparison among BSS, SBS and DCBS with model parameter values summarized in Table 2. Details of the simulation settings are as follows:

Setting H (computation time comparison). In scenario H, p = 20, q = 1, $m_0 = 2$, $t_1 = \lfloor \frac{T}{3} \rfloor$, $t_2 = \lfloor \frac{2T}{3} \rfloor$. The auto-regressive coefficients are chosen to have the same simple 1-off diagonal structure as in Scenario A.1 as shown in the top left panel of Figure 2 with repeated entries -0.6, 0.6, and -0.6 off the main diagonal. The sample size for scenarios H.1 through H.6 are T = 1000, 2000, 3000, 4000, 5000, 6000, respectively.

| | | | - | | | 0 | |
|-----|-------|----|--------------|-----------------------------------------------------------------------------------------------------------|-------|--------------|--|
| Sim | T | p | AR order q | block size b_n | m_0 | AR structure | |
| H.1 | 1,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |
| H.2 | 2,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |
| H.3 | 3,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |
| H.4 | 4,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |
| H.5 | 5,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |
| H.6 | 6,000 | 20 | 1 | $\lfloor n^{\frac{1}{2}} \rfloor$, $\lfloor n^{\frac{2}{5}} \rfloor$, $\lfloor n^{\frac{1}{3}} \rfloor$ | 2 | simple | |

Table 2: Details of model parameters for simulation settings H.

The average computation time over 100 replicates for simulation setting H is plotted in Figure 2. BSS with large block sizes is the fastest method overall, while DCBS is the slowest one. It is worth noting that BSS with small block sizes remains faster than both SBS and DCBS, while its estimation accuracy and selection rate are the best over all these three methods as explained in simulation setting G. In this numerical experiment, the reduction in computation time in BSS (small block size) compared to SBS and DCBS is around 30%



Figure 2: Average computational time for the BSS (large, medium, small block sizes), SBS and DCBS methods .

and 55%, respectively, while BSS with medium and large block sizes achieves even higher time reduction.

References

- Basu, S. and G. Michailidis (2015). Regularized estimation in sparse high-dimensional time series models. *The Annals of Statistics* 43(4), 1535–1567.
- Cho, H. (2016). Change-point detection in panel data via double cusum statistic. *Electronic Journal of Statistics* 10(2), 2000–2038.
- Cho, H. and P. Fryzlewicz (2015). Multiple-change-point detection for high dimensional time series via sparsified binary segmentation. *Journal of the Royal Statistical Society:* Series B (Statistical Methodology) 77(2), 475–507.
- Friedman, J., T. Hastie, and R. Tibshirani (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software 33*(1), 1.
- Hartigan, J. A. and M. A. Wong (1979). Algorithm as 136: A k-means clustering algorithm. Journal of the Royal Statistical Society. Series C (Applied Statistics) 28(1), 100–108.
- Loh, P.-L. and M. J. Wainwright (2012, 06). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. Ann. Statist. 40(3), 1637–1664.
- Lütkepohl, H. (2005). New introduction to multiple time series analysis. Springer Science & Business Media.
- Safikhani, A. and A. Shojaie (2020). Joint structural break detection and parameter estimation in high-dimensional non-stationary var models. *Journal of American Statistical Association (Theory and methods)*, To Appear.
- Schwarz, G. et al. (1978). Estimating the dimension of a model. The annals of statistics 6(2), 461-464.
- Zou, H., T. Hastie, and R. Tibshirani (2007). On the "degrees of freedom" of the lasso. *The Annals of Statistics 35*(5), 2173–2192.