

Appendix

A.1 Deviation of the Lower Bound

Given the assumption of the Markov property and fixed state distributions we can first split the upper bound into two integrals given by

$$\begin{aligned}
& \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\frac{p(\boldsymbol{\tau}, \boldsymbol{\theta})}{p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta})} \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau} \\
&= \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\frac{p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t)}{q(\boldsymbol{\theta})} \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau} \\
&\quad + \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\frac{p(s_1) \prod_{t=1}^T p(s_{t+1} | a_t, s_t)}{p(s_1) \prod_{t=1}^T p(s_{t+1} | a_t, s_t) \pi'(a_t | s_t)} \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau},
\end{aligned} \tag{1}$$

where $\pi'(a_t | s_t)$ is the old policy. Simplifications and the fixed nature of $\pi'(a_t | s_t)$ leads then directly to the proposed simple form given by

$$\begin{aligned}
&= \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\frac{p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t)}{q(\boldsymbol{\theta})} \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau} \\
&\quad - \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\prod_{t=1}^T \pi'(a_t | s_t) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau} \\
&= \iint p_{\text{old}}(\boldsymbol{\tau})q(\boldsymbol{\theta}) \log \left(\frac{p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t)}{q(\boldsymbol{\theta})} \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\theta} d\boldsymbol{\tau} \\
&\quad + \text{const.}
\end{aligned} \tag{2}$$

A.2 Applying the Idea of Variational Bayes

Given the simplified form from A.1 we can apply the techniques of Variational Bayes for the estimation of the parameters Θ . In this section we will derive the approximated distributions step by step from the original integral. For simplicity the parameters are labelled with Θ_i and we assume that the factorization $q(\Theta) = \prod_i q(\Theta_i)$ holds. The derivations are then given by

$$\begin{aligned}
& \int \int p_{\text{old}}(\boldsymbol{\tau}) q(\Theta) \log \left(\frac{p(\Theta) \prod_{t=1}^T \pi(a_t | \Theta, s_t)}{q(\Theta)} \right) p(r = 1 | \boldsymbol{\tau}) d\Theta d\boldsymbol{\tau} \\
&= \int \prod_i q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \left(\log \left(p(\Theta) \prod_{t=1}^T \pi(a_t | \Theta, s_t) \right) - \sum_k \log q_k(\Theta_k) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta \\
&= \int_{\Theta_j} \int_{\Theta_{-j}} q_j(\Theta_j) \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \left(\log \left(p(\Theta) \prod_{t=1}^T p(a_t | \Theta, s_t) \right) - \sum_k \log q_k(\Theta_k) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j} d\Theta_j \\
&= \int_{\Theta_j} \int_{\Theta_{-j}} q_j(\Theta_j) \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \log \left(p(\Theta) \prod_{t=1}^T p(a_t | \Theta, s_t) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j} d\Theta_j \\
&\quad - \int_{\Theta_j} \int_{\Theta_{-j}} q_j(\Theta_j) \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \left(\sum_{k \neq j} \log q_k(\Theta_k) + \log q_j(\Theta_j) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j} d\Theta_j \\
&= \int_{\Theta_j} q_j(\Theta_j) \int_{\Theta_{-j}} \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \log \left(p(\Theta) \prod_{t=1}^T p(a_t | \Theta, s_t) \right) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j} d\Theta_j \\
&\quad - \int_{\Theta_j} q_j(\Theta_j) \int p_{\text{old}}(\boldsymbol{\tau}) \log q_j(\Theta_j) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_j + \text{const.}
\end{aligned} \tag{3}$$

Omitting the constant term we find that the maximization of this formula is given when the inner parts of the two separated integrals over Θ_j are equal. Thus, we can find that the maximization is given by

$$\begin{aligned}
& \int p_{\text{old}}(\boldsymbol{\tau}) \log q_j(\Theta_j) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} = \int_{\Theta_{-j}} \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \log \prod_{t=1}^T \pi(a_t, \Theta | s_t) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j} \\
& \Leftrightarrow \log q_j(\Theta_j) = \left(\int p_{\text{old}}(\boldsymbol{\tau}) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} \right)^{-1} \int_{\Theta_{-j}} \prod_{i \neq j} q_i(\Theta_i) \int p_{\text{old}}(\boldsymbol{\tau}) \log \prod_{t=1}^T \pi(a_t, \Theta | s_t) p(r = 1 | \boldsymbol{\tau}) d\boldsymbol{\tau} d\Theta_{-j},
\end{aligned} \tag{4}$$

where the last line is the final approximation.

A.3 Parameter Distribution for \mathbf{M}

The distribution of $q_{\mathbf{M}}(\mathbf{M})$ can be derived by

$$\begin{aligned}
\log q_{\mathbf{M}}(\mathbf{M}) &= \mathbb{E}_{-\mathbf{M}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \log \left(p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{M}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(\log \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t) \right. \right. \right. \\
&\quad \left. \left. \left. + \log \prod_{m=1}^M \mathcal{N} \left(\mathbf{m}^{(m)} | \mathbf{m}_{\text{old}}^{(m)}, \sigma^2 \mathbf{I} \right) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{M}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(\sum_{m=1}^M \sum_{j=1}^{D_m} \sum_{t=1}^T \log \pi(a_{tj}^{(m)} | \boldsymbol{\theta}, s_t) + \right. \right. \right. \\
&\quad \left. \left. \left. \sum_{m=1}^M \sum_{j=1}^{D_m} \mathcal{N} \left(\mathbf{m}_{j,:}^{(m)\top} | \mathbf{m}_{\text{old}_{j,:}}^{(m)\top}, \sigma^2 \mathbf{I} \right) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{M}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(\sum_{m=1}^M \sum_{j=1}^{D_m} \sum_{t=1}^T \right. \right. \right. \\
&\quad \left. \left. \left. \left(-\frac{1}{2} \left(\left(a_{t,j}^{(m)} - \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t \right) - \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right)^{\top} \left(\text{Tr} \left(\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} \right) \tilde{\tau}_m^{-1} \mathbf{I} \right)^{-1} \right. \right. \right. \\
&\quad \left. \left. \left. \cdot \left(\left(a_{t,j}^{(m)} - \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t \right) - \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right) \right) \right) \right. \right. \\
&\quad \left. \left. \left. - \sum_{m=1}^M \sum_{j=1}^{D_m} \frac{1}{2} \left(\mathbf{m}_{j,:}^{(m)} - \mathbf{m}_{\text{old}_{j,:}}^{(m)} \right) \sigma^{-2} \left(\mathbf{m}_{j,:}^{(m)} - \mathbf{m}_{\text{old}_{j,:}}^{(m)} \right)^{\top} \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{M}} \left[\sum_{m=1}^M \sum_{j=1}^{D_m} \mathbb{E}_{p(\boldsymbol{\tau})} \left[R(\boldsymbol{\tau}) \left(\sum_{t=1}^T \text{Tr} \left(\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} \right)^{-1} \tilde{\tau}_m \right. \right. \right. \\
&\quad \left. \left. \left. \cdot \left(\boldsymbol{\Phi}^{\top} \mathbf{m}_{j,:}^{(m)\top} \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} - 2 \left(a_{t,j}^{(m)} - \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t \right)^{\top} \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&\quad - \frac{1}{2} \left(\mathbf{m}_{j,:}^{(m)\top} \mathbf{m}_{j,:}^{(m)} - 2 \mathbf{m}_{j,:}^{(m)\top} \mathbf{m}_{\text{old}_{j,:}}^{(m)} \right) \sigma^{-2} \left. \right]
\end{aligned} \tag{5}$$

A.4 Parameter Distribution for \mathbf{W}

The parameter distribution $q_{\mathbf{W}}(\mathbf{W})$ can be found with

$$\begin{aligned}
\log q_{\mathbf{W}}(\mathbf{W}) &= \mathbb{E}_{-\mathbf{W}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \log \left(p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t|\boldsymbol{\theta}, s_t) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{W}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau}) (\right. \right. \\
&\quad \sum_{t=1}^T \sum_{m=1}^M \sum_{j=1}^{D_m} \log \mathcal{N} \left(a_{t,j}^{(m)} | \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t + \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi}, \text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) \tilde{\tau}_m^{-1} \right) \\
&\quad \left. \left. + \sum_{m=1}^M \sum_{j=1}^{D_m} \mathcal{N} \left(\mathbf{w}_{j,:}^{(m)\top} | \mathbf{0}, \bar{\boldsymbol{\alpha}}_{m,K} \right) \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\mathbf{W}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \left(\sum_{t=1}^T \sum_{m=1}^M \sum_{j=1}^{D_m} -\frac{1}{2} \text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T)^{-1} \tilde{\tau}_m \right. \right. \right. \\
&\quad \cdot \left(\left(a_{t,j}^{(m)} - \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right) - \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t \right)^{\top} \left(\left(a_{t,j}^{(m)} - \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right) - \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t \right) \\
&\quad \left. \left. \left. - \sum_{m=1}^M \sum_{j=1}^{D_m} \frac{1}{2} \mathbf{w}_{j,:}^{(m)\top} \bar{\boldsymbol{\alpha}}_{m,K} \mathbf{w}_{j,:}^{(m)} \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \right. \tag{6} \\
&= \mathbb{E}_{-\mathbf{W}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \left(\sum_{t=1}^T \sum_{m=1}^M \sum_{j=1}^{D_m} -\frac{1}{2} \text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T)^{-1} \tilde{\tau}_m \right. \right. \right. \\
&\quad \cdot \left(-2 \left(a_{t,j}^{(m)} - \mathbf{m}_{j,:}^{(m)} \boldsymbol{\Phi} \right) \mathbf{w}_{j,:}^{(m)} \tilde{\mathbf{z}}_t + \mathbf{w}_{j,:}^{(m)} + \tilde{\mathbf{z}}_t \tilde{\mathbf{z}}_t^{\top} \mathbf{w}_{j,:}^{(m)\top} \right) \\
&\quad \left. \left. \left. - \sum_{m=1}^M \sum_{j=1}^{D_m} \frac{1}{2} \mathbf{w}_{j,:}^{(m)\top} \bar{\boldsymbol{\alpha}}_{m,K} \mathbf{w}_{j,:}^{(m)} \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \right. \\
&= \mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \left(-2 \sum_{t=1}^T \frac{\left(a_{t,j}^{(m)} - \mathbb{E}_{\mathbf{m}} [\mathbf{m}_{j,:}^{(m)}] \boldsymbol{\Phi} \right) \mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t]^{\top} \mathbf{w}_{j,:}^{(m)\top}}{\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) \mathbb{E}_{\tilde{\tau}} [\tilde{\tau}_m^{-1}]} \right. \right. \\
&\quad \left. \left. \mathbf{w}_{j,:}^{(m)} \left(\sum_{t=1}^T \frac{\mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t \tilde{\mathbf{z}}_t^{\top}]}{\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) \mathbb{E}_{\tilde{\tau}} [\tilde{\tau}_m^{-1}]} + \bar{\boldsymbol{\alpha}}_{m,K} \right) \mathbf{w}_{j,:}^{(m)\top} \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1}
\end{aligned}$$

A.5 Distribution of Z

The distribution $q_{\tilde{\mathbf{Z}}}(\tilde{\mathbf{Z}})$ can be derived with

$$\begin{aligned}
\log q_{\tilde{\mathbf{Z}}}(\tilde{\mathbf{Z}}) &= \mathbb{E}_{-\tilde{\mathbf{z}}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \log \left(p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t | \boldsymbol{\theta}, s_t) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\tilde{\mathbf{z}}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(\sum_{t=1}^T \sum_{m=1}^M \log \mathcal{N}(a_t^{(m)} | \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t + \mathbf{M}^{(m)} \boldsymbol{\Phi}, \text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T) \tilde{\tau}_m^{-1} \mathbf{I}) \right. \right. \right. \\
&\quad \left. \left. \left. + \sum_{t=1}^T \log \mathcal{N}(\tilde{\mathbf{z}}_t | \mathbf{0}, \text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T) \mathbf{I}) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\tilde{\mathbf{z}}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(-\frac{1}{2} \sum_{t=1}^T \sum_{m=1}^M \text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T)^{-1} \tilde{\tau}_m \left(-2(a_t^{(m)} - \mathbf{M}^{(m)} \boldsymbol{\Phi})^T \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t \right. \right. \right. \right. \\
&\quad \left. \left. \left. \left. + \tilde{\mathbf{z}}_t^T \mathbf{W}^{(m)T} \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t \right) - \frac{1}{2} \text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T)^{-1} \tilde{\mathbf{z}}_t^T \tilde{\mathbf{z}}_t \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1 | \boldsymbol{\tau}) \left(-\frac{1}{2} \sum_{t=1}^T \left(-2 \left(\sum_{m=1}^M \frac{(a_t^{(m)} - \mathbb{E}_{\mathbf{M}}[\mathbf{M}^{(m)}] \boldsymbol{\Phi})^T \mathbf{W}^{(m)}}{\text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T) \mathbb{E}_{\tilde{\tau}}[\tilde{\tau}_m]^{-1}} \right) \tilde{\mathbf{z}}_t \right. \right. \right. \\
&\quad \left. \left. \left. + \tilde{\mathbf{z}}_t^T \left(\sum_{m=1}^M \frac{\mathbb{E}_{\mathbf{W}}[\mathbf{W}^{(m)T} \mathbf{W}^{(m)}]}{\text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T) \mathbb{E}_{\tilde{\tau}}[\tilde{\tau}_m]^{-1}} \right) \tilde{\mathbf{z}}_t + \frac{\tilde{\mathbf{z}}_t^T \tilde{\mathbf{z}}_t}{\text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T)} \right) \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1 | \boldsymbol{\tau})])^{-1}
\end{aligned} \tag{7}$$

A.6 Distribution of Tau

The distribution $q_{\tilde{\tau}}(\tilde{\tau})$ of the precision can be found with

$$\begin{aligned}
& \log q_{\tilde{\tau}}(\tilde{\tau}) \\
&= \mathbb{E}_{-\tilde{\tau}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \log \left(p(\boldsymbol{\theta}) \prod_{t=1}^T \pi(a_t|\boldsymbol{\theta}, s_t) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\tilde{\tau}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \log \prod_{t=1}^T \pi(a_t|\boldsymbol{\theta}, s_t) + \log \mathcal{G}(\tilde{\tau}|a^{\tilde{\tau}}, b^{\tilde{\tau}}) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\tilde{\tau}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \left(\sum_{t=1}^T \log \pi(a_t|\boldsymbol{\theta}, s_t) + \sum_{m=1}^M \log \mathcal{G}(\tilde{\tau}_m|a^{\tilde{\tau}}, b^{\tilde{\tau}}) \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \mathbb{E}_{-\tilde{\tau}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \left(\left(\sum_{m=1}^M \sum_{t=1}^T \log \left(|\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) \tilde{\tau}_m^{-1} \mathbf{I}|^{-\frac{1}{2}} \right) \right. \right. \right. \right. \right. \\
&\quad - \frac{1}{2} \left(\left(a_t^{(m)} - \mathbf{M}^{(m)} \boldsymbol{\Phi} \right) - \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t \right)^T \left(\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) \tilde{\tau}_m^{-1} \mathbf{I} \right)^{-1} \left(a_t^{(m)} - \mathbf{M}^{(m)} \boldsymbol{\Phi} \right) - \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t \right) \left. \left. \left. \left. \left. \right) \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \sum_{m=1}^M \frac{1}{2} D_m T \log \tilde{\tau}_m + \sum_{m=1}^M \left(\log(\tilde{\tau}_m) (a^{\tilde{\tau}} - 1) - b^{\tilde{\tau}} \tilde{\tau}_m \right) + \\
&\quad \mathbb{E}_{-\tilde{\tau}} \left[\mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \sum_{m=1}^M \sum_{t=1}^T \frac{\tilde{\tau}_m}{\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T)} \left(a_t^{(m)T} a_t^{(m)} - 2a_t^{(m)T} \mathbf{M}^{(m)} \boldsymbol{\Phi} \right. \right. \right. \\
&\quad + 2\tilde{\mathbf{z}}_t^T \mathbf{W}^{(m)T} \mathbf{M}^{(m)} \boldsymbol{\Phi} - 2a_t^{(m)T} \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t + \boldsymbol{\Phi}^T \mathbf{M}^{(m)T} \mathbf{M}^{(m)} \boldsymbol{\Phi} \\
&\quad \left. \left. \left. + \tilde{\mathbf{z}}_t^T \mathbf{W}^{(m)T} \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t \right) \right] \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1} \\
&= \sum_{m=1}^M \frac{1}{2} D_m T \log \tilde{\tau}_m + \sum_{m=1}^M \left(\log(\tilde{\tau}_m) (a^{\tilde{\tau}} - 1) - b^{\tilde{\tau}} \tilde{\tau}_m \right) + \\
&\quad \mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \sum_{m=1}^M \sum_{t=1}^T \frac{\tilde{\tau}_m}{\text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T)} \left(a_t^{(m)T} a_t^{(m)} - 2a_t^{(m)T} \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)}] \boldsymbol{\Phi} \right. \right. \\
&\quad + 2\mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t]^T \mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)}]^T \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)}] \boldsymbol{\Phi} - 2a_t^{(m)T} \mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)}] \mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t] + \boldsymbol{\Phi}^T \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)T} \mathbf{M}^{(m)}] \boldsymbol{\Phi} \\
&\quad \left. \left. + \mathbb{E}_{\tilde{\mathbf{z}}, \mathbf{W}} [\tilde{\mathbf{z}}_t^T \mathbf{W}^{(m)T} \mathbf{W}^{(m)} \tilde{\mathbf{z}}_t] \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1}
\end{aligned} \tag{8}$$

If we apply now the quadratic expectation

$$\mathbb{E} [\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr} (\mathbf{A} \text{Cov} [\mathbf{x}]) + \mathbb{E} [\mathbf{x}]^T \mathbf{A} \mathbb{E} [\mathbf{x}], \tag{9}$$

then we get

$$\begin{aligned}
&= \sum_{m=1}^M \frac{1}{2} D_m T \log \tilde{\tau}_m + \sum_{m=1}^M \left(\log(\tilde{\tau}_m) (a^{\tilde{\tau}} - 1) - b^{\tilde{\tau}} \tilde{\tau}_m \right) - \\
&\quad \frac{1}{2} \sum_{m=1}^M \tilde{\tau}_m \mathbb{E}_{p(\boldsymbol{\tau})} \left[p(r = 1|\boldsymbol{\tau}) \sum_{t=1}^T \text{Tr}(\boldsymbol{\Phi}\boldsymbol{\Phi}^T)^{-1} \left(a_t^{(m)T} a_t^{(m)} - 2a_t^{(m)T} \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)}] \boldsymbol{\Phi} \right. \right. \\
&\quad + 2\mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t]^T \mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)}]^T \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)}] \boldsymbol{\Phi} - 2a_t^{(m)T} \mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)}] \mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t] + \boldsymbol{\Phi}^T \mathbb{E}_{\mathbf{M}} [\mathbf{M}^{(m)T} \mathbf{M}^{(m)}] \boldsymbol{\Phi} \\
&\quad \left. \left. + \text{Tr} (\mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)T} \mathbf{W}^{(m)}] \text{Cov}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t]) + \mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t]^T \mathbb{E}_{\mathbf{W}} [\mathbf{W}^{(m)T} \mathbf{W}^{(m)}] \mathbb{E}_{\tilde{\mathbf{z}}} [\tilde{\mathbf{z}}_t] \right) \right] (\mathbb{E}_{p(\boldsymbol{\tau})} [p(r = 1|\boldsymbol{\tau})])^{-1}
\end{aligned} \tag{10}$$