# Supplementary Material for 'Action-Priority Driven Policy Optimization for Flexible Job Shop Scheduling Problem'

Wenjun Zheng[1], Weilin Cai[2], Wei Qu[3], Jianfeng Mao[*]

[1,2,3]:{wenjunzheng,weilincai,weiqu}@link.cuhk.edu.cn

[*]: jfmao@cuhk.edu.cn; corresponding author

## I. IMPORTANT NOTATIONS

This paper introduces various notations for problem descriptions, property definitions, formal proofs, and related purposes. In this subsection, we highlight the key notations that are frequently used throughout the chapters:

| | | | |
|---|---|---|---|
| $\mathcal{J}$: | The set for all jobs. | $\mathcal{M}$: | The set for all machines. |
| $\mathbb{m}$: | The number of jobs. | $\mathbb{m}$: | The number of machines. |
| $J_j$: | The $j$-th job. | $M_m$: | The $m$-th machine. |
| $\mathbb{k}_j$: | # operations of job $j$. | $O_{j,i}$: | $i$-th operation of job $j$. |
| $\mathcal{M}_{j,i}$: | Machines that can process $O_{j,i}$. | $\tilde{p}_{m,j,i}$: | Processing time of $M_m$ on $O_{j,i}$. |
| $p_{m,j,i}$: | Remaining processing time of $M_m$ on $O_{j,i}$. | $p^*$: | Max processing time among all $\tilde{p}_{m,j,i}$. |
| $\hat{p}$: | Min remaining processing time in $\mathbf{W}_t$. | $C_{\max}$: | The makespan of a schedule. |
| $C_{\max}^*$: | Optimal makespan in our model. | $C_{\max}^{\mathrm{TI}^*}$: | Optimal makespan in time-indexed model. |
| $t$: | Stage index. | $\mathcal{H}_t$: | Set of jump stages before stage $t$. |
| $s_t$: | State at stage $t$. | $o_t$: | Observation at stage $t$. |
| $\pi_{\mathrm{PTB}}$: | Processing time-based policy. | $\mathcal{T}^{\mathrm{TI}}$: | Trajectory set in time-indexed system. |
| $\mathcal{T}^s$: | Trajectories attained by $\pi^s$. | $\mathcal{T}^o$: | Trajectories attained by $\pi^o$. |
| $|\tau|$: | # stages in trajectory $\tau$. | $\mathbf{J}_t$: | Job info matrix at stage $t$. |
| $\mathbf{M}_t$: | Machine info vector at stage $t$. | $\mathbf{W}_t$: | Working info set at stage $t$. |
| $\mathcal{S}$: | Set of all states. | $\mathcal{O}$: | Set of all observations. |
| $\mathcal{A}(s)$: | Feasible actions given state $s$. | $\mathcal{A}(o)$: | Feasible actions given observation $o$. |
| $\mathcal{A}^r(o)$: | Reduced action space given $o$. | $\Pi$: | Policy set over states. |
| $\hat{\Pi}$: | Policy set over observations. | $\Pi^f$: | Policy set with state aggregation. |
| $V^*(s)$: | Optimal value of $s$ under state policy. | $V^\dagger(s)$: | Optimal value of $s$ under observation policy. |
| $\mathcal{P}$: | A specific FJSP problem. | $\mathrm{Pr}$: | Probability that something happens. |
| $\gamma$: | Discount factor. | $R(s,a)$: | Reward given $s$ and $a$. |
| $\mathbb{R}$: | Set of real numbers. | $\pi \rightsquigarrow \tau$: | Policy $\pi$ follows trajectory $\tau$. |

Please note that these notations are explained in detail in their respective sections. Not all variables are included here, as **some intermediate variables are intentionally omitted**. This subsection serves as a centralized reference for quick access to commonly used notations.

## II. MODEL EXPLANATION

### A. Detailed Version of System Dynamics

In this subsection, we introduce the system dynamics in details. Since the FJSP problem is deterministic, the next state $s_{t+1} = (\mathbf{J}_{t+1}, \mathbf{M}_{t+1}, \mathbf{W}_{t+1}, \mathbb{T}_t)$ is uniquely determined given the current state $s_t = (\mathbf{J}_t, \mathbf{M}_t, \mathbf{W}_t)$ and action $a_t = (m, j, i)$. The state update follows three steps: Step (1) is an **Intermediate Information Storage**: In this step, intermediate job, machine, and working information are created to store the state after applying action $a$. Here Step (2) is for **Checking Available Actions**: Since, in our model, a stage represents a decision point when an available machine can be assigned to an operation, we need to check whether any available actions remain after executing $a$. If no actions are available, we proceed to Step (3), i.e. **Processing Progression**: The system advances the processing stage, during which some operations are completed, and some machines are released. After this update, we again check for available actions, leading back to Step (2) if any exist. The details for each step are shown below:

(1) First, obtain intermediate job information $\mathbf{J}$, machine information $\mathbf{M}$, working information $\mathbf{W}$ and current system time $\mathbb{T}$. The newly assigned machine $M_m$ for operation $O_{j,i}$ and its remaining processing time are added to the work information:

$$\mathbf{W} = \mathbf{W}_t \cup \{(m, j, i, (\mathcal{P})_{m,j,i})\}$$

The first column of $\mathbf{J}_t$ remains unchanged:

$$(\mathbf{J})_{*,1} = (\mathbf{J}_t)_{*,1}$$

The status of job $J_j$ is updated from idle to in progress:

$$(\mathbf{J})_{j',2} = \begin{cases} 1, & \text{if } j' = j, \\ (\mathbf{J}_t)_{j',2}, & \text{otherwise.} \end{cases}$$

The status of machine $M_m$ is changed from idle to busy:

$$(\mathbf{M})_{m'} = \begin{cases} 1, & \text{if } m' = m, \\ (\mathbf{M}_t)_{m'}, & \text{otherwise.} \end{cases}$$

The system time remains the same: $\mathbb{T} = \mathbb{T}_t$.

(2) If based on $\mathbf{J}, \mathbf{M}$ and $\mathcal{P}$, there is no available action, then goes to Step (3), otherwise let $\mathbf{J}_{t+1} = \mathbf{J}, \mathbf{M}_{t+1} = \mathbf{M}, \mathbf{W}_{t+1} = \mathbf{W}, \mathbb{T}_{t+1} = \mathbb{T}$.

(3) Since there is no action available, the system should precede the processing progress. Let $\hat{p}$ represent the minimum remaining processing time, defined as:

$$\hat{p} = \min \{ p_{m,j,i} \mid (m, j, i, p_{m,j,i}) \in \mathbf{W} \}.$$

Then the system time should updated as $\hat{\mathbb{T}} = \mathbb{T} + \hat{p}$. The set $\hat{\mathbf{W}}$ is then updated as follows:

$$\hat{\mathbf{W}} = \{ (m, j, i, p_{m,j,i} - \hat{p}_t) \mid p_{m,j,i} > \hat{p}_t \}.$$

This update process removes machine-operation pairs where the remaining processing time is exactly $\hat{p}$, while reducing the remaining time for all other pairs accordingly. We define the set of machines that are released as:

$$\mathcal{M}^{\text{released}} = \{ m \mid (m, j, i, p_{m,j,i}) \in \mathbf{W}, p_{m,j,i} = \hat{p} \}.$$

Similarly, the set of jobs that need to be updated is given by:

$$\mathcal{J}^{\text{released}} = \{ j \mid (m, j, i, p_{m,j,i}) \in \mathbf{W}, p_{m,j,i} = \hat{p} \}.$$

Additionally, let

$$\mathcal{J}^{\text{done}} = \{ j \mid (\mathbf{J}_t)_{j,1} = \Bbbk_j \}$$

denote the set of jobs that have already completed their final operation. The job status is updated as follows for the first column:

$$(\hat{\mathbf{J}})_{j,1} = \begin{cases} 0, & \text{if } j \in \mathcal{J}^{\text{released}} \cap \mathcal{J}^{\text{done}}, \\ (\mathbf{J})_{j,1} + 1, & \text{if } j \in \mathcal{J}^{\text{released}} \setminus \mathcal{J}^{\text{done}}, \\ (\mathbf{J})_{j,1}, & \text{if } j \notin \mathcal{J}^{\text{released}}. \end{cases}$$

- If $j \in \mathcal{J}^{\text{released}} \cap \mathcal{J}^{\text{done}}$, the job has completed its final operation.
- If $j \in \mathcal{J}^{\text{released}} \setminus \mathcal{J}^{\text{done}}$, the job has finished its current operation but still has remaining steps.

For the second column, the update rule is:

$$(\hat{\mathbf{J}})_{j,2} = \begin{cases} 0, & \text{if } j \in \mathcal{J}^{\text{released}}, \\ (\mathbf{J})_{j,2}, & \text{if } j \notin \mathcal{J}^{\text{released}}. \end{cases}$$

The machine status is updated as follows:

$$(\mathbf{M}_{t+1})_m = \begin{cases} 0, & \text{if } m \in \mathcal{M}^{\text{released}}, \\ (\mathbf{M})_m, & \text{if } m \notin \mathcal{M}^{\text{released}}. \end{cases}$$

This ensures that machines in $\mathcal{M}^{\text{released}}$ are marked as available, while others retain their previous status. Let $\mathbf{J} = \hat{\mathbf{J}}, \mathbf{M} = \hat{\mathbf{M}}, \mathbf{W} = \hat{\mathbf{W}}, \mathbb{T} = \hat{\mathbb{T}}$, then goes to Step (2).

### B. Model Illustration

## III. PROOFS IN DETAILS

### A. Proof of Lemma 1 (with the original lemma provided)

**Lemma 1.** *Given a policy $\pi$ and its corresponding action-priority vector $\Phi$, suppose $G(\tau') > V^\pi(s_0')$, then for the new policy $\pi'$ and its corresponding $\Phi'$ constructed according to (8), we have following relations:*

- $\pi(a|s) = \pi'(a|s), \forall a \neq a_0', a \in \mathcal{A}(s)$, *if $a_0'$ is invalid for state $s$.*
- $\pi'(a|s) < \pi(a|s), \forall a \neq a_0', a \in \mathcal{A}(s)$, *if $a_0'$ is valid for state $s$. More specifically, if $s$ is the state that satisfies $\mathcal{A}(s) \subseteq \mathcal{A}(s_0')$, then $\pi'(a|s) \leq \delta\pi(a|s)$ with equality holds when $\mathcal{A}(s) = \mathcal{A}(s_0')$.*

*For the case $G(\tau') < V^\pi(s_0')$, the first result still holds, while the second result changes the inequality $\leq$ to $\geq$.*

*Proof.* Here, we provide the proof for the case $G(\tau') > V^\pi(s_0')$. The proof for the converse case is similar.

- If $\tilde{a}_0$ is invalid for state $s$, then $\mathbb{I}(s, \tilde{a}_0) = 0$.

$$\begin{aligned} \pi'(a|s) &= \frac{\phi_a'}{\sum_j \mathbb{I}(s,j)\phi_j'} = \frac{\phi_a'}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j'} \\ &= \frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\delta\phi_j} = \frac{\delta\phi_a}{\delta \sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j} \end{aligned}$$
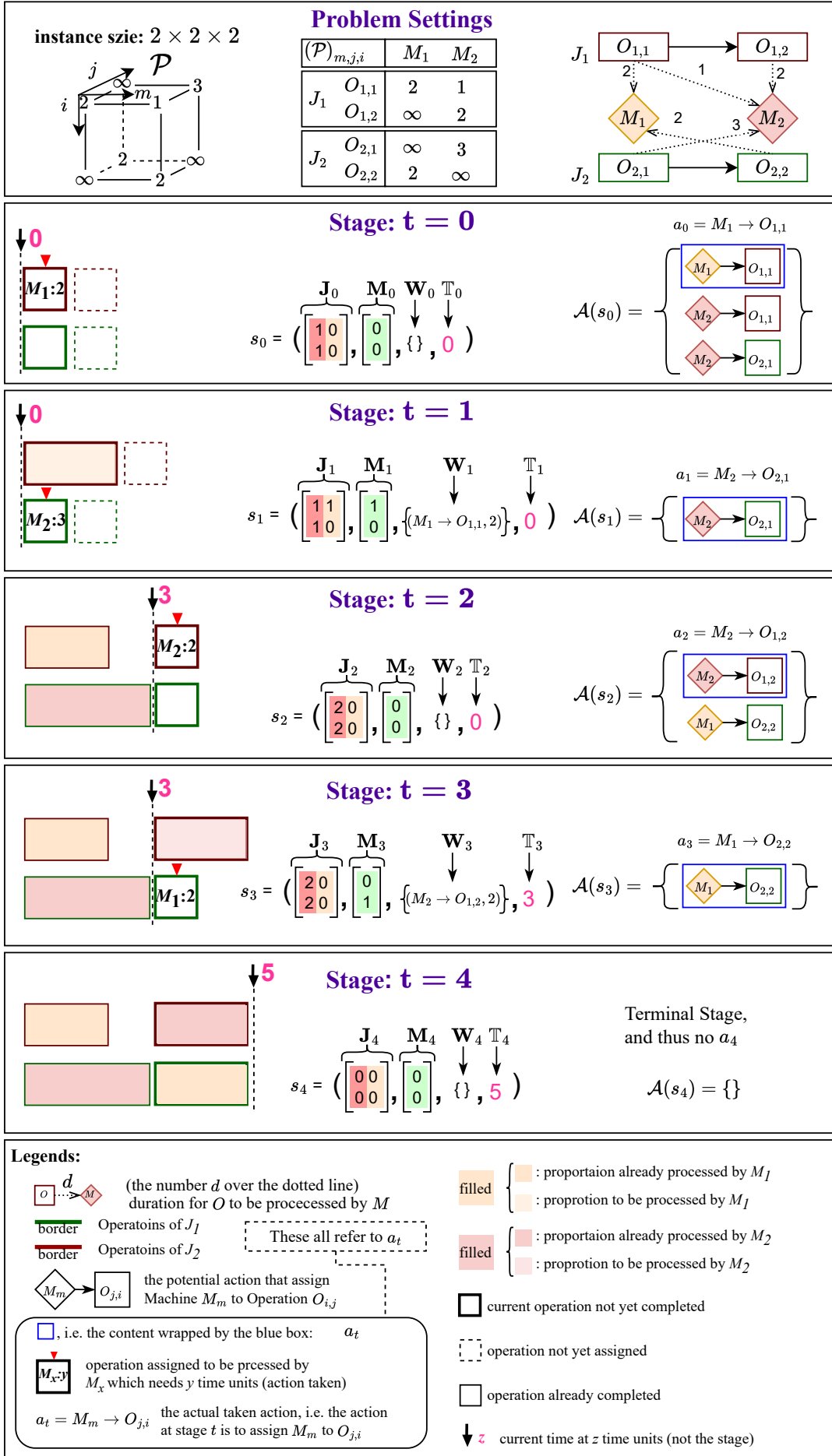
# Problem Settings

**instance szie:** $2 \times 2 \times 2$

$\mathcal{P}$

| $(\mathcal{P})_{m,j,i}$ | | $M_1$ | $M_2$ |
|---|---|---|---|
| $J_1$ | $O_{1,1}$ | 2 | 1 |
| | $O_{1,2}$ | $\infty$ | 2 |
| $J_2$ | $O_{2,1}$ | $\infty$ | 3 |
| | $O_{2,2}$ | 2 | $\infty$ |

$J_1$: $O_{1,1} \to O_{1,2}$

$J_2$: $O_{2,1} \to O_{2,2}$

$M_1$ $M_2$

## Stage: $t = 0$

$M_1{:}2$

$$s_0 = \left( \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix}}_{\mathbf{J}_0}, \underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\mathbf{M}_0}, \underbrace{\{\}}_{\mathbf{W}_0}, \underbrace{0}_{\mathbb{T}_0} \right)$$

$a_0 = M_1 \to O_{1,1}$

$$\mathcal{A}(s_0) = \left\{ \begin{array}{c} M_1 \to O_{1,1} \\ M_2 \to O_{1,1} \\ M_2 \to O_{2,1} \end{array} \right\}$$

## Stage: $t = 1$

$M_2{:}3$

$$s_1 = \left( \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}}_{\mathbf{J}_1}, \underbrace{\begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\mathbf{M}_1}, \underbrace{\{(M_1 \to O_{1,1}, 2)\}}_{\mathbf{W}_1}, \underbrace{0}_{\mathbb{T}_1} \right)$$

$a_1 = M_2 \to O_{2,1}$

$$\mathcal{A}(s_1) = \left\{ M_2 \to O_{2,1} \right\}$$

## Stage: $t = 2$

$M_2{:}2$

$$s_2 = \left( \underbrace{\begin{bmatrix} 2 & 0 \\ 2 & 0 \end{bmatrix}}_{\mathbf{J}_2}, \underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\mathbf{M}_2}, \underbrace{\{\}}_{\mathbf{W}_2}, \underbrace{0}_{\mathbb{T}_2} \right)$$

$a_2 = M_2 \to O_{1,2}$

$$\mathcal{A}(s_2) = \left\{ \begin{array}{c} M_2 \to O_{1,2} \\ M_1 \to O_{2,2} \end{array} \right\}$$

## Stage: $t = 3$

$M_1{:}2$

$$s_3 = \left( \underbrace{\begin{bmatrix} 2 & 0 \\ 2 & 0 \end{bmatrix}}_{\mathbf{J}_3}, \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_{\mathbf{M}_3}, \underbrace{\{(M_2 \to O_{1,2}, 2)\}}_{\mathbf{W}_3}, \underbrace{3}_{\mathbb{T}_3} \right)$$

$a_3 = M_1 \to O_{2,2}$

$$\mathcal{A}(s_3) = \left\{ M_1 \to O_{2,2} \right\}$$

## Stage: $t = 4$

$$s_4 = \left( \underbrace{\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{J}_4}, \underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\mathbf{M}_4}, \underbrace{\{\}}_{\mathbf{W}_4}, \underbrace{5}_{\mathbb{T}_4} \right)$$

Terminal Stage, and thus no $a_4$

$$\mathcal{A}(s_4) = \{\}$$

**Legends:**

$O \xrightarrow{d} M$ (the number $d$ over the dotted line) duration for $O$ to be proccessed by $M$

green border: Operatoins of $J_1$

red border: Operatoins of $J_2$

These all refer to $a_t$

$M_m \to O_{j,i}$ the potential action that assign Machine $M_m$ to Operation $O_{i,j}$

blue box, i.e. the content wrapped by the blue box: $a_t$

$M_x{:}y$ operation assigned to be prcessed by $M_x$ which needs $y$ time units (action taken)

$a_t = M_m \to O_{j,i}$ the actual taken action, i.e. the action at stage $t$ is to assign $M_m$ to $O_{j,i}$

filled: : proportaion already processed by $M_1$ / : proprotion to be processed by $M_1$

filled: : proportaion already processed by $M_2$ / : proprotion to be processed by $M_2$

□ current operation not yet completed

⊡ operation not yet assigned

□ operation already completed

↓ $z$ current time at $z$ time units (not the stage)

Fig. 1. Model Illustration

$$= \frac{\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j} = \frac{\phi_a}{\sum_j \mathbb{I}(s,j)\phi_j}$$
$$= \pi(a|s)$$

- If $\tilde{a}_0$ is valid for state $s$, then $\mathbb{I}(s, \tilde{a}_0) = 1$.

$$\pi'(a|s) = \frac{\phi'_a}{\sum_j \mathbb{I}(s,j)\phi'_j} = \frac{\phi'_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi'_j + \phi'_{\tilde{a}_0}}$$
$$= \frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\delta\phi_j + \phi_{\tilde{a}_0} + (1-\delta)\sum_{j \neq \tilde{a}_0} \mathbb{I}(\tilde{s}_0,j)\phi_j} < \frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\delta\phi_j + \delta\phi_{\tilde{a}_0}}$$
$$= \frac{\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j + \phi_{\tilde{a}_0}} = \frac{\phi_a}{\sum_j \mathbb{I}(s,j)\phi_j}$$
$$= \pi(a|s)$$

Now, if $\mathcal{A}(s) \subseteq \mathcal{A}(\tilde{s}_0)$, we have $\sum_{j \neq \tilde{a}_0} \mathbb{I}(\tilde{s}_0,j)\phi_j \geq \sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j$ since $\phi_i > 0$. Therefore,

$$\frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\delta\phi_j + \phi_{\tilde{a}_0} + (1-\delta)\sum_{j \neq \tilde{a}_0} \mathbb{I}(\tilde{s}_0,j)\phi_j} \leq \frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\delta\phi_j + \phi_{\tilde{a}_0} + (1-\delta)\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j}$$
$$= \frac{\delta\phi_a}{\sum_{j \neq \tilde{a}_0} \mathbb{I}(s,j)\phi_j + \phi_{\tilde{a}_0}} = \delta\pi(a|s)$$

That is $\pi'(a|s) \leq \pi(a|s)$. ∎

*B. Proof of Lemma 2 (with the original lemma provided)*

**Lemma 2.** *Suppose $\mathcal{A}_m(s'_0)$ satisfying **CD3** and $\tau \in \mathcal{T}$, then $t_\tau^+ \geq 2$.*

*Proof.* Since $\tau \in \mathcal{T}$, by the definition of $\mathcal{T}$, we have $\tau = (s'_0, a_0, \cdots)$, where either $a_0 \in \mathcal{A}_m(s'_0) \cup \mathcal{A}_{(j,i)}(s'_0)$ or $a_0 \notin \mathcal{A}_m(s'_0) \cup \mathcal{A}_{(j,i)}(s'_0)$. If $a_0 \in \mathcal{A}_m(s'_0) \cup \mathcal{A}_{(j,i)}(s'_0)$, then by **CD3**, we directly obtain $t_\tau^+ \geq 2$. Otherwise, if $a_0 \notin \mathcal{A}_m(s'_0)$, suppose $a_0 = (m,j,i)$ and $a'_0 = (m',j',i')$. Since $m \neq m'$, and $j \neq j'$, $i \neq i'$, it follows that after executing action $a_0$, action $a'_0$ is still available. Hence, we conclude that $t_\tau^+ \geq 2$. ∎

*C. Proof of Lemma 3 (with the original lemma provided)*

**Lemma 3.** *Suppose $\tau \in \mathcal{T}_1$, then $a'_0 \notin \mathcal{A}_{t_\tau^+}(\tau)$. If $\pi'$ is updated according to updating rule (8), then*

$$P_\tau^{\pi'} < \frac{\delta^2(1-\delta\eta)\sum_j \mathbb{I}_j(s)\phi_j}{(1-\eta)(\delta\sum_j \mathbb{I}_j(s)\phi_j + (1-\delta)\sum_j \mathbb{I}(s'_0,j)\phi_j)} P_\tau^\pi, s \in \tau$$

*where $P_\tau^{\pi'}$ represents the probability of $\tau$ given policy $\pi'$. $\pi'$ is the new policy and $\pi$ is the old one.*

*Proof.* Suppose $\tau = (s'_0, a_0, \ldots, a_{t_\tau^+ - 1}, s_{t_\tau^+}, \ldots)$. By the definition of $\mathcal{T}_1$, we have $a_t \neq a'_0, \forall 0 \leq t \leq t_\tau^+ - 1$. Then, by Lemma 1, it follows that $\pi'(a_0|s'_0) = \delta\pi(a_0|s'_0)$ and $\pi'(a_t|s_t) < \delta\pi(a_t|s_t)$ for $1 \leq t \leq t_\tau^+ - 1$.

In the worst case, if action $a'_0$ is selected at some state $s$ in $\tau$, then

$$\pi(a'_0|s) = \frac{\phi_{a'_0}}{\sum_{j \neq a'_0} \mathbb{I}_j(s)\phi_j + \phi_{a'_0}}$$

and

$$\pi'(a'_0|s) = \frac{\phi_{a'_0} + (1-\delta)\sum_{j \neq a'_0} \mathbb{I}_j(s'_0)\phi_j}{\delta\sum_{j \neq a'_0} \mathbb{I}_j(s)\phi_j + \phi_{a'_0} + (1-\delta)\sum_{j \neq a'_0} \mathbb{I}_j(s'_0)\phi_j}.$$

Again, by Lemma 1, once action $a'_0$ is selected, we have $\pi'(a|s) = \pi(a|s)$. Thus, for some $s \in \tau$, it follows that

$$p_\tau^{\pi'} < \frac{\delta^{t_\tau^+}[\phi_{a'_0} + (1-\delta)\sum_{j \neq a'_0} \mathbb{I}_j(s'_0)\phi_j]}{\delta\sum_{j \neq a'_0} \mathbb{I}_j(s)\phi_j + \phi_{a'_0} + (1-\delta)\sum_{j \neq a'_0} \mathbb{I}_j(s'_0)\phi_j} \cdot \frac{\sum_{j \neq a'_0} \mathbb{I}_j(s)\phi_j + \phi_{a'_0}}{\phi_{a'_0}} p_\tau^\pi, \quad \text{for some } s \in \tau$$
$$= \frac{\delta^{t_\tau^+}(1-\delta\eta)\sum_j \mathbb{I}_j(s)\phi_j}{(1-\eta)(\delta\sum_j \mathbb{I}_j(s)\phi_j + (1-\delta)\sum_j \mathbb{I}_j(s'_0)\phi_j)} p_\tau^\pi.$$

Since, by Lemma 2, we have $t_\tau^+ \geq 2$, it follows that

$$p_\tau^{\pi'} < \frac{\delta^2(1-\delta\eta)\sum_j \mathbb{I}_j(s)\phi_j}{(1-\eta)(\delta\sum_j \mathbb{I}_j(s)\phi_j + (1-\delta)\sum_j \mathbb{I}_j(s'_0)\phi_j)} p_\tau^\pi.$$
∎

*D. Proof of Lemma 4 (with the original lemma provided)*

**Lemma 4.** *Suppose $\tau \in \mathcal{T}$, $a'_0 \in \mathcal{A}_{t^+_\tau}(\tau)$ and $a'_0$ is executed at stage $t$, we know $0 \leq t < t^+_\tau$. Then,*

$$\mathrm{P}^{\pi'}_\tau < \frac{\delta^t(1 - \delta\eta)}{1 - \eta}\mathrm{P}^\pi_\tau$$

*Proof.* By definition, the set $\mathcal{A}_{t^+_\tau}(\tau)$ consists of the actions taken in trajectory $\tau$ before stage $t^+_\tau$. Since $a'_0 \in \mathcal{A}_{t^+_\tau}(\tau)$, we can express $\tau$ as

$$\tau = (s'_0, a_0, \ldots, s_t, a'_0, \ldots, a_{t^+_\tau - 1}, s_{t^+_\tau}, \ldots).$$

Then, by Lemma 1, it follows that $\pi'(a_0|s'_0) = \delta\pi(a_0|s'_0)$ and $\pi'(a_{t'}|s_{t'}) < \delta\pi(a_{t'}|s_{t'})$ for $1 \leq t' \leq t - 1$. At state $s_t$, we have

$$\pi(a'_0|s_t) = \frac{\phi_{a'_0}}{\sum_j \mathbb{I}_j(s_t)\phi_j}$$

and

$$\pi'(a'_0|s_t) = \frac{\phi_{a'_0} + (1 - \delta)\sum_{j \neq a'_0} \mathbb{I}_j(s'_0)\phi_j}{\delta\sum_{j \neq a'_0}\mathbb{I}_j(s_t)\phi_j + \phi_{a'_0} + (1 - \delta)\sum_{j \neq a'_0}\mathbb{I}_j(s'_0)\phi_j}.$$

From Lemma 1, once action $a'_0$ is taken, we have $\pi'(a|s_{t'}) = \pi(a|s_{t'})$ for $t' > t$. Thus, we obtain

$$p^{\pi'}_\tau < \frac{\delta^t(\phi_{a'_0} + (1 - \delta)\sum_{j \neq a'_0}\mathbb{I}_j(s'_0)\phi_j)}{\delta\sum_{j \neq a'_0}\mathbb{I}_j(s_t)\phi_j + \phi_{a'_0} + (1 - \delta)\sum_{j \neq a'_0}\mathbb{I}_j(s'_0)\phi_j} \cdot \frac{\sum_j \mathbb{I}_j(s_t)\phi_j}{\phi_{a'_0}}p^\pi_\tau$$

$$< \frac{\delta^t(\delta\phi_{a'_0} + (1 - \delta)\sum_j \mathbb{I}_j(s'_0)\phi_j)}{\phi_{a'_0}}p^\pi_\tau$$

$$= \frac{\delta^t(1 - \delta\eta)}{1 - \eta}p^\pi_\tau.$$

Such relaxation exactly match the required inequality in Lemma 4. ∎

*E. Proof of Lemma 5 (with the original lemma provided)*

**Lemma 5.** *Suppose that $\tilde{\tau} := (s'_0, a'_0, \tilde{s}_1, \tilde{a}_1, \ldots, \tilde{a}_{t^+_{\tilde{\tau}} - 1}, \tilde{s}_{t^+_{\tilde{\tau}}}, \ldots) \in \mathcal{T}_3$, and the set $\{a'_0, \tilde{a}_1, \ldots, a_{t^+_{\tilde{\tau}} - 1}\}$ contains all the actions taken before the first change in system time within trajectory $\tilde{\tau}$. Next, define $\tilde{\tau}^{(t)}$ as the trajectory whose first $t^+_{\tilde{\tau}}$ actions match the set $\{a'_0, \tilde{a}_1, \ldots, a_{t^+_{\tilde{\tau}} - 1}\}$, and satisfy $\tilde{\tau}^{(t)}|\tilde{s}_{t^+_{\tilde{\tau}}} = \tilde{\tau}|\tilde{s}_{t^+_{\tilde{\tau}}}$. By the existence of the **Permutation Irrelevance Property** (Proposition 1), the equality $\tilde{\tau}^{(t)}|\tilde{s}_{t^+_{\tilde{\tau}}} = \tilde{\tau}|\tilde{s}_{t^+_{\tilde{\tau}}}$ is well-defined.*

*The order of the first $t^+_{\tilde{\tau}}$ actions in $\tilde{\tau}^{(t)}$ is given by*

$$a^{(t)}_{t'} = \begin{cases} a_1, & \text{if } t' = 0, \\ a'_0, & \text{if } t' = t, \\ a_{t'+1}, & \text{if } 0 < t' < t, \\ a_{t'}, & \text{if } t' > t, \end{cases} \tag{1}$$

*where $0 \leq t' \leq t^+_{\tilde{\tau}} - 1$.*

*Then, an equivalent expression for the set $\mathcal{T}_2$ is given by*

$$\mathcal{T}_2 = \bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t^+_{\tilde{\tau}} - 1)}\}.$$

*Similarly, the set $\mathcal{T}_4$ can be expressed in the same form. In other words, $\mathcal{T}_4$ is a special case of this relation, since $\mathcal{T}_5$ contains only a single element, $\tau'$. Specifically, we have*

$$\mathcal{T}_4 = \{\tau'^{(0)}, \ldots, \tau'^{(t^+_{\tau'} - 1)}\}.$$

*where $\tau'^{(t)}$ for $0 \leq t \leq t^+_{\tau'} - 1$ is defined similarly, satisfying $\tau'^{(t)}|s'_{t^+_{\tau'}} = \tau'|s'_{t^+_{\tau'}}$.*

*Proof.* We first derive an equivalent expression for the set $\mathcal{T}_2$. Recall its definition:

$$\mathcal{T}_2 = \left\{\tau \mid \exists\tilde{\tau} \in \mathcal{T}_3, \text{ s.t. } \mathcal{A}_{t^+_{\tilde{\tau}}}(\tau) = \mathcal{A}_{t^+_{\tilde{\tau}}}(\tilde{\tau}), \tau \neq \tilde{\tau}, \tau \notin \mathcal{T}_3, \tau \in \mathcal{T}\right\}.$$

For any $\tau \in \mathcal{T}_2$, there exists some $\tilde{\tau} \in \mathcal{T}_3$ such that $\mathcal{A}_{t^+_{\tilde{\tau}}}(\tau) = \mathcal{A}_{t^+_{\tilde{\tau}}}(\tilde{\tau})$. This implies that the action set before the first system change is identical for both trajectories. Since $\tau \neq \tilde{\tau}$, the sequence of actions before stage $t^+_{\tilde{\tau}}$ must be a permutation of the set $\mathcal{A}_{t^+_{\tilde{\tau}}}$. Consequently, $\tau \in \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t^+_{\tilde{\tau}} - 1)}\}$. Thus, we conclude that

$$\mathcal{T}_2 \subseteq \bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t^+_{\tilde{\tau}} - 1)}\}.$$

Conversely, suppose $\tau \in \bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$. Then, there must exist some $\tilde{\tau} \in \mathcal{T}_3$ such that $\tau \in \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$. By the definition of the set $\{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$ (see equation (1)), we have

$$\mathcal{A}_{t_{\tilde{\tau}}^+}(\tau) = \mathcal{A}_{t_{\tilde{\tau}}^+}(\tilde{\tau}) \quad \text{and} \quad \tau \neq \tilde{\tau}.$$

This implies $\tau \in \mathcal{T}_2$. Thus, $\bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\} \subseteq \mathcal{T}_2$. Combining both directions, we can then conclude that $\mathcal{T}_2 = \bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$. The proof of the equivalence relation for the set $\mathcal{T}_4$ is similar and follows the same structure as the proof above. ∎

### F. Proof of Proposition 1 (with the original proposition provided)

**Proposition 1** (Permutation Irrelevant). *Suppose $(s_0, a_0, \cdots, s_t, a_t, s_{t+1})$ is a state-action sequence with $\mathbb{T}_0 = \mathbb{T}_t < \mathbb{T}_{t+1}$, meaning that the system times of states $s_0$ and $s_t$ are the same and both are less than that of $s_{t+1}$. Then, consider any permutation of the action set $\{a_0, \ldots, a_t\}$, which forms a new ordered action sequence. If we execute these actions in the new given order, the resulting state will be exactly $s_{t+1}$.*

*Proof.* It's equivalent to prove the existence of some function (or more generally, some mapping) $h$ that maps the input $(s_0, \{a_1, a_2, \cdots, a_t\})$ to $s_{t+1}$, i.e. we shall have $s_{t+1} = h(s_0, \{a_1, a_2, \cdots, a_t\})$.

As we write the action as $a = (m, j, i)$ and the calculation of $\mathbf{W}$ actually adds the term $(m, j, i, (\mathcal{P})_{m,j,i})$, we may rewrite such term as $(a, (\mathcal{P})_a)$. Denote the set $\underline{\mathbf{W}} \doteq \mathbf{W}_0 \cup \{(a_0, (\mathcal{P})_{a_0}), (a_1, (\mathcal{P})_{a_1}), \cdots, (a_t, (\mathcal{P})_{a_t})\}$, the processing time $\tilde{p} \doteq \min\{(\mathcal{P})_a | (a, (\mathcal{P})_a) \in \underline{\mathbf{W}}\}$, the action set $\tilde{\mathcal{A}} \doteq \{a \mid (\mathcal{P})_a = \tilde{p}\}$, and the sets $\overline{\mathbf{W}} \doteq \underline{\mathbf{W}} \setminus \{(a, (\mathcal{P})_a) | a \in \tilde{\mathcal{A}}\}$, $\tilde{\mathbf{W}} \doteq \overline{\mathbf{W}} \setminus \mathbf{W}_0$. We also define the following auxiliary sets:

- $\tilde{\mathcal{J}} \doteq \{j | (m, j, i) \in \tilde{\mathcal{A}}\}$
- $\hat{\mathcal{J}} \doteq \{j | (m, j, i, (\mathcal{P})_{m,j,i}) \in \overline{\mathbf{W}}\}$
- $\bar{\mathcal{J}} \doteq \{j \mid (\mathbf{J}_0)_{j,1} = \Bbbk_j\}$
- $\tilde{\mathcal{M}} \doteq \{m | (m, j, i, (\mathcal{P})_{m,j,i}) \in \tilde{\mathbf{W}}\}$
- $\hat{\mathcal{M}} \doteq \{m | (m, j, i) \in \tilde{\mathcal{A}}\} \cap \{m | (m, j, i, (\mathcal{P})_{m,j,i}) \in \mathbf{W}_0\}$

Now, we can calculate $s_{t+1}$ in as follows according to the system dynamics:

- $(\mathbf{J}_{t+1})_{j,1} = \begin{cases} 0, & \text{if } j \in \tilde{\mathcal{J}} \cap \bar{\mathcal{J}}; \\ (\mathbf{J}_0)_{j,1} + 1, & \text{if } j \in \tilde{\mathcal{J}} \setminus \bar{\mathcal{J}}; \\ (\mathbf{J}_0)_{j,1}, & \text{otherwise.} \end{cases}$

- $(\mathbf{J}_{t+1})_{j,1} = \begin{cases} 0, & \text{if } j \in \tilde{\mathcal{J}}; \\ 1, & \text{if } j \in \hat{\mathcal{J}}. \end{cases}$

- $(\mathbf{M}_{t+1})_m = \begin{cases} 0, & \text{if } m \in \hat{\mathcal{M}}; \\ 1, & \text{if } j \in \tilde{\mathcal{M}}. \end{cases}$

- $\mathbf{W}_{t+1} = \{(a, (\mathcal{P})_a - \tilde{p}) | (a, (\mathcal{P}_a \in \overline{\mathbf{W}})\}$;

- $\mathbb{T}_{t+1} = \mathbb{T}_0 + \tilde{p}$.

Here the second part of the inputs of $h$ is some set $\{a_0, a_1, \ldots, a_t\}$ rather than some sequence. Since a set is identical no matter how the elements are permuted, with $s_{t+1} = h(s_0, \{a_0, a_1, \ldots, a_t\})$, we can tell that the permutation irrelevant property exists.

A potential issue is that when the sequence $(a_0, a_1, \cdots, a_t)$ is reordered, the proof above requires that these actions are still feasible at each time step. This is true: for $t_x, t_y \in \{0, 1, \cdots, t\}$, suppose $a_{t_x} = (m_{t_x}, j_{t_x}, i_{t_x})$ and $a_{t_y} = (m_{t_y}, j_{t_y}, i_{t_y})$, as long as $t_x \neq t_y$, we shall have $m_{t_x} \neq m_{t_y}$ and $(j_{t_x} \neq j_{t_y})$. Since the assignments at each time step are mutually exclusive, the feasibility can be guaranteed even if the action sequence is permuted.

Further explanation is shown in Figure 2. ∎

### G. Proof of Proposition 2 (with the original proposition provided)

**Proposition 2.** *For any $\delta \in (0, 1)$, condition **CD1** is satisfied, i.e.,*

$$\frac{\delta^2(1 - \delta\eta)}{(1 - \eta)\left(\delta + (1 - \delta)\frac{\sum_j \mathbb{I}(s_0', j)\phi_j}{\sum_j \mathbb{I}_j(s)\phi_j}\right)} < \frac{1 - (1 - \delta\eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}{1 - (1 - \eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}, \quad \forall s \neq s_0'.$$

*Proof.* Denote $\sigma = \frac{\sum_j \mathbb{I}(s_0', j)\phi_j}{\sum_j \mathbb{I}_j(s)\phi_j} > 0$, then

$$\frac{\delta^2(1 - \delta\eta)}{(1 - \eta)(\delta + (1 - \delta)\frac{\sum_j \mathbb{I}(s_0', j)\phi_j}{\sum_j \mathbb{I}_j(s)\phi_j})} < \frac{1 - (1 - \delta\eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}{1 - (1 - \eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}$$

$$\Leftrightarrow \frac{\delta^2(1 - \delta\eta)}{(1 - \eta)((1 - \sigma)\delta + \sigma)} < \frac{1 - (1 - \delta\eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}{1 - (1 - \eta)\mathrm{P}^\pi_{\tau' | (s_0', a_0')}}$$
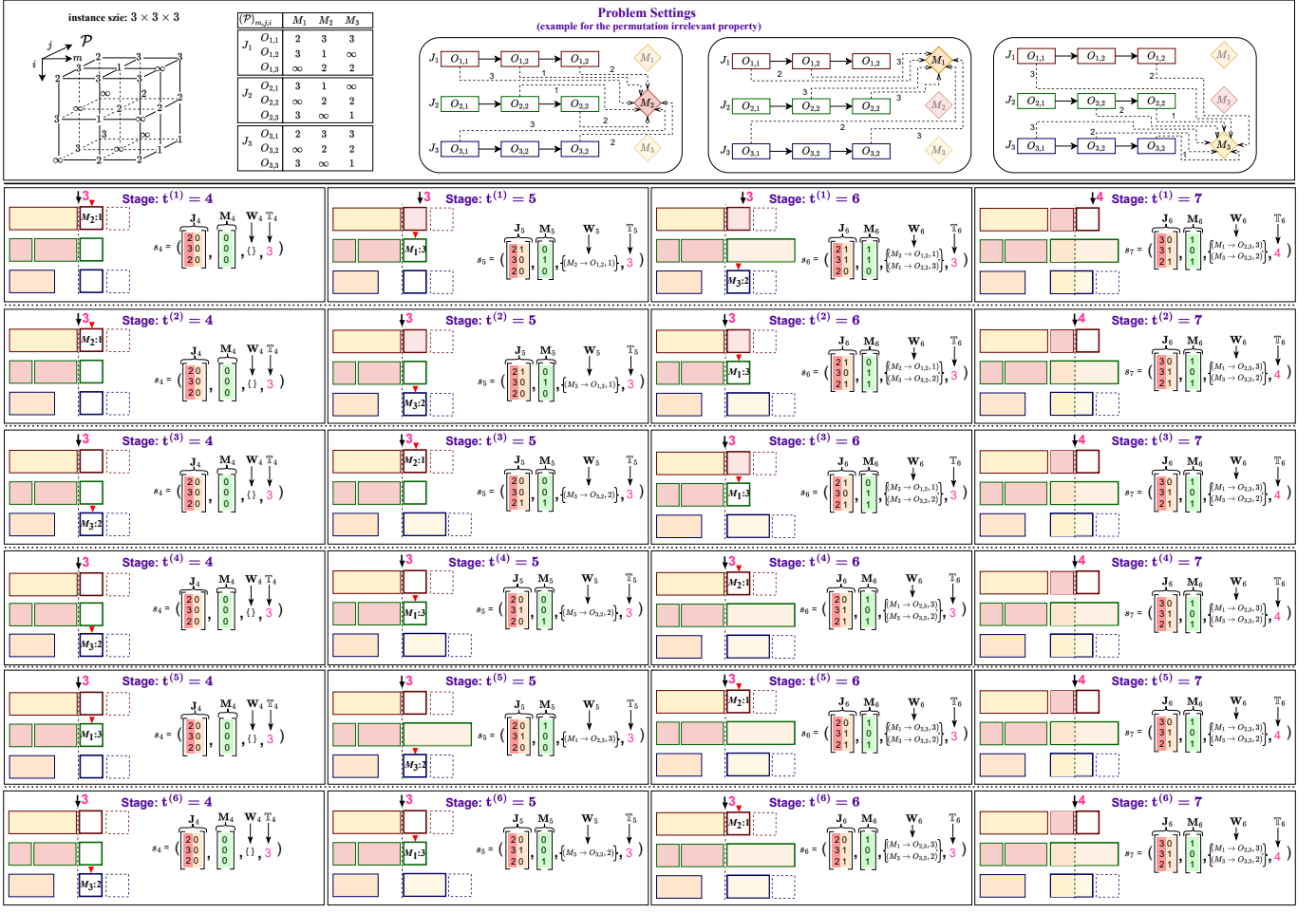
Fig. 2. Permutation Irrelevant Property Explanation

$$\Leftrightarrow \delta^2(1-\delta\eta)\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} < ((1-\sigma)\delta+\sigma)(1-(1-\delta\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')})$$

$$\Leftrightarrow \delta^2(1-\delta\eta)\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} < (1-\sigma)\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}\delta^2 + [(1-\sigma)(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')})+\sigma\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}]\delta + \sigma(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')})$$

$$\Leftrightarrow \left[\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} - (1-\sigma)\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}\right]\delta^2 - [(1-\sigma)(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')})+\sigma\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}]\delta - \sigma(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')}) <$$

$$\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta}\eta\delta^3$$

Let $f_1(\delta) = \left[\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} - (1-\sigma)\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}\right]\delta^2 - [(1-\sigma)(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')})+\sigma\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}]\delta - \sigma(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')})$

and $f_2(\delta) = \frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta}\eta\delta^3$. We find that $f_1(1) = f_2(1)$ and $f_1(0) = -\sigma(1-\mathrm{P}^\pi_{\tau'|(s_0',a_0')}) < 0$. Since

$$\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} - (1-\sigma)\eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')} < \frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta} - \eta\mathrm{P}^\pi_{\tau'|(s_0',a_0')}$$

$$= \frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')} - \eta(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta}$$

$$= \frac{1-(1-\eta^2)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta}$$

$$> \frac{1-(1-\eta^2)}{1-\eta} = \frac{\eta^2}{1-\eta} > 0$$

Hence, $f_1(\delta)$ is a convex quadratic function. Since

$$\frac{1-(1-\eta)\mathrm{P}^\pi_{\tau'|(s_0',a_0')}}{1-\eta}\eta > 0,$$

it follows that, as shown in Figure 3, we have $f_1(\delta) < f_2(\delta), \quad \forall \delta \in (0,1)$. ∎
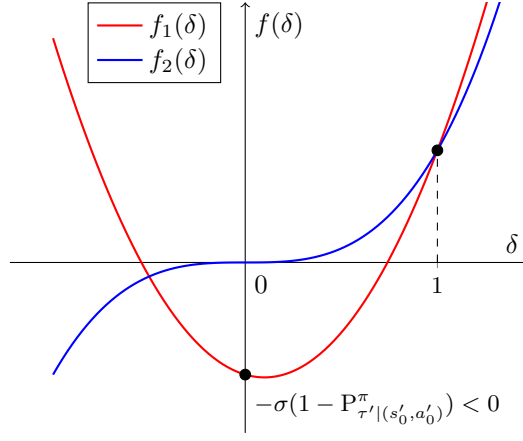
Fig. 3. $f_1(\delta), f_2(\delta)$

## H. Proof of Proposition 3 (with the original proposition provided)

**Proposition 3.** *If* $1 - \frac{\eta}{1-\eta} > \alpha + \beta$, *where*

$$\alpha = \frac{\eta \mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}}{1 - (1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}}, \quad \beta = \frac{1}{1-\eta}\frac{\eta}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}}\frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_1',j)\phi_j},$$

*then there exist a $\epsilon > 0$, such that when $\delta \in (\epsilon, 1)$, **CD2** is satisfied.*

**Remark 1.** *By the definition of*

$$\eta = \frac{\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_0',j)\phi_j},$$

*as well as the trajectory $\tau'$, states $s_0'$ and $s_1'$, the values of $\alpha$ and $\beta$ can be easily computed. Thus, the condition*

$$1 - \frac{\eta}{1-\eta} > \alpha + \beta$$

*can be efficiently verified.*

*Proof.* Given $\tau \in \mathcal{T}_3, \mathcal{T}_5$, let $g_1(\delta) = \sum_{t=1}^{t_\tau^+ - 1} \delta^t \frac{1 - \delta\eta}{1-\eta} \mathrm{P}^{\pi}_{\tau^{(t)}}$ and

$$g_2(\delta) = \delta \left( \sum_{t=1}^{t_\tau^+ - 1} \frac{\eta \mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}\mathrm{P}^{\pi}_{\tau^{(t)}}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}} + \frac{\eta \mathrm{P}^{\pi}_{\tau|(s_0',a_0')}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}} \right) \sum_{t=1}^{t_\tau^+ - 1} \frac{(1 - \mathrm{P}^{\pi}_{\tau'|(s_0',a_0')})\mathrm{P}^{\pi}_{\tau^{(t)}}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}} - \frac{\eta \mathrm{P}^{\pi}_{\tau|(s_0',a_0')}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}}$$

The following theorem is established in Spivak, M.'s book [1], and we will use it directly:

**Theorem** (for Proposition 3). *Let $g_1(\delta)$ and $g_2(\delta)$ be continuously differentiable functions on $(0,1]$. If*

$$g_1(1) = g_2(1) \quad and \quad g_1'(1) > g_2'(1),$$

*then there exists $\epsilon > 0$ such that*

$$g_1(\delta) < g_2(\delta), \quad \forall \delta \in (\epsilon, 1).$$

Since $g_1(1) = g_2(1)$ and both functions are continuously differentiable, it suffices to show that $g_1'(1) > g_2'(1)$. If this holds, then for $\delta$ close to 1, we have $g_1(\delta) < g_2(\delta)$.

The derivatives at $\delta = 1$ are given by

$$g_1'(1) = \sum_{t=1}^{t_\tau^+ - 1} \left( t - \frac{\eta}{1-\eta} \right) \mathrm{P}^{\pi}_{\tau^{(t)}}$$

and

$$g_2'(1) = \sum_{t=1}^{t_\tau^+ - 1} \frac{\eta \mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}\mathrm{P}^{\pi}_{\tau^{(t)}}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}} + \frac{\eta \mathrm{P}^{\pi}_{\tau|(s_0',a_0')}}{1-(1-\eta)\mathrm{P}^{\pi}_{\tau'|(s_0',a_0')}}.$$

Observing the dependency on $t_\tau^+$, we see that when $t_\tau^+$ increases by 1,

$$g_1'(1) \quad \text{increases by} \quad (t_\tau^+ - \frac{\eta}{1-\eta})\mathrm{P}^{\pi}_{\tau^{(t_\tau^+)}},$$

whereas

$$g_2'(1) \quad \text{increases by} \quad \frac{\eta \mathrm{P}_{\tau'|(s_0',a_0')}^{\pi} \mathrm{P}_{\tau^{(t_\tau^+)}}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}.$$

By Lemma 2, we know that $t_\tau^+ \geq 2$ and since the increase in $g_1'(1)$ is greater than that of $g_2'(1)$, therefore, we only need to consider the base case $t_\tau^+ = 2$, where

$$g_1'(1) = \left(1 - \frac{\eta}{1-\eta}\right)\mathrm{P}_{\tau^{(1)}}^{\pi}$$

and

$$g_2'(1) = \frac{\eta \mathrm{P}_{\tau'|(s_0',a_0')}^{\pi} \mathrm{P}_{\tau^{(1)}}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}} + \frac{\eta \mathrm{P}_{\tau|(s_0',a_0')}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}.$$

Now, given that $t_\tau^+ = 2$, suppose $\tau = (s_0', a_0', s_1, a_1, s_2, \dots)$, then we have $\tau^{(1)} = (s_0', a_1, s_1^{(1)}, a_0', s_2, \dots)$. Thus,

$$\mathrm{P}_{\tau|(s_0',a_0')}^{\pi} = \pi(a_1|s_1)\mathrm{P}_{\tau|(s_1,a_1)}^{\pi} = \frac{\pi(a_1|s_1)\pi(a_1|s_0')\pi(a_0'|s_1^{(1)})\mathrm{P}_{\tau|(s_1,a_1)}^{\pi}}{\pi(a_1|s_0')\pi(a_0'|s_1^{(1)})} \tag{2}$$

$$= \frac{\pi(a_1|s_1)\mathrm{P}_{\tau^{(1)}}^{\pi}}{\pi(a_1|s_0')\pi(a_0'|s_1^{(1)})} \tag{3}$$

$$= \mathrm{P}_{\tau^{(1)}}^{\pi} \frac{\phi_{a_1}}{\sum_j \mathbb{I}(s_1,j)\phi_j} \frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\phi_{a_1}} \frac{\sum_j \mathbb{I}(s_1^{(1)},j)\phi_j}{\phi_{a_0'}} \tag{4}$$

$$= \mathrm{P}_{\tau^{(1)}}^{\pi} \frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\phi_{a_0'}} \frac{\sum_j \mathbb{I}(s_1^{(1)},j)\phi_j}{\sum_j \mathbb{I}(s_1,j)\phi_j} \tag{5}$$

$$= \mathrm{P}_{\tau^{(1)}}^{\pi} \frac{1}{1-\eta} \frac{\sum_j \mathbb{I}(s_1^{(1)},j)\phi_j}{\sum_j \mathbb{I}(s_1,j)\phi_j} \tag{6}$$

$$< \mathrm{P}_{\tau^{(1)}}^{\pi} \frac{1}{1-\eta} \frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_1,j)\phi_j} \tag{7}$$

The last inequality follows from the fact that $\mathcal{A}(s_1^{(1)}) \subset \mathcal{A}(s_0')$. Thus, we obtain

$$g_2'(1) = \frac{\eta \mathrm{P}_{\tau'|(s_0',a_0')}^{\pi} \mathrm{P}_{\tau^{(1)}}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}} + \frac{\eta \mathrm{P}_{\tau|(s_0',a_0')}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}$$

$$< \left(\frac{\eta \mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}} + \frac{1}{1-\eta} \frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_1,j)\phi_j}\right)\mathrm{P}_{\tau^{(1)}}^{\pi}$$

$$= (\alpha + \beta)\mathrm{P}_{\tau^{(1)}}^{\pi}$$

Thus, if $1 - \frac{\eta}{1-\eta} > \alpha + \beta$, then

$$g_2'(1) < (\alpha+\beta)\mathrm{P}_{\tau^{(1)}}^{\pi} < \left(1 - \frac{\eta}{1-\eta}\right)\mathrm{P}_{\tau^{(1)}}^{\pi} = g_1'(1).$$

This establishes the case for $t_\tau^+ = 2$. For $t_\tau^+ > 2$, the result follows from the previous discussion. Finally, by applying Theorem for Proposition 3, we complete the proof. ∎

*I. Proof of Theorem 1 (with the original theorem provided)*

**Theorem 1** (Policy Improvement-Version I). *Given a policy $\pi$ corresponding action-priority vector $\Phi$ and a trajectory $\tau' := (s_0', a_0', s_1', ...)$. Assume $a_0' = (m, j, i)$, let $\mathcal{A}_m(s_0') \subseteq \mathcal{A}(s_0')$ denote the set of feasible actions where the machine is fixed as $M_m$. Similarly, let $\mathcal{A}_{(j,i)}(s_0') \subseteq \mathcal{A}(s_0')$ denote the set of feasible actions where the job and operation are fixed as $J_j$ and $O_{j,i}$, respectively. Suppose $G(\tau') > V^\pi(s_0')$, then we have new policy $\pi'$ and its corresponding $\Phi'$ where*

$$\phi_a' = \begin{cases} \phi_a + (1-\delta)\sum_{a' \neq a} \mathbb{I}(s_0', a')\phi_{a'}, & \text{if } a = a_0' \\ \delta\phi_a, & \text{otherwise} \end{cases} \tag{8}$$

*and $\delta \in (0,1)$. Here the policy $\pi'$ and its corresponding $\Phi'$ updated according to (8), will perform better, that is $V^{\pi'}(s_0') > V^\pi(s_0')$ if $\delta$ satisfies **CD1**, **CD2** and $\mathcal{A}_m(s_0')$ satisfies **CD3**:*

**CD1** $\quad \dfrac{\delta^2(1-\delta\eta)}{(1-\eta)(\delta + (1-\delta)\frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}_j(s)\phi_j})} < \dfrac{1 - (1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}, \forall s \neq s_0'$

**CD2** $\quad \sum_{t=1}^{t_\tau^+ - 1} \delta^t \frac{1-\delta\eta}{1-\eta}\mathrm{P}_{\tau^{(t)}}^{\pi} < \delta\left(\sum_{t=1}^{t_\tau^+ - 1} \frac{\eta\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}\mathrm{P}_{\tau^{(t)}}^{\pi}}{1-(1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}} + \frac{\eta\mathrm{P}_{\tau|(s_0',a_0')}^{\pi}}{1-(1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}\right) + \sum_{t=1}^{t_\tau^+ - 1} \frac{(1-\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi})\mathrm{P}_{\tau^{(t)}}^{\pi}}{1-(1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}} - \frac{\eta\mathrm{P}_{\tau|(s_0',a_0')}^{\pi}}{1-(1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^{\pi}}, \forall \tau \in \mathcal{T}_3, \mathcal{T}_5$

**CD3** *For any $a \in \mathcal{A}_m(s_0') \cup \mathcal{A}_{(j,i)}(s_0')$, after executing action $a$, the system time of the resulting state remains the same as $\mathbb{T}_0'$, which is the system time of $s_0'$.*

*Proof.* First, the value function $V^\pi(s_0')$ can be expressed as

$$V^\pi(s_0') = \sum_{\tau \in \mathcal{T}_1} \mathrm{P}_\tau^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_2} \mathrm{P}_\tau^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_3} \pi(a_0'|s_0')\mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_4} \pi(a_0'|s_0')\mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi G(\tau')$$

Since $G(\tau') > V^\pi(s_0')$, we have

$$G(\tau') > \frac{\sum_{\tau \in \mathcal{T}_1} \mathrm{P}_\tau^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_2} \mathrm{P}_\tau^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_3} \pi(a_0'|s_0')\mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_4} \pi(a_0'|s_0')\mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau)}{1 - \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \tag{9}$$

Since in our modeling, once an action is executed, it will never be valid again, and together with Lemma 1, then

$$\mathrm{P}_{\tau|(s_0',a_0')}^{\pi'} = \mathrm{P}_{\tau|(s_0',a_0')}^\pi \tag{10}$$

Since

$$\begin{aligned}
\pi'(a_0'|s_0') &= \frac{\phi_{a_0'}'}{\sum_j \mathbb{I}(s_0',j)\phi_j'} \\
&= \frac{\phi_{a_0'} + (1-\delta)\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\delta\phi_j + \phi_{a_0'} + (1-\delta)\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j} \\
&= \frac{\phi_{a_0'} + (1-\delta)\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j + \phi_{a_0'}} = \frac{\phi_{a_0'} + (1-\delta)\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_0',j)\phi_j} \\
&= \pi(a_0'|s_0') + \frac{(1-\delta)\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_0',j)\phi_j}
\end{aligned}$$

Since $\eta = \frac{\sum_{j \neq a_0'} \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}(s_0',j)\phi_j}$, then

$$\pi'(a_0'|s_0') = \pi(a_0'|s_0') + (1-\delta)\eta \tag{11}$$

and

$$\pi(a_0'|s_0') = \frac{\phi_{a_0'}}{\sum_j \mathbb{I}(s_0',j)\phi_j} = 1 - \eta \tag{12}$$

By (9), (10), (11) and (12),

$$\begin{aligned}
V^{\pi'}(s_0') - V^\pi(s_0') &= \sum_{\tau \in \mathcal{T}_1} (\mathrm{P}_\tau^{\pi'} - \mathrm{P}_\tau^\pi)G(\tau) + \sum_{\tau \in \mathcal{T}_2} (\mathrm{P}_\tau^{\pi'} - \mathrm{P}_\tau^\pi)G(\tau) + \sum_{\tau \in \mathcal{T}_3} (1-\delta)\eta \mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) \\
&\quad + \sum_{\tau \in \mathcal{T}_4} (1-\delta)\eta \mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + (1-\delta)\eta \mathrm{P}_{\tau'|(s_0',a_0')}^\pi G(\tau') \\
&> \sum_{\tau \in \mathcal{T}_1} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta)\eta \mathrm{P}_{\tau'|(s_0',a_0')}^\pi}{1 - \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi - \mathrm{P}_\tau^\pi \right] G(\tau) + \sum_{\tau \in \mathcal{T}_2} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta)\eta \mathrm{P}_{\tau'|(s_0',a_0')}^\pi}{1 - \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi - \mathrm{P}_\tau^\pi \right] G(\tau) \\
&\quad + \sum_{\tau \in \mathcal{T}_3} \left[ (1-\delta)\eta + \frac{(1-\delta)\eta(1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi}{1 - \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \right] \mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_4} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta)\eta \mathrm{P}_{\tau'|(s_0',a_0')}^\pi}{1 - \pi(a_0'|s_0')\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi - \mathrm{P}_\tau^\pi \right] G(\tau) \\
&= \sum_{\tau \in \mathcal{T}_1} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi \right] G(\tau) + \sum_{\tau \in \mathcal{T}_2} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi \right] G(\tau) \\
&\quad + \sum_{\tau \in \mathcal{T}_3} \left[ \frac{(1-\delta)\eta}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \right] \mathrm{P}_{\tau|(s_0',a_0')}^\pi G(\tau) + \sum_{\tau \in \mathcal{T}_4} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi \right] G(\tau)
\end{aligned}$$

Since $G(\tau) < 0$, **CD1** and Lemma 3, we have

$$\begin{aligned}
\sum_{\tau \in \mathcal{T}_1} \left[ \mathrm{P}_\tau^{\pi'} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \mathrm{P}_\tau^\pi \right] G(\tau) &> \sum_{\tau \in \mathcal{T}_1} \left[ \frac{\delta^2(1-\delta\eta)\sum_j \mathbb{I}_j(s)\phi_j}{(1-\eta)(\delta\sum_j \mathbb{I}_j(s)\phi_j + (1-\delta)\sum_j \mathbb{I}(s_0',j)\phi_j)} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \right] \mathrm{P}_\tau^\pi G(\tau) \\
&= \sum_{\tau \in \mathcal{T}_1} \left[ \frac{\delta^2(1-\delta\eta)}{(1-\eta)(\delta + (1-\delta)\frac{\sum_j \mathbb{I}(s_0',j)\phi_j}{\sum_j \mathbb{I}_j(s)\phi_j})} + \frac{(1-\delta\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi - 1}{1 - (1-\eta)\mathrm{P}_{\tau'|(s_0',a_0')}^\pi} \right] \mathrm{P}_\tau^\pi G(\tau) \\
&\geq 0
\end{aligned}$$

By Lemma 5, we have

$$\mathcal{T}_2 = \bigcup_{\tilde{\tau} \in \mathcal{T}_3} \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$$

and

$$G(\tilde{\tau}^{(0)}) = \cdots = G(\tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}) = G(\tilde{\tau}).$$

Furthermore, for any $\tilde{\tau}^{(t)} \in \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}$ with $1 \le t \le t_{\tilde{\tau}}^+ - 1$, it follows from (1) that the action $a_0'$ is taken at stage $t$. Then, by Lemma 4, we obtain

$$\mathrm{P}_{\tau}^{\pi'} G(\tilde{\tau}^{(t)}) = \mathrm{P}_{\tau}^{\pi'} G(\tilde{\tau}) < \frac{\delta^t (1 - \delta\eta)}{1 - \eta} \mathrm{P}_{\tau}^{\pi} G(\tilde{\tau}).$$

Hence, we have

$$\sum_{\tau \in \mathcal{T}_2} \left[ \mathrm{P}_{\tau}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tau}^{\pi} \right] G(\tau) + \sum_{\tau \in \mathcal{T}_3} \left[ \frac{(1 - \delta)\eta}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \right] \mathrm{P}_{\tau|(s_0', a_0')}^{\pi} G(\tau) \tag{13}$$

$$= \sum_{\tilde{\tau} \in \mathcal{T}_3} \left( \sum_{\tilde{\tau}^{(t)} \in \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}} \left[ \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi} \right] G(\tilde{\tau}^{(t)}) + \left[ \frac{(1 - \delta)\eta}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \right] \mathrm{P}_{\tilde{\tau}|(s_0', a_0')}^{\pi} G(\tilde{\tau}) \right) \tag{14}$$

$$= \sum_{\tilde{\tau} \in \mathcal{T}_3} \left( \sum_{\tilde{\tau}^{(t)} \in \{\tilde{\tau}^{(0)}, \ldots, \tilde{\tau}^{(t_{\tilde{\tau}}^+ - 1)}\}} \left[ \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi} \right] + \left[ \frac{(1 - \delta)\eta}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \right] \mathrm{P}_{\tilde{\tau}|(s_0', a_0')}^{\pi} \right) G(\tilde{\tau}) \tag{15}$$

$$> \sum_{\tilde{\tau} \in \mathcal{T}_3} \left( \sum_{t=1}^{t_{\tilde{\tau}}^+ - 1} \left[ \frac{\delta^t (1 - \delta\eta)}{1 - \eta} \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi} - \frac{1 - (1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tilde{\tau}^{(t)}}^{\pi} \right] + \left[ \frac{(1 - \delta)\eta}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \right] \mathrm{P}_{\tilde{\tau}|(s_0', a_0')}^{\pi} \right) G(\tilde{\tau}) \tag{16}$$

$$> 0 \tag{17}$$

The last inequality follows from **CD2**. Similarly, for the set $\mathcal{T}_4$, applying Lemma 4 and Lemma 5, we obtain

$$\sum_{\tau \in \mathcal{T}_4} \left[ \mathrm{P}_{\tau}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tau}^{\pi} \right] G(\tau) = \sum_{\tau'^{(t)} \in \{\tau'^{(0)}, \ldots, \tau'^{(t_{\tau'}^+ - 1)}\}} \left[ \mathrm{P}_{\tau'^{(t)}}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tau'^{(t)}}^{\pi} \right] G(\tau'^{(t)})$$

$$= \sum_{t=1}^{t_{\tau'}^+ - 1} \left[ \mathrm{P}_{\tau'^{(t)}}^{\pi'} + \frac{(1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi} - 1}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tau'^{(t)}}^{\pi} \right] G(\tau')$$

$$> \sum_{t=1}^{t_{\tau'}^+ - 1} \left[ \frac{\delta^t (1 - \delta\eta)}{1 - \eta} \mathrm{P}_{\tau'^{(t)}}^{\pi} - \frac{1 - (1 - \delta\eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}}{1 - (1 - \eta)\mathrm{P}_{\tau'|(s_0', a_0')}^{\pi}} \mathrm{P}_{\tau'^{(t)}}^{\pi} \right] G(\tau')$$

$$> 0$$

The last inequality follows from **CD2**. Consequently, we have proven that $V^{\pi'}(s_0') > V^{\pi}(s_0')$. $\blacksquare$

## IV. NUMERICAL RESULTS

### REFERENCES

[1] M. Spivak, *Calculus*. Cambridge University Press, 2006.