

APPENDIX A  
PROOF OF THEOREM 1

The Stage III resource allocation problem is a linear programming problem with the objective of minimizing operational cost  $C(\mathbf{x})$  subject to constraints (3), (4), (9b), and (9c).

Given the cost structure where  $c^{TH} < c^{TL}$  and  $c^{IH} > c^{IL}$ , the optimal allocation strategy is to:

- 1) Assign as much training demand as possible to high-performance GPUs (cheaper for training).
- 2) Assign as much inference demand as possible to low-performance GPUs (cheaper for inference).

The optimal allocation follows directly from this cost-minimization principle:

- When  $\sum_{n=1}^N g_n^T \leq G^H$ , all training demand can be allocated to high-performance GPUs:  $x_n^{TH*} = g_n^T$ .
- When  $\sum_{n=1}^N g_n^T > G^H$ , high-performance GPUs are fully utilized and training demand is proportionally allocated:  $x_n^{TH*} = g_n^T \frac{G^H}{\sum_{n=1}^N g_n^T}$ .
- When  $\sum_{n=1}^N g_n^I \leq G^L$ , all inference demand can be allocated to low-performance GPUs:  $x_n^{IL*} = g_n^I$ .
- When  $\sum_{n=1}^N g_n^I > G^L$ , low-performance GPUs are fully utilized and inference demand is proportionally allocated:  $x_n^{IL*} = g_n^I \frac{G^L}{\sum_{n=1}^N g_n^I}$ .

The remaining allocations are determined by the demand constraints:

$$\begin{aligned} x_n^{TL*} &= g_n^T - x_n^{TH*}, \\ x_n^{IH*} &= g_n^I - x_n^{IL*}. \end{aligned}$$

If  $\sum_{n=1}^N g_n^T + \sum_{n=1}^N g_n^I > G^H + G^L = G$ , the problem is infeasible due to insufficient GPU capacity.

APPENDIX B  
PROOF OF THEOREM 2

To solve the Problem (8), due to the  $P_n^T(g_n^T)$  is only determined by  $g_n^T$  and  $P_n^I(g_n^I)$  is only determined by  $g_n^I$ . We can solve the problem by optimize  $g_n^T$  and  $g_n^I$ , respectively.

*A. For the training service*

Client  $n$  solves

$$\max_{g_n^T \geq 0} P_n^T(g_n^T) = \gamma_n^T \cdot \log(1 + g_n^T/\theta_n) - p^T \cdot g_n^T. \quad (25)$$

Taking the first-order derivative,

$$\frac{\partial P_n^T}{\partial g_n^T} = \gamma_n^T \cdot \frac{1/\theta_n}{1 + g_n^T/\theta_n} - p^T = \frac{\gamma_n^T}{\theta_n + g_n^T} - p^T. \quad (26)$$

Setting the derivative to zero,

$$\frac{\gamma_n^T}{\theta_n + g_n^T} = p^T. \quad (27)$$

Solving for  $g_n^T$ ,

$$\theta_n + g_n^T = \frac{\gamma_n^T}{p^T}, \quad (28)$$

$$g_n^T = \frac{\gamma_n^T}{p^T} - \theta_n = \theta_n \left( \frac{\gamma_n^T}{p^T \theta_n} - 1 \right). \quad (29)$$

Considering  $g_n^T \geq 0$ , we require:

$$\frac{\gamma_n^T}{p^T} - \theta_n \geq 0 \Rightarrow p^T \leq \frac{\gamma_n^T}{\theta_n}. \quad (30)$$

Therefore, the optimal training demand is

$$g_n^*(p^T) = \theta_n \left[ \frac{\gamma_n^T}{p^T \theta_n} - 1 \right]_+ = \mathbf{1}(p^T \leq \frac{\gamma_n^T}{\theta_n}) \theta_n \left( \frac{\gamma_n^T}{p^T \theta_n} - 1 \right). \quad (31)$$

*B. For the inference service*

Client  $n$  solves

$$\max_{g_n^I \geq 0} P_n^I(g_n^I) = -\gamma_n^I \cdot \frac{1}{\left( \frac{g_n^I}{d_n} - \lambda_n \right)^2} - p^I \cdot g_n^I. \quad (32)$$

Taking the first-order derivative,

$$\frac{\partial P_n^I}{\partial g_n^I} = \gamma_n^I \cdot \frac{2}{d_n \left( \frac{g_n^I}{d_n} - \lambda_n \right)^3} - p^I. \quad (33)$$

Setting the derivative to zero,

$$\gamma_n^I \cdot \frac{2}{d_n \left( \frac{g_n^I}{d_n} - \lambda_n \right)^3} = p^I. \quad (34)$$

Solving for  $g_n^I$ ,

$$\left( \frac{g_n^I}{d_n} - \lambda_n \right)^3 = \frac{2\gamma_n^I}{d_n p^I}, \quad (35)$$

and we denote the solution  $g_n^I$  of (35) by  $g_n^{I(1)}$  and

$$g_n^{I(1)} = d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right). \quad (36)$$

Back to (33), we identify the positive and negative of it. We first give the conclusion: The (33) has a singularity  $g_n^{I(2)} = d_n \lambda_n$ , where  $\frac{\partial P_n^I}{\partial g_n^I}(g_n^{I(2)}) = +\infty$  and  $\frac{\partial P_n^I}{\partial g_n^I}(g_n^{I(2)}) = -\infty$ . The  $\frac{\partial P_n^I}{\partial g_n^I}$  is negative within  $[0, g_n^{I(2)}]$ , positive within  $(g_n^{I(2)}, g_n^{I(1)})$  and negative within  $(g_n^{I(1)}, +\infty)$ . Thus, we have  $P_n^I(g_n^I)$  is decreasing within  $[0, g_n^{I(2)}]$ , increasing within  $(g_n^{I(2)}, g_n^{I(1)})$  and decreasing within  $(g_n^{I(1)}, +\infty)$ . So the maximum of  $\frac{\partial P_n^I}{\partial g_n^I}$  may be  $g_n^I = 0$  or  $g_n^I = g_n^{I(1)}$ . This is proved in Appendix C.

We define a threshold price  $p_n^{I\dagger}$  to identify which one is the maximum point: when  $p^I \leq p_n^{I\dagger}$ , the maximum point is  $g_n^{I*} = g_n^{I(1)}$ , otherwise, the maximum point is  $g_n^{I*} = 0$ .

Substitute the two points into  $P_n^I(g_n^I)$  and the threshold price  $p_n^{I\dagger}$  is obtained by solving

$$P_n^I(0) = P_n^I(g_n^{I(1)}(p^{I\dagger})), \quad (37)$$

which is

$$p_n^{I\dagger} d_n \lambda_n - \frac{3}{\sqrt[3]{4}} (p_n^{I\dagger} d_n)^{2/3} (\gamma_n^I)^{1/3} = \frac{-\gamma_n^I}{\lambda_n^2}. \quad (38)$$

Continue simplifying (38) as

$$\frac{\gamma_n^I}{\lambda_n^2} = \frac{3}{2^{2/3}} (\gamma_n^I)^{1/3} (d_n p^I)^{2/3} + p^I d_n \lambda_n. \quad (39)$$

Solving this equation and we have

$$p_n^{I\dagger} = \frac{\gamma_n^I}{4d_n\lambda_n^3}. \quad (40)$$

Thus, when  $p^I \leq p_n^{I\dagger}$ , the optimal point is  $g_n^{I*}(p^I) = g_n^{I(1)} = d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right)$ , otherwise  $g_n^{I*}(p^I) = 0$ . We can reexpress it as

$$\begin{aligned} g_n^{I*}(p^I) &= \begin{cases} d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right), & \text{if } p^I \leq p_n^{I\dagger}, \\ 0, & \text{otherwise,} \end{cases} \\ &= \mathbf{1}(p^I \leq p_n^{I\dagger}) d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right). \end{aligned} \quad (41)$$

### APPENDIX C PROOF OF THE MONOTONICITY OF $\frac{\partial P_n^I}{\partial g_n^I}$

For ease of description, we denote  $x = \frac{g_n^I}{d_n} - \lambda_n$ , then the derivative can be written as

$$\frac{\partial P_n^I}{\partial g_n^I} = \gamma_n^I \cdot \frac{2}{d_n x^3} - p^I. \quad (42)$$

where  $x$  varies with  $g_n^I$ . The sign of the derivative depends on the sign and magnitude of  $x$ .

- 1) When  $x < 0$ , i.e.,  $g_n^I < d_n \lambda_n$ , there is  $x^3 < 0$ , and  $\frac{2}{d_n x^3} < 0$ , thus  $\gamma_n^I \cdot \frac{2}{d_n x^3} - p^I < 0$ . The original function  $P_n^I$  decreases as  $g_n^I$  increases.
- 2) When  $x > 0$ , i.e.,  $g_n^I > d_n \lambda_n$ , there is  $x^3 > 0$ , so  $\frac{2}{d_n x^3} > 0$ , thus  $\gamma_n^I \cdot \frac{2}{d_n x^3} > 0$ . The derivative  $\frac{\partial P_n^I}{\partial g_n^I} = \gamma_n^I \cdot \frac{2}{d_n x^3} - p^I$ , and its sign depends on the comparison between  $\gamma_n^I \cdot \frac{2}{d_n x^3}$  and  $p^I$ . Define the critical point  $x_c = \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3}$ , corresponding to

$$g_c = d_n \lambda_n + d_n x_c = d_n \lambda_n + d_n \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3}. \quad (43)$$

- When  $x < x_c$ , i.e.,  $g_n^I < g_c$ , there is  $\gamma_n^I \cdot \frac{2}{d_n x^3} > p^I$ , so the derivative  $\frac{\partial P_n^I}{\partial g_n^I} > 0$ . The original function  $P_n^I$  increases as  $g_n^I$  increases.
  - When  $x > x_c$ , i.e.,  $g_n^I > g_c$ , there is  $\gamma_n^I \cdot \frac{2}{d_n x^3} < p^I$ , so the derivative  $\frac{\partial P_n^I}{\partial g_n^I} < 0$ . The original function  $P_n^I$  decreases as  $g_n^I$  increases.
- 3) When  $x = 0$ , i.e.,  $g_n^I = d_n \lambda_n$ , the denominator in the derivative expression becomes zero, so the derivative is undefined.

In one word, the  $\frac{\partial P_n^I}{\partial g_n^I}$  is negative within  $[0, g_n^{I(2)}]$ , positive within  $(g_n^{I(2)}, g_n^{I(1)})$  and negative within  $(g_n^{I(1)}, +\infty)$ . Thus, we have  $P_n^I(g_n^I)$  is decreasing within  $[0, g_n^{I(2)}]$ , increasing within  $(g_n^{I(2)}, g_n^{I(1)})$  and decreasing within  $(g_n^{I(1)}, +\infty)$ . So the maximum of  $\frac{\partial P_n^I}{\partial g_n^I}$  may be  $g_n^I = 0$  or  $g_n^I = g_n^{I(1)}$ .

### APPENDIX D PROOF OF EQUATION (17)

The cost function  $C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p})))$  is derived by substituting the optimal Stage III allocation into the cost function and considering all possible cases of resource utilization:

- 1) Case 1:  $D^T(p^T) \leq G^H$  and  $D^I(p^I) \leq G^L$ .

All training on high-performance GPUs, all inference on low-performance GPUs and

$$C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) = c^{TH} D^T(p^T) + c^{IL} D^I(p^I). \quad (44)$$

- 2) Case 2:  $D^T(p^T) > G^H$  and  $D^I(p^I) \leq G^L$ .

High-performance GPUs fully utilized for training, excess training on low-performance GPUs, all inference on low-performance GPUs and

$$C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) = c^{TH} G^H + c^{TL} (D^T(p^T) - G^H) + c^{IL} D^I(p^I). \quad (45)$$

- 3) Case 3:  $D^T(p^T) \leq G^H$  and  $D^I(p^I) > G^L$ .

All training on high-performance GPUs, low-performance GPUs fully utilized for inference, excess inference on high-performance GPUs and

$$C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) = c^{TH} D^T(p^T) + c^{IL} G^L + c^{IH} (D^I(p^I) - G^L). \quad (46)$$

- 4) Case 4:  $D^T(p^T) > G^H$  and  $D^I(p^I) > G^L$ .

Infeasible due to insufficient resources.

The indicator functions in Equation (17) precisely capture these four cases.

### APPENDIX E DERIVATION OF $\mathcal{C}_{p,q}^T$ AND $\mathcal{C}_{p,q}^I$

In subregion  $Q_{p,q}$  with  $p^T \in [a_p^T, a_{p+1}^T]$  and  $p^I \in [a_q^I, a_{q+1}^I]$ , the active client sets are determined by the price thresholds:

For training services, client  $n$  has positive demand when

$$p^T \leq \frac{\gamma_n^T}{\theta_n}. \quad (47)$$

Since  $p^T \in [a_p^T, a_{p+1}^T]$ , the active training clients are

$$\mathcal{C}_{p,q}^T = \left\{ n \in \mathcal{N} : \frac{\gamma_n^T}{\theta_n} \geq a_p^T \right\}. \quad (48)$$

For inference services, client  $n$  has positive demand when

$$p^I \leq p_n^{I\dagger}. \quad (49)$$

The active inference clients are

$$\mathcal{C}_{p,q}^I = \left\{ n \in \mathcal{N} : p_n^{I\dagger} \geq a_q^I \right\}. \quad (50)$$

These sets remain constant within each subregion  $Q_{p,q}$ .

### APPENDIX F PROOF OF EQUATION (21)

The additional critical points  $a_p'^T$  and  $a_q'^I$  are derived from the resource capacity constraints.

For training services, the critical point where  $D^T(p^T) = G^H$  is found by solving:

$$\sum_{n \in \mathcal{C}_{p,q}^T} \theta_n \left( \frac{\gamma_n^T}{p^T \theta_n} - 1 \right) = G^H. \quad (51)$$

And we denote the solution  $p^T$  of (51) by  $a'^T_p$  and

$$a'^T_p = \frac{\sum_{n \in \mathcal{C}_{p,q}^T} \gamma_n^T}{G^H + \sum_{n \in \mathcal{C}_{p,q}^T} \theta_n}. \quad (52)$$

For inference services, the critical point where  $D^I(p^I) = G^L$  is found by solving

$$\sum_{n \in \mathcal{C}_{p,q}^I} d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right) = G^L \quad (53)$$

And we denote the solution  $p^I$  of (53) by  $a'^I_q$  and

$$a'^I_q = \left( \frac{\sum_{n \in \mathcal{C}_{p,q}^I} d_n^{2/3} (2\gamma_n^I)^{1/3}}{G^L - \sum_{n \in \mathcal{C}_{p,q}^I} d_n \lambda_n} \right)^3 \quad (54)$$

## APPENDIX G PROOF OF EQUATION (22)

In each subregion  $Q'_{p,q}$ , the indicator functions in Equation (17) become constants, allowing us to express the cost function as a linear combination of the demand functions

$$C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) = A_{p,q} D_{p,q}^T(p^T) + B_{p,q} D_{p,q}^I(p^I) + E_{p,q}, \quad (55)$$

where the coefficients  $A_{p,q}$ ,  $B_{p,q}$ , and  $E_{p,q}$  depend on which case of the piecewise cost function applies in subregion  $Q'_{p,q}$ .

The revenue function remains

$$R(\mathbf{p}, \mathbf{g}^*(\mathbf{p})) = p^T D_{p,q}^T(p^T) + p^I D_{p,q}^I(p^I). \quad (56)$$

The specific values of  $A_{p,q}$ ,  $B_{p,q}$ , and  $E_{p,q}$  depend on which of the four cases in Equation (17) applies in subregion  $Q'_{p,q}$ <sup>3</sup>:

Case 1:  $D_{p,q}^T(p^T) \leq G^H$  and  $D_{p,q}^I(p^I) \leq G^L$ .

$$\begin{aligned} C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) &= c^{TH} D_{p,q}^T(p^T) + c^{IL} D_{p,q}^I(p^I), \\ A_{p,q} &= c^{TH}, \\ B_{p,q} &= c^{IL}, \\ E_{p,q} &= 0. \end{aligned}$$

Case 2:  $D_{p,q}^T(p^T) > G^H$  and  $D_{p,q}^I(p^I) \leq G^L$ .

$$\begin{aligned} C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) &= c^{TH} G^H + c^{TL}(D_{p,q}^T(p^T) - G^H) + c^{IL} D_{p,q}^I(p^I) \\ &= c^{TL} D_{p,q}^T(p^T) + c^{IL} D_{p,q}^I(p^I) + (c^{TH} - c^{TL}) G^H, \\ A_{p,q} &= c^{TL}, \\ B_{p,q} &= c^{IL}, \\ E_{p,q} &= (c^{TH} - c^{TL}) G^H. \end{aligned}$$

<sup>3</sup>Here, the Case 2 - 3 do not imply that the entire subregion satisfies the constraints (7b), as  $D_{p,q}^T(p^T) + D_{p,q}^I(p^I) > G$  may still happen.

Case 3:  $D_{p,q}^T(p^T) \leq G^H$  and  $D_{p,q}^I(p^I) > G^L$ .

$$\begin{aligned} C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) &= c^{TH} G^H + c^{TL}(D_{p,q}^T(p^T) - G^H) + c^{IL} D_{p,q}^I(p^I) \\ &= c^{TH} D_{p,q}^T(p^T) + c^{IH} D_{p,q}^I(p^I) + (c^{IL} - c^{IH}) G^L, \\ A_{p,q} &= c^{TH}, \\ B_{p,q} &= c^{IH}, \\ E_{p,q} &= (c^{IL} - c^{IH}) G^L. \end{aligned}$$

Case 4:  $D_{p,q}^T(p^T) > G^H$  and  $D_{p,q}^I(p^I) > G^L$ .

$$\begin{aligned} C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) &= +\infty, \\ A_{p,q} &= +\infty, \\ B_{p,q} &= +\infty, \\ E_{p,q} &= +\infty. \end{aligned}$$

The revenue function remains

$$R(\mathbf{p}, \mathbf{g}^*(\mathbf{p})) = p^T D_{p,q}^T(p^T) + p^I D_{p,q}^I(p^I). \quad (57)$$

Therefore, the profit function becomes

$$\begin{aligned} F_{p,q}(\mathbf{p}) &= R(\mathbf{p}, \mathbf{g}^*(\mathbf{p})) - C(\mathbf{x}^*(\mathbf{g}^*(\mathbf{p}))) \\ &= p^T D_{p,q}^T(p^T) + p^I D_{p,q}^I(p^I) \\ &\quad - A_{p,q} D_{p,q}^T(p^T) - B_{p,q} D_{p,q}^I(p^I) - E_{p,q}. \end{aligned}$$

This reformulation eliminates the indicator functions and provides a continuous objective function within each subregion.

## APPENDIX H PROOF OF THEOREM 3

Within each subregion  $Q'_{p,q}$ , we analyze the monotonicity properties of the objective function  $F_{p,q}(\mathbf{p})$ :

A. For  $p^I$ :

The  $p^I$  related part of the objective function is

$$F_{p,q}(p^I) = p^I D_{p,q}^I(p^I) - B_{p,q} D_{p,q}^I(p^I), \quad (58)$$

$$\text{where } D_{p,q}^I(p^I) = \sum_{n \in \mathcal{C}_{p,q}^I} d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n p^I} \right)^{1/3} \right).$$

Let  $\Lambda = \sum_{n \in \mathcal{C}_{p,q}^I} d_n \lambda_n$  and  $\Upsilon = \sum_{n \in \mathcal{C}_{p,q}^I} d_n^{2/3} (2\gamma_n^I)^{1/3}$ , then

$$D_{p,q}^I(p^I) = \Lambda + \Upsilon p^{I-1/3}. \quad (59)$$

Substituting to the (58), we have

$$\begin{aligned} F_{p,q}^I(p^I) &= p^I (\Lambda + \Upsilon p^{I-1/3}) - B_{p,q} (\Lambda + \Upsilon p^{I-1/3}) \\ &= \Lambda p^I + \Upsilon p^{I-2/3} - B_{p,q} \Lambda - B_{p,q} \Upsilon p^{I-1/3}. \end{aligned} \quad (60)$$

Taking the derivative with respect to  $p^I$ ,

$$\frac{\partial F_{p,q}^I}{\partial p^I} = \Lambda + \frac{2}{3} \Upsilon p^{I-4/3} + \frac{1}{3} B_{p,q} \Upsilon p^{I-4/3}. \quad (61)$$

Now we analyze each term

- $\Lambda > 0$  (sum of positive terms).
- $\frac{2}{3} \Upsilon p^{I-4/3} > 0$  (since  $\Upsilon > 0$  and  $p^I > 0$ ).
- $\frac{1}{3} B_{p,q} \Upsilon p^{I-4/3} > 0$  (since  $B_{p,q} > 0$ ).

All three terms are strictly positive for  $p^I > 0$ . Therefore,

$$\frac{\partial F_{p,q}^I}{\partial p^I} > 0 \quad \text{for all } p^I > 0. \quad (62)$$

This proves that  $F_{p,q}(p)$  is strictly increasing with respect to  $p^I$  within each subregion  $Q'_{p,q}$ . Therefore, the optimal  $p^I$  is at the upper boundary of the subregion

$$p_{p,q}^{I*} = a_{q+1}^I. \quad (63)$$

### B. For $p^T$ :

The  $p^T$  related part of the objective function is

$$F_{p,q}(p^T) = p^T D_{p,q}^T(p^T) - A_{p,q} D_{p,q}^T(p^T) \quad (64)$$

where

$$\begin{aligned} D_{p,q}^T(p^T) &= \sum_{n \in \mathcal{C}_{p,q}^T} \theta_n \left( \frac{\gamma_n^T}{p^T \theta_n} - 1 \right) \\ &= \frac{1}{p^T} \sum_{n \in \mathcal{C}_{p,q}^T} \gamma_n^T - \sum_{n \in \mathcal{C}_{p,q}^T} \theta_n. \end{aligned} \quad (65)$$

Let  $\Gamma = \sum_{n \in \mathcal{C}_{p,q}^T} \gamma_n^T$  and  $\Theta = \sum_{n \in \mathcal{C}_{p,q}^T} \theta_n$ , then

$$D_{p,q}^T(p^T) = \frac{\Gamma}{p^T} - \Theta. \quad (66)$$

Substituting to the (64), we have

$$\begin{aligned} F_{p,q}^T(p^T) &= p^T \left( \frac{\Gamma}{p^T} - \Theta \right) - A_{p,q} \left( \frac{\Gamma}{p^T} - \Theta \right) \\ &= \Gamma - \Theta p^T - \frac{A_{p,q} \Gamma}{p^T} + A_{p,q} \Theta. \end{aligned} \quad (67)$$

Taking the derivative with respect to  $p^T$ ,

$$\frac{\partial F_{p,q}^T}{\partial p^T} = -\Theta + \frac{A_{p,q} \Gamma}{(p^T)^2}. \quad (68)$$

The  $F_{p,q}^T(p^T)$  is not monotonic with  $p^T$ .

### 1) Derivation of $p'^T$ (critical point for $p^T$ ):

Setting the derivative to zero to find the critical point,

$$-\Theta + \frac{A_{p,q} \Gamma}{(p^T)^2} = 0. \quad (69)$$

And we have

$$p^T = \sqrt{\frac{A_{p,q} \Gamma}{\Theta}}. \quad (70)$$

This gives us the critical point  $p'^T = \sqrt{\frac{A_{p,q} \sum_{n \in \mathcal{C}_{p,q}^T} \gamma_n^T}{\sum_{n \in \mathcal{C}_{p,q}^T} \theta_n}}$ .

Next, we check the second derivative

$$\frac{\partial^2 F_{p,q}^T}{\partial (p^T)^2} = -\frac{2A_{p,q} \Gamma}{(p^T)^3} < 0 \quad (\text{since } A_{p,q} > 0, \Gamma > 0, p^T > 0). \quad (71)$$

Thus,  $p'^T$  is indeed a maximum point of the objective function. Solving this equation gives the critical point  $p'^T$ . The optimal solution depends on the position of  $p'^T$  relative to the subregion boundaries  $[a_p^T, a_{p+1}^T]$ .

The three cases in the theorem cover all possible scenarios:

- 1) If  $p'^T \leq a_p^T$ , the function is decreasing in the subregion, so the optimal is at the lower boundary  $a_p^T$ .
- 2) If  $p'^T \in [a_p^T, a_{p+1}^T]$ , the critical point is feasible and optimal.
- 3) If  $p'^T \geq a_{p+1}^T$ , the function is increasing in the subregion, so the optimal is at the upper boundary  $a_{p+1}^T$ .

After that, we identify the optimal  $p^T$  in each subregion. However, we still not consider the constraint (7b). For a subregion  $Q'_{p,q}$ , there may be some points satisfied the constraint but the others not. Thus, we need check the impact of (7b) to the closed-form solution.

Since the optimal  $p^I$  is the upper bound of  $p^I$  in the subregion, considering the  $D_{p,q}^I(p^I)$  is deceasing to  $p^I$ , if there exists feasible points in this subregion, the upper bound is also the best choice to satisfy the constraint. Thus, the solution of  $p^I$  is not affected by the constraint (7b). Next, we will identify the constraint's impact to the  $p^T$ .

### 2) Derivation of $b_p^T$ (constraint boundary for $p^T$ ):

The resource constraint from (7b) is

$$\sum_{n=1}^N (g_n^T(\mathbf{p}) + g_n^I(\mathbf{p})) \leq G. \quad (72)$$

Substituting the  $D_{p,q}^T(p^T)$  and  $D_{p,q}^I(p^I)$  to it, we have

$$D_{p,q}^T(p^T) + D_{p,q}^I(a_{q+1}^I) \leq G. \quad (73)$$

We already have  $p_{p,q}^{I*} = a_{q+1}^I$  from the monotonicity analysis. Substituting and we have

$$D_{p,q}^T(p^T) + D_{p,q}^I(a_{q+1}^I) \leq G. \quad (74)$$

Using the expression for  $D_{p,q}^T(p^T)$ ,

$$\left( \frac{\Gamma}{p^T} - \Theta \right) + D_{p,q}^I(a_{q+1}^I) \leq G. \quad (75)$$

Solving for  $p^T$ , we have

$$p^T \geq \frac{\Gamma}{G + \Theta - D_{p,q}^I(a_{q+1}^I)}. \quad (76)$$

Thus, the constraint boundary is

$$b_p^T = \frac{\sum_{n \in \mathcal{C}_{p,q}^T} \gamma_n^T}{G + \sum_{n \in \mathcal{C}_{p,q}^T} \theta_n - D_{p,q}^I(a_{q+1}^I)}, \quad (77)$$

where  $D_{p,q}^I(a_{q+1}^I) = \sum_{n \in \mathcal{C}_{p,q}^I} d_n \left( \lambda_n + \left( \frac{2\gamma_n^I}{d_n a_{q+1}^I} \right)^{1/3} \right)$ .

In a word, to satisfy (7b), the solution  $p^T$  in each subregion must be no less than  $b_p^T$ .

### 3) Optimal solution selection:

Based on the position of  $p'^T$  relative to the subregion boundaries  $[a_p^T, a_{p+1}^T]$  and the constraint boundary  $b_p^T$ , we rewrite three cases:

- Case 1:  $p'^T < a_p^T$ .

The critical point is below the subregion, so the function is decreasing in the subregion. The optimal is at the lower boundary, but must satisfy the constraint

$$p_{p,q}^{T*} = \max(a_p^T, b_p^T). \quad (78)$$

- Case 2:  $a_p^T \leq p'^T < a_{p+1}^T$ .

The critical point is within the subregion. The optimal is the critical point, but must satisfy the constraint

$$p_{p,q}^{T*} = \max(p'^T, b_p^T). \quad (79)$$

- Case 3:  $a_{p+1}^T \leq p'^T$ .

The critical point is above the subregion, so the function is increasing in the subregion. The optimal is at the upper boundary

$$p_{p,q}^{T*} = a_{p+1}^T. \quad (80)$$

In this case, we don't need to consider  $b_p^T$  explicitly because we've already verified the feasibility of the subregion (At least the point  $(a_{p+1}^T, a_{q+1}^I)$  satisfy the constraints, which means there must be  $b_p^T \leq a_{p+1}^T$ ), which ensures that  $a_{p+1}^T$  satisfies the constraint.

This completes the proof of Theorem 3.