

Appendix for the Greedy Scheduling

1 Greedy-based Solution

To address such an NP-hard problem, the request manager in **Apt-Serve** uses a greedy-based approximate solution. Its core idea is to allocate memory occupation iteratively to candidate requests based on their marginal gain. Specifically, this involves continuously assigning additional memory usage to the request that provides the highest marginal gain per unit of memory consumption, thereby progressively approaching the optimal solution. For simplicity, we omit the superscript e denoting the iteration in the following illustration.

For a given request i , denote its current memory usage as \bar{m}_i , if $\bar{m}_i = 0$, a marginal memory usage increase Δm_i to it involves assigning half of its maximum memory requirement $\frac{m_i}{2}$, indicating that request i is scheduled using hidden cache. Similarly, if $\bar{m}_i = \frac{m_i}{2}$, the marginal memory usage increase Δm_i involves assigning the other half of its maximum memory requirement, indicating that request i is now scheduled using KV cache, given that it was previously scheduled with hidden cache. If $\bar{m}_i = m_i$, assigning further memory does not provide any additional gain in value. Therefore, for a given request i , its marginal gain in value θ_i based on its current memory usage \bar{m}_i , is as follows:

$$\theta_i = \begin{cases} \frac{(p_i - (|W| + |R|)\alpha m_i) - 0}{\frac{m_i}{2} - 0} = \frac{2p_i}{m_i} - 2(|W| + |R|)\alpha, & \text{if } \bar{m}_i = 0; \\ \frac{p_i - (p_i - (|W| + |R|)\alpha m_i)}{m_i - m_i/2} = 2(|W| + |R|)\alpha, & \text{if } \bar{m}_i = \frac{m_i}{2}; \\ 0, & \text{if } \bar{m}_i = m_i. \end{cases}$$

Note that for a given request i to be scheduled, there may be cases where its hidden cache usage imposes too much of a negative impact on other requests. This is evident when the marginal gain to the overall scheduling value from allocating just enough memory to hold its hidden cache (from $\bar{m}_i = 0$ to $\bar{m}_i = \frac{m_i}{2}$) is smaller than which by directly assigning the exact memory space to hold its KV cache (from $\bar{m}_i = 0$ to $\bar{m}_i = m_i$). In such cases, its hidden cache usage is avoided, thus its marginal gain is refined as follows:

$$\theta_i = \begin{cases} \frac{p_i - 0}{\frac{m_i}{2} - 0} = \frac{2p_i}{m_i}, & \text{if } \bar{m}_i = 0; \\ 0, & \text{if } \bar{m}_i = m_i. \end{cases}$$

Assume the total number of candidate schedules (how much marginal memory usage increase to which requests) is n , all the possible marginal gain $\theta_1, \theta_2, \dots, \theta_n$

Algorithm 1 Greedy-based Scheduling by Apt-Serve.

Input: The candidate request set U ; memory constraint M ; the candidate schedule set Υ (Eq. 1).

Output: Final scheduling decisions S (Eq. 2).

```
1: Initialize scheduling decisions  $S \leftarrow \{(x_i, y_i) = (0, 0) | \forall i \in U\}$ ;
2: Initialize the current memory usage  $\bar{M} \leftarrow \{\bar{m}_i = 0 | \forall i \in U\}$ ;
3: Initialize the overall scheduling values  $V \leftarrow 0$ ;
4: Sort  $\Upsilon$  based on the marginal gain in a decreasing order and derive  $\hat{\Upsilon} \leftarrow \{(\hat{\theta}_j, \hat{r}_j, \Delta\hat{m}_j, \hat{m}_{\hat{r}_j}) | j = 1, 2, \dots, n\}$ ;
5: for  $j = 1, 2, \dots, n$  do
6:   if  $\text{sum}(\bar{M}) + \Delta\hat{m}_j \leq M$  then
7:      $\bar{m}_{\hat{r}_j} \leftarrow \bar{m}_{\hat{r}_j} + \Delta\hat{m}_j$  ▷ Update memory usage for  $\hat{r}_j$ .
8:      $V \leftarrow V + \hat{\theta}_j \cdot \Delta\hat{m}_j$  ▷ Update the overall values.
9:     if  $x_{\hat{r}_j} = 1$  then ▷  $\hat{r}_j$  scheduled with hidden cache before.
10:       $y_{\hat{r}_j} \leftarrow 0$  ▷ Schedule  $\hat{r}_j$  with KV cache now.
11:     else ▷ First-time schedule for  $\hat{r}_j$ .
12:       if  $\Delta\hat{m}_j \neq m_{\hat{r}_j}$  then ▷ Hidden cache allowed for  $\hat{r}_j$ .
13:          $y_{\hat{r}_j} \leftarrow 1$  ▷ Schedule  $\hat{r}_j$  with hidden cache.
14:       end if
15:     end if
16:      $x_{\hat{r}_j} \leftarrow 1$  ▷ Schedule  $\hat{r}_j$ .
17:   else ▷ Keep adding would violate memory constraint  $M$ .
18:      $V' \leftarrow \hat{\theta}_j \cdot \Delta\hat{m}_j$ 
19:      $S' \leftarrow \{(x_i, y_i) | \forall i \in U, (x_i, y_i) = (0, 0) \text{ except for } i = \hat{r}_j \text{ where } (x_{\hat{r}_j}, y_{\hat{r}_j}) = (1, 0)\}$ 
20:     Break
21:   end if
22: end for
23: if  $V < V'$  then ▷ Double check to prevent the extreme case.
24:    $S \leftarrow S'$ 
25: end if
```

associated with each candidate request can be pre-recomputed. Then, we can derive the candidate schedule set Υ described as follow:

$$\Upsilon = \{(\theta_j, r_j, \Delta m_j, m_{r_j}) \mid j = 1, 2, \dots, n\}. \quad (1)$$

For the j -th candidate schedule in the set Υ , θ_j represents its marginal gain in value. r_j corresponds to the request associated with the j -th candidate schedule. Δm_j indicates the marginal increase in memory usage required by this schedule, while m_{r_j} denotes the maximum memory required by request r_j for the j -th candidate schedule. With the candidate schedule set Υ derived, together with the candidate request set U and the memory constraint M , it is fed to the greedy-based scheduling process outlined in Algorithm 1 as input, to derive the final scheduling decisions S as follow:

$$S = \{(x_i, y_i) \mid i \in U\}. \quad (2)$$

The complexity of Algorithm 1 is $\mathcal{O}(n \log n)$, resulting from sorting the candidate schedule set Υ at the beginning (Line 4). When increasing memory usage for a certain request \hat{r}_j causes a memory constraint violation, the resultant gain towards the overall scheduling value V' is still recorded (Lines 18-19) and compared with the existing overall scheduling value V to derive the final scheduling decisions (Lines 23-25). This double-check procedure is designed to prevent the extreme case where a request with a large absolute value but a small marginal gain in value is omitted during scheduling, resulting in an overall scheduling value far from optimal.

2 Approximation Guarantee

We further develop the approximation guarantee for the proposed greedy-based solution (Algorithm 1), which is summarized in Theorem 1.

Theorem 1. *For the hybrid-cache-based scheduling problem, let Φ_{ALG} denote the overall scheduling value produced by the greedy solution (Algorithm 1), and Φ_{OPT} denote the overall scheduling value produced by the optimal solution. We have:*

$$\frac{\Phi_{OPT}}{\Phi_{ALG}} \leq 2. \quad (3)$$

Proof. Consider an LP-relaxed variation of the original hybrid-cache-based scheduling problem, where $0 \leq x_i \leq 1$. In Algorithm 1, combining the first solution S with a simple tweak to the alternative solution S' — by setting the scheduling decision variable x'_{r_l} of the only scheduled request r_l to $\frac{M - \text{sum}(\bar{M})}{\Delta \bar{m}_{r_l}}$ — can result in the optimal solution for the LP-relaxed variation of the original hybrid-cache-based scheduling problem. This is because increasing the memory usage for the request with the highest marginal gain in value continuously until all the available memory is used up (achievable due to $0 \leq x_i \leq 1$), ensures that every unit of additional memory usage contributes the maximum possible value gain, thus resulting in the optimal solution Φ_{OPTLP} for the LP-relaxed variant of the problem (which can be proved by contradiction). Consequently, $V + V' = V + \hat{\theta}_{r_l} \cdot \Delta \hat{m}_{r_l} > V + \hat{\theta}_{r_l} \cdot \Delta \hat{m}_{r_l} \cdot \frac{M - \text{sum}(\bar{M})}{\Delta \bar{m}_{r_l}} = \Phi_{OPTLP}$.

Then, since Algorithm 1 essentially chooses the solution corresponding to the larger value between V and V' , it can be deduced that the final overall scheduling value Φ_{ALG} obtained by Algorithm 1 satisfies $\Phi_{ALG} = \max(V, V') \geq \frac{V + V'}{2} \geq \frac{\Phi_{OPTLP}}{2}$, which can be easily proved by contradiction. Finally, since the solution space of the hybrid-cache-based scheduling problem is a subset of that for its LP-relaxed variant with $0 \leq x_i \leq 1$, we have $\Phi_{OPTLP} \geq \Phi_{OPT}$, thus $\Phi_{ALG} \geq \frac{\Phi_{OPT}}{2}$, which ends the proof. \square