

1 Appendix

Algorithm 1 Layer Redistribution with Migration Coordinator

```
1: Input: Current layer assignment  $L_{curr}$ , Adjusted layer assignment  $L_{adj}$ 
2: Output: Efficient migration with minimal inference disruption
3: Initialization Phase:
4:   Pre-allocate memory for potential reassigned layers
5:   Preload weights for all candidate layers
6: On Scheduler Trigger of Redistribution:
7:    $M \leftarrow \text{IdentifyMigratedLayers}(L_{curr}, L_{adj})$ 
8:   for each layer  $l$  in  $M$  do
9:     if  $l$  is to be sent then
10:      Wait for computation of  $l$  to complete on source stage
11:      Asynchronously send KV cache of  $l$  to target stage
12:     else if  $l$  is to be received then
13:      Asynchronously receive KV cache
14:     end if
15:   end for
16: During Inference Execution:
17:   for each layer  $l$  assigned to target stage do
18:     if  $l \in M$  then
19:       Wait until KV cache for  $l$  is received
20:     end if
21:     Execute forward computation for  $l$ 
22:   end for
```
