

569 **A Hyper-Parameter Search**

570 In this section, we showcase the hyper-parameter grid search we performed for **layer stacking** and
 571 **selective backpropagation**. For **layer dropping**, we include the search results in Figure 5a.

572 **A.1 Layer stacking: When To Stack**

573 Figure 7 shows that **layer stacking** is relatively insensitive to different stacking RST hour times
 574 $\{(2.4, 7.2), (3, 7.2), (7.2, 9.6)\}$, with $\{(3, 7.2), (7.2, 9.6)\}$ performing about the same and $(2.4, 7.2)$
 575 slightly worse. We choose $(3, 7.2)$, as it matches the same $\frac{\text{stacking step}}{\text{all training steps}}$ ratio as proposed by Gong
 576 et al. [27].

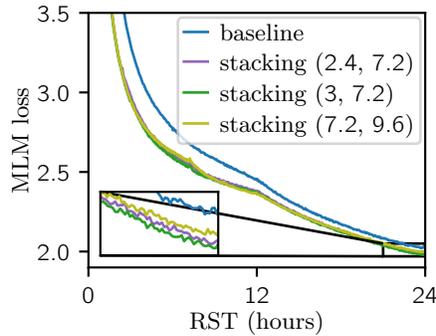


Figure 7: **Layer stacking** grid search: we tune the times at which the model is doubled. “stacking (a, b) ” indicates that the model was doubled in size once at a hours and then again at b hours, measured using RST.

577 **A.2 Selective backpropagation: Selectivity Scale**

578 We tune the selectivity scale β , where $\beta \in \{1, 2, 3\}$. Jiang et al. [35] use 33% and 50% selectivity in
 579 their experiments, which approximately corresponds to $\beta = \{1, 2\}$, respectively. We find that the
 580 larger the β value, the worse the pre-training performance. Note that the higher the β value, the more
 581 forward passes **selective backpropagation** needs to perform in order to collect enough samples for a
 582 backward pass, which decreases the total number of parameter update steps within the RST budget.
 583 For the experiments in Section 4.4, we chose $\beta = 1$, as it consistently achieves the best performance.

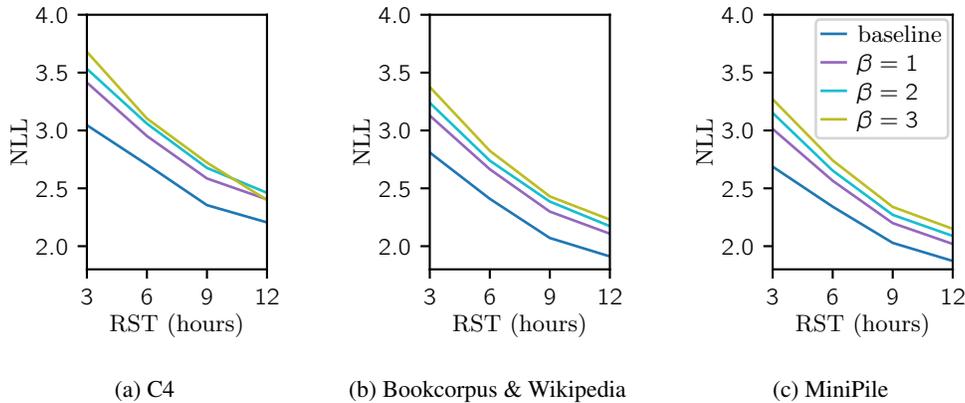


Figure 8: **Selective backpropagation** grid search: we tune the β hyperparameter. Each plot shows the validation loss over time during training for the given dataset.