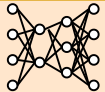


## P-pruning

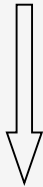
Pre-trained  
Language Model



Sampled  
Task Dataset



FLOPs  
Constraint



Compressed Model

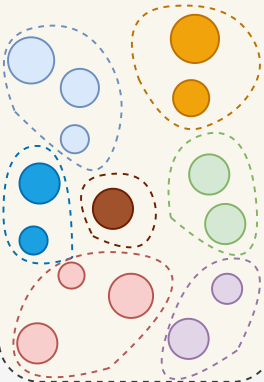


Fine-tuned Model



## Module Clustering

neuron clustering



dimension mapping

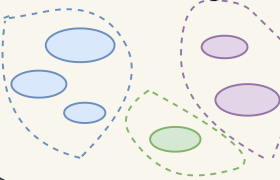
head  $i$



head  $j$

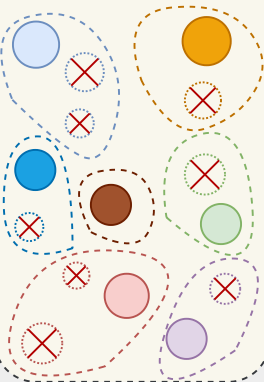


head clustering



## Centroid Selection

neuron selection



head selection

