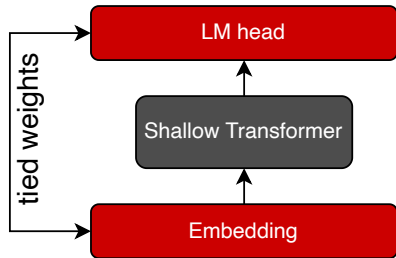


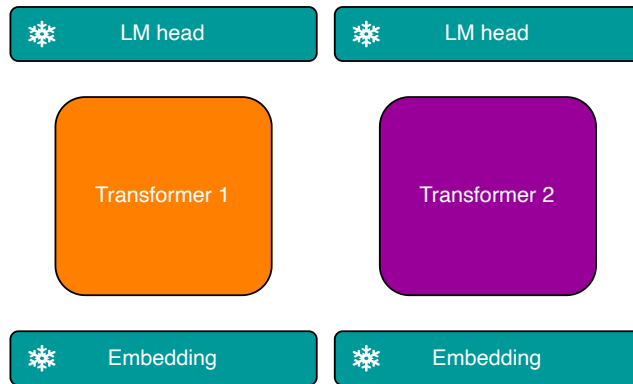
### STEP 1:

Train Embedding & LM head (tied weights) using shallow transformer



### STEP 2:

Asynchronously train models using these Embeddings



### STEP 3:

Stack the models  
(w/ or w/o common training)

