STEP 1:
Train Embedding & LM head (tied weights) using shallow transformer
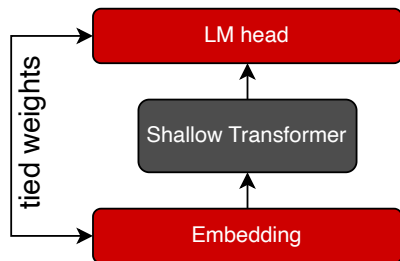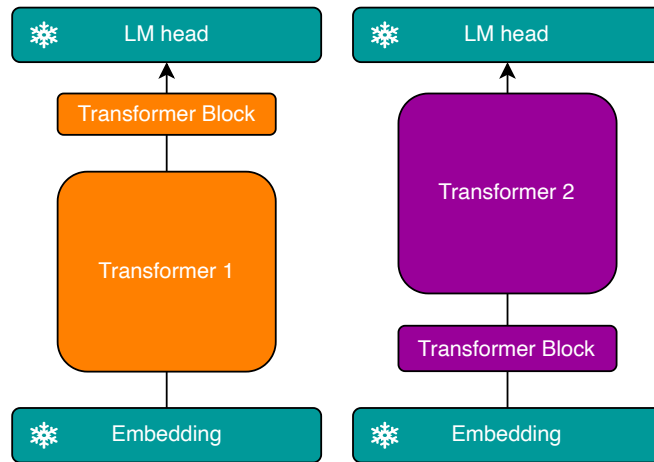
LM head
tied weights
Shallow Transformer
Embedding

STEP 2:
Train models asynchronously on trained Embedding weights.
Treat last layer of model 1 and first layer of model 2 seperately.

LM head
Transformer Block
Transformer 1
Embedding

LM head
Transformer 2
Transformer Block
Embedding

STEP 3:
Stack the models.
Throw away first layer of model 2.
Use last layer of model 1 only for decoding early.

LM head
Transformer 2
LM head
Transformer Block
Transformer 1
Embedding