

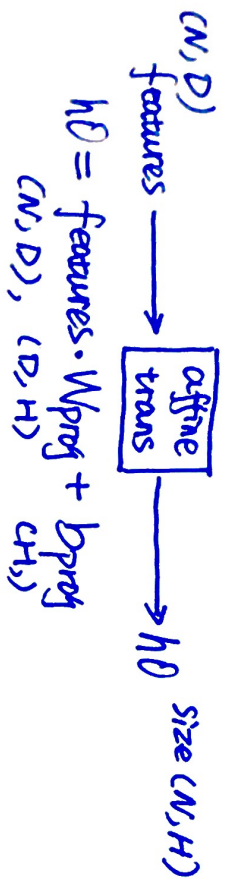
# Computational process for 'Word Embedding':

- given images input 'features':  $(N, D)$   $N$ : number of images
- 'captions' with size  $(N, T)$ . Each row here is a vector, each element in a vector is a number, correspond to a word using (idx-to-word).

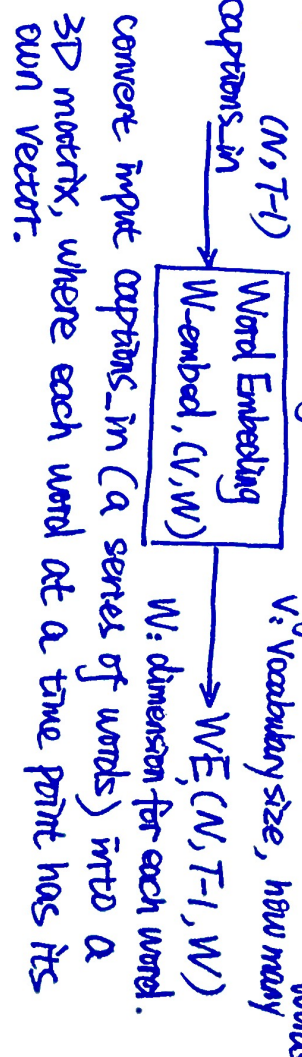
$T$  means there  $T$  number total of words in each word vector.

- captions\_in = captions[:, :-1], size  $(N, T-1)$
- captions\_out = captions[:, 1:] , size  $(N, T-1)$
- mask = captions\_out != self.\_null

(1) • Transfer input features into hidden inputs.

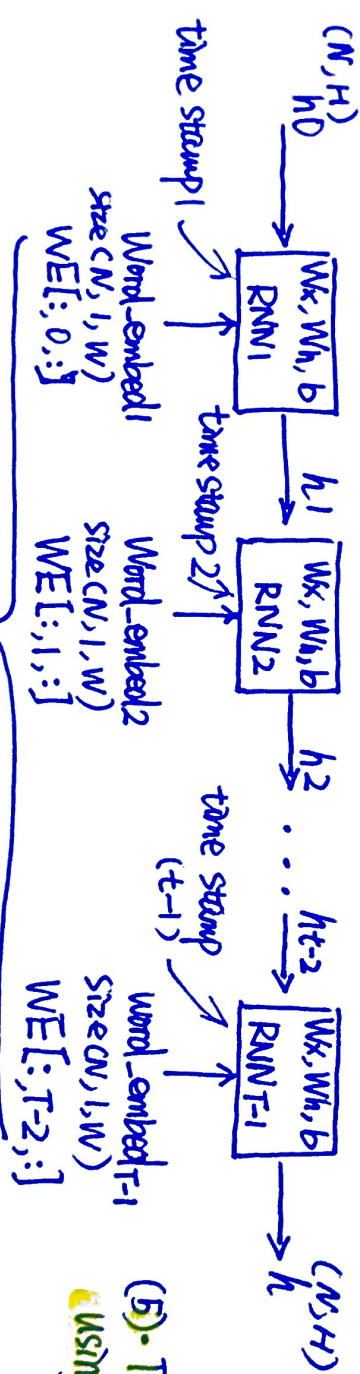


(2) • Use Word embedding to transfer captions\_in

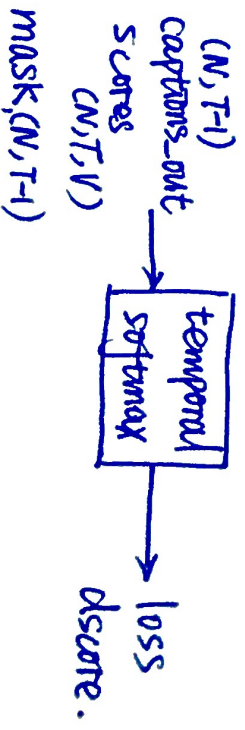


(3) • Recurrent Neural Network (RNN) process

- use either Vanilla RNN or '1st in RNN',  $W_k, W_h, b$  might have different size, depends which RNN.



(5) • Temporal Softmax to compute loss using captions\_out



(4) • Temporal affine transform to compute scores of all words at each time stamp.

