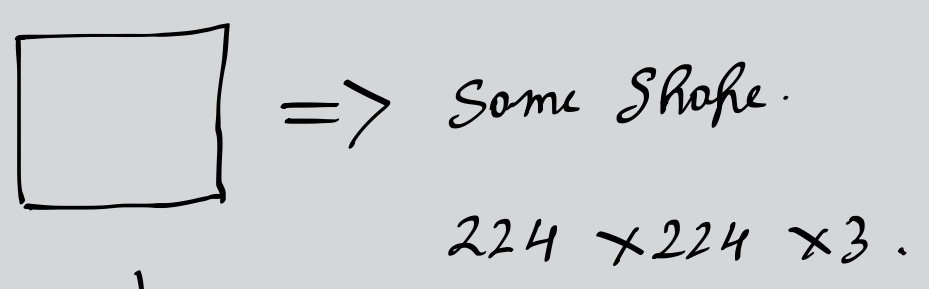


Image



Reshaping.

Batch Size - 32/64.

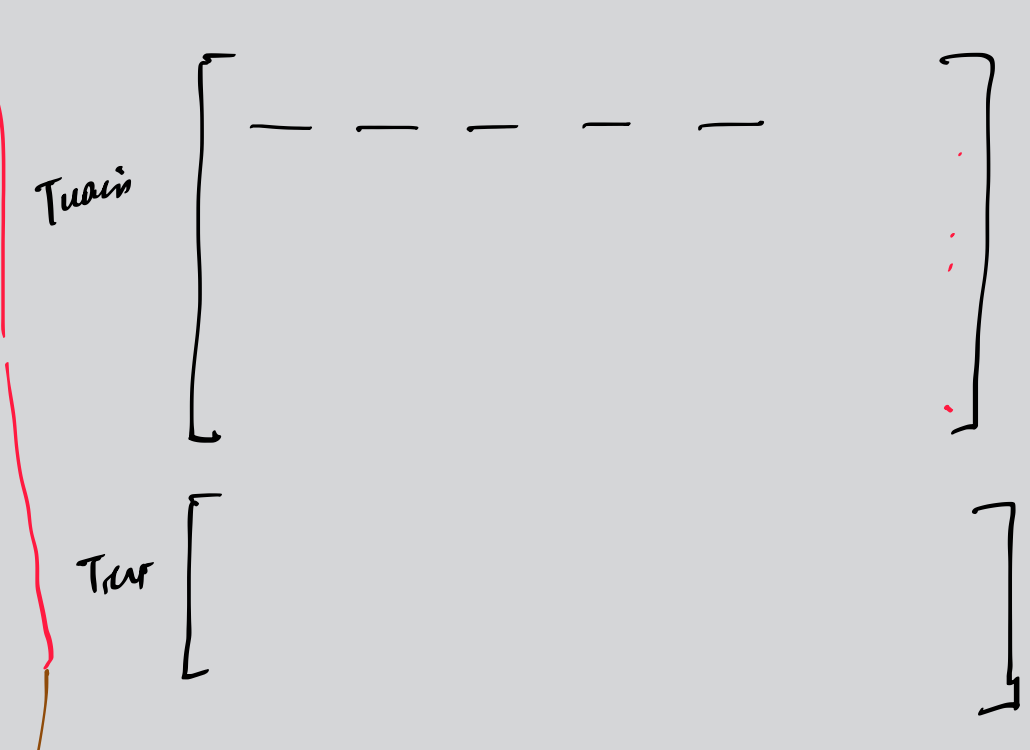
$32 \times 224 \times 224 \times 3$

$32 \times 3 \times 224 \times 224$   
 $B \times C \times H \times W$

Normalization.

Zero Mean | Unit Variance.

Text



1. In Train Dataset.
2. Iterate over the dataset.
3. Tokenize using Spacy Tokenizer.  
I am a student.  $\leftarrow$  lowercase  
 $\Rightarrow ["I", "am", "a", "student"]$
4. Vocab list = []  $\leftarrow$  Initially Empty.

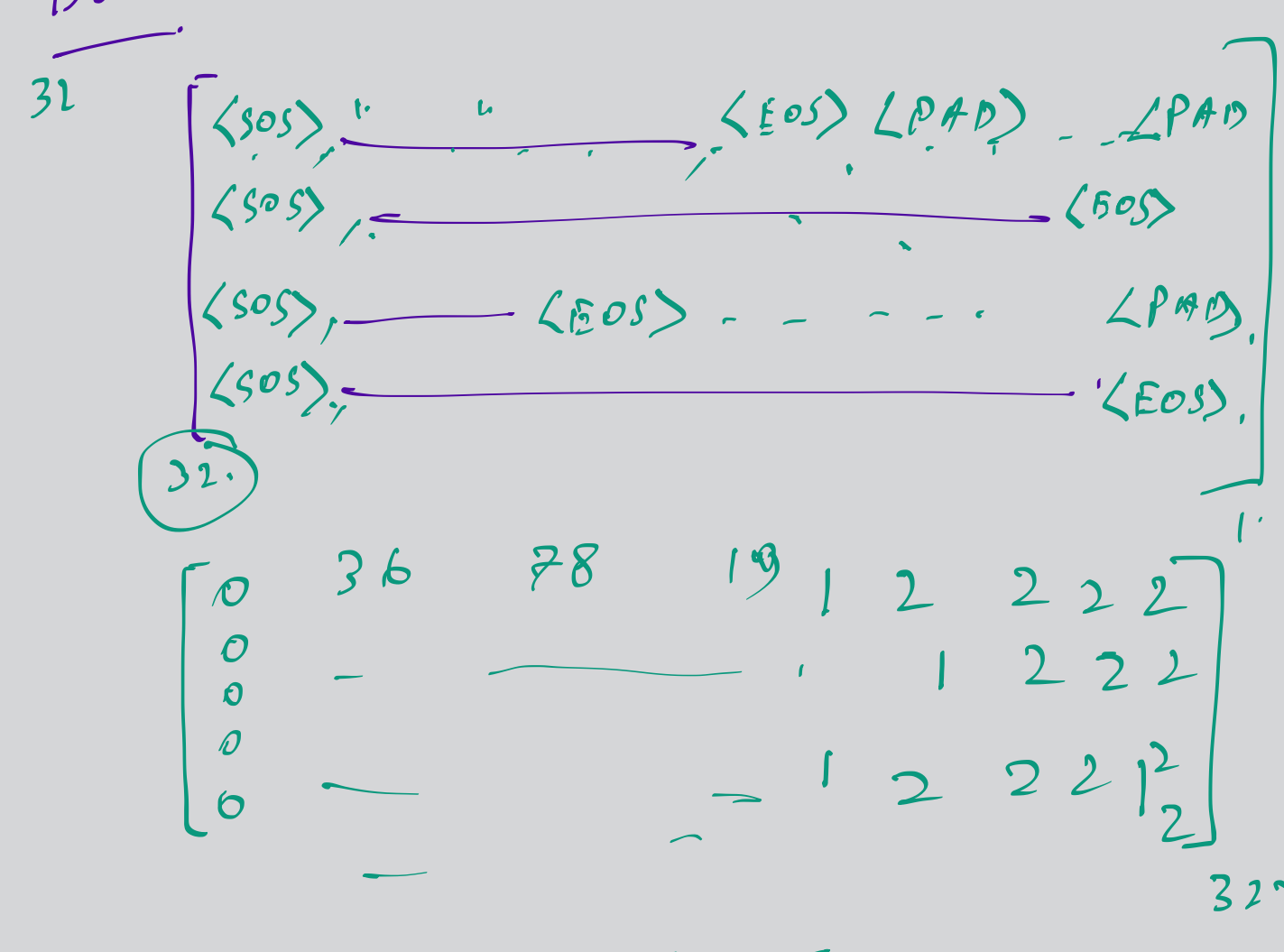
```
vocab_list = []
for sent in train - test:
    sent = sent.lower()
    sent = sent.tokenize() ["I", "am", ...]
    vocab_list.extend(sent)
```

```
vocab_list = set(vocab_list).
vocab_list = list(vocab_list).
vocab_list.extend(["<SOS>", "<EOS>", "<PAD>", "<UNK>"])
```

Key	Value
<SOS>	0
<EOS>	1
<PAD>	2
<UNK>	3
...	...
:	:

```
word2int = {}
for i, w in enumerate(vocab_list):
    word2int[w] = i
int2word = {}  $\leftarrow$  Reverse of word2int
```

Batch



1. max length of sentence in a batch.  
 $= 15, 20$

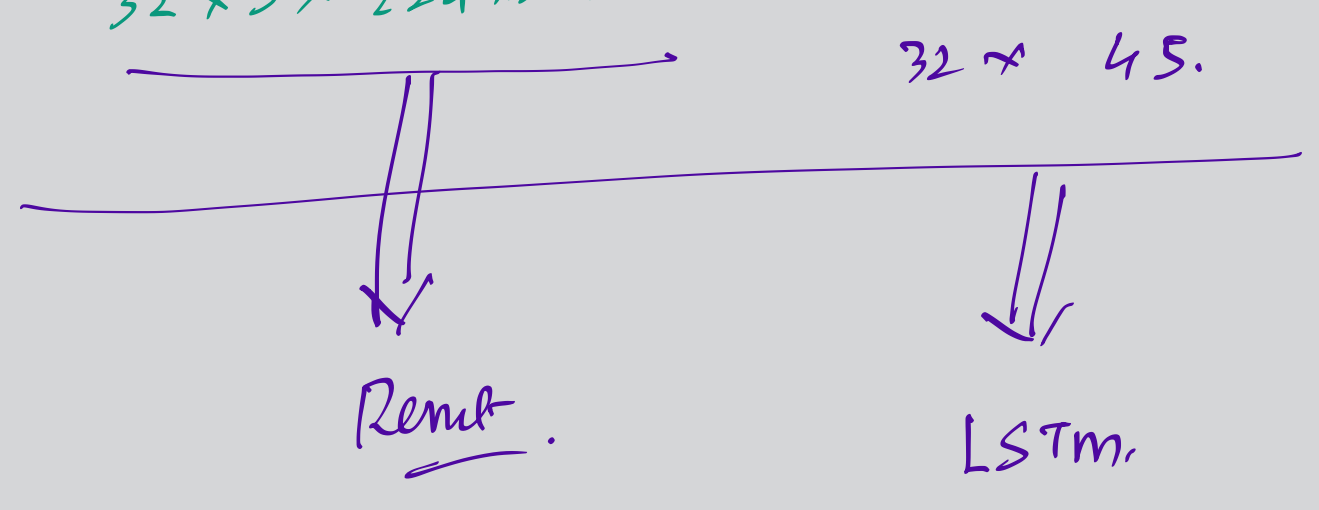
Target Tensor

Src Tensor

$32 \times 3 \times 224 \times 224$

Tgt Tensor

$32 \times 20$   
 $32 \times 45$



Img  $\rightarrow$  Resnet  $\rightarrow$  512 dim.

