#### Attention based models in End-to-End ASR

Exploration of Attention in ESPNET toolkit

Shreekantha Nadig

November 9, 2018

International Institute of Information Technology - Bangalore

#### Table of contents



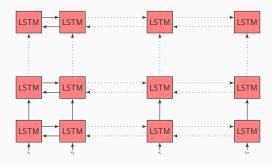
- 1. Introduction
- 2. No Attention [Equal Attention?]
- 3. Dot product Attention
- 4. Additive Attention
- 5. Location Aware Attention
- 6. 2D Location Aware Attention
- 7. Location Aware Recurrent Attention
- 8. Coverage mechanism Attention
- 9. Coverage mechanism location aware Attention
- 10. Multi-Head Attention
- 11 Multi Head Dot Product Attention

# Introduction

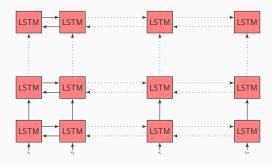




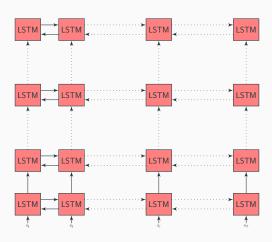




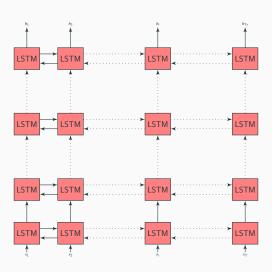




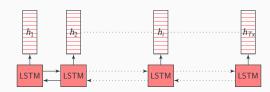






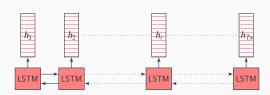




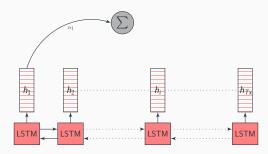




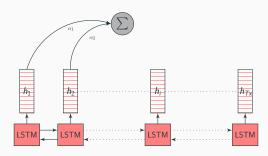




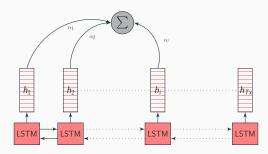




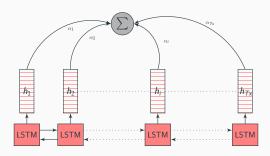




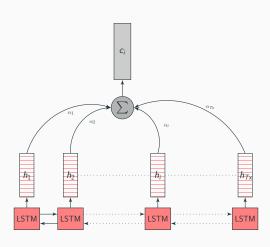




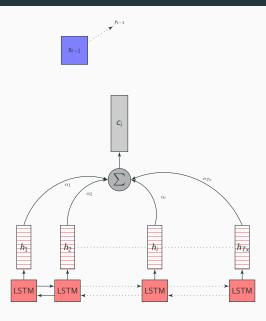




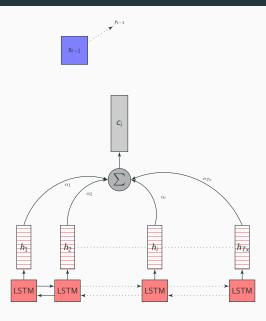




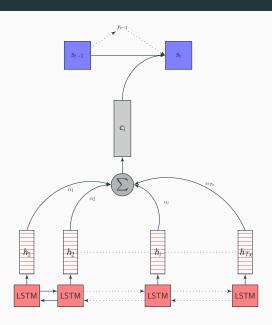




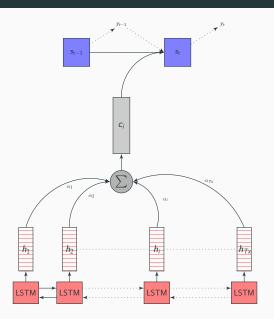








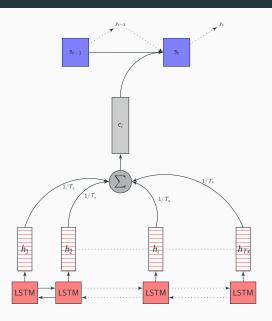




No Attention [Equal Attention?]

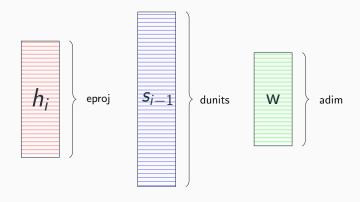
## No Attention [Equal Attention?]





## **Dimensions of representations**

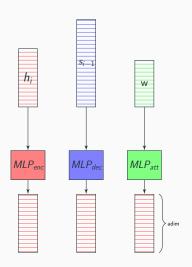




Mostly  $eproj \neq dunits \neq adim$ 

## Matching the dimensions of representations

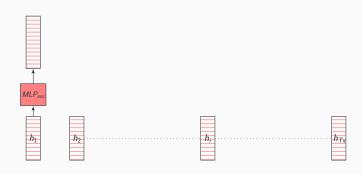




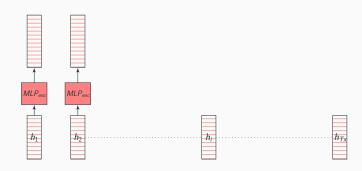


			_
- la	- La	- L	_
111	112	h;	+7
			_

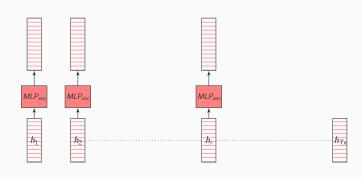




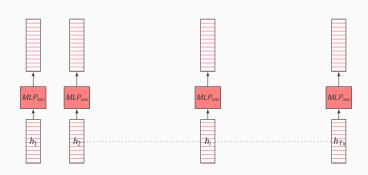




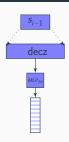


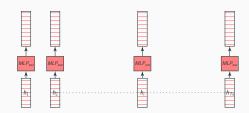




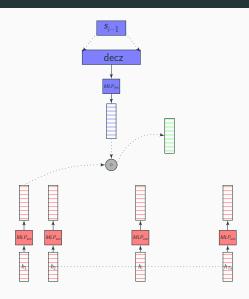




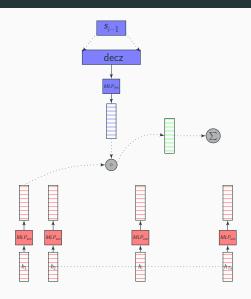




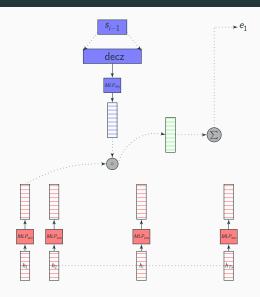




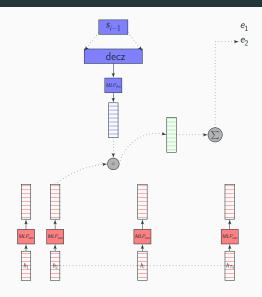




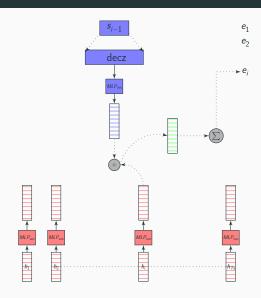




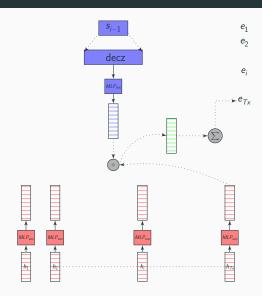






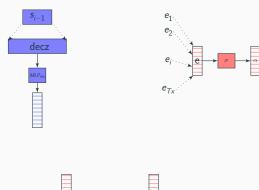


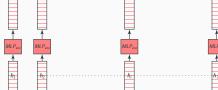




### **Dot product Attention**

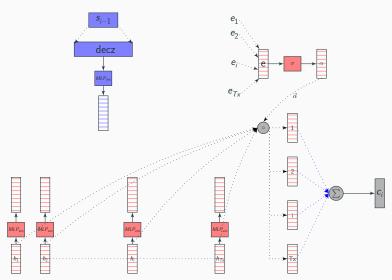






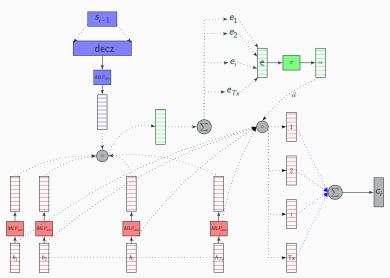
### **Dot product Attention**



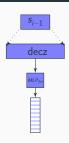


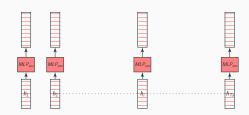
## **Dot product Attention - Full picture**



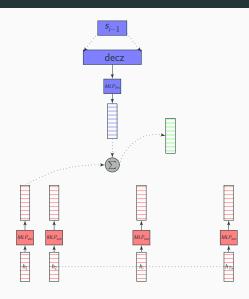




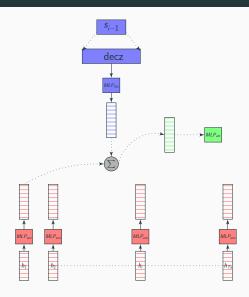




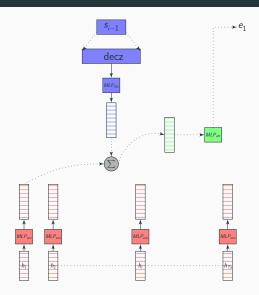




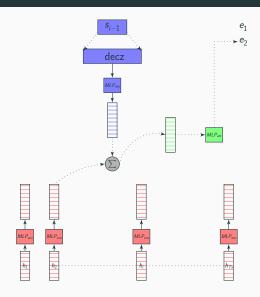




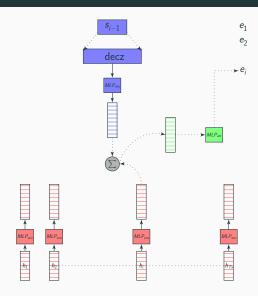




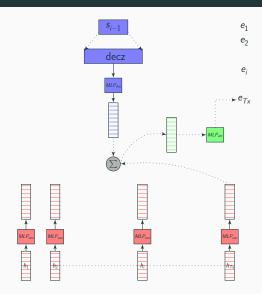




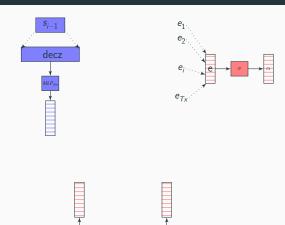




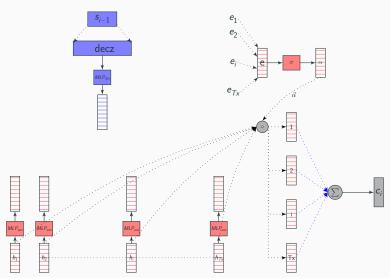






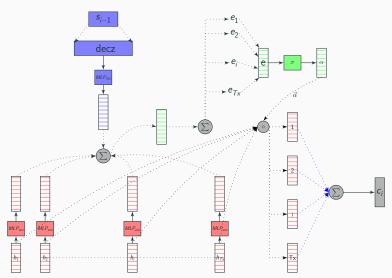




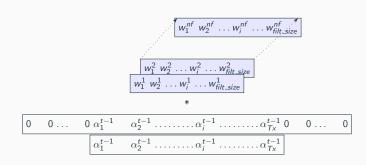


### Additive Attention - Full picture









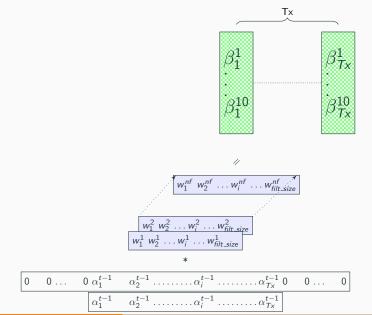


$$\begin{bmatrix} \beta_1^1 & \beta_2^1 & \dots & \beta_{i}^1 & \dots & \beta_{T_X}^1 \\ & & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\$$



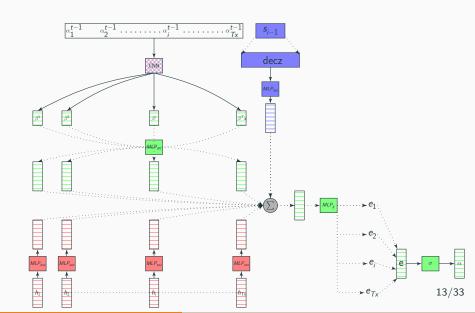
$$\begin{bmatrix} \beta_1^2 & \beta_2^2 & \dots & \beta_i^2 & \dots & \beta_{Tx}^2 \\ & & & & & & & \\ w_1^2 & w_2^2 & \dots & w_i^2 & \dots & w_{filt\_size}^2 \\ & & & & * \\ \\ \hline 0 & 0 & \dots & 0 & \alpha_1^{t-1} & \alpha_2^{t-1} & \dots & \alpha_i^{t-1} & \dots & \alpha_{Tx}^{t-1} & 0 & 0 & \dots & 0 \\ \end{bmatrix}$$



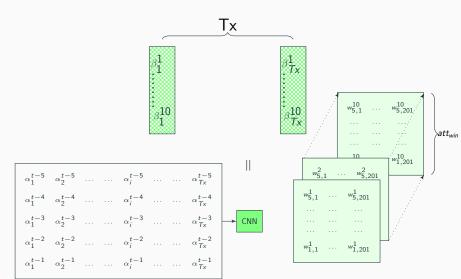


#### **Location Aware Attention - Full picture**



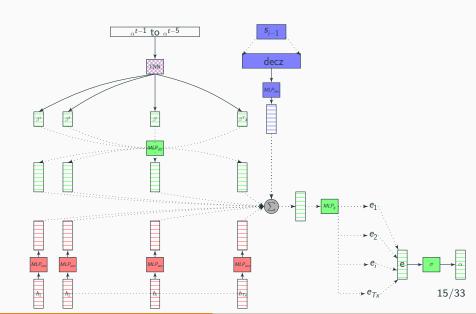






#### 2D Location Aware Attention - Full picture

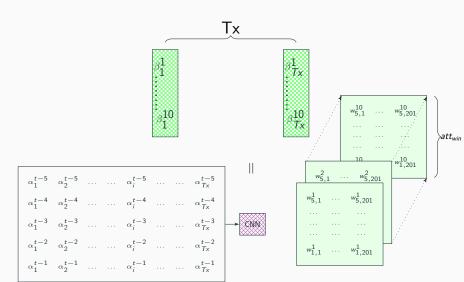




**Location Aware Recurrent Attention** 

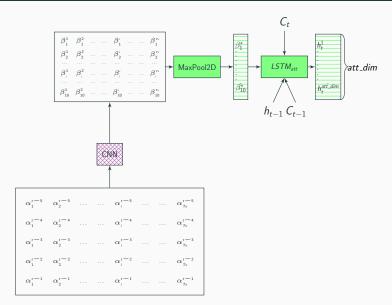
#### **Location Aware Recurrent Attention**





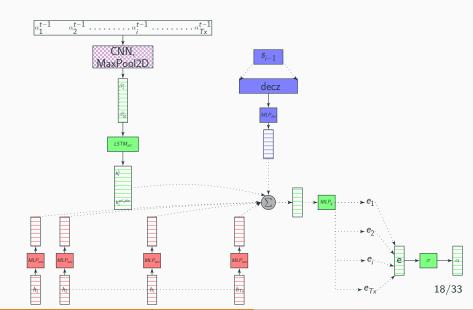
#### **Location Aware Recurrent Attention - weights**





#### **Location Aware Recurrent Attention - Full picture**

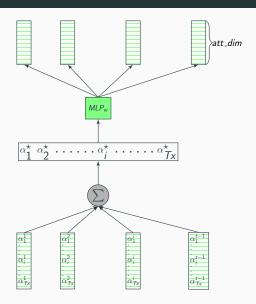




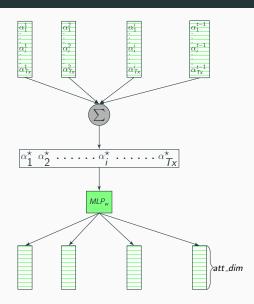


- Text summarization Seq-to-Seq models
- Not reliable in producing factual details correctly
- Extend the standard seq-to-seq attention models
  - Hybrid pointer-generator network
  - Coverage

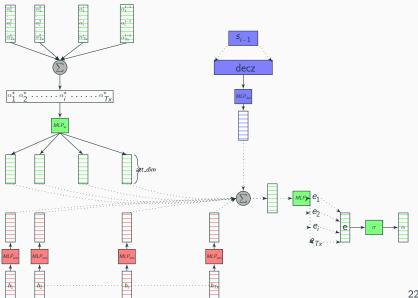




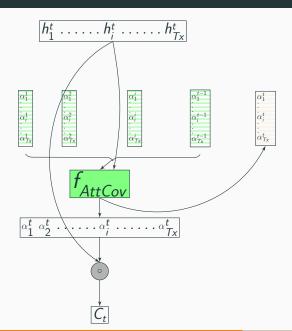










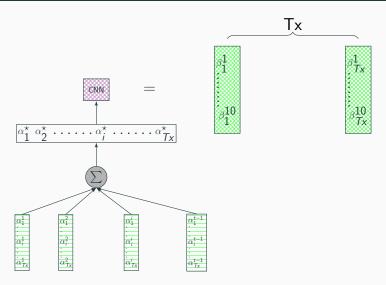


Coverage mechanism location aware

**Attention** 

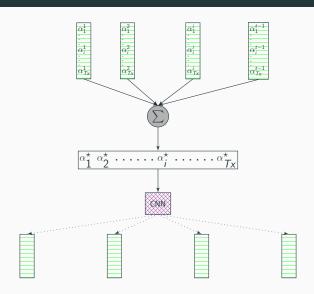
#### Coverage mechanism location aware Attention





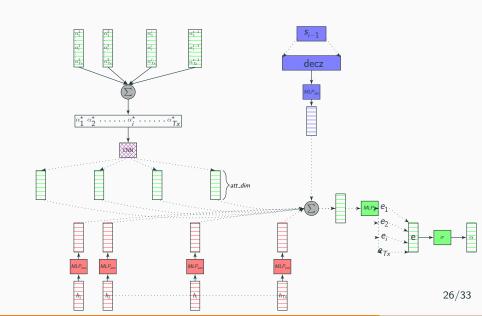
#### Coverage mechanism location aware Attention





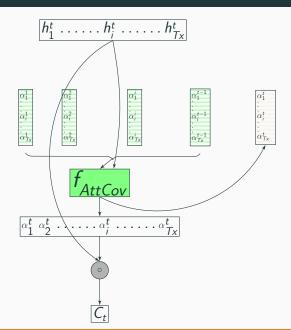
### Coverage mechanism location aware Attention





### Coverage mechanism location aware Attention



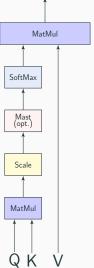


**Multi-Head Attention** 

## Scaled Dot product

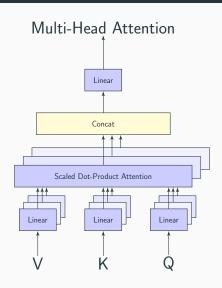






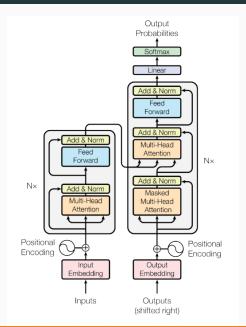
### **Multi-Head Attention**





### The Transformer - model architecture





## **Positional Encoding - Example**





$$PE_{(pos,2_i)} = sin(pos/1000^{2i/d_{model}})$$

$$PE_{(pos,2_{i+1})} = cos(pos/1000^{2i/d_{model}})$$

### **Pending Discussion**



- Pointer Generator Attention network
- Multi Head location based Attention
- Multi Head multi resolution location based Attention

# **Questions?**

Multi Head Dot Product Attention

Multi Head location based Attention

Multi Head multi resolution location

based Attention