

# A Generalized Language Model in Tensor Space



Lipeng Zhang<sup>1</sup>, Peng Zhang<sup>1,\*</sup>, Xindian Ma<sup>1</sup>, Shuqin Gu<sup>1</sup>, Zhan Su<sup>1</sup>, Dawei Song<sup>2</sup>

<sup>1</sup> College of Intelligence and Computing, Tianjin University, China

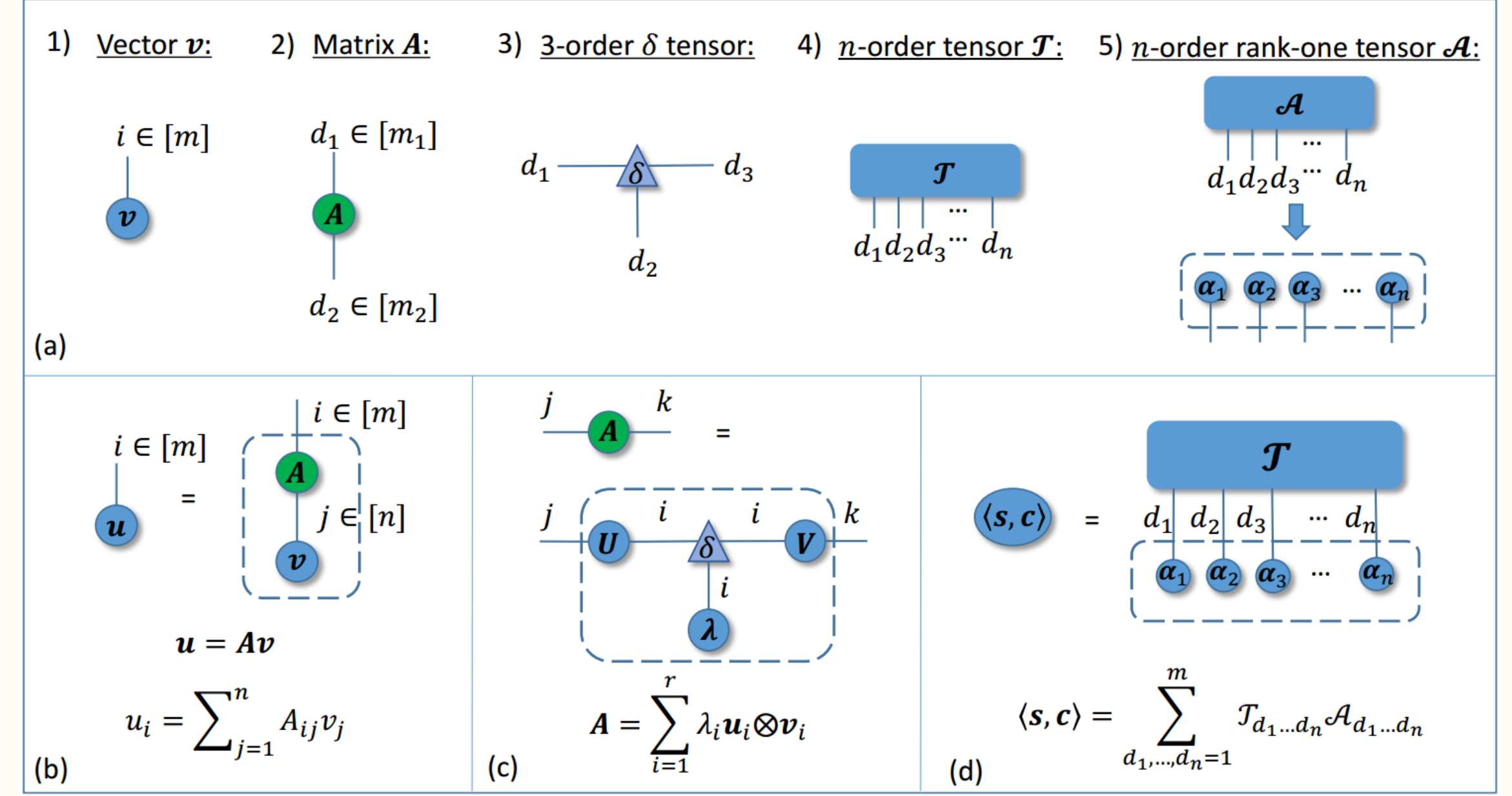
<sup>2</sup> School of Computer Science and Technology, Beijing Institute of Technology, China

## 1. Introduction

### Motivation

- ❑ The existing methods usually adopt relatively low-order tensors, which have limited expressive power in modeling language.
- ❑ We propose a language model based on relatively high-order tensor representation—Tensor Space Language Model (TSLM).
- ❑ Challenges
  - ◆ To derive an effective solution for such high order representation;
  - ◆ To demonstrate such a solution is a general approach for language modeling;
  - ◆ To solve that such a high-order tensor contains exponential magnitude of parameters.

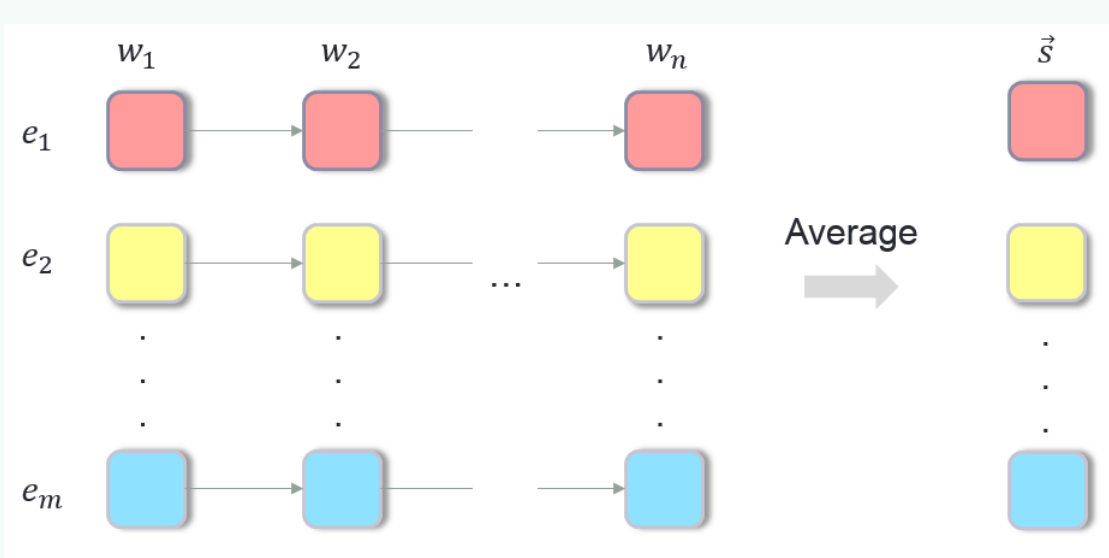
## 2. Tensor Network



## 3. Tensor Space Language Model and The Generalization

### TSLM Basic Representation

Hypotheses: A sentence has  $n$  words. Each word has  $m$  semantic meanings.



The sentence still has  $m$  meanings.

- ❑ How to represent a single word

$$\mathbf{w}_i = \sum_{d_i=1}^m \alpha_{i,d_i} \mathbf{e}_{d_i}$$

- ❑ How to represent a sentence

$$\mathbf{s} = \mathbf{w}_1 \otimes \cdots \otimes \mathbf{w}_n = \sum_{d_1, \dots, d_n=1}^m \mathcal{A}_{d_1 \dots d_n} \mathbf{e}_{d_1} \otimes \cdots \otimes \mathbf{e}_{d_n}$$

- ❑ Assume that each sentence  $s_i$  appears with a probability  $p_i$ . We can denote the corpus as:

$$\mathbf{c} = \sum_i p_i \mathbf{s}_i = \sum_{d_1, \dots, d_n=1}^m \mathcal{T}_{d_1 \dots d_n} \mathbf{e}_{d_1} \otimes \cdots \otimes \mathbf{e}_{d_n}$$

- ❑ The sentence probability:

$$p(\mathbf{s}) = \langle \mathbf{s}, \mathbf{c} \rangle = \sum_{d_1, \dots, d_n=1}^m \mathcal{T}_{d_1 \dots d_n} \mathcal{A}_{d_1 \dots d_n}$$

### The Generalization of N-Gram Language

- ❑ **Claim 1:**

In our TSLM, when we set the dimension of vector space  $m = |V|$  and each word  $w$  as an one-hot vector, the probability of sentence  $s$  consist of words  $d_1, \dots, d_n$  in vocabulary is the entry  $\mathcal{T}_{d_1 \dots d_n}$  of tensor  $\mathcal{T}$ .

- ❑ **Claim 2:**

In our TSLM, we define the word sequence  $w_1^i = (w_1, w_2, \dots, w_i)$  with length  $i$  as:

$$\mathbf{w}_1^i = \mathbf{w}_1 \otimes \cdots \otimes \mathbf{w}_i \otimes \mathbf{1}_{i+1} \otimes \cdots \otimes \mathbf{1}_n$$

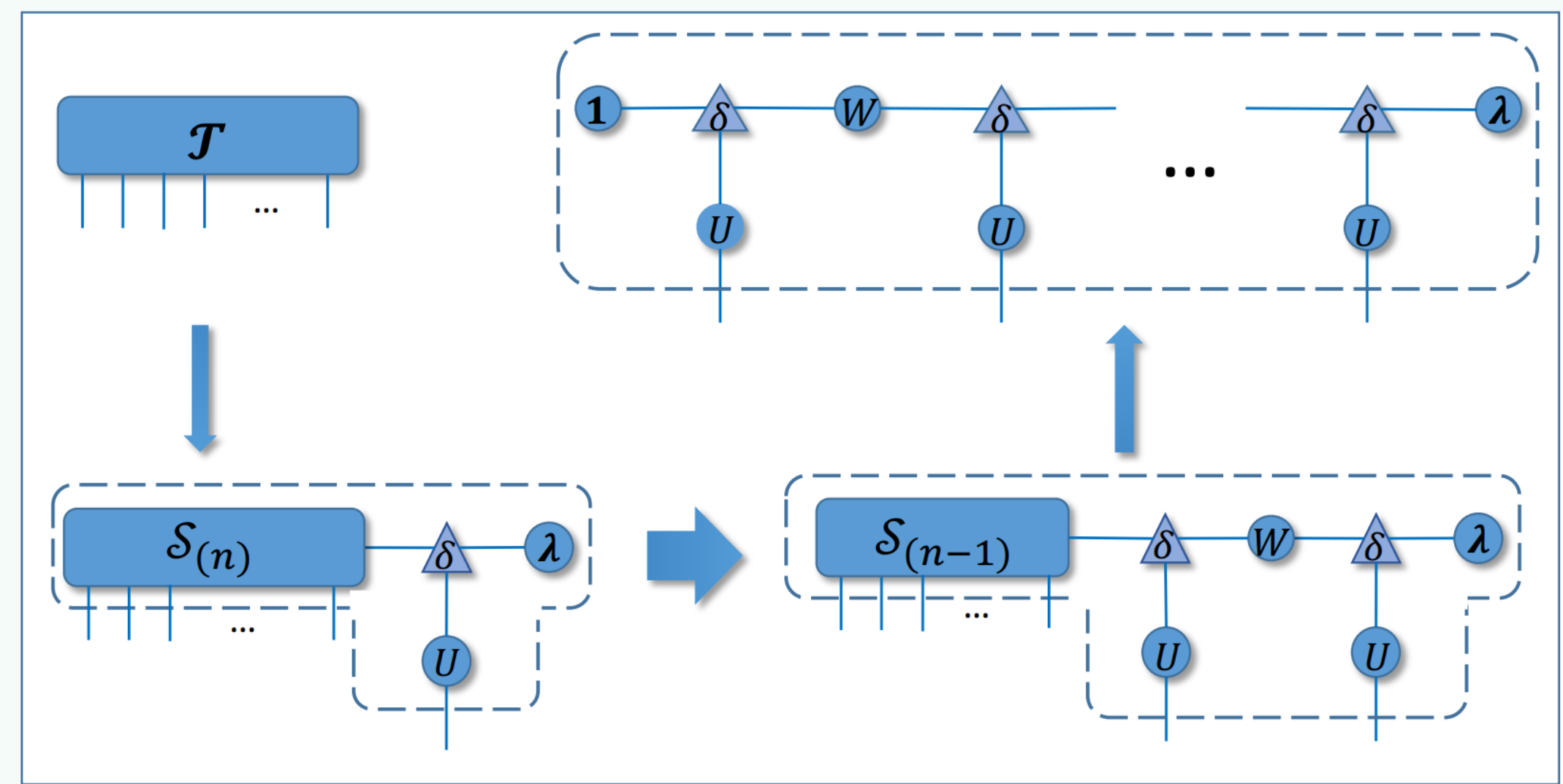
Which means that the sequence  $w_1^i$  is padded via using full one vector  $\mathbf{1}$ . Then, the probability  $p(w_1^i) = \langle \mathbf{w}_1^i, \mathbf{c} \rangle$ .

- ❑ In TSLM, the conditional probability  $p(w_i | w_1^{i-1})$  can be computed as:

$$p(w_i | w_1^{i-1}) = \frac{p(w_1^i)}{p(w_1^{i-1})} = \frac{\langle \mathbf{w}_1^i, \mathbf{c} \rangle}{\langle \mathbf{w}_1^{i-1}, \mathbf{c} \rangle}$$

## 4. Deriving Recursive Language Modeling Process from TSLM

### Recursive Tensor Decomposition



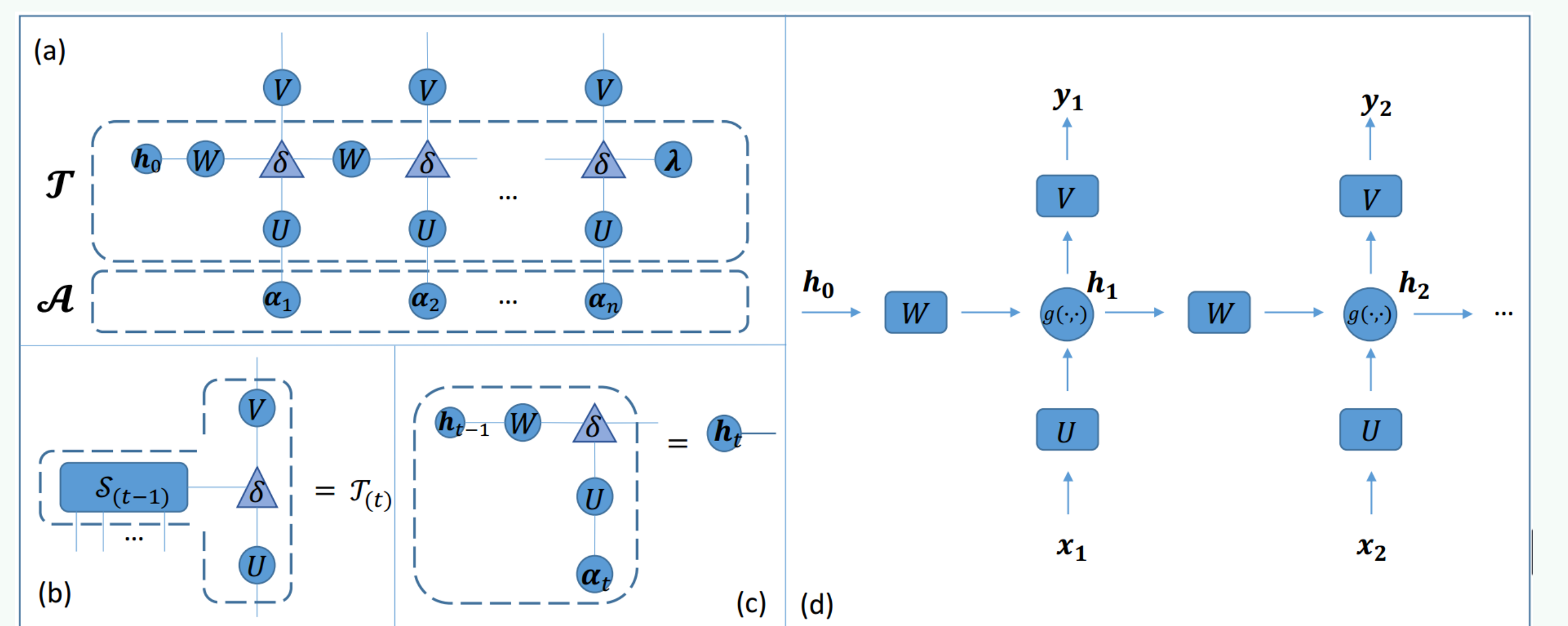
- ❑ Decomposition :

$$\mathcal{T} = \sum_{i=1}^r \lambda_i \mathcal{S}_{(n),i} \otimes \mathbf{u}_i$$

$$\mathcal{S}_{(n),k} = \sum_{i=1}^r W_{k,i} \mathcal{S}_{(n-1),i} \otimes \mathbf{u}_i$$

$$\dots$$

$$\mathcal{S}_{(1)} = \mathbf{1}$$



- ❑ Denote :  $\mathbf{h}_0 = W^{-1} \mathbf{1}$

- ❑ Computing  $\mathbf{h}_t$  recursively :

$$\mathbf{h}_t = W \mathbf{h}_{t-1} \odot U \alpha_t$$

- ❑ Constructing a tensor mapping to vocabulary by a matrix  $V \in \mathbb{R}^{r \times |V|}$ :

$$\mathcal{T}_{(t),k} = \sum_{i=1}^r V_{k,i} \mathcal{S}_{(t-1),i} \otimes \mathbf{u}_i$$

- ❑ Therefore, computing the conditional probability recursively:

$$p(w_t | w_1^{t-1}) = \text{softmax}(\mathbf{y}_t)$$

$$\mathbf{y}_t = V \mathbf{h}_t$$

$$\mathbf{h}_t = g(W \mathbf{h}_{t-1}, U \alpha_t)$$

$$g(\mathbf{a}, \mathbf{b}) = \mathbf{a} \odot \mathbf{b}$$

## 5. Empirical Evaluation and Conclusion

### Experiment Results

Model	PTB				WikiText-2			
	Hidden size	Layers	Valid	Test	Hidden size	Layers	Valid	Test
KN-5(Mikolov and Zweig 2012)	-	-	-	141.2	-	-	-	-
RNN(Mikolov and Zweig 2012)	300	1	-	124.7	-	-	-	-
LSTM(Zaremba, Sutskever, and Vinyals 2014)	200	2	120.7	114.5	-	-	-	-
LSTM(Grave, Joulin, and Usunier 2016)	1024	1	-	82.3	1024	1	-	99.3
LSTM(Merity et al. 2017)	650	2	84.4	80.6	650	2	108.7	100.9
RNN†	256	1	130.3	124.1	512	1	126.0	120.4
LSTM†	256	1	118.6	110.3	512	1	105.6	101.4
TSLM	256	1	<b>117.2</b>	<b>108.1</b>	512	1	<b>104.9</b>	<b>100.4</b>
RNN+MoS†(Yang et al. 2018)	256	1	88.7	84.3	512	1	85.6	81.8
TSLM+MoS	256	1	<b>86.4</b>	<b>83.6</b>	512	1	<b>83.9</b>	<b>81.0</b>

### Conclusion

- ❑ We propose a novel language model, named Tensor Space Language Model, aiming to consider high-order dependencies of words via tensors and tensor networks.
- ❑ We prove that TSLM is a generalization of the n-gram language model.
- ❑ We can derive a recursive calculation of conditional probability for language modeling via tensor decomposition in TSLM.