

# Automatic Text Summarization using closed-Pattern-Infused Edit Similarity (PIESim)

閉合模式融合於編輯相似度之自動文本摘要

---

台大統計碩士學位學程 陳育婷

指導教授：蔡政安 教授 許聞廉 教授

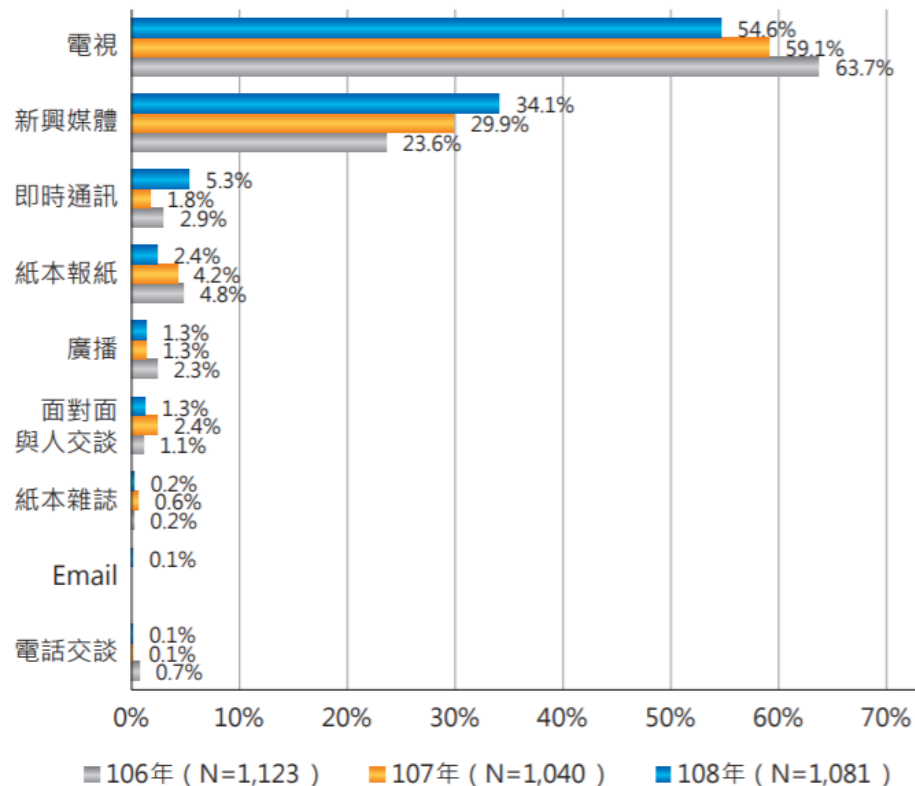
2020/ 07/ 27

# Table of Content

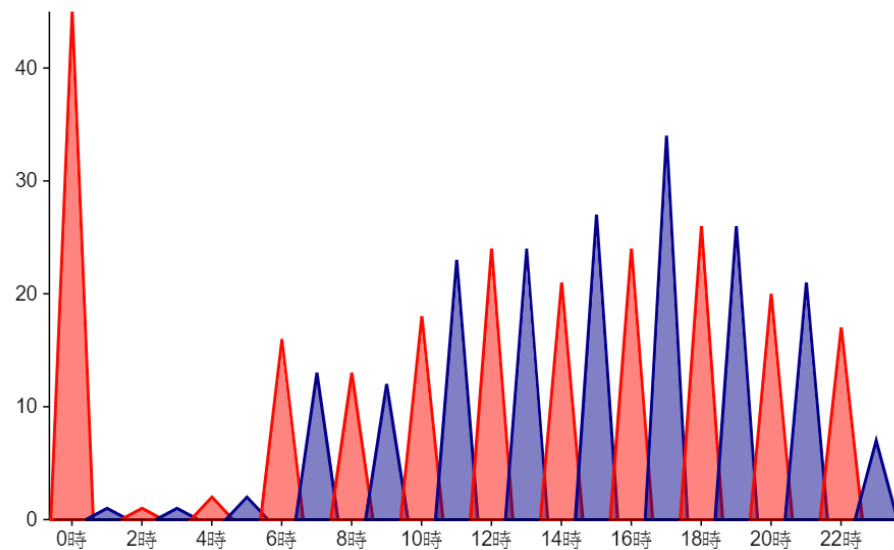
1. Introduction and Aims
2. Literature Review
3. Background Study
4. The Proposed PIESim Model
5. Experiments and Results
6. Different Settings Analysis
7. Error Analysis & User Interaction
8. Conclusion

# Introduction & Aims

主要透過哪個管道獲得新聞資訊



108/6/1- 7/31 National Communication Commission (NCC) Digital Convergence Survey with 95% C.I.



Number of streamed News on 2019/04/07 of Apple Daily News

- ✓ Preliminary Goal : A precise and revisable automatic text summarization (ATS) system can ease the burden of news publishers and benefit all kinds of human

# Introduction & Aims

- ✓ Preliminary Goal : design an extractive, precise, and revisable summarization system with compression techniques

## More Deeper Application of NLP primo.ai.com

Group 1	Group 2	Group 3
Cleanup, Tokenization	Information Retrieval and Extraction (IR)	Machine Translation
Stemming	Relationship Extraction	Automatic Summarization/ Paraphrasing
Lemmatization	Named Entity Recognition (NER)	Natural Language Generation
Part of Speech Tagging	Sentiment Analysis/Sentence Boundary Disambiguation	Reasoning over Knowledge Based
Query Expansion	World sense and Disambiguation	Quation Answering System
Parsing	Text Similarity	Dialog System
Topic Segmentationand Recognition	Coreference Resolution	Image Captioning & other Multimodel Tasks
Morphological Degmentation (Word/Sentences)	Discourse Analysis	



- abstract
- compress
- extract

### Sources of Summaries

- A. Aries and W. K. Hidouci, "Automatic text summarization: What has been done and what has to be done,"

# Introduction & Aims

What's  
Important  
?

No systems can cover all aspects. Considering from general and query aspects is more useful in real life

Incorporate Query Aspect

Use of  
Enormous  
Data

Most non DL or ML approaches do not use information outside the document

Structure can consider new info.

Vector-  
based  
Represen-  
-tation

Term-based : Lack Semantic Information  
Deep learning : Hard to revise  
Not conform to human process

Non-vector rep. with semantic info.

✓ Goal : design an extractive, precise, non-vector-based , revisable, query & new-info.-considered ATS system with compression techniques

# Introduction & Aims

precise

Performance is the **best** among large categories of ATS systems on **Chinese** and **English** dataset.

non-vector-based & revisable

**Non-vector based** but preserves **semantic information**; **thus** can be intuitively revised

query & new-info.-considered

Structure is **easy to absorb** new info. We propose a novel memory similarity criterion to learn from **training data**

innovative

The first to **combine pattern and edit distance** in ATS and to propose a new **PIESim** similarity measure

# Literature Review – Genres of ATS

Pros: exploit shared

information between units,  
more coherent

Cons: constructed from the  
target document only

Pros : More complex and  
automatic

Cons : Needs corpora to  
learn the rules, model itself is  
not easy to revise.

Stati-  
stical

Pros: simple, explanation and revision

Cons: many use word-independent  
features

Graph-  
based

Pros : usually more  
accurate

Lingu-  
-istic

Cons : restricted to  
languages and domains,  
more laboring work

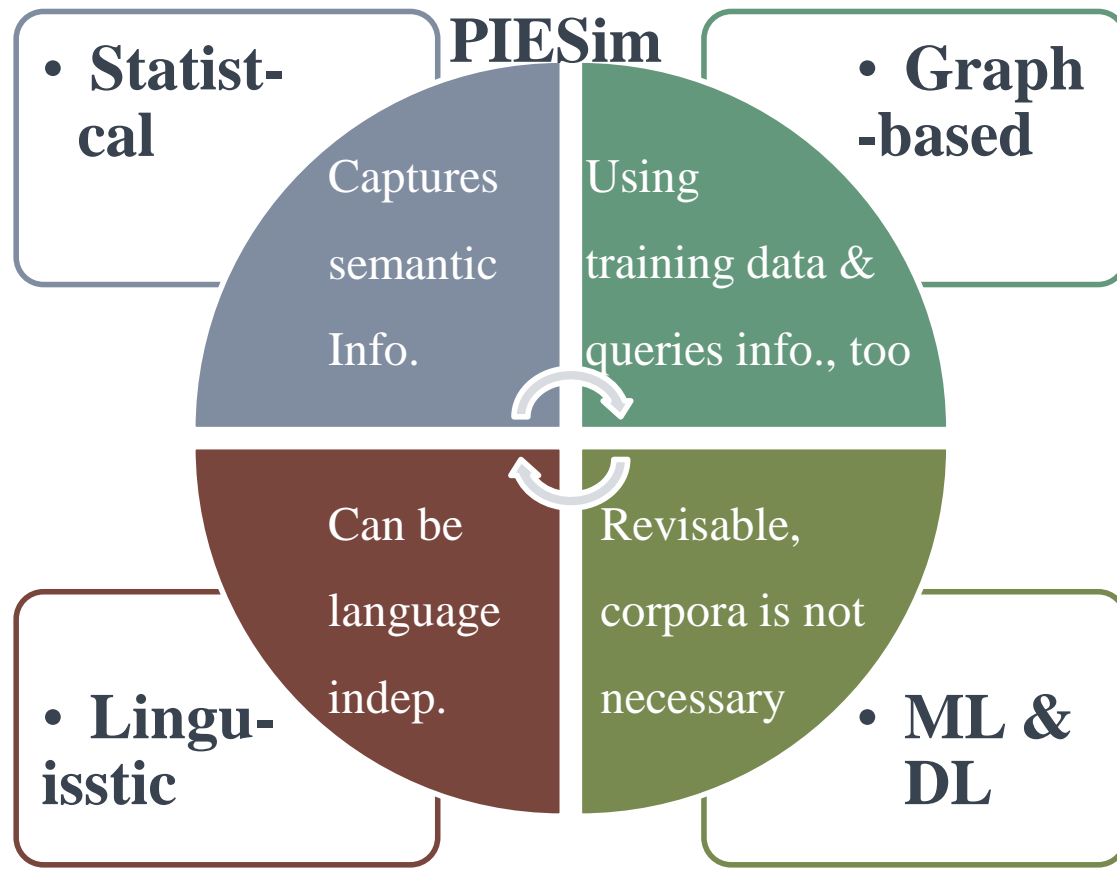
ML

Pros : Fully automatic. Perform well  
when corpora is enough

DL

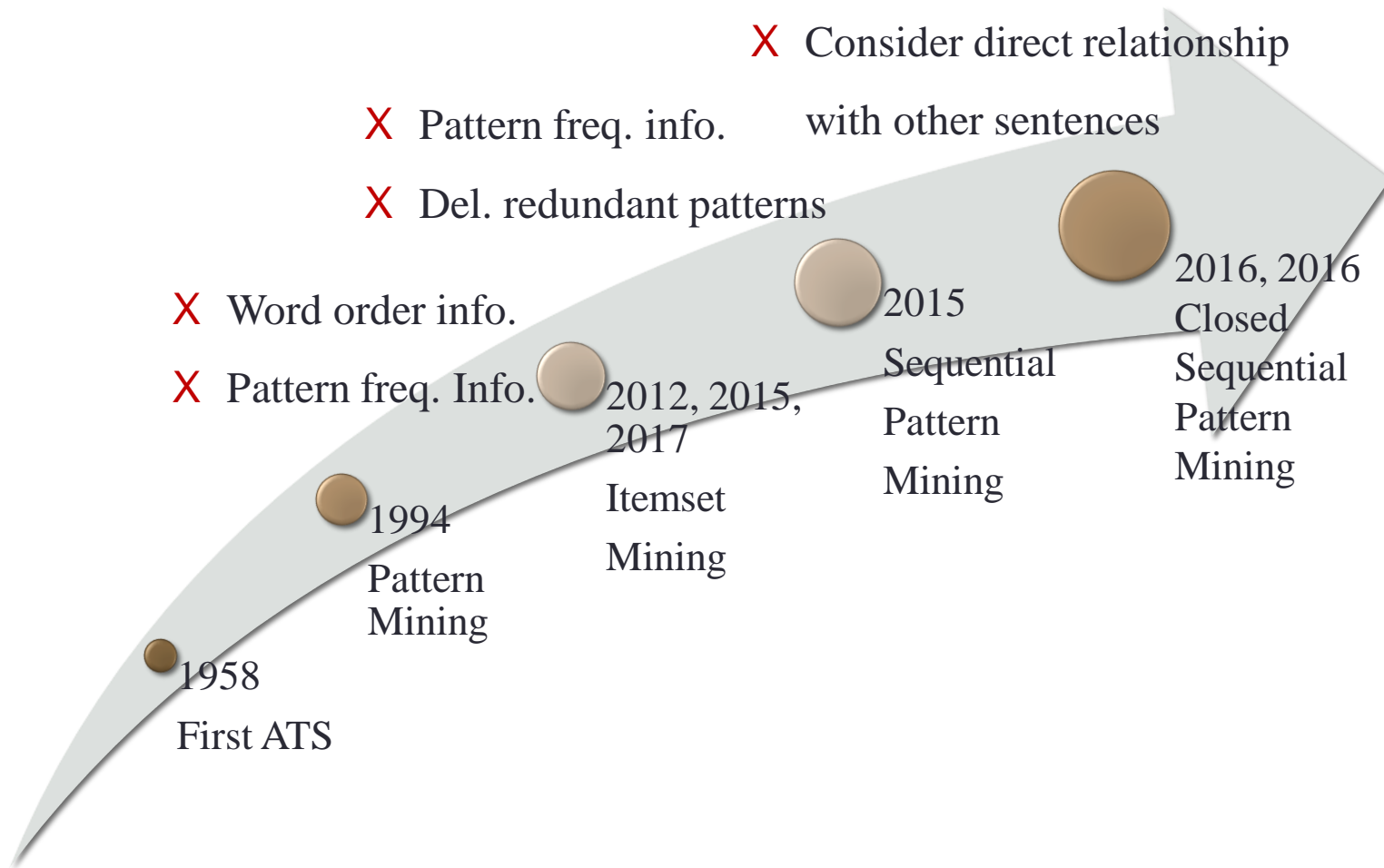
Cons : Explanation and revision  
itself is not easy to revise. Often non  
grammatically correct.

# Literature Review – Genres of ATS





# Literature Review - Pattern based ATS



# Literature Review – Edit Dist. in ATS

## Preprocessing Tool

- graph-based to cluster similar-phrase node or string-based to pre-filter sent.s

## Sole Sim. Measurement

### Pattern-based

- weighted Edit Distance
  - consider contextual similarity, and with intuitive explanation and revision ability

## Similarity

### Measurement

- graph-based or term-based weighted edit dist.

# Background Study – Categorization of ATS

Type		Description	
Extractive	Abstractive	Sentence selection	Rephrasing, reinterpreting concepts
Single document	Multiple documents	Summarization of single document	Summarization of multiple document
Generic	Query-Focused	Preserve general important parts	Consider only user preferences
Supervised	Unsupervised	Use training data	Not using training data

# Background Study – Sequential Pattern Mining

A subfield of pattern mining. Discover interesting (sequential) patterns in a (sequential) database

- $I : I = \{i_1; i_2; \dots i_m\}, i_k: \text{item } k$ . A set of all items.
- *itemset*  $X : X \subseteq I$ .
- *sequence*  $S : S = \langle I_1, I_2, \dots, I_n \rangle$ . An ordered itemset list such that  $I_k \subseteq I$ .
- *sequence*  $S_a = \langle A_1, A_2, \dots, A_m \rangle$  *is contained/ subsequence* in/of *sequence*  $S_b = \langle B_1, B_2, \dots, B_n \rangle$  :  
$$S_a \sqsubseteq S_b \text{ if } \exists 1 \leq i_1 \leq i_2 \leq \dots \leq n \text{ s.t. } \forall A_k \in S_a, A_k \subseteq B_{i_k}$$
- *sequential database*  $SDB : SDB = \langle S_1, S_2, \dots, S_n \rangle$ . A list of sequences.
- *support*  $\text{sup}(S_a) : \text{sup}(S_a) = |S| S_a \sqsubseteq S \wedge S \in SDB|$  . Number of sequences in a  $SDB$  containing  $S_a$

# Background Study – Sequential Pattern Mining

SI D	Sequence	Sequential Patterns (minsup = 2)	Closed Sequential Patterns (minsup = 2)
1	$\langle \{a\}, \{b, c\}, \{d, e\} \rangle$	$\langle \{c\} \rangle, \langle \{d\} \rangle, \langle \{e\} \rangle$	$\langle \{c\}, \{d\} \rangle, \langle \{c\}, \{e\} \rangle, \langle \{d\}, \{e\} \rangle$ $\{d, e\}$
2	$\langle \{c\}, \{d, e\} \rangle$	$\langle \{c\}, \{d\} \rangle, \langle \{c\}, \{e\} \rangle, \langle \{d\}, \{e\} \rangle$	
3	$\langle \{e, f\} \rangle$		

A Sequential Database Example with  $I = \{a, b, c, d, e, f\}$

1. Order Matters :

$$\langle \{b\}, \{d, e\} \rangle \sqsubseteq S_1 = \langle \{a\}, \{b, c\}, \{d, e\} \rangle \text{ but } \langle \{d, e\}, \{b\} \rangle \not\sqsubseteq S_1 = \langle \{a\}, \{b, c\}, \{d, e\} \rangle$$

2. Sequential Patterns/ Frequent Subsequences  $FS$ :


$$FS = \{S_{\text{sub}} \mid |S_{\text{sub}} \sqsubseteq S \wedge S \in SDB| \geq \text{minsup}\}.$$

3. Closed Sequential Patterns  $CS$  :

$$CS = \{S_a \mid S_a \in FS \wedge \nexists S_b \in FS \text{ st. } S_a \sqsubset S_b \wedge \text{sup}(S_a) = \text{sup}(S_b)\}.$$

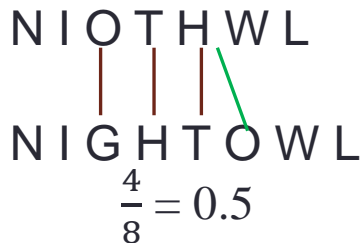
*lossless representation, largest subsequences common to sets of sequences*

# Background Study – Edit Distance

- Distance Measure between two units (usually strings) : counting number of edit operations needed to transform one string to another
- Examples : 

Levenshtein  
Distance

N I O T H W L  
N I G H T O W L

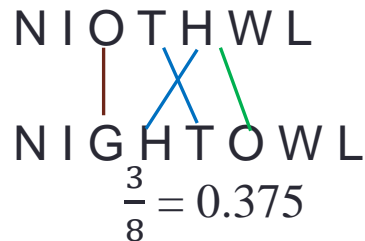


$\frac{4}{8} = 0.5$

Substitution  
Insertion  
Deletion

Damerau–Levenshtein  
Distance

N I O T H W L  
N I G H T O W L



$\frac{3}{8} = 0.375$

Substitution  
Insertion Deletion  
Transposition  
( 2 adjacent units)

Longest Common  
Subsequence Distance

N I O T H W L  
N I G H T O W L

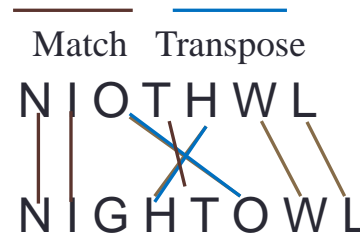


$\frac{5}{8} = 0.625$

Insertion  
Deletion

# Background Study – Jaro Similarity

Transposition



$$1 - \frac{1}{3} \left( \frac{7}{7} + \frac{7}{8} + \frac{7-1}{7} \right) = 0.2898$$

$$sim_j = \frac{1}{3} \left( \frac{|m|}{|s_1|} + \frac{|m|}{|s_2|} + \frac{|m| - \frac{|t|}{2}}{|m|} \right), dist_j = 1 - sim_j$$

$$m = \langle (t_{1i}, t_{2j}) | t_{1i} = t_{2j} \wedge |p_{1i} - p_{2j}| \leq \left\lfloor \frac{\max(|s_1|, |s_2|)}{2} \right\rfloor - 1 \wedge i \notin \bigcup_{k=1}^{i-1} k \rangle$$

$$= \langle (u_{1k}, u_{2k}) \in m \wedge \exists 1 \leq i_1 \leq i_2 \leq \dots \leq l \text{ s.t. } \forall (u_{1k}, u_{2k}), u_{1k} = t_{1i_k}, u_{2k} = t_{2i_k} \rangle$$

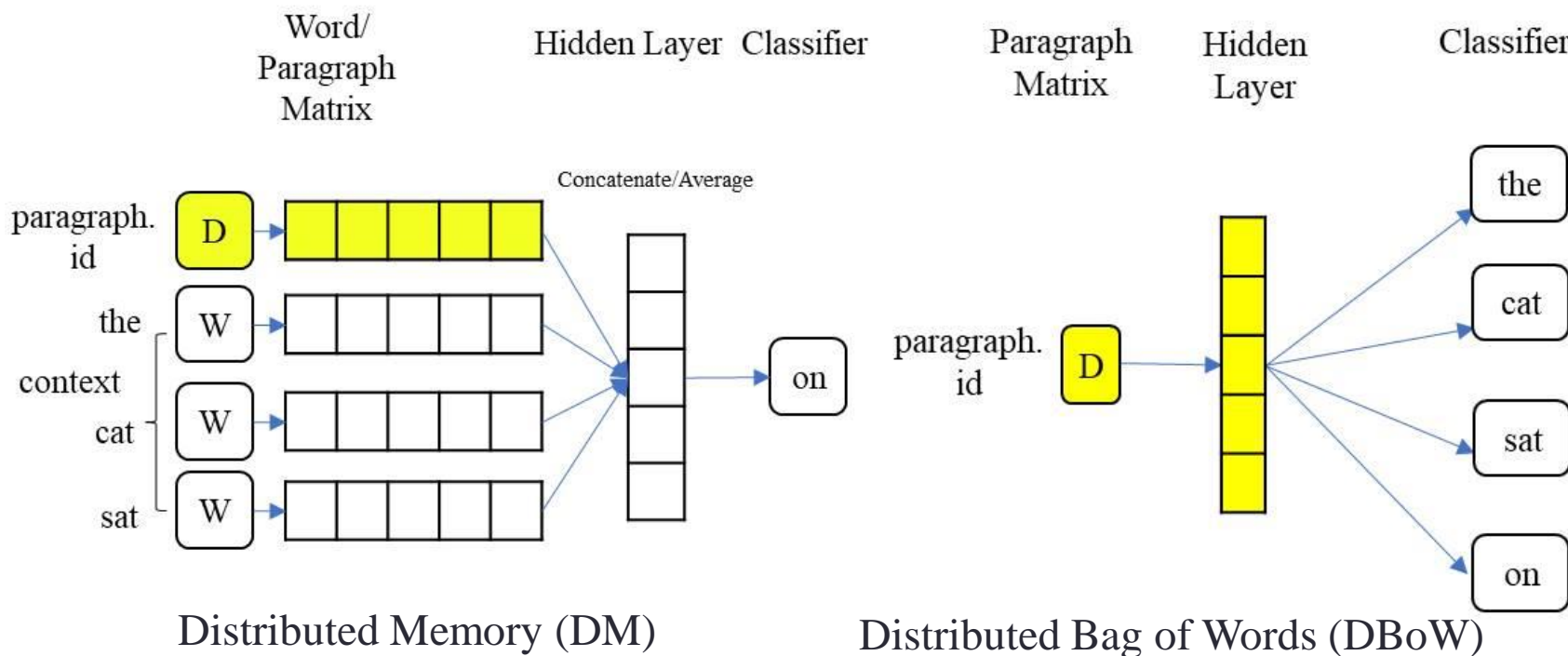
$$t = \langle (u_{1k}, u_{2k}) | (u_{1k}, u_{2k}) \in m_{order} \wedge u_{1k} \neq u_{2k} \rangle$$

$p_{li}$ : position of  $t_{li}$  in  $s_l$

## ● Perform Well

- W. Cohen, P. Ravikumar, and S. Fienberg, “A comparison of string metrics for matching names and records”
- G. Recchia and M. M. Louwerse, “A Comparison of String Similarity Measures for Toponym Matching”

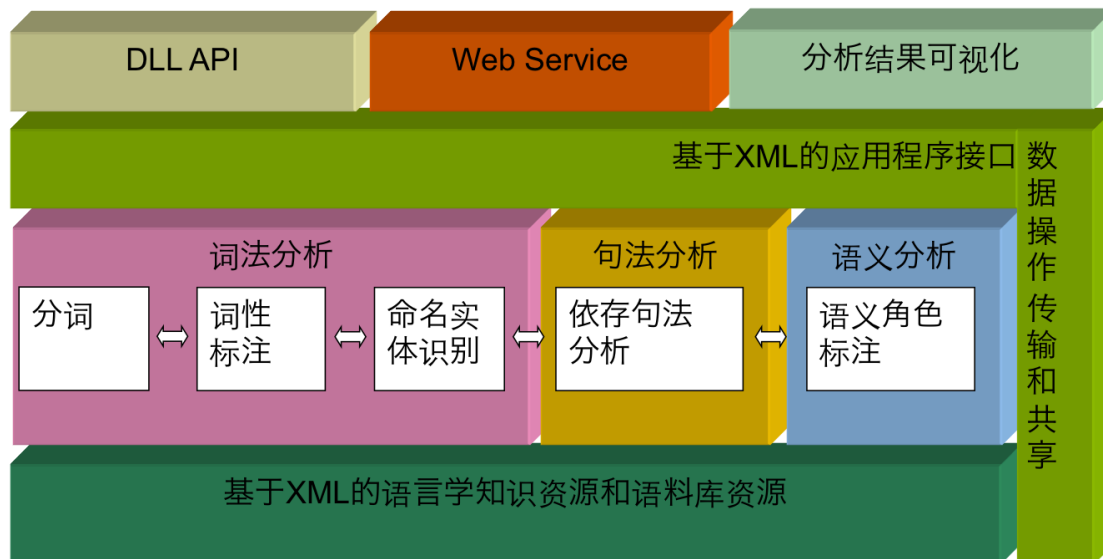
# Background Study – Paragraph Embedding



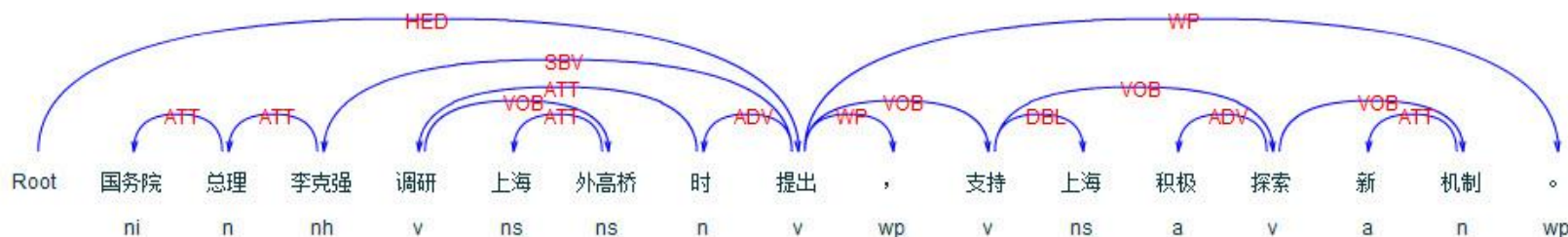
- ✓ Learn contextual information by predicting words from their context in the model
- ✓ Indirectly captures semantics of the paragraph



# Background Study – Tools & Knowledge for Sent. Compression



Framework of Language Technology Platform (LTP) provided by HIT-SCIR ( LTP official page)



Example of a Dependency Parsing Structure ( LTP official page)

# Background Study – Rouge-n

- Full name : Recall-Oriented Understudy for Gisting Evaluation
- To compute the number of units (n-grams) in both the system's summary and the reference one and calculate the recall

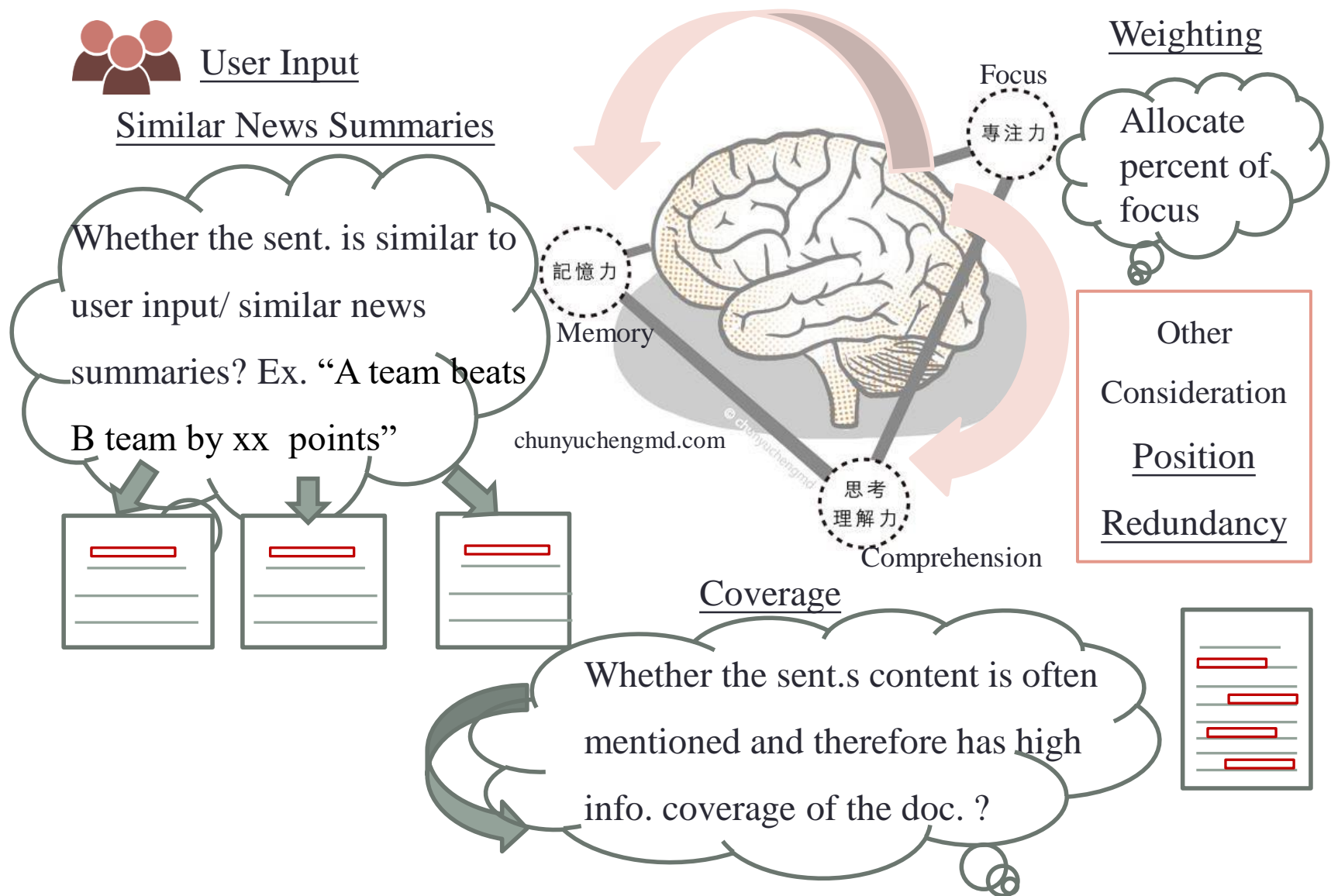
$$\checkmark \text{ Rouge - } n = \frac{\sum_{S_{ri} \in \text{Summ}_{ref}} \sum_{S_{rij} \in S_{ri}} |\{ngram | ngram \in S_{rij} \wedge ngram \in S_{si}\}|}{\sum_{S_{ri} \in \text{Summ}_{ref}} \sum_{S_{rij} \in S_{ri}} |\{ngram | ngram \in S_{rij}\}|}$$

*Summ<sub>ref</sub>*: set of reference summaries, *Summ<sub>sys</sub>* : set of system summaries  
*S<sub>rij</sub>*: jth summary in ith summary set in *Summ<sub>ref</sub>*, *S<sub>si</sub>*: ith summary in *Summ<sub>sys</sub>*

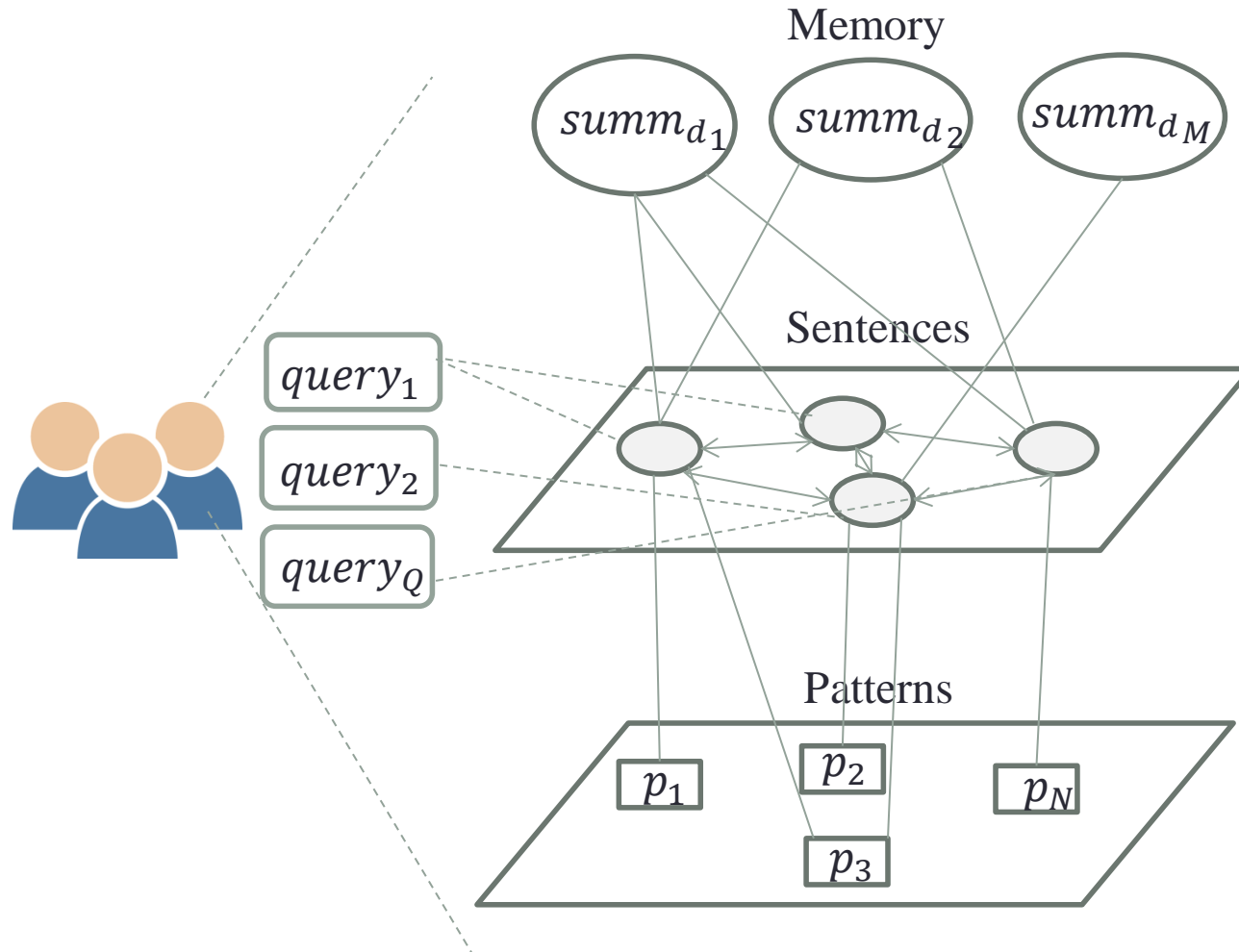
- ✓ *Precision - n* : replacing denominator with total counts of n-grams in system summaries

$$\checkmark F1 : \frac{1}{\frac{1}{2}(\text{Rouge-n} + \text{Precision-n})}$$

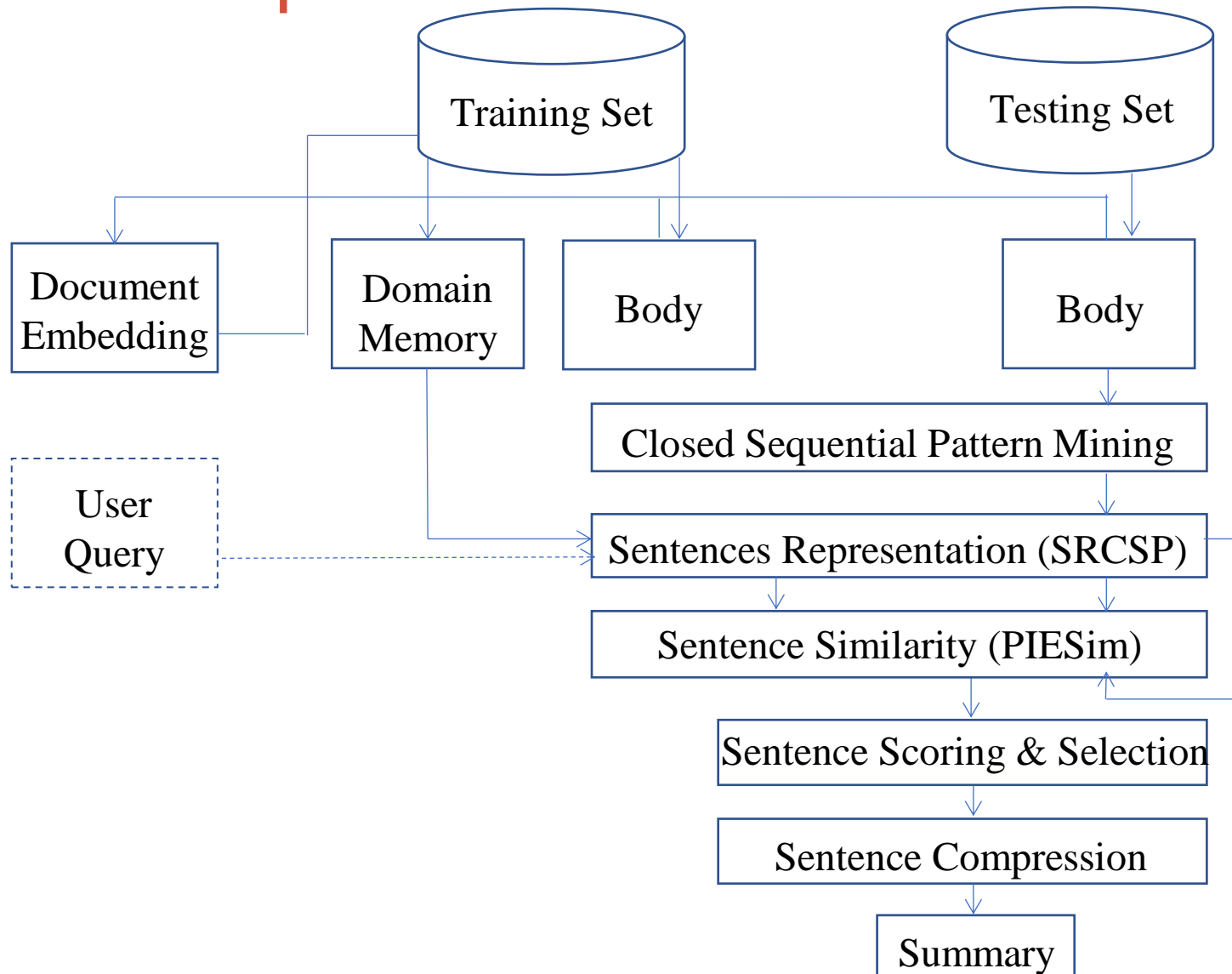
# The Proposed PIESim Model – Analogy to Human



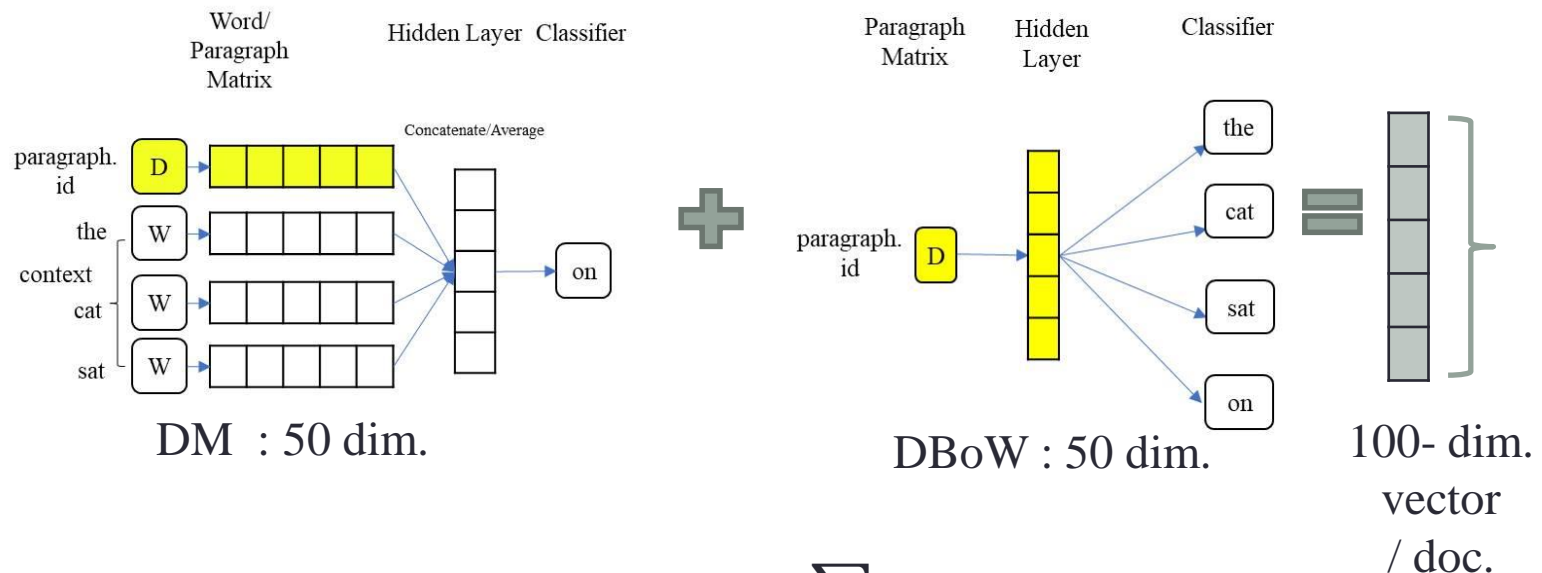
# The Proposed PIESim Model – Structure



# The Proposed PIESim Model – WorkFlow



# The Proposed PIESim Model – Memory Retrieval



$$memory_d = \{s_i | i \in \operatorname{argmax}_{D' \subset D, |D'|=10} \sum_{j \in D \setminus d} \cos(embed_d, embed_j)\}$$

$$\cos(embed_d, embed_j) = \frac{embed_d \cdot embed_j}{|embed_d| |embed_j|}$$

$s_i$ :  $i$ th document's summary,  $embed_j$ :  $j$ th document's embedding,

$$D = \{i | i \in 1, 2, \dots, |\text{training doc.s}|\}$$

# The Proposed PIESim Model – Sent. Rep. using Closed Sequential Patterns (SRCSP)

**Definition 1** Let  $p_i = \langle t_1, t_2, \dots, t_m \rangle$  be a closed sequential pattern with  $\text{sup}(p_i)$ ;  $t_j$  is  $j$ th term and  $m$  is count of terms in  $p_i$ . A “pattern-weight-pair set” of a pattern is

$$pw(p_i) = \{(t_1, w), (t_2, w), \dots, (t_m, w)\}, w = \text{sup}(p_i), t_j \in p_i$$

**Definition 2** Let  $pw(p_i) = \{(t_1, w_i), (t_2, w_i), \dots, (t_m, w_i)\}$  and  $pw(p_j) =$

$\{(s_1, w_j), (s_2, w_j), \dots, (s_n, w_j)\}$  be two pattern-weight-pair sets, the interchangeable composition operation  $\oplus$  of them is formulated below:

$$\begin{aligned} pw(p_i) \oplus pw(p_j) = & \{(t_i, w_i + w_j) | \{(t_i, w_i) | (t_i, w_i) \in pw(p_i) \wedge (s_j, w_j) \in pw(p_j) \wedge t_i = s_j\} \\ & \cup \{(t_i, w_i) | (t_i, w_i) \in pw(p_i), t_i \notin \bigcup_{j=1}^n s_j\} \\ & \cup \{(s_j, w_j) | (s_j, w_j) \in pw(p_j), s_j \notin \bigcup_{i=1}^m t_i\} \end{aligned}$$

# The Proposed PIESim Model - SRCSP

**Definition 3** A sentence representation of a sentence sequence  $s = \langle t_1, t_2, \dots, t_m \rangle$  is the function of ordered composition operation of pattern-weight pairs from patterns contained in the sentence.

$$\begin{aligned}
 & s_{pw} = \{pw(p_i) | p_i \sqsubseteq s, p_i \in CS_d\} \\
 & \quad \quad \quad CS_d : \text{closed sequential patterns of a document } d \\
 & \text{term-based rep.} \left\{ \begin{aligned} & sentrep(s) = pw_{s1} \oplus pw_{s2} \oplus pw_{sk}, \quad pw_{si} \in s_{pw} \\ & sentrep_{term}(s) \\ & \quad = \langle (a_i, w_i) | s_{pw} \neq \emptyset \wedge (a_i, w_i) \in sentrep(s) \wedge \exists 1 \leq i_1 \leq i_2 \leq \dots \\ & \quad \leq m \text{ s.t. } \forall (a_k, w_k), a_k = t_{ik} \wedge a_i \notin \bigcup_{j=1}^{i-1} a_j \rangle \\ & \quad \cup \langle (a_i, 1) | s_{pw} = \emptyset \wedge a_i \in s \wedge \exists 1 \leq i_1 \leq i_2 \leq \dots \leq n \text{ s.t. } \forall a_k, a_k = t_{ik} \rangle \end{aligned} \right. \\
 & \text{char-based rep.} \left\{ \begin{aligned} & s_c = \langle c_{11}, \dots, c_{1|t_1|}, c_{10}, c_{21}, \dots, c_{m-1|t_{m-1}|}, c_{(m-1)0}, c_{m1} \dots c_{m|t_m|} \rangle \\ & f(c_{ij}) = \begin{cases} w_i, j \neq 0 \\ \frac{1}{8}, j = 0 \end{cases}, \quad \begin{matrix} t_i : \text{term } i \text{ in } s, c_{ij} : \text{character } j \text{ in term } i \text{ (include space)} \\ c_{ij} \in s_c, w_i \in \langle w_k | (t_k, w_k) \in sentrep_{term}(s) \rangle \end{matrix} \\ & sentrep_{char.}(s) = \langle (s_{ij_k}, f(s_{ij})_k) | s_{ij} \in s_c \rangle \end{aligned} \right.
 \end{aligned}$$



# The Proposed PIESim Model – SRCSP Examples

**Def. 1**  $p_i = \{read, book\}$  with support 2  $\rightarrow pw(p_i) = \{(read, 2), (book, 2)\}$

**Def. 2**  $pw(p_j) = \{(buy, 3), (book, 3)\} \rightarrow pw(p_i) \oplus pw(p_j) = \{(book, 5), (buy, 3), (read, 2)\}$

Def. 3	Original Sentence	Contained CS	SRCSP
English	I, usually, buy, a, book, < from, the, bookstore, > and, read, the, book, at, night	$\{(buy, 3), (book, 3)\}$ $\{(read, 2), (book, 2)\}$	$\langle (buy, 3), (book, 5), (read, 2) \rangle$
Chinese	< 我, 通常, 從, 書店, 買, 書本, > 然後, 在, 晚上, 讀, 書本 >	$\{(買, 3), (書本, 3)\},$ $\{(讀, 2), (書本, 2)\}$	$\langle (買, 3), (‘ ‘, \frac{1}{8}) (書, 5), (本, 5), (‘ ‘, \frac{1}{8}), (讀, 2) \rangle$

- *English or lang.s with few char.s*:, matching in unit of character will be misleading
- *Chinese or lang.s with many char.s*: 50000- 100000 characters. Character match actually implies semantic relation in many cases. Ex. “腳踏車(bicycle)” & “汽車(car)”  
Will be inappropriate sometimes, gain is larger than loss of info. In our experiments

# The Proposed PIESim Model – PIESim Measure

$$s_{rep1} = sentrep_{term(character)}(s_1), s_{rep2} = sentrep_{term(character)}(s_2)$$

$$PIESim(s_1, s_2) = \frac{1}{3} \left( \frac{m_w}{\sum_{\{w_i | (t_i, w_i) \in s_{rep1}\}} w_k} + \frac{m_w}{\sum_{\{w_j | (t_j, w_j) \in s_{rep2}\}} w_k} + \frac{m_w - \frac{|t|}{2}}{m_w} \right)$$

$$m = \langle (c_{1i}, c_{2j}) | c_{1i} = c_{2j} \wedge |p_{1i} - p_{2j}| \leq \left\lfloor \frac{\max(|s_{rep1}|, |s_{rep2}|)}{2} \right\rfloor - 1 \wedge i \notin \bigcup_{k=1}^{i-1} k \rangle$$

$$m_w = \sum_{\substack{\{(w_{1i}, w_{2j}) | (c_{1i}, c_{2j}) \in m \wedge (c_{1i}, w_{1i}) \in s_{rep1} \wedge \\ (c_{2j}, w_{2j}) \in s_{rep2}\}}} \frac{w_{1i} + w_{2j}}{2 m_{order}} = \langle (u_{1k}, u_{2k}) \in m \wedge \exists 1 \leq i_1 \leq i_2 \leq \dots \leq l \rangle$$

$\bar{u}_{1k}^1, u_{2k}) |$   
*s. t.*  
 $\forall (u_{1k}, u_{2k}), u_{1k} = c_{1i_k}, u_{2k} = c_{2i_k}$

$$t = \langle (u_{1k}, u_{2k}) | (u_{1k}, u_{2k}) \in m_{order} \wedge u_{1k} \neq u_{2k} \rangle$$

$c_{li}$ : unit  $i$  in  $s_l$ ,  $p_{li}$ : position of  $c_{li}$  in  $s_l$

Sim. Measure	SRCSP of s1 s2	Sim. Computation
PIESim <u>match</u>	{(buy, 3), (book, 5), (read, 2)}	$\frac{1}{3} \left( \frac{\frac{5+5}{2} + \frac{2+2}{2}}{10} = 7 + \frac{7}{13} + \frac{7-1}{7} \right)$ $\approx 0.7$
Jaro Similarity <u>transposition</u>	{(I, 3), (love, 3), ( <u>read</u> , 2), ( <u>book</u> , 5)}	$\frac{1}{3} \left( \frac{2}{4} + \frac{2}{3} + \frac{2-1}{2} \right) \approx 0.55$

# The Proposed PIESim Model – Scoring Selection

$$score(s_i) = \{ (1 - \beta) * \left[ \alpha * \frac{cov(s_i)}{\sum_j cov(s_j)} + (1 - \alpha) * \frac{mem(s_i)}{\sum_j mem(s_j)} \right] + \underbrace{\beta * \frac{div(s_i)}{\sum_j div(s_j)}}_{\text{Diversity (- Redundancy)}} * \underbrace{pos_i}_{\text{Positional Weight}} \}$$

Coverage
Memory Sim.

**Diversity (- Redundancy) Positional Weight**

$$cov(s_i) = \frac{\sum_{s \in D} PIESim(s_i, s) - 1}{|D| - 1}, mem(s_i) = \frac{\sum_{h \in memory_d} PIESim(s_i, h)}{|memory_d|}$$

$$div(s_i) = \min(\{1 - PIESim(s_k, s_i) | s_k \in S\})$$

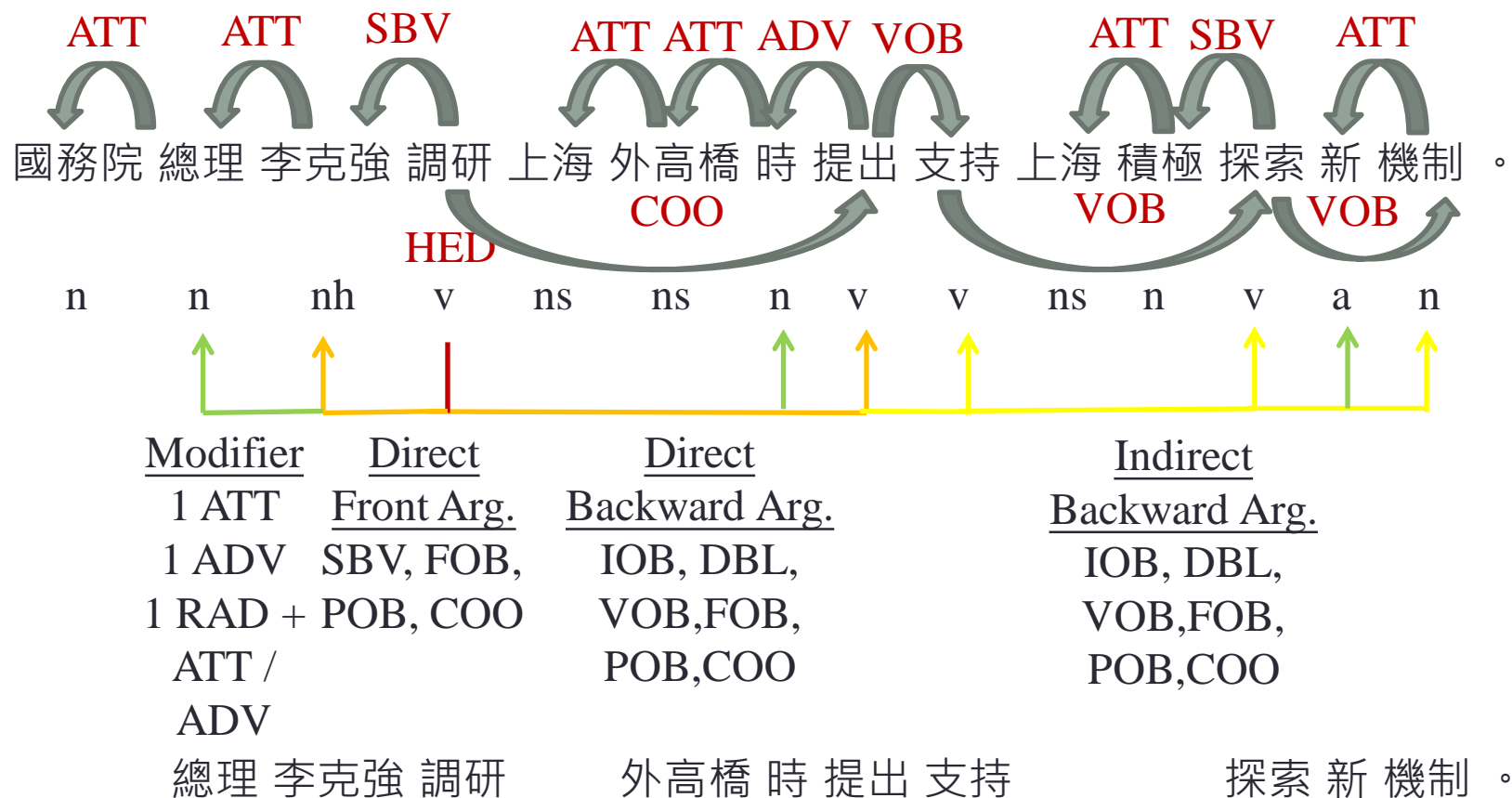
$$pos_i = 1 - \frac{(i-1)}{|\{s | s \in \{sent.s \text{ in document } d\}\}|}, \quad h \in memory_d, s \in \{sentences \text{ in } d\},$$

$$pos_i^j = 1 - \frac{(i-1)}{|\{s_i^j | s_i^j \in \{sent.s \text{ in document } d_j\}\}|}$$

$s_i$  :  $i$ th sentence of  $d$ ,  $d$ : document  $d$   $S$ : summary set of  $d$

✓ greedily select sent. with highest score iteratively until we reach the limits

# The Proposed PIESim Model – Compression (Chinese)







# Experiments and Results – Chinese

- 2017 -2018 492997 UDN news : Select 2 sent. or 1 sent. > 30 char.s



summary stat	med.	std.
word/ hl.	10	2.34
word/ sent. in bd.	7	5.14
char./ word in bd.	2	0.93

Size	UDN
Train	443700
Test	49297
Median	~188
Size	LCSTS
Train	2407551
Test	1106
Median	~60

	Statistical based
	Graph based
	Pattern based
	DNN

System	PIESim	PIESim <sub>comp.</sub>	ILP	Submod.1	Submod.2
Rouge-1 F1	<b>0.3137</b>	0.2753	0.2766	0.2647	0.2646
Rouge-2 F1	<b>0.1836</b>	0.1520	0.1574	0.1467	0.1467
System	Luhn	SumBasic	LSA	TextRank	LexRank
Rouge-1 F1	0.2459	0.2706	0.2083	0.2385	0.2083
Rouge-2 F1	0.1352	0.1508	0.1061	0.1330	0.1061
System	Reduction	PatSum	Seq2Seq	SeqAtt.	LCSTS Seq / SeqAtt.
Rouge-1 F1	0.2385	0.3001	0.1933	0.2607	0.215/ 0.089
Rouge-2 F1	0.1330	0.1775	0.0732	0.1117	0.299/ 0.174

# Experiments and Results - English

- Task 2 from benchmark Dataset from DUC 2004 Conference
- 50 clusters with 10 docs. each, No training data – multi. doc.s / unsupervised ATS
- ✓ Pattern Mining : 10-doc.s / Position Weight : Start from 1 in each doc. / Len. Limit : 665 bytes

<div>Sent. Split (DUC script)</div> <div>↓</div> <div>Remove Stop Words (Cornell)</div> <div>↓</div> <div>Stem Words (Porter Stemmer)</div>		Rouge 2			Rouge 4		
		Precision	Recall	F1	Precision	Recall	F1
	<b>PIESim</b>	<b>0.0943</b>	<b>0.0953</b>	<b>0.0947</b>	<b>0.0168</b>	<b>0.0170</b>	<b>0.0169</b>
	Pat Sum	0.0988	0.0990	0.0986	0.0167	0.0166	0.0166
	Freq. Itemset	0.0864	0.0869	0.0866	0.0135	0.0136	0.0135
	Weight. Itemset	0.0916	0.0904	0.0909			
	Best Peer	0.092	0.091	0.091	0.015	0.015	0.015
	Human A	0.088	0.092	0.090	0.010	0.009	0.010
	Human B	0.096	0.091	0.092	0.013	0.013	0.013
	Human C	0.102	0.094	0.098	0.012	0.011	0.012
	Human D	0.106	0.100	0.102	0.010	0.010	0.010
	Human E	0.099	0.094	0.097	0.012	0.011	0.012
	Human H	0.105	0.101	0.103	0.013	0.012	0.012

# Different Settings Analysis – Pattern Usage

\*, \*\*, \*\*\* represent two-sided-paired-sample t-tests with  $\alpha = 0.05, 0.01, 0.001$  significance level

Tested under 25 combinations of coverage  $\alpha$  and redundancy par.s  $\beta$  0.1, 0.3, 0.5, 0.7, 0.9 on UDN

	Mean Rouge 1	Mean Rouge 2	Mean	Mean	Best Rouge 1
	F1 Difference	F1 Difference	Rouge 1 F1	Rouge 2 F1	F1 / Rouge 2 F2
Use pattern or not – Base : Testing doc. SRCSP w. min. occur. 2 + Sent.s in Memory's Hl.s					
No Pattern	0.0007	-0.0023	0.2958	0.1678	0.3085/ 0.1756
Train/ Test doc. Patterns	-0.0013***	-0.0010***	0.2937	0.1692	0.3062/ 0.1787
<b>Test doc. Pattern</b>	<b>base</b>	<b>base</b>	<b>0.2951</b>	<b>0.1702</b>	<b>0.3067/0.1789</b>
Pattern Variant– Base : Testing doc. SRCSP with min. occurrence 2					
Occurrence 3	0.0013	0.0011	0.2965	0.1714	0.3047/ 0.1777
All Occurrences	-0.00017***	- 0.00015***	0.2949	0.1700	0.3067/0.1789
Text Format – Base : Original Testing doc. SRCSP with min. occurrence 2					
Map Ehownet & NER	-0.0111***	-0.0067***	0.2840	0.1634	0.3015/ 0.1756
Map Ehownet	-0.009***	-0.006***	0.2853	0.1632	0.3012/ 0.1744

# Different Settings Analysis – PIESim Variant

	Mean Rouge 1 F1 Difference	Mean Rouge 2 F1 Difference	Mean Rouge 1 F1	Mean Rouge 2 F1	Best Rouge 1 F1/ Rouge 2 F1
w. or w/o. PIESim-weights information ( space weight = 1)					
<b>PIESim Weight</b>	<b>base</b>	<b>base</b>	<b>0.3012</b>	<b>0.1747</b>	<b>0.3118/ 0.1821</b>
No Weight	-0.0060***	-0.0045***	0.2951	0.1702	0.3067/ 0.1789
w. or w/o. word level information ( space weight = 1)					
Not Sep. by Spaces	-0.0065**	-0.0067***	0.2946	0.1680	0.3043/ 0.1742
Matched by words	-0.0200***	-0.0115***	0.2812	0.1632	0.2870/ 0.1666
Space Weight					
1/2	0.0015***	0.0011***	0.3027	0.1758	0.3128/ 0.1829
1/4	0.0020***	0.0015***	0.3032	0.1762	0.3133/ 0.1833
1/8	0.0022***	0.0016***	0.3034	0.1764	0.3137/ 0.1836
0	-0.0077***	-0.0069***	0.2934	0.1678	0.3033/ 0.1739



# Different Settings Analysis – Sent. Rep./ Sim. Measure/ Scoring Variant

	Mean Rouge 1 F1 Difference	Mean Rouge 2 F1 Difference	Mean Rouge 1 F1	Mean Rouge 2 F1	Best Rouge 1 F1/ Rouge 2 F1
<b>PIESim</b>	<b>base</b>	<b>base</b>	<b>0.3034</b>	<b>0.1764</b>	<b>0.3137/ 0.1836</b>
Sent. Rep. & Sim. measurement Variants					
term-freq. infused edit sim.	-0.0088***	-0.0055***	0.2946	0.1709	0.3036/0.1772
tfidf	-0.0165***	-0.0112***	0.2869	0.1652	0.3102/0.1812
sent embedding	-0.0334***	-0.0216***	0.2700	0.1548	0.2905/0.1682
Sent. Scoring					
PageRank	0.0008	0.0011	0.3042	0.1775	0.3121/0.1828

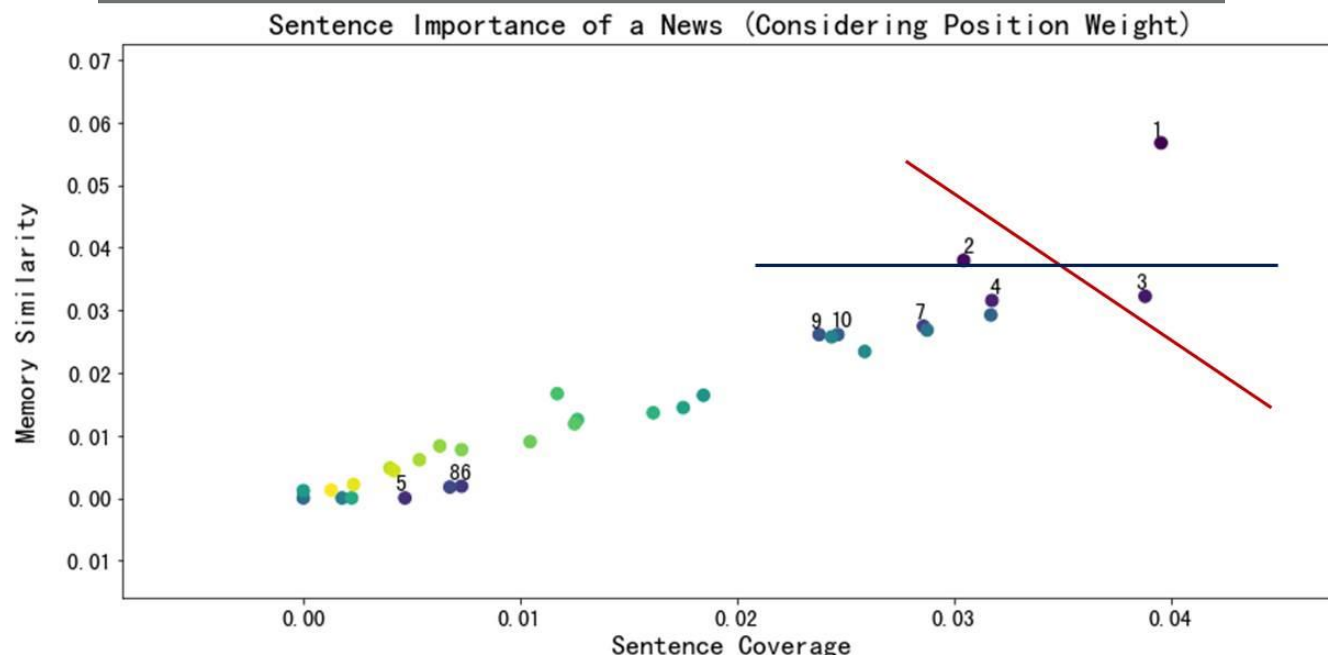
The word embedding is obtained from the open source FastText [62] trained on Wikipedia  
The sent. embedding is the average of all available word embedding in the sent.

# Error Analysis – Effectiveness of Coverage

懷特 ( 4108 ) 研發上市新藥獲癌症醫界肯定<sup>1</sup>，黃耆經萃取分離及高度純化、研發成功的黃耆多醣注射劑<sup>2</sup>，經證實六成以上患者治療後可緩解癌因性疲憊症<sup>3</sup>，懷特4月營收858萬元，較上月減少21.57%...

Headline :懷特新藥獲肯定，六成患者治療後可緩解癌因性疲憊症。

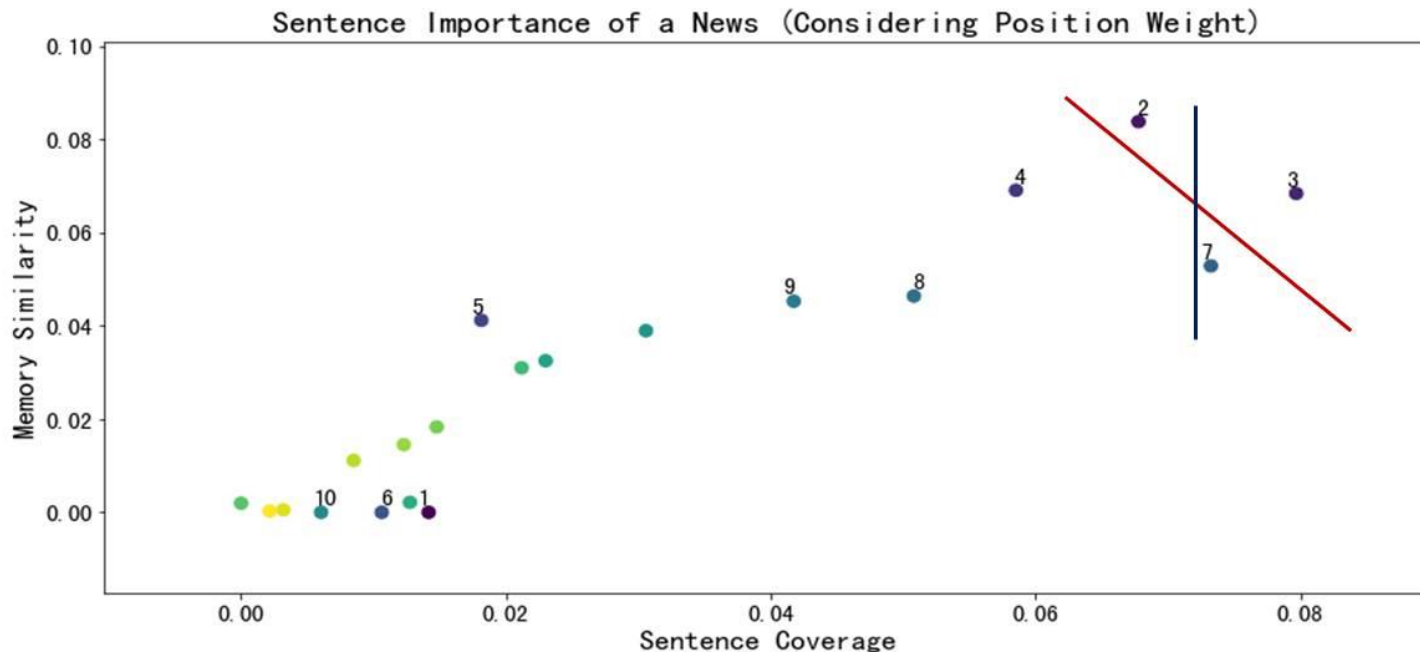
Newly invented medicine by Phytohealth is praised. Cancer-related fatigue of near 60% patients can be eased after the treatment.



# Error Analysis – Effectiveness of Memory Sim.

因應瑪莉亞颱風襲台，台鐵上午9時宣布今（10）日16時前<sup>2</sup>，全線各級列車照常行駛<sup>3</sup>，中午12時會公布16時以後的列車行駛情形，台鐵表示，自發布海上颱風警報起至解除海上颱風警報止，購買上述期間內各級列車車票之旅客<sup>7</sup>，可自乘車日起一年內，持未經使用之車票至各  
Headline:台鐵16時前全線列車照常行駛。

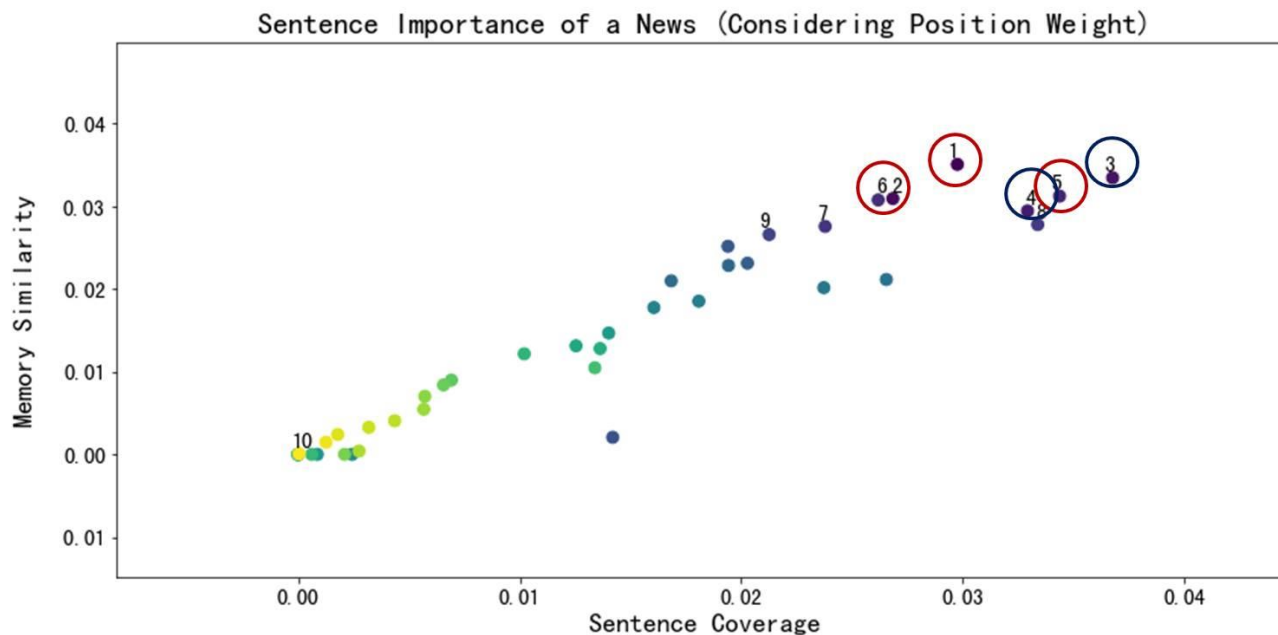
Taiwan railway announced the railway system would operate as usual before 4 p.m.



# Error Analysis – Revise from Coverage

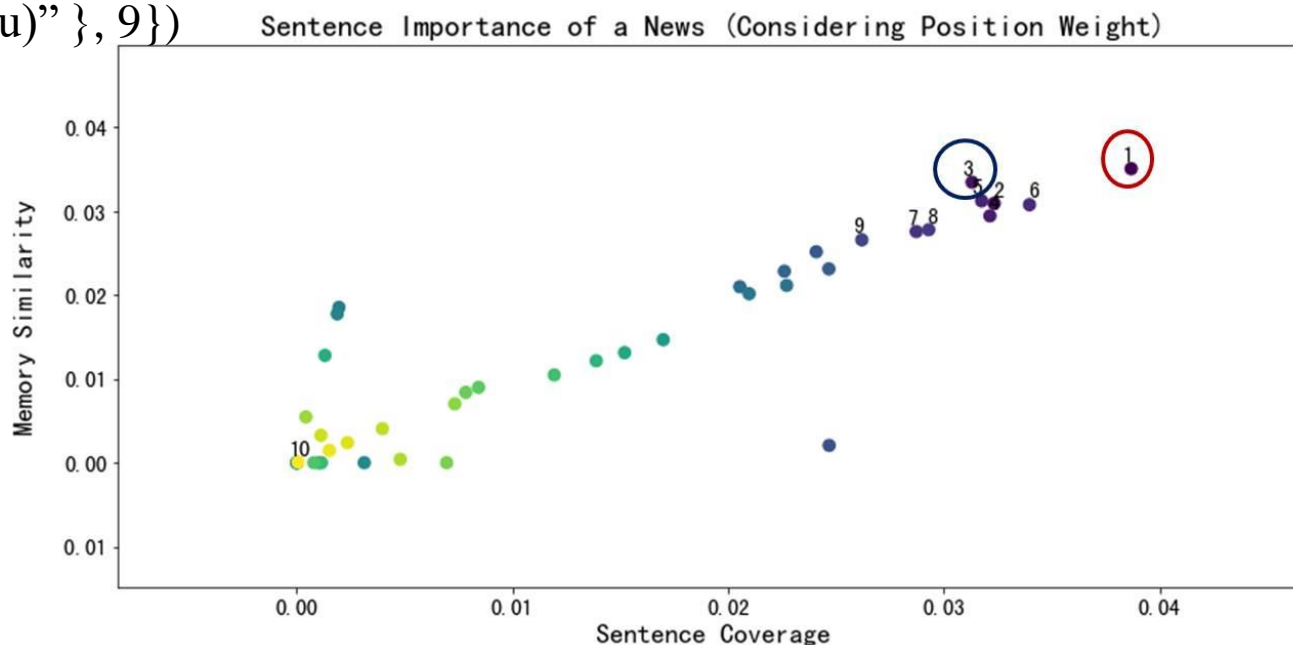
台視、東森「鐘樓愛人」端午連假仍加緊拍攝中<sup>1</sup>，藝人比莉神秘探班，  
為兒子周湯豪加油打氣，劇中演員林帥甫也帶著粽子慰勞劇組人員<sup>4</sup>，  
還爆料周湯豪會畫專屬Q版人物給他<sup>5</sup>，端午節連假「鐘樓愛人」雖馬<sup>6</sup>  
不停蹄拍攝，但也有不少藝人到場探班...The TV series “Love,  
Headline: 「鐘樓愛人」端午趕拍，周湯豪手繪Q版人物。

The TV series “Love, Timeless” kept filming during the Dragon boat festival. Nick Chou drew exclusive cutie version characters.



# Error Analysis – Revise from Coverage

▶ Examine	▶ Revise from Patterns	▶ Repeat same SRCSP PIESim, Scoring process
Wrong sentences'	Delete Patterns :	
Dominant Patterns :	((“周湯豪” }, 10), ((“林帥甫”}, 9)	
((“周湯豪(Nick	Retain Patterns : ((“周湯豪”, “林帥甫”},3)	
Chou)” }, 10), ((“林帥甫	((“劇(TV series)”, “中(In)”, “林帥甫”,} 2)	
(Lin, Shuai-Fu)” }, 9}}		



# Error Analysis – Revise from Memory Sim.

搶先布局5G商機<sup>1</sup>，遠傳電信董事長徐旭東今(12)日宣布，正式成立「遠傳5G先鋒隊」<sup>3</sup>，攜手工研院、愛立信、國內供應商，打造車聯網創新基地<sup>5</sup>，現場展示兩款自駕車，展現台灣5G應用實力。

Headline :徐旭東宣布成立「5G先鋒隊」主攻車聯網。

Douglas Hsu, announced to establish “Far East vanguard” focusing on Internet of Vehicle.

▶ Limitation : Coverage

▶ Change 4 least sim. memory

“組建車聯網 (Build IoV)”

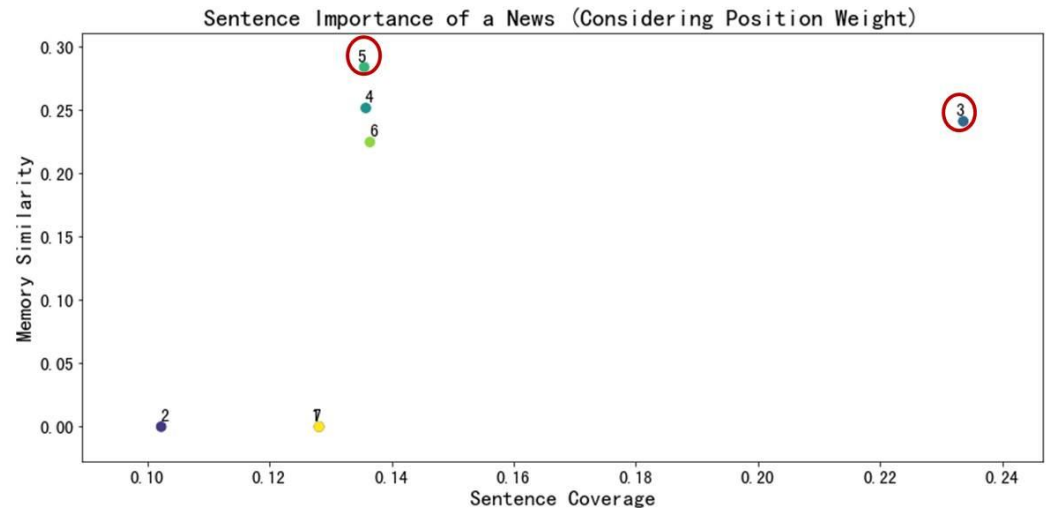
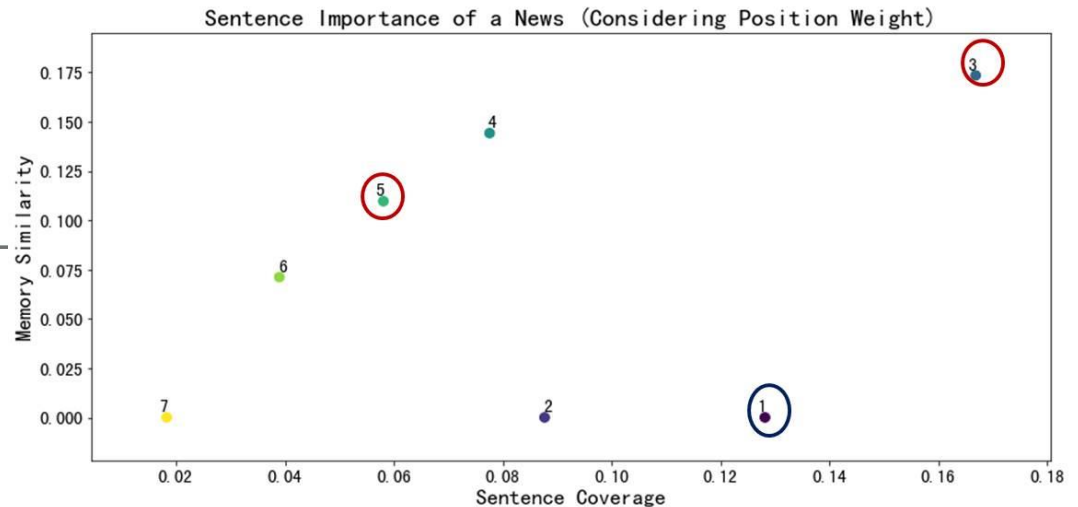
“5G推動車聯網

(5G accelerates IoV)” etc

▶ Short Text

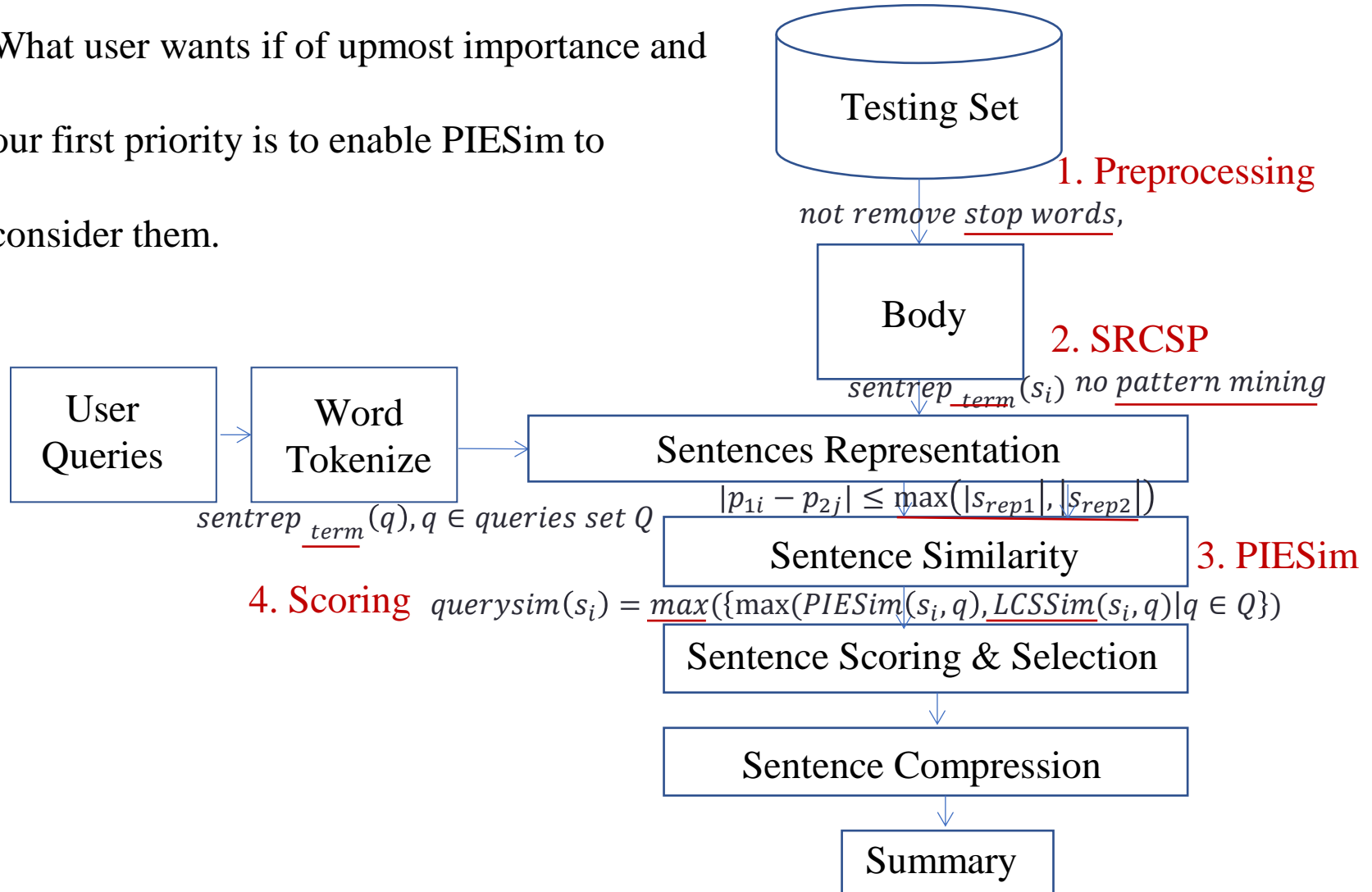
Remove position aspects

▶ Repeat same SRCSP  
PIESim, Scoring process



# User Interaction – Work Flow

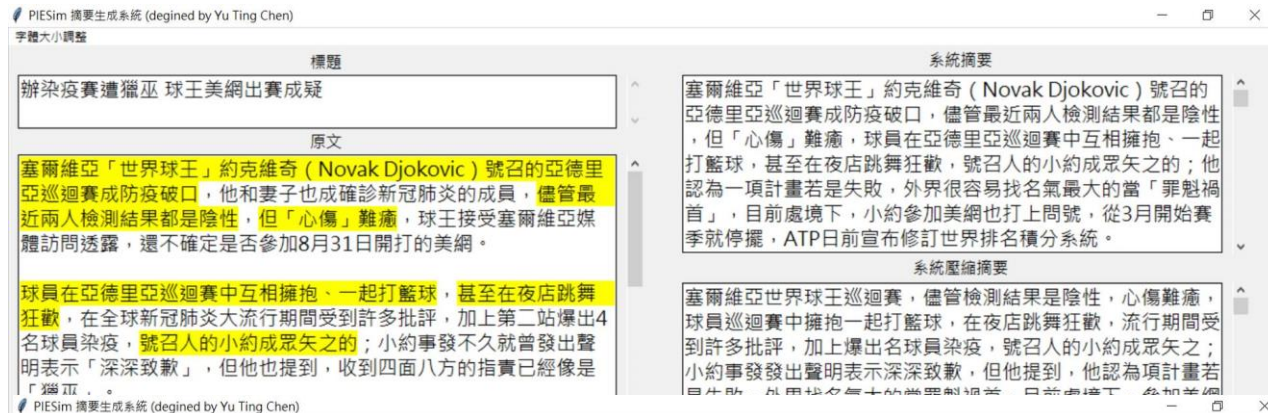
- What user wants if of upmost importance and our first priority is to enable PIESim to consider them.



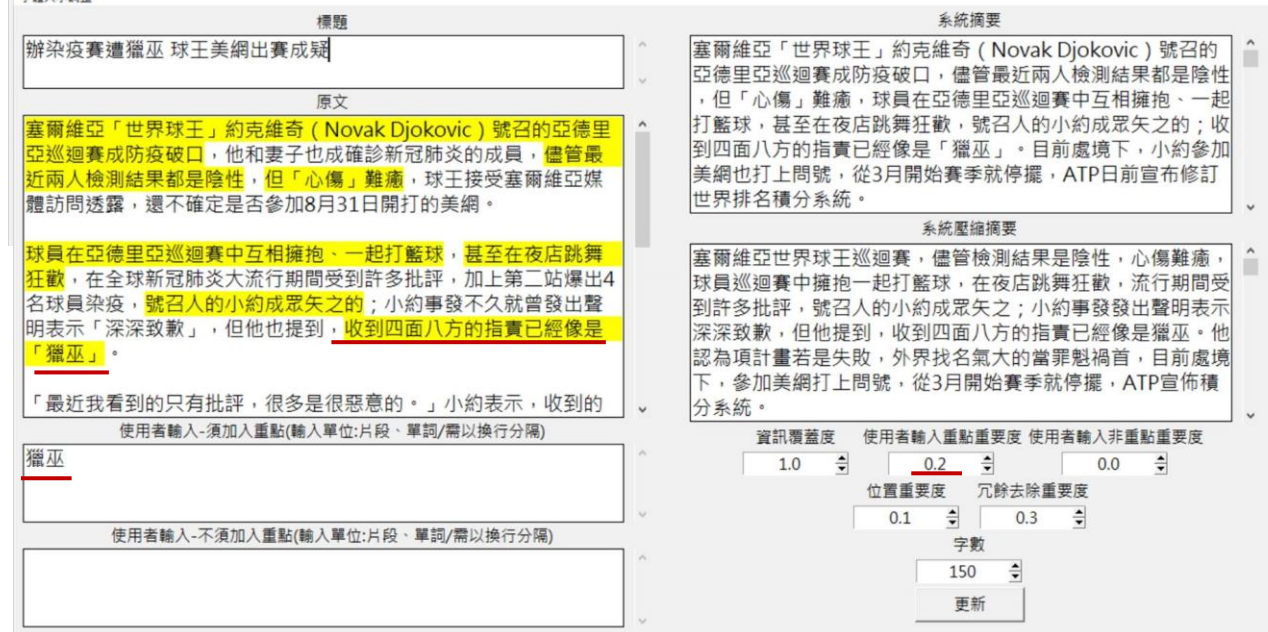


# User Interaction – Interface Demo

- Before Query

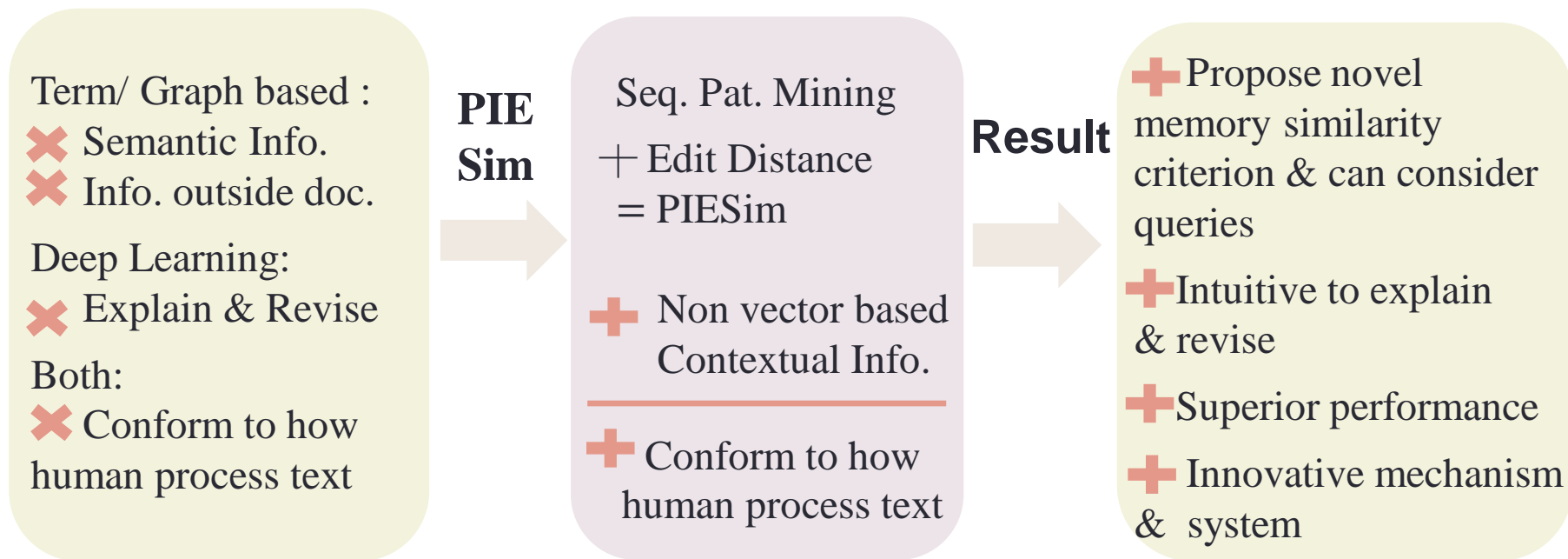


- After Query





# Conclusion & Contribution

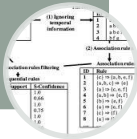


The interaction GUI can be downloaded from <https://rb.gy/cagef4>. (big since it contains many models & modules). Currently work on **Windows 8/10 64 bytes**. Unzip it and click the **PIESim.exe** in the folder Make your input text **clean will produce a better-quality summary**. Also, since we tested it on UDN, we recommend you to copy news from <https://udn.com/news/index>, but all formatted articles should work.

# Future Work and Prospect

- Sequential Rule Mining

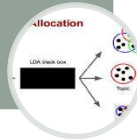
# Cause affect relation



CMRules paper

- Topic Modeling or Clustering

## Topical Information

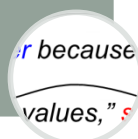


## Topic modeling of weather and climate condition paper

- ❑ General and Query-focused
- ❑ Extractive and Abstractive (Chinese)
- ❑ Single document and Multi-document
- ❑ Unsupervised and Supervised

- Linguistic Features ex. co-ref. resolution

# Linguistic Techniques



Stanford NLP



bilalqambrani.blogspot

More works using pattern mining and edit dist. In NLP and AI tasks



www.moea.gov.tw

## PIESim as both the ATS system and the similarity measurement

Thanks for Your Listening !