Rethinking Coherence Modeling: Synthetic vs. Downstream Tasks

Coherence Analysis

- Coherence models distinguish a coherent text from incoherent ones
- Proposed with the aim of applying them in text generation and ranking - e.g., machine translation, essay scoring, dialog systems, summarization, ...
- But most work on coherence modeling **ignores** downstream applications



Motivation

- Common evaluation: synthetic discrimination tasks
- Unclear if existing models are capable of modeling coherence of generated text from downstream applications

Goal: Assess coherence models on standard discrimination tasks and compare against performance on downstream use-cases.

Standard Training and Synthetic Tasks

Training:

- WSJ news dataset
- Standard pairwise training setup: coherent vs. incoherent document ranking
- 20 random permutations (incoherent documents) for training and testing

Global Discrimination Tasks

- **Standard**: Incoherent document is a random permutation of sentences
- . **Inverse**: Incoherent document has a fully reversed sentence order

s ₀	"The House voted to boost the federal minimum wage for the first time since early 1981, casting a solid 382-37 vote for a compromise measure backed by President Bush."
s ₁	"The vote came after a debate replete with complaints from both proponents and critics of a substantial increase in the wage floor."
s ₂	"Advocates said the 90-cent-an-hour rise , to \$ 4.25 an hour by April 1991 , is too small for the working poor , while opponents argued that the increase will still hurt small business and cost many thousands of jobs."
S 3	"But the legislation reflected a compromise agreed to on Tuesday by President Bush and Democratic leaders in Congress, after congressional Republicans urged the White House to bend a bit from its previous resistance to compromise."
s ₄	"So both sides accepted the compromise , which would lead to the first lifting of the minimum wage since a four-year law was enacted in 1977 , raising the wage to \$ 3.35 an hour from \$ 2.65." (a) Positive sample data
s ₄	"So both sides accepted the compromise , which would lead to the first lifting of the minimum wage since a four-year law was enacted in 1977 , raising the wage to \$ 3.35 an hour from \$ 2.65."
S 3	"But the legislation reflected a compromise agreed to on Tuesday by President Bush and Democratic leaders in Congress , after congressional Republicans urged the White House to bend a bit from its previous resistance to compromise."
s ₀	"The House voted to boost the federal minimum wage for the first time since early 1981, casting a solid 382-37 vote for a compromise measure backed by President Bush."
s ₂	"Advocates said the 90-cent-an-hour rise , to \$ 4.25 an hour by April 1991 , is too small for the working poor , while opponents argued that the increase will still hurt small business and cost many thousands of jobs."
s ₁	"The vote came after a debate replete with complaints from both proponents and critics of a substantial increase in the wage floor."
	(b) Negative sample data

Model	Standard	Inverse
EGRID	81.60	75.78
NEURALEGRID	84.36	83.94
LEXNEUEGRID	88.51	88.13
TRANSMODEL	91.77	99.62
UNIFIEDMODEL	93.19	96.78

Table 1. Results: Accuracies of the coherence models in the **Global Discrimination** task.

Tasnim Mohiuddin*, Prathyusha Jwalapuram*, Xiang Lin*, and Shafiq Joty* +

Nanyang Technological University

Salesforce Research Asia

Downstream Tasks

Machine Translation Coherence

- Reference translations more coherent than system translations at document level. Test if models score reference higher
- . User study: get pairwise coherence ranks of translations from different systems. Check model agreement with humans

Model	Acc. (%)	AC1 Agr
EGRID	51.75	0.80
NEURALEGRID	54.75	0.77
LexNeuEGrid	49.34	0.76
TRANSMODEL	48.67	0.77
UnifiedModel	43.36	0.78

Candidate A:
"A Generation I
demanding poli
Witnesses say t
1 ~
Candidate B:
"A Nesil Protest

Table 2. Machine Translation Results

 \Rightarrow Models like TRANSMODEL and UNIFIEDMODEL with strong global discrimination performance do not perform well

Summarization

- **Abstractive:** User study for coherence ranking of summaries from abstractive summarization systems. Check model agreements with humans
- **Extractive:** Human coherence ratings for extractive summaries from Document Understanding Conference. Check model agreements with human ranks converted from ratings

Models	Abs. Agr.	Ext. Agr.
EGRID	0.71	0.52
NEURALEGRID	0.68	0.70
LexNeuEGrid	0.71	0.57
TRANSMODEL	0.55	0.38
UNIFIEDMODEL	0.68	0.35

Freddie Gray, 25, died in police custody 15 days ago after he was arrested on a weapons charge. His family said his voice box was crushed and his neck snapped before he slipped into a coma. Hundreds of protesters peacefully rallied on the streets of Baltimore on Saturday against the alleged police role in Gray's death.

Table 3. Text Summarization Results

 \Rightarrow Similar pattern; models with high performance in synthetic tasks perform poorly

^{*}Equal Contribution



Eric garner's family and other members of families united for justice will attend gray's funeral. Gray was arrested April 12 and died a week later from a sever spinal cord injury. Three white house officials will also attend gray's funeral



useful task-specific features

Machine Translation Setup

- negative documents using WMT data

Model	Acc. (%)	AC1 Agr.
EGRID	48.74	0.797
NEURALEGRID	52.58	0.76
LexNeuEGrid	56.84	0.795
TRANSMODEL	57.65	0.751
UnifiedModel	77.35	0.828

Discussion:

- downstream applications
- coherence issues more subtle than permuted text
- sentence by sentence unlikely to be out of order

Conclusions:

Please scan the QR code for more information and resources:





Task-specific Training

Investigate whether a change in the training setup might help models learn more

Train coherence models with reference text as positive and system translations as

Report accuracy and agreement on the same test data as previous experiment

 \Rightarrow Performance improves across models; UNIFIEDMODEL improves from 43.36% to 77.35%, high agreement of 0.83

Conclusions

Models trained on permuted sentences may not be learning features useful for

Models based on synthetic tasks may be overfitting on the these tasks, failing to find

Only considering incoherence from permuted documents may be a **poor**

approximation of real-world coherence problems *e.g.*, MT output often produced

• Standard training paradigm may need reform to build more generalizable models Standard evaluation needs reform to be more indicative of real-world performance

Resources

