

N-best Response-based Analysis of Contradiction-awareness in Neural Response Generation Models

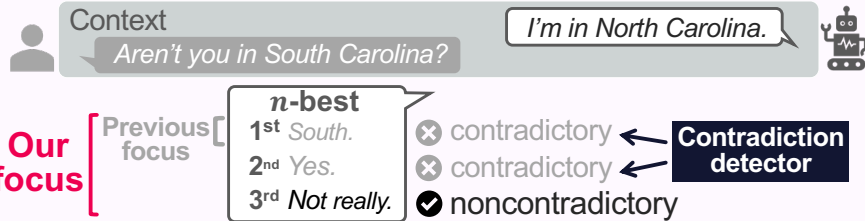
Shiki Sato¹, Reina Akama^{1,3}, Hiroki Ouchi^{2,3}, Ryoko Tokuhisa¹, Jun Suzuki^{1,3}, Kentaro Inui^{1,3}

¹Tohoku University ²Nara Institute of Science and Technology ³RIKEN

Overview

- We analyze contextual **contradiction-awareness** of response generation models focusing on **consistency of n -best candidates**
- Beam search has limitation** on avoiding contradiction and **unlikelihood training reduce** it

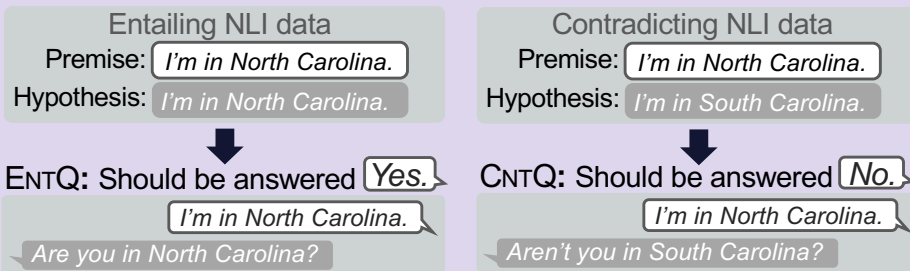
Background: Avoiding contradiction with contradiction detector



All candidates affect whether final output is contradictory

Need to observe all candidates' consistency

Proposed framework



Yes-no questions allow for **clear judge** of responses

All n -bests should have **✓** to always avoid contradiction

Certainty = $\frac{\text{Number of } n\text{-bests having at least one } \checkmark}{\text{Number of } n\text{-bests}}$

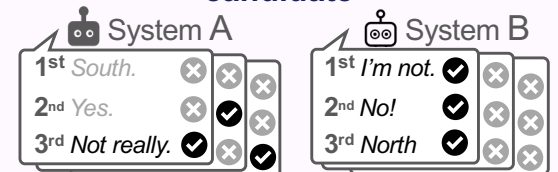
✓'s proportion in n -best should be **high** to provide many available candidates **Leads to better output**

Variety = $\text{Average} \left(\frac{\text{Number of } \checkmark \text{ in the } n\text{-best}}{\text{Size of } n\text{-best having at least one } \checkmark} \right)$

1. Prepare stimulus inputs



2. Detect noncontradictory candidate



yes-no classifier Categorize into yes/no

3. Compute scores

Certainty
3/3 = 1.00

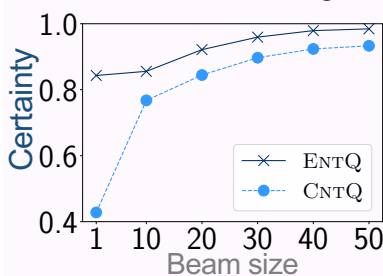
Variety
(.33+.33+.33)/3 = 0.33

Certainty
1/3 = 0.33

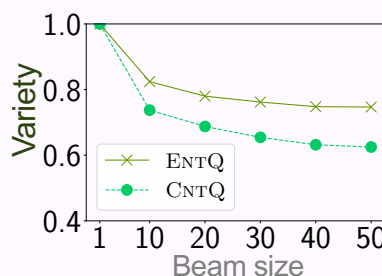
Variety
(1.00)/1 = 1.00

Experiments: Proposed framework reveals n -best's properties

Using **beam search**,



Increases as beam size increases



Decreases as beam size increases

trade-off

Using **newer techniques**,

Technique	Certainty		Variety	
	ENTQ	CNTQ	ENTQ	CNTQ
Beam search	.856	.768	.824	.737
Diverse beam search [Vijayakumar+16]	.999	.981	.758	.478
Nucleus sampling [Holtzman+20]	1.00	.994	.755	.462
Unlikelihood training [Li+20]	.910	.937	.969	.968

Unlikelihood training **improves both scores**

Settings

Generation model: Blender 3B [Roller+21] (See our paper for results of DialoGPT [Zhang+20])

Inputs: 2,000 ENTQ/CNTQ from Multi-Genre NLI [Williams+18], Yes-no classifier: RoBERTa [Liu+19] fine-tuned on Circa [Louis+20]