

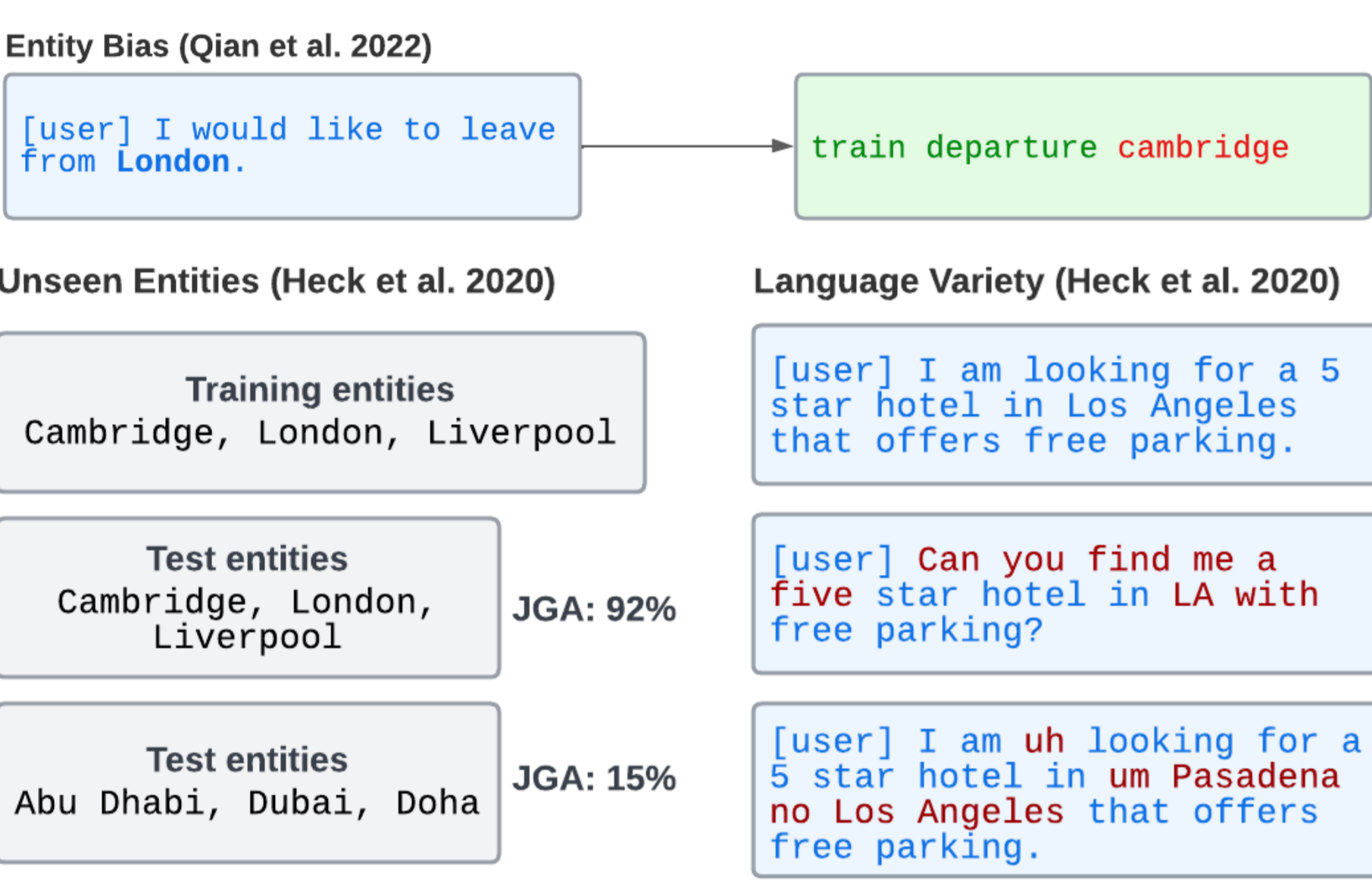
Has it been difficult to compare your dialogue state tracking model beyond accuracy on an in-distribution test set?

Conduct comprehensive dialogue state tracking diagnostics to discover strengths and weaknesses and overlooked opportunities for improvement.

Know Thy Strengths: Comprehensive Dialogue State Tracking Diagnostics

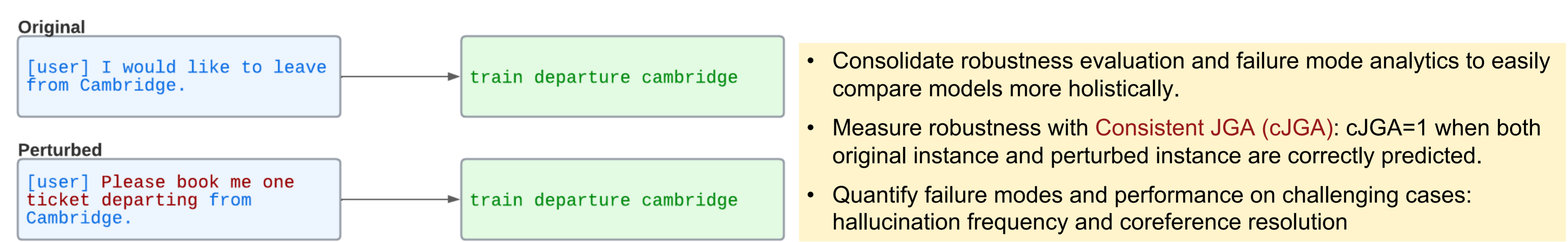
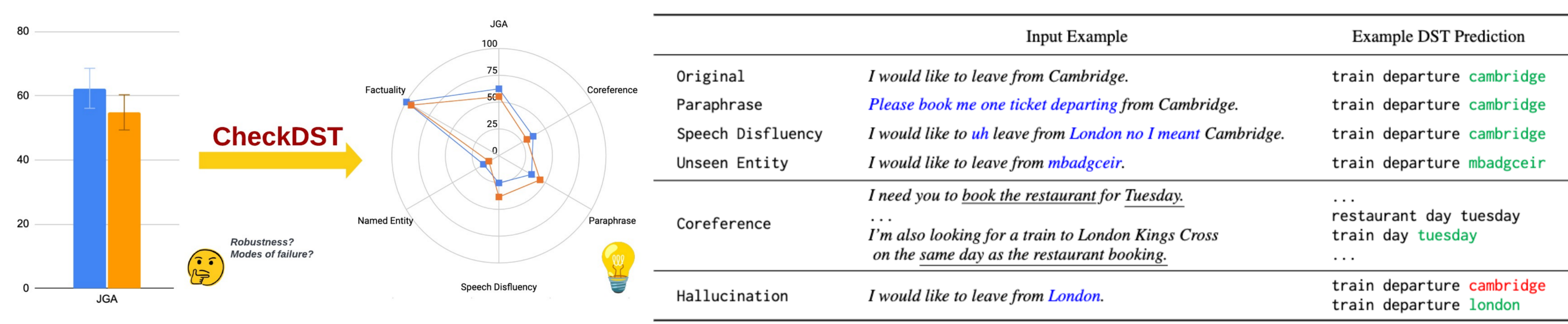
Hyundong Cho, Chinnadhurai Sankar, Christopher Lin, Shahin Shayandeh, Asli Celikyilmaz, Jonathan May, Ahmad Beirami

Problem: sparse and uncoordinated comparisons beyond joint goal accuracy

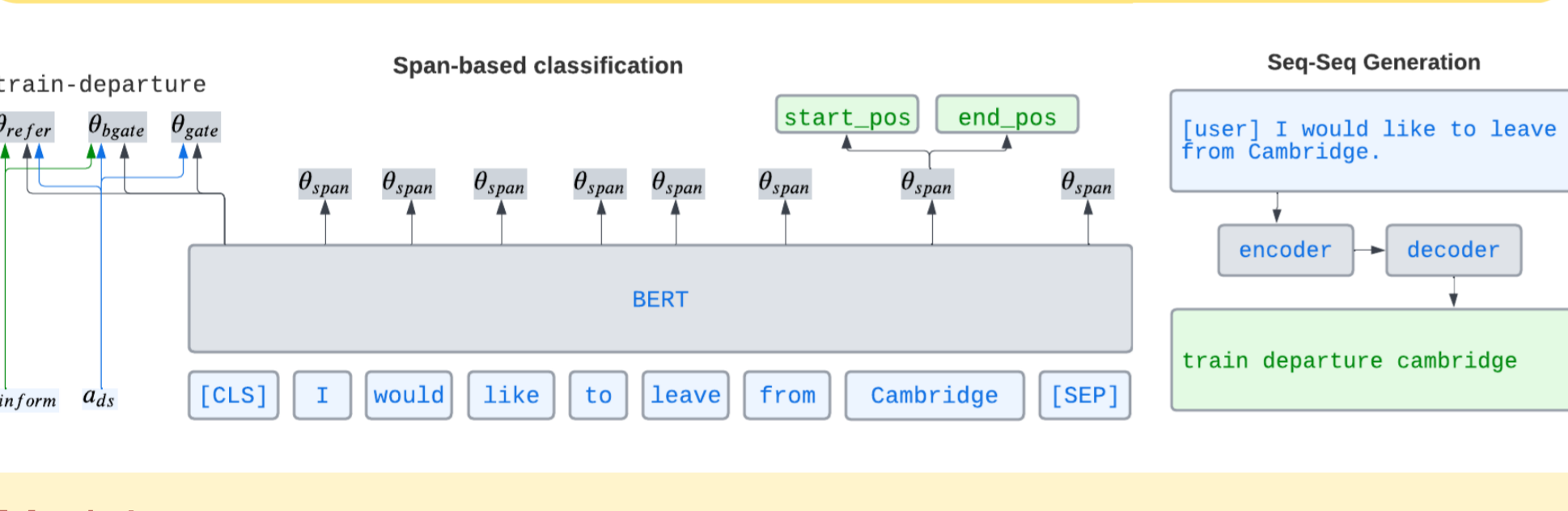


- Joint goal accuracy (JGA) on an in-distribution test set does not capture enough information as dialogue state tracking (DST) models are brittle.
- This is a known problem and previous work have done individual analysis beyond JGA. But they are sparse and uncoordinated, making it difficult to holistically compare strengths and weaknesses of DST models.

Solution: CheckDST
Toolkit that facilitates consolidated robustness evaluation and failure mode analytics



Experimental setup



Models

- Span-based classification models (SCLS)
- SimpleTOD-style Seq-Seq generation models (GEN)

Dataset

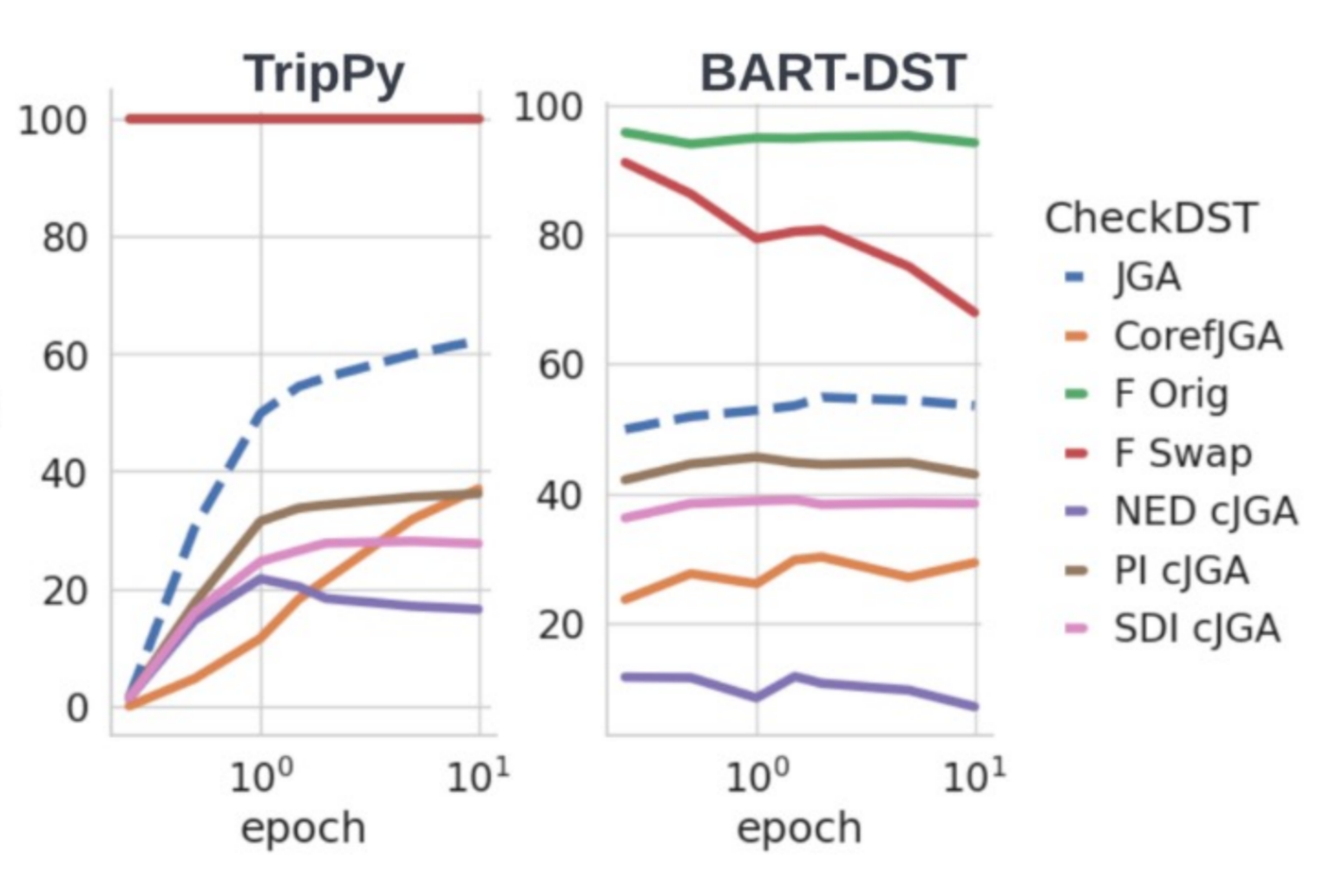
- Original train/val/test: MultiWOZ 2.3 (Han et al. 2020)
- Paraphrase & speech disfluency perturbations from LAUG (Liu et al. 2021)
- Named entity replacements (Huang et al. 2021)

Metrics

- Paraphrase invariant cJGA (PI cJGA)
- Speech disfluency invariant cJGA (SDI cJGA)
- Named Entity Directional cJGA (NED cJGA)
- CorefJGA: JGA on cases that require coreference resolution
- Factuality (F): 1 - hallucination frequency

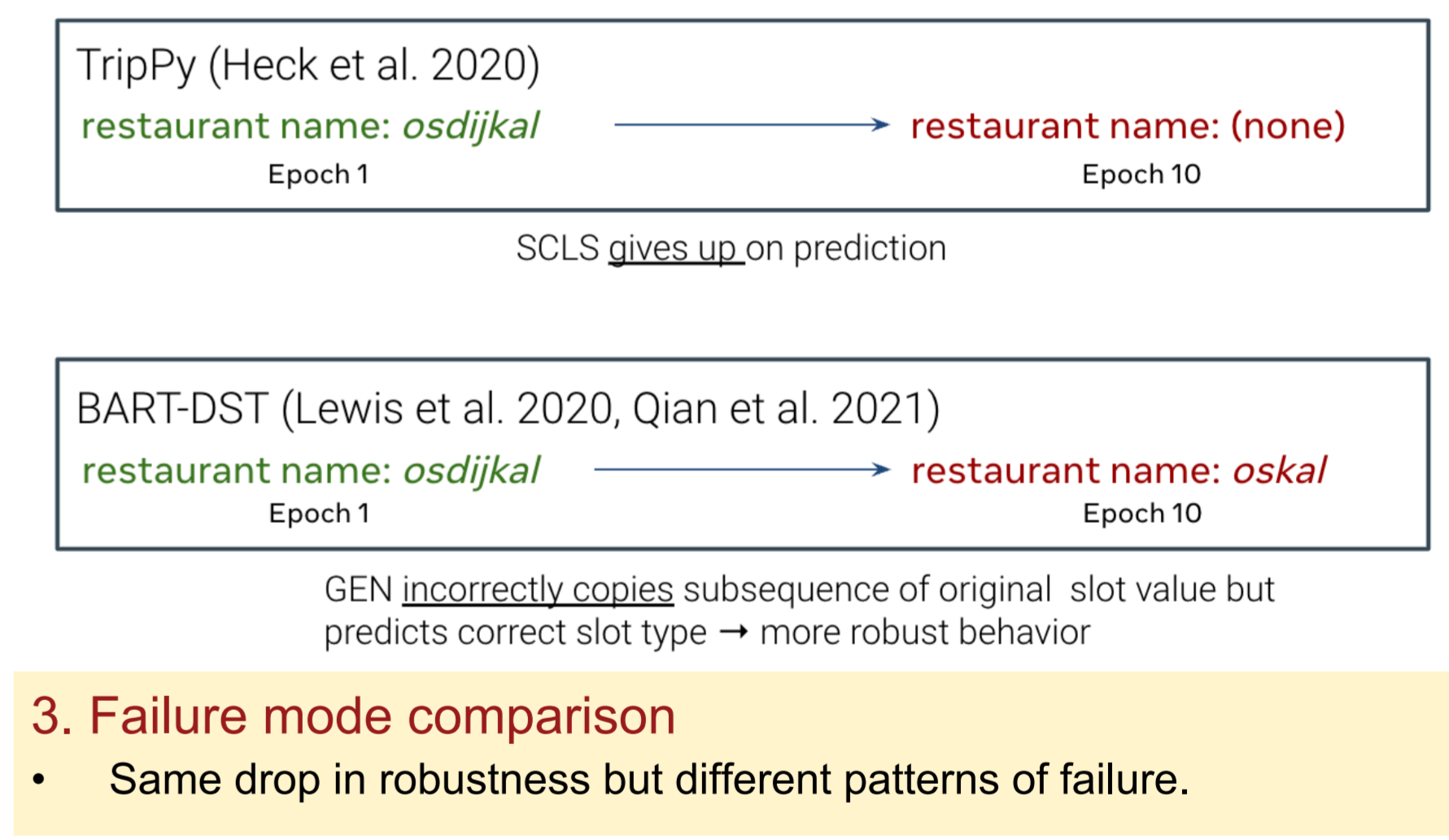
CheckDST quickly quantifies relative strengths and weaknesses between different models and also the same model at different points of training

Model	JGA	CorefJGA	PI cJGA	SDI cJGA	NED cJGA	F _{orig}	F _{swap}
SCLS	TripPy (2020)	62.3 _{0.2}	37.0 _{0.7}	36.2 _{0.1}	27.7 _{0.5}	16.5 _{0.4}	100 ₀
	ConvBERT-DG (2020)	62.0 _{0.1}	36.6 _{0.7}	35.9 _{0.2}	29.5 _{0.4}	16.2 _{0.1}	100 ₀
GEN	BART-DST (2020)	52.5 _{0.2}	25.5 _{1.2}	43.0 _{0.8}	36.5 _{1.3}	9.8 _{0.7}	75.4 _{1.8}
	SOLOIST (2021a)	54.8 _{0.4}	30.4 _{1.1}	44.5 _{0.6}	38.5 _{0.5}	10.7 _{0.4}	81.1 _{1.6}
	MUPPET-DST (2021)	54.9 _{0.4}	29.9 _{1.9}	45.4 _{0.4}	39.1 _{0.7}	7.8 _{1.4}	68.4 _{3.8}

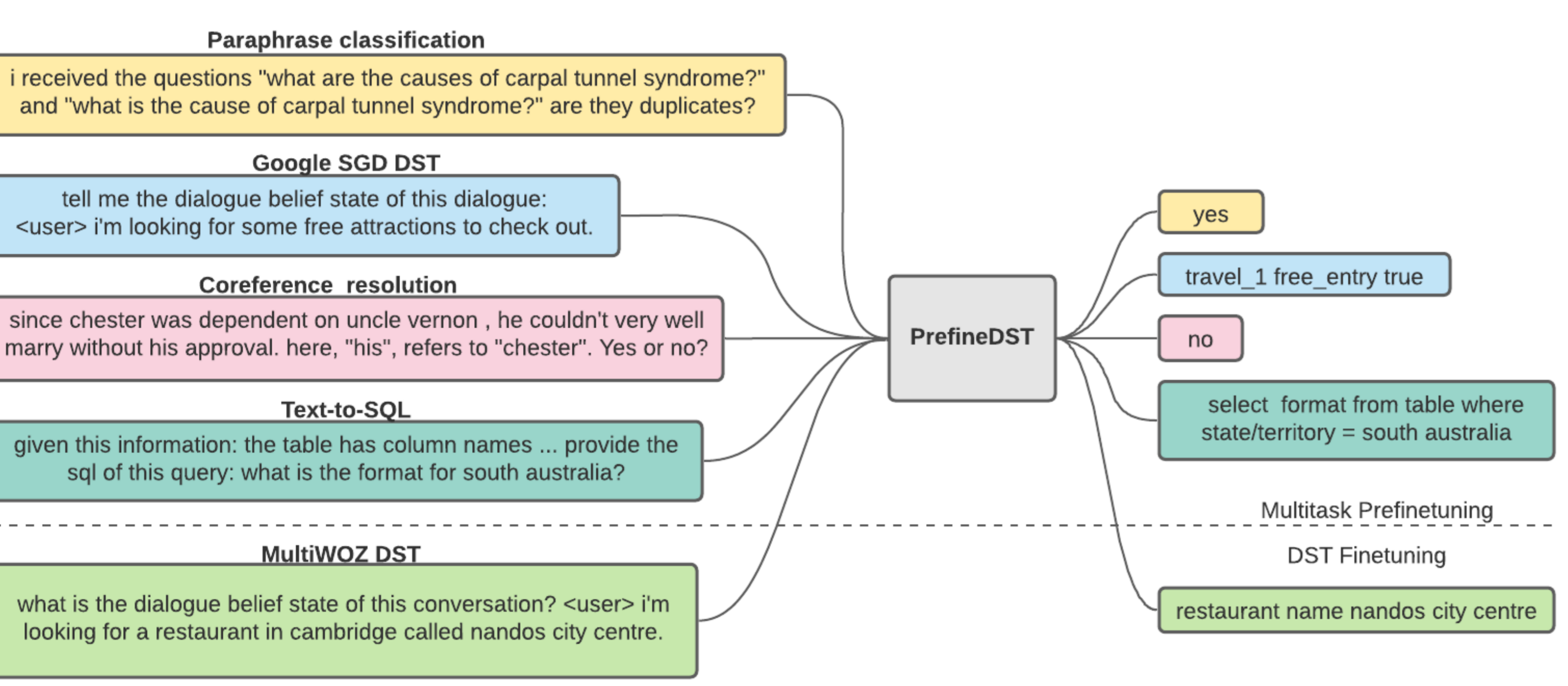


2. Intra-model comparison

- ↑JGA ≠ ↑Robustness
- Both SCLS and GEN models share same trend of **dropping robustness as training progresses**, despite gains in JGA.
- This trend starts even from the start for GEN models and thus encourages few-shot / zero-shot approaches.



CheckDST can guide holistic performance gains! PrefineDST targets exposed weaknesses with pre-finetuning to attain comprehensive improvements.



Model	JGA	CorefJGA	PI cJGA	SDI cJGA	NED cJGA	F _{orig}	F _{swap}
BART-DST	52.5 _{0.2}	25.5 _{1.2}	43.0 _{0.8}	36.5 _{1.3}	9.8 _{0.7}	94.8 _{0.2}	75.4 _{1.8}
SOLOIST	54.8 _{0.4}	30.4 _{1.1}	44.5 _{0.6}	38.5 _{0.5}	10.7 _{0.4}	94.9 _{0.1}	81.1 _{1.6}
MUPPET-DST	54.9 _{0.4}	29.9 _{1.9}	45.4 _{0.4}	39.1 _{0.7}	7.8 _{1.4}	94.6 _{0.1}	68.4 _{3.8}
PrefineDST	55.7_{0.3}	30.5_{0.5}	46.1_{1.0}	41.8_{1.0}	11.0_{0.5}	95.2 _{0.2}	76.7 _{1.4}

Average Δ from Best ↓

TripPy (Heck et al., 2020)	4.55
ConvBERT-DG (Mehri et al., 2020)	4.57
SimpleTOD (Hosseini-Asl et al., 2020)	6.93
SOLOIST (Peng et al., 2021a)	4.98
MUPPET-DST (Aghajanyan et al., 2021)	5.40
PrefineDST	3.92

- CheckDST exposes failure modes and brittleness, **guiding the development of more robust DST models.**
- PrefineDST is pre-finetuned with tasks that intuitively addresses these weaknesses**, such as paraphrase understanding and various parsing tasks that require correctly copying spans from the input.
- PrefineDST shows most well-rounded performance** and shows promise towards robust DST models.

