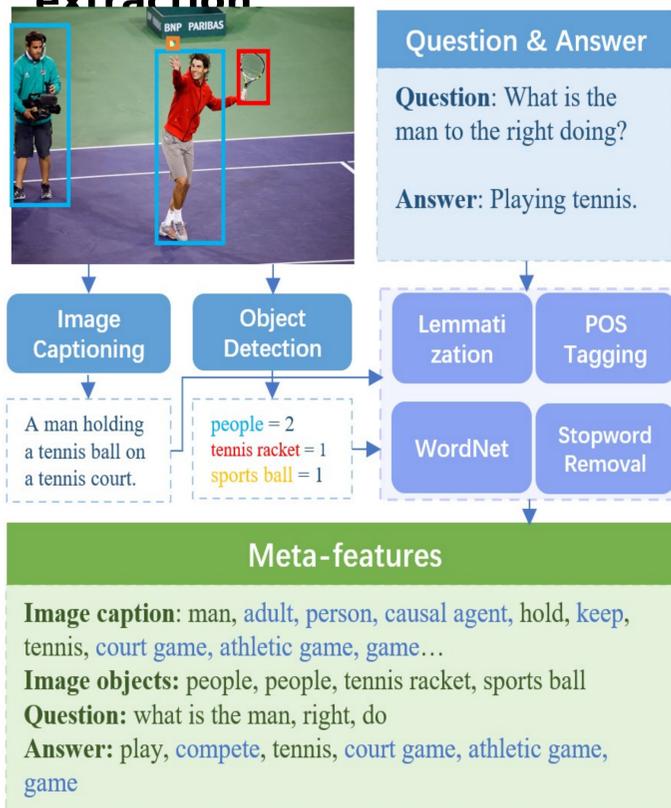


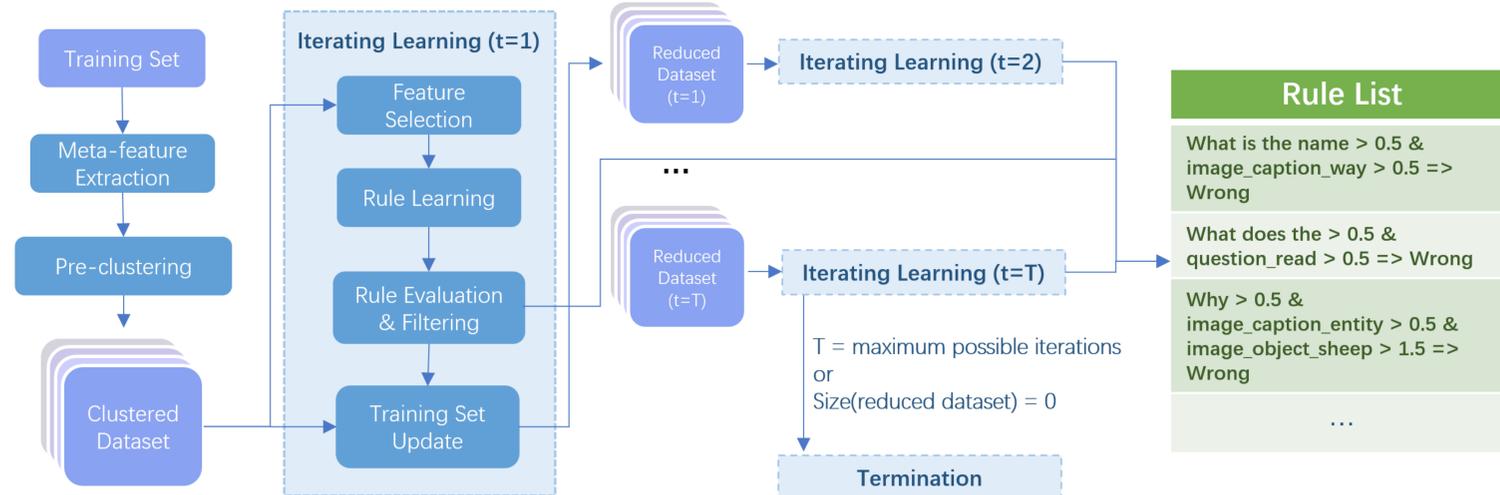
Objectives

- Devise a general method to extract representative rules describing and prediction errors of deep learning models
- Improve models based on the discovered insights

Step 1: Meta-feature extraction



Step 2: Rule Extraction



Step 3: Analyze the Rules

	VQA – OCR	CSQA
Rules	<p>Question type: what Question contains letter Error rate: 82.7% Coverage: 0.404%</p> <p>Question type: what Question contains website Error rate: 92.1% Coverage: 0.0658%</p>	<p>Rule: Answer contains: American State Error rate: 70% Coverage: 3.85%</p>
Examples found	 <p>Question: What letters are on the umbrellas? Image caption: A group of people holding umbrellas Expected: wwf Prediction: white</p>  <p>Question: What is the name of the street? Image caption: A green street sign on the side of a street Expected: flaming lips alley Prediction: main street</p>	<p>Examples</p> <p>Where are the most famous BBQ steakhouses in America? Gold: Texas Prediction: Kansas city Negative Answers: building, Kansas city, Maine, falling down</p> <p>The tourist entered Mammoth cave, what state were they in? Gold: Kentucky Prediction: West Virginia Negative Answers: West Virginia, Rocky Hills, Scotland, Canyon</p> <p>Rule: Answer contains: kitchen Error rate: 100% Coverage:</p> <p>Example</p> <p>From where would you normally take a cup when you're about to get a drink? Gold: kitchen cabinet Prediction: water fountain Negative Answers: dishwasher, water fountain, sand box, toilet</p> <p>Rule: Question contains: chordate, where Error rate: 100%</p> <p>Examples</p> <p>Where does a wild bird usually live? Gold: countryside Prediction: sky Negative Answers: cage, sky, desert, windowsill</p> <p>Where is a bird likely to make its home? Gold: forest Prediction: nest Negative Answers: nest, roof, leaves, sky</p>

4. Improve the Models

	dev	test-dev	test-std
ViLBERt	68.75	69.47	69.64
Ours	69.44	69.64	69.82

Performance comparison on VQA v2.0

Model	Acc
RoBERTa	72.1
Ours	72.45

Performance comparison on Commonsense QA

Conclusion

We presented a novel pipeline that helps automate the error analysis process by learning interpretable rules that characterize the type of mistakes that a system makes. We demonstrated the ability to "close the loop" and use the insight gained from some of the induced rules to make modest improvements to ViLBERt and RoBERTa. These simple but effective approaches have the potential to be applied in production environment shortening the iterative update cycle of models.

Acknowledgments

This research was supported by the DARPA XAI program under a grant from AFRL. We would like to thank Bill Ferguson and his colleagues at Raytheon for inspiring our work.