# Interpretability for Black-Box Seq2Seq Models
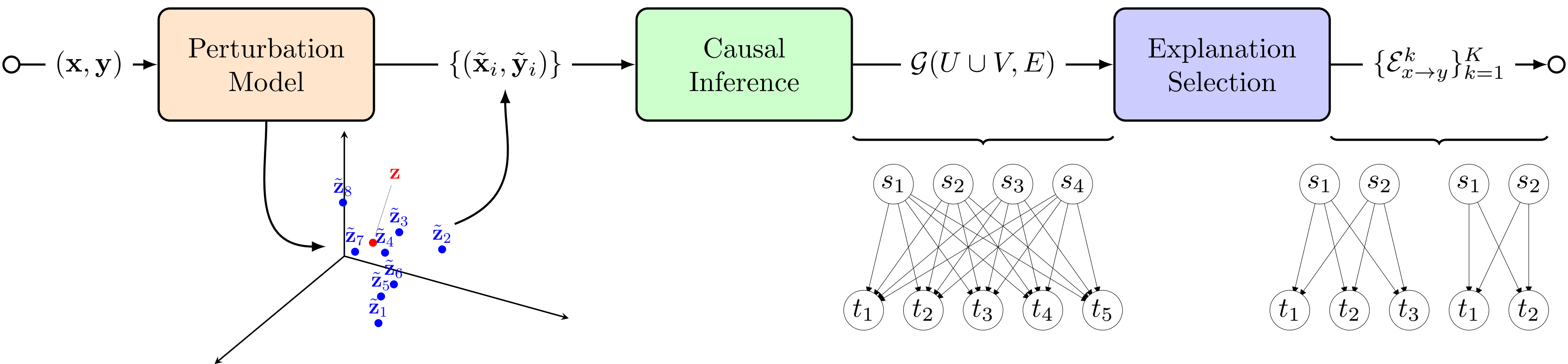
AM+ Jaakkola, *EMNLP'17*

- **Perturb:** Encode -> perturb -> decode -> query

- **Infer:** Logistic regression to infer causal dependencies

- **Select:** Partition dependency graph into *explanation chunks*

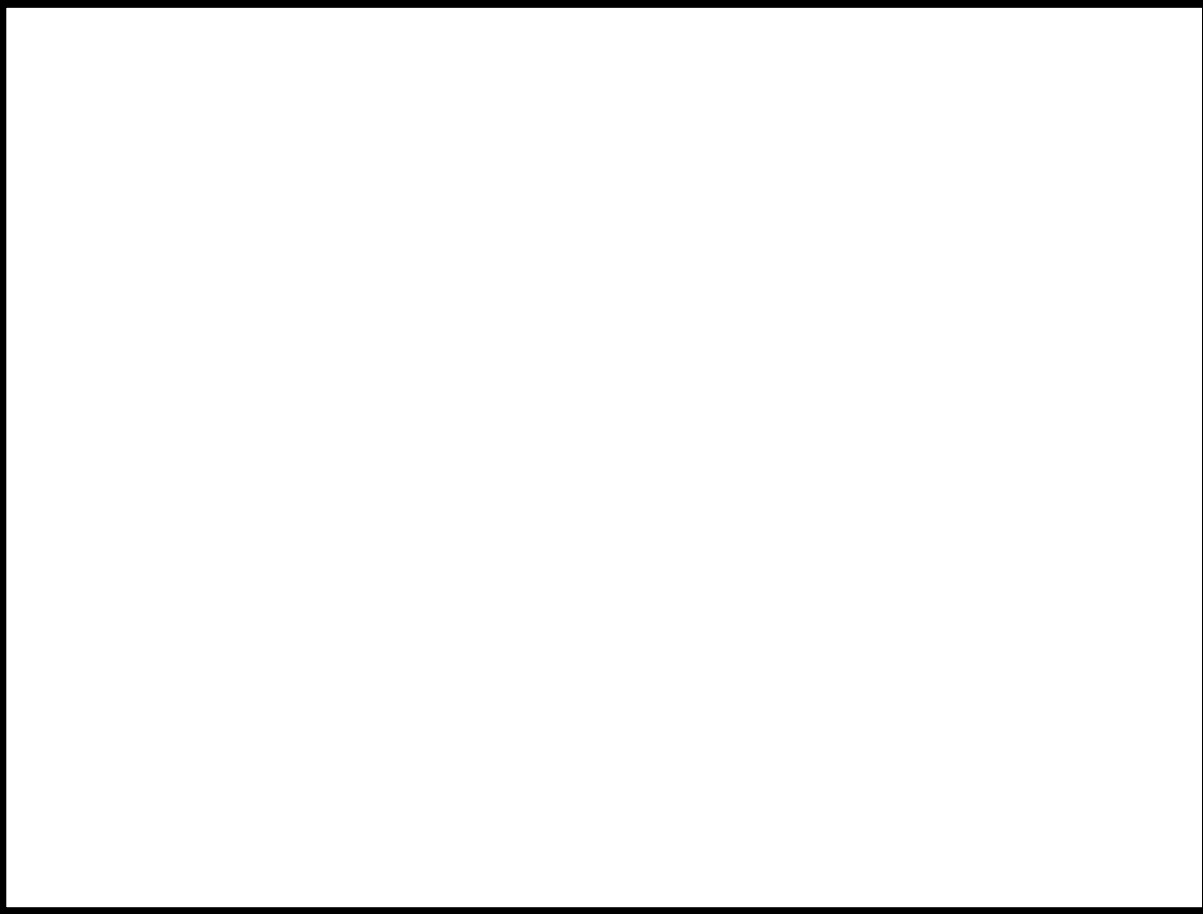# Interpretability for Black-Box Seq2Seq Models

AM+ Jaakkola, *EMNLP'17*



- **Perturb:** Encode -> perturb -> decode -> query

- **Infer:** Logistic regression to infer causal dependencies

- **Select:** Partition dependency graph into *explanation chunks*

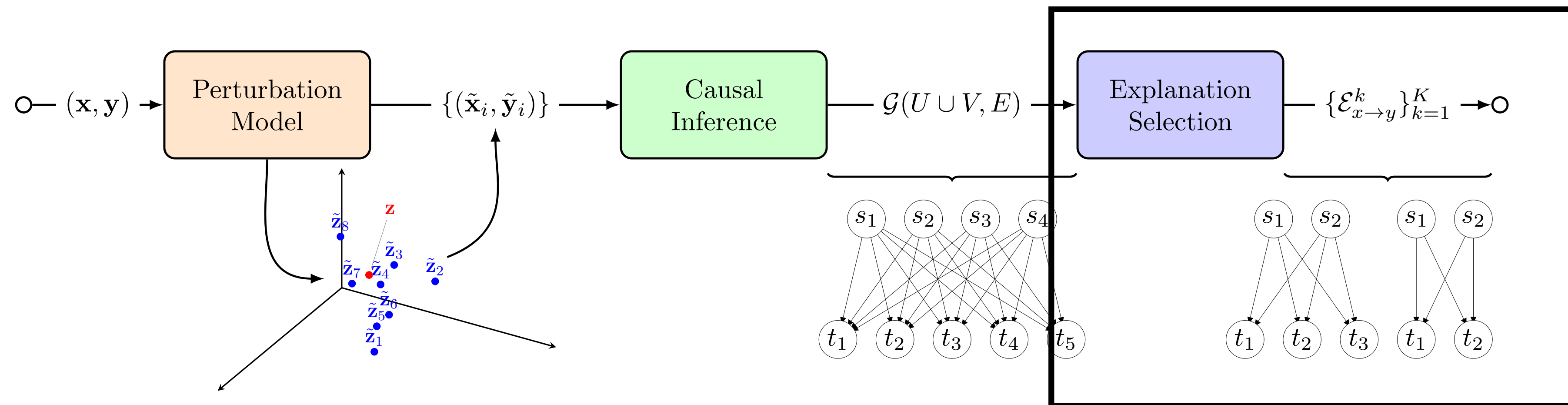# Interpretability for
# Bias Detection in Machine Translation