KNOWING WHEN TO STOP: EVALUATION AND VERIFICATION OF CONFORMITY TO OUTPUT-SIZE SPECS

Chenglong Wang¹, Rudy Bunel², Krishnamurthy Dvijotham³, Po-Sen Huang³, Edward Grefenstette⁴, Pushmeet Kohli³ ¹University of Washington, ²University of Oxford, ³Deepmind, ⁴Facebook AI Research

PROBLEM

We study the <u>termination problem</u> of <u>variable length</u> <u>computing</u> models (Image2Text, Seq2Seq)

Q: Given a model **M**, a sample input **x**, does the model terminates in **K** steps for all inputs $\mathbf{x'} = \mathbf{x} + \boldsymbol{\delta}$?



A: A <u>new testing algorithm</u> and <u>the first verification</u> <u>algorithm</u> to eval robustness of sequence models



MOTIVATION

Achieving Computational Robustness

- ► ML as a service
- Adaptive-time computation model
- Understanding and Debugging Models

Discover abnormal model behaviors (e.g., privacy) Canonical specification for testing variable-length models

ROBUSTNESS OBJECTIVE



Scalable but No Guarantee

$$V^n \mid \sum_{i=1} \mathbb{1}[x_i = x'_i] \le \delta \cdot n\}$$

Complete but Expensive



1. Perturb test inputs randomly



2. Perturb test input with PGD attack



names names names names names names names eos

DeepMind

EXPERIMENTS