# GSA EULA Machine Learning
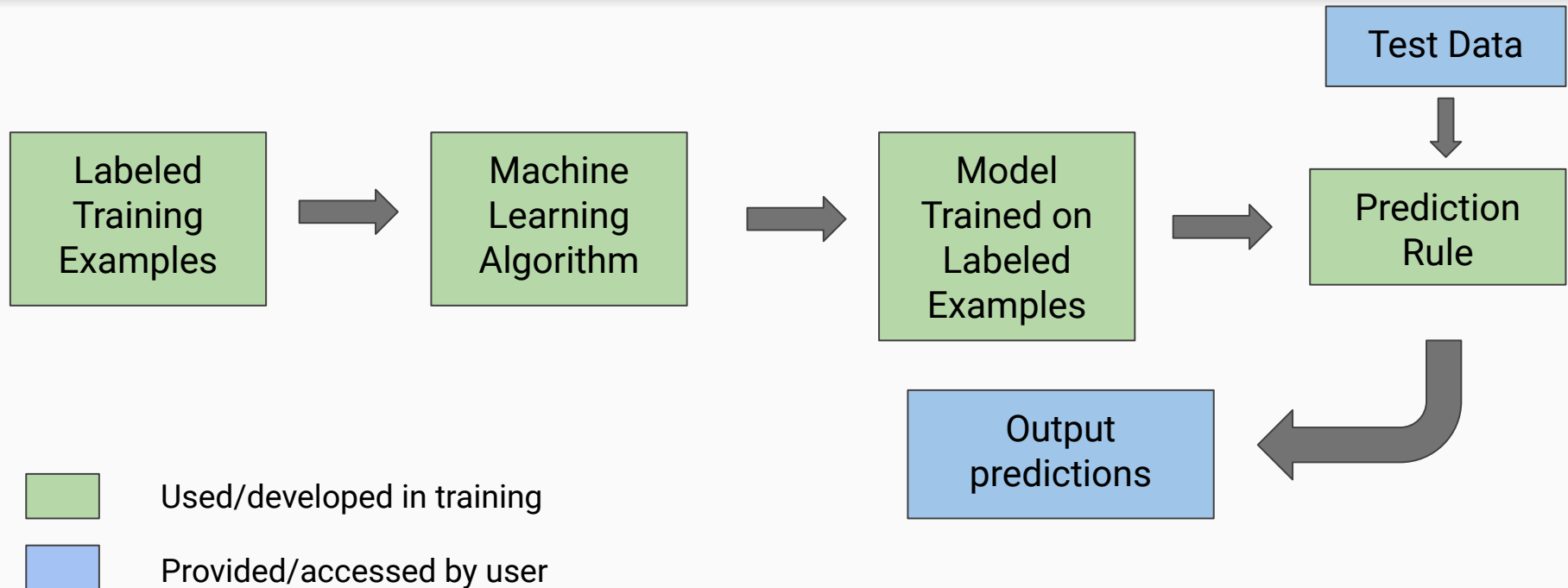
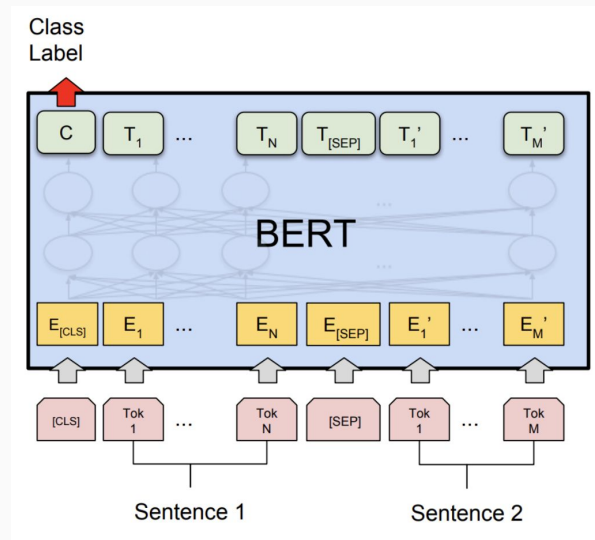Rohan Narain + Andrew Kim

# The Problem

- End User License Agreements must adhere to specific regulations when involving software and products sold to the government.
- Some End User License Agreements contain clauses that may violate regulations, and it can be very difficult to find these clauses.
- By leveraging Artificial Intelligence/Machine Learning, one can detect the presence of clauses that violate regulations and figure out which clauses violate them.

# Supervised Machine Learning Process



Labeled Training Examples → Machine Learning Algorithm → Model Trained on Labeled Examples → Prediction Rule

Test Data → Prediction Rule

Prediction Rule → Output predictions

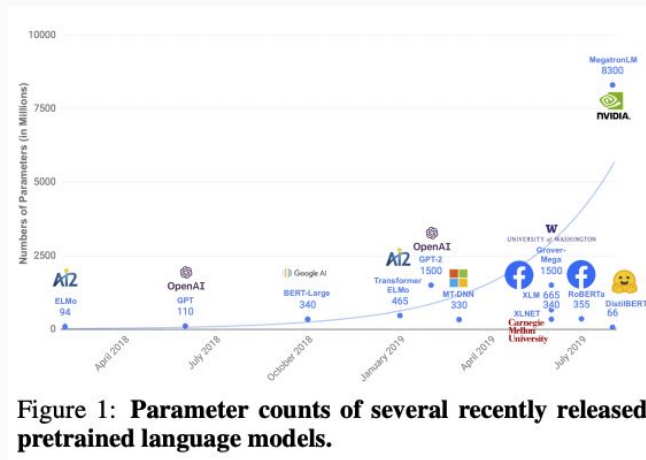Used/developed in training

Provided/accessed by user

# The Model: BERT

- BERT is a transformer model--a deep neural net that efficiently extracts **contextual information** from every part of the input sentence
- Transformer models are the **state-of-the-art** in many natural language processing tasks, including text classification
- Models released by research company HuggingFace are pre-trained on massive text corpus--using a pre-trained model and fine-tuning it is called **transfer learning**

# Scaling BERT: DistilBERT

- In 2019, HuggingFace released DistilBERT, a version of BERT with only 60% of the parameters, but **maintains 95% of the performance** of the original BERT
- DistilBERT makes using BERT in a production system very easy



Figure 1: **Parameter counts of several recently released pretrained language models.**

# Data Format

- Sample data was given in the following format:
  - Three columns: Clause ID, Clause Text, Classification
  - Classification: 0 (acceptable) or 1 (unacceptable)
- Data can be provided by the user as either a Microsoft Word Document in .docx format, or in Portable Document Format (PDF).

# The Solution: EULA-ML

# How Does It Work?

# Model Performance and Explanation

- Weighted F-1 score of 0.87 / 1
- Brier score of 0.09 / 1
- Baseline model scores -- F-1: 0.73, Brier: 0.54
- Non-neural models: XGBoost, SVM, Random Forests…
  - Best F-1: 0.78, Best Brier: 0.18

# Areas of Improvement

- Parsing clauses is imperfect--currently using a simple heuristic, but the formatting of word documents and PDFs can be tough to work with.
- Model retraining
- Adding in explainability
- UI enhacement for sorting clauses by character length, probabilities, predictions (perhaps a table)
- Allow multiple users to provide their feedback on re-labeling our predictions (several stakeholders may have differing opinions)
- Better client side validation of file input (where upload of files happens)

# Thank you!