Data Driven

Training Data

"text":

"The party was absolutely fantastic!"

Instruction Data

"instruction": "Convert to a negative sentiment.",
"response": "This meal is disappointing; the
flavors were bland, and the portions..."

Human Feedback

"It's okay to lie if it helps you get ahead.": 1
"Honesty is always the best policy.": 5

External Knowledge

Theme	Text
Finance	The financial market is experiencing significant volatility, leading to
Sports	The team played an outstanding match, securing a decisive victory that

Model Driven

Classifier, Scorer

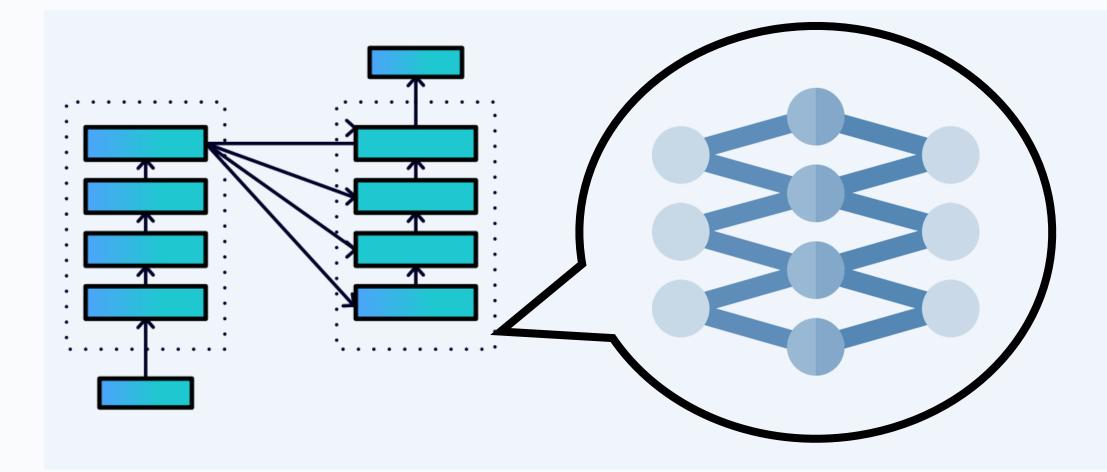
"I love this product": 0.9103,

"This is the worst experience ever.": 0.1734

Class-conditional Language Model

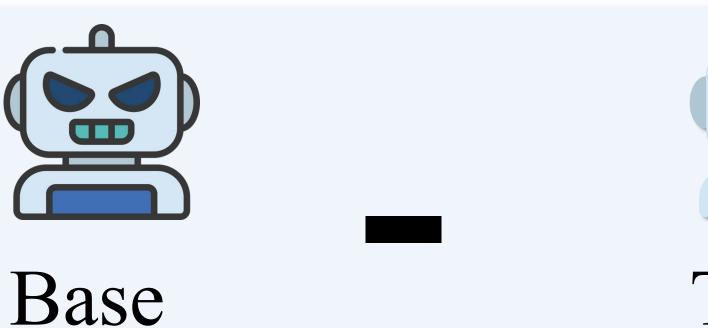
"class": "positive",
 "generated_text": " I love this product."
"class": "negative",
 "generated_text": " I hate this product."

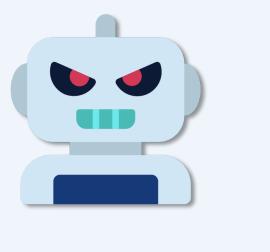
Model Module

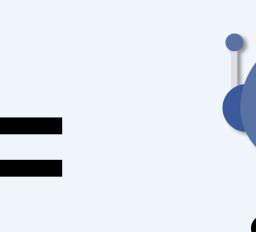


Task-specific module

Model Itself







Toxic

Safe