

Human Centered AI



## **Meet the Founders of Anote**



### Data Scientist at Deloitte Applied Al

Electrical & Computer Engineering | Computer Science @ Cornell University





### Software Engineer at Google for 5+ years

Financial Engineering | Computer Science @ Washington University



## **Presentation Overview**

### **Overview of Anote's Capabilities:**

- A end to end MLOps platform that enables you to obtain the best large language model for your data.
- A evaluation framework to compare zero shot LLMs like GPT, Claude, Llama3 and Mistral, with fine tuned LLMs that are trained on your domain specific training data (via supervised, unsupervised and RLHF fine tuning).
- A data annotation interface to convert raw unstructured data into an LLM ready format, and incorporate subject matter expertise into your training process to improve model accuracies.
- An integration layer to route the best LLM into your own, on premise, private chatbot, via our fine tuning SDK.

### Product Demos of Two Use Cases:

- Generative AI Red Teaming for Chatbots
- Automated CUI Tagging for Text Classification

Notable Customers S&P Global Marubeni Weill Cornell Harvard Medical School US Navy M2 Compliance

#### 5 + Text Data Solutions Platform TASK TYPES CHAT HISTORY Ľ File Uploader Document content infe 🦯 前 Hello, I am your financial assistant, upload a file to get started. MODEL SELECTION Al Engineering Studen 🦯 前 OpenAl what is anote's solution? Text Data Solutions Pla 🦯 🛅 Your own fine-tuned model key: Chat 407 / 亩 Anote's solution consists of three fully functional products and three main steps: Model key / 面 Chat 374 1) Label Text Data - This entails classifying text, extracting entities, and UPLOADED FILES / 亩 Chat 373 answering questions on documents. 2) Fine Tune Model - Allows you to run your PitchDeck\_Anote.pdf / 亩 Chat 372 Fine Tuned LLM locally with their model inference API. 3) Private Chatbot -Al Startup Panacea So 🧳 前 PrivateChatbotDeck.pdf Enables you to chat with your documents while keeping your data private and / 亩 Chat 153 secure. The core technology is a fine-Anote -而 tuning library that leverages state-of-the-Solution.pdf / 亩 Chat 152 art few-shot learning to make high-quality predictions with a few labeled samples. Healthcare Summariza 🥂 💼 The fine-tuning can be done in three ways: 1) Unsupervised Fine Tuning -Chat 19 / 亩 SOURCES Ask your document a question 1 Ť Answer Assistance Sei 🧪 前 PitchDeck\_Anote.pdf V

🖉 Private Chatbot

/**?**?} ~

Download Private Version

# **The Problems**

### **Model Hallucinations:**

As the volume of intelligence reports and classified documents increases, and questions become highly domain-specific, the model often hallucinates, leading to inaccurate or misleading responses.

Ex. "The military operation involved 23,700 personnel" vs. "The military operation involved 2,370 personnel."

### Lack of Domain Specificity:

When asked questions about specific national security domains, such as intelligence gathering or defense strategies, the model's responses are often too vague or imprecise.

Ex. ""The strategy was effective in safeguarding national interests" vs. a detailed explanation of the defense measures used and their outcomes."

### Adversarial Vulnerabilities:

Al systems used in national security are vulnerable to adversarial attacks such as data poisoning, where malicious data corrupts the training process, or evasion attacks, where inputs are altered to deceive the Al into making incorrect predictions. *Ex. "The Al system correctly classified this intelligence report" vs. "The Al system misclassified this report due to adversarial changes in the data."* 

# **Use Case 1 - Generative AI Red Teaming**

### Anote co-led the U.S.-wide red-teaming challenge, NIST's ARIA pilot evaluation of LLM risks:

The first phase of this challenge tested and evaluated the robustness, security, and ethical implications of cutting-edge AI systems through adversarial testing. Participants identified exploits in three proxy scenarios, embodying issues of data exfiltration, bias, and hallucination.

The second phase of this challenge is an in-person exercise to be held alongside <u>CAMLIS</u>, that will include a red team evaluation using generative AI models. The in-person exercise will use the "Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile (NIST AI 600-1)."

https://ai-challenges.nist.gov/aria

## **Our Solution**



#### STEP 1

#### Label Text Data

Classify Text, Extract Entities, and Answer Questions on Documents with LLMs

#### STEP 2

#### **Fine Tune Model**

Run your fine-tuned LLMs locally with our Model Inference API

#### STEP 3

#### **Build Your Own Private Chatbot**

Chat with your documents with LLMs while keeping your data private/secure

### **Data Labeler**

State of the art few shot learning to make high quality predictions with a few labeled samples.

| CSV |
|-----|
|     |

Supervised Fine-tuning Fine-tune your model on your labeled data

| ſ٩ |
|----|
|    |

Unsupervised Fine-tuning Fine-tune your model from your raw unstructured documents



#### RLHF / RLAIF Actively improve your models from

Actively improve your models from human / AI feedback

#### Upload

Create a new text based dataset



#### Annotate

As you annotate a few edge cases, the model actively learns to predict the rest.



#### Customize

Add the categories, entities or questions you care about

| ANOTE                     |          |  |   | <b>A</b> |
|---------------------------|----------|--|---|----------|
| 4- Bank                   | O Upload |  |   | Nest +   |
| Categories                |          |  |   | Dented   |
| sadress #                 | •        |  |   |          |
| IF humilated THEN sadross | •        |  |   |          |
| arger B                   | +        | ( and the contract   | - annua   |          |
| If grouply THEN anger     |          | I care p) from favoling an increasion in an diamond forgative part in<br>actions) increasing who cares and is people |   |          |
| tore #                    |          |  |   |          |
| suprise B                 |          | Lare over failing restrige about the Replace Laid Investment<br>the presents   |   |          |
| har B                     |          |  |   |          |
| I COLLAPSE A              | 4        |  |   |          |
| + All Category            | •.       | An some familing a title burdented lately weeks sure why that w  | mu  |          |
| Add Labeling Fund         | ***      |  |   |          |
|                           |          |  | 1 1 1   |          |
|                           |          |  | And a state of the second s |          |

#### Download

Download the resulting labels as a CSV. Export fine tuned model as an API endpoint

| Inert Intern                       |                          | Maidwin Dauttou           | e Peste Disettant |
|------------------------------------|--------------------------|---------------------------|-------------------|
|                                    |                          |                           |                   |
|                                    | Decorrectation           | • 1=0                     |                   |
|                                    | Inground                 | Coconentation             |                   |
|                                    | Pasture                  | · inprovement             |                   |
|                                    | 144                      | Pastan                    |                   |
|                                    | • **                     | • 1==1                    |                   |
| ingkenent multi languaga support   | <ul> <li>Test</li> </ul> | • he                      | 0.0947            |
|                                    | Constantiation           | <ul> <li>test.</li> </ul> |                   |
| Fis performence degradation listes | • fa                     | Deconantation             |                   |
|                                    | · Feetre                 | •                         |                   |
|                                    |                          | • to:                     |                   |

## **Private Chatbot**

Your accurate private enterprise AI assistant

Upload  $\uparrow$ Upload your financial documents (10-Ks, 10-Qs, Earnings Calls Transcripts)

Chat

Ask Questions on your documents with models like GPT, Claude, Llama2 or Mistral

99 answer is correct to mitigate hallucinations

Evaluate Get Citations for answers, and ensure the



### **Evaluation Results**





| Evaluation Metric Scores <sup>①</sup> |                  |           | <sub>4</sub> 7 |
|---------------------------------------|------------------|-----------|----------------|
| Score                                 | Fine-tuned Model | A\ Claude | 🚳 Open Al      |
| Cosine Similarity Score               | 0.778            | 0.821     | 0.621          |
| Rouge-L                               | 0.824            | 0.901     | 0.780          |
| LLM Evaluation Score                  | 0.824            | 0.901     | 0.780          |
| + ADD A EVALUATION SCORE              |                  |           |                |

## the best model for your data



Each "weak model" runs asynchronously, though to the user they appear to run instantaneously

## **Differentiators**

#### ACCURATE PREDICTIONS

Fine tuning and enhanced RAG for more accurate and tailored predictions. Our AI models actively learns and rapidly improve from SMEs.



#### ACCURATE CITATIONS

Accurate sources (page number, chunk of text, important features) to explain the models predictions and mitigate hallucinations.

#### COMPREHENSIVE CAPABILITIES

Supervised, Unsupervised and RLHF / RLAIF fine tuning for classifying text, extracting entities, answering questions, and chatting with documents



#### EASY TO USE

Accessible UI similar to ChatGPT, and simple SDK for developers where you can input a fine tuned model for improved results.

#### PRIVATE VERSION



Ēq

On premise enterprise-grade solution using Llama2 and GPT4All to leverage LLMs on your unstructured documents while keeping your data local and secure.



#### **EVALUATION FRAMEWORK**

Robust evaluation framework with metrics like Ragas Rouge-L, Cosine Similarity and Answer Relevance to show fine tuned model performance improvements

## Why It Matters

### Before

Manually Labeling their data themselves in a spreadsheet

Tedious, Time Consuming, Costly

Manual Iterative Relabeling

Given a raw unstructured documents, such as a <u>10-k</u> or earnings call transcript, you can't get answers to the questions right if trying to extract info, where accuracy really matters.

Not accurate and largely manual extraction

Sub-optimal analytics for critical business decisions

### After / Anote

State of the art few shot learning to make high quality predictions with a few labeled samples.

Less time, less expensive, higher accuracy

Rapid flexibility for changing business requirements

After a few interventions, we go from 10 questions right, to 15 questions right, to 20 questions right, to enable insights that were otherwise impossible to obtain.

Higher accuracy for raw unstructured documents

New insights that otherwise were not obtainable

# **Use Case 2 - The CUI Problem**

The DOD generates millions of documents (PDFs, PPTX, DOCX, CSV, XLSX, EMAILS) per month that need to be tagged for Controlled Unclassified Information (CUI) content.

**These documents contain sensitive information**, so it is crucial to tag these documents with correct CUI categories to prevent highly classified documents from falling into the wrong people's hands.

https://www.dodcui.mil/CUI-Registry-New/

https://www.challenge.gov/?challenge=cui

## The Problem with the Current Approach



## **Operational Resilience**

**If documents are under classified**, people with not appropriate privilege levels have access to information that they are not allowed to access, which is a major security risk.

**If documents are over classified,** DOD personnel that should be able to collaborate and work on solving critical national security problems are not able to have access, which causes massive inefficiencies in productivity, communications, collaboration and getting things done.

Not solving the CUI problem is a national security risk, and will directly affect our ability to win future wars in the digital age, where information sharing is mission critical.

## Time and Cost Savings

| Team                  | DOD Current Process | Anote's Proposal |
|-----------------------|---------------------|------------------|
| Number of People      | 80                  | 5                |
| Cost Per Person       | \$70,000/year       | \$70,000 / year  |
| Cost of Product POC   | \$0                 | \$X              |
| Total Time to Label   | 6 months            | 10 days          |
| Total Costs Per Month | \$2,800,000         | \$175,000 + \$X  |

## **Requirements for a CUI Solution**



## **Our Solution**





2

HOW TO ACCURATELY TAG CUI CONTENT

### HOW TO IDENTIFY MISLABELED CUI TAGS

### 3 HOW TO SCALE A ROBUST CUI SOLUTION

| Iels: Natan's Fine-tune | d Model X Setfit X |       |           |        |         |
|-------------------------|--------------------|-------|-----------|--------|---------|
|                         |                    |       |           |        |         |
| Classification Report   | Metrics O          |       |           |        |         |
| Model                   | MPC Accuracy       |       | Precision | Recall | Support |
| PROCURE                 | 0.978              | 0.778 | 0.821     | 0.621  | 10      |
| PRVCY                   | 0.924              | 0.824 | 0.901     | 0.780  | 10      |
| LEI                     | 0.846              | 0.946 | 0.702     | 0.924  |         |
| нітн                    | 0.945              | 0.776 | 0.765     | 0.924  | 10      |
|                         |                    |       |           |        |         |



**Evaluation Results** 

| Document Name | Human Label | Model Prediction |
|---------------|-------------|------------------|
| Doc1.txt      | PRVCY       | HLTH             |
| Doc2.pdf      | HLTH        | LE               |
| Doc5.pptx     | нітн        | PROCURE          |
| Doc23.csv     | PROCURE     | LEI              |

# **Statement of Work**

## Phase 1 - MVP

1. Predict categories on test documents from within the predict table view

2. Training the models on training documents to improve the performance via supervised fine tuning

- 3. Evaluate the models performance to see how different models perform on the testing documents
- 4. Baseline recomposition for standard PDFs, DOCx and PDF files

### Phase 2 - Production Ready / Integration

1. Data labeler - incorporate subject matter expertise, and implement the RLHF training to make the models more accurate.

- 2. Core user permissions projects, datasets, models, admin and annotator roles
- 3. Expand to more categories / questions and more documents
- 4. Integrate the fine tuned model with the Private Chatbot for on-premise deployment

BENCHMARKING Q&A MODELS TALK Feb 2024 | Katherine Jijo

## **Thank You!**

Anote is a startup in New York City, helping make artificial intelligence more accessible. We believe there is a massive gap between the tremendous power of AI models, and the everyday tasks that people care about.

### Contact Info



https://anote.ai/



<u>nvidra@anote.ai</u>



https://www.linkedin.com/company/anote-ai/





BENCHMARKING TEXT CLASSIFICATION TALK

FEW SHOT LEARNING TED TALK April 2023 | Natan Vidra



#### IMPROVING RETRIEVAL FOR Q&A TALK

Feb 2024 | Spurthi Setty

Jan 2024 | Katherine Jijo



### FINE TUNING OF LLMS TALK

Jan 2024 | Spurthi Setty