# eduworks

# Threats and Opportunities for Leveraging LLMs

Eduworks Corporation

BENGAL Proposers Day

# Company Summary

- Small U.S. business est. 2001
- Core capabilities in AI, Machine Learning

- Serving multiple USG clients in DOD, DHS, NSF

- Significant experience designing, building, testing and evaluating LLMs

- Research into building and deploying AI for DARPA

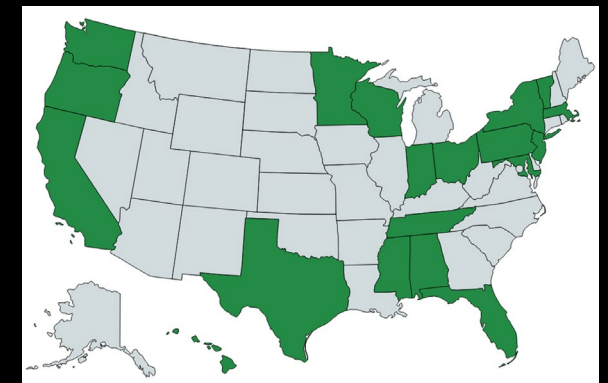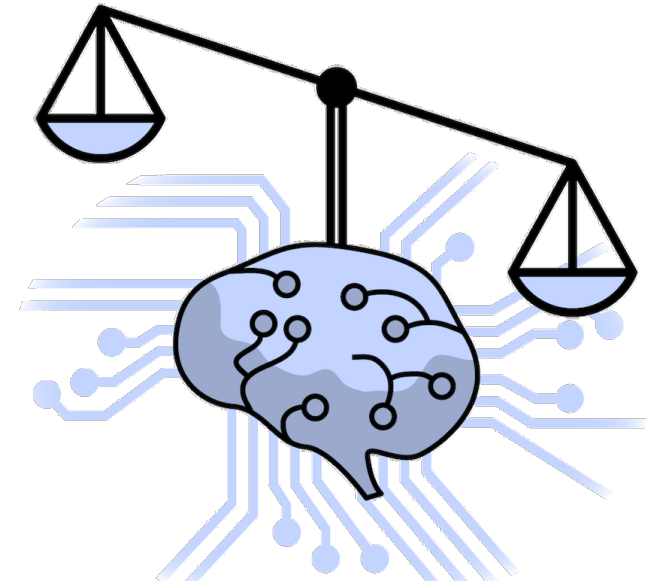- Cutting-edge NSF work on LLM bias, threat and vulnerability detection and mitigation
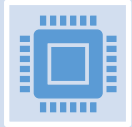
# Goals

- Rapidly summarize/contextualize information relevant to IC

- Avoid biases and toxic or erroneous outputs

- Create novel approaches to address LLM vulnerabilities

# Threats

Injection attacks while training LLMs: nearly impossible to detect with 100% certainty

Inference attacks: can be mitigated with extensive testing, but full protection requires human in the loop

LLMs trained on US/Western data will have biases that cloud responses to foreign threats

Hallucinations, especially from COTS LLMs, will produce responses that appear legitimate but are entirely falsified
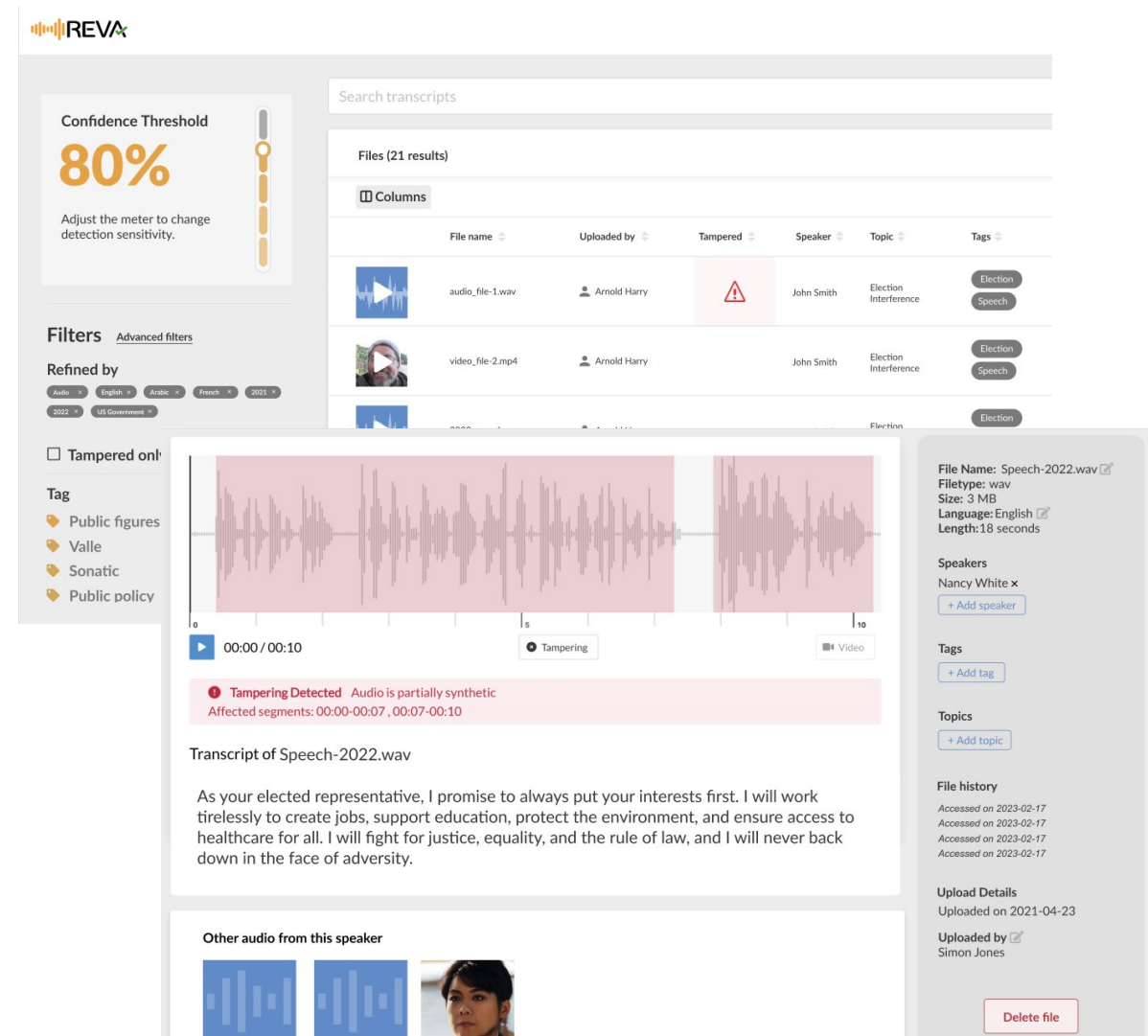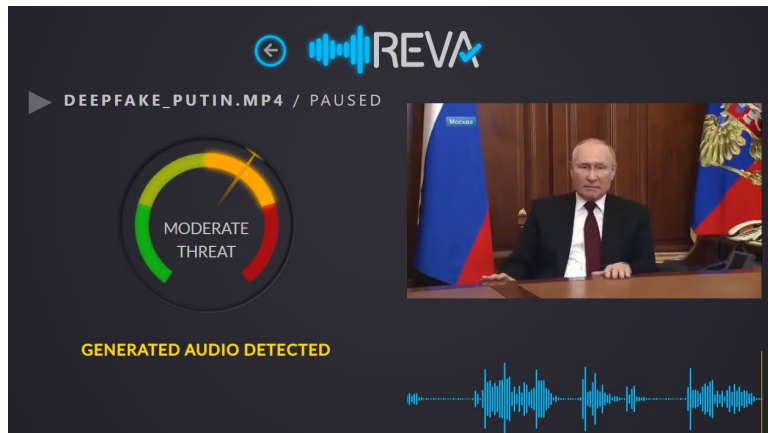
# Current DARPA Work: Countering Deepfakes

## Summary

- Developed by Eduworks for DARPA under SBIR 2017-2024
- Detects GAN-generated speech deepfakes in real time
- Can be used forensically to analyze large media collections
- Uses ML ensemble trained on a large, varied dataset
- Robust to degraded audio conditions (telephone, VOIP)

## Benefits

- Gives analysts better SA than vulnerable legacy techniques
- API can readily provide REVA services to any client system
- Powerful job aid to rapidly triage collected media
- Localization, Diarization, Topic Extraction, Language ID

# Current NSF Work: Safe LLMs

## Summary

- AI alignment of skills with job postings
- Extracts skills and relationships in unstructured text.
- Prioritized list of most relevant skills for a job title or description.
- Scores alignment between KSAs and a course materials.
- Finds and replaces Company Identifiable Information (CII)
- Developing de-biasing techniques for LLM training and tuning

## Benefits

- **Bias Mitigation** via pre-training our LLMs enables mitigation and bias measurement at all levels
- Bias testing pipelines built into each model design
- **Counterfactual Data Augmentation** replaces biased terms with neutral equivalent or multidirectional expansion
- **Bias regularization** debiases embedding during LLM training by minimizing projection of neutral words on relevant axes

# Summary and Benefits

## Use-case focused LLMs:

- Custom LLMs via pretraining, fine tuning, and expanded architecture
- Models that offer confidence scores
- RLHF options to improve accuracy of confidence measures

## Anticipated Benefits:

- Dramatically reduce hallucination
- Minimize costly retraining
- Provide access to real-time data
- Integrate into existing infrastructure

# Thank You

**Contact: Dr. Benjamin Bell**
**COO, Eduworks Corporation**
**benjamin.bell@eduworks.com**