# Safe, Secure and Trustworthy AI...

# Safe, Secure and Trustworthy AI

- AI here and now
- Robustness
- Privacy and Security
- What about generative AI?
- Transparency and ecosystems

**Brookhaven**
National Laboratory

# It's less about Skynet and HAL-9000..

.. than about how AI is used today.

- How AI models are built
- How they are used
- How people *want* to use them
- What to expect from AI
- What you may not consider

# Robustness and Consistency



AI models that "work" may still fail in the real world

- Error sources: noise, class overlap, shift in data distribution, **adversarial attacks**
  - Some attacks require knowledge of model details (*white-box*), other don't (*black-box*)
- Robust training of AI models can *reduce* susceptibility to attacks

**Adversarial attacks on AI model** (*Kotyan 2023*)



**Robust AI training to counter adversarial attacks** (*Hong Wang, BNL 2024*)

# Confidence and Uncertainty

Closely related to robustness (confidence in AI responses, consistency)

Two types of uncertainty:

- **Aleatoric:** data uncertainty (e.g. from measurements), induces bias, propagates through model, *irreducible* without additional data

- **Epistemic**: knowledge/modeling uncertainty, *may be reduced* with improved modeling

Techniques such as Bayesian inference may directly estimate AI uncertainty & robustness

- Challenging to extend to large models



**Data (aleatoric) and modeling (epistemic) uncertainty**
(*Abdar 2021*)



**Bayesian inference for uncertainty quantification**
(*Sanket Jantre, BNL 2023*)

# Private and Secure AI Training and Inference

- Training on private/sensitive data poses challenge
  - Training data leakage
  - Moving sensitive data to AI compute
- Privacy-preserving AI
  - **Differential privacy (DP)** introduces stochastic noise to mask individual data samples
  - BNL demonstrated first distributed DP
- Secure AI training and inference
  - **Fully Homomorphic Encryption (FHE)** enables training without ever exposing secure data



**Extracting training data from deployed model**
*(Carlini 2021)*



**Using FHE to train secure AI models on encrypted data**

# Safety, Security and Trust for Generative AI

Large Language Models (LLMs) and other generative Foundation Models (FMs) pose additional challenges

- Biases and alignment issues
  - May be mitigated *or reinforced* by **RLHF (Reinforcement Learning from Human Feedback)**
  - e.g., personification
- Hallucinations and Verification
  - Imperfect memorization, pseudo-reasoning
  - May be mitigated in part by **Chain-of-Thought (CoT)** reasoning, **self-critique**, **multi-LLM ensembles**, and leveraging external resources, e.g., **Retrieval-Augmented Generation (RAG)**
- Other risks: data poisoning, prompt injections

**Step 0:** LLM

**Step 1:** supervised task tuning

**Step 2:** sample, train proxy reward model

**Step 3:** optimize LLM with RL reward
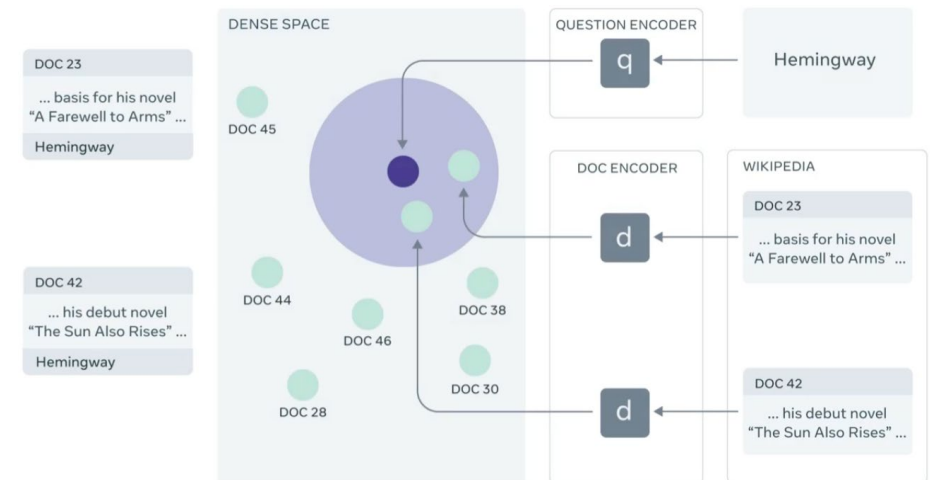
"reinforcement learning from human feedback"

**RLHF is used to align LLM output with human reviewer expectations**



**RAG can map queries and knowledge resources to common embedding space, LLM retrieves relevant context to inject in prompt or response** *(Lewis 2020)*

Brookhaven
National Laboratory

# Transparency and AI Ecosystems

Interpretable and **Explainable AI (XAI)** continues to grow in relevance

- Feature and activation **visualizations** enable per-instance inspection
- Model-agnostic **interpretation** techniques (e.g. LIME, ICE, ALE, SHAP) generally aim to identify feature contributions

Trust, safety, and security of AI ultimately comes down to **use**

- Risks exist throughout **AI ecosystems and lifecycle**



**Layer-wise Relevance Propagation (LRP) visualization for deep NNs**
(*Huang 2021*)



**The AI Lifecycle**
(*Castro 2023*)



**Visualizing feature importance in AI climate model**
(*Wei Xu, BNL 2021*)

**Brookhaven** National Laboratory

# Thank you

csoto@bnl.gov

## References

- Kotyan, Shashank. "A reading survey on adversarial machine learning: Adversarial attacks and their understanding." *arXiv preprint arXiv:2308.03363* (2023).
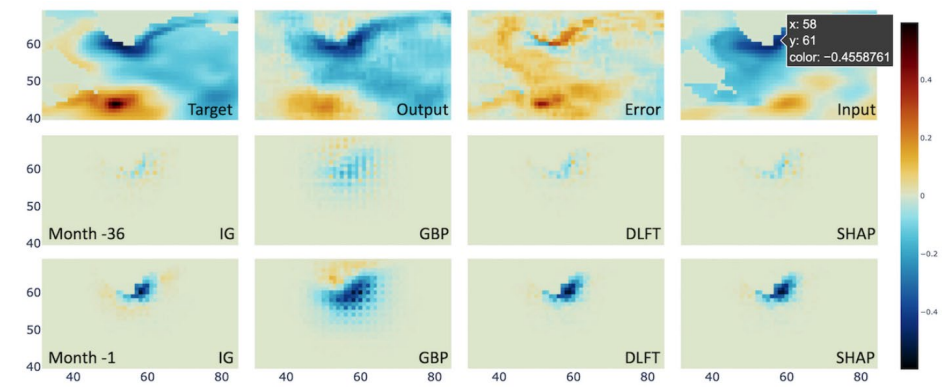
- Wang, Hong, et al. "Exploring robust features for improving adversarial robustness." *IEEE Transactions on Cybernetics* (2024).

- Abdar, Moloud, et al. "A review of uncertainty quantification in deep learning: Techniques, applications and challenges." *Information fusion* 76 (2021): 243-297.

- Jantre, Sanket, et al. "Learning active subspaces for effective and scalable uncertainty quantification in deep neural networks." *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024.

- Carlini, Nicholas, et al. "Extracting training data from large language models." *30th USENIX Security Symposium (USENIX Security 21)*. 2021.

- Huang, Xinyi, et al. "A Visual Designer of Layer-wise Relevance Propagation Models." *Computer Graphics Forum*. Vol. 40. No. 3. 2021.

- Xu, Wei, et al. "Feature importance in a deep learning climate emulator." *Modeling Oceans and Climate Change Workshop at ICLR* 2021.

- Castro, Leyla Jael, et al. "Lifecycle for FAIR Machine Learning." *technical report*, 2023

**Brookhaven** National Laboratory