# Blake Bullwinkel

✉ blakebullwinkel@gmail.com   🌐 blakebullwinkel.com   ⚙ GitHub   in LinkedIn   🎓 Scholar

## EDUCATION

| | |
|---|---|
| **Harvard University** | Cambridge, MA |
| M.S. in Data Science. GPA 3.95/4 | *May 2022* |
| **Williams College** | Williamstown, MA |
| B.A. in Mathematics, Chinese. GPA 3.83/4 (*cum laude*) | *June 2020* |
| **University of Oxford** | Oxford, UK |
| Attended as part of the selective, year-long Williams-Exeter Program at Oxford (WEPO). | *June 2019* |

## PROFESSIONAL EXPERIENCE

**Microsoft**                                                                                        Redmond, WA
*AI Safety Researcher II, AI Red Team*                                           *Jan 2024–Present*
• Developing novel attacks against generative AI systems to uncover safety and security vulnerabilities.
• Active contributor to PyRIT ⬈, an open-source project that automates AI red teaming techniques.
• Led red teaming of the Phi-3 language models including Phi-3-mini, small, medium and MoE.
• Researched gradient-based data exfiltration attacks against Copilots with jailbreak filters.

*Data & Applied Scientist*                                                                 *Aug 2022–Dec 2023*
• Introduced a method to classify performance bugs and customer incidents using text embeddings.
• Built a pipeline to detect and prioritize kernel-mode memory leaks across the Azure fleet.

**Harvard University**                                                                          Cambridge, MA
*Teaching Fellow*                                                                                  *Feb–May 2022*
• Assisted professors in teaching of CS 109b: Advanced Topics in Data Science, a course focused on non-linear statistical methods and deep learning models, including CNNs, RNNs, LSTMs, GANs, and transformers.

## RECENT RESEARCH

**B Bullwinkel** et al. Lessons From Red Teaming 100 Generative AI Products. *Microsoft BlueHat 2024. NeurIPS Workshop on Red Teaming GenAI, 2024.*

**B Bullwinkel** et al. Phi-3 Safety Post-Training: Aligning Language Models with a "Break-Fix" Cycle. *Arxiv 2024.*

**B Bullwinkel** et al. PyRIT: A Framework for Security Risk Identification and Red Teaming in Generative AI Systems. *CAMLIS 2024.*

**B Bullwinkel** et al. Using Large Language Models for Humanitarian Frontline Negotiation: Opportunities and Considerations. *ICML Workshop on the Next Generation of AI Safety, 2024.*

**B Bullwinkel***, D Randle*, P Protopapas, D Sondak. DEQGAN: Learning the Loss Function for PINNs with Generative Adversarial Networks. *ICML Workshop on AI for Science, 2022.*

**B Bullwinkel**, K Grabarz, L Ke, S Gong, C Tanner, J Allen. Evaluating the Fairness Impact of Differentially Private Synthetic Data. *ICML Workshop on Theory and Practice of Differential Privacy, 2022.*

## HONORS AND AWARDS

| | |
|---|---|
| **CES Infinite Mindset Partnership Award** for leading Phi-3 language model red teaming (Microsoft) | *2024* |
| **Quality Stars Award** for building a novel memory leak detection pipeline for Azure (Microsoft) | *2023* |
| **Certificate of Distinction in Teaching** based on student ratings (Harvard University) | *2022* |
| **IACS Student Scholarship** to support data science thesis research (Harvard University) | *2021* |
| **Goldberg Prize in Mathematics** for the best senior mathematics colloquium (Williams College) | *2020* |
| **Linen Prize in Chinese** for achieving distinction in Chinese (Williams College) | *2020* |

## SKILLS

| | |
|---|---|
| **Programming** | Python, R, HTML/CSS, JavaScript, SQL, KQL |
| **Libraries** | NumPy, Pandas, SciPy, Scikit-Learn, PyRIT, HuggingFace, PyTorch, TensorFlow |
| **Platforms** | Azure, AWS, Docker, Linux, Windows |
| **Language** | Working proficiency in written and spoken Chinese (Mandarin) |