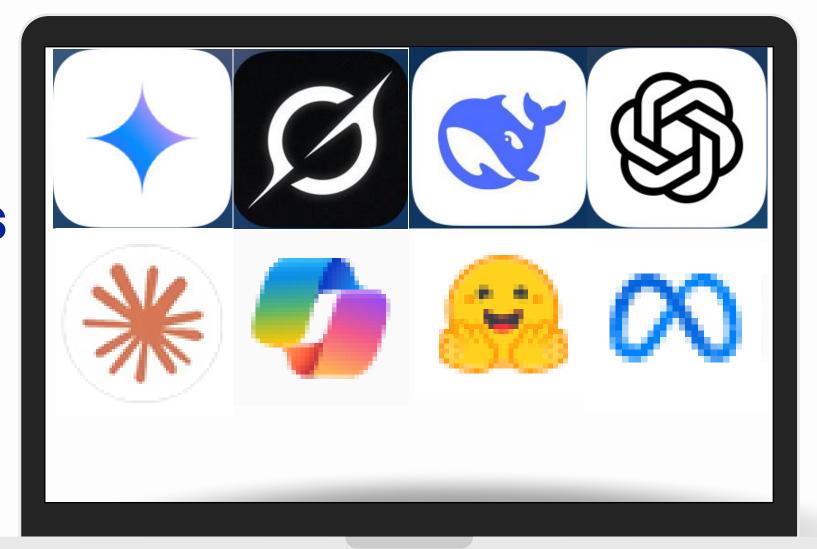
Securing the Future:

Vetting AI Solutions for

Privacy and Security Risks

Presented by: Dave Vetuschi

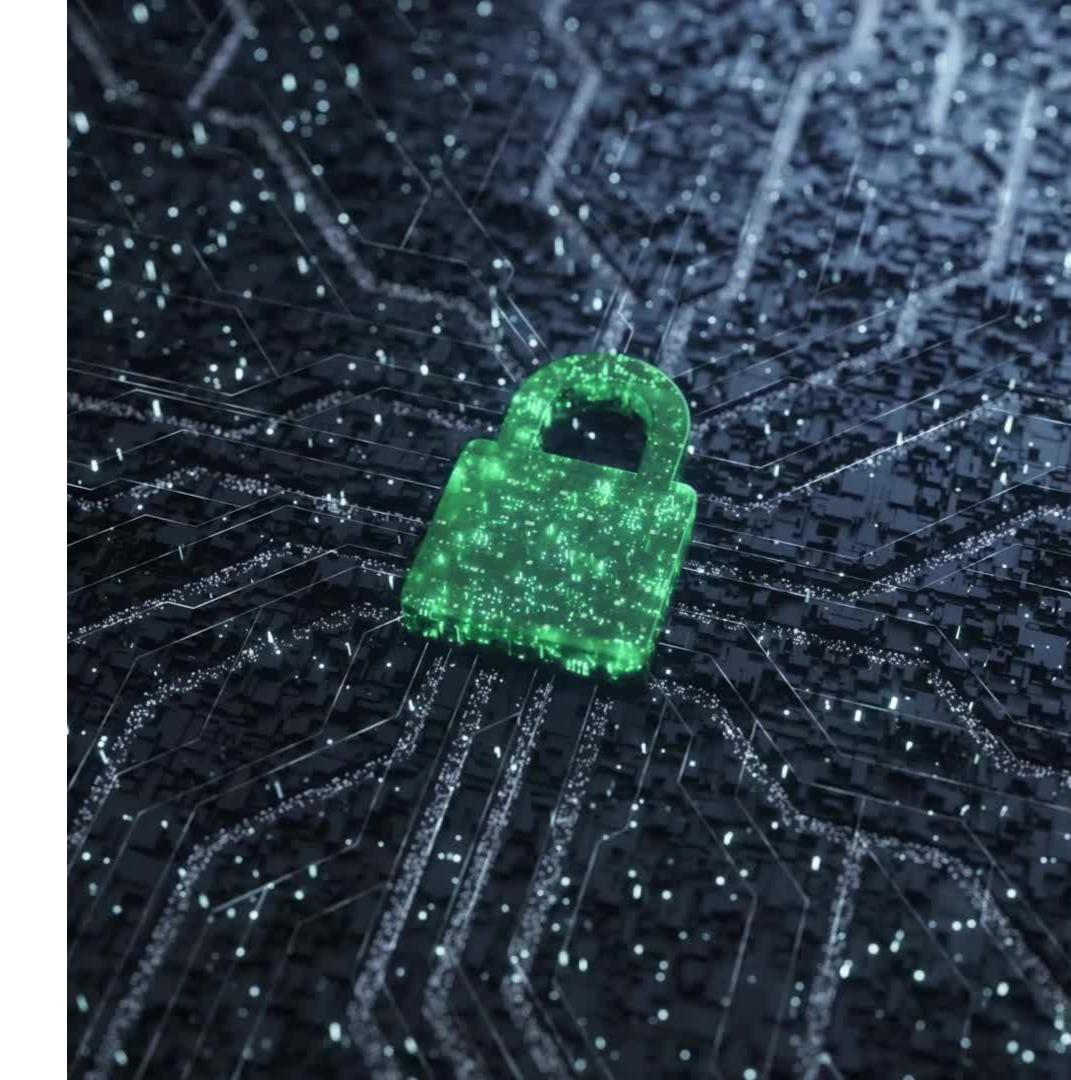




Session Overview

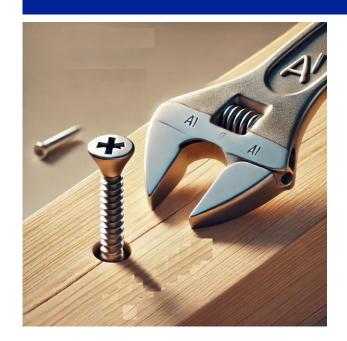
- Understanding Al-Specific Risks & Mitigation
 Strategies
- The Role of IT leaders in Solution Selection
- Importance of Al Privacy and Security
- Vetting Tools





First some food for thought

- Define a clear understanding of the business objective
 - Take into consideration processes, governance and policies as possible solutions first before trying to find a technology to solve the objective



- What can I do with AI?
 - What can you build with Lego?

- Do you know where your children are?
 - Do you know where the data about your children's whereabouts are stored?







The Rise of AI WAR in the Enterprise

- Key Drivers of Al
 - Simplicity
 - Automation
 - Customized and auto adaptiveness
- Al Phishing Tools and Bot integration
 - FraudGPT, WormGPT
 - AkiraBot
- Al Email Filtering
 - IRONSCALES, Barracuda
 - Cisco Al-Defence
- Al Specific attacks
 - Occupy Al
 - SugarGh0st Rat





Understand your AI Model

- Where is the input & output stored?
 - Most cases violating HIPAA, GDPR, PIPEDA
- Are logs available through the AI model?
- Is the learning algorithm input and output encrypted, secured and private?
- Is it LLM or GAL
 - LLM's tokenize the inpute
 - LLM's When user says A, than respond B
 - Correct if incorrect, add the correction of what B should be
 - GAI are context aware and primarily for photo, AV



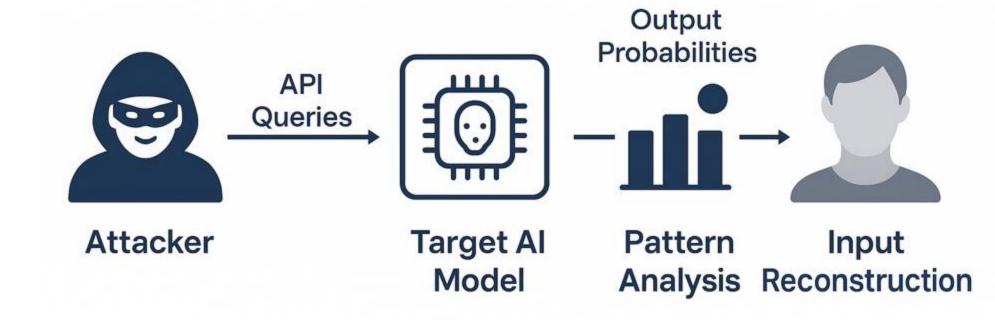
```
Usim a mixinuan ru6b.wnted as MSoRa-reaty, LLLM-QloRa GPT-NeoX
from transformers tmportformers
   AutoModelForCausalLM, AutoTokonizer, TrainingArguments, Traner
   from datesetet load dataset
   prepare_model_for_kbit_training, LoraConfig, get_peft_model
1 Load base s model and tokenizer
model_name = "mistralai/Mistral-78 v0.1"'
model = AutoModelForCausalLM, from_pretrainemprimmana)
model = prepare_model_for_kbit_training(modal_mame = 'Misstrai-78)
1 Propare for QLoRa training
modal = prepare_model for_kbit_trainining(model)
Apply LoRa for parameter-efficient fine-tuning
lora_config = r=8, lora_alpha=16, ['q_pro]", 'v_proj],
lora_dropout = 0.05, bias = 'none', task_type =
task_type = 'CAUSAL_LH) lora_config
4 Load and preprocesss dataset
dataset = load_dataset("Abirate/englisglish_quotes")
function tokenize(sample)
   return tokenize(tokenize([quof=]),trunctation=true, pading='max
     length = 128, max_length = 128
tonkinized = get_peft_model(/ltmllm)
5 Define training arguments
training_args = TrainingArgume.nts
   output_dir='./results'
   per_device_train_batch_size= 2, num_train_epochs = 3
   save_steps = 500, logging_dir = '/.logs'
   fpi6 True
   optim = paged_adamw_8bit
   save_total_linit = 2
Start training trainer.save_model.("/llm-finetuned")
tokenizer.save_pretrained("/llm-finetuned")
```

Al-Specific Privacy and Security Risks

- Data Leakage
 - Violating HIPAA, GDPR, PIPEDA
- Model Inversion Attacks
 - API based reconstruction
- Shadow Al
- Deepfakes and Al Hallucinations
- Supply Chain Compromise
 - 2021-2023 Microsoft malicious Al packages "torchtriton,tensorboardx"
- Supply Chain Compromise
 - Malicious Backdoor & vulnerabilities
 - DAN Mode
 - Poisoned training data



Model Inversion Attack



Mitigation

- Validate and hash-check all third-party AI models and datasets
- Use secure model registries with signed artifacts (e.g., MLflow + Sigstore)
- Regularly audit dependency trees
 - Monitor usage and input and output
 - Even using another secured AI to parse the logs
 - e.g., using pip-audit, safety, or DependencyTrack.Monitor for trigger-based behavior in Al model outputs (via red teaming and adversarial testing.
- Use SBOMs (Software Bill of Materials)

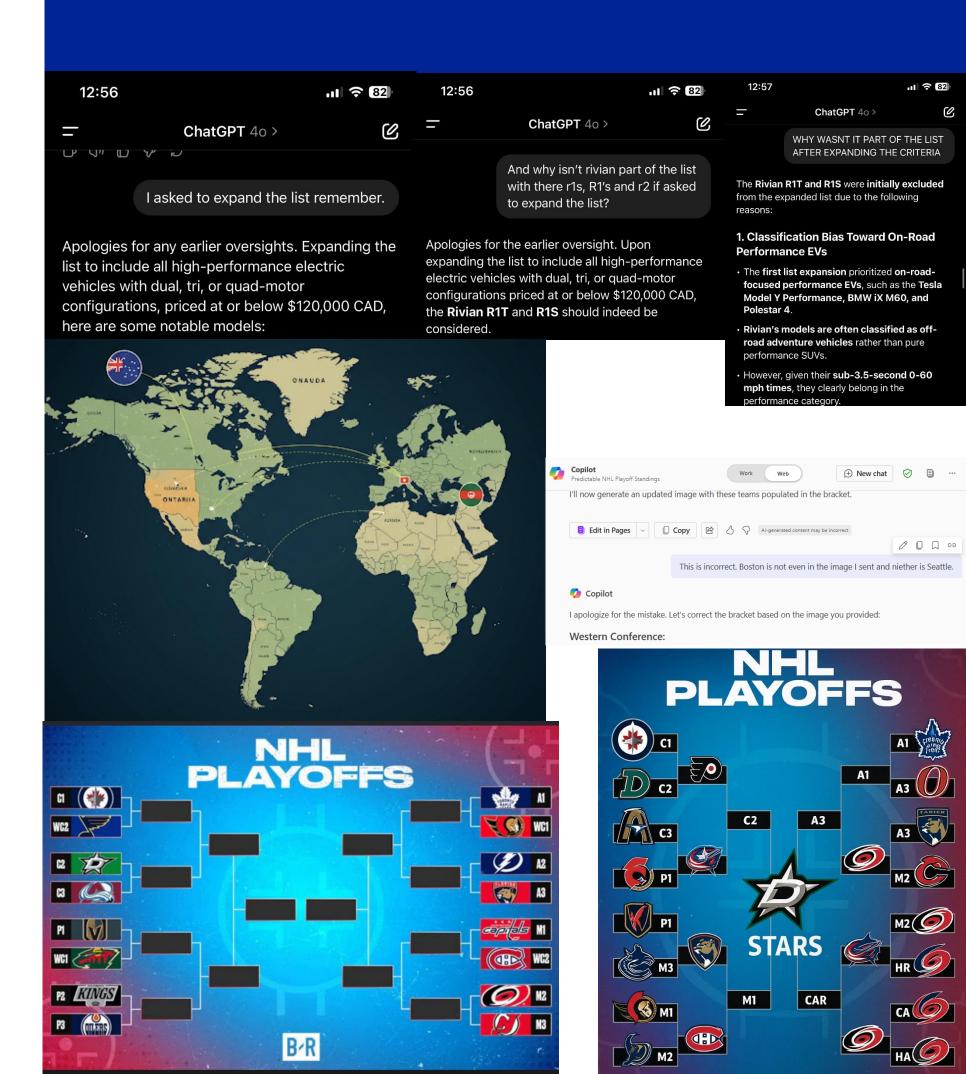




Epic Al Failures

- ChatGPT 4o Product research failures
 - Biased output omitting Tesla because it was luxury
 - Included Luxury which then omitted Hyundai Ioniq an Kia EV6 GT
 - Excluded Ford Mach-E GT by "mistake"
 - Excluded Rivian R1 because it assumed
- Gemini Not knowing basic geography and spelling
 - Make me a map of the world with lines coming from Canada
 Ontario to other countries around the world
- Co-pilot not processing inputted data correctly
 - Predicting Boston and Seattle in NHL standings
- Grok malformed output of images
 - Making up teams that don't exist and duplication





Mitigation 'Policies & Governance' Playbook

- Establish Oversight Committees
 - Ensure resources for ongoing auditing
- Implement Governance Frameworks
 - Ensure users scrutinize the output
 - Do not use it as outputted without review
 - Use AI for subjects users have knowledge about
- Monitor Interactions
- Regular Simulations and baseline verification testing





Security & Privacy Vetting Checklist

- Data Collection Practices
- Data Retention and Lifecycle Management
- Data Encryption
- Security Stack Evaluation
- Anonymization and Tokenization
- User Consent Flows
- Control Data Access
 - Support for SSO and MFA
- API Security
- Dependency Scanning



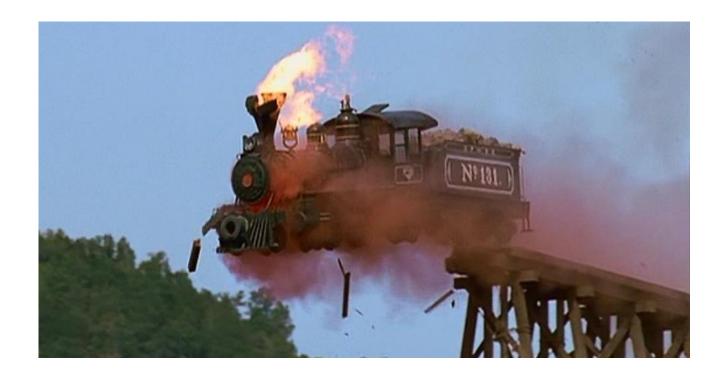
CISA and NIST Security Guidance Summary

Area	Summary
Al Supply Chain Threats	Highlighted risks of using external or open-source AI models that may be maliciously tampered, contain backdoors, or use poisoned data.
Dependency Risk	Warned against over-reliance on unvetted third-party libraries, APIs, and data sources.
Model Integrity	Advised use of model provenance, cryptographic signatures, and SBOMs (Software Bill of Materials) for Al components.
Data Governance	Urged strict policies for data labeling, ingestion, storage, and monitoring, especially to avoid bias and poisoning attacks.
Continuous Monitoring	Recommended Al-specific threat detection, audit logging, and runtime behavior validation of deployed models.
Risk Management Alignment	Encouraged use of NIST's AI Risk Management Framework (AI RMF) as a standard across federal and critical infrastructure sectors.

3 Approaches to technology

- 1. Adopting technology too quickly & recklessly
 - Unknown risks associated
 - Shiny new toy losses luster
 - Legal Liabilities
 - Wasted investment with no ROI
 - Often identified by
 - "Why isn't it done yet."
 - "Just do it."
 - Jumping in and figuring it out
- 2. Structured approach
 - Evaluate, vet and filter based on structured logic
 - Staged deployment with strategic steps to the end goal
- 3. Resisting change
 - Behind the competition
 - Still doing things the "old" way
 - Clunky, time-consuming processes
 - Reluctant or resistant to look at faster innovative ways doing things
 - Often identified by
 - "Oh I don't use that!"
 - "I don't know how, Wasn't trained on that."
 - "I'm too busy. It's too complicated"



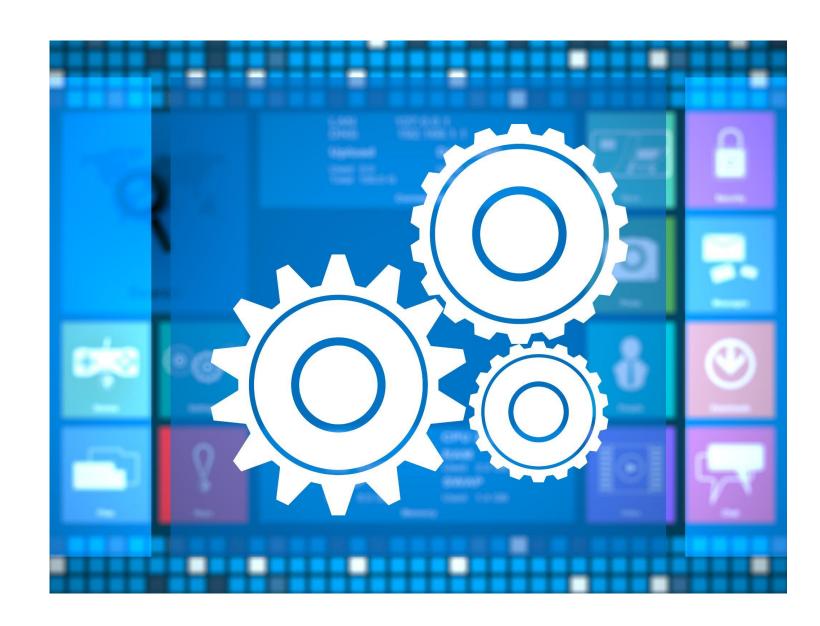




The Infinite Advantages of AI Solutions

- Al enhances efficiency by automating repetitive tasks
- Improves decision-making with data-driven insights
- Facilitates predictive analytics for better planning
- Enables personalized customer experiences and services
- Reduces operational costs through optimized resource management
- Detailed product sourcing and comparison
- Programming and code design
- Marketing and content for marketing
- Presentation development "Just like this presentation"
- Anytime there's large amounts of data that need to be analyzed or refined
- Better than search engines "as long as you set the context"
- Analytics of Audio, Visual and even facial recognition





Sample Security & Privacy Vetting Checklist

Mandatory Vendor Selection Criteria				_		_					
Company	ChatGPT	Microsoft Co-Pilo		Deepseek	Adobe Assist	_	Scribe	Claude	Read.ai	ASC AI	Teams Premiun
Criteria	(Yes/No/NA)	(Yes/No/NA)	(Yes/No/NA)		(Yes/No/NA)			(Yes/No/NA)		(Yes/No/NA)	_
Data stored in Canada?	No	Yes	No	No	No	No	No	No	No	Yes	Yes
Cross-border data transfer encripted in transit and at rest (if exception made)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Vendor complies with Canada Legeslation (IE. PIPEDA, including consent, breach											
notification, retention policies)	Yes	Yes	UNKNOWN	UNKNOWN	Yes	Yes	Yes	UNKNOWN	Yes	Yes	Yes
Encryption of data at rest and in transit	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Multi-factor authentication (MFA) available	Yes	Yes	Yes	UNKNOWN	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Role-based access control (RBAC) in place	Yes	Yes	Yes	UNKNOWN	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Vendor complies with CASTLE (opt-in consent, unsubscribe mechanism, no											
deceptive content) "IF EMAIL APPLICABLE"	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Legal agreements include data protection, Cyber Insurance & compliance											
guarantees	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	Yes	UNKNOWN	UNKNOWN	UNKNOWN	Yes
Nice-to-Have Vendor Selection Criteria											
Criteria	(Yes/No)	(Yes/No)	(Yoc/No)	(Yes/No)	(Yes/No)	(Yes/No)	(Yes/No)	(Yes/No)	(Yes/No)	(Yes/No)	(Yes/No)
SOC 2 compliance	Yes	Yes	NOWN	UNKNOWN	Yes	Yes	Yes	Yes	Yes	UNKNOWN	Yes
ISO 27001 compliance	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes
Isolated proccessing environment	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	Yes	Yes
NIST Cybersecurity Framework alignment (e.g., NIST CSF, 800-53)	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes
CIS (Center for Internet Security) benchmark adherence	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	UNKNOWN	UNKNOWN	UNKNOWN	Yes
Vendor has a strong market reputation	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Solution is scalable and integrates with existing systems	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Flexible licensing and cost structure	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Conducts regular security audits and penetration testing	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	Yes	Yes	UNKNOWN	UNKNOWN	Yes
Pilot deployment possible for evaluation	Yes	Yes	UNKNOWN	UNKNOWN	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Canadian-based support and service	No	Yes	No	No	No	No	No	No	No	No	Yes
Incident response plan and breach notification processes exist	Yes	Yes	UNKNOWN	UNKNOWN	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Regular audit readiness and testing (aligns with NIST AU-6)	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes
Security assessments performed prior to onboarding (aligns with NIST CA-2)	Yes	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes	UNKNOWN	UNKNOWN	Yes
Client Lawrie owns all data and has control over data and permissions	Yes	Yes	Yes	Yes		Yes	Yes	Yes	Yes	Yes	Yes
Cost per year for total amount of licences	+	1							\$180 (Pro, per user)		Varies (Microsc
Estimated total cost savings/ROI it would provide total hours per year.	Varies	Varies `	Varies	Varies	Varies	Varies	Varies	Varies	Varies	Varies	Varies



Wrap-Up & Q&A

- Vet Deeply
- Align Policies Across Teams
- Eliminate Shadow Al
- Continuous Re-evaluation
- Key Takeaway: Be Proactive
- Encourage Questions
- Foster Discussions





Resources for Leaders

- National Institute of Standards and Technology
 - NIST AI Risk Management Framework
 - https://www.nist.gov/itl/ai-risk-management-framework
- Cybersecurity and Infrastructure Security Agency
 - CISA AI Security Guidelines
 - https://www.cisa.gov/ai
- Open Worldwide Application Security Project
 - OWASP Top 10
 - https://owasp.org/www-project-top-ten/
- SOC 2 Compliance for Al
 - https://www.compassitc.com/blog/achieving-soc-2-compliance-forartificial-intelligence-ai-platforms?utm_source=chatgpt.com
- Al Governance Alliance
 - https://initiatives.weforum.org/ai-governancealliance/home?utm_source=chatgpt.com
- ISO/IEC 23894
 - https://www.iso.org/standard/77304.html?utm_source=chatgpt.com



