

# AI のリスク軽減と AI による サイバーセキュリティの変革



# AIの脅威の状況

データポイズニング・  
トレーニング攻撃

インファレンス  
攻撃

回避攻撃

モデル 誠実さ・  
バックドア

抽出攻撃

コンプライアンス

システミック  
インフラ  
ストラクチャー

抽出・統合問題

# AI のリスク軽減 – 主要なリスクと軽減策

## コンプライアンスのリスク

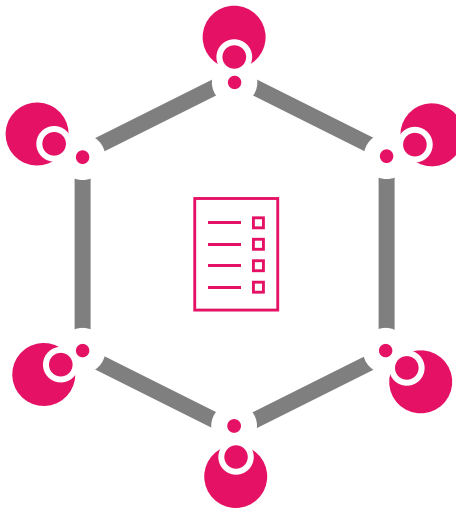
信頼できないソースから収集されたデータには、コンプライアンス違反につながる PII や機密データが含まれる可能性

### データの整合性

データセキュリティと環境制御が弱いため、AI ライフサイクル中にデータの整合性が損なわれる可能性

### モデルリスク

脆弱なモデル コード、敵対的な条件に耐えるために不適切にトレーニングされたモデル、トレーニングと幻覚、バイアスに対する過剰適合モデル



### AI/ML セキュリティ ポリシー

AI/ML のセキュリティ制御を管理する文書化されたポリシー (例 - モデルのトレーニング、モデルのデプロイメント、データ収集)

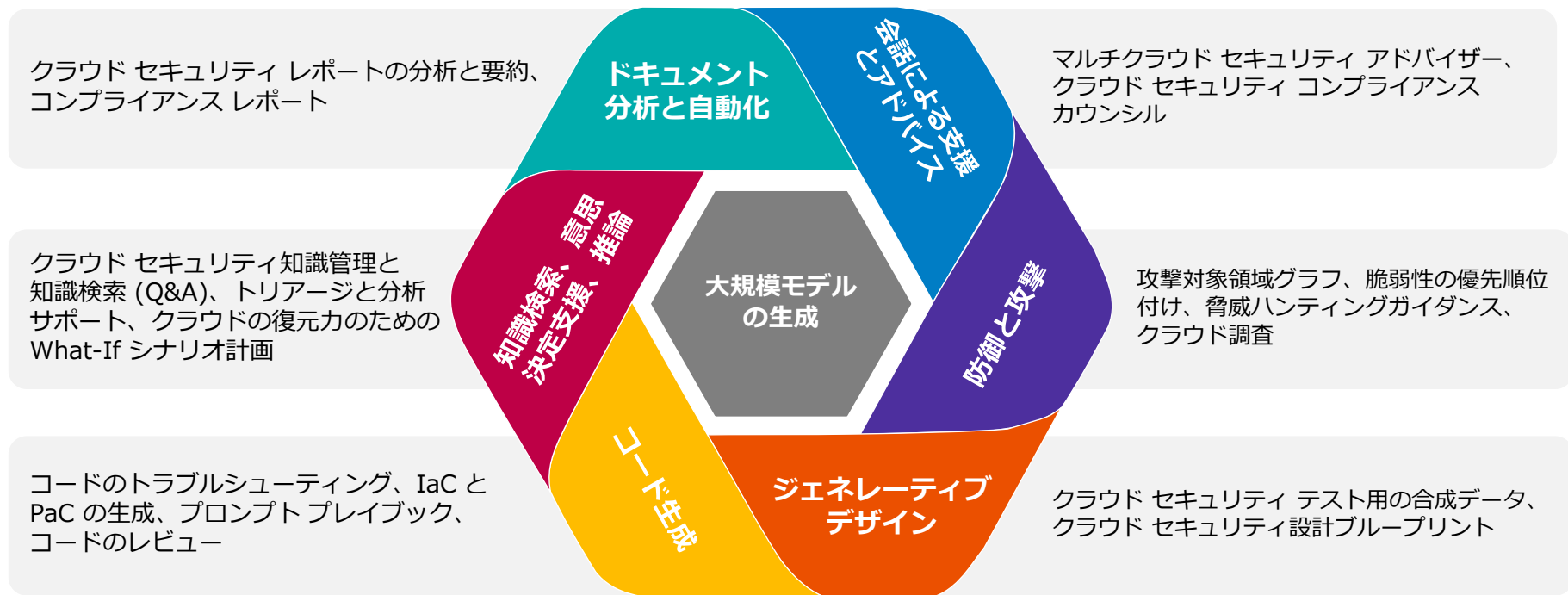
### 音声データ管理 練習

信頼できるソースからのデータの利用、データの分類とラベル付け、適切な保護とデータプライバシーの確保

## 信頼できる AI モデル

起こり得る状況 (最悪の場合、破損したデータセット) に基づいてトレーニングされたモデル、モデルコード用の安全な SDLC、モデルのバージョン管理

# ジェネレーティブ AI：機会と主なユースケース





Thank You