# Multi-Agent Collaboration in Incident Response with Large Language Models

**Zefang Liu** (Georgia Institute of Technology)
Contact: liuzefang@gatech.edu

## Overview

- **Incident response (IR)** is essential for cybersecurity, requiring **quick decision-making and coordination**.
- **Large Language Models (LLMs)** can serve as **intelligent agents** to enhance collaboration and efficiency.
- **LLM-based multi-agent collaboration** in **cybersecurity** is explored using **Backdoors & Breaches**, a tabletop game designed for IR training.
- Different **team structures (centralized, decentralized, hybrid)** are analyzed to evaluate their impact on IR effectiveness.

## Backdoors & Breaches

- A **tabletop game** that simulates **real-world cybersecurity incidents**.
- Players take on roles as **incident captain** and **defenders**, working together to **uncover and mitigate attack vectors**.
- The game includes different card types:
  - **Attack Cards:** Represent stages of a cyberattack (e.g., initial compromise, pivot and escalate, command and control (C2) and exfiltration, and persistence).
  - **Procedure Cards:** Defensive strategies used to detect and counter threats.
  - **Inject Cards:** Unexpected events that introduce new challenges.
- The goal is to **reveal all hidden attack cards** within limited turns through strategic decision-making.



**Initial Compromise**     **Pivot and Escalate**     **C2 and Exfil**
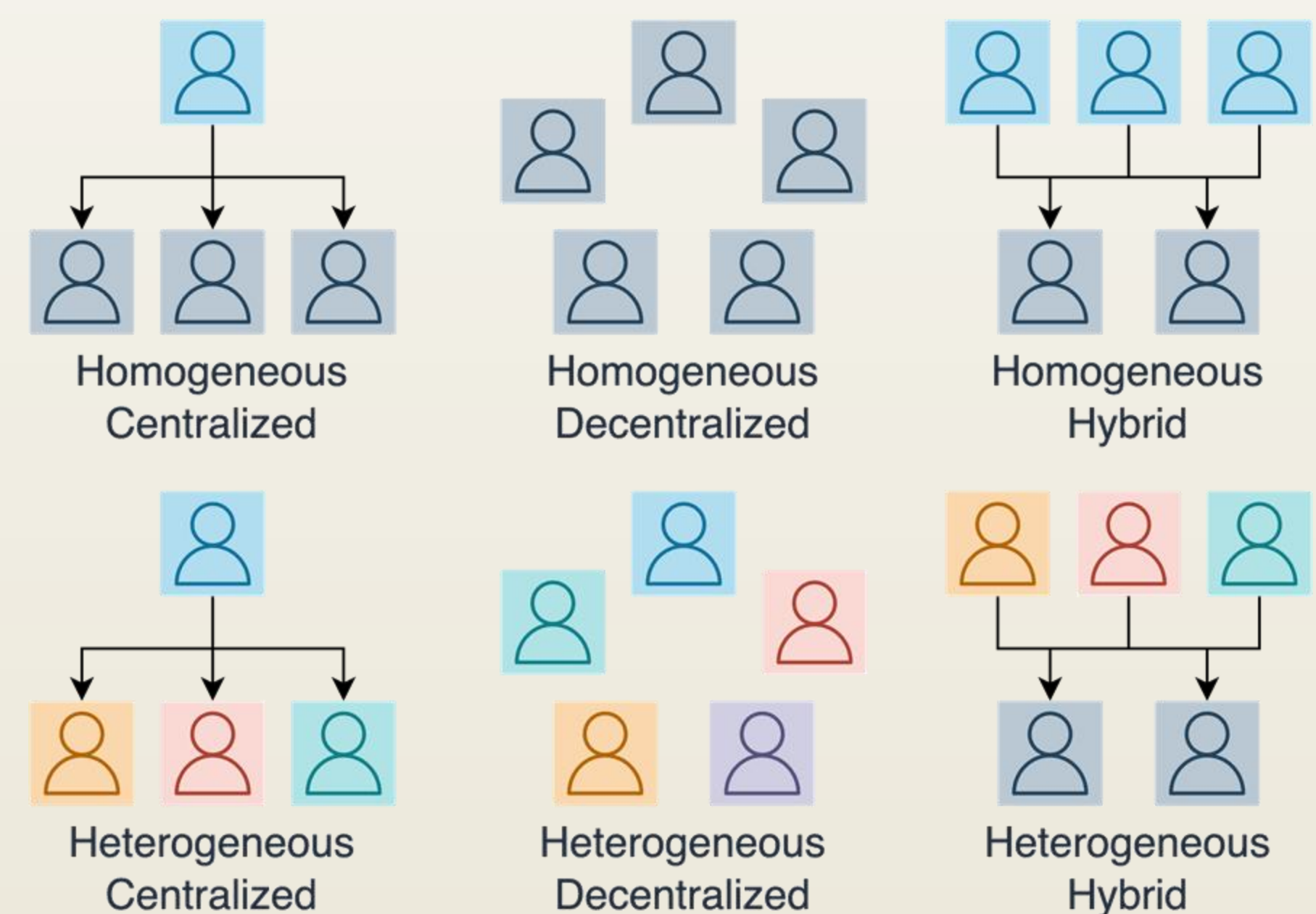


**Persistence**     **Procedure**     **Inject**

## Experimental Setup

- **LLM-based agents:** Implemented using **AutoGen**, with structured roles and interactions.
- **Team structures:**
  - **Centralized:** Leadership-driven decision-making.
  - **Decentralized:** Equal decision-making across all agents.
  - **Hybrid:** Mix of leadership and collaboration.
- **Evaluation metrics:**
  - **Success rate** in uncovering all attack stages.
  - **Failure patterns** across different team structures.

## Team Structures



Homogeneous Centralized — Homogeneous Decentralized — Homogeneous Hybrid

Heterogeneous Centralized — Heterogeneous Decentralized — Heterogeneous Hybrid

## Experimental Results

- **Centralized teams** perform well due to clear leadership but may struggle with adaptability.
- **Decentralized teams** leverage diverse expertise but can have coordination challenges.
- **Hybrid teams** balance structured leadership with flexibility, leading to strong performance.
- **LLM-based agents** facilitate IR processes by assisting in decision-making and coordination.

| Team | Success | Failure | Pentest | Invalid |
|------|---------|---------|---------|---------|
| Homo-Cen | 14 | 1 | 2 | 3 |
| Heter-Cen | 13 | 3 | 3 | 1 |
| Homo-Dec | 13 | 1 | 1 | 5 |
| Hetero-Dec | 12 | 3 | 3 | 2 |
| Homo-Hyb | 14 | 3 | 2 | 1 |
| Hetero-Hyb | 13 | 1 | 2 | 4 |

## Case Studies

- **Homogeneous Centralized:** Over-reliance on high-modifier procedures led to poor adaptability.
- **Heterogeneous Centralized:** Struggled with prioritization and aligning expert inputs.
- **Homogeneous Decentralized:** Slow decision-making and overuse of standard procedures.
- **Heterogeneous Decentralized:** Lack of coordination led to missed attack indicators.
- **Homogeneous Hybrid:** Misprioritized investigations delayed threat detection.
- **Heterogeneous Hybrid:** Expertise misalignment caused early attack stages to be overlooked.

## Conclusion & Future Work

- **LLMs** demonstrate strong potential in **multi-agent collaboration** in **incident response** and **cybersecurity**.
- **Future directions** includes:
  - Improving adaptability of LLMs for unpredictable cyber threats.
  - Extending simulations to real-world cybersecurity environments.
  - Exploring human-LLM hybrid teams for incident response.

## References

Young, Jacob, and Farshadkhah, Sahar. "Backdoors & Breaches: Using a Tabletop Exercise Game to Teach Cybersecurity Incident Response." *Proceedings of the EDSIG Conference ISSN*. Vol. 2473. 2021.
Wu, Qingyun, et al. "AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation." *ICLR 2024 Workshop on Large Language Model Agents*.
Liu, Zefang. "Multi-Agent Collaboration in Incident Response with Large Language Models." *AAAI 2025 Workshop on Multi-Agent AI in the Real World*.

SCAN ME