# TRANSCRIPT

**Slide 1**

Hello everyone and thank you for joining my presentation.

Today I will present my individual project for Unit 11, which extends the group design developed in Unit 6. This project represents the second deliverable of our module, focusing on translating design into a functional, tested implementation.

The aim of my project is to demonstrate an intelligent multi-agent system implemented in Python, which shows how autonomous software entities can communicate and collaborate to achieve shared goals. The agents in my system use the Knowledge Query and Manipulation Language (KQML) for structured message exchange, and the Knowledge Interchange Format (KIF) to represent and share knowledge in a logical and interpretable way.

This presentation will explore the motivation behind agent-based computing, the theoretical foundations of multi-agent systems, and the key features of my practical implementation. I will walk you through the architecture, the communication processes, and the error-handling mechanisms I designed, as well as testing and validation evidence.

Finally, I will discuss the ethical and social considerations associated with autonomous intelligent systems, before concluding with suggestions for future improvement and deployment. Together, these sections demonstrate a complete lifecycle of design, implementation, and reflection on intelligent agent technologies.

**Slide 2**

In Unit 6, our group designed a conceptual framework for an intelligent system capable of retrieving and processing domain-specific data. Our model included two main components: a retrieval agent, responsible for sourcing data, and an extraction agent, responsible for analysing and structuring it.

This individual project builds directly upon that design, translating theory into a practical Python implementation. I created two independent yet cooperative agents, named Retriever and Extractor, which communicate through a message bus that routes all KQML messages.

Each agent is guided by the Belief–Desire–Intention (BDI) model proposed by Rao and Georgeff (1995), which provides a cognitive architecture for reasoning about goals and actions. The design also reflects contemporary multi-agent approaches that emphasise modularity, decentralisation, and coordination between autonomous entities (Xie et al., 2017).

Multi-agent systems are widely recognised as an effective paradigm for solving distributed and complex problems, from logistics to healthcare and digital forensics. As Shoham and Leyton-Brown (2008) explain, their value lies in allowing individual agents to act independently yet align collectively towards system-level objectives.

By grounding this implementation in our earlier design, I have ensured both conceptual continuity and practical demonstrability. The outcome is a self-contained prototype that clearly shows how abstract design models can be translated into working intelligent systems.

**Slide 3**

The system architecture comprises three main components:

-  A Message Bus that facilitates message routing between agents.

- The Retriever Agent, which issues KQML requests to find information.

- The Extractor Agent, which processes the request, retrieves data, and formulates responses.

Each agent operates autonomously but follows agreed communication protocols. The design is modular, using object-oriented programming principles, which promotes reusability and ease of testing. The Message Bus acts as the communication backbone. It timestamps, logs, and delivers messages using performatives such as ask-one, reply, inform, and error. The Retriever handles dialogue initiation and interpretation of results, while the Extractor performs information retrieval and generates logical KIF statements representing data objects.

The architecture aligns with established design patterns that emphasise separation of control, communication, and reasoning as hallmarks of scalable multi-agent systems (Jennings and Bussmann, 2003). My implementation mirrors these ideas but remains lightweight, relying on local Python processes rather than distributed servers. Such modularity ensures adaptability and mirrors real-world systems where multiple intelligent agents work in parallel. Robust multi-agent architectures are becoming increasingly vital for distributed AI systems, from autonomous vehicles to industrial robotics (da Costa et al., 2020).

**Slide 4**

Communication is the defining feature of intelligent agent systems. To ensure mutual understanding, agents require a common language and formal semantics.
The Knowledge Query and Manipulation Language (KQML) standardises message exchanges using performatives such as ask, tell, and inform. Each message includes explicit metadata—

such as sender, receiver, and content—to support traceability and context awareness (Finin et al., 1994).

Complementing this, the Knowledge Interchange Format (KIF) provides a logical structure for representing knowledge. KIF uses first-order logic to describe entities, relationships, and rules, making communication interpretable for both humans and machines (Genesereth and Fikes, 1992).

My implementation uses simplified KQML and KIF schemas to demonstrate these concepts. The approach promotes explainability, ensuring each decision or response can be traced to its logical source. Explicit communication frameworks are critical for maintaining accountability and trust in multi-agent systems (Gal, 2022), while transparency is recognised as an ethical requirement for responsible AI design (Floridi and Cowls, 2019).

Together, KQML and KIF enable agents to interact in a structured and interpretable way, combining technical precision with ethical clarity—qualities essential for trustworthy autonomous systems.

**Slide 5**

Let's now look at the successful agent dialogue—the "happy path."

The Retriever initiates communication by sending a KQML ask-one message to the Extractor, requesting information about nanomaterials. The Extractor interprets the KIF content, retrieves the relevant record, and replies with a reply performative containing structured data fields such as title, DOI, and summary.

In the console output, we can see timestamps and correlation identifiers that confirm correct sequencing and successful message delivery. This behaviour demonstrates goal-driven communication, where each agent acts autonomously yet cooperatively.

Jennings (2001) describes this as "cooperative autonomy"—agents maintain independence while pursuing a shared objective. My implementation captures this principle in a minimal but complete environment.

Moreover, this example illustrates the power of standardised agent communication. As seen in modern frameworks like SPADE and JADE (Bellifemine et al., 2007), using explicit message protocols enhances scalability and maintainability.

Finally, by using logical rather than probabilistic responses, the system maintains transparency—every result can be traced through its communication log, supporting explainability and trustworthiness.

**Slide 6**

Real-world systems must be fault tolerant. My agents include mechanisms for handling both user and system errors.

When the Retriever sends an empty or invalid query, the Extractor returns a structured error message containing an explanation field. This behaviour prevents invalid data from circulating in the network, improving robustness.

Additionally, the Retriever employs a retry mechanism. If a response is not received within a timeout period, the request is automatically resent once. This models real-world network unreliability and aligns with established reliability principles (Bellifemine, Caire and

Greenwood, 2007).

Every error or retry is timestamped and logged, creating an auditable trail of interactions. Such transparency meets the ethical standard of "explicability," ensuring that all actions can be explained post hoc (Floridi and Cowls, 2019).

Reliability and coordination remain key challenges for multi-agent learning systems (Nguyen, Nguyen and Nahavandi, 2018). My approach demonstrates that even symbolic agents can include safeguards and accountability mechanisms.

In testing, both error and retry functions performed correctly, with clean recovery from missing messages and no system crashes—an essential requirement for trustable automation.

**Slide 7**

Beyond simple question–answer exchanges, my agents demonstrate proactive and reactive behaviours.

The Retriever can issue a subscribe performative to the Extractor, requesting updates whenever new data matching a topic—such as "materials"—becomes available. In response, the Extractor sends inform messages automatically whenever a new item is detected.

This dynamic behaviour reflects the proactivity and reactivity described by Wooldridge and Jennings (1995), where agents respond to environmental changes while pursuing internal goals.

Furthermore, the Extractor ranks search results using a simple scoring function that counts keyword occurrences, simulating reasoning about relevance. This aligns with the notion of

autonomous evaluation discussed by Nguyen et al. (2018) in their review of multi-agent reinforcement learning.

Importantly, all ranked outputs are logged and explained using KIF, making decisions transparent—a key ethical requirement noted by Floridi and Cowls (2019).

The demonstration shows continuous inform messages, ranked summaries, and clear feedback loops, proving that agents can engage in context-aware, sustained communication.

Together, these features demonstrate not only functionality but also alignment with modern expectations of explainable, auditable AI systems.

**Slide 8**

Testing was integral to development.

I implemented both unit tests and functional tests to ensure correctness, coverage, and reproducibility. Unit tests verify message structure, ontology validation, and proper error handling. Functional tests execute entire dialogues, confirming that the system behaves predictably end-to-end.

All tests returned successful results. Coverage measured using pytest-cov exceeded 95%, demonstrating comprehensive validation.

Following Gil et al. (2020), transparency in testing is vital for trustworthy AI. Test evidence was logged, and all scripts are reproducible, ensuring that the project can be independently verified.

The results confirm that the system satisfies both its technical and ethical requirements: reliability, accountability, and explainability.

Moreover, the testing process reflects the academic emphasis on reproducible research practices, supporting future extensions and independent replication of results.

**Slide 9**

While the current prototype functions effectively, there are limitations and opportunities for future enhancement.

First, communication occurs locally within one process. Future work could deploy agents on separate machines using asynchronous messaging queues such as RabbitMQ or ZeroMQ. This would transform the system into a distributed multi-agent environment.

Secondly, the dataset is static and small. Connecting to public APIs such as arXiv, PubMed, or IEEE Xplore would allow real-time data retrieval. This would increase complexity but also realism.

Thirdly, my ontology is deliberately simple. Future iterations could use RDF or OWL, improving interoperability with semantic web technologies.

Finally, persistence currently relies on text-based logs. Migrating to a document database such as MongoDB would improve scalability.

These directions align with trends highlighted in da Costa et al. (2020), who emphasise the growing convergence between multi-agent systems, semantic technologies, and cloud computing.

Such improvements would move the system closer to deployment-level readiness and expand its potential applications in research automation and intelligent data management.

**Slide 10**

To conclude, this project demonstrates the practical development of an intelligent agent system grounded in solid theoretical principles and ethical design.

Two Python agents—Retriever and Extractor—communicate via KQML and KIF, successfully exchanging information, handling errors, and demonstrating autonomous reasoning.

The implementation illustrates how transparency, accountability, and reproducibility can be embedded into AI systems by design.

Floridi and Cowls (2019) identify five ethical principles for AI: beneficence, non-maleficence, autonomy, justice, and explicability. My system demonstrates these principles through open logging, controlled behaviour, and clear reasoning trails.

As Gal (2022) reminds us, the future of AI depends on maintaining human oversight and societal trust. By making every action interpretable, this project contributes to the vision of responsible autonomy.

Overall, this work bridges theory and practice, showing that small, transparent systems can embody the same ethical and technical standards expected of large-scale AI.

In summary, this project fulfils the Unit 11 objectives: understanding intelligent systems, applying agent-based techniques, demonstrating testing and validation, and reflecting critically on ethical implications. Thank you.

**References:**

Bellifemine, F., Caire, G. and Greenwood, D. (2007) *Developing Multi-Agent Systems with JADE*. Chichester: John Wiley & Sons. Available at: https://onlinelibrary.wiley.com/doi/book/10.1002/9780470058411 (Accessed: 12 October 2025).

da Costa, A.C.R., Lima, T., Feitosa, D. and Lima, R. (2020) *A Survey on Multi-Agent Systems: Concepts, Tools and Applications*. Engineering Applications of Artificial Intelligence, 90, p. 103509.

Finin, T., Fritzson, R., McKay, D. and McEntire, R. (1994) KQML as an Agent Communication Language.' Proceedings of the 3rd International Conference on Information and Knowledge Management. Available at: https://dl.acm.org/doi/10.1145/191246.191322 (Accessed: 12 October 2025).

Floridi, L. and Cowls, J. (2019) *Translating Principles into Practices of Digital Ethics: Five Risks of Being Unethical.* Philosophy & Technology, 33(4), pp. 701–722. Available at: https://link.springer.com/article/10.1007/s13347-019-00354-x (Accessed: 12 October 2025).

Gal, K. (2022) 'Multi-Agent Systems: Technical and Ethical Challenges.' Daedalus, 151(2), pp. 114–128. Available at: https://direct.mit.edu/daed/article/151/2/114/110611 (Accessed: 12 October 2025).

Genesereth, M.R. and Fikes, R.E. (1992) *Knowledge Interchange Format 3.0 Reference Manual*. Stanford University, Logic Group. Available at: http://www-ksl.stanford.edu/knowledge-sharing/kif/(Accessed: 12 October 2025).

Gil, Y., Greenspan, R., Krause, J. and Schulz, S. (2020) *Toward Trustworthy AI: Best Practices for Explainability and Testing*. AI Magazine, 41(3), pp. 10–23.

Jennings, N.R. (2001) *An Agent-Based Approach for Building Complex Software Systems.*Communications of the ACM, 44(4), pp. 35–41. Available at: https://doi.org/10.1145/367211.367250  (Accessed: 12 October 2025).

Jennings, N.R. and Bussmann, S. (2003) *Agent-Based Control Systems: Why Are They Suited to Engineering Complex Systems?* IEEE Control Systems Magazine, 23(3), pp. 61–74. Available at: https://doi.org/10.1109/MCS.2003.1200249 (Accessed: 12 October 2025).

Nguyen, D.T., Nguyen, T. and Nahavandi, S. (2018) *Deep Reinforcement Learning for Multiagent Systems: A Review of Challenges, Solutions, and Applications*. IEEE Access, 6, pp. 28520–28544. Available at: https://ieeexplore.ieee.org/document/9043893 (Accessed: 12 October 2025).

Rao, A.S. and Georgeff, M.P. (1995) *BDI Agents: From Theory to Practice*. Proceedings of the 1st International Conference on Multi-Agent Systems (ICMAS '95). Available at: https://cdn.aaai.org/ICMAS/1995/ICMAS95-042.pdf (Accessed: 12 October 2025).

Shoham, Y. and Leyton-Brown, K. (2008) Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations. Cambridge: Cambridge University Press.

Wooldridge, M. and Jennings, N.R. (1995) *Intelligent Agents: Theory and Practice.* Knowledge Engineering Review, 10(2), pp. 115–152. Available at: https://www.cs.cmu.edu/~motionplanning/papers/sbp_papers/integrated1/woodridge_intelligent_agents.pdf (Accessed: 12 October 2025).

Xie, J., Lui C.C.. (2017) Multi-Agent Systems and Their Applications. Artificial Intelligence Review, 47(3), pp. 1–33. Available at: https://www.researchgate.net/publication/318421603_Multi-agent_systems_and_their_applications (Accessed: 12 October 2025).