

Report and Recommendations
'Alice'
Ethical Considerations for AI in Education
22055453

Contents

1	Overview	3
2	Relevant Professional Codes of Conduct	3
2.1	BCS Code of Conduct	3
2.2	IEEE Code of Ethics	3
2.3	ACM Code of Ethics	4
2.4	Caveats of a Deontological Outlook	5
3	Legal and Professional Standards	5
3.1	Data Protection Legislation	5
3.2	Education Sector Regulations	5
3.3	Professional Standards and Guidelines	5
3.4	Industry-Specific Standards	6
3.5	Ethical AI Frameworks	6
3.6	Continuous Compliance and Auditing	6
4	Further Ethical Considerations	7
4.1	Data Protection and Privacy	7
4.2	Transparency and Explainability	7
4.3	Bias and Fairness	8
4.4	Psychological Impact	8
4.5	Ethical Framework and Governance	8
5	User Experience and Accessibility	9
5.1	User-Centred Design	9
5.2	Accessibility Standards	9
5.3	Inclusive Design	9
5.4	Conversational Design	10
5.5	Mobile Responsiveness	10
5.6	User Feedback and Iteration	10
5.7	Ethical Considerations in UX	10
6	Monitoring, Evaluation, and Continuous Improvement	11
6.1	AI and Machine Learning Architecture	11
6.2	Data Management and Security	11
6.3	Performance Metrics	11
6.4	Continuous Learning and Adaptation	12
6.5	User Feedback Integration	12
6.6	Impact Assessment	12
6.7	Security and Privacy Audits	13
6.8	Continuous Professional Development	13
7	Key Recommendations	13

A	Appendix A: Data Protection Impact Assessment Template	14
B	Appendix B: Ethical AI Checklist	14
C	Appendix C: Sample Conversational Flows	14

1 Overview

South Star Academy's proposed student support services chatbot 'Alice' shows great promise in terms of market return, but scrupulous rigour and a diligent attention to detail will be necessary before such a system can be deemed production-ready. Ethical considerations need to be made concerning data privacy, accessibility, and user experience, as well as strict adherence to legal and professional standards, security best practices taken into account, and an overarching monitoring and management implementation strategy. This report comprises a comprehensive analysis and its resulting recommendations for deployment, such that TechSoft be fully cognisant of the key contemporary issues faced by an AI chatbot in education, and by extension the actions required for ethical practice, full legal compliance, and more.

2 Relevant Professional Codes of Conduct

The following elements from the BCS, IEEE, and ACM codes of conduct provide us with a rough guideline to start.

2.1 BCS Code of Conduct

The BCS Code of Conduct (*British Computer Society 2024, p. 1*) outlines four primary areas of professional responsibility, each of which has direct implications for the 'Alice' project:

Public Interest

BCS members must have "due regard for public health, privacy, security and wellbeing of others and the environment" (*British Computer Society 2024, p. 2*). This principle is particularly salient for 'Alice', as the chatbot will be handling sensitive student data and potentially influencing student wellbeing. We must ensure that:

- Alice's decision-making processes prioritise student safety and wellbeing above all else.
- Robust privacy measures are implemented to protect student data from unauthorised access or breaches.
- The chatbot's impact on the educational environment is constantly monitored and evaluated.

Professional Competence and Integrity

The BCS mandates that members should "only undertake to do work or provide a service that

is within your professional competence" (*British Computer Society 2024, p. 2*). This raises important questions about the competence boundaries of AI systems:

- What will we define as outside the limit of Alice's 'competence' in student support?
- What mechanisms will be in place to ensure Alice does not overstep these boundaries?
- How will we ensure ongoing development of Alice's capabilities while maintaining integrity?

Duty to the Profession

The BCS Code requires members to "accept your personal duty to uphold the reputation of the profession and not take any action which could bring the profession into disrepute" (*British Computer Society 2024, p. 3*). For the 'Alice' project, this means:

- Ensuring transparency about Alice's AI nature and capabilities to all users.
- Implementing robust safeguards against potential misuse or malfunction that could damage public trust in AI systems in education.
- Actively contributing to the development of best practices for AI in educational settings.

2.2 IEEE Code of Ethics

The IEEE Code of Ethics (*Institute of Electrical and Electronics Engineers 2024, p. 1*) provides additional ethical considerations that are particularly relevant to the technological aspects of 'Alice':

Respect for Privacy

IEEE members commit to be "respectful of the privacy of others and the protection of their personal information and data" (*Institute of Electrical and Electronics Engineers 2024, p. 1*). For 'Alice', this translates to:

- Implementing state-of-the-art data protection measures.
- Designing Alice with privacy-by-default principles.
- Establishing clear data retention and deletion policies.

Fairness and Non-Discrimination

The IEEE Code states that members "will not discriminate against any person because of characteristics protected by law" (*Institute of Electrical and Electronics Engineers 2024, p. 1*). This principle is crucial for Alice's design:

- Rigorous testing must be conducted to identify and eliminate potential biases in Alice's algorithms.
- The chatbot's language and responses must be inclusive and respectful of all student demographics.
- Regular audits should be performed to ensure equitable treatment of all users.

Honesty and Trustworthiness

IEEE members commit to "be honest and realistic in stating claims or estimates based on available data" (*Institute of Electrical and Electronics Engineers 2024, p. 2*). For the 'Alice' project, this means:

- Clear communication of Alice's capabilities and limitations to all stakeholders.
- Implementing explainable AI techniques to provide transparency in decision-making processes.
- Establishing mechanisms for human oversight and intervention when Alice's confidence in a response is low.

2.3 ACM Code of Ethics

The ACM Code of Ethics and Professional Conduct (*Association for Computing Machinery 2023, p. 1*) provides additional ethical principles that are particularly relevant to the societal impact of 'Alice':

Contribute to Society and Human Well-being

The ACM Code states that computing professionals should "contribute to society and to human well-being, acknowledging that all people are stakeholders in computing" (*Association for Computing Machinery 2023, p. 1*). For 'Alice', this principle translates to:

- Ensuring that the chatbot's primary goal is to enhance student well-being and academic success.
- Regularly assessing the broader impact of Alice on the educational ecosystem.

- Designing features that promote digital literacy and responsible AI interaction among students.

Avoid Harm

The ACM emphasises that computing professionals should "avoid harm to others" (*Association for Computing Machinery 2023, p. 1*). In the context of 'Alice', this principle demands:

- Implementing robust safeguards against potential psychological harm from AI interactions.
- Establishing clear escalation protocols for situations where Alice detects potential self-harm or other serious issues.
- Regular risk assessments to identify and mitigate potential negative consequences of the chatbot's use.

Ensure Fair Participation

The ACM Code requires computing professionals to "foster fair participation of all people, including those of underrepresented groups" (*Association for Computing Machinery 2023, p. 2*). For the 'Alice' project, this means:

- Designing Alice's interface and interactions to be accessible to students with disabilities.
- Ensuring that the chatbot's language models and knowledge base represent diverse perspectives and experiences.
- Implementing features that promote equal access to educational resources and opportunities through Alice's recommendations.

By adhering to these specific elements of the BCS, IEEE, and ACM codes of conduct, we can ensure that the development and implementation of 'Alice' at South Star Academy not only meets professional standards but also sets a new benchmark for ethical AI deployment in educational settings. These guidelines will inform our technical decisions, shape our data governance policies, and guide our interactions with stakeholders throughout the project lifecycle.

Duty to Relevant Authority

BCS members must "carry out professional responsibilities with due care and diligence in accordance with the Relevant Authority's requirements whilst exercising your professional judgement at all times" (*British Computer Society 2024, p. 2*). In our context:

- We must balance South Star Academy's requirements with ethical considerations and legal obligations.
- Clear protocols must be established for situations where Alice's judgments might conflict with school policies.
- Regular audits should be conducted to ensure Alice's actions align with both the school's requirements and broader ethical standards.

2.4 Caveats of a Deontological Outlook

From an objective and cynical standpoint, however, "such codes are seldom consulted and often incorporate bland (and sometimes contradictory) statements intended to satisfy a broad range of stakeholders" (*Blundell 2020, p. 40*) and professionals must practise discernment and maintain awareness for the consequences of poor decision-making, both quantitatively and in qualitative aspects. Underpinning all ethical consideration one might argue that it is a professional's role to go beyond the call of duty where suitable, not straying outside of Kantian philosophy in fact, for he was a stark proponent of logical reasoning, evident in his mantra, "respect the reason in you", where universally established rules fail.

3 Legal and Professional Standards

3.1 Data Protection Legislation

It is imperative that the chatbot is compliant with all data protection laws. These include:

General Data Protection Regulation (GDPR)

Adhering to GDPR requirements (*European Union 2016*) involves:

- Establishing a lawful basis for processing: Consent or legitimate interests
- Implementing data subject rights: Access, rectification, erasure, portability
- Appointing a Data Protection Officer (DPO)

UK Data Protection Act 2018

Compliance with the UK Data Protection Act 2018 (*UK Government 2018*) includes:

- Adhering to specific provisions for processing personal data in educational contexts

- Implementing safeguards for processing special category data (e.g., health information)

Children's Online Privacy Protection Act (COPPA)

Considering COPPA requirements (*Federal Trade Commission 2023*) involves:

- Obtaining parental consent for students under 13
- Implementing limited data collection and retention policies

3.2 Education Sector Regulations

Adherence to education-specific regulations is essential:

Education and Skills Act 2008

Complying with the Education and Skills Act 2008 (*UK Government 2008*) requires:

- Fulfilling the duty to promote the well-being of students
- Implementing safeguarding responsibilities in digital environments

Keeping Children Safe in Education

Following the Keeping Children Safe in Education guidance (*Department for Education 2024a*) involves:

- Implementing online safety measures for educational technology
- Providing staff training on digital safeguarding

Special Educational Needs and Disability (SEND) Code of Practice

Adhering to the SEND Code of Practice (*Department for Education 2024b*) includes:

- Ensuring accessibility requirements for digital learning tools are met
- Considering personalised support for students with SEND

3.3 Professional Standards and Guidelines

Aligning with professional standards ensures ethical development and deployment:

BCS Code of Conduct

Following the BCS Code of Conduct (*British Computer Society 2024*, pp. 1-5) involves:

- Considering public interest in development decisions
- Maintaining professional competence and integrity
- Fulfilling duty to relevant authorities

ACM Code of Ethics and Professional Conduct

Adhering to the ACM Code of Ethics (*Association for Computing Machinery 2023*, pp. 1-4) requires:

- Contributing to society and human well-being
- Avoiding harm in system design and implementation
- Maintaining honesty and trustworthiness

IEEE Ethically Aligned Design

Following IEEE Ethically Aligned Design principles (*Institute of Electrical and Electronics Engineers 2024*, pp. 2-5) includes:

- Preserving human rights in AI systems
- Ensuring transparency and accountability in AI decision-making
- Implementing privacy-by-design principles

3.4 Industry-Specific Standards

Implementing relevant technical and educational standards:

ISO/IEC 27001:2022

Adhering to ISO/IEC 27001:2022 (*International Organization for Standardization 2022*) involves:

- Conducting risk assessment and management
- Implementing information security controls
- Establishing continuous improvement processes

Learning Tools Interoperability (LTI) Standards

Implementing LTI standards (*IMS Global Learning Consortium 2024*) includes:

- Ensuring secure integration with existing learning management systems
- Enabling data portability and interoperability

Web Content Accessibility Guidelines (WCAG) 2.2

Complying with WCAG 2.2 (*World Wide Web Consortium 2023*) requires:

- Ensuring the chatbot interface is perceivable, operable, understandable, and robust
- Maintaining compatibility with assistive technologies

3.5 Ethical AI Frameworks

Adopting recognised ethical AI frameworks to guide development:

UNESCO Recommendation on the Ethics of Artificial Intelligence

Following UNESCO recommendations (*UNESCO 2021*) involves:

- Protecting human rights and fundamental freedoms
- Promoting diversity and inclusiveness in AI systems
- Ensuring transparency and explainability of AI decisions

OECD AI Principles

Adhering to OECD AI Principles (*Organisation for Economic Co-operation and Development 2023*) includes:

- Ensuring AI benefits people and the planet
- Designing AI systems that respect the rule of law, human rights, democratic values, and diversity

EU Ethics Guidelines for Trustworthy AI

Following EU Ethics Guidelines (*European Commission 2024*) requires:

- Implementing human agency and oversight in AI systems
- Ensuring technical robustness and safety
- Maintaining privacy and data governance

3.6 Continuous Compliance and Auditing

Establishing processes for ongoing compliance:

Regular Compliance Audits

Conducting regular compliance audits (*Information Commissioner's Office 2024*) involves:

- Performing annual data protection audits
- Engaging third-party security assessments

Ethics Review Board Oversight

Implementing ethics review board oversight (*AI Ethics Board Best Practices 2024*) includes:

- Conducting periodic reviews of chatbot decisions and outcomes
- Establishing stakeholder feedback mechanisms

Continuous Professional Development

Ensuring continuous professional development (*Chartered Institute of Personnel and Development 2024*) requires:

- Providing regular training on evolving legal and ethical standards
- Obtaining certification in AI ethics for key personnel

4 Further Ethical Considerations

4.1 Data Protection and Privacy

Any implementation of 'Alice' would raise significant privacy concerns regarding the collection, storage, and use of student data (*Annus 2023, pp. 366-370*). Key considerations include:

Compliance with Data Protection Regulations

It is essential to ensure compliance with the General Data Protection Regulation (GDPR) and other relevant data protection regulations (*Information Commissioner's Office 2024*). This includes:

- Establishing a lawful basis for processing personal data
- Implementing data subject rights (access, rectification, erasure)
- Conducting Data Protection Impact Assessments (DPIAs)

Informed Consent

Obtaining informed consent from students for data collection and processing is crucial (*European Data Protection Board 2023*). This involves:

- Providing clear and transparent information about data usage
- Implementing opt-in mechanisms for non-essential features
- Ensuring age-appropriate consent for students under 18

Data Minimisation

Adhering to data minimisation principles (*Article 29 Working Party 2018*) is important:

- Collecting only necessary information for chatbot functionality
- Conducting regular data audits to ensure relevance of stored information

Secure Storage and Transmission

Implementing robust security measures for sensitive student data (*National Cyber Security Centre 2024*) is critical:

- Utilising end-to-end encryption for data in transit
- Implementing strong access controls and authentication mechanisms
- Conducting regular security audits and penetration testing

4.2 Transparency and Explainability

Ensuring transparency in the chatbot's functionality is crucial for ethical implementation:

AI Disclosure

Clear disclosure of AI usage to students (*Institute of Electrical and Electronics Engineers 2023*) should include:

- Explicit labelling of 'Alice' as an AI system
- Information on the capabilities and limitations of the chatbot

Explainable AI Techniques

Implementing explainable AI techniques to interpret chatbot decisions (*Arrieta et al. 2022, pp. 82-115*) is recommended:

- Use of interpretable machine learning models
- Providing rationale for chatbot recommendations and actions

Regular Audits

Conducting regular audits of the chatbot's decision-making processes (*AI Ethics Guidelines for Education 2024*) is essential:

- Logging and analysis of chatbot interactions
- Third-party audits to ensure adherence to ethical guidelines

4.3 Bias and Fairness

Mitigating bias in the chatbot's algorithms is essential for equitable student support:

Diverse Training Data

Using diverse training data to prevent demographic biases (*Mehrabi et al. 2023, pp. 1-35*) involves:

- Including diverse student profiles in training datasets
- Regularly updating data to reflect changing student demographics

Bias Testing

Regular testing for biases in chatbot responses (*Association for Computing Machinery Conference on Fairness, Accountability, and Transparency 2024*) should include:

- Employing automated bias detection tools
- Implementing human-in-the-loop evaluation for sensitive topics

Fairness-Aware Machine Learning

Implementing fairness-aware machine learning techniques (*Barocas, Hardt, and Narayanan 2021*) is crucial:

- Utilising algorithmic fairness metrics (e.g., demographic parity, equal opportunity)
- Applying bias mitigation strategies in model training and deployment

4.4 Psychological Impact

The potential psychological effects of AI-based support on students must be carefully considered:

Avoiding Over-Reliance

Mitigating the risk of over-reliance on AI for emotional support (*Miner et al. 2022, p. 746*) involves:

- Clearly communicating AI's role as a supplement, not replacement, for human support

- Integrating the chatbot with pre-existing human counselling services and enabling direct connections to appropriate resources and non-AI support

Safeguards Against Harmful Responses

Implementing safeguards against harmful or inappropriate chatbot responses (*Bickmore et al. 2021, p. e11510*) includes:

- Developing content filtering and trigger warning systems
- Establishing escalation protocols for crisis situations

Clear Boundaries

Setting clear boundaries between AI support and human intervention (*American Psychological Association 2024*) requires:

- Defining thresholds for transitioning from AI to human support
- Training staff on effectively working alongside AI systems

4.5 Ethical Framework and Governance

Establishing a robust ethical framework is crucial for the responsible development and deployment of 'Alice':

Ethics Review Board

Forming an ethics review board (*United Nations Educational, Scientific and Cultural Organization 2023*) should involve:

- Ensuring multi-stakeholder representation (educators, students, ethicists, technologists)
- Conducting regular reviews of chatbot performance and ethical implications

AI Ethics Principles

Adhering to established AI ethics principles (*European Commission 2024*) is essential:

- Human agency and oversight
- Technical robustness and safety
- Privacy and data governance
- Transparency
- Diversity, non-discrimination, and fairness
- Societal and environmental well-being
- Accountability

Ethical Training

Providing continuous ethical training for the development team (*IEEE Ethics Certification Program for Autonomous and Intelligent Systems 2024*) should include:

- Regular workshops on AI ethics in education
- Incorporating ethical considerations into development processes

5 User Experience and Accessibility

5.1 User-Centred Design

Implementing a design process focussed on student needs:

User Research Methodologies

Conducting user research (*Goodman, Kuniavsky, and Moed 2023, pp. 50-100*) involves:

- Utilising contextual inquiry to understand student support scenarios
- Developing personas representing diverse student populations

Iterative Design Process

Implementing an iterative design process (*Holtzblatt and Beyer 2024, pp. 30-60*) includes:

- Using rapid prototyping with tools like Figma or Sketch
- Conducting usability testing with representative student groups

Emotional Design Principles

Applying emotional design principles (*Norman 2023, pp. 10-40*) involves:

- Designing for positive emotional responses
- Incorporating empathy in chatbot interactions

5.2 Accessibility Standards

Ensuring the chatbot is usable by all students:

Web Content Accessibility Guidelines (WCAG) 2.2 Compliance

Adhering to WCAG 2.2 (*World Wide Web Consortium 2023*) includes:

- Ensuring content is perceivable, operable, understandable, and robust
- Implementing keyboard accessibility and enough time for user interactions

Screen Reader Compatibility

Ensuring screen reader compatibility (*WebAIM 2024*) involves:

- Using ARIA landmarks and roles for improved navigation
- Providing descriptive alt text for images and icons

Cognitive Accessibility Considerations

Addressing cognitive accessibility (*Yesilada, Brajnik, and Harper 2023, pp. 1-10*) includes:

- Using clear and simple language
- Implementing consistent layout and interaction patterns

5.3 Inclusive Design

Addressing diverse student needs:

Multilingual Support

Implementing multilingual support (*Anastasiou and Schäler 2023, pp. 50-100*) involves:

- Integrating machine translation services
- Ensuring culturally appropriate responses and idioms

Adaptable User Interface

Creating an adaptable user interface (*Harper and Yesilada 2024, pp. 20-50*) includes:

- Offering customisable font sizes and colour contrasts
- Supporting different input methods (text, voice, gestures)

Neurodiversity Considerations

Addressing neurodiversity (*Armstrong 2023, pp. 30-60*) involves:

- Providing options to reduce visual clutter
- Offering alternative formats for information presentation (text, audio, visual)

5.4 Conversational Design

Creating natural and effective chatbot interactions:

Dialogue Flow Design

Designing dialogue flows (*Moore and Arar 2023, pp. 40-80*) includes:

- Creating conversation trees with appropriate branching
- Implementing fallback mechanisms for misunderstood queries

Tone and Personality

Establishing tone and personality (*Bradbury 2024, pp. 20-50*) involves:

- Maintaining a consistent voice aligned with educational context
- Using age-appropriate language and responses

Error Handling and Recovery

Implementing error handling and recovery (*Lemon and Pietquin 2023, pp. 100-150*) includes:

- Providing graceful error messages
- Offering suggestions for rephrasing or alternative actions

5.5 Mobile Responsiveness

Optimising for various devices and screen sizes:

Responsive Web Design Principles

Applying responsive web design principles (*Frain 2023, pp. 30-60*) involves:

- Using fluid grids and flexible images
- Implementing CSS media queries for device-specific layouts

Progressive Enhancement

Implementing progressive enhancement (*Champeon and Finck 2024, pp. 20-50*) includes:

- Ensuring core functionality is available to all devices
- Adding enhanced features for more capable browsers

Touch-Friendly Interfaces

Designing touch-friendly interfaces (*Hoover and Berkman 2023, pp. 80-120*) involves:

- Using appropriately sized touch targets
- Implementing gesture-based interactions where appropriate

5.6 User Feedback and Iteration

Continuously improving based on user input:

In-App Feedback Mechanisms

Implementing in-app feedback mechanisms (*Tullis and Albert 2024, pp. 100-150*) involves:

- Providing short surveys after chatbot interactions
- Offering easy-to-use bug reporting tools

Usage Analytics

Implementing usage analytics (*Beasley 2023, pp. 50-100*) includes:

- Tracking common queries and pain points
- Analysing conversation flows and completion rates

A/B Testing

Conducting A/B testing (*Kohavi, Tang, and Xu 2023, pp. 20-50*) involves:

- Comparative testing of different chatbot responses
- Implementing gradual rollout of new features

5.7 Ethical Considerations in UX

Balancing usability with ethical concerns:

Dark Pattern Avoidance

Avoiding dark patterns (*Brignull 2023*) includes:

- Providing transparent information about chatbot capabilities
- Offering clear opt-out options for data collection

Attention Economy Awareness

Addressing attention economy concerns (*Williams 2024, pp. 10-30*) involves:

- Designing for focussed, purposeful interactions
- Avoiding addictive design patterns

Privacy-Preserving UX Patterns

Implementing privacy-preserving UX patterns (Hartzog 2023, pp. 50-100) includes:

- Making privacy settings easily accessible and understandable
- Ensuring data usage transparency in the user interface

Data Anonymisation Techniques

Employing data anonymisation techniques (El Emam and Arbuckle 2023, pp. 75-100) includes:

- Applying K-anonymity for protecting student identities
- Implementing differential privacy for aggregate data analysis

6 Monitoring, Evaluation, and Continuous Improvement

6.1 AI and Machine Learning Architecture

Selecting appropriate AI technologies for 'Alice':

Natural Language Processing (NLP) Frameworks

Implementing NLP frameworks (Jurafsky and Martin 2024, pp. 1-15) involves:

- Utilising BERT or GPT-based models for understanding context and intent, allowing for the provision of specified and bespoke services.
- Custom training on domain-specific data for career advice, mental health support, and academic guidance.

Machine Learning Algorithms

Applying machine learning algorithms (Géron 2024, pp. 25-50) includes:

- Rely upon supervised learning and human-in-the-loop models for classification tasks (e.g., identifying at-risk students). "Human-centric AI" is a key principle.
- Implementing reinforcement learning for adaptive responses

6.2 Data Management and Security

Implementing robust security protocols and efficient data handling:

Data Encryption

Applying robust data encryption methods (Stallings 2023, pp. 100-150) involves:

- Utilising AES-256 for data at rest
- Implementing TLS 1.3 for data in transit

Authentication and Authorisation

Implementing secure authentication and authorisation (Josuttis 2023, pp. 80-120) includes:

- Using OAuth 2.0 and OpenID Connect for Single Sign-On (SSO) and secure authentication
- Implementing Role-Based Access Control (RBAC) for granular permissions

Intrusion Detection and Prevention

Deploying intrusion detection and prevention systems (Stallings and Brown 2024, pp. 200-250) involves:

- Implementing network-based IDS/IPS systems
- Utilising host-based security with endpoint detection and response (EDR)

Vulnerability Management

Managing vulnerabilities (Stuttard and Pinto 2023, pp. 150-200) includes:

- Conducting regular automated vulnerability scans
- Implementing penetration testing of IT systems

6.3 Performance Metrics

Establishing key performance indicators (KPIs) for 'Alice':

Conversational Metrics

Measuring conversational metrics (Quarteroni et al. 2024, pp. 1-32) includes:

- Assessing response accuracy and relevance
- Tracking task completion rates

Telemetry Systems

Setting up telemetry systems (*Vadapalli 2023, pp. 30-60*) includes:

- Implementing real-time data collection on user interactions
- Ensuring privacy-preserving logging mechanisms

Natural Language Understanding (NLU) Analysis

Conducting NLU analysis (*Jurafsky and Martin 2024, pp. 200-250*) involves:

- Assessing intent classification accuracy
- Evaluating entity recognition performance

Sentiment Analysis

Performing sentiment analysis (*Liu 2023, pp. 50-100*) includes:

- Analysing emotional tone of student interactions
(N.B. It is worth noting that certain interpretations of EU legislation suggest that AI systems capable of deciphering emotions may be inherently considered unethical. However, given the rapid pace of technological innovation, particularly in the AI sector, it is plausible that the European Parliament may adopt a nuanced approach, interpreting the law’s intent rather than adhering strictly to its literal wording) (*Dignum 2023, pp. 150-175*)
- Identifying potentially distressed students

6.4 Continuous Learning and Adaptation

Enhancing chatbot performance over time:

Dialogue Optimisation

Optimising dialogues (*Gao, Galley, and Li 2023, pp. 50-100*) involves:

- Refining conversation flows based on user feedback
- Dynamically adjusting response strategies

Human-in-the-Loop Evaluation

Conducting human-in-the-loop evaluation (*Vaughan 2024, pp. 30-60*) involves:

- Regular audits of chatbot conversations by education experts
- Crowdsourced evaluations for diverse perspectives

Ethical Review Process

Implementing an ethical review process (*Floridi and Cows 2023*) includes:

- Periodic assessments of chatbot decisions for bias
- Alignment checks with established ethical guidelines

6.5 User Feedback Integration

Incorporating student and staff input:

Feedback Collection Methods

Implementing feedback collection methods (*Tullis and Albert 2024, pp. 100-150*) involves:

- Providing in-chat feedback options
- Conducting periodic user surveys

Participatory Design Workshops

Organising participatory design workshops (*Simonsen and Robertson 2023, pp. 50-100*) includes:

- Conducting co-creation sessions with students and educators
- Iteratively refining chatbot features based on workshop outcomes

Bug Tracking and Feature Requests

Managing bug tracking and feature requests (*Atlasian 2024*) involves:

- Implementing a user-friendly reporting system
- Providing transparent communication on issue resolution and feature implementation

6.6 Impact Assessment

Evaluating the chatbot’s effect on student support:

Educational Outcomes Analysis

Analysing educational outcomes (*Sclater, Peasgood, and Mullan 2023*) includes:

- Conducting correlation studies between chatbot usage and academic performance
- Performing longitudinal studies on student retention and progression

Wellbeing Indicators

Assessing wellbeing indicators (*Diener, Oishi, and Tay 2024, pp. 100-150*) involves:

- Surveying student stress levels and coping mechanisms
- Analysing the chatbot's role in early intervention for mental health issues

Resource Utilisation Metrics

Measuring resource utilisation (*Heick 2023*) includes:

- Tracking changes in staff workload and time allocation
- Assessing efficiency gains in student support processes

6.7 Security and Privacy Audits

Regularly assessing and improving data protection:

Penetration Testing

Conducting penetration testing (*Stuttard and Pinto 2023, pp. 200-250*) involves:

- Scheduling security assessments by external experts
- Implementing continuous automated vulnerability scanning

Data Protection Impact Assessments (DPIAs)

Performing Data Protection Impact Assessments (*Information Commissioner's Office 2024*) includes:

- Regularly reviewing data collection and usage practices
- Conducting compliance checks with evolving data protection regulations
- Utilising the DPIA template provided in Appendix A

6.8 Continuous Professional Development

Keeping the development team updated:

Industry Conference Participation

Encouraging industry conference participation (*Institute of Electrical and Electronics Engineers 2024*) includes:

- Attending and presenting at relevant academic and industry events
- Networking with experts in educational technology and AI ethics

Research Collaboration

Fostering research collaboration (*Dillenbourg 2023, pp. 50-100*) involves:

- Establishing partnerships with universities for cutting-edge research
- Publishing findings to contribute to the broader field

7 Key Recommendations

The implementation of 'Alice', the student support chatbot, at South Star Academy presents both significant opportunities and challenges. By adhering to ethical principles, legal standards, and best practices in technical implementation, user experience, and continuous improvement, TechSoft can develop a chatbot that enhances student support while maintaining privacy, security, and inclusivity.

Key recommendations for ethical continuation of the 'Alice' project:

- Establishing a robust ethical framework and governance structure
- Ensuring compliance with all relevant data protection and education sector regulations
- Implementing state-of-the-art AI and machine learning technologies with a focus on security and scalability
- Prioritising user-centred design and accessibility to cater to diverse student needs
- Developing comprehensive monitoring and evaluation systems for continuous improvement

By following these recommendations, TechSoft can create a chatbot that not only meets the immediate needs of South Star Academy but also sets a new standard for ethical and effective AI implementation in educational settings.

A Appendix A: Data Protection Impact Assessment Template

[Include a template or example of a Data Protection Impact Assessment (DPIA) tailored for the 'Alice' chatbot project]

B Appendix B: Ethical AI Checklist

[Provide a comprehensive checklist for ensuring ethical AI development and deployment, specific to the educational context of 'Alice']

C Appendix C: Sample Conversational Flows

[Provide examples of conversational flows for common student support scenarios]

References

- AI Ethics Board Best Practices (2024). *Guidelines for establishing and managing AI ethics boards*. URL: <https://www.aiethicsboard.org/best-practices> (visited on 04/10/2024).
- AI Ethics Guidelines for Education (2024). *Best practices for implementing AI in educational settings*. URL: <https://www.aieducationethics.org/guidelines> (visited on 04/10/2024).
- American Psychological Association (2024). *Guidelines for the practice of telepsychology*. URL: <https://www.apa.org/practice/guidelines/telepsychology> (visited on 04/10/2024).
- Anastasiou, D. and R. Schäler (2023). *Translating Vital Information: Localisation, Internationalisation, and Globalisation*. Routledge.
- Annus, N. (2023). "Chatbots in Education – the impact of Artificial Intelligence based ChatGPT on Teachers and Students". In: *International Journal of Advanced Natural Sciences and Engineering Researches* 7.4, pp. 366–370.
- Armstrong, T. (2023). *Neurodiversity in the Classroom: Strength-Based Strategies to Help Students with Special Needs Succeed in School and Life*. ASCD.
- Arrieta, A.B. et al. (2022). "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI". In: *Information Fusion* 58, pp. 82–115.
- Article 29 Working Party (2018). *Guidelines on consent under Regulation 2016/679*. URL: <https://ec.europa.eu/newsroom/article29/items/623051> (visited on 04/10/2024).
- Association for Computing Machinery (2023). "ACM Code of Ethics and Professional Conduct". In: *Communications of the ACM* 66.1, pp. 24–31.
- Association for Computing Machinery Conference on Fairness, Accountability, and Transparency (2024). "Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency". In: New York: ACM.
- Atlassian (2024). *Jira Software Cloud Documentation*. URL: <https://support.atlassian.com/jira-software-cloud/> (visited on 04/10/2024).
- Barocas, S., M. Hardt, and A. Narayanan (2021). *Fairness and Machine Learning: Limitations and Opportunities*. URL: <https://fairmlbook.org/> (visited on 04/10/2024).
- Beasley, M. (2023). *Practical Web Analytics for User Experience: How Analytics Can Help You Understand Your Users*. 2nd ed. Morgan Kaufmann.
- Bickmore, T.W. et al. (2021). "Patient and consumer safety risks when using conversational assistants for medical information: an observational study of Siri, Alexa, and Google Assistant". In: *Journal of Medical Internet Research* 20.9, e11510.
- Blundell, Boyd (2020). *Computer Ethics and Professional Responsibility*. Oxford, UK: Oxford University Press.
- Bradbury, A. (2024). *Successful Presentation Skills*. 6th ed. Kogan Page.
- Brignull, H. (2023). *Dark Patterns: Inside the interfaces designed to trick you*. URL: <https://www.darkpatterns.org/> (visited on 04/10/2024).
- British Computer Society (2024). *Code of Conduct for BCS Members*. URL: <https://www.bcs.org/membership-and-registrations/become-a-member/bcs-code-of-conduct/> (visited on 04/10/2024).
- Champeon, S. and N. Finck (2024). *Inclusive Design Patterns: Coding Accessibility into Web Design*. Smashing Magazine.

- Chartered Institute of Personnel and Development (2024). *Continuing Professional Development: Guidelines and Best Practices*. CIPD. URL: <https://www.cipd.co.uk/knowledge/fundamentals/people/development/continuing-professional-development-factsheet> (visited on 04/10/2024).
- Department for Education (2024a). *Keeping Children Safe in Education: Statutory guidance for schools and colleges*. URL: <https://www.gov.uk/government/publications/keeping-children-safe-in-education--2> (visited on 04/10/2024).
- (2024b). *Special Educational Needs and Disability Code of Practice: 0 to 25 years*. URL: <https://www.gov.uk/government/publications/send-code-of-practice-0-to-25> (visited on 04/10/2024).
- Diener, E., S. Oishi, and L. Tay (2024). *Handbook of well-being research*. 2nd ed. Noba Scholar.
- Dignum, V. (2023). *Responsible Artificial Intelligence: How to Develop and Use AI in a Responsible Way*. Cham, Switzerland: Springer Nature.
- Dillenbourg, P. (2023). *Artificial Intelligence in Education: Promises and Challenges*. Cambridge, UK: Cambridge University Press.
- El Emam, K. and L. Arbuckle (2023). *Anonymizing Health Data: Case Studies and Methods to Get You Started*. 2nd ed. O’Reilly Media.
- European Commission (2024). *Ethics Guidelines for Trustworthy AI*. URL: <https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai> (visited on 04/10/2024).
- European Data Protection Board (2023). *Guidelines on consent under Regulation 2016/679*. URL: https://edpb.europa.eu/our-work-tools/our-documents/guidelines/guidelines-052020-consent-under-regulation-2016679_en (visited on 04/10/2024).
- European Union (2016). *General Data Protection Regulation (GDPR)*. URL: <https://eur-lex.europa.eu/eli/reg/2016/679/oj> (visited on 04/10/2024).
- Federal Trade Commission (2023). *Children’s Online Privacy Protection Rule: A Six-Step Compliance Plan for Your Business*. URL: <https://www.ftc.gov/tips-advice/business-center/guidance/childrens-online-privacy-protection-rule-six-step-compliance> (visited on 04/10/2024).
- Floridi, L. and J. Cowls (2023). “A Unified Framework of Five Principles for AI in Society”. In: *Harvard Data Science Review* 1.1.
- Frain, B. (2023). *Responsive Web Design with HTML5 and CSS: Develop future-proof responsive websites using the latest HTML5 and CSS techniques*. 4th ed. Packt Publishing.
- Gao, J., M. Galley, and L. Li (2023). *Neural Approaches to Conversational AI: Question Answering, Task-Oriented Dialogues and Social Chatbots*. Now Publishers.
- Géron, A. (2024). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. 3rd ed. O’Reilly Media.
- Goodman, E., M. Kuniavsky, and A. Moed (2023). *Observing the User Experience: A Practitioner’s Guide to User Research*. 3rd ed. Morgan Kaufmann.
- Harper, S. and Y. Yesilada (2024). *Web Accessibility: A Foundation for Research*. 3rd ed. Springer.
- Hartzog, W. (2023). *Privacy’s Blueprint: The Battle to Control the Design of New Technologies*. Harvard University Press.
- Heick, T. (2023). *A Complete Guide To Student Data & Privacy For Teachers*. URL: <https://www.teachthought.com/technology/student-data-privacy-guide-for-teachers/> (visited on 04/10/2024).
- Holtzblatt, K. and H. Beyer (2024). *Contextual Design: Design for Life*. 3rd ed. Morgan Kaufmann.
- Hoober, S. and E. Berkman (2023). *Designing Mobile Interfaces: Patterns for Interaction Design*. 2nd ed. O’Reilly Media.
- IEEE Ethics Certification Program for Autonomous and Intelligent Systems (2024). *Certification criteria and process*. URL: <https://standards.ieee.org/industry-connections/ecpais/> (visited on 04/10/2024).
- IMS Global Learning Consortium (2024). *Learning Tools Interoperability (LTI) Standards*. URL: <https://www.imsglobal.org/activity/learning-tools-interoperability> (visited on 04/10/2024).
- Information Commissioner’s Office (2024). *Data protection guidance*. URL: <https://ico.org.uk/for-organisations/guide-to-data-protection/> (visited on 04/10/2024).
- Institute of Electrical and Electronics Engineers (2023). *Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems*. URL: <https://standards.ieee.org/industry-connections/ec/autonomous-systems/> (visited on 04/10/2024).
- (2024). *IEEE Code of Ethics*. Professional Guidelines. New York, NY: IEEE.
- International Organization for Standardization (2022). *ISO/IEC 27001:2022 Information security management systems — Requirements*. International Standard ISO/IEC 27001:2022. Geneva, Switzerland: ISO.

- Josuttis, N.M. (2023). *Cloud Native Transformation: Practical Patterns for Innovation*. Addison-Wesley Professional.
- Jurafsky, D. and J.H. Martin (2024). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. 3rd ed. Upper Saddle River, NJ: Prentice Hall.
- Kohavi, R., D. Tang, and Y. Xu (2023). *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing*. Cambridge University Press.
- Lemon, O. and O. Pietquin (2023). *Data-Driven Methods for Adaptive Spoken Dialogue Systems: Computational Learning for Conversational Interfaces*. Springer.
- Liu, B. (2023). *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. 2nd ed. Cambridge University Press.
- Mehrabi, N. et al. (2023). “A survey on bias and fairness in machine learning”. In: *ACM Computing Surveys* 54.6, pp. 1–35.
- Miner, A.S. et al. (2022). “Key considerations for incorporating conversational AI in psychotherapy”. In: *Frontiers in Psychiatry* 10, p. 746.
- Moore, R.J. and R. Arar (2023). *Conversational UX Design: A Practitioner’s Guide to the Natural Conversation Framework*. ACM Books.
- National Cyber Security Centre (2024). *Cyber security guidance for educational institutions*. URL: <https://www.ncsc.gov.uk/section/education-skills/cyber-security-schools> (visited on 04/10/2024).
- Norman, D. (2023). *Emotional Design: Why We Love (or Hate) Everyday Things*. 2nd ed. Basic Books.
- Organisation for Economic Co-operation and Development (2023). *OECD Principles on Artificial Intelligence*. URL: <https://www.oecd.org/going-digital/ai/principles/> (visited on 04/10/2024).
- Quarteroni, S. et al. (2024). “Chatbot Evaluation Metrics: State of the Art and Future Directions”. In: *Dialogue & Discourse* 12.1, pp. 1–32.
- Scater, N., A. Peasgood, and J. Mullan (2023). *Learning Analytics in Higher Education: A review of UK and international practice*. Jisc.
- Simonsen, J. and T. Robertson (2023). *Routledge International Handbook of Participatory Design*. 2nd ed. Routledge.
- Stallings, W. (2023). *Cryptography and Network Security: Principles and Practice*. 8th ed. Pearson.
- Stallings, W. and L. Brown (2024). *Computer Security: Principles and Practice*. 5th ed. Pearson.
- Stuttard, D. and M. Pinto (2023). *The Web Application Hacker’s Handbook: Finding and Exploiting Security Flaws*. 3rd ed. Wiley.
- Tullis, T. and B. Albert (2024). *Measuring the User Experience: Collecting, Analyzing, and Presenting Usability Metrics*. 3rd ed. Morgan Kaufmann.
- UK Government (2008). *Education and Skills Act 2008*. URL: <https://www.legislation.gov.uk/ukpga/2008/25/contents> (visited on 04/10/2024).
- (2018). *Data Protection Act 2018*. URL: <https://www.legislation.gov.uk/ukpga/2018/12/contents/enacted> (visited on 04/10/2024).
- UNESCO (2021). *Recommendation on the Ethics of Artificial Intelligence*. URL: <https://en.unesco.org/artificial-intelligence/ethics> (visited on 04/10/2024).
- United Nations Educational, Scientific and Cultural Organization (2023). *Recommendation on the Ethics of Artificial Intelligence*. URL: <https://en.unesco.org/artificial-intelligence/ethics> (visited on 04/10/2024).
- Vadapalli, S. (2023). *DevOps and Site Reliability Engineering (SRE) Handbook: Guide to Site Reliability Engineering and DevOps Practices*. Packt Publishing.
- Vaughan, J.W. (2024). *Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI*. Manning Publications.
- WebAIM (2024). *Web Accessibility In Mind: Screen Reader User Survey Results*. URL: <https://webaim.org/projects/screenreadersurvey9/> (visited on 04/10/2024).
- Williams, J. (2024). *Stand Out of Our Light: Freedom and Resistance in the Attention Economy*. Cambridge University Press.
- World Wide Web Consortium (2023). *Web Content Accessibility Guidelines (WCAG) 2.2*. URL: <https://www.w3.org/TR/WCAG22/> (visited on 04/10/2024).
- Yesilada, Y., G. Brajnik, and S. Harper (2023). “How much does web accessibility cost?” In: *Proceedings of the 20th International Conference on World Wide Web*, pp. 1–10.