

# L'intelligenza artificiale negli archivi: soluzioni in campo e casi concreti

📅 10 giugno 2025 - ore 11:30

Termine iscrizioni: 09 giugno 2025 ore 14:00



L'intelligenza artificiale è sempre più presente anche nel dibattito sugli archivi e sulla gestione documentale. Ma cosa è davvero realtà e cosa narrazione?

Il webinar vuole offrire un momento di riflessione sulla crescente crucialità dell'intelligenza artificiale nel governo delle fonti, soprattutto in relazione alla garanzia della loro autenticità e affidabilità, ruolo questo, da sempre affidato all'archivista.

L'obiettivo è fornire, per quanto possibile, uno sguardo critico e aggiornato sull'impiego concreto dell'IA nel settore archivistico, un'occasione per distinguere tra narrazione e pratica e per capire dove stia andando davvero l'innovazione negli archivi.

L'evento sarà inoltre occasione per offrire una panoramica su alcune applicazioni già operative e su soluzioni offerte dal mercato, sia in Italia che all'estero, nonché per raccontare un caso di applicazione pratica, in cui l'IA è stata impiegata per la valutazione dei formati dei file, in particolare elaborando l'indice di interoperabilità descritto nell'Allegato 2 alle Linee guida sulla formazione, gestione e conservazione del documento informatico di Agid.

Il webinar si inserisce all'interno di un ciclo previsto dall'Accordo di collaborazione sulla conservazione tra **Formez PA e Regione Emilia-Romagna (Parer)**.

Il termine per le iscrizioni al webinar è **lunedì 9 giugno 2025 - 14:00, salvo esaurimento posti**.

## ☰ Programma

**11.30 - 11.35** Introduzione (**Alessandra Cornero**, Formez PA)

**11.35 - 11.55** Intelligenza artificiale e archivi: a che punto siamo, oggi (**Giorgia Di Marcantonio**, Università degli Studi di Napoli, l'Orientale)

**11.55 - 12.15** Classificare e fascicolare i documenti con l'AI: l'esperienza del progetto InterPares Trust (**Stefano Allegrezza**, Università degli Studi di Macerata)

**12.15 - 12.35** Calcolo automatico dell'indice di interoperabilità dei formati (**Marianna Tascone**, Polo archivistico e gestione documentale della Regione Emilia-Romagna- ParER)

**12.35 - 13.00** Domande e risposte

## Partecipanti

857 iscritti a questo evento

## Contatti

✉ [ufficio.competenzadocumentale@formez.it](mailto:ufficio.competenzadocumentale@formez.it)

## Progetto di riferimento

Attività istituzionali

## Ulteriori informazioni

Argomento: *Conservazione, Gestione documentale, Intelligenza artificiale*

Politica di registrazione: *Senza approvazione*

## 🔔 Rimani aggiornato

Accedi per rimanere aggiornato

# Classificare e fascicolare i documenti con l'AI: l'esperienza del progetto InterPares Trust AI

Stefano Allegrezza  
Università degli Studi di Macerata  
[stefano.allegrezza@unimc.it](mailto:stefano.allegrezza@unimc.it)

# Agenda

---

1. Introduzione
2. La ricerca condotta sull'impiego dell'AI per la classificazione e fascicolazione nell'ambito del progetto InterPares Trust AI
3. Risultati
4. Considerazioni finali

# **1. Introduzione**

# La classificazione dei documenti informatici

## [Linee guida AgID, par. 3. Classificazione dei documenti informatici]

La **classificazione** ha il fine di organizzare logicamente tutti i documenti amministrativi informatici prodotti o ricevuti da un ente nell'esercizio delle sue funzioni. L'attività di classificazione si avvale del **piano di classificazione** che mappa, su più livelli gerarchici, tutte le funzioni dell'ente.

La **classificazione** è un'**attività obbligatoria** nel sistema di gestione informatica dei documenti dell'AOO e si applica a tutti i documenti prodotti e acquisiti dalla stessa AOO sottoposti o meno alla registrazione di protocollo, ai sensi degli articoli 56 e 64, comma 4, del TUDA. Le informazioni relative alla classificazione nei casi dei documenti amministrativi informatici costituiscono parte integrante dei metadati previsti per la formazione dei documenti medesimi.

Nel sistema di gestione informatica dei documenti dell'AOO l'**attività di classificazione** guida la formazione dell'archivio mediante il **piano di organizzazione delle aggregazioni documentali**.

## [Linee guida AgID, All. 1 «Glossario dei termini e degli acronimi»]

**Piano di classificazione:** «Struttura logica che permette di organizzare documenti e oggetti digitali secondo uno schema desunto dalle funzioni e dalle attività dell'amministrazione interessata».

**[Piano di classificazione per gli archivi dei comuni italiani]:** «Il piano di classificazione o titolario è il sistema preconstituito di partizioni astratte, gerarchicamente ordinate (dal generale al particolare), fissate sulla base dell'analisi delle funzioni dell'ente, al quale deve ricondursi la molteplicità dei documenti prodotti, per organizzarne la sedimentazione ordinata».



**Senza la classificazione non si produce un archivio ma un insieme di documenti senza alcuna relazione tra loro, che rende difficile lavorare e inefficiente la ricerca**

# La formazione delle aggregazioni documentali

[Linee guida AgID, par. 3.3 «Aggregazioni documentali informatiche»]

La Pubblica Amministrazione documenta la propria attività tramite funzioni del sistema di gestione informatica dei documenti finalizzate alla produzione, alla gestione e all'uso delle **aggregazioni documentali informatiche**, corredate da opportuni metadati, così come definiti nell'allegato 5 “Metadati” alle presenti Linee guida.

[Linee guida AgID, All. 1 «Glossario dei termini e degli acronimi»]

**Piano di organizzazione delle aggregazioni documentali:**

«Strumento integrato con il sistema di classificazione a partire dai livelli gerarchici inferiori di quest'ultimo e finalizzato a individuare le tipologie di aggregazioni documentali (tipologie di serie e tipologie di fascicoli) che devono essere prodotte e gestite in rapporto ai procedimenti e attività in cui si declinano le funzioni svolte dall'ente».



# La domanda di ricerca

## Esempio n. 1

In molte amministrazioni pubbliche e aziende private i **documenti non sono né classificati né aggregati**.

In altri casi, le aggregazioni documentali **non sono ben create**, con il risultato di un numero incontrollato di documenti **non ordinati**, non collocati nella cartella corretta e **difficili da trovare**.

In molti casi **mancano i metadati**, necessari per garantire l'affidabilità, l'attendibilità, la qualità e la sostenibilità delle valutazioni e delle acquisizioni.

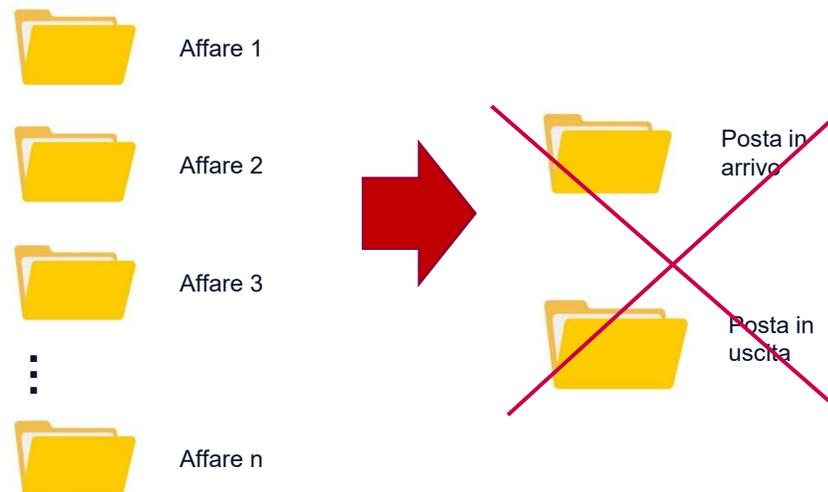
Nonostante i progressi compiuti nelle varie tecnologie a supporto della gestione dei documenti, il supporto software per queste attività rimane limitato.



## Esempio n. 2

**La gestione della posta elettronica** è diventata **una delle attività che richiedono più tempo** sia nel settore pubblico che nelle aziende private e nelle attività personali.

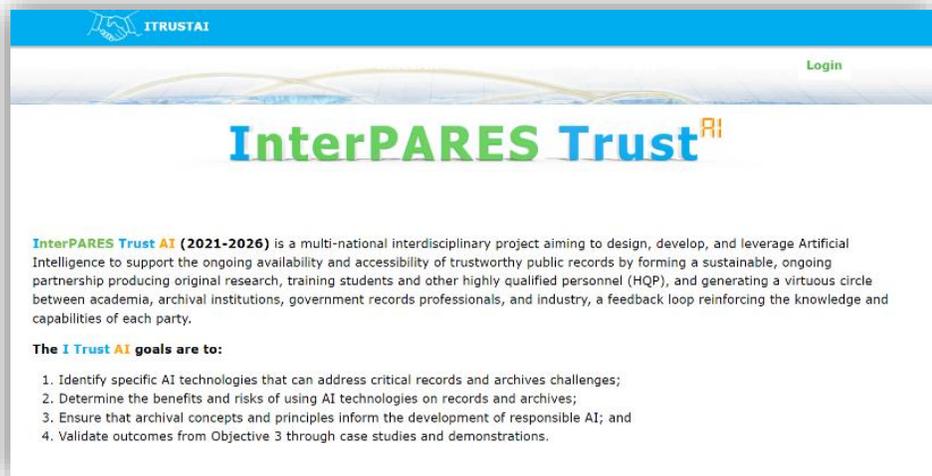
Le e-mail sono spesso gestite come singoli record senza alcun legame con altre e-mail e **non sono classificate né organizzate in aggregazioni archivistiche (cartelle)**; inoltre non sono collegate e classificate nel sistema di gestione documentale del soggetto produttore.



È possibile utilizzare gli **strumenti di IA** per **classificare** i documenti e creare (o ricreare) le **aggregazioni documentali** e individuare i **metadati**?

## **2. La ricerca condotta sull'impiego dell'AI per la classificazione e fascicolazione**

# Il Progetto InterPARES Trust AI e il gruppo di lavoro CU05



**InterPARES Trust AI (2021-2026)** è un progetto **interdisciplinare multinazionale** che mira a:

- identificare **tecnologie specifiche di intelligenza artificiale** in grado di affrontare le sfide critiche relative ai documenti e agli archivi;
- determinare i **benefici e i rischi dell'uso delle tecnologie AI su documenti e archivi**;
- garantire che i **concetti e i principi archivistici** siano tenuti in conto nello sviluppo di un'IA responsabile;
- convalidare i risultati attraverso **studi di caso**.

<https://interparestrustai.org/>

The role of AI in identifying or reconstituting archival aggregations of digital records and enriching metadata schemas

CU05

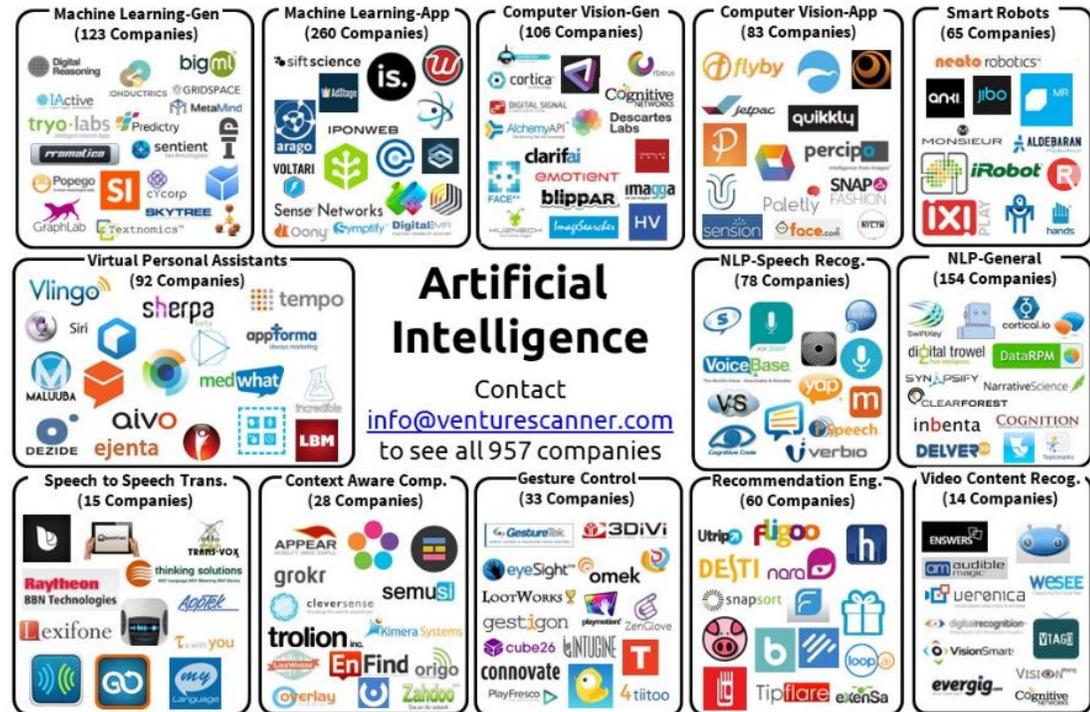
**Stefano Allegrezza** (co-chair) (**Università di Macerata**)  
**Mariella Guercio** (co-chair) (Associazione nazionale archivistica italiana - **ANAI**)  
**Massimiliano Grandi** (Associazione nazionale archivistica italiana - **ANAI**)  
**Bruna La Sorda** (Associazione nazionale archivistica italiana - **ANAI**)  
**Francesca Magnoni** (North Atlantic Treaty Organization - **NATO**)  
**Maria Mata Caravaca** (International Centre for the Study of the Preservation and Restoration of Cultural Property - **ICCROM**)  
**Samir Musa** (Historical Archives of European Union – **HAEU**)  
**Luís-Esteve Casellas Serra**, **Municipality of Girona – Spain** (connection with AA01 “Employing AI for Retention & Disposition in Digital Information and Recordkeeping Systems (DIRS)”)  
**Nicola Di Matteo** (**Halifax University**, Canada)  
**Eleonora Luzi** (Associazione nazionale archivistica italiana - **ANAI**)  
**Gianni Penzo Doria** (**Università dell’Insubria**)  
**Gabriele Bezzi** (Associazione nazionale archivistica italiana - **ANAI**)  
**Stefano Delli Ponti**, **Roberto D’Ippolito** (**SIAV**)  
**Marianna Tascone** (**PAR-ER, Regione Emilia-Romagna**)

# La domanda di ricerca

È possibile utilizzare gli **strumenti di IA** al fine di **classificare** i documenti, creare (o ricreare) le **aggregazioni documentali**, l'individuazione dei **metadati**?



Ci sono **migliaia di aziende** che dichiarano di utilizzare l'IA. **Centinaia** di loro dichiarano di utilizzare le tecnologie AI sul campo o gli ERMS/EDMS.



Quali **tecnologie AI** potrebbero essere utili per l'esecuzione automatica o semi-automatica di operazioni archivistiche quali:

- la **classificazione dei documenti**?
- la creazione delle **aggregazioni documentali**?
- l' **integrazione dei metadati**?
- la gestione delle **aggregazioni** nel caso della posta elettronica?



Come verificare le funzionalità degli strumenti proposti da queste aziende e come verificare se questi strumenti funzionano davvero?

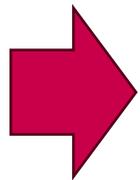
# Fase 1: Identificazione delle aziende di IA

Identificazione di un gruppo iniziale di **300 aziende** interessanti per lo studio, ovvero aziende che sviluppano soluzioni basate su tecnologie legate all'IA e sono pertinenti all'ambito dello studio CU05.

Strumenti per la costruzione dell'elenco:

- **ricerche dirette su Internet** utilizzando parole chiave e stringhe di testo;
- **le risorse e le conoscenze messe a disposizione dai professionisti** (Alan Pelz-Sharpe, Andrew Warland, James Lappin, Jenny Bunn e Paul Young)

Nr.	Company	Location	Website	Rating	Notes
1	a3doc Wolters Kluwer	Alphen aan den Rijn, The Netherlands	<a href="http://www.a3doc.com/experiencia-clientes-gest-ion-documental-cloud-colaborativa.html">http://www.a3doc.com/experiencia-clientes-gest-ion-documental-cloud-colaborativa.html</a>	1	a3doc Wolters Kluwer has developed a3doc cloud, an application for document management which is able to classify and store automatically corporate documents. There is no much information in the webpage describing the product. The company is based in the Netherlands, but the webpage is in Spanish.
2	Abbyy	California, United States	<a href="https://www.abbyy.com/">https://www.abbyy.com/</a>	1	The main headquarters of Abbyy are Milpitas, California, USA. Abbyy mainly develops application for data capture, but one of their software products is ABBYY Vantage, software for Intelligent Document Processing. ABBYY Vantage uses AI-driven technologies to process "documents of any kind—structured, semi-structured, or unstructured, and all type of data including machine printed, hand printed, barcodes, signatures, and check boxes". In ABBYY Vantage "Trained skills can be quickly designed to understand and extract information from all types of documents". "Once skills are deployed, Vantage then monitors, measures, and analyses performance of all your deployed skills, creating new learning models—so you can continuously improve and move your automation to the next level." ABBYY "Vantage skills are continuously getting smarter and more accurate over time, as new document variations and statistical data is collected during human-in-the-loop review". "The easy-to-use, no-code Vantage platform can be utilized to set up and train Document, Classification, and Process Skills for just about any document type and flow" - For this information see <a href="https://www.abbyy.com/vantage/">https://www.abbyy.com/vantage/</a>
3	ActiveNav	Reston, Virginia, USA	activenav.com	1	ActiveNav is based in US, UK and Australia. They seem to deal mainly with automated data mapping and automated data classification. As it is the case for Automated Intelligence, you do not find a specific mention of records and archives, but clearly the services they offer (or at least say they can offer) also concern archives and records management.
4	Acumatica - Webiplex - PairSoft	Washington, USA (Acumatica) Florida, USA (PairSoft)	<a href="https://www.acumatica.com/media/2017/04/DocuPeak-for-Acumatica.pdf">https://www.acumatica.com/media/2017/04/DocuPeak-for-Acumatica.pdf</a>	1	DocuPeak is an application developed by Webiplex - that has been recently purchased by PairSoft, a company based in Florida, USA. DocuPeak has been built to be integrated with the Enterprise Resource Planning platform created by Acumatica, whose main headquarters is in Seattle, Washington, USA. DocuPeak is powered by Robotic Process Automation technologies. This is the introductory description of DocuPeak: "Robotic Process Apps built on the DocuPeak cloud platform streamline operations 'before and including data entry', from automated data extraction and data entry, to approval workflows and document lifecycle management to entirely electronic forms-based applications". Some features of DocuPeak are "Leverage Smart Document Recognition (SDR) to extract key data from documents, without templates, automating document indexing and transactional data entry"; "DocuPeak's Rules Engine to automatically route documents such as AP invoices through a predefined role based review and approval process, including notifications and escalation procedures"; "Integrate all documents to the associated transaction within Acumatica for synchronization of key document data and instant retrieval"; "Create a secure, compliant document management environment, maintaining an audit trail of all document-based activity, including check-in, check-out and version control on changed documents".
5	Ademero	Florida, United States	<a href="https://www.ademero.com/">https://www.ademero.com/</a>	1	The company is based in Lakeland, Florida, USA. Ademero is a Document Management Software platform. Thanks to its AI capabilities, Ademero can "intelligently identify, categorize and process accounts payable invoices or any other low or high volume paper entering Content Central" - see <a href="https://www.ademero.com/document-management-software/">https://www.ademero.com/document-management-software/</a> . Ademero includes 2 modules, one for document scanning / capture, and another one for document management after the capture of the document. "Capture Software lays the foundation for an automated onboarding process by both classifying and indexing documents, then working hand-in-hand with a Document Management System and your other business software solutions, streamlines document and office workflows. "Intelligent" or fully-automatic solutions in this category means that your staff simply scan in paper documents at your office multifunction printer (mfp) and the software uses optical character recognition (OCR) to turn that scanned image into a digital version of that document, then handles classifying and indexing each document before handing it off to your DMS or other software applications."As to the Document Management module, one of its most interesting features is "the logical and consistent folder and file building it provides". "You just upload your document and the system will handle filing it away so that you or anyone who has permission to access those documents can locate it quickly by navigating logically named folders based on the standards you require".
6	Adlib	Burlington, Ontario, Canada	<a href="http://www.adlibsoftware.com">www.adlibsoftware.com</a>	3	Adlib is a Canadian company which has produced "Adlib Elevate". "Adlib Elevate" is a File Analytics platform to automate discovery, extraction and classification of vital data from complex documents to streamline data-intensive processes and accelerate process automation. Adlib Elevate is one of the 5 suppliers AI-based application for records management assessed by The UK National Archives.
7	Aida Cloud	Turin, Italy	<a href="https://www.aidacloud.com/home">https://www.aidacloud.com/home</a>	0	AIDA has been developed by Technology & Cognition LAB, based in Turin, Italy. AIDA is a document retrieval application that uses AI "to recognise any type of document and, depending on the user's needs, extract all the information needed, with a simple learning process that requires no technical knowledge". Basically it retrieves the information you need, organizes it and makes it available for users - see <a href="https://doc.aidacloud.com/aida/">https://doc.aidacloud.com/aida/</a> Users define document types, and AIDA through Intelligent Document Analysis manages to recognise such document



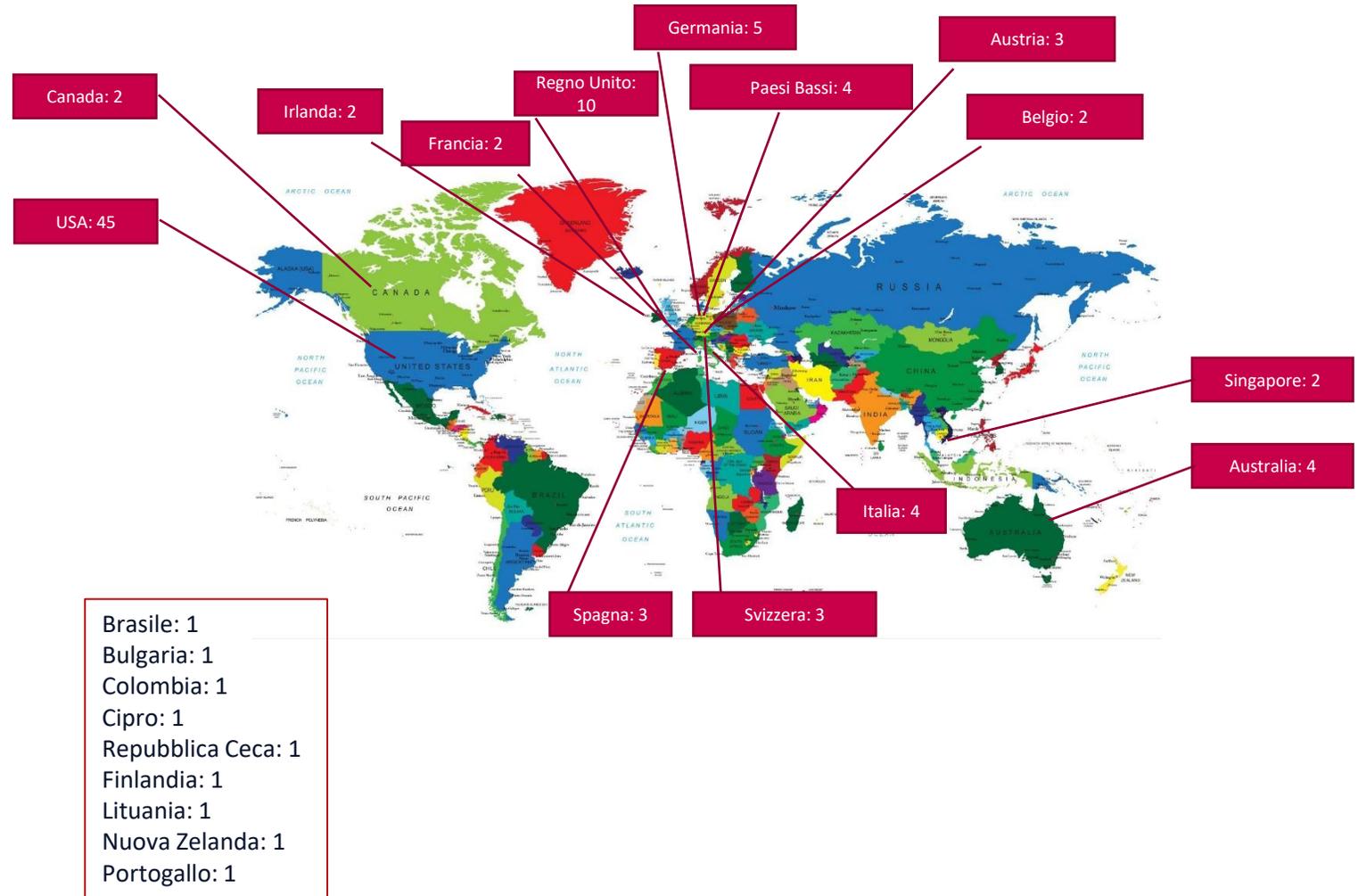
L'elenco non è né esaustivo né definitivo, ma è un punto di partenza.

# Fase 1: 100 aziende - distribuzione geografica

L'elenco è stato successivamente limitato a **100 aziende**

Le caratteristiche dei loro **software di intelligenza artificiale** sono state analizzate in base alle **informazioni disponibili sui loro siti web** con riferimento a:

- in cui l'azienda dichiara che la **gestione dei documenti** è uno degli obiettivi delle sue applicazioni di intelligenza artificiale;
- **evidenza di interesse** per gli aspetti del **records management** e degli **archives** (anche se in alcuni casi ciò non è dichiarato apertamente ma si può solo intuire dal contenuto del sito web).



# Fase 1: Identificazione delle aziende di IA

Poiché non era possibile intervistare tutte le **100 aziende**, dall'elenco iniziale abbiamo selezionato un elenco di **28 aziende** sulla base di:

- il loro **portafoglio**
- il loro coinvolgimento diretto nel **campo della gestione dei documenti**
- la loro **conformità** ai quadri normativi e agli standard rilevanti per il settore
- **la reputazione** generale dell'azienda.

1	Microsoft	Washington, DC, USA	<a href="https://www.microsoft.com/en-gb">https://www.microsoft.com/en-gb</a>
2	Iron Mountain	Boston, Massachusetts, USA	<a href="http://www.ironmountain.com">www.ironmountain.com</a>
3	Adlib	Burlington, Ontario, Canada	<a href="http://www.adlibsoftware.com">www.adlibsoftware.com</a>
4	Castlepoint	Canberra, Australia	<a href="http://www.castlepoint.systems">www.castlepoint.systems</a>
5	Gimmel	Texas, USA	<a href="https://www.gimmel.com/">https://www.gimmel.com/</a>
6	Quest-it	Siena, Italia	<a href="http://www.quest-it.com">www.quest-it.com</a>
7	Gruppo Adapting	Valencia, Spagna	<a href="https://www.adapting.com/en/">https://www.adapting.com/en/</a>
8	Hyland	Westlake, Ohio, USA	<a href="https://www.hyland.com/en">https://www.hyland.com/en</a>
9	Stratagemma	Aurora, Colorado, USA	<a href="http://www.stratagemgroup.com">www.stratagemgroup.com</a>
10	Aluma	Cambridge, Regno Unito e New York, USA	<a href="https://aluma.io/">https://aluma.io/</a>
11	Collabware	Washington, DC, USA	<a href="http://collabware.com">collabware.com</a>
12	Ephesoft	Irvine, California, USA	<a href="https://ephesoft.com/">https://ephesoft.com/</a>
13	Leggi-Coop	Innsbruck, Austria	<a href="https://readcoop.eu/transkribus/">https://readcoop.eu/transkribus/</a>
14	Punto di registrazione	Sydney, Australia	<a href="http://www.recordpoint.com">www.recordpoint.com</a>
15	Software Prisma	California, USA	<a href="https://prismsoftware.com/">https://prismsoftware.com/</a>
16	Sistema esperto	Modena, Italia	<a href="https://www.expert.ai/">https://www.expert.ai/</a>
17	GRMGestione dei documenti	New Jersey, USA	<a href="https://www.grmdocumentmanagement.com/">https://www.grmdocumentmanagement.com/</a>
18	Grooper	Oklahoma, USA	<a href="https://www.bisok.com/intelligent-document-processing/">https://www.bisok.com/intelligent-document-processing/</a>
19	Ripcord	Hayward, California, USA	<a href="http://www.ripcord.com">www.ripcord.com</a>
20	Cortical	New York, USA	<a href="http://www.cortical.io">www.cortical.io</a>
21	AmyGB.ai	Mumbai, India	<a href="http://www.amygb.ai">www.amygb.ai</a>
22	Bizamica	Pune, India	<a href="http://www.bizamica.com">www.bizamica.com</a>
23	Docxflow	Popayán, Colombia	<a href="https://www.docxflow.com/">https://www.docxflow.com/</a>
24	IA Gleemático	Singapore	<a href="https://gleematic.com/">https://gleematic.com/</a>
25	Soluzioni aziendali SBK	São Bernardo do Campo, San Paolo, Brasile	<a href="http://www.sbkbs.com.br">www.sbkbs.com.br</a>
26	Datacentrix	Johannesburg, Sudafrica	<a href="http://www.datacentrix.co.za">www.datacentrix.co.za</a>



Anyz (Norvegia)  
DXC (Italia)

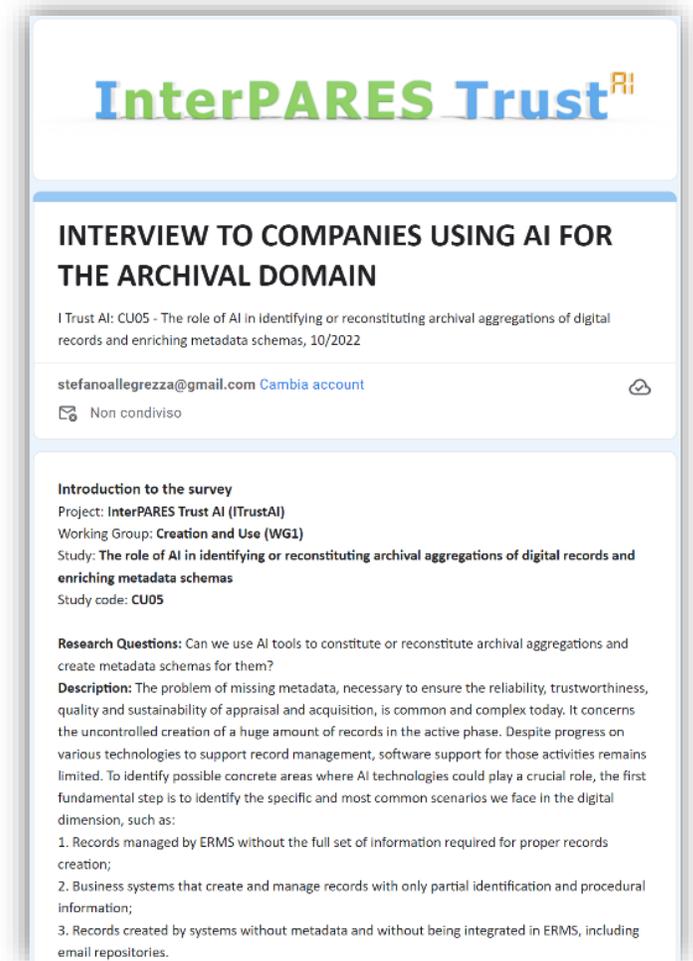
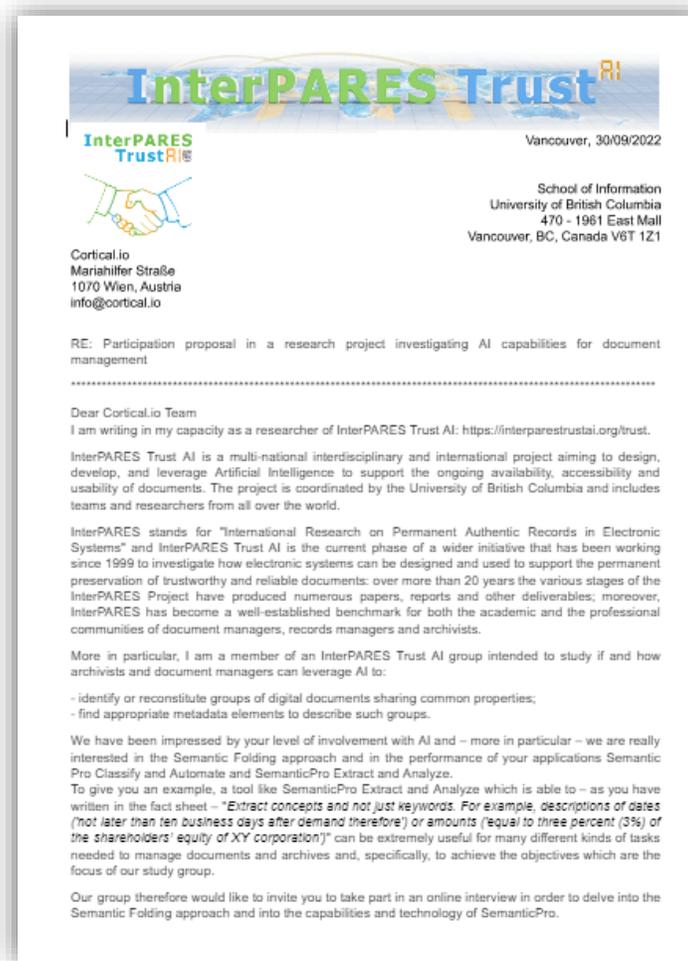
## Fase 2. Questionario e interviste

Per raccogliere informazioni più precise, abbiamo preparato un **questionario** molto dettagliato, volto a raccogliere sistematicamente le informazioni per un'adeguata valutazione delle domande.

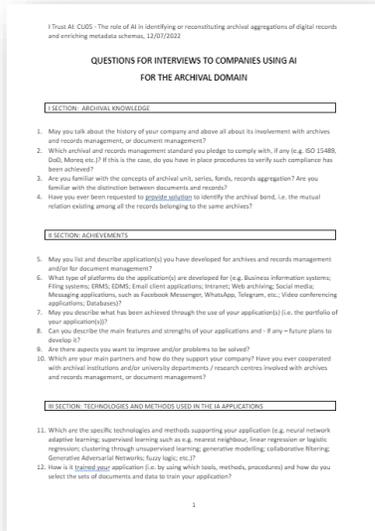
Abbiamo inviato alle **28 aziende** una **lettera di invito** ufficiale (in **inglese**, **spagnolo** o **portoghese**, a seconda della lingua preferita dall'azienda) a partecipare all'indagine.

Il questionario è stato spiegato oralmente durante un **incontro preliminare** con il personale di gestione delle informazioni e gli ingegneri informatici.

Successivamente, le aziende hanno compilato il questionario disponibile su **Google Forms**.



# Fase 2. Questionario e interviste



25 questions

I SECTION

achievements

II SECTION

specific capabilities (for recordkeeping and email systems)

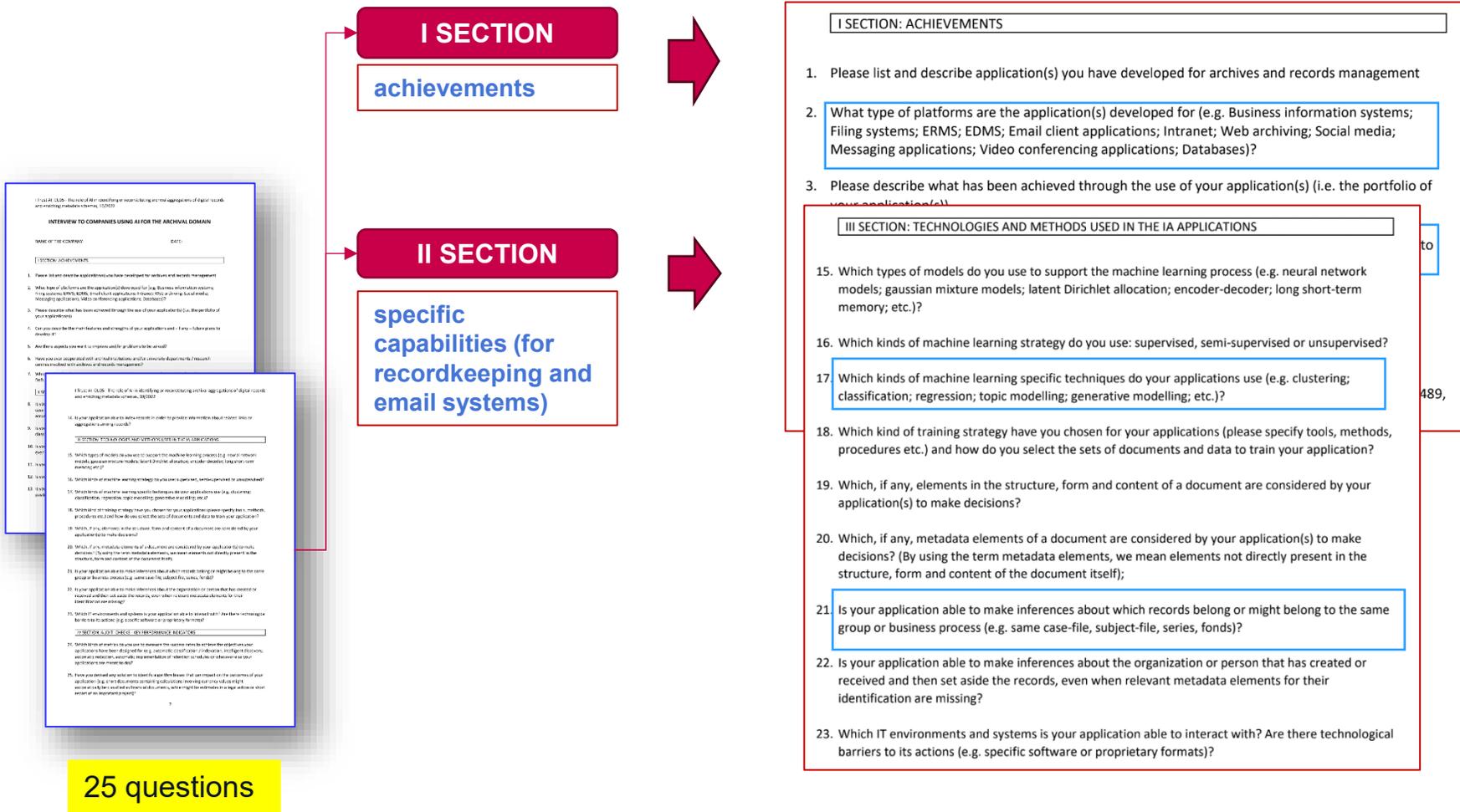
III SECTION

technologies and methods used in the IA applications

IV SECTION

audit-checks -- key performance indicators

# Fase 2. Questionario e interviste



25 questions

1. List of QUB: This tool is for identifying recordkeeping and applications of digital records and managing metadata systems. (2022)

**INTERVIEW TO COMPANIES USING AI FOR THE ARCHIVAL DOMAIN**

NAME OF THE COMPANY: \_\_\_\_\_ DATE: \_\_\_\_\_

CONTACT PERSON: \_\_\_\_\_

1. Please list and describe the applications you have developed for archives and records management.

2. What types of datasets from the applications did you use for (e.g. Business information systems; ERP systems; CRM; HRMS; email client applications; Intranet; Web archiving; Social media; Messaging applications; Video conferencing applications; Databases)?

3. Please describe what has been achieved through the use of your application(s) (i.e. the portfolio of your applications).

4. Can you describe the main features and strengths of your application(s) and if any – features you consider to be novel?

5. Are there aspects you want to improve and/or add on to be novel?

6. Have you ever cooperated with any third institutions or other entities (universities/research centers) to develop or enhance your application(s)?

7. Why?

8. How?

9. What data?

10. What type of data?

11. What type of data?

12. What type of data?

13. What type of data?

14. Is your application able to automatically generate metadata about records (e.g. subject, keywords, classification, etc.)?

**III SECTION: TECHNOLOGIES AND METHODS USED IN THE IA APPLICATIONS**

15. Which types of models do you use to support the machine learning process (e.g. neural network models; gaussian mixture models; latent Dirichlet allocation; encoder-decoder; long short-term memory; etc.)?

16. Which kinds of machine learning strategy do you use: supervised, semi-supervised or unsupervised?

17. Which kinds of machine learning specific techniques do your applications use (e.g. clustering; classification; regression; topic modelling; generative modelling; etc.)?

18. Which kind of training strategy have you chosen for your applications (please specify tools, methods, procedures etc.) and how do you select the sets of documents and data to train your application?

19. Which, if any, elements in the structure, form and content of a document are considered by your application(s) to make decisions?

20. Which, if any, metadata elements of a document are considered by your application(s) to make decisions? (By using the term metadata elements, we mean elements not directly present in the structure, form and content of the document itself);

21. Is your application able to make inferences about which records belong or might belong to the same group or business process (e.g. same case-file, subject-file, series, fonds)?

22. Is your application able to make inferences about the organization or person that has created or received and then set aside the records, even when relevant metadata elements for their identification are missing?

23. Which IT environments and systems is your application able to interact with? Are there technological barriers to its actions (e.g. specific software or proprietary formats)?

**II SECTION: SPECIFIC CAPABILITIES (FOR RECORDKEEPING AND EMAIL SYSTEMS)**

4. Can you describe the main features and strengths of your application(s) and if any – features you consider to be novel?

5. Are there aspects you want to improve and/or add on to be novel?

6. Have you ever cooperated with any third institutions or other entities (universities/research centers) to develop or enhance your application(s)?

7. Why?

8. How?

9. What data?

10. What type of data?

11. What type of data?

12. What type of data?

13. What type of data?

14. Is your application able to automatically generate metadata about records (e.g. subject, keywords, classification, etc.)?

15. Which types of models do you use to support the machine learning process (e.g. neural network models; gaussian mixture models; latent Dirichlet allocation; encoder-decoder; long short-term memory; etc.)?

16. Which kinds of machine learning strategy do you use: supervised, semi-supervised or unsupervised?

17. Which kinds of machine learning specific techniques do your applications use (e.g. clustering; classification; regression; topic modelling; generative modelling; etc.)?

18. Which kind of training strategy have you chosen for your applications (please specify tools, methods, procedures etc.) and how do you select the sets of documents and data to train your application?

19. Which, if any, elements in the structure, form and content of a document are considered by your application(s) to make decisions?

20. Which, if any, metadata elements of a document are considered by your application(s) to make decisions? (By using the term metadata elements, we mean elements not directly present in the structure, form and content of the document itself);

21. Is your application able to make inferences about which records belong or might belong to the same group or business process (e.g. same case-file, subject-file, series, fonds)?

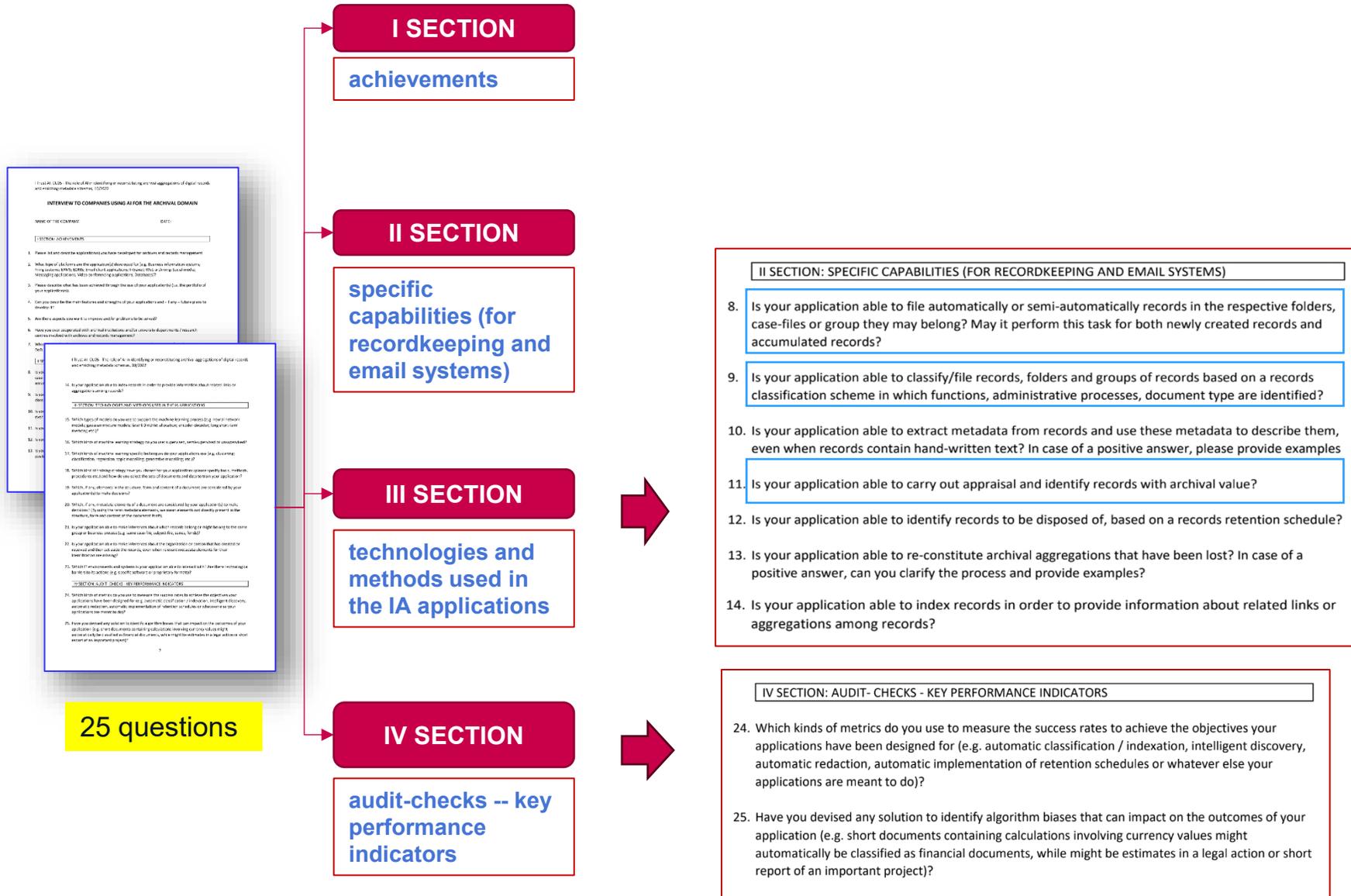
22. Is your application able to make inferences about the organization or person that has created or received and then set aside the records, even when relevant metadata elements for their identification are missing?

23. Which IT environments and systems is your application able to interact with? Are there technological barriers to its actions (e.g. specific software or proprietary formats)?

25

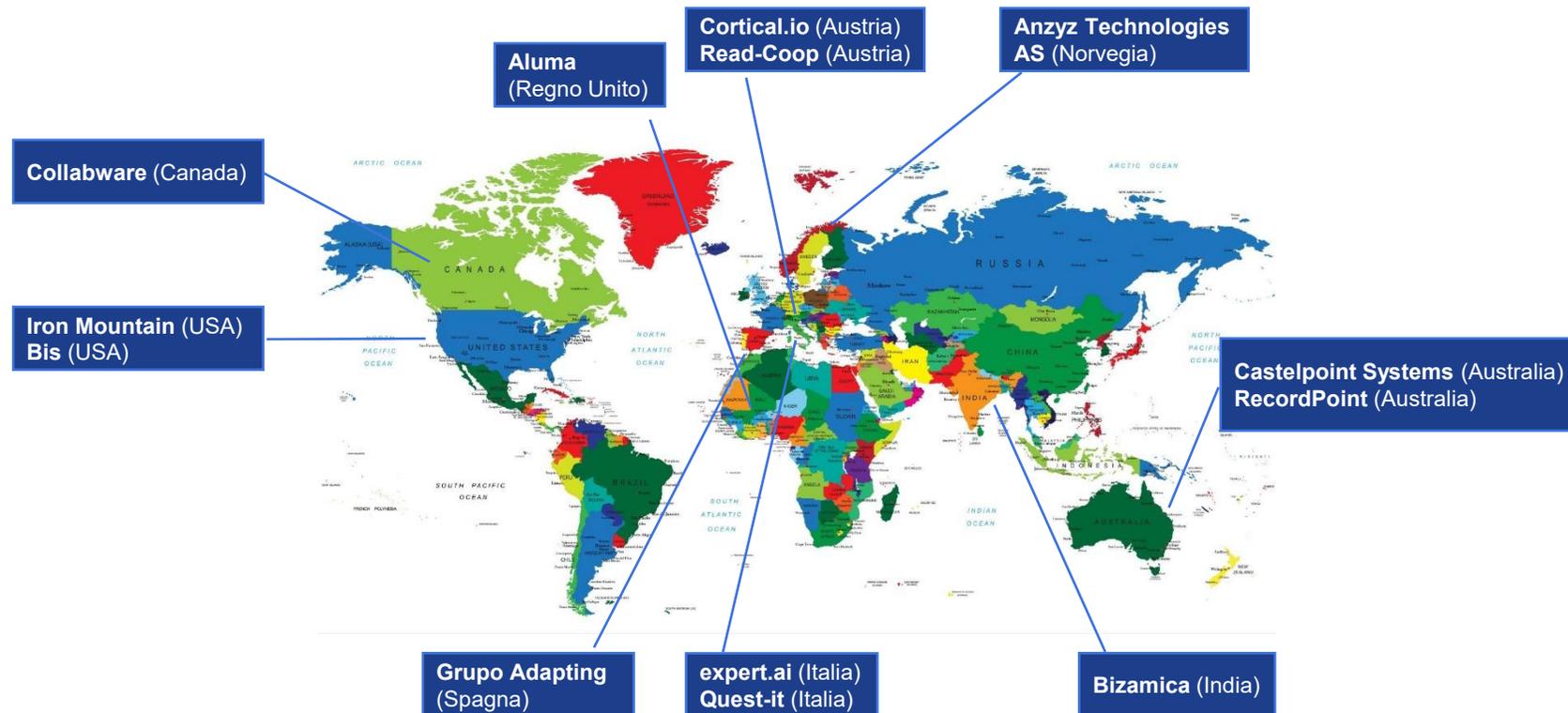
to  
489,

# Fase 2. Questionario e interviste



## Fase 2. Aziende che hanno risposto al sondaggio

Hanno risposto **13 aziende**:



# Il rapporto finale



<b>Study Title</b>	The role of AI in identifying or reconstituting archival aggregations of digital records and enriching metadata schemas
<b>Working group code</b>	Creation and Use: CU05
<b>Document type</b>	Final report
<b>Status</b>	Final version; Public
<b>Version</b>	10.0
<b>Writers</b>	Stefano Allegrezza, Mariella Guercio, Maria Mata Caravaca, Massimiliano Grandi, Bruna La Sorda
<b>Date</b>	November 1, 2023

1

[https://interparestrustai.org/assets/public/dissemination/Report-CU05-Survey-of-the-Companies\\_v121.pdf](https://interparestrustai.org/assets/public/dissemination/Report-CU05-Survey-of-the-Companies_v121.pdf)

**Stefano Allegrezza** (co-presidente) (Università di Bologna - Istituto di Ricerca per l'Intelligenza Artificiale Centrata sull'Uomo-ALMA AI)



**Mariella Guercio** (co-presidente) (Associazione nazionale archivistica italiana - ANAI)



**Maria Mata Caravaca** (Centro Internazionale per lo Studio della Conservazione e del Restauro dei Beni Culturali - ICCROM)



**Massimiliano Grandi** (Associazione nazionale archivistica italiana - ANAI)



**Bruna La Sorda** (Associazione nazionale archivistica italiana - ANAI)



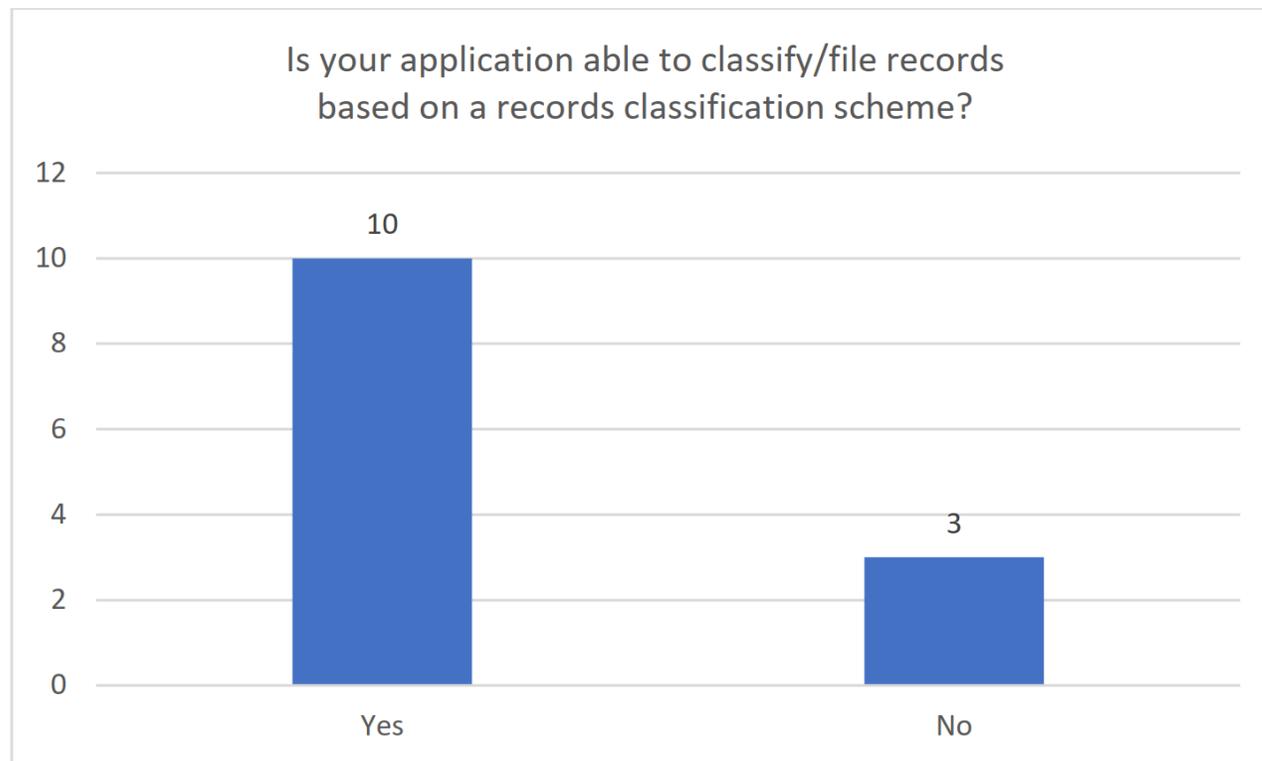
## **3. Risultati**

## Alcuni risultati: la classificazione

La **classificazione automatica dei documenti** sulla base di un piano di classificazione è offerta da **quasi tutte le aziende che hanno risposto ai questionari (10 su 13)**, comprese quelle non specificamente coinvolte nella gestione di archivi e documenti.

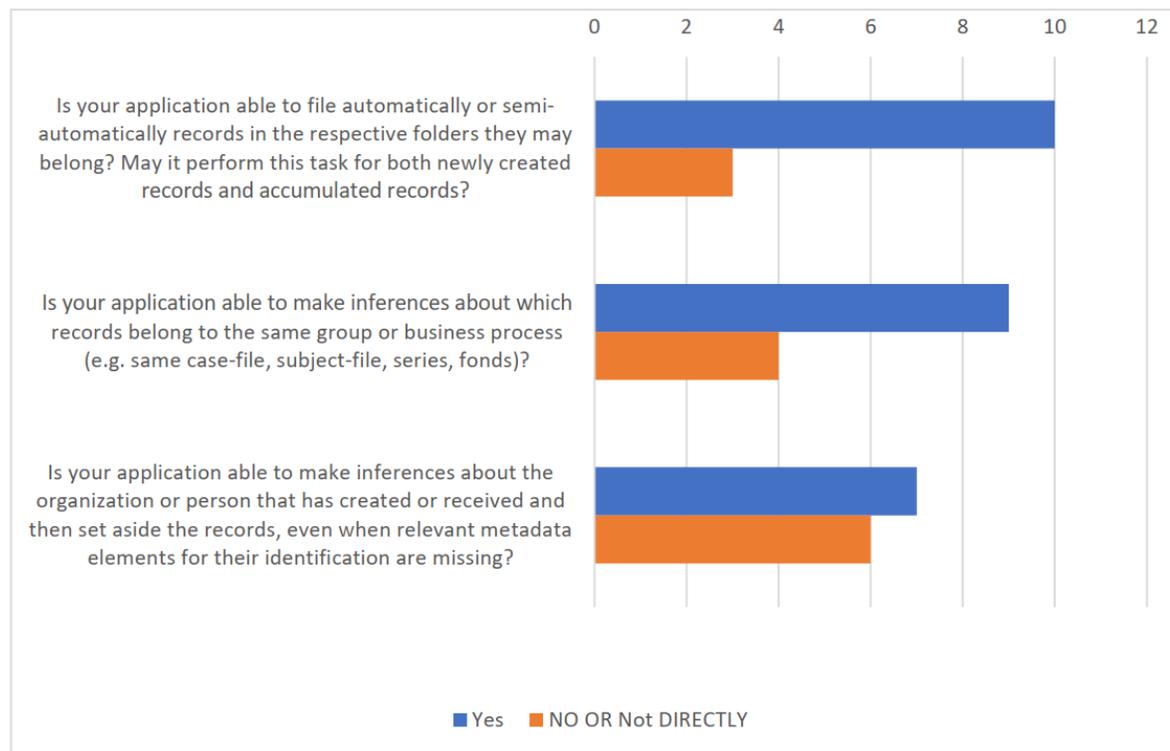
Tra le aziende le cui **applicazioni non prevedono la classificazione automatica**, una si occupa principalmente del riconoscimento del testo manoscritto e le altre due si concentrano sull'indicizzazione e sull'estrazione di elementi di metadati.

Alcune di queste affermano che le loro applicazioni **potrebbero essere addestrate** per la classificazione dei documenti.

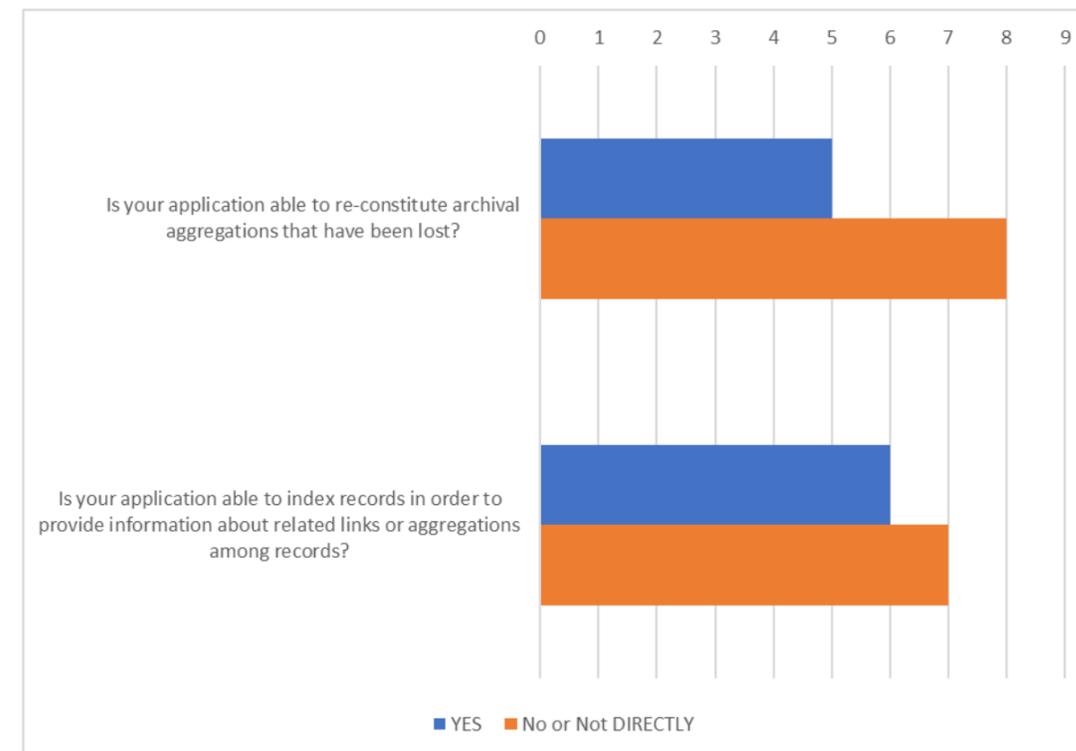


Classificazione

## Alcuni risultati: le aggregazioni documentali e il vincolo



Fascicolazione



Ricostituzione del vincolo archivistico

# Soluzioni tecnologiche

## Techniques and Analysis Models

The companies listed a wide range of different **analysis models**: **24 different entries** - the 5 most recurring are:

**Neural Network Models** (4 companies)

**Support Vector Machines** (4 companies)

**Decision Trees** (3 companies)

**Random Forests** (3 companies)

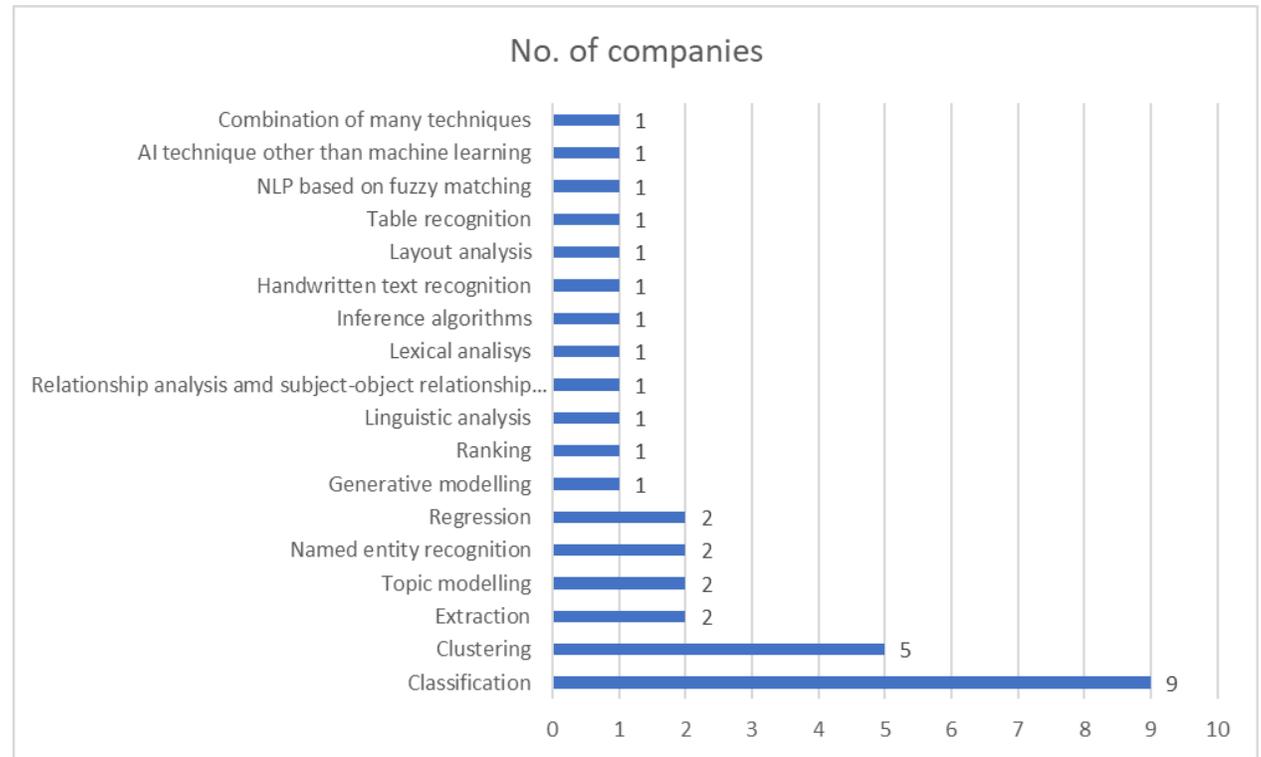
**LSTM - Long-Short Term Memory** (3 companies)

As to **the types of techniques** featured in the products of the companies - the 2 most recurring are:

**Classification** (9 companies)

**Clustering** (5 companies)

This is no surprising as the companies were selected because of their expertise at least in document management. However **18 different types of techniques** overall have been reported

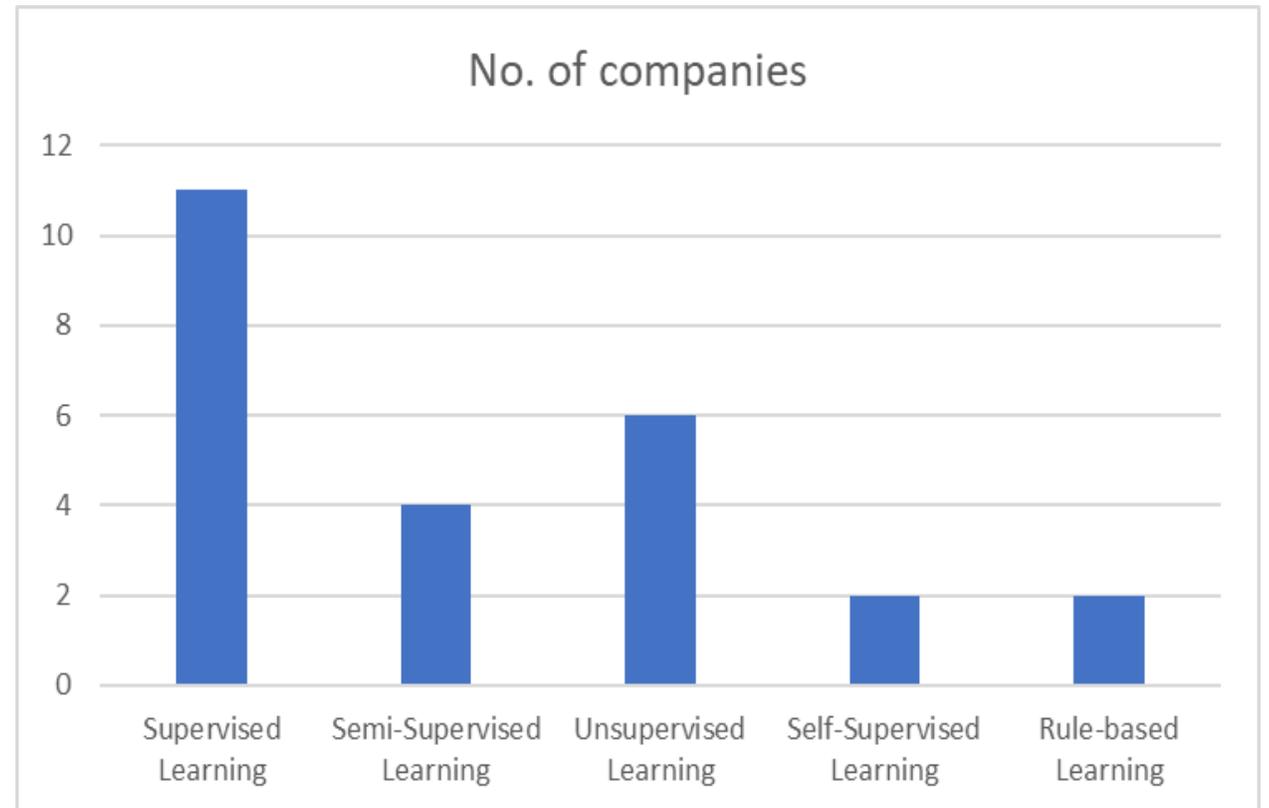


# Soluzioni tecnologiche

## Training strategies

The companies declared they use a combinations of different training strategies:

- **Supervised Learning:** 11 companies
- **Semi-Supervised Learning:** 4 companies
- **Unsupervised Learning:** 6 companies
- **Self-Supervised Learning :** 2 companies
- **Rule-based Learning :** 2 companies



## Conformità alle norme

---

Le aziende dichiarano che i loro prodotti sono conformi alle seguenti norme:

- **ISO 15489** (Records management);
- **ISO 16175** (Information and documentation — Processes and functional requirements for software for managing records);
- **ISO 23081-1:2017** (Information and documentation — Records management processes - Metadata for records);
- **ISO 30301:2019** (Information and documentation — Management systems for records — Requirements);
- **ISO/IEC 27001** (Information security management systems);
- **MoReq 2010** (Modular Requirements for records systems);

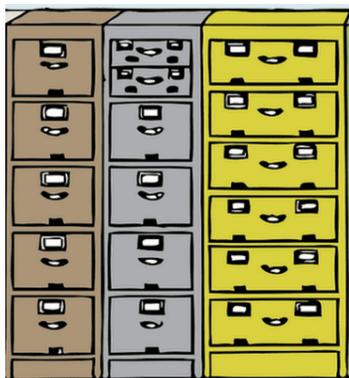
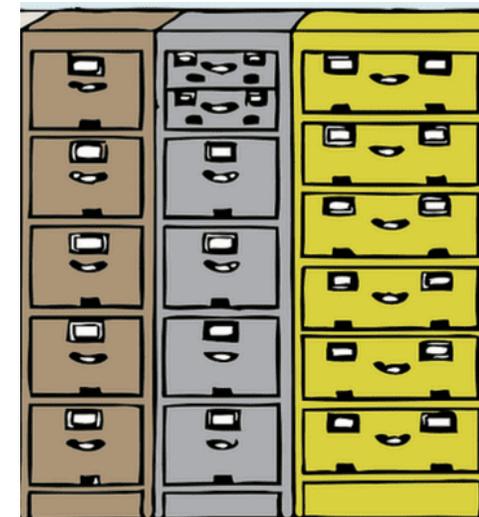
## Risultati dell'indagine dal punto di vista archivistico

Per quanto riguarda la **classificazione dei documenti**, tutte le aziende intervistate hanno dichiarato di:

- aver sviluppato soluzioni basate su tecnologie di intelligenza artificiale per l'**indicizzazione** e/o la **classificazione** di documenti/dati **strutturati, semi-strutturati e non strutturati**,

Il ruolo dei **metadati presenti o dedotti** è sempre al centro di tutte le risposte.

L'informazione sulla **tipologia dei documenti** - quando disponibile - è spesso considerata un'altra componente cruciale per il successo dell'applicazione delle tecniche di IA ai documenti.



Interessante notare che:

- **solo un'azienda** ha dichiarato che la sua piattaforma **può essere addestrata**, grazie a uno specifico set di dati, per generare autonomamente (senza intervento umano) etichette e tag relativi a qualsiasi **schema di classificazione**,
- in **tutti gli altri casi** **l'intermediazione umana è considerata insostituibile** se si vogliono ottenere risultati affidabili.

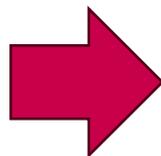
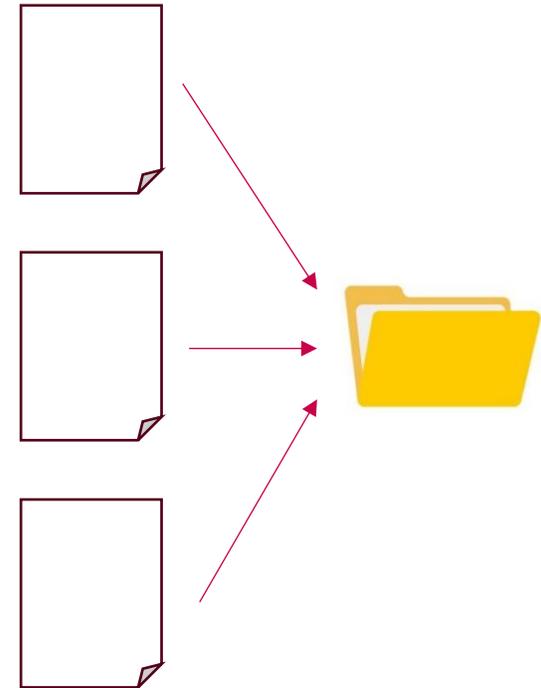
## Risultati dell'indagine dal punto di vista archivistico

Per quanto riguarda la **creazione delle aggregazioni documentali** (o la ricostituzione), le promesse di automatizzazione non sono molto incoraggianti, poiché questa possibilità è **limitata a casi molto specifici**, come ad esempio:

- quando esistono delle **specifiche ben definite sulle tipologie di documenti**
- quando vi sono informazioni sulla **struttura originale e ciò consente l'individuazione di relazioni funzionali tra i documenti**

L'aggregazione **automatica o semi-automatica** basata sul contenuto del documento è **solo suggerita** e di solito deve comunque essere **supportata dalla convalida da parte dell'utente** o da **regole disponibili al momento della creazione**.

Anche la **ri-costituzione del vincolo archivistico** - quando perduto o non esplicitamente definito - è riconosciuta come un'**attività complessa**, difficilmente realizzabile senza l'aiuto significativo da parte dell'uomo e/o informazioni descrittive.



Nella maggior parte dei casi le funzionalità **non sono già completamente sviluppate**, ma **in via di sviluppo** e ciò implica **maggiori investimenti (sostenuti dal mercato?)**.



## **4. Considerazioni finali**

## Considerazioni finali

---

La maggior parte delle aziende intervistate ha dimostrato:

- di **comprendere la complessità dell'ambiente e delle funzioni archivistiche**
- di essere consapevoli dell'**importanza dei metadati originali acquisiti nelle attività correnti del soggetto produttore**, sia nel caso in cui si tratti di classificare automaticamente i documenti sia nel caso si tratti di creare le aggregazioni documentali

L'indagine testimonia che la **complessità delle funzioni archivistiche** non può essere **facilmente ridotta** e ad un approccio automatico, ma solo **supportata** dalle tecnologie AI attraverso l'**intermediazione degli archivisti**.



## Future work

---



**Casi di  
studio**



NATO (North Atlantic  
Treaty Organization)



**ParER**

Polo archivistico dell'Emilia-Romagna

PAR-ER (Polo archivistico  
Regione Emilia Romagna)

## Considerazioni finali

---

- Dalla intelligenza artificiale, essendo basata su **modelli probabilistici** e non deterministici (come gli expert systems degli anni 80-90), non ci si può aspettare che classifichi i documenti **in maniera corretta al 100%** (esattamente come non ci si può aspettare che individui sempre correttamente l'immagine di un gatto tra le tante che le vengono sottoposte). Quindi gli errori (**hallucinations**) ci saranno sempre.
- Questo, però, è esattamente quello che avviene con gli esseri umani: anche gli **esseri umani** commettono degli errori. Ricordiamoci che l'uomo non è una macchina e quindi a volte sbaglia!
- Quello che ci si può ragionevolmente aspettare è che migliorando l'addestramento dell'AI (ad es. aumentando la dimensione dei corpora documentari) il numero di errori diminuisca sempre più e migliori sempre più **l'accuratezza**.
- Al momento la verifica umana sembra necessaria, ma sarà sempre così?
- Sì, fino a quando la precisione dell'intelligenza artificiale non sarà comparabile o addirittura superiore a quella dell'essere umano (e in taluni campi già è così) ed allora le operazioni di **classificazione** e di creazione delle **aggregazioni documentali** potranno (forse) essere svolte integralmente dall'intelligenza artificiale senza intervento umano.

Per saperne di più...



## Research Dissemination

Research Dissemination is authored by InterPARES Trust AI researchers and research assistants. They will be listed below on an ongoing basis.

Search

by Author:

by Text in Citation:

### Publications

- [ 9 ] [Books \(Including Chapters\)](#)
- [ 23 ] [Articles - Refereed](#)
- [ 16 ] [Conference Proceedings - Refereed](#)
- [ 26 ] [ITrustAI Research Documents](#)
- [ 2 ] [Social Media](#)

### Presentations

- [ 117 ] [Lectures, Workshops, and Seminars](#)
- [ 187 ] [Conferences \(Symposia, Sessions, Panels, Papers\)](#)

### Public Relations

- [ 9 ] [Articles and Reviews](#)
- [ 4 ] [Broadcast Interviews](#)
- [ 1 ] [Reports](#)
- [ 1 ] [Photos and Videos](#)

### Education

- [ 1 ] [Tutorials](#)
- [ 12 ] [Curriculum Materials](#)

[https://interparestrustai.org/trust/research\\_dissemination](https://interparestrustai.org/trust/research_dissemination)

**Grazie per l'attenzione!**

Stefano Allegrezza  
Università di Macerata (ITALIA), Dipartimento di Studi umanistici  
[stefano.allegrezza@unimc.it](mailto:stefano.allegrezza@unimc.it)