*Universal Bibliographic Control at the crossroads: the challenges of unifying IFLA bibliographic standards*
*18–19 August 2023*
*KBR, Brussels, Belgium*

# From Bias to Transparency : Ethical Imperatives in AI-Based Library Cataloging

**Hannes Lowagie**
Bibliographic Information Agency, KBR, Brussels, Belgium
E-mail address: Hannes.Lowagie@kbr.be

**Abstract:**

*In today's digital age, the need for a cohesive system to organize global knowledge is evident. Universal bibliographic control, long a scholarly interest, faces challenges due to diverse cultural perspectives. This article proposes an innovative approach, using AI and linked subject indexing, to create standardized cataloguing. The approach integrates diverse local classifications into a unified knowledge graph through AI, celebrating cultural diversity. Ethical considerations address bias and fair representation. This interconnected system benefits researchers, librarians, and users by promoting cross-cultural knowledge sharing, fostering collaboration, and enriching our understanding while preserving cultural uniqueness.*

**Introduction**

In today's digital age, the vast repository of knowledge available worldwide necessitates a cohesive and efficient system for organizing and accessing information. Universal bibliographic control, with its aim to create a standardized approach to cataloging, has long been a subject of scholarly interest. However, achieving this goal in a global context presents a formidable challenge due to the rich tapestry of cultural perspectives and diverse interpretations of concepts. To address this, we propose an innovative approach that capitalizes on the power of Artificial Intelligence (AI) and linked subject indexing.

In his 2022 article on 'Universal bibliographic control in the digital ecosystem', Mauro Guerrini contemplates in the last phrase the future possibilities brought about by alternative technologies such as machine learning and artificial intelligence (AI) in the realms of metadata and authority control. He highlights the potential usefulness of AI as a tool that complements the cataloguer's judgment, emphasizing its role in enhancing the fundamental intellectual activity in cataloging (Guerrini, 2022).

Today, only one year later, I am convinced there are indeed tools available to create comprehensive bibliographic records, with subject indexing as a crucial role part in this endeavor. Next week in Rotterdam, I will present a method I developed that uses AI to detect metadata on a page and creates the bibliographic record. One component of this flow involves text classification for subject indexing. We have trained a model using an Excel file consisting of two columns: the text itself and the corresponding "tags" to be added. Other automatic subject systems, like annif, need a same kind of training data with on the one hand the texts, on the other the subject terms.

## Automatic subject indexing using AI

Before I elaborate on the idea, first a small word about the possible use of AI in subject indexing. Since the development of Annif, a lot of libraries discovered the potential of AI. In KBR, we used a Microsoft category classification model that can also be trained to 'tag' text based on a self-trained model. Using NLP-techniques (national language processing), it is possible to train a AI-model that can, based on the training data, attribute terms to a new text.

In the beginning of training such a model, the results has to be monitored and the model has to be refined. Consider a training set for automated subject indexing consisting of three texts.

| | |
|---|---|
| The integration of AI and Internet of Things has transformed environmental monitoring. Through sensor networks and **data analytics**, AI can predict **natural disasters**, monitor air quality, and manage resources efficiently, ushering in a new era of smart **environmental** management. | Attributed:<br>Environmental Monitoring<br>Internet of Things |
| The widespread adoption of Internet of Things devices has raised concerns about **data privacy**. AI-driven analysis of personal data collected by these devices necessitates stringent regulations to protect user **privacy** and prevent unauthorized access. | Attributed:<br>Data Privacy<br>Internet of Things |
| In the realm of modern technology **data analytics** has paved the way for remarkable advancements. This synergy enables us to predict and mitigate the impact of **natural disasters** with unprecedented accuracy. By harnessing data from diverse sources, we gain crucial insights into the changing environment. However, this integration also raises concerns about **data privacy**. As we delve into the intricacies of **environmental** data collection, the safeguarding of personal data **privacy** becomes a paramount consideration. | Suggested:<br>Internet of Things |

This case highlights the potential pitfalls of automated subject indexing. The system may form associations based on keywords without comprehending the larger context, leading to inaccurate tags. In this scenario, the wrong suggestion of the "Internet of Things" tag can misrepresent the content, hinder search accuracy, and create confusion for users relying on the indexing system.

To enhance the accuracy of automated subject indexing and mitigate such errors, the following solutions can be implemented:

**Diverse Training Data**: Incorporate a more diverse range of training data that covers a wider variety of topics and contexts. By including a broader range of texts that are representative of

different subjects, themes, and regions, the system can learn to make more accurate associations and avoid overgeneralization.

**Fine-tuning and Customization**: Fine-tune the pre-trained NLP models on domain-specific data or tailor them to the specific needs and characteristics of the cataloging domain. This customization can help the models better capture the specific nuances and terminology relevant to the subject indexing process.

**Human-in-the-Loop Approach**: Incorporate human review and validation in the automated subject indexing process. Instead of solely relying on the output of the automated system, involve human experts who can review and verify the suggested tags, correcting any inaccuracies or false associations. This hybrid approach ensures a balance between automation and human expertise.

**Regular Evaluation and Monitoring**: Continuously evaluate and monitor the performance of the automated subject indexing system. Conduct periodic assessments to identify and rectify any biases, errors, or inconsistencies that may avaluation helps to refine the models, improve their accuracy, and adapt them to evolving cataloging needs.

**Transparency and Documentation**: Foster transparency in the subject indexing process by documenting the data sources, algorithms, and decision-making criteria used in the automated system. This documentation enables experts and stakeholders to better understand the reasoning behind the suggested tags and address any potential biases or inaccuracies.

By implementing these solutions, the accuracy and reliability of automated subject indexing can be improved, reducing the likelihood of suggesting wrong tags and ensuring a more accurate representation of the content's subject matter.

## Combining models for international collaboration

During the development of this model, I thought, is every country going to develop its own model and create its own training data, or can we collaborate here? And if so, how?

My first idea was to gather the expertise of cataloguers from different nations to build one big dataset and a comprehensive model capable of classifying texts in a detailed and accurate manner.

However, as I delved deeper into this concept, I soon realized its limitations. Although we all inhabit the same world, our interpretations of real-world objects and concepts differ significantly. Even basic concepts like 'poverty,' 'youth,' ('juvenile literature') can carry various meanings and interpretations across different cultural contexts.

Reaching universal agreement on how these concepts should be conceptualized remains elusive. Embracing these diverse interpretations is essential, as it fosters a vibrant and dynamic global community. It is crucial to avoid imposing a singular conceptualization as the 'true' one, as the value lies in acknowledging and appreciating the multiplicity of perspectives and interpretations.

In KBR, to train our model, we used data from data.bnf.fr, which provides book summaries related to specific Rameau subject terms through SPARQL queries from data.bnf. In that view,

it is important to take into account is the fact all the data used for training the model KBR created originates from the manual efforts of French cataloguers, resulting in an inherent bias with a French and European perspective. It is not a bad thing, but something to take into account.

Nonetheless, I remain convinced that the potential of AI can contribute to bolstering the concept of universal bibliographic control, particularly in the realm of universal classification. Here, the term 'universal' does not imply 'uniform' or 'identical' but rather signifies a commitment to honoring the diversity and distinctiveness of each culture and society, and thus for each interpretation of a concept.

That is also why, in the realm of subject indexing, it is crucial to recognize and embrace the diversity of cultural perspectives and interpretations that exist worldwide.
If we look at the IFLA National Bibliographic Register, we see that, while DDC is the most used, still 20% uses its own classification scheme, 60% own subject vocabulary.

While aiming for a more comprehensive and inclusive approach, we must appreciate the value of local classifications and subject vocabularies that capture the unique viewpoints of each culture. This recognition of diverse conceptualizations enhances the authenticity and accuracy of the subject indexing process, resulting in a classification system that truly reflects the global tapestry of knowledge. Each nation's library classification system reflects its unique cultural heritage, historical context, and intellectual traditions. Local classifications capture nuanced concepts that may be overlooked in broader, standardized schemes, thereby ensuring a more accurate and contextually relevant portrayal of the works. At the same time, we strive for an inclusive approach that transcends geographical boundaries and fosters collaboration among libraries worldwide. By embracing the value of both local and global contributions, we can create a unified system that celebrates the multiplicity of cultural viewpoints and promotes universal accessibility to knowledge.

**MAIN QUESTION**
But how can we achieve this, a universal way of organizing and accessing information with respect to each culture's different interpretations resulting in an improved discoverability of nation's publications?

I think the solution lies in the combination of AI-techniques with a creation of a knowledge graph that maps local classification and subject indexing methods.

> Side note: Today I use the terms "classification" and "subject indexing" interchangeably because our idea is using firstly a very general classification, such as 'history', 'social sciences', 'applied sciences', and then adding subject terms linked to FAST or Rameau, for instance, for the category 'History': 'Middle Ages,' 'Vikings,' and 'Codicolgy.' This allows us to combine both approaches and utilize subject terms for classifying publications. We believe this approach is the best option for two reasons. First, it enables us to use automatic subject indexing, starting with the detection of the general category and then utilizing a model trained with terms related to that category. Second, it provides great flexibility, allowing us to easily add subject terms without the need to modify the thesaurus. This flexibility enables us to respond effectively to new subjects found in Belgian books.

While it remains highly theoretical and perhaps even challenging, impossible, but despite its complexities, I believe it can be interesting to present it, and its worth discussing it at the round table. The idea is to create a unified and interconnected knowledge graph by combining and mapping different local classifications through SKOS links. Rather than striving for homogeneity, this model embraces the diversity of cultural perspectives and interpretations of concepts. The links between the different concepts of a classification can be created manually, but they can also be created (automatically?) by comparing the classifications given to the same text by every nation.

For example, a book about 17 years old: if Country A gives the classification "Childhood" , a link 'is related to' between "Childhood" and "young adulthood' from country B And Childhood and Adolescence from country B
By creating such a knowledge graphs that maps different local and global classifications, we acknowledge and celebrate the richness of various knowledge systems.

Today I will not talk too much about the technical aspects, because I am not a specialist, but the ideai is to create a dataset with texts for each classification, that can be used to train a AI subject indexing tool (such as Annif). This can be done by importing summaries and texts from different nation's publications outputs, together with the attributed classifications, into a common file. To begin this process, libraries and information institutions can collaborate to gather summaries and full-text materials from their respective national publications. These materials can include books, articles, research papers, and other written works. A first step for this endeavor, can be found in the output of an ongoing research project at KBR: BELTRANS. This project focuses on comparing various translations of Belgian books, and our data manager has diligently exported and analyzed the same publications from different catalogues, including KBR, Bnf, BL, and KBNL. The advantage of these translated books is that we can compare classifications and subject indexing from the same works at different institutions.
Leveraging the classification information assigned to these books offers a promising starting point for our larger dataset. By extending this method to encompass an even broader range of publications, we can effectively establish meaningful relationships between the diverse classification schemes employed by different libraries. This groundwork is instrumental in creating a unified and interconnected classification system, enhancing universal bibliographic control for a wider array of works.
*This process harmonizes data from different sources, facilitating the creation of a richer and more diverse dataset. By combining knowledge from various nations, libraries can enrich their own subject classifications, while also contributing to the larger global endeavor of achieving universal bibliographic control. The benefits of this approach extend beyond the enhancement of individual library collections, enabling greater cross-cultural information retrieval and knowledge dissemination worldwide.*

Later, the imported texts can then be enriched, through the established mapping, with each countries own classification term. If needed, this attribution (mapping) can be refined. That way we create a dataset that can be used to train its own AI-model, with your own classification, but that not only includes your own nation's publications, but that is enriched with global publications. One of the key advantages of importing summaries and texts from diverse sources into a common file is the creation of a richer dataset. By pooling together information from various national publications, the classification system gains access to a wide array of subject domains, cultural perspectives, and linguistic variations. This diversity allows for a more inclusive representation of global knowledge and cultural heritage. Additionally, the dataset becomes more extensive, covering topics that may be underrepresented in individual national

collections. As a result, researchers and users worldwide can benefit from a broader scope of information, enabling them to explore interdisciplinary connections and gain a deeper understanding of different subjects.

To summarize, a possible workflow can be:
- **Initial Data Import**: As libraries begin to import their summaries and classifications into the common file, an initial dataset is formed with each library's subject terms.
- **Identifying Common Concepts**: Libraries can then analyze their subject terms and identify common concepts shared across different classifications. For example, if one library classifies a book under "Youth Literature," while another uses "Children's Literature," they can identify these as representing the same concept of literature for young readers.
- **Mapping Creation**: Based on the identified common concepts, libraries can collaboratively create a mapping table or dictionary that links subject terms from different classifications. This mapping should indicate the relationships between terms, such as 'same as,' 'is related to,' 'narrower,' or 'broader.'. This mapping should identify which tags or terms from one scheme correspond to the same or similar concepts in another scheme. For example, if one library uses the term "adolescents" to represent youth and another library uses "young people," establish a mapping that connects the two concepts. These concepts should also be mapped with LOD sources as wikidata or FAST subject terms. Based on these relationships a Knowledge graphs can be created. Knowledge graphs are powerful tools that represent information as nodes and edges, forming a network of interconnected concepts. By constructing a knowledge graph based on subject terms and their relationships, libraries can establish links between different classifications and navigate through related concepts.
- **Incorporating AI Techniques**: AI technologies, such as natural language processing and machine learning, can be employed to assist in identifying similar subject terms and suggesting potential mappings. These techniques can help automate the process and ensure a more comprehensive linkage. AI techniques play a crucial role in the implementation of the proposed model for enhancing universal bibliographic control through linked subject indexing. On technique is that of word embeddings. Word embeddings are dense vector representations of words in a continuous vector space. These embeddings capture semantic relationships between words, enabling AI models to recognize semantically similar terms even if they are expressed differently. By using word embeddings, libraries can compare subject terms from different classifications and find potential matches based on semantic similarity. By clustering terms based on their semantic similarities, libraries can identify potential mappings and linkages between clusters representing similar concepts.
- **Sharing Mapping Information**: The mapping table should be shared with all participating libraries to ensure consistency and uniformity. Each library can then update its classification system to include the mappings for cross-referencing.
- **Continuous Refinement**: As more libraries contribute their data to the common file, the mappings can be continuously refined and expanded to accommodate new subject terms and concepts.

By adopting this approach, libraries can collaboratively create a network of interconnected subject terms, enhancing the discoverability and accessibility of information across different classifications. Over time, the establishment of mappings and linkages will enable libraries to achieve a more unified and inclusive classification system that reflects a broader, global perspective.

## Challenges and limitations

The implementation of the proposed model for enhancing universal bibliographic control through linked classification and subject indexing comes with several challenges and limitations that must be carefully addressed. One significant challenge is the potential existence of conflicts in interpretations between different classification schemes. As libraries from various countries contribute their own classifications based on their unique cultural perspectives, variations in terminology and conceptualizations may arise. By adding different classifications for the same text, links between the classification schemes can be detected and different interpretations visualised. I will give an example later.

Another challenge lies in the ambiguities. Some subject terms may have multiple meanings or interpretations, leading to uncertainty in their alignment. Strategies for resolving these ambiguities involve thorough analysis, contextual understanding, and consultation with subject matter experts. Establishing a transparent and participatory process for resolving ambiguities can help build trust and ensure the integrity of the interconnected classification system.

Additionally, the scalability of the model is a crucial consideration. As the dataset grows with contributions from numerous libraries, managing and updating the mappings between subject terms can become complex. Developing efficient algorithms and technologies for continuous data integration and mapping updates is essential to maintain the accuracy and relevance of the unified system.

To address these challenges, international collaborations and partnerships are essential. Libraries, researchers, and stakeholders must work together to create standardized guidelines, establish a shared ontology, and implement transparent processes for validating and updating mappings. By addressing these challenges and limitations proactively, the proposed mapping can become a powerful tool for promoting universal bibliographic control and facilitating cross-cultural knowledge sharing in the global information landscape.

## Advantages

The process of importing data from multiple sources encourages collaboration and knowledge exchange among libraries and institutions. It fosters a sense of shared responsibility for maintaining and enriching the classification system, promoting a collective effort towards universal bibliographic control. Moreover, the collaborative nature of this approach allows for continuous updates and refinements, reflecting the ever-evolving landscape of knowledge production and cultural expressions. With each contribution, the classification system becomes more robust, accommodating new perspectives, and refining its ability to categorize and link subject terms effectively.

The interconnected classification system also offers a multitude of benefits for researchers, librarians, and users worldwide.

For **researchers**, this system provides an extensive and interconnected knowledge base that transcends borders and language barriers. With access to a diverse range of subject terms and linked concepts, researchers can explore interdisciplinary connections, gain deeper insights into various topics, and make informed cross-cultural comparisons. The unified classification system enhances the efficiency and accuracy of information retrieval, allowing researchers to discover relevant resources more effectively and conduct comprehensive and nuanced studies.

**Librarians** also reap substantial advantages from the interconnected system. By incorporating multiple local classifications into a unified framework, librarians can streamline cataloging processes, minimize duplication efforts, and maintain consistent and standardized metadata across international collections. This system facilitates seamless cooperation among libraries, fostering an environment of collaboration and knowledge exchange. Librarians can better serve their users by offering broader and more diverse collections, enriching the library experience and supporting users in their pursuit of knowledge. By pooling their expertise and contributing local classifications, libraries become part of a global knowledge-sharing network. This collaborative approach not only enriches the dataset but also fosters a sense of shared responsibility in preserving and promoting cultural heritage and knowledge from various regions. As a result, the unified system becomes a powerful tool for building bridges of understanding across different cultures and promoting a more inclusive and interconnected global information landscape.

**Users** worldwide also benefit greatly from the interconnected classification system. The unified approach ensures that users can access a wider array of resources beyond their local library, enriching their understanding of diverse cultures, ideas, and perspectives. With linked subject terms and interconnected concepts, users can navigate through the global knowledge network, discovering related materials and gaining a deeper appreciation of the interconnectedness of knowledge. Ultimately, the interconnected classification system empowers researchers, librarians, and users alike, fostering a more inclusive, collaborative, and accessible global information environment.

One of the most significant advantages of the unified system is its potential to foster cross-cultural information retrieval and knowledge sharing. With the establishment of links between subject terms from different classifications, users can navigate through the interconnected web of concepts, uncovering related materials from different libraries and cultural contexts. This cross-cultural exchange promotes a deeper understanding of diverse viewpoints and enables users to access a wide range of resources that would otherwise remain isolated within individual library catalogs.

With subject terms linked across various classifications, users from different cultural backgrounds can effortlessly navigate and explore information in a manner that aligns with their unique contexts and preferences. This system promotes a more inclusive and equitable knowledge-sharing environment, where users can access materials relevant to their interests, regardless of the originating library's classification.

By breaking down language barriers through linked subject terms, the unified system promotes efficient cross-cultural communication and understanding. Users seeking information in their native language can easily find resources classified in other languages, fostering global knowledge exchange and collaboration. This interconnected approach empowers users to delve into diverse cultural contexts and gain insights from various viewpoints, enriching their intellectual exploration and encouraging intercultural dialogue.

Overall, the unified system revolutionizes cross-cultural information retrieval and knowledge sharing by transcending linguistic and cultural barriers, facilitating meaningful interactions among diverse perspectives, and promoting a global exchange of ideas. It fosters an inclusive, accessible, and interconnected information environment that empowers users to embrace the richness and diversity of global knowledge.

**Ethical considerations**
Because today we talk about the ethics of cataloging, I also want to talk about some ethical considerations of the idea I just presented.

I see some benefits:
- **Cultural sensitivity** : Respect for diverse cultural perspectives is essential. When integrating local classifications into a global framework, it is crucial to avoid cultural appropriation or misrepresentation. Ensuring that local communities have a say in how their knowledge is represented and used can mitigate potential issues.
- **Diversity and Inclusion**: Linking classification schemes ensures that the training set considers diverse perspectives, cultures, and languages, providing equitable access to information for all users. Avoid biases by incorporating a wide range of materials and representing different subject areas and regions.
- **Bias and Fair Representation**: Local classifications may inadvertently contain biases based on historical, social, cultural or political factors. When merging them into a global system, it is vital to critically examine and address any biases to ensure fair and equitable representation of all perspectives.

I will elaborate on this using a well-known example in our regions, namely "Zwarte Piet," translated as "Black Pete,". He is a traditional character associated with the Dutch holiday of Sinterklaas (St. Nicholas). Everyone from my age grew up with this black-faced helper, with big red lips, curly hair. "Zwarte Piet" has long been a beloved and integral part of the Sinterklaas tradition, representing a figure who assists Sinterklaas in delivering gifts to children. Why the black color: one of the explanations was that he delivered gifts through the chimney which made him a black color

However, in recent years this character has been a subject of controversy due to its depiction, which we (or many people) nowadays perceive as perpetuating racial stereotypes and blackface imagery. The character's portrayal has raised concerns about racial insensitivity and cultural appropriation.

So on one hand, the Dutch classification may view "Zwarte Piet" as a cherished part of their heritage and folklore, emphasizing the character's positive role as a jovial helper during the holiday season. On the other hand, from a global perspective, using a subject term "Zwarte Piet" could be seen as perpetuating racial stereotypes and promoting cultural insensitivity. It may conflict with broader principles of inclusivity and respect for diverse cultural perspectives, particularly those that are critical of the character's representation. Addressing this ethical concern requires engaging in open dialogues with diverse communities to understand their perspectives and concerns regarding "Zwarte Piet." Through inclusive consultations and comparisons, librarians can ensure that the global classification system thoughtfully reflects the complex and nuanced views on this issue. If a cataloguer attributes a possible harmfull concept of "Zwarte Piet", it can be 'flagged' by other communities, raising possible issues with that term. In the meantime it can also give context to the user (who still will be using search terms als 'zwarte piet'). Ultimately, the ethical approach to handling "Zwarte Piet" in the classification system lies in striking a balance between recognizing the character's cultural significance within the Dutch context and being attentive to the concerns raised by 'other' individuals and communities who find the depiction offensive or hurtful. By doing so, the global classification system can uphold principles of fairness, respect, and cultural sensitivity in its representation of this contentious cultural element. Linking classification schemes will speed up the detection and possible cleanup of possible harmful terms. It can also be reassuring to a black man living here who finds the depiction of Black Pete problematic that he can see that the "local" term "Black Pete" is flagged by other communities.

**Conclusion**

The proposed approach for enhancing universal bibliographic control revolves around the integration of linked subject indexing and AI classification. The proposed approach aims to create a unified and interconnected knowledge graph by combining different local classifications through SKOS links. Rather than striving for homogeneity, this model embraces the diversity of cultural perspectives and interpretations of concepts. By importing summaries and texts into a common file and allowing each library to contribute its own classification, we acknowledge and celebrate the richness of various knowledge systems. The establishment of links between subject terms, such as 'same as,' 'is related to,' 'narrower,' and 'broader,' fosters a harmonious network of interconnected concepts. AI plays a pivotal role in recognizing patterns and relationships across datasets, thereby facilitating the creation of a comprehensive and enriched knowledge graph.

However, it is crucial to address potential ethical considerations, such as bias and fair representation. Careful examination and resolution of biases in local classifications ensure that the interconnected graph respects the uniqueness of each culture and avoids imposing a singular worldview.

The advantages of this interconnected system are manifold. Researchers, librarians, and users worldwide benefit from a vast and diverse pool of knowledge resources, promoting cross-cultural information retrieval and knowledge sharing. This model fosters global collaboration and enriches our understanding by embracing different perspectives. As we continue to refine and expand the interconnected knowledge graph, we move closer to achieving universal bibliographic control while preserving and celebrating the cultural diversity that makes our world so vibrant and unique.

*--In this article about the use of AI, I have to admit that OpenAI's chatbot ChatGPT played a role in improving the clarity and coherence of the content by rephrasing and polishing my initial draft.--*