**ETD-MS v2.0: A Proposed Extended Standard for Metadata of Electronic Theses and Dissertations**

Lamia Salsabil[1], Jian Wu[1], William A. Ingram[2], Edward Fox[3]

[1] Computer Science, Old Dominion University

[2] University Libraries, Virginia Polytechnic Institute and State University

[3] Computer Science, Virginia Polytechnic Institute and State University

**Abstract**

The growth of Electronic Theses and Dissertations (ETDs) in academic repositories requires comprehensive and robust schemas for compliance with the FAIR (Findability, Accessibility, Interoperability, and Reusability) principles. Dublin Core and ETD-MS v1.1 were established as the metadata standards for general scholarly documents and ETDs. We identified several gaps between the existing schemas and the need to represent ETDs comprehensively toward a better digital service. The content-level data, including objects comprising ETDs, become increasingly crucial to facilitate the development of machine learning models to mine scientific knowledge from ETDs, and scholarly big data services in general. By organizing content-level data into a new schema, this paper addresses the critical need for enhancing the expressiveness and depth of metadata for ETDs. The schema proposed includes a Core Component building on the existing ETD-MS v1.1 schema, and an Extended Component that captures objects, their provenance, and user interactions for ETDs. The schema covers 28 entities with a total of 160 metadata fields. To demonstrate applicability, we implemented the schema using MySQL and populated it with data derived from 1,000 ETDs collected from U.S. university libraries. This work provides a comprehensive and flexible approach that addresses the limitations of existing standards by enabling the description of content-level data, laying the groundwork for integrating advanced AI techniques into academic repositories.

*Keywords:* ETD, metadata, schema, FAIR

**Introduction**

Digital repositories dedicated to Electronic Theses and Dissertations (ETDs) have experienced tremendous growth (Uddin et al. 2021), creating a vast and valuable resource for researchers. To maximize the impact of this resource and provide more and better services, it is essential to adhere to the FAIR (Findable, Accessible, Interoperable, and Reusable) principles. Existing ETD metadata schemas, including Dublin Core (DC) and ETD-MS v1.1 (Hickey et al., 2021), have been used in building digital

repositories, search engines, and application programming interfaces (APIs). For example, DC and ETD-MS v1.1 are both used to build Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH) (Lagoze et al., 2003) APIs for institutional repositories, which are capable of supporting multiple metadata formats, including both the mandatory DC and specialized ETD Metadata Object Description Schema (MODS) (Gartner, 2003). These OAI-PMH APIs allow repositories to meet both general interoperability requirements and specific informational needs. However, existing ETD schemas are insufficient to fulfill all aspects of the FAIR principles because of three major gaps. The first is that existing ETD schemas do not provide complete and detailed representations of ETDs. For example, the "dc.rights" element in ETD-MS v1.1 (Hickey et al., 2021) only allows three values regarding accessibility but does not support the rich semantics of rights information. Another example is that ETD-MS v1.1 contains "dc.format" that describes standard MIME types available, which assumes an ETD has a single format. Certain ETDs include multiple files with various MIME types. The second gap is that existing ETD schemas lack metadata elements that describe parts of an ETD, such as chapters, figures, and tables. In the following sections, we will call them objects because they each have attributes. The third gap is the lack of metadata elements that have been added from sources other than those responsible for handling the ETD submission. Such sources might be human (e.g., authors, advisors, library catalogers, or ProQuest catalogers). All such metadata should be accompanied by provenance information that explains how the metadata arose. For example, if we generate a chapter summary, the metadata needs to point to the process or model and other details that led to such a summary, e.g., Ingram et al., (2024).

This paper proposes an extended metadata standard to overcome these gaps and enhance the findability and reusability of ETDs. The new schema will extend the existing metadata standard by incorporating new metadata elements to provide complete and detailed descriptions of the ETD document, the objects comprising it, their provenance, and relationships between ETDs and between objects to enhance the *findability and reusability* of ETD at document and content levels. In this work,

we introduce the concept of "content-level metadata," which refers to metadata that describes the parts of an ETD, including various types of objects, their attributes and provenance, and relationships between objects. Due to space constraints, figures and tables are accessible online by URL identifiers in footnotes.

**Related Work**

Metadata is essential to the structure and functionality of digital library systems, particularly for ETDs. Managing ETD metadata presents ongoing challenges due to advances in technology and diverse practices. Effective metadata management is essential to comply with the FAIR principles, as demonstrated by Santos et al. (2023) with the FAIR Data Point system.

The DC standard is widely recognized for its simplicity and role in ensuring basic interoperability across digital collections. Comprising 15 elements, the Simple DC provides a foundational framework for metadata management. Its core principles—simplicity, standardization, and interoperability—make it suitable for broad application. However, DC's general-purpose nature limits its capability to address the specific and detailed needs of different digital resources (Park et al., 2015). To overcome these limitations, certain repositories extended the Simple DC through qualifiers, adapting it to represent metadata covering scholarly communication, department affiliations, and grant details. For instance, Ho et al. (2019) developed guidelines to ensure comprehensive and consistent metadata in their repository by aligning it with both DC and MODS schemas. This approach balances the simplicity of DC with the depth of MODS, by specifying mandatory, recommended, and optional metadata elements.

ETD-MS v1.1, which includes 22 metadata elements, was designed to meet the unique needs of ETDs. This schema extends the DC to provide detailed descriptions specific to ETDs. Significant efforts have been made to update and enhance ETD metadata in institutional repositories. To ensure consistent and compatible metadata, the Texas Digital Library (TDL) developed a standardized ETD MODS schema, which was subsequently simplified into DC format for compatibility with DSpace (Potvin et al., 2015). Thompson et al. (2019) focused on improving ETD metadata in the University of Houston Libraries'

repositories to align with the TDL Descriptive Metadata Guidelines for ETDs, Version 2.0. Repositories using the OAI-PMH protocol can support various metadata standards, including both the Simple DC and specialized formats such as TDL's ETD MODS. This approach allows repositories to meet general interoperability requirements while accommodating specific informational needs. The OAI-PMH specification mandates that repositories provide records in the Simple DC format to ensure basic interoperability. In contrast, the TDL's ETD schema offers a more detailed metadata format designed for comprehensive ETD descriptions, including the fields for Author, Thesis Advisor, and Committee Member information, as well as unique elements such as Author Identifier (e.g., ORCID) and Embargo. However, all existing ETD metadata standards fall short in representing content-level metadata and or do not provide a comprehensive description of ETDs, such as missing detailed rights information.

The Open Annotation Core (OAC) Data Model was an early effort to incorporate content, associations between related resources, and provenance for web annotations (Hunter et al., 2010). Open Annotations can easily be shared between platforms with sufficient richness of expression to satisfy complex requirements. Although the full model supports additional functionality, such as enabling semantic annotations and embedding content, the OAC data model was designed for web resources in general. A concrete schema for ETDs still needs to be developed.

ETD content, including objects comprising a hierarchical structure, contains rich information about domain knowledge, scientific findings, technical details, and author profile. Recent efforts have been made to mine the content of ETDs using state-of-the-art natural language processing and computer vision tools. For example, Choudhury et al. (2024) introduced multimodal models for classifying ETD pages. Kahu et al. (2021) proposed YOLO-based methods for detecting and extracting figures and tables. Topic modeling, using methods such as Latent Dirichlet Allocation (LDA), was used to identify core themes of ETDs (Lamba et al., 2019; Aboelnaga et al., 2024). Sequence-to-sequence models and Pointer-Generator networks were employed to generate ETD summaries (Ingram et al., 2019; Ahuja et

al., 2018). Large Language Models (LLMs) and fine-tuned pre-trained models were used to generate chapter-level classification labels and summaries (Bipasha, 2024). Concept maps were used to visualize ETD content (Richardson et al., 2008). Object detection techniques, such as ETD-OD (Ahuja et al., 2022), identify and categorize non-textual elements in ETDs. Automatic classification techniques were also applied to chapters (Banerjee, 2024). ETD objects need metadata to describe their attributes and provenance, which will be addressed by our proposed schema.

## Methodology

Our proposed metadata standard, called ETD-MS v2.0 contains a Core Component and an Extended Component.  The Core Component aims to address the need for a *comprehensive* and *detailed* description of document-level data of ETDs by introducing new entities and attributes. The Extended Component was aimed at addressing the needs of describing complex ETD content.

### Core Component Development

We followed a top-down approach for the Core Component development. We analyzed 500 ETDs sampled from a collection of over 500,000 ETDs downloaded from 114 U.S. university libraries (Uddin et al., 2021). The 500 ETDs cover 350 STEM and 150 non-STEM majors, including 469 doctoral, 27 master's, and 5 bachelor's degrees, and were published between 1945 and 1990. Each ETD was accompanied by an XML file containing metadata obtained from the OAI-PMH API or by scraping the HTML landing page. We applied three types of modification to build the Core Component. First, we started with high-level attributes of ETDs in ETD-MS v1.1 and converted certain attributes into entities if they could be described by the second-level attributes. For example, we converted the dc.rights into an entity called "Rights", which included attributes describing the full text of the copyright statement, the rights type, and the date when the rights was applied. Second, we supplemented attributes or entities that were not covered by ETD-MS v1.1. For example, ETD-MS v1.1 did not include elements to describe references, so we added the "References" entity describing a reference cited by an ETD. The attributes of "References"

include reference_text, author, title, year, and venue. Finally, we observed that certain attributes in ETD-MS v1.1 are only applicable for special ETDS, so we created a new entity that provides a generalized description that is applicable to all ETDs we inspected. For example, the ETD-MS v1.1 contains "dc.format", which assumes an ETD has a single format. However, certain ETDs were published in multiple PDFs and multimedia files. We thus designed the "ETD_File" entity with attributes including the file description, generation method (scanned or born-digital), checksum, and MIME types.

**Core Component Description**

The Core Component consists of 10 entities with 73 metadata fields. The "ETDs" entity includes key metadata describing an ETD such as a unique identifier, title, author, committee chair, year of publication, and abstract. The other entities can be classified into two categories.

Category C.1: describing the ETD document itself, including "Rights", "ETD_file", "Subjects", "ETD_classes", and "ETD_topics".

Category C.2: describing the ETD's relationship with other documents, such as "References", "ETD-ETD_neighbors", "Collections", and "Collection_topics". For example, if a collection consisting of 2,000 scanned ETDs was built for a text summarization project. The "Collections" entity will contain the collection name, a text description, and a list of ETD IDs belonging to this collection. "Collection_topics" represents topics identified for each collection. For example, a collection of scanned ETDs may have topics "Neural Networks" and "LLM," which are then stored in "Collection_topics" and linked to individual ETDs in "ETD_topics." If this collection is classified as "Computer Science" according to ProQuest, the "Computer Science" class should be recorded in the "ETD_classes" entity.

Table 1[1] provides the details of the Core Component. Compared with ETD-MS v1.1 and DC, the hierarchical structure of the Core Component provides a deeper description of ETDs with attributed entities instead of flat attributes.

---

[1] https://github.com/lamps-lab/ETDMiner/blob/master/ETD-MS-v2.0/ETD-MS-v2.0-Schema.pdf

**Extended Component Development**

The Extended Component was developed using a bootstrap approach. The development started with identifying the objects (e.g., chapters, figures, tables) and their attributes, such as text, object_classes, and object summaries. Then we added entities to describe object provenance. For example, the "Classifications" entity represents a classification system (e.g., ProQuest Subject Categories) that provides a list of subject categories to categorize ETDs and their objects. The "Classification_entries" entity represents specific subject categories, such as "Library Science" and "Web Studies", and the "Classifiers" entity represents models to classify an ETD or its objects relative to a Classification. Next, we added an entity called object-object_neighbors to describe the object-object relations. Finally, we added several entities to describe user interactions with ETDs, such as "Users", "User_classes", and "User-user_neighbors."

**Extended Component Description**

The Extended Component includes interconnected entities that represent objects within ETDs, their provenance, relations, and user interactions. The Extended Component comprises 18 entities, with a total of 87 attributes. The four categories are shown below.

Category E.1: describing objects including "Object," "Text," "Object_classes," "Object_summaries," "Object metadata," and "Object_topics." The "Objects" entity represents objects detected (e.g., figures and tables) or generated (e.g., summaries) from ETDs, while "Object_classes" represents a list of classes describing the text objects of an ETD. For example, a table in an ETD that displays statistics of employee satisfaction surveys is categorized as "Statistics." This category is recorded in the "Object_classes" entity for the corresponding "table" object. "Object_summaries" represents summaries of objects, and "Object_topics" represents topics associated with text objects. "Object-object_neighbors" represents similar objects within an ETD.

Category E.2: describing provenance including "Classifications," "Classification_entries," "Classifiers", "Topic_models," and "Summarizers." Note that these provenance entities can be used for objects in the Extended Component and ETDs in the Core Component. For example, the "ETD_classes" entity represents the assignment of subject categories from an instance in the "Classifications" entity. In addition, the "Topic_models" is an entity that represents details about topic modelings, and "Summarizers" is an entity that represents details about various summarization models.

Category E.3: describing user interactions including "Users," "User_classes," "User_queries," "User_queries_clicks," "User_topics," and "User-user_neighbors." "User_classes" represents a user's subject category such as "Information Science." "User_queries" represents queries made by individual users in a search engine, and "User_topics" represents topics of interest to a user.

Table 2[1] provides a detailed overview of the Extended Component. Figure 1[2] illustrates the entities and their relations in the Core and Extended Components.

## Implementation

**Proof of Concept Database Implementation**

We implemented the proposed schema using a MySQL database to demonstrate its applicability. The process involved extracting document-level and content-level data and metadata from a subset of 1,000 ETDs selected from our collection of over 500,000 ETDs, with no overlaps with the 500 ETDs used for developing the schema. The 1,000 ETDs were selected from Year 2005 to 2019, sourced from 50 distinct U.S. universities. The database consists of 28 tables each mapping to an entity in the schema. From the raw data collected by parsing XML files obtained from OAI-PMH or HTML files crawled from university library sitemaps, we extracted 17 metadata fields, such as title, year, author, and advisor. The HTML files provide additional metadata such as copyright details. We inserted NULL values when the metadata was unavailable. To derive data using AI-based methods, we employed several state-of-the-art

---

[2] https://github.com/lamps-lab/ETDMiner/blob/master/ETD-MS-v2.0/ERD.png

models. Specifically, the text of born-digital ETDs was extracted using PyMuPDF[3]. Pytesseract, a wrapper for Google's Tesseract-OCR engine (Google, 2006), was used to extract text from scanned ETDs. We developed a specialized prompt[4] to classify ETDs into ProQuest subject categories (ProQuest, 2023) using GPT-3.5. The prompt defines the role of the AI assistant for this task and provides detailed instructions for analysis and output formatting. The classification leverages the first 19,900 tokens of an ETD as input and assigns a ProQuest subject category to an input. We used the Fine-Tuned T5 Small (Raffel et al., 2020) and Pegasus (Zhang et al., 2020) for text summarization, which were stored in the "Summarizers" and "Object_summaries" tables. For tables such as "Object_topics", "Topic_models", "Objects", and "Object-object_neighbors", we used LDA, LDA2Vec, and BERT (Devlin et al., 2018) to extract topics. We used Convolutional Neural Networks and YOLOv7 (Wang et al., 2023) to extract figures and tables. We populated user-related tables such as "User_classes", "User_queries", "user-User_neighbors", "Users", and "User_queries_clicks" with dummy data. Automatically populating the database with 1,000 ETD entries took about 11 minutes on a virtual machine with 32 CPUs, 125 Gigabytes of memory, and Hard Disk Drives, indicating reasonable scalability. Most metadata fields successfully map to our new schema, with exceptions discussed in the next section.

**Mapping Between Simplified Dublin Core, ETD-MS v1.1, and the Core Component of ETD-MS v2.0**

The new schema enhances FAIR principles by improving ETD data findability through content-level metadata in digital library search engines, and reusability through collection subsets used for individual projects. Interoperability is reduced due to the introduction of new fields. To mitigate this issue, we map fields in the new schema to equivalent fields in the existing metadata schema.

Our proposed schema is compatible with the Dublin Core and ETD-MS v1.1. Table 3[5] shows the mapping across the three schemas. The mapping was performed by aligning fields with the closest

---

[3] https://pymupdf.readthedocs.io/en/latest/
[4] https://github.com/lamps-lab/ETDMiner/blob/master/ETD-MS-v2.0/ETD_classification_prompt.py
[5] https://github.com/lamps-lab/ETDMiner/blob/master/ETD-MS-v2.0/Mapping-ETD-MSV2.2.pdf

interpretations. For example, the "dc.rights" field from Dublin Core and ETD-MS v1.1 is mapped to "ETDs.owner_and_statement" in ETD-MS v2.0. In certain cases, we could only map an entity/attribute in ETD-MS v2.0 to ETD-MS v1.1 or DC. For example, "ETDs.discipline" in ETD-MS v2.0 maps to "thesis.degree.discipline" in ETD-MS v1.1 but does not map to any field in DC. The "ETDs.References.reference_text" in ETD-MS v2.0 maps to "dc.relation" in DC but does not map to entity/attribute in ETD-MS v1.1.

### Discussion and Conclusion

One limitation of the proposed schema is that it was developed based on a corpus of 500 ETDs which may not fully capture the metadata of ETDs beyond the scope of selection. When implementing the schema using an independently selected ETD sample, we found additional fields that were not incorporated into our schema. For example, certain ETDs included metadata fields specifying different dates when the ETD was officially added to the repository and when it became accessible. Our schema includes only a single "year" field to represent the year of publication. In addition, a fraction of ETDs contains a "Peer-reviewed" field, indicating whether the document has undergone peer review, which was not included in our schema. In the future, we will refine the schema by incorporating more fields to make it more comprehensive. We will also collect feedback from ETD users.

In summary, the expansion of ETDs in academic repositories necessitates a more comprehensive and fine-granular metadata standard to comply with FAIR principles. Our proposed schema extended existing standards by providing a more complete, detailed description, and integrating content-level metadata. This ETD schema addresses critical gaps between existing ETD metadata and the growing need to mine and study knowledge embedded in ETD content for various downstream tasks.

# References

Aboelnaga, A. A., Sivakumar, A., Narla, J., Dasu, P. U., Seetharaman, R., Bhaskar, S., & Srinivas, S. S. (2024). Final Report CS 5604: Information Storage and Retrieval.

Ahuja, A. (2023). Analyzing and Navigating Electronic Theses and Dissertations. PhD defense 12 June 2023, Virginia Tech, Computer Science, with [demonstration of HTML results](demonstration of HTML results) and [related user manual](related user manual), http://hdl.handle.net/10919/115817, winner of 2023 Innovative ETD Award (awarded by NDLTD).

Ahuja, A., Devera, A., & Fox, E. A. (2022, November). Parsing electronic theses and dissertations using object detection. In Proceedings of the first Workshop on Information Extraction from Scientific Publications (pp. 121-130). https://aclanthology.org/2022.wiesp-1.14

Ahuja, N., Bansal, R., Ingram, W. A., Jude, P., Kahu, S., & Wang, X. (2018). Big Data Text Summarization: Using Deep Learning to Summarize Theses and Dissertations. http://hdl.handle.net/10919/86406

Banerjee, B. (2024, June). Improving access to ETD elements through chapter categorization and summarization. Virginia Tech, Computer Science, doctoral dissertation. https://hdl.handle.net/10919/120890

Choudhury, M. H., Salsabil, L., Ingram, W. A., Fox, E. A., & Wu, J. (2024, March). ETDPC: A Multimodality Framework for Classifying Pages in Electronic Theses and Dissertations. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 38, No. 21, pp. 22878-22884). https://doi.org/10.48550/arXiv.2311.04262

da Silva Santos, L. O. B., Burger, K., Kaliyaperumal, R., & Wilkinson, M. D. (2023). FAIR data point: a FAIR-oriented approach for metadata publication. Data Intelligence, 5(1), 163-183. https://doi.org/10.1162/dint_a_00160

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional

transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

https://doi.org/10.48550/arXiv.1810.04805

Gartner, R. (2003). MODS: Metadata object description schema. JISC Techwatch report TSW, 03-06.

Google (2006). Tesseract-OCR.

https://github.com/tesseract-ocr/tesseract?tab=readme-ov-file

Hickey, T., Pavani A., Suleman, H. ETD-MS v1.1: Metadata standard for electronic theses and

dissertations. https://ndltd.org/wp-content/uploads/2021/04/etd-ms-v1.1.html

Ho, Jeannette, and Charity Kay Stokes. "Metadata Guidelines for Digital Resources at Texas A&M

University Libraries." (2019). https://hdl.handle.net/1969.1/175368

Hunter, J., Cole, T., Sanderson, R., & Van de Sompel, H. (2010). The open annotation collaboration: A data

model to support sharing and interoperability of scholarly annotations.

Ingram, W. A., Wu, J., Kahu, S. Y., Manzoor, J. A., Banerjee, B., Ahuja, A., ... & Fox, E. A. (2024). Building

datasets to support information extraction and structure parsing from electronic theses and

dissertations. International Journal on Digital Libraries, 1-22.

https://link.springer.com/article/10.1007/s00799-024-00395-4

Ingram, W. A., Banerjee, B., & Fox, E. A. (2019). Summarizing ETDs with deep learning. Cadernos BAD, 1,

46-52. https://doi.org/10.48798/cadernosbad.2014

Kahu, S. Y., Ingram, W. A., Fox, E. A., & Wu, J. (2021). Scanbank: A benchmark dataset for figure extraction

from scanned electronic theses and dissertations. arXiv preprint arXiv:2106.15320.

https://doi.org/10.48550/arXiv.2106.15320

Koster, L., & Woutersen-Windhouwer, S. (2018). FAIR Principles for Library, Archive and Museum

Collections: A proposal for standards for reusable collections. Code4Lib Journal, (40).

Lagoze, C., & Van de Sompel, H. (2003). The making of the open archives initiative protocol for metadata

harvesting. Library hi tech, 21(2), 118-128.

Lamba, M., & Madhusudhan, M. (2019). Mapping of ETDs in ProQuest dissertations and theses (PQDT)

global database (2014-2018). Cadernos BAD, 1, 169-182.

https://doi.org/10.5281/zenodo.3599788

Osman, R., AMK, Y. I., & Abrizah, A. (2023). Metadata matters: evaluating the quality of Electronic Theses

and Dissertations (ETDs) descriptions in Malaysian institutional repositories. Malaysian Journal of

Library and Information Science, 28(1), 109-125. https://doi.org/10.22452/mjlis.vol28no1.7

Park, J. R., Brenza, A., & Lu, C. (2015). A comparative analysis of metadata best practices and guidelines:

issues and implications. International Journal of Metadata, Semantics and Ontologies, 10(4),

240-260. https://doi.org/10.1504/IJMSO.2015.074751

Potvin, S., Thompson, S., Rivero, M., Long, K., Lyon, C., & Park, K. (2015). Texas Digital Library Descriptive

Metadata Guidelines for Electronic Theses and Dissertations, Version 2.0.

http://hdl.handle.net/2249.1/68437

ProQuest. UMI subject categories guide. ProQuest. Retrieved from

https://pq-static-content.proquest.com/collateral/media2/documents/umi_subjectcategoriesgui

de.pdf

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the

limits of transfer learning with a unified text-to-text transformer. Journal of machine learning

research, 21(140), 1-67.

Richardson, R., & Fox, E. A. (2008). Using concept maps in NDLTD as a cross-language summarization tool

for computing-related ETDs.

Richardson, L. (n.d.). *Beautiful Soup documentation*. Beautiful Soup

4.https://beautiful-soup-4.readthedocs.io/en/latest/

Rushing, A., Koenig, J., Mitchell, A., Moen, W., Strawn, T., & Thomale, J. (2008). Texas Digital Library

      Descriptive Metadata Guidelines for Electronic Theses and Dissertations, Version 1.0. Prepared

      for and published by the Texas Digital Library.

Tani, A., Candela, L., & Castelli, D. (2013). Dealing with metadata quality: The legacy of digital library

      efforts. Information Processing & Management, 49(6), 1194-1205.

      https://doi.org/10.1016/j.ipm.2013.05.003

Thompson, S., Liu, X., Duran, A., & Washington, A. (2019). A case study of ETD metadata remediation at

      the University of Houston libraries.

      https://journals.ala.org/index.php/lrts/article/view/6764/9320

Uddin, S., Banerjee, B., Wu, J., Ingram, W. A., & Fox, E. A. (2021, December). Building A large collection of

      multi-domain electronic theses and dissertations. In 2021 IEEE International Conference on Big

      Data (Big Data) (pp. 6043-6045). IEEE. https://doi.org/10.1109/BigData52589.2021.9672058

Wang, C. Y., Bochkovskiy, A., & Liao, H. Y. M. (2023). YOLOv7: Trainable bag-of-freebies sets new

      state-of-the-art for real-time object detectors. In Proceedings of the IEEE/CVF conference on

      computer vision and pattern recognition (pp. 7464-7475).

Wani, J. A. (2019). Open access electronic thesis and dissertation repositories: An assessment. Library

      Philosophy and Practice, 0_1-11. https://digitalcommons.unl.edu/libphilprac/2528/

Zhang, J., Zhao, Y., Saleh, M., & Liu, P. (2020, November). Pegasus: Pre-training with extracted

      gap-sentences for abstractive summarization. In international conference on machine learning

      (pp. 11328-11339). PMLR.