# Mathematical Information Retrieval

## Search and Question Answering

**Richard Zanibbi**
Department of Computer Science
Rochester Institute of Technology
rxzvcs@rit.edu

**Behrooz Mansouri**
Department of Computer Science
University of Southern Maine
behrooz.mansouri@maine.edu

**Anurag Agarwal**
School of Mathematics & Statistics
Rochester Institute of Technology
axasma@rit.edu

arXiv:2408.11646v1 [cs.IR] 21 Aug 2024

# Contents

# Mathematical Information Retrieval

Richard Zanibbi[1], Behrooz Mansouri[2] and Anurag Agarwal[1]

[1] *Rochester Institute of Technology; {rxzvcs,axasma}@rit.edu*
[2] *University of Southern Maine; behrooz.mansouri@maine.edu*

ABSTRACT

Mathematical information is essential for technical work, but its creation, interpretation, and search are challenging. To help address these challenges, researchers have developed multimodal search engines and mathematical question answering systems. This book begins with a simple framework characterizing the information tasks that people and systems perform as we work to answer math-related questions. The framework is used to organize and relate the other core topics of the book, including interactions between people and systems, representing math formulas in sources, and evaluation. We close with some key questions and concrete directions for future work. This book is intended for use by students, instructors, and researchers, and those who simply wish that it was easier to find and use mathematical information.

# Preface

**This is an early draft of the book.** All chapters are now complete, but are likely to be revised in a later version.

If you have any comments questions, please feel free to contact the authors at [rxzvcs@rit.edu](mailto:rxzvcs@rit.edu).

Best wishes,

Richard Zanibbi
Behrooz Mansouri
Anurag Agarwal

# 1

---

## Sources, Tasks, and Search

---

We often pause to search when something that we read, watch, or hear prompts questions that we want answers to. We then go about finding answers using additional sources of information: some already exist, some are created in response to requests (*e.g.,* emails or search results), and some are created to record and organize what we find along with partial or complete answers to our questions.

In this way, information sources are the backbone of our information ecosystems. The sources available in our information ecosystem place a hard limit upon which questions we can answer. In addition to the information content in a source, its terminology, notation, writing style, and other factors determine the amount of information one can recover from a source, and how accurately and completely. This is a key reason why math instructors that communicate well are so highly regarded: they help us more easily understand topics by *how* they speak, write, and present exercises. Through course materials, lectures, and conversations, these instructors provide multiple sources tailored to their students' level of understanding and communication style.

Outside the classroom, we still often find ourselves in need of mathematical information. It might be as simple as finding a formula to

convert temperatures in Fahrenheit to Celsius, or the formula associated with a name (*e.g.,* inverse document frequency). Or the goal may be more complex, such as understanding a proof of the *sensitivity conjecture.*

As we look for answers, in addition to collecting sources that we find, we will in some way annotate and organize them to identify and apply pertinent information, *e.g.,* to find other sources, choose different search terms, execute suggested exercises, and make notes about partial answers to our questions. The effort needed for these tasks depends largely on the content and presentation in the sources that we access, including search results and discussion online and in-person with helpful people. To save time, we will often annotate sources and create additional sources of our own (*e.g.,* bookmarking a web page, placing notes in a file, or highlighting a PDF document).

In this book, when we speak about **sources**, we are referring to individual documents, recordings (*e.g.,* videos) or other artifacts that contain information directly. Libraries and other people are of course also information sources, in the sense that they can provide information, but here we use 'sources' to refer to records of specific information.

To help organize our study of mathematical information retrieval, in this chapter we introduce a framework for information tasks based on sources, and the different ways that we retrieve, analyze, and use sources to answer mathematical questions. This information task framework is built upon two main ideas:

1. Search begins, progresses, and ends with sources.
2. Tasks other than search are often needed to find information.

The key components of the framework are:

- **information needs** that individuals have,
- **sources of information** that we search, consult, and create,
- information **tasks** performed to address information needs, and
- their roles in **search algorithms and user-interfaces**.

In the next section, we consider how these components interact when

we have a mathematical question that we wish to answer.

## 1.1 When and where do we search?

Some short answers to this question are (1) when we have a question, and (2) wherever is easiest. While not very satisfying, these answers are basically correct. Search is generally performed as part of some larger information task, and not for its own sake.[1] This motivates finding quick paths to answers.

However, technical subjects such as mathematics can be complex. Finding and understanding information on math may require multiple activities, such as web search, reading sources (*e.g.,* Wikipedia pages and textbooks), taking notes, talking to instructors or colleagues, and doing exercises. As a result, when retrieving technical material on math and other specialized topics (*e.g.,* law, chemistry, music history), it is helpful to understand how search interacts with other information tasks.

To illustrate, consider the more general problem of *sensemaking*, which learning about detailed mathematical topics is closely related to.[2] In sensemaking, we construct a conceptual understanding of a topic with many sources, usually along with communicating this understanding. Common examples include writing a school term papers on an unfamiliar topic (*e.g.,* applications of category theory), or summarizing a complex historical event from multiple news reports.

Sensemaking tasks are challenging because information must be found in multiple sources, but also because this information must be analyzed, compared, and integrated. These thinking activities often require most of the effort for sensemaking. To manage these thinking tasks, we record plans, notes, and outlines to organize our work. These working documents may be checked repeatedly as we work, and as we write our final summary. They are themselves important information sources that provide the scaffolding needed to focus and ultimately complete work on a sensemaking task.

To further illustrate information tasks that complement search, imagine taking handwritten notes on eigenvectors as described in an

---

[1]A fact that is both important and humbling for IR researchers.

[2]See (Hearst, 2009) for an overview of early research on sensemaking.

algebra textbook. The notes allow us to annotate this source with our own observations, and record them for reference at a later time. The analysis and insights in the notes come from applying information that we know and find. These notes communicate a new information source to a specialized audience: ourselves.

For our notes to be useful, we organize them. Perhaps this is a purple sticky note that we attach to a monitor to consult later in the evening, after dinner. Or, perhaps we use a paper notebook with separate sections for different subjects, along with other organizational devices (*e.g.,* sticky notes acting as bookmarks). We might instead be using a tablet computer, which also provides handwriting recognition to convert the notes to computer-searchable data (*e.g.,* using `Ctrl-f`).

The information tasks above are distinct from a basic search task where we submit a query, post a question, or send an email to obtain new information sources. However, it turns out that search engines implement variations of the same basic tasks described above: they need to *index*, *communicate*, *annotate*, and *apply* information in sources to be effective. For example, we organize sources when we arrange sticky notes by topic and color on a wall, or construct an *inverted index* mapping words or formulas to their document locations: these are both forms of *indexing*. As another example, search engines produce Search Engine Result Pages (SERPs) summarizing documents matching a query, and question answering systems or AI 'bots' produce answers. These retrieval system outputs and our notes are *communications* creating new information sources.

Making notes on a passage requires us to *apply* information to create an *annotation*: additional information associated with the passage. In turn, if those notes were handwritten on a tablet computer, a system converting these to text and LATEX for math applies information captured in an algorithm, annotating the notes themselves. We end up with a hierarchy of annotations: the notes annotate a passage, while a recognition algorithm annotates the notes.

In our framework we will distinguish different source types, based largely on what information tasks they are primarily used for. More specifically, we distinguish:

1. available sources on a topic including search queries and results,
2. information added to sources (*annotation*), and
3. structures and organizations created for search (*indexing*)

Getting back to our motivating question, when we have identified a mathematical *information need*, we generally start with questions, and hope to end with one or more information sources that we feel address or ideally answer those questions (*i.e., relevant* sources for the information need). Where we search is motivated by the types of sources we expect to find from places online and/or the physical world (*e.g.,* conversations and post-its). Unless we are casually browsing resources on a topic, the places and order in which these sources are found will generally reflect attempts to reduce our time and effort.[3] Relevant sources are often of different types: perhaps a passage in a web page along with a SERP page, an answer from an online AI system, an email from a friend, and a green sticky note on your monitor.

From this perspective, math-aware search engines and question answering systems are important tools, but only one among many resources for finding math information, and only a small part of what happens when we search.

## 1.2 Information task framework

While we focus in this book on information retrieval using computers, we wish to address sources in their broadest sense here. Not all sources are text documents, and not all sources are recorded in documents. Consider an informal conversation about Bayesian decision theory in the hallway, or observing that there are no clouds in the sky: often, your only record of important information sources is your own memory.

In addition to textbooks, technical papers, and web pages, in recent years the types of resources used to locate mathematical information has grown to include substantial amounts of video (Davila *et al.*, 2021) and audio, *e.g.,* for course lectures, tutorials, and technical talks. Community Question Answering sites, and direct question answering is provided by

---

[3]Information foraging theory (Pirolli and Card, 1999) suggests we evolved to gather and consume information similar to food, governed by cost-benefit analyses.

resources such as Math Stack Exchange,[4] Wolfram Alpha, and large language models such as the Generative Pre-trained Transformer (GPT).

It is worth noting that when we have found or produced information we want to share or reuse, we usually produce a source of information ourselves. For example, when we have an answer to a math homework question, we create a physical or digital document, so that this can be checked by ourselves and graded by our instructor. If we found a helpful video while doing the homework, we might share it in a text message, which is itself a form of 'micro-source.'

**Differences in Information Sources.** Especially when we include information obtained directly from our environment along with modern computing and communication devices, information sources may come in many forms. Sources vary in the dimensions listed below, among others.

- immediacy, *e.g.,* having a conversation vs. reading a transcription
- authorship, *e.g.,* human, machine generated, environment
- interactivity, *e.g.,* a human/chabot conversation vs. a document
- audience, *e.g.,* grade school students vs. math professors
- modality, *e.g.,* text, video, audio, or a web page combining these
- purpose, *e.g.,* textbook, search results, or a search index
- structure, *e.g.,* free text in a sticky note, vs. a book with chapters
- length
- formality, *e.g.,* proof vs. text message
- style, *e.g.,* how concepts and examples are communicated
- correctness, *e.g.,* correct vs. incorrect definition or proof
- completeness, *e.g.,* partial vs. full search index

For the *immediacy* of a source, we are referring to whether the source comes directly from observing a person's environment (*e.g.,* through conversation, experimentation, travel, etc.) or is recorded, as in a document or audiovisual recording. Particularly with the advent of large language models, authorship and correctness are important concerns. In many cases LLMs and other sources may appear credible but are incorrect.

---

[4]https://math.stackexchange.com

Knowing how a source is created can help us determine how trusting we should be of it when we are uncertain about validity (*e.g.,* from the perceived expertise of an author or system).

The intended audience and style of a source are also critical concerns. They determine the prerequisite knowledge needed to decide whether a source is relevant and to interpret and use information in the source.

**Information task types.** The primary task types we will use for people and systems come from common descriptions of sensemaking and simpler information tasks: retrieving, analyzing, and synthesizing. For this framework, these **tasks are considered at the source level**. For example, *analyzing* a source refers to analysis that produces additional recorded information for a source (*e.g.,* in a note) rather than reading and interpreting a source, which does not produce a new observable artifact.

We subdivide each of these into two subtasks based on how sources are created and used, producing six tasks in total.

1. **Retrieve**
   - *Query* to request sources of information
   - *Consult* and interpret available sources, *examining* and *navigating* within and across sources
2. **Analyze**
   - *Annotate* sources with additional information, *e.g.,* notes, add formula locations
   - *Index* sources by organizing them for retrieval
3. **Synthesize**
   - *Apply* available information that we know, have in available sources, or is encoded in algorithms, etc.
   - *Communicate* information by creating new sources

The *apply* task is critical, and used for all other tasks (*e.g.,* creating queries, consulting or creating information sources, or generating annotations and/or indexes). It is distinct because people often apply information without producing observable sources. For example, recognizing what a variable represents does not involve creating an an-

⋮
TermID, 62
Term normalization, 30
Term-partitioned index, 70
Terms. *See also* Queries
    BIM ranking function, deriving,
       207
    defined, 3, 21
    function notations, xi
    partitioning, 416
    statistical properties of, 82
    tree-structured dependencies, 213
    vectors, weighting and, 113
Term weighting. *See* Weighting
Test data, 237
⋮

**Figure 1.1:** Index Excerpt from "Introduction to Information Retrieval" (2008) by Manning, Raghavan, and Schütze.

notation, index, or source outside of our own minds. The *apply* task also identifies an important commonality between thought and computation (*e.g.,* algorithms): both apply information but using differing levels of formality, flexibility, and automation.

Not all information needs require queries. If we have a helpful document describing the inverse document frequency on our laptop, we may simply consult it to review previously highlighted passages. This locating of items in an available source or across available sources through references and links is known as *navigation*, which is distinct from submitting a *query* to a system or person to find new sources. As another example of using navigation to satisfy an information need, in some case we me may simply use the contents of available sources directly (*e.g.,* copying-and-pasting into an online form).

The process of analyzing a source and recording a map for use in retrieval is known as *indexing.* Consider Figure 1.1, where a book index provides a map for the book, so that a reader can quickly *navigate* to parts of the book discussing 'Terms,' for example. Contrast this *subject index* with the index used in a traditional term-based search engine, which provide a much simpler map known as a *concordance* recording where specific terms appear in documents (Duncan, 2021). While these

**Figure 1.2:** Information Task Framework: The Source Jar. The jar contains source 'marbles' that we initially have to start, and as we work we add, create, annotate and organize the sources in the jar. *Apply* has a dashed line because it is used in other tasks.

different indices are both used for retrieval, they differ in their scales (one document vs. a collection) and intended audiences (human reader vs. search algorithm). Other forms of indexing are less formal, such as collecting and organizing notes on different sticky notes for easier use.

As discussed earlier, we distinguish tasks for analyzing sources in terms of organizing them for use and retrieval (*indexing*), and adding information to sources (*annotation*). Annotations are often used in indexing sources, such as adding formula locations for PDF documents.

**Source Jar Framework.**  To put sources and the tasks used to create them in a more intuitive relationship, Figure 1.2 visualizes our task framework as a jar of sources with a lid. The jar contains immediately available sources as marbles in the jar. Each marble has an identifying color and shape. The source marbles contain information of different types, and may refer to other sources inside and outside of the jar. Sources that are directly available are either with us, or inside the jar.

Stickers on the jar identify the information task types we can perform, and the lid is labeled with the need it is used to address. When we find or create a new information source, we add a marble to the jar. If a

new source *annotates* another source, we place it in a container with the source it describes inside the jar (*e.g.,* using a small plastic box). *indexing* produces a marble containing a description of which sources it organizes, and how. We take source marbles and containers out of the jar to use them. When they are no longer useful, we return them to the jar. It is also possible to lose sources when the jar is accidentally left open and 'spilled.'[5]

We can imagine having a shelf of these jars for different information needs. To reuse or get additional information for a need we worked on previously, we open a jar from our shelf.[6] For a new information need we create a jar, adding any potentially useful sources from our other jars. When we stop working, we may select any final sources for use later, and then close the lid.

This informal jar model is intended to roughly capture how people experience working with information in a simple way. It captures observable sources and observable task actions. We tend to move from source to source, performing tasks of some specific types with a goal in mind. We are often unaware of why we performed tasks in a particular order, and so this is not represented explicitly, other than as marbles moving in and out of the jar through time.

## 1.3   Information needs and search strategies

**Information Needs**   When searching for information on math, what we need to find will vary from finding definitions for terminology, math symbols, formulas, operational knowledge such as proof techniques, applications of mathematics (*e.g.,* information retrieval models), resources for instructors, and detailed information on mathematical spaces, theorems, etc.

Example 1 illustrates information needs that different audiences may be seeking to address using the same query, along with a list of sources that might be used to address their needs. These needs vary from finding definitions to exploring sophisticated relationships between

---

[5] *e.g.,* 'the dog ate it,' 'my internet is down,' or 'I know it's here...somewhere.'

[6] *e.g.,* 'Wait; I forgot one of the types of category theory applications I wanted to discuss in my paper from my notes...'

Example 1: Differing Information Needs

**Query:** What does $a^2 + b^2 = c^2$ represent and how is it useful?

*Students* might use this query to learn the Pythagoras theorem, and perhaps find an example demonstrating the theorem, and a possible proof.

*Educators* may have similar interests to students, but may seek additional resources on how to teach this result.

*Researchers* can have very different interests than the other audiences. They may be interested in one of more of the following:

- For a mathematician: Is this true in a general metric space and/or a Hilbert space?
- For a physicist: How is it linked to the probability assignments in quantum mechanics?
- For an IR expert: How is it related to probability assignments in a Hilbert space used in describing interaction for information retrieval?

**Possible Relevant Sources:** web-pages (MSE, Brilliant.org etc.), YouTube videos, online lecture slides, text documents (*e.g.,* digital books (OpenStax, LibreTexts), articles, online notes (MIT OpenCourseWare)).

the generalization of the theorem in different mathematical structures (Hilbert spaces) and applications in other fields (quantum mechanics).

As a result, the types of sources needed by each audience differ dramatically, but the initial (admittedly vague) query is identical: the *query intent* differs for these audiences. For math information needs, we have found it important to consider information needs both in terms of the desired information, as well as who is searching. Also shown in the example are relevant sources that might be used to address these needs from the different audiences, whether wholly or in part.

For a broader sense of the types of mathematical information needs users have online, Table 1.1 illustrates information needs for math organized by Broder's taxonomy of needs/intents behind web search queries (Broder, 2002). While some question the usefulness of the transactional class in Broder's model, for math, the transactional class is a useful distinction. For example, a user may be looking to *refind* a web page they used to enter formulas in LaTeX (i.e., a navigational intent). Or, they may instead be looking to find such a web page for the first time, thereby looking to interact/transact with as-yet unknown websites (i.e., a transactional intent).

Within the *informational* needs class, a distinct subclass of *computational/operational* information needs exist. These include needs to evaluate or simplify a formula, or to produce a proof for a statement using logical operations. It was useful to distinguish questions that were seeking *concepts, proofs,* and *computation* for the ARQMath shared tasks (Mansouri *et al.*, 2022b) that we discuss in Chapter 3.

In our work we have found it useful to consider math information tasks in two dimensions, based on the type of information need as shown in Table 1.1, and the user's mathematical background. More formally, we have a space/set of mathematical information needs $N$ defined by a Cartesian product of possible information needs $(T)$ypes and user/audience $(B)$ackgrounds $(N \in T \times B)$. How these types and backgrounds interact is illustrated in Example 1.

**Search strategies.**    For a given information need, it helps to think about *strategies* that might be used to satisfy it. We can sketch these in strategy 'jar' diagrams as seen in the panel labeled Strategy 1. The

**Table 1.1:** Examples of Mathematical Information Needs within Broder's Taxonomy (Broder, 2002)). A user's math background is another dimension.

| **Navigational** | Find a specific source ('known item' retrieval) |
|---|---|
| | Web page (*e.g.,* for formula entry) |
| | Document (*e.g.,* Book, Technical Paper) |
| | YouTube or Khan Academy Video |
| | Podcast |
| | |
| **Transactional** | Find online resources for use/interaction |
| | Formula entry |
| | Evaluating and plotting a formula |
| | Simplification of a formula |
| | Interactive theorem proving |
| | |
| **Informational** | Find information for a topic or question |
| | *Sub-categories*: computation, concepts, and proofs |
| | How to compute an expression (*e.g.,* integral) |
| | Symbol and operation definitions (*e.g.,* $\zeta$, $\binom{n}{k}$) |
| | Concept name(s) associated with a formula |
| | When is a function not differentiable? |
| | Who was Gauss? |
| | Proof drafts for P = NP |

diagrams identify an information need, initial queries, expected tasks, and where relevant sources might be found. We can imagine beginning a new search by writing the information need on the lid, putting already available source 'marbles' in the jar, and then writing planned tasks on the jar labels. For readability, we use informal descriptions for the three main task types along with initial queries. For this discussion, in our examples strategies we identify multiple relevant sources already known to the authors, some or none of which a searcher might actually know when they begin searching. An actual list of possibly relevant sources may be much shorter, or empty.

Let us first consider search strategies that might be used by undergraduate students, for learning how to complete a square, and to change the base of a logarithm (Strategy 1 and Strategy 2). In both examples, two queries that might be used are given, and the *Synthesis* tasks clarify

---

**Strategy 1: (Student) Completing the Square**

**Retrieve:**

**Query:** completing the square $OR$ $ax^2 + bx + c = (\star)^2 + \text{constant}$?

Search using the text query or possibly the symbolic query; $(\star)$ is a wildcard for any subexpression. Identify where the general method can be found, and examine the proof of the result.

**Analyze:** Mark-up/bookmark sources to identify useful information. Use a notebook to summarize key details found in sources. Save examples for different cases, *e.g.,* $a, b > 0$ and $c \leq 0$.

**Synthesize:** Solve an integration problem on paper, such as $\displaystyle\int \frac{1}{x^2 + 4x + 3}\, dx$.

**Possible Relevant Sources:** web-pages (Math Stack Exchange (MSE), Brilliant.org etc.), YouTube videos, online lecture slides, text documents (*e.g.,* digital books (OpenStax, LibreTexts), articles, online notes (MIT OpenCourseWare)), Online tutorial sites (Khan Academy, Magoosh), online answer engines (WolframAlpha), online calculators (SymboLab).

---

the specific information need: the source they want to produce. In the first example, this involves completing an exercise on paper, and in the second example, obtaining a value from a calculator. Note that for the queries containing formulas, students might find it difficult or be unable to express the formulas in queries using a standard text query box, particularly if they are unfamiliar with LaTeX.

Now let's consider more advanced information needs for researchers. The researchers may be interested in following progress on an old conjecture (*e.g.,* Riemann Hypothesis). Or, they may be interested in learning about a new possible proof of the problem, or perhaps they were unfamiliar with the problem but are curious to know more about it. Strategy 3 seeks information and a proof for a problem that was posed in 1994. It became a major unresolved question in mathematical computer science until 2019, when Hao Huang solved it.

As another example, imagine that a researcher encounters a technical statement for the sensitivity conjecture, but which does not name it. They want to know the status of the statement, and if there are associated results they can use in their own work. Here the searcher only wants to learn the conjecture's name, properties, and proofs for later reference. The strategy from Strategy 3 needs to be altered, as reflected in Strategy 4. In this second case, the researcher has a document

## Strategy 2: (Student) Log Base Change

**Retrieve:**
   **Query:** log base change *OR* how to convert $\log_b x$ to $\log_c x$?

   The student may use the text or symbolic query. Find sources giving the conversion rule with general bases.

**Analyze:** Markup sources and note down where relevant sources are located in a list (*e.g.,* in a text file). Save some special cases like converting $\log_{10} x$ to $\ln x$.

**Synthesize:** They use this to compute $\log_4 13$ on a calculator as the $\boxed{\log}$ button on most calculators only represents $\log_{10}(\cdot)$.

---

**Possible Relevant Sources:** web-pages (MSE, Brilliant.org etc.), YouTube videos, online lecture slides, text documents (*e.g.,* digital books, articles, online notes (MIT OpenCourseWare)), Online tutorial sites (Khan Academy), online answer engines (WolframAlpha), online calculators (SymboLab), online databases for definitions and theorems (ProofWiki).

## Strategy 3: (Researcher) Sensitivity Conjecture

**Retrieve:**
   **Query:** What is **Sensitivity Conjecture**? Has it been proven?

   Find papers/books defining the conjecture and providing proofs.

**Analyze:** Since the conjecture is very technical, retrieved material is annotated with sources where terminology in the conjecture can be comprehended. An index (graph) is made capturing the chronological account of progress on the proof.

**Synthesize:** Results and the methods for proving this conjecture are used for similar problems, and new articles/material are created to disseminate the findings.

---

**Possible Relevant Sources:** online encyclopedias (Wikipedia, Wolfram MathWorld), online Q&A sites (MathOverflow.net, AoPS, sciencedirect.com), YouTube videos, online lecture notes, text documents (*e.g.,* digital books, research articles), online science & math magazines (Quanta Magazine), online math databases (Cornell's mathematics library, zbMATH Open, $\mathcal{AMS}$: Math Reviews)

**Strategy 4: (Researcher) Unknown Conjecture**

**Retrieve:**

**Query:** Any set $H$ of $2^{n-1} + 1$ vertices of the $n-$cube contains a vertex with at least $\sqrt{n}$ neighbors in $H$.

The search is done using a textual query with LaTeX for the formulas. Related papers/books are collected and consulted for theorem definitions and proofs.

**Analyze:** Retrieved sources are annotated with links to other sources where terminology used can be comprehended. Highlight the name of the statement when it is found.

**Synthesize:** Create document summarizing the theorem name and key details, with cites/links to key sources found. Include link to a file directory on a laptop where additional notes in text and LaTeX files can be found, if any.

**Possible Relevant Sources:** online encyclopedias (Wikipedia, Wolfram MathWorld), online Q&A sites (MathOverflow.net, AoPS, sciencedirect.com), YouTube videos, online lecture notes, text documents (*e.g.,* digital books, research articles), online science & math magazines (Quanta Magazine), online math databases (Cornell's mathematics library, zbMATH Open, $\mathcal{AMS}$: Math Reviews)

summarizing the key findings and where sources may be found.

As yet another example, let us consider a situation where a researcher needs to define unfamiliar notation that they encounter in a paper. As shown in Strategy 5, for this information need the researcher is simply collecting and utilizing relevant sources, rather than generating a new source of their own. If successful, they will learn that this is the *Hypergeometric function & series*, and if necessary use this name later to find identities associated with

We also note that the interpretation of the majority of mathematical expressions is *context-dependent*, i.e., the same formula may refer to different concepts in different contexts. A student executing Strategy 6 will end up with multiple interpretations, which might represent:

- the distributive law: $\pi(m + n) = \pi m + \pi n$, or
- the value of the *prime-counting function* that counts the number of primes less than or equal to $m + n$.

The more general property of a single object signifying multiple entities is known as *polysemy*, such as the word 'apple' being used to represent both a food and a company, and often poses challenges for both information retrieval and natural language processing.

## Strategy 5: (Researcher) Unfamiliar Notation

**Retrieve:**

 **Query:** What is $_2F_1(a, b; c; z)$?

 The search is done as a symbolic query using LaTeX for the formula. Find places where similar notation is defined and help disambiguate this specific instance.

**Analyze:** Retrieved material is highlighted where the same or similar notation is used.

**Synthesize:** Use the definition to help understand another source being read.

---

**Possible Relevant Sources:** online encyclopedias (Wikipedia, Wolfram MathWorld), online Q&A sites (MathOverflow.net, AoPS, sciencedirect.com), YouTube videos, online lecture notes, text documents (*e.g.,* digital books, research articles), online science & math magazines (Quanta Magazine), online math databases (Cornell's mathematics library, zbMATH Open, $\mathcal{AMS}$: Math Reviews)

## Strategy 6: (Student) Unfamiliar Notation

**Retrieve:**

 **Query:** $\pi(m + n)$

 The student does not know LaTeX and so a text query is used with a unicode symbol for $\pi$.

**Analyze:** Highlight and make notes on where sources for the definition is found.

**Synthesize:** Create a web post identifying helpful sources for other students working on a group project.

---

**Possible Relevant Sources:** online resources (Wikipedia, Wolfram MathWorld), online Q&A sites (MathOverflow.net, AoPS, sciencedirect.com), YouTube videos, online lecture notes, text documents (*e.g.,* digital books, research articles), online encyclopedia such as OEIS.

**User studies and use cases.** There are a small number of papers examining math retrieval online. We know of just one study examining user behaviors when using a standard text-based search engine for math (Mansouri *et al.*, 2019b). Query logs from a Persian general-purpose search engine were used. Compared to the general case, search sessions for math topics were typically longer with more query refinements (*i.e.,* changing queries to try and improve results) and were less successful. Queries were also longer and more varied more than queries overall. In another interesting study, posts to threads in an online math Community Question Answering (CQA) site were studied (MathOverflow[7]). The authors identified patterns in the collaborative actions they exhibit (*e.g.,* providing information, clarifying a question, revising an answer) and their impact on the final solution quality (Tausczik *et al.*, 2014).

Earlier work considered use cases for math-aware search in a study of mathematics graduate students and faculty (Zhao *et al.*, 2008). Surprisingly the participants did not find formula search was useful overall, perhaps because they generally knew the names of entities they wanted to search on. The study also points out that the type of a source is an important relevance factor (*e.g.,* exercises vs. code). Another analysis of expert use cases is also available (Kohlhase and Kohlhase, 2007), in which formula search was studied using the MathWebSearch tool.

## 1.4 Retrieval systems

Figure 1.3 provides an overview of retrieval system interactions with people, and the specific sub-tasks from the 'jar' framework that they perform. Unlike the freely interacting tasks of the 'jar' model, retrieval systems generally perform information tasks in a fixed order, shown by arrows in Figure 1.3. The figure shows two main information flows for the collection of sources that a retrieval system uses.

1. **Index construction (offline).** Information passes from the sources at top and flows to the bottom-right, as sources are annotated with additional information, and then used to compile

---

[7]https://mathoverflow.net

**Figure 1.3:** Information Tasks in Retrieval Systems (Backend). Arrows show the flow of information. All tasks in Figure 1.2 other than *Apply* are shown.

    a searchable index of patterns. The collection index is precomputed before the system is used for retrieval.

2. **Retrieval (online).** Submitted queries are annotated and then matched against patterns in the index, returning one or more matching sources. The collection is generally consulted for passages, bibliographic data, and other contents when generating the result returned to the user.

**Consulting sources.**   Search engines that match queries to contents in sources are a type of filter. A standard search result is useful precisely because it contains sources with patterns of information shared with the query, omitting all other sources.

    The implementation of *consult* tasks that access sources is important for both index construction and retrieval, and is another way that sources are filtered in a retrieval system. Source contents shown in search results directly impact our impression of which returned sources are promising. Source contents used for index construction define the available patterns

for matching queries to sources.

For example, omitting high frequency terms from queries and sources that do not signify a topic (*i.e.,* skipping *stop words* such as 'the') can greatly reduce index sizes and increase retrieval speed, but at the risk of performing poorly on queries using these terms; a classic example is the phrase 'to be or not to be' from Shakespeare's play Hamlet, which is instantly recognizable but composed entirely of stop words.

For math-aware search, a similar decision would be omitting tokens and strings representing formulas (*e.g.,* in LaTeX source files). Limitations on what can be consulted includes formulas in PDF documents, which are usually not represented explicitly (Shah *et al.*, 2021). This and other missing information can be addressed by annotating sources.

**Annotation and indexing.**    In direct contrast to filtering performed in the *consult* step, we will also *annotate* sources with additional information. This extra information can be used to add patterns for matching sources in the index, or to add information to retrieval results.

For example, some neural net-based techniques such as SPLADE automatically add words that do not appear in a source to the inverted index (Formal *et al.*, 2021a).[8] These additional terms are synonyms and other words appearing in similar contexts within a training collection. For math, a simple example is adding additional representations for formulas in sources, such as generating Content MathML for operator trees corresponding to formulas represented in LaTeX or Presentation MathML, allowing formulas to be searched using both formula appearance and operation structure.

From the information obtained through *consulting* and *annotating* documents, an index of patterns for matching queries is produced. This can take different forms, but is generally one or a combination of:

1. *inverted indexes* that map entities to sources and source locations (*e.g.,* tokens or paths in graphs for math formulas), and
2. *embedding spaces* mapping entities to points in a vector space, where entities with more similar contexts across a collection are

---

[8]This augmentation is also applied to queries. Query annotations are called a *query expansion* when they add tokens or other patterns for matching additional sources in the collection.

closer (*e.g.,* words, sentences, and formulas).

Embedding vectors have their own dictionary mapping vectors to specific sources or source locations (*e.g.,* when search is done on text *passages* or individual math formulas). This allows sources matched in vectors to be consulted when communicating results to users.

In indexes, sources are referenced using a compact reference scheme to reduce the index size. For example, we might use integer tuples such as *(document, page, \*region)*, with *\*region* having a different number of integers depending on the content type. For example, for text data *\*region* might provide the starting character position for a word, or the starting and ending character positions for a sentence or other text excerpt. For math formulas, formulas expressed in L<sup>A</sup>T<sub>E</sub>X or MathML can be indexed similarly to text using character positions. However, for a PDF file the *\*region* might instead be the formula's location on the identified page, as given by a page number and the top-left and bottom-right $(x, y)$ coordinates for a *bounding box*.

**Retrieval: Querying sources and communicating system results.** We query a collection using the collection index and an annotated query containing additional terms and/or an embedding vector. Search using inverted indexes is referred to as *sparse* retrieval, while search using embedding spaces is referred to as *dense retrieval*, based on the underlying vector representations for each. In particular, term vectors representing the presence of words or formula structures in a document are mostly zeros. In general, sparse retrieval models such as BM25 that use tokens or other source contents directly for lookup are faster (Robertson and Zaragoza, 2009), but dense retrieval models such as ColBERT (Khattab and Zaharia, 2020a) are more effective (Wang *et al.*, 2023; Giacalone *et al.*, 2024). Some retrieval models use dense models to improve sparse models *e.g.,* SPLADE, mentioned earlier.

The improved effectiveness of dense retrieval models is partly from additional context used in defining patterns, *e.g.,* using the words referring to and surrounding a formula to represent a formula in a pattern vs. the formula alone. The use of a vector space also provides more holistic and flexible pattern matching, *e.g.,* finding source vectors with the most similar angles to a query vector, rather than matching

query formula tokens individually to vocabulary entries in an inverted index. These help bridge the *vocabulary problem* discussed in the next subsection.

How the final result of a query is *communicated* (generated) can vary substantially, and often makes use of query and source annotations.In a traditional search engine, specific sources are matched in the index for the *query* task, with the index comprised of some combination of inverted indexes and embedding spaces. Source contents are then used to generate a result in the *communicate* task, using sources and source locations matched in the index. However, for a generative question answering or retrieval system, the result of the *query* task may be a single vector capturing the similarity of patterns in the query to patterns within sources of a very large collection produced using a neural network. This vector is then used in the *communicate* task as a starting point for generating the response, for example using a second recurrent neural network trained on the collection, possibly along with additional information from the original collection of sources (*e.g.,* with references to specific sources). Some are used to generate a list of ranked sources directly (Zeng *et al.*, 2024), ultimately producing an *extractive* search result summary based on source contents.

Other recent systems such as Google's recently released AI search assistant produce *abstractive* summaries of retrieved sources, which summarize matching sources but without limiting the summary to contents found in the matched sources or their annotations.

**System design.**    System designers and IR researchers are interested in the efficiency and effectiveness of a retrieval system. As seen in Figure 1.3, these are observed in live systems through query and user interaction logs. For experiments, system results are computed using simulated user interactions for a fixed set of queries, and relevance scores for sources, along with a description of the information needs associated with each query. Designers and researchers also make use of additional tools for evaluation, some of which we discuss in Chapter 3.
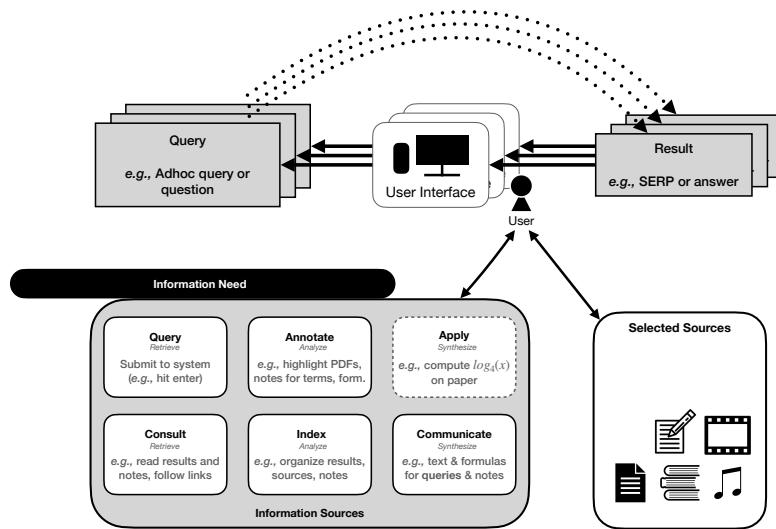
## 1.5   User-system interactions and interfaces

User interfaces play a very important role in mathematical information retrieval. In addition to executing queries and returning results, how queries are entered, how results are returned, and how other information tasks in Figure 1.2 are supported can help speed up or even limit a person's efforts to find and use information.

We next present a user-centered view of retrieval systems in math information tasks. We then share some key challenges for retrieval system interaction, along with interface designs aiming to address them.

**Interfaces-in-the-task-loop.**   Figure 1.4 illustrates a student working to change the base of a logarithm (*i.e.,* Strategy 2) using multiple retrieval systems. At bottom-left of Figure 1.4 is a jar holding sources the student had on hand when they started searching, new sources they find or create, along with other linked sources, *e.g.,* by web link, citation, or mention. In addition to these sources, their queries, results from queries, and handwritten notes (*e.g.,* from converting bases by hand) are also found in the jar. Of these available sources, the ones currently being used are at bottom-right of Figure 1.4.

Some selected sources partially or fully answer the student's needs, but others do not, such as sources later deemed not relevant. Other selected sources might exercise knowledge such as shown for the *Apply* task in Figure 1.4, or come from other tasks such as annotating and indexing sources of interest. Some selected sources may be added even after finding answers, perhaps because they provide a different perspective, or have a presentation that is easier to understand.

For more complex tasks, such as for the sensitivity conjecture retrieval strategies described in Strategy 3 and Strategy 4, we may even see the focus of our selected sources drift. In the *berry picking* model of retrieval (Bates, 1989), people see their queries and information needs change as they search and learn. Particularly for unfamiliar topics, our needs and queries may change dramatically as our understanding does (Belkin, 1980). For example, this is likely to happen when a person explores unfamiliar concepts associated with unfamiliar notation. In our jar model, information need changes involve changing the jar lid

**Figure 1.4:** Interacting with Multiple Retrieval Systems (Frontend). Each dotted arrow represents a retrieval system backend (see Figure 1.3). Sources currently used to address the information need are shown in a separate container at bottom right.

label, perhaps using an orange sticky note placed over the original description.[9]

What we have in Figure 1.4 is a person generating, selecting. and using sources for needs that may change as they work. Retrieval systems are a part of this process, but not the focus.

**Interaction challenges.** All systems embody design decisions and biases. Naturally, no one retrieval system will be ideal for all queries or subjects (*e.g.,* mathematics). However, users often have challenges in search that are more cognitive than system-related. These are important considerations in creating usable systems, particularly for search interfaces (Hearst, 2009; White, 2016; Holmes *et al.*, 2019).

Norman identifies two 'gulfs' that limit human task performance (Norman, 1988). Broadly speaking, for retrieval systems the two main categories of interaction challenges are with expressing queries (a *gulf of execution*) and interpreting results (a *gulf of evaluation*). In both

---

[9]sticky notes: a versatile information tool in this chapter and in life.

cases, for unfamiliar topics, the user may be unable to formulate an effective query or interpret results reliably precisely because of what they don't know, or because their understanding is incorrect (*i.e.,* their *Anomalous State of Knowledge* (Belkin, 1980)).

A common cause of a gulf of execution in query formulation is the *vocabulary problem*, where the terms/patterns a person uses for search differ from those used to index sources. For example, in one study undergraduate students were challenged while trying to define the binomial coefficient '$\binom{n}{k}$' (Wangari *et al.*, 2014). Because of this notation-based vocabulary problem, the students' were unable to find a definition using standard text search. When allowed to enter the expression by hand with automatic translation to LaTeX, they found definitions using the same text search engines.

People sometimes also encounter a *gulf of evaluation* when trying to identify relevant information in search results. Aside from missing relevant items in results due to the vocabulary problem, an important factor here is how retrieval results are presented to a user. For example, (Reichenbach *et al.*, 2014) report statistically significant differences in the ability of participants to identify relevant sources in SERPs when excerpts present formulas as raw LaTeX vs. rendered formulas. Additional gulfs occur when selected excerpts are not relevant for a need, or a person lacks the math background to understand a result.[10]

Learning new terminology and notation while searching allows a user to extend their patterns used to express queries and identify relevant results, bridging these gulfs of execution and evaluation. Some of these new patterns might be recorded explicitly in a source, *e.g.,* recording an unfamiliar notation for eigenvectors on a blue sticky note.

**Query input: math-aware search bars.**    For the most part, math-aware search bars differ in how they include formulas. Perhaps the simplest design is for users to enter *both* text terms and formulas as text. An early example is the Digital Library of Mathematical Functions[11] which accepts LaTeX commands for formulas along with text terms in queries

---

[10]impatience, inattention, mental strain, and tiredness are also factors here.
[11]https://dlmf.nist.gov

**Figure 1.5:** MathDeck query entry, formula chips, and cards (Diaz *et al.*, 2021). Chips can be dragged, edited, and combined. Editing may be done using raw LaTeX, or a combination of operations, chips, handwriting, and LaTeX using the canvas at center. Formula cards (bottom left) contain chips, titles and descriptions. New cards can be created by users, and searched by formula & title (video: https://www.youtube.com/watch?v=XfXQhwIQlbc).

(Miller and Youssef, 2003). The more recent Approach Zero system[12] system uses MathQuill[13] to render LaTeX formulas as they are typed in the search bar, and allows writing lines and argument positions to be reached with arrow keys rather than LaTeX commands (*e.g.,* for superscripts and fraction denominators).

To avoid remembering many names for operations and symbols, or to avoid unfamiliar LaTeX or other syntax for creating formulas, usually a palette of buttons with images for symbols and operations accompanies the search bar. Buttons add formula elements including operation structures (*e.g.,* fractions, integrals, and radicals) and symbols not found on a keyboard (*e.g.,* greek letters such as $\zeta$ (*zeta*)). Query bars with palettes often display formulas in a structured editor like those in document editors (*e.g.,* Word). Early examples of prototypes with symbol/operation palettes include MathWebSearch (Kohlhase and Prodescu, 2013) and MIAS (Sojka *et al.*, 2018).

---

[12]https://approach0.xyz
[13]http://mathquill.com

As another way to reduce the effort and expertise required for formula entry, some search bars also support *multimodal* formula entry. Multimodal query editors allow formulas to be uploaded from images or entered using handwriting in addition to standard keyboard and mouse-based entry. There are also multimodal tools such as Detexify[14], which looks up LaTeX commands for symbols drawn using a tablet or mouse (Kirsch, 2010). In addition to search, recognizing math in handwriting and images has been used for interactive computer algebra systems and other applications, and is an active area of research dating back to the 1960's (Zanibbi and Blostein, 2012; Truong *et al.*, 2024).

An example of a search bar with multimodal formula entry is the MathDeck system[15] (Diaz *et al.*, 2021). As seen at top-left in Figure 1.5, a text search box can be used to enter words and LaTeX for formulas. Formulas can also be added from a visual formula editor shown at center-left in Figure 1.5, and using formula 'chips' with embedded LaTeX (*e.g.,* blue oval at right of the query text box). Like MathQuill and structured formula editors, MathDeck renders a formula as it is entered, but with more flexible subexpression selection and entry. MathDeck's query and formula entry interface is designed to:

1. support text entry; natural for text, and one can easily type 'x + 2', or copy-and-paste LaTeX with small changes (*e.g.,* $a \to x$)
2. provide symbol palettes to help enter symbols and structures
3. provide handwriting input for those who prefer it, and to avoid searching palettes for symbols & structures
4. support formula reuse in chips; chips can be used in editing, and can be exported/shared as images with LaTeX metadata
5. construct formulas interactively using a structured editor, with larger formulas easily built up from smaller pieces.

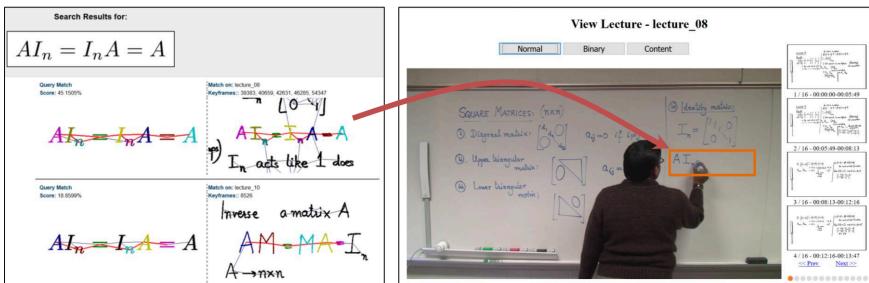Other multimodal query entry interfaces have similar design goals, most commonly to support image and keyboard/mouse input.

Other familiar ways to reduce query and formula entry effort are query suggestions and query autocompletion. Their helpfulness is related to the principle of *recognition over recall*: it is usually easier to recognize

---

[14]http://detexify.kirelabs.org/classify.html
[15]https://mathdeck.org

**Figure 1.6:** Tagent-V Formula Search Results (left) and Video Player Supporting Navigation (right)(Davila and Zanibbi, 2018). A rendered LaTeX formula is used to search handwritten math symbols recognized in a video (Davila and Zanibbi, 2017b). Here the user has clicked on the '=' of a matched formula on the whiteboard, and this advances the video to where it is first drawn (video: https://www.youtube.com/watch?v=gn24qo1MLN0).

something we see than describe the same thing from memory (Hearst, 2009). As a simple example, a query autocompletion might include a concept whose name but not formula we can remember, and allow us to quickly select a query containing both.

**Query response: communicating results.** To illustrate the communication of retrieval results, we'll use a system for visual search that uses an inverted index. Figure 1.6 shows handwriting in math lecture videos being queried with a LaTeX-generated formula image. The inverted index uses pairs of symbols (*e.g.*, $(I, n)$, $(=, A)$) as the vocabulary for lookup. Before searching, the query image is annotated with a graph containing nodes for symbols and edges with angles between adjacent symbols.

The inverted index is queried by looking up *all* adjacent query symbol pairs, to find their occurrences in the video collection. Each entry in the *posting list* for a queried symbol pair (*e.g.*, $(I, n)$) refers to an edge connecting the same symbols drawn in a video. Before indexing, videos are annotated with keyframes of drawn symbols that overlap in time along with an adjacency graph for each keyframe. Each edge in a keyframe graph is added to the posting list for its pair of symbols, as a *posting* containing a unique identifier for the edge, its keyframe graph, and video. In Figure 1.6, zoomed-in video keyframe graphs are shown

in results at left, and keyframe thumbnails are shown at far right.

It is actually the drawn symbol keyframes annotated on videos that are searched. Keyframes are scored by the similarity of matched adjacency subgraphs to the query graph, based on the similarity of matched symbols (nodes) and their angles (edges). In the results shown in Figure 1.6, symbols matched in a query/keyframe have the same color, and matched graph edges are red. To avoid missing symbols due to recognition errors, symbol pairs are indexed using all combinations of possible labels for each symbol. This is how $n$ matches $M$ in the second match shown.[16]

Let's consider the results in Figure 1.6 more closely, with associated system tasks in Figure 1.3. The query result shows the top-2 matching videos, and not keyframes. To generate this view, the keyframe ranking from *querying* the index is restructured as a video ranking, with videos ranked by best keyframe match. Also, the annotated query and videos have been *consulted* to produce the graph matches shown for each video. The videos are consulted again for annotations not used for indexing: clicking the mouse on a symbol in a result keyframe makes the video player jump to where the symbol starts to be drawn.

How the search results filter and present the videos is motivated by tasks users carry out (see Figure 1.4). For example, having symbols linked to frames can help people *consult* videos by quickly navigating to where a formula is drawn and discussed. Results rank videos rather than keyframes to make the search results more concise and easier to consult. The *communication*, *annotation*, *indexing*, and *querying* tasks can also be supported from search results. In MathDeck formulas in results can be used directly for search, selected for editing or export, or annotated in a card with a title and description. Cards are also automatically indexed in a 'deck' searchable by formula or title.

Showing matching graphs in results is more helpful for designers than users; simple bolding or highlighting is more common. In contrast, MathDeck's search results highlight matched query words and formulas

---

[16]A variation of adding tokens to queries and documents to increase possible matches in an inverted index. Symbol similarity is computed from all label probabilities assigned to each symbol. Tangent-v has also been used to search formula collections using unique symbol labels in PDF (Davila *et al.*, 2019).

located in PDF documents (*e.g.,* for papers from the ACL Anthology (Amador *et al.*, 2023)).[17]

We've used just two systems here to illustrate search results that rank sources, and how they interact with human information tasks. However, results from other systems have different types. Some systems plot, simplify, and/or perform requested operations on formulas, or provide solutions for math problems posed in text and/or formulas directly (*e.g.,* using Wolfram Alpha or a math-aware chatbot). In these cases the response is an answer to a (possibly inferred) question, rather than a ranked list of sources. These are designated as *question answering* systems, and interactions with chatbots addressing math queries are a type of *conversational search* where clarifying questions and additional information may be provided or received in multiple rounds of query/result interactions.

Chatbots using Large Language Models (*LLMs*) have proven intriguing and useful in some instances, but there is an increasing awareness of issues related to the validity of responses and other substantive concerns (Bender *et al.*, 2021). However, any retrieval result is only an information source – understanding and verifying any source requires additional work. Related to this, in Community Question Answering platforms (CQAs), many posts request clarification of a question, or clarify/correct posted answers and comments.[18] This illustrates how human responses to math queries also often contain misunderstandings, ambiguities and errors.

Regardless of the result type, how information is chosen and presented in results is important. It has a real impact on the usefulness of the result as a source of information, and on how tasks other than consulting the result itself are supported. In many cases usability testing can be used to check the effectiveness of result presentations and other interface design elements, and to discover refinements and alternatives.[19]

---

[17]ACL Anthology search demonstration: https://drive.google.com/file/d/1fbiMyHtlfEYUJvmrbZsWzhfL0X_zo9-t/view

[18]Including the classic, 'my correction to @your correction.'

[19]*e.g.,* parts of the MathDeck were usability tested (Dmello, 2019a; Nishizawa, 2020; Diaz, 2021), which led to substantial improvements.

**Figure 1.7:** ScholarPhi system showing definition for math symbols found within the same PDF paper (Head *et al.*, 2021). To assist skimming for details, text other than for definitions of a selected formula is greyed out. Source: https://www.youtube.com/watch?v=yYcQf-Yq8B0

**Supporting tasks for individual sources.**  Programs used to view/consult sources can also help with the user tasks illustrated in Figure 1.4. A nice example is the ScholarPhi system shown in Figure 1.7 (Head *et al.*, 2021). Reading formulas can be challenging, as symbols may be defined throughout a paper. ScholarPhi provides annotations decorating a selected formula/subexpression, providing symbol definitions in-place. Definitions are linked to where they appear, and text not associated with a selection is greyed out.

To produce the definition views in ScholarPhi, sources need to be annotated with formulas, symbols, and definition locations, and then definitions need to be linked with associated entities where they appear in the paper (*i.e.,* symbols or subexpressions). Definition segmentation and linking entities are performed with natural language processing techniques. The original prototype identified math symbols with LaTeX source files used to generate PDFs, simplifying formula detection.

A second example is the keyframe list at right in Figure 1.6. When viewing a video, all keyframes for handwritten content are available in a thumbnail list. Keyframes can be selected, and individual symbols clicked on to jump to where it is first drawn in the video (similar to the search results). This requires annotating video sources with generated

keyframes produced using computer vision techniques.

Both ScholarPhi and Tangent-v require generating additional information using automated inference (*i.e.,* AI), and their usefulness is limited by the accuracy and scalability of the methods employed. However, we believe that this is an important future direction for mathematical information retrieval, because the content and organization of mathematical sources can be complex. Particularly for non-expert users, mature versions of these techniques may be very helpful.

A number of well-known formats were devised or augmented to support detailed annotation with links and tags, including TIFF, PDF, and XML. Unfortunately detailed 'semantic' annotation has proven difficult at scale despite significant efforts. Some possible reasons include the time required to create sources before annotations, the diversity of information needs (*e.g.,* which information do we annotate?), attaching large annotations makes files large and unwieldy, and overall progress in scalable AI has been slower than many anticipated.

As AI continues to improve, creating source annotations to support examining and navigating math sources and other user information tasks within UIs seems likely to be beneficial. Perhaps application-specific annotations such as used in ScholarPhi, Tangent-v, and MathDeck are a good starting point.

## 1.6   Summary

This chapter provided a brief overview of mathematical information retrieval. A simple 'source jar' framework capturing how people's math information needs interact with information sources was presented. We showed how retrieval strategies can be sketched within a larger context using the framework, and how the information retrieval, analysis, and synthesis tasks in the model are performed by both people and systems, with some important differences. The chapter concludes with a characterization of challenges that people encounter interacting with retrieval systems, and examples of user interface elements designed to address those challenges for math-aware search.

In the next chapter we consider sources themselves more closely – in particular, how formulas are represented in sources.

# 2

## Representing Formulas and Text

In this chapter we focus our attention on how formulas and text are represented in sources and search indices used by retrieval systems. Figure 2.1 shows the tasks from the system overview in Figure 1.3 that we will discuss. These tasks include consulting sources for information in formulas and text, annotating formulas and text with additional information, and then constructing and organizing patterns for search in the index.

For context the query annotation and querying tasks are shown in Figure 2.1, but greyed-out. Annotations added to sources are often also added to queries, so that the correct pattern type is used to search the index. For example, for a dense retrieval model for formulas, both formulas in sources and queries are converted to vectors before identifying the formulas that are most similar to the query formula.

To help establish the roles of formulas and text and their relationships within sources, let's consider a simple but slippery question. What does a formula tell us, exactly?

**Figure 2.1:** Tasks from Figure 1.3 for Representing Formulas and Text in an Index. For context, query annotation and querying the index are shown greyed-out.

## 2.1    What information does a formula carry?

Let's start addressing our question using the *Inverse Document Frequency (IDF)*. IDF is used directly and indirectly in influential sparse retrieval models, including variants of TF-IDF (*Term Frequency*-Inverse Document Frequency) and BM25 (Robertson and Walker, 1994; Robertson and Zaragoza, 2009). IDF is used to score query terms that appear in sources. Its utility comes from a simple but powerful idea: query terms appearing in fewer documents are rarer and thus *more specific*, and so should be given higher weight when ranking sources. For example, the term 'BM25' is predominantly found in sources on information retrieval, while the term 'weight' is used for many topics and in multiple senses, including the heaviness of an object and computing scalar factors for numeric values. When scoring terms from a query containing 'BM25' and 'weight', IDF produces a higher score for 'BM25' reflecting its narrower usage (*i.e., specialization* in associated topics).
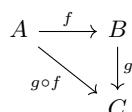
In Example 2 we see a definition for IDF, below which is the defi-

nition's structure as a sequence of words and formulas. Formulas have visual structure representable using a Symbol Layout Tree (SLT) as shown for the *idf* formula. The SLT represents the placement of symbols on writing lines using the spatial relationships shown in Table 2.1.[1] In the example adjacent symbols are shown using horizontal lines, and other relationships including *subscript* and *below* with directed arrows. This SLT represents the *idf* and *log* function names in single nodes, with their characters grouped into *tokens*. SLTs for variables $N$, $t_i$, and $n_i$ are subtrees, with one node for $N$, and two nodes for $t_i$ and $n_i$. Table 2.1 also includes *containers* that represent operators and their associated writing lines compactly. We will encounter these later in the chapter.

Taken together, the excerpt is representable as a directed graph with formulas in SLTs connected to the formula nodes.[2] Given this, a reader tries to recover the represented information. An important part of this process is identifying unstated or assumed information. These gaps are filled using information known to the reader, or available elsewhere in the same source or additional sources. As a simple example, a person identifies unstated information when using context to distinguish 'sin $hx$' (sine of angle $hx$) from 'sinh $x$' (hyperbolic sine of $x$), perhaps because of a typo. How operations and arguments are grouped needs to be correctly identified for a correct SLT representation.

In Example 2 one piece of missing information is the hierarchy of operations and arguments represented in the *idf* formula. While a reader is unlikely to think about this consciously, interpreting the formula essentially involves converting the SLT to an Operator Tree (OPT). In OPTs variables and other arguments appear at the leaves, with operations above the leaves in internal nodes. While SLTs are

---

[1]Diagrams and other graphics are frequently used in math, but outside of our discussion here; *e.g.,* commutative diagrams can be expressed as a matrix-like SLT container, but are really directed graphs with nodes/edges labeled by formulas:

$$A \xrightarrow{f} B$$
$$g \circ f \searrow \quad \downarrow g$$
$$C$$

[2]SLTs are hidden in the Example 2 middle panel for easier reading.

Example 2: Inverse Document Frequency (*IDF*)

**Excerpt from (Robertson, 2004):**

... Assume there are $N$ documents in the collection, and that term $t_i$ occurs in $n_i$ of them ... the measure proposed by Sparck Jones, as a weight to be applied to term $t_i$, is essentially

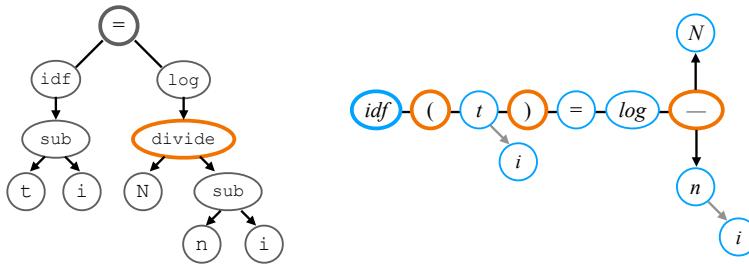$$idf(t_i) = \log \frac{N}{n_i} \tag{1}$$

...

**Word and formula sequence:**

Assume there are $N$ documents in the collection, and that term $t_i$ occurs in $n_i$

$$idf(t_i) = \log \frac{N}{n_i} \tag{1}$$

**Equation (1) operations and appearance:**

**Left:** Operator tree (OPT) with hierarchy of operations and arguments.
**Right:** Symbol layout tree (SLT) placing symbols on writing lines.

oriented left-right to reflect reading order, OPTs are oriented vertically
to reflect operation order. The order of operations is bottom-up in the
OPT, with precedence decreasing as we move away from the leaves,
*e.g.,* '=' is applied last for the *idf* example.

Operations appear directly above their arguments in an operator
tree. If an operator's arguments have different roles (*e.g.,* $N$ and $n_i$ in
$\frac{N}{n_i}$) they are represented in a fixed left-right order below the operator.
When argument order does not affect the result such as for '=' or '+',
a left-right order is chosen arbitrarily, usually reflecting the reading
order in the SLT. In the example OPT, ordered operator arguments
are shown with arrows, and the unordered arguments for '=' are shown
with lines. A `sub` operator is used to represent subscripted variable
names within the SLT. Orange nodes indicate where SLT symbol nodes
are renamed or removed in the OPT. This includes replacing a fraction
line by `divide` in the OPT, and removing parentheses because OPT
argument edges make them redundant.

For search, it is worth considering what is fixed and what varies in
OPTs and SLTs:

- In SLTs, spatial relationships are limited to those in Table 2.1,
  while symbols appearing on writing lines can vary greatly.
- In OPTs, an operation hirearchy with arguments at the leaves is
  defined. Like SLTs, operation and argument symbols vary but some
  operations are implied and missing in SLTs, *e.g.,* $xy$ represents
  $x \times y$.

We can improve search by reducing variation in OPTs and SLTs, partic-
ularly when one formula has multiple representations. We will discuss
this further in Section 2.2.

In addition to the underlying operations represented in formulas,
there are some other types of information that formulas provide. For-
mulas represent or *imply* mathematical properties, and provide visual
information in their presentation. Also, like all parts of language their
uniqueness is another form of information, as the IDF formula suggests.

**Information from formulas and text.**   Despite its brevity, the excerpt
in Example 2 contains a fair amount of represented and implied infor-
mation. But what information is represented separately by formulas and

**Table 2.1:** SLT Spatial Relationships (8 total) and Containers. Relationships identify writing lines around a symbol/token. 'Adjacent at left' is excluded because of reading order: writing lines are read left-to-right.

| Relationship | Examples |
|---|---|
| Inside a symbol | $\sqrt{N}$ |
| Adjacent on writing line | $idf(t) = \log \frac{\phantom{}}{\phantom{}}$ |
| *(at right of symbol)* | |
| Subscript & Superscript | $t_i \quad N^2 \quad \int_{-\infty}^{+\infty} \quad \sum_{i=1}^{n}$ |
| *(diagonally at right above/below symbol)* | |
| *Prefix Superscript & Subscript | $_2^n C \quad {}^{235}U$ |
| *(diagonally at left above/below symbol)* | |
| Above & Below | $\frac{N}{n_i} \quad \binom{n}{2} \quad \lim_{n\to\infty} \sum_{i=1}^{n}$ |
| *(directly above/below symbol)* | |

| **Containers: *combine relationships*** | **Examples** |
|---|---|
| *Fraction* | $\frac{N}{n_i}$ |
| (above + below fraction line) | |
| *nth-Root* | $\sqrt[3]{8}$ |
| (above + inside radical symbol) | |
| *Matrix* and *Choice* | $\begin{bmatrix} 1 & 0 \\ 0 & x^2 \end{bmatrix} \quad \binom{n}{2}$ |
| *(adjacent for brackets/rows, below for cols)* | |
| *Note: nested SLTs within cells* | |

*Aside from specific communities and topics (*e.g.,* chemistry), these are rare.

text, and what information do they present together? The following is a structured summary of information in the excerpt, including uncertain details we'll come back to later.

*Variables and Functions*

- Variable $N$ refers to a quantity, similar to a noun. $N$ acts like a common noun (*e.g., street*) rather than proper noun (*e.g., Bath Road*) because the collection is unspecified.
- Variables $t_i$ and $n_i$ are also common noun-like, with by a shared reference to $i$ indicating *any* term in a vocabulary. For example, $(t_i, n_i)$ could be (‘IDF’, 7) or (‘weight’, 11).
- log represents the set of logarithmic functions with unspecified base, and so variable/common noun-like.
- $idf(t_i)$ is a *specific* function applied to term $t_i$ making *idf* constant and akin to a proper noun.

- Definitions
  - Variables in text with implied types: natural 'counting' numbers ($\mathbb{N}$) for $N$ and $n_i$, character strings for $t_i$.
  - *idf* function in Equation (1).
  - The full definition includes textual variable definitions and the *idf* function definition in Equation (1).

*Operations*

- Division ($-\div-$), application ( $idf(\cdot)$ , $\log\cdot$), and equivalence ($=$) act on variables and subexpressions (like verbs).
- Definitions
  - *Missing* in text and formulas of the passage.

*Formulas in Text*

- Variable and function names act as single-word nouns ($N$ vs. *Sunday*) and noun phrases ($t_i$ vs. *'any given Sunday'*).
- Equation (1) acts as a sentence clause with a subject ($idf$), subject verb ($=$), and predicate (the definition).

*Text Only*

- Provides the purpose of *idf*: defining weights for terms.
- Provides context: formula is similar to its original proposal by Spärck Jones (Jones, 1972).

*\*Unspecified Details*

- Base of the log function
- Specific collection that $N$ and $n_i$ refer to
- Specific vocabulary (*i.e.,* index terms/entries) $t_i$ refers to
- Operator definitions

Notice the different communication roles that both words and formulas take in this excerpt. This includes naming/reference, actions and properties, definitions, and discourse:

1. Nouns, variables, and function names refer to objects.
2. Verbs and operators give actions on objects and object properties.
3. Definitions are provided in text and formulas, with text adding information to formulas and vice-versa.
4. Discourse: *idf* is referenced in the larger paper discussion, and the text provides additional information on the formula's origin.

The excerpt also provides an example for our discussion here.

Let's next consider an alternative definition for the *idf* function that uses no mathematical symbols or names at all:

> The *inverse document frequency* for a term is defined as the number of collection documents divided by the number of documents containing the term, which is then converted to a logarithmic value.

This seems simple enough. But we lose some useful things when we remove the math notation, as summarized below.

*Compact Reference:* Variable and function name reference is more efficient than reusing descriptions, *e.g.,* referring to $N$ vs. 'the number of collection documents.'

*Visibility:* Formulas are italicized and use distinct symbols (*e.g.,* operators and greek letters), and in addition to appearing in sentences (*inline*) they may ve indented and offset (*displayed*). Equation (1) in Example 2 is displayed.

*Compact Structure:* Formulas define relationships more compactly than text (*e.g.,* the *idf* formula vs. the description above). This also helps visualize mathematical concepts. In the *idf* formula, the fraction visualizes a 'flipped' (inverse) percentage of sources. As another example, the *distributive property* from algebra is easily expressed as $x(y + z) = xy + xz$.

*Abstraction and Generality:* Formulas identify relationships that appear in many different contexts, *e.g., idf* can be applied to formulas as well as text. Abstractions in formulas help identify, define, and name these commonalities, simplifying analysis and discussion. For example, formulas can be repurposed by redefining variables and redefining operations.

- Example 3 provides an example of applying a formula in different contexts just by redefining variables.
- We can redefine variables and operations to alter the distributive property above for vectors. To do this, multiplication ($\times$) is replaced by $\cdot$ (scalar product) and $+$ is replaced by

vector addition. To distinguish vectors from individual values, we'll also make the variable names bold. This gives: $\mathbf{x} \cdot (\mathbf{y} + \mathbf{z}) = \mathbf{x} \cdot \mathbf{y} + \mathbf{x} \cdot \mathbf{z}$.

**Missing Information.**    Having seen the helpful abstractions that math formulas provide, let's further discuss the unspecified details for the excerpt in Example 2. It turns out that these are all *deliberate* omissions, to support abstraction in the formulas and for clearer communication:
- By not specifying a collection, $N$ is defined for any collection.
- $t_i$ and $n_i$ refers to any term and the number of documents it appears in for any specific term, collection, and vocabulary.
- Logarithm values increase with input values (*i.e.,* they are *monotonic*) so any base suffices. Omitting the base emphasizes this.
- The operators used are common in computing, and so they are assumed rather than given to save space and reader effort.

Omitting information from formulas is useful for generalization and brevity, provided that the reader is clear about what is omitted and why. The value of omission holds for text as well. When we discuss a person, place or thing, we skip shared information and unrelated details; we can find it frustrating or amusing when others forget to.[3] Choosing what to omit in formulas and text is informed by:
1. The context and focus of discussion, *i.e.,* items discussed earlier and current topic
2. The background of the audience

These two factors largely dictate what can and *should* be left out.

**Rareness and topic specificity.**    The IDF formula is related to *entropy* from information theory (Shannon, 1948), but is not fully compatible with it due the different probabilities involved (Robertson, 2004). Entropy measures uncertainty for events with discrete outcomes such as flipping coins or the next word or formula that will appear in a passage. It quantifies how surprising an outcome is, and equivalently how much probabilistic *information* knowing the outcome provides.

---

[3] *e.g.,* 'oversharing', like telling a stranger your favorite color while booking a dentist appointment.

---

**Example 3: A Decay Function**

This is an exponential decay function for quantity $N$ after a number of timesteps $t$:

$$N(t) = N_0 e^{-\lambda t} \tag{2}$$

$N(t)$ is the quantity at time $t$, $N_0$ the initial quantity (when $t = 0$), and $\lambda$ is a decay constant. This same formula can be used in multiple contexts, simply by changing $N_0$ and $\lambda$. For example,

**Financial analyst:** retirement fund balance after $t$ months assuming monthly payouts and a fixed interest rate.
**Chemical engineer:** rate of a chemical reaction.
**Physicist:** discharge of the potential contained in a capacitor.

---

For example, we are less surprised when someone correctly predicts a coin toss than guesses the hidden top card in a shuffled deck of 52 cards. If the deck is shuffled 10 times, and after getting the first card wrong, the correct card is guessed correctly 9 times, we rightly suspect cheating because of how unlikely (surprising) this is.[4]

The entropy $H$ for event $X$ with possible outcomes $\mathcal{X}$ is given by

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x) \tag{2.1}$$

where $\log p(x)$ is the log probability of one possible outcome. These outcome values are combined in a weighted average using the likelihood for each alternative ($p(x)$). Because the maximum probability is 1.0, all outcome probabilities are smaller than standard base values such as 2 or $e$, producing negative logarithm values, and a positive entropy value after the final sum and negation. The maximum entropy for $n$ alternatives is obtained by a purely random (*i.e.*, *uniform*) distribution when outcomes are equally likely with $p(x) = 1/n$, such as (fairly) guessing the top card in a shuffled deck.

Rarer terms are more *informative* in this probabilistic sense. The

---

[4]With proper shuffling, 9 correct guesses in a row has an estimated likelihood of $\frac{1}{52}^9 = 3.39 \times 10^{-16}$. Here the initial error was made because the signal for identifying the card had not been perfected as two adults play the game to befuddle children.

most frequent words are articles, prepositions, pronouns/references, and other syntactic 'latticework' for text and their likelihood is dramatically higher than other words (*e.g.,* the, and, your etc.), and are often referred to as *stopwords.* Across languages, word frequencies sorted in decreasing order follow a Zipfian distribution, with an inverse exponential-shaped curve. Rarer words and phrases in the long 'tail' of the curve often signify specific properties or topics (*e.g.,* Zipf, inverse, 'information retrieval'), but can also arise from spelling mistakes, for example.

So another form of information that formulas carry are their uniqueness or rarity within a passage or collection. For example, the *idf* formula is more likely to be used in sources from computer science on information retrieval, natural language processing, and data mining than other topics. In contrast, the small formulas $x^2$ and $x$ will be found so frequently across disciplines that we can consider it a 'latticework' part of communication similar to *the* (*i.e.,* a 'stopformula').

## 2.2 Formula representations

SLTs and OPTs are helpful abstractions when working with the appearance and operations of formulas. We use the information representable in SLTs and OPTs whenever we work directly with mathematical definitions and operations, or indirectly with applications of math. For example, we do this when using calculators, playing games, reading formulas, talking with others about what a formula represents, or when implementing a formula in a computer program.

As discussed earlier, math formulas in documents help with abstraction by capturing properties and patterns found in multiple contexts. Redefining variables and operations in formulas allows their represented information to change without altering their appearance. For familiar symbols like $x$ or $+$, definitions are usually omitted for brevity.

When implementing an OPT in a computer program, there are additional constraints. In code variables and operations *must* be uniquely defined before they can be used. This limits a flexibly-defined OPT to exactly *one* interpretation. Definitions can be omitted from program code when they are provided by the language or imported libraries, but all symbols must be defined for a running program to work. Definitions

must be fixed, but can be replaced, *e.g.,* by importing a different library.

Representing SLTs in code is more straightforward than for OPTs. SLTs primarily define the placement of symbols on writing lines along with formatting information (*e.g.,* spacing, fonts and sizes). The LaTeX math syntax is an example of a commonly-used SLT representation. LaTeX supports defining symbols and layout structures in macros and libraries.

**Representing OPTs and SLTs in code.**   Example 4 provides computer code representing the OPTs and SLTs for the *idf* formula shown in Example 2 that demonstrate differences between representations for OPTs and SLTs in computer programs.

Syntax for OPTs and SLTs can both be similar to formula appearance on a page as seen for the Python and LaTeX excerpts in Example 4. However, operator *prefix* syntax showing operators before arguments match OPT structure directly, as seen in the Lisp examples. To compute values from an OPT in code, we require unique definitions for operators and arguments, which require decisions about data structures and algorithms as discussed below for the implementation of subscripts.

SLT code representations are simpler than OPT code representations. They identify where and how to draw symbols, without any concern for what the symbols represent. For example, in LaTeX commands such as \sum, \choose, and \frac helpfully suggest math operations to make the syntax easier, but the commands only define where to place symbols (*i.e.,* these are containers as shown in Table 2.1).

Let's now consider the code examples in more detail, starting with the OPT implementations. While our example OPT can be interpreted as defining an *idf* function, for simplicity we'll use a slightly different interpretation for testing function values. In the code we assume that:

- variables N, n, t, i, n_i and t_i have defined values,
- *idf* is defined as a function in code, and
- = (Lisp) and == (Python) test for equal values.

The Lisp and Python OPT implementations test if the *idf* function produces the value obtained by applying operations in the definition.

In the top row Lisp represents the OPT using indentation to visualize levels in the tree. All Lisp operations are grouped with their arguments

---

Example 4: *idf* Formula in Code (OPT & SLT)

**Lisp: prefix operations**

```
( =
    ( idf
        ( sub t i ))
    ( log
        ( / N ( sub n i ))))
```

**Lisp + subscripts in names**

```
( =
    ( idf t_i )
    ( log
        ( / N n_i ))
```

**Python:** infix operations + subscripts in names

```
idf( t_i ) = log( N / n_i )
```

**LaTeX:** *appearance as SLT (writing lines)*

```
idf( t_i ) = \log \frac{ N }{ n_i }
```

**OPTs**: Lisp and Python, **SLT**: LaTeX

---

using bracketed expressions, while individual variables appear without brackets. A *prefix* syntax is used, putting operations before arguments in groups. The =, / and **sub** operators have the form

> ( op arg1 arg2 )

while the log and *idf* functions have the form

> ( op arg1 ).

The first Lisp expression represents the *idf* OPT using the same symbols other than replacing **divide** by the / Lisp operator. The operation node hierarchy is represented by nesting arguments in brackets.

Choosing how to represent the subscripts is surprisingly subtle. In the top-left Lisp expression subscripts are defined using the **sub** operator in the OPT. However, **sub** and its arguments need specific definitions if this is to run in a program. Let's start with:

- **sub** is a standard array lookup,
- **t** is an array of strings for vocabulary terms,
- **n** is an array of counts for terms (*same* order as **t**), and
- **i** is an index position for a term.

If i = 2, ( sub t i ) might return 'weight', ( sub n i ) returns 15 as the number of documents where 'weight' appears, and the **idf** function receives a string argument as ( idf "weight"). This **idf** function *must* refer to data not shown in the OPT, *e.g.,* pre-computed values in a dictionary, or a dictionary of (term, count) pairs used to compute IDF values in the function.

Now consider the top-right Lisp excerpt, where subscripts are represented in variable names. This makes the code shorter, but also less general. The code can compute values for exactly one term, rather than use the index $i$ to access any vocabulary term. Depending on other code in the program, one or the other definition may be more helpful.

An equivalent program in Python syntax is shown that includes *infix* operators, with arguments to the left and right of $==$ and the division operator (/).[5] When operations appear before arguments, it is outside at left of a bracketed list. The Python code certainly *looks* closer to the original formula than the Lisp. Small differences in appearance between the original formula and Python code include the fraction being written left-right, and adding required parentheses for the log operation.

The final example is LATEX representing an SLT for the appearance of the formula. LATEX commands begin with a backslash, and the fraction is represented using the \frac{}{} command. Only the parentheses ( '(' and ')') that actually appear in the drawn formula are included in the LATEX. Subscripts are represented with operators to capture writing line changes, with '_' being the LATEX subscript operator rather than part of variable names. There are commands in LATEX that modify the interpretation of '_' and the superscript operator '^'. For example, the \displaystyle command changes argument positions from sub/superscript to above and below for $\sum_{i=1}^{n}$ in Table 2.1.

**Canonicalization and MathML.** Example 5 visualizes some common ways to identify properties of variables and operations in OPTs. At middle we number each unique variable from left-to-right in the OPT, starting from 1. This type of numbering for entities is known as an *enumeration*. '2' for $i$ is repeated because it appears twice in the expression.

We can use enumerations to capture variable placement while ignoring specific variable names. For example, using variable enumeration the *pythagorean theorem* expressed as
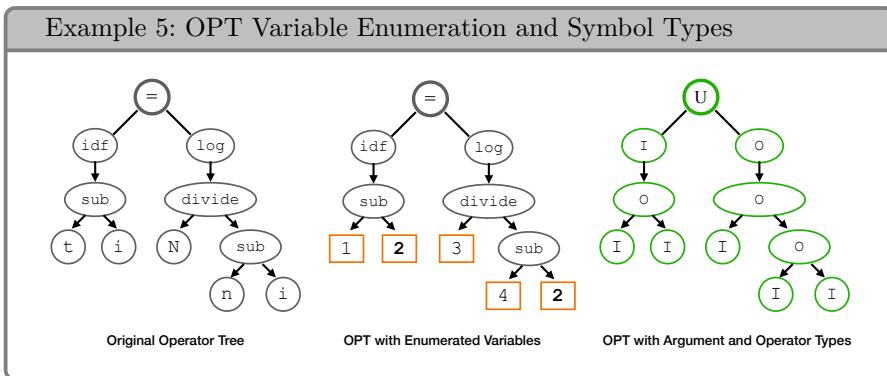
$$x^2 + y^2 = z^2 \tag{2.2}$$

and

$$a^2 + b^2 = c^2 \tag{2.3}$$

---

[5]Technically, `log` should be `math.log`.

Example 5: OPT Variable Enumeration and Symbol Types

Original Operator Tree     OPT with Enumerated Variables     OPT with Argument and Operator Types

both have the form

$$\boxed{1}^2 + \boxed{2}^2 = \boxed{3}^2. \tag{2.4}$$

Enumerations are helpful for *unifying* the variables, making the two versions of the pythagorean theorem identical. We do this by substituting variables (*e.g.,* $\{x/a,\ b/y,\ c/z\}$) or representing variables with enumeration numbers instead of variable names. Unification is useful in many contexts, including theorem proving and formula search. It is possible to enumerate operations in the internal nodes as well, but this is less common for retrieval applications.

Variable enumeration removes distinctions between formulas with identical operations but different variables. Applying transformations like this that generalize formulas and reduce the number of unique formulas is known as *canonicalization.* Canonicalization often involves *normalizing* symbols and structures as well. A common normalization is re-ordering variable names alphabetically for operations where the order of arguments is unimportant, *e.g.,* to have both $x + y$ and $y + x$ represented by $x + y$.

Both OPTs and SLTs may be canonicalized using enumeration, normalizations, and other transformations. Without canonicalization it is harder to identify the same information in different formulas. However, using too much canonicalization can remove meaningful differences between formulas. A common compromise is to use multiple formula representations. For example, we might search OPTs using enumerated variables, and then refer back to the original OPTs to rank those

containing query variable names higher.

Types are also often used for annotation and canonicalization, as shown at right in Example 5. Here we have an OPT with nodes replaced by an assigned *type* for each symbol. *I* indicates identifiers (*i.e.,* names) for variables and operations/functions like the `idf` function identifier. Note that the identifier `idf` must represent an operation: it appears at an internal node and not a leaf in the OPT. In our example, other predefined mathematical operations are assigned one of two types: *O* for operations with ordered arguments (`sub`, `log`, `divide`) and *U* for operations with unordered arguments (`=`). Among other uses, types can be used as constraints when unifying symbols including operators with the same argument structure in formulas.

Let's now talk about file encodings that are used to represent SLTs and OPTs, along with annotations including types. There are many formula file encodings available, but here we focus on the XML-based Mathematical Markup Language (MathML[6]) which has been used in a number of evaluation benchmarks for math-aware search. Being XML-based, the syntax for representing formulas is much closer in structure to the prefix syntax of Lisp than the Python syntax with infix operations in Example 4. All MathML commands have a start and end tag, each containing a list of tags. The syntax is roughly

```
<tag-command [attribute list]>
    <argument-1>
        ...
    </argument-1>
    ...
    <argument-n>
        ...
    </argument-n>
</tag-command>
```
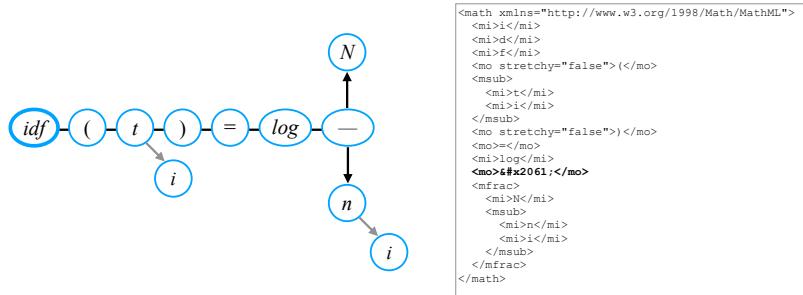
This is used both for SLT representation in *Presentation MathML*, and OPT representation in *Content MathML*. For single symbols tags can be combined with the symbol itelf, *e.g.,* for the `<log\>` operator in Content MathML.

---

[6]https://www.w3.org/Math

Example 6: MathML for IDF Formula (LaTeXML[a] Output)

The MathML shown was produced by LaTeXML[a] using the LaTeX in Example 4. *idf* is not defined; $i$, $d$, and $f$ are treated as variables.

**Symbol Layout Tree in *Presentation MathML***



```
<math xmlns="http://www.w3.org/1998/Math/MathML">
  <mi>i</mi>
  <mi>d</mi>
  <mi>f</mi>
  <mo stretchy="false">(</mo>
  <msub>
    <mi>t</mi>
    <mi>i</mi>
  </msub>
  <mo stretchy="false">)</mo>
  <mo>=</mo>
  <mi>log</mi>
  <mo>&#x2061;</mo>
  <mfrac>
    <mi>N</mi>
    <msub>
      <mi>n</mi>
      <mi>i</mi>
    </msub>
  </mfrac>
</math>
```

Other than adjacency, relationships are represented in tags (*e.g.*, `<msub>` for subscript, `<frac>` for fraction). Unicode x2061 is inserted by LaTeXML representing the *function application* of log to $N/N_i$.

**Operator Tree in *Content MathML***



```
<math xmlns="http://www.w3.org/1998/Math/MathML">
<apply>
  <eq/>
  <apply>
    <times/>
    <ci>i</ci>
    <ci>d</ci>
    <ci>f</ci>
    <msub>
      <ci>t</ci>
      <ci>i</ci>
    <msub>
  </apply>
  <apply>
    <log/>
    <apply>
      <divide/>
      <ci>N</ci>
      <msub>
        <ci>n</ci>
        <ci>i</ci>
      </msub>
    </apply>
  </apply>
</apply>
</math>
```

Nested `<apply>` tags represent the nesting of operations in the OPT. Each `<apply>` tag contains a function (*e.g.*, log) or relation (*e.g.*, =) and its arguments. With *idf* interpreted as three variables, a multiplication replaces the application of *idf* to $t_i$.

_____

[a]https://math.nist.gov/~BMiller/LaTeXML

Some MathML commands, such as `<msub>` for subscripts in OPTs and SLTs have a fixed number of arguments. Others like `<apply>` in Content MathML and `<math>` or `<mrow>` representing writing lines in Presentation MathML have variable length argument lists. The `apply` operation is close to the bracketed operation groups in Lisp: the first argument is an operation followed by its arguments.

MathML also defines types for arguments, including `<mi>` and `<ci>` for variable identifiers, `<mn>` and `<cn>` for numbers, and `<mi>` and `<mo>` for operators in SLTs/Presentation MathML (*e.g.,* `<mi>log</mi>`). Defined operations in OPTs/Content MathML have their own predefined tags, and so `log` appears as `<log/>` and not `<ci>log</ci>`. The LaTeXML tool used to produce Example 6 is aware that `\log` is an operator, and inserts an *invisible* node in the Presentation MathML to capture the application to the fraction, using the Unicode value `x2061` in hexadecimal. This symbol does not appear when this formula is rendered (*e.g.,* in a web page by MathJax[7]).

*idf* was not defined as an operation when translating the LaTeX in Example 4 to MathML. As a result this is broken up into three adjacent variables in the SLT, $i$, $d$, and $f$. In both the Presentation (SLT) and Content (OPT) MathML excerpts, these are treated as variable identifiers. Because of this change, in the Content MathML the letters are multiplied with each other and $t_i$.

Because a fixed set of definitions is required to convert formulas, for large collections inconsistencies such as those seen in Example 4 are common. There is a Content MathML `<cerror>` tag for unrecognized symbols or structures seen frequently in large-scale conversions as well.[8] Fortunately, for search if these interpretations not intended by their authors and 'errors' are consistent, they still provide patterns that are useful in search, provided that both formulas in queries and in sources are converted in the same way.

Getting back to canonicalization, we often restructure automatically generated OPTs and SLTs in MathML before indexing to simplify the tag structures. Common examples include removing `<cerror>` tags, and

---

[7]https://www.mathjax.org/

[8]Of course, some formulas are actually incorrect.

replacing nested `<mrow>` tags by a single `<mrow>` or top-level `<math>` tag in Presentation MathML, *e.g.,* to avoid:
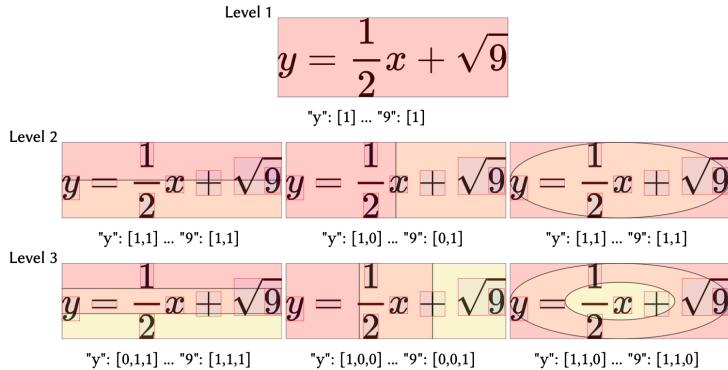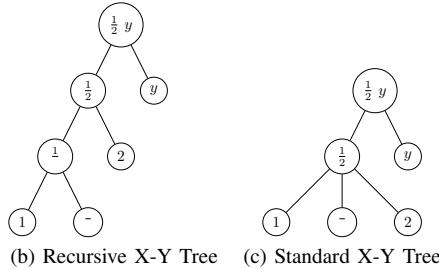
```
<mrow>i <mrow>d <mrow>f ... </mrow></mrow></mrow>
```

Similarly, for Content MathML it is also common to 'flatten' some operation sequences to a single operation node. For example, we can flatten multiplications to a single multiplication (represented by `times`) without changing the original meaning, because multiplication is an unordered operation (*e.g.,* $a \times b \times c = c \times b \times a$). Both of these transformations for `<mrow>` and `<times>` are applied in Example 6, flattening a chain of `<times>` nodes into one node in the Content MathML (OPT), and removing `<mrow>` tags entirely from the Presentation MathML (SLT).

**Other formula representations.** SLTs and OPTs are probably the most natural representations for formulas, giving appearance by symbol placement and the represented operation hierarchy. However, in many source types including images, videos, and PDF files, the location of formulas are unavailable, as well as SLTs or OPTs. In PDF formulas are usually represented in embedded pixel images (*i.e.,* raster images) or PDF drawing instructions for characters and lines (*i.e.,* vector graphics, such as in SVG or PDF-rendered LaTeX).

For these source types, we need to annotate sources with formula locations and/or create an isolated formula collection. Once formulas have been located or placed in individual images, we need to either (1) annotate formulas with SLTs/OPTs using recognition systems, (2) embed formula images directly as vectors for dense retrieval, or (3) create visual representations that can be computed for raster (*e.g.,* PNG) or vector (*e.g.,* PDF) images. Image embeddings and visual formula representations provide patterns distinct from SLTs. In this regard they can be viewed as complementary rather than strictly as substitutes or replacements for SLTs.

Both formula detection and recognition are active research topics in the document recognition community, and image embedding and image retrieval continue to be explored extensively in the information retrieval, computer vision, medical computing, and multimedia research

---

Example 7: Region-Based Spatial Formula Representations

**Pyramidal Histogram of Characters (PHOC)**

Level 1

$$y = \frac{1}{2}x + \sqrt{9}$$

"y": [1] ... "9": [1]

Level 2

$$y = \frac{1}{2}x + \sqrt{9} \quad y = \frac{1}{2}x + \sqrt{9} \quad y = \frac{1}{2}x + \sqrt{9}$$

"y": [1,1] ... "9": [1,1]        "y": [1,0] ... "9": [0,1]        "y": [1,1] ... "9": [1,1]

Level 3

$$y = \frac{1}{2}x + \sqrt{9} \quad y = \frac{1}{2}x + \sqrt{9} \quad y = \frac{1}{2}x + \sqrt{9}$$

"y": [0,1,1] ... "9": [1,1,1]        "y": [1,0,0] ... "9": [0,0,1]        "y": [1,1,0] ... "9": [1,1,0]

**XY-Cut Trees**



(b) Recursive X-Y Tree    (c) Standard X-Y Tree

---

literatures. Here we will consider two examples of visual-spatial representations that have been used as alternatives to SLTs for formula search; others are certainly possible, and an interesting direction for future work.

For vector images with explicit symbol locations (*e.g.,* in SVG or PDF) or where isolated symbols in images have been recognized, we can use a *spatial* formula representation. We saw an example of this with the line-of-sight graphs over symbols used to search handwritten and typeset math in the previous chapter. An example of a region-based spatial represenation is the Pyramidal Histogram of Characters (PHOC), which represents symbols in a fixed set of recursively partitioned regions.

PHOC was originally created for retrieving words in handwritten text (Almazán *et al.*, 2014) but can be generalized in a straight-forward way for representing two-dimensional structures like formulas (Langsenkamp *et al.*, 2022; Avenoso *et al.*, 2021; Amador *et al.*, 2023). For a formula, the whole formula is the first region. Additional regions are defined by splitting this formula region into 2, 3, or more regions horizontally, vertically, and concentrically using ellipses or rectangles. A bit vector is created for each symbol, using a '1' to mark each region that the symbol appears in. Each formula is then represented by a very long vector that concatenates this bit vector for every symbol in the vocabulary. A much more compact representation defines vectors only for symbols appearing in a formula, as illustrated in Example 7.

For raster (pixel-based) images symbol locations are unknown, but we can still capture visual structure without OCR using an *XY-cut tree.* XY-cut trees partition touching pixel groups (connected components) by recursively cutting at vertical and horizontal gaps in an image, either in strict alternation or using the largest gap in either direction (Ha *et al.*, 1995; Nagy and Seth, 1984). The resulting graphs are related to SLTs: cutting vertically tends to partition the image at subexpressions along a writing line, and cutting horizontally tends to separate writing lines, *e.g.,* above and below a fraction. Symbols can be recognized or features computed from sub-images at nodes (Zanibbi and Yu, 2011). This was one of the earliest spatial representations used for OCR of math formulas (Okamoto and Miao, 1991).[9] XY-trees could also be used when symbols are known (*e.g.,* for PDF), as it is simply a region partitioning method.

There are also alternative symbolic representations that can be used for search, such as substitution indexing trees originally developed to support automated theorem proving (Graf, 1995b). Substitution indexing trees can be used to group OPTs or SLTs with shared structure using a hierarchy of symbol and operation replacements and enumerated variables (Kohlhase and Sucan, 2006; Schellenberg *et al.*, 2012). These trees represent a set of possible operation sequences that produce a

---

[9]Some additional steps are needed to remove radicals to access their subexpressions, and to identify where writing lines are subscripted or superscripted.

concrete formula at the leaves, with operations performing replacements top-down in an OPT.

Whatever formula representations are used, it is common to embed them in vectors for use in dense retrieval models. Generally speaking, formula embedding vectors are obtained by translating an image, text, or graph into a fixed-length vector. Formula images may be converted to vectors, many methods convert LaTeX and/or MathML token sequences to vectors, and some convert SLT and/or OPT graphs directly to vectors using Graph Neural Networks (GNN). For example, a graph convolutional neural network can be applied to SLTs in Presentation MathML generated from LaTeX formulas in arXiv papers (Pfahler and Morik, 2020).

## 2.3  Formula annotation and indexing

Figure 2.1 illustrates the main tasks for indexing sources for retrieval. When we talk about indexing, we're actually referring to a process that consults and annotates sources with additional information, and then creates a collection index. The collection index contains data structures used to search the collection of sources using *patterns* generated from sources and their annotations. More concretely, indexing involves:

1. consulting source text, formulas and sub-expressions and generating dictionaries for fast lookup and analysis (*i.e.,* source content annotations),
2. adding new information to sources using additional dictionaries, *e.g.,* formula locations, formula representations, connecting formulas to descriptions and variables to types (*i.e., additional* source annotations), and
3. generating inverted index and/or dense vector index files from source contents and additional annotations.

In this section we first discuss indexing formulas in isolation, and then how text annotations can be added and included in indexing. We'll then briefly talk about indexing sources as whole, including all formulas and text in the next section.

**Preliminary book-keeping: sources and formulas.** Whatever method we use for formula search, retrieval systems need to access (1) individual sources and (2) the location of formulas within sources. We need to detect formula locations in sources if they are not explicitly identified. For example, videos and PDF documents generally do not identify formula locations while LaTeX and MathML demarcate formulas using known symbols and tags.

We generally use integer identifiers for sources and formulas because of their fixed small size. These integer identifiers are faster to process than variable-length character strings, and their use avoids duplicating sources in the index. For example, given a collection of LaTeX documents we might produce two dictionaries:

1. **Source formula locations**
   `formula-int-id` → (`source-int-id, start-char-index,`
        `end-char-index`)
2. **Files for sources**
   `source-int-id` → `file-location`

To generate a search result containing LaTeX for a formula, we use three lookups. The first uses the formula id to obtain the source and excerpt containing the formula, the second retrieves the file location using the source id, and the third obtains LaTeX commands from the excerpt within the file after it is loaded. The dictionaries save space through *indirection* (*i.e.,* ids referring to other ids). For example, within the index the file names for sources *only* appear in the second dictionary – all other references to files are made using source ids.

Where formula locations are detected, a modified formula location dictionary is needed. An example are the Page-Region-Object tables used for ACL anthology PDFs in the MathDeck system (Amador *et al.*, 2023). Each detected formula is assigned an identifier and annotated with its source, page, and page location.

1. **\*Detected PDF source formula locations**
   `formula-int-id` → (`source-int-id, page-int-id,`
        `page-region-bb, obj-int-id-list`)

`page-region-bb` is two $(x, y)$ coordinates giving top-left and bottom-right corners for a page region (*i.e.,* a *bounding box*). A secondary dictionary holds the object type, value, and bounding box for each

`obj-int-id:`

> **\*Formula object index**
> `obj-int-id → (type, value, obj-region-bb)`

For raster images, formulas have only one object, the raster image. For formulas containing drawing instructions, there may be multiple objects and types (*e.g.,* character/line/image) and values (*e.g.,* character/line endpoints/raster matrix).

So far our formula dictionaries represent only the locations of formulas in sources. This means that formula representations must be retrieved from files. It is convenient to add additional dictionaries for formula representations, ideally with manageable added space. We will also need to annotate sources with representations we want not provided in the sources, *e.g.,* OPT is normally not included in LaTeX files or HTML pages. For example, with three representations, we could produce one dictionary for each and use a single formula id to get each representation separately:
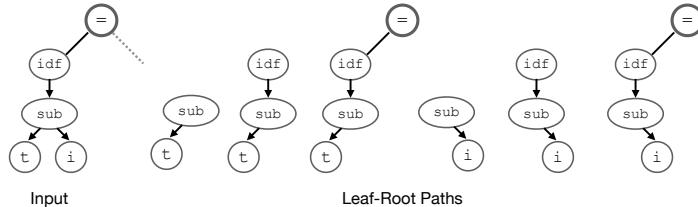
3. **SLT id:** `formula-int-id → SLT-int-uid`
4. **OPT id:** `formula-int-id → OPT-int-uid`
5. **LaTeX id:** `formula-int-id → LaTeX-int-uid`

Notice that each dictionary entry is another identifier; here `-uid` represents *unique identifier.* These refer to three more dictionaries (dictionaries 6-8) holding *unique* SLT, OPT, and LaTeX representations, *e.g.,* `LaTeX-int-uid → latex-string`. Unique formulas are identified after canonicalization to ignore unimportant differences (*e.g.,* LaTeX comments). This avoids repeating common expressions like $x^2$ thousands or millions of times. We can now retrieve a representation by formula id in two lookups: first for the unique id, and then for the representation (*e.g.,* LaTeX string). We can also produce formula dictionaries for additional representation types and/or canonicalization methods.

The choice of how we canonicalize formula representations is important, as this impacts which formulas are treated as identical, and how abstract/generalized formula representations are (*e.g.,* from variable enumeration and token type information). After this we can produce the building-block dictionaries for sources and formulas such as those above, which provide the foundation for formula search index construction.

Example 8: Indexing Paths in SLT and OPT

**Operator Tree Paths:**

Input · Leaf-Root Paths

**Symbol Layout Tree Paths:**

SLT Path Length 1 (All Edges)

SLT Path Length 2 (Partial Edge List)

**Indexing formula sub-structures.** We can also index parts of formulas. Example 8 show paths for OPTs and SLTs of different lengths. For the OPT the paths are rooted to variables $t$ and $i$ in leaf nodes, ending at the root of the OPT. This guarantees that the indexed paths belong to complete subexpressions in the OPT. For the SLT, we use a more generic approach, constructing all paths of lengths 1 and 2. Other approaches to generating paths or substructures have also been used for SLTs and OPTs, *e.g.,* complete subtrees rather than paths.

Let's again use dictionaries to assign an identifier to each unique path seen in our formula collection, representing paths as token tuples (*e.g.,* ($t$, sub) with unique id 7). The next two dictionaries define a

*vocabulary* of paths for SLTs and OPTs in our collection.

9. **SLT path vocab:** `SLT-path-tuple` → `SLT-path-int-uid`

10. **OPT path vocab:** `OPT-path-tuple` → `OPT-path-int-uid`

We then create two more dictionaries for posting lists identifying unique SLTs and OPTs containing a path.

11. **SLTs with path:**  `SLT-path-int-uid` → [ `SLT-int-uid`, ... ]

12. **OPTs with path:**  `OPT-path-int-uid` → [ `OPT-int-uid`, ... ]

With these new dictionaries we can identify all formulas whose SLT contains a specific path tuple with 3 lookups. The first returns the unique id for the SLT tuple, and the second gets a list of `SLT-int-uid` for unique SLTs. We then do reverse lookups for each `SLT-int-uid` in dictionary 3 to collect all matching formulas from their `formula-int-id`. OPT paths are handled in the same manner. This organization allows us to quickly access unique representations separately from all formulas with duplicates, while still connecting unique formula representations with all of their occurrences in a collection.

As a simpler example, to index LATEX tokens, we just perform text indexing at the formula level. Just as for text indexing, we produce a unique integer id for each token – equivalently a token tuple of length 1. We can then map each token to a posting list of unique LATEX formulas:

14. **LATEX token vocab:** `LATEX-token` → `LATEX-token-int-uid`

15. **LATEX with token:** `LATEX-token-int-uid` → [ `LATEX-int-uid`, ... ]

We can do the same type of 3-stage lookup used for SLT paths above to find formulas containing a token, using a reverse lookup on dictionary 5 in the final step. For larger retrieval units than single tokens, it is generally better to use SLT paths rather than LATEX code excerpts, because this captures writing line structure more flexibly.

**Retrieval unit granularities and index types.** Along with the formula representations and canonicalization, the size or *granularity* of the unit(s) used in retrieval is very important. Intuitively, patterns for whole formulas represent *more* information, while formula parts represent *specific* information. Which retrieval unit is more helpful depends upon the query type, how units are represented as patterns, and how similarity is measured. Similarity in sparse vector spaces is computed by accumulating weighted matches from inverted index lookups, *e.g.,* using

IDF (Zobel and Moffat, 2006).[10] Vector similarity is instead geometric, and usually computed from angles between vectors (*e.g.,* cosine similarity) or a line between vector endpoints (*e.g.,* Euclidean distance).

For example, for sparse retrieval we can use our unique SLT representation dictionary and a second dictionary mapping canonicalized SLTs to unique SLT ids, *i.e.,* `SLT-canon` → `SLT-int-uid` (*i.e.,* a unique SLT vocabulary). Unfortunately, querying returns only formulas with the same unique SLT id; results are empty if the formula is not in the vocabulary. This will perform poorly if we want to find variations of the *idf* formula using a formula query, for example. However, searching using SLT paths as described above matches formulas different from the query formula, because we search multiple paths for a single query, and each path need only match part of a formula. Because sparse retrieval is essentially dictionary lookup, approximate matching for query lookups must come from either canonicalization (*e.g.,* variable enumeration), or query expansion.

The choice of pattern(s) to embed in vectors for dense retrieval is just as important as for sparse retrieval. For example, imagine that we embed formula SLTs in vectors and store these in a matrix (*i.e.,* our index of 'dense' vectors). A query SLT vector will return many results sorted by increasing distance from the query vector. However, many similar formulas may not be *idf*-related because the query pattern represents the formula as a whole. Including smaller patterns such as paths that include the function name *idf* will be helpful. For dense retrieval canonicalization is still often helpful *e.g.,* normalizing `<mrow>` tags in Presentation MathML to avoid false distinctions.

There is a clear difference in purpose here with associated trade-offs: inverted indexes catalogue the exact locations of specific patterns in formula representations, while embedded spaces represent similarity through geometric differences between transformed (vectorized) patterns. Construction-wise, inverted indexes are obtained through dictionaries such as described above, where ultimately the key for each index is some type of integer (*e.g.,* for source, formula, unique formula, or

---

[10]Sparse retrieval is generally faster because query patterns are looked up in dictionaries, rather than embedded by a neural net followed by finding neighbors in the embedding space.

unique path). We often use dictionaries when producing dense indexes as well, including enumerating pattern vocabularies for use in indexing (*e.g.,* paths or token integer ids), and to map embedding vectors to their original formula/source.

**Constructing embedding spaces for dense retrieval.**    Unlike the graph or text-based formula representations organized within inverted indexes, for dense retrieval these formula patterns are transformed into embedding vectors. A standard approach is having a neural network play an imitation game that is a classification task. In the game we hide tokens (*i.e.,* mask them) in an input sequence (*e.g.,* LaTeX), or node/edge labels in an input graph (*e.g.,* an OPT or SLT), and the network produces likelihoods for every alternative in a vocabulary. During training, this game is played repeatedly for multiple iterations over a training data set, with network weights updated to improve estimates as defined by a loss function that scores the game, such as the *cross-entropy loss.*[11] Masking and other types of *self-supervised* learning capture a *language model* reflecting how objects are associated through their likelihood of appearing together and in particular ways. The often large amounts of computation required is somewhat confusingly referred to as *pre-training* because network weights are not optimized for retrieval directly.[12]

To further improve dense retrieval performance, additional imitation games known as *learning to rank* tasks are run using the network that require generating ranks for individual sources (*point-wise*), comparing two items at a time (*pair-wise*) or entire rankings (*list-wise*). Network output layers are replaced to play a new game. Weights are again updated to try and improve a chosen metric; as a simple example, for ranking two items at a time this might be the cross-entropy loss with

---

[11]Related to our earlier discussion of IDF and entropy, the cross-entropy loss $\ell_{CE}$ is interesting. For hidden label $l_x$ with correct label $x$ we have $\ell_{CE}(l_x) = -\log p(x)$, where $p(x)$ is the correct label probability estimate. This is exactly the same quantity for alternatives seen in the entropy equation (Equation 2.1). Note that small correct label probabilities produce very large losses.

[12]'Pre-trained' language models often produce surprisingly strong retrieval baselines. For example, symbols or sub-expressions that look quite different may have similar vectors if they are often used in similar contexts, because when masking these, one or the other will be more likely than other alternatives.

four alternatives: (1) item 1 before item 2, and (2) item 2 before item 1, (3) equally relevant, (4) both non-relevant. This is more generically know as *fine-tuning* of network weights. Learning to rank requires test collections where sources *relevant* to specific queries have been identified, as described in the next chapter; normally relevance labeling is at most partially automated (Faggioli *et al.*, 2024). As a result, data available for learning-to-rank is often much smaller than for language model learning. Data for pre-training is often created by randomly hiding labels or applying other manipulations without human involvement. After training is finished, output layers for last game played are removed, and the remaining network layers kept for use in embedding formulas.

Some final steps include computing and storing embedding vectors for collection formulas in a matrix, along with additional matrices or dictionaries storing references to our previously built dictionary keys, *e.g.,* `formula-int-id` values for individual formulas in sources, or `OPT-path-int-uid` for unique OPT paths in the collection. Multiple dense indexes can be created for different representations, canonicalizations, and retrieval units (*e.g.,* formula vs. path vs. LaTeX token).

**Textual Formula Annotations.** As discussed at the beginning of this chapter, formulas provide useful abstraction through the ability to redefine their symbols. However, this means that the symbols of a formula in isolation are often ambiguous, as well as what specifically a formula represents. We can address this by contextualizing the formulas through annotating them with references to surrounding text and formulas. We previously saw a visualization of this in Figure 1.7, where symbol definitions are shown directly around the symbols of a formula in a paper.

Math Entity Linking (MEL) connects math formulas to surrounding context, and includes the related task of resolving *co-references* (*i.e.,* capturing multiple references to the same entity). Context may include descriptions for symbols and entire formulas, other formulas defining symbols in a formula, and external sources (*e.g.,* linking formulas to Wikipedia pages). Example 9 shows a possible annotation of Equation 1 in the Example 2 excerpt. Aside from annotating symbols with descriptions, by resolving the *anaphoric* reference from 'them'

Example 9: Math entity linking and coreference resolution

... Assume there are $N$ documents in the collection, and that term $t_i$ occurs in $n_i$ of them ... the measure proposed by Sparck Jones, as a weight to be applied to term $t_i$, is essentially

$$idf(t_i) = \log \frac{N}{n_i} \tag{1}$$

- Individual variable and function name entities indicated by filled colored boxes
- Entity descriptions shown with underlines matching entity box color
- The *anaphora* of *them* referring to *N documents* shown with an unfilled blue box

to '$N$ documents', we can include text references identifying that $n_i$ represents a document count for term $t_i$, and that the *idf* formula is a 'weight to be applied to the term $t_i$.' Beyond this small example, we should note that symbols like $x$ and $\lambda$ are frequently re-defined within a single paper, leading to multiple definitions (Asakura *et al.*, 2022). This complicates the task of *coreference resolution*, where multiple references to the same mathematical symbol or entity need to be identified and disambiguated when symbols are redefined (Ito *et al.*, 2017).[13]

The majority of early work on MEL was rule-based approaches, due to limited data for training machine learning models. One of the earliest textual formula annotations linked math expressions to their corresponding Wikipedia page (Kristianto *et al.*, 2016b); unfortunately not all math expressions have Wikipedia pages, and context provided in the document where a formula appears is likely more relevant and/or accurate for the formula. Another early system annotated formulas with descriptions and relationships to other formulas in dependency graphs (Kristianto *et al.*, 2017). Textual descriptions are extracted using an SVM-based model to link description nodes to formulas and symbols (Kristianto, Aizawa, *et al.*, 2014). References between formula nodes are captured through structural matching of formula sub-expressions.

---

[13]The Math Identifier-oriented Grounding Annotation Tool (MioGatto) (Asakura *et al.*, 2021) provides a tool for annotating different roles for formulas and symbol, linking identifiers to pre-defined math concepts extracted from the document.

Later systems including MathAlign (Alexeeva *et al.*, 2020) focused on textual annotations within the documents where formulas appear. There has also been work on automated variable typing, where pre-defined mathematical types (*e.g.,* integer, real) are assigned to variables in mathematical formulas using sentences containing descriptions where a symbol appears (Stathopoulos *et al.*, 2018).

More recently there was the *SymLink* shared task at SemEval 2022 (Lai *et al.*, 2022). Symlink requires extracting math symbols with their textual descriptions from LATEX source files collected from the arXiv. The main task requires this to be performed within a LATEX paragraph. First, all text spans (contiguous excerpts) containing math symbols and descriptions are identified, and then symbols are matched with their descriptions. The dataset provides more than 31,000 entities and 20,000 relation pairs. This available data allowed several BERT-based models to be proposed for the task.

As with everything to be used for retrieval, these textual descriptions and their links with formulas and formula symbols need to be compiled in tables. As a simple example, we might represent formula text descriptions using source and excerpt range (span) information, as we did for formulas in LATEX files earlier, and where specific formula elements (objects) are described by text and/or another formula:

- **Formula annotations:** `formula-int-id` → `text-span-id`
- **Text spans:** `text-span-id` → `(source-int-id,`
      `start-char-index, end-char-index)`
- **Symbol annotations:**
   `obj-int-id` → `(text-span-id, formula-int-id)`

If needed, representing formulas within text spans can be done using minor changes to dictionary entires, and the types of formula dictionaries described previously. From here we can produce inverted indices using these linked dictionaries directly, or using smaller parts (*e.g.,* cutting up long text descriptions). For dense retrieval models, we use one or more of the resulting pattern indexes, *e.g.,* embedding symbols/formulas with their associated descriptions, and separately creating a dense vector index for textual descriptions on their own.

Aside from annotation and indexing for retrieval, textual formula annotations can be used for other downsteam information tasks. For

example, using formula symbol identifier descriptions as features for automatically generating Mathematics Subject Classification (MSC) subject codes (Schubotz *et al.*, 2016). MSC is a collaboratively-produced hierarchical classification scheme used to identify subject codes for papers in math journals. Recent math-aware search engines have also explored using annotated formulas as their collection, including the math entity cards in MathDeck (Dmello, 2019b) that connect formulas to titles and descriptions from Wikipedia. Another retrieval system, MathMex (Durgin *et al.*, 2024) indexes formulas that appear with their textual descriptions in a document.

## 2.4   Indexing formulas and text

Having discussed indexing formulas and their annotations, let's now briefly discuss indexing formulas and text together. Example 10 presents a definition of IDF from Wikipedia. Notice that our definition has not changed in any significant way, however, the formula used to represent *idf* has changed. The source collection now appears explicitly as a parameter $D$, and we no longer use $i$ to refer to any term in a vocabulary. The number of documents containing a term, previously variable $n_i$ is now an expression for the number of documents containing the term, with the term noted as simply $t$, and not $t_i$. Here is an example of a different formula for the same *idf* concept, with different variables and operations.

Consider the HTML excerpt in Example 10, with some parts hidden by [...]. We have a section title (`<h3>`), paragraph (`<p>`) and an SLT formula given as Presentation MathML with the original LaTeX used to generate this given as a `<math>` tag attribute and `<annotation>` tag value. The tag `<semantics>` can be is misleading here – this signifies an SLT definition rather than an OPT definition.

We can apply the same types of formula annotations discussed earlier, *e.g.,* generating additional formula representations, indexes for parts of the formula, connecting text descriptions to symbols, and so on. However, we also have opportunities to break the text into different retrieval unit patterns. Titles often identify the topic of a text section, and we can index this separately to provide search of title specifically,

---

Example 10: IDF Definition from Wikipedia tf-idf Page

> **Inverse document frequency**
>
> The **inverse document frequency** is a measure of how much information the word provides, i.e., how common or rare it is across all documents. It is the logarithmically scaled inverse fraction of the documents that contain the word (obtained by dividing the total number of documents by the number of documents containing the term, and then taking the logarithm of that quotient):
>
> $$\mathrm{idf}(t, D) = \log \frac{N}{|\{d : d \in D \text{ and } t \in d\}|}$$
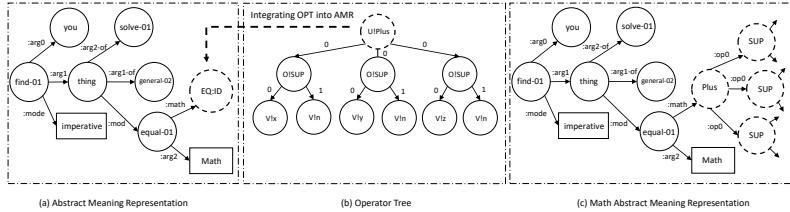
### HTML with Presentation MathML + LaTeX

```
<h3 id="Inverse_document_frequency">Inverse document frequency</h3>
[...]
<p>The <b>inverse document frequency</b> is a measure of how much
    information the word provides, i.e., how common or rare it is
    across all documents. It is the <a href="/wiki/
    Logarithmic_scale" title="Logarithmic␣scale">logarithmically
    scaled</a> inverse fraction of [...]
</p>
[...]
<math xmlns="http://www.w3.org/1998/Math/MathML"
      alttext="{\displaystyle␣\mathrm␣{idf}␣[...]">
  <semantics>
    <mrow class="MJX-TeXAtom-ORD">
      <mstyle displaystyle="true" scriptlevel="0">
        <mrow class="MJX-TeXAtom-ORD">
          <mi mathvariant="normal">i</mi>
          <mi mathvariant="normal">d</mi>
          <mi mathvariant="normal">f</mi>
        </mrow>
        <mo stretchy="false">(</mo>
        <mi>t</mi>
        <mo>,</mo>
        <mi>D</mi>
        <mo stretchy="false">)</mo>
    [...]
    </mrow>
    <annotation encoding="application/x-tex">
      {\displaystyle \mathrm {idf} (t,D)=\log {\frac {N}{|\{d:d\in D
          {\text{ and }}t\in d\}|}}}
    </annotation>
  </semantics>
</math>
```

*Accessed July 31, 2024 (https://en.wikipedia.org/wiki/Tf-idf)

Example 11: Augmenting AMR Trees with Operator Trees

Abstract Meaning Representation (AMR) tree, with inserted operator tree for formula in the query "Find $x^n + y^n + z^n$ general solution."

(a) Abstract Meaning Representation     (b) Operator Tree     (c) Math Abstract Meaning Representation

or to support weighting titles higher in retrieval scores. Also, we have a link to the *Logarithmic scale* article, and we can use the article name to annotate text in the tag, and/or follow the link to annotate the tagged phrase with text from that article.[14]

To compile content for indexing, it is often helpful to index formulas first, and then replace formulas by their generated identifiers in text passages. For example, the MathML excerpt could be added to a Presentation MathML formula index, and we could then replace the formula tags with the index identifier, *e.g.,* `<math [...] </math>` becomes `EQ::42` for `formula-int-id` 42. If instead want to use the HTML directly including the MathML tags for indexing, we can use an approach similar to that for LaTeX discussed earlier. A dictionary is used to assign integers to unique tokens, and then a secondary dictionary represents each file as a token integer sequence.[15]

In contrast to treating formulas as token sequences, we can also convert text to graphs containing formulas (*e.g.,* OPTs or SLTs). An Abstract Meaning Representation (AMR) graph represents the meaning of text using a hierarchical representation of subject, objects, actions, and attributes (Langkilde and Knight, 1998). An example is shown in Example 11 for a query that includes a formula. AMR graphs have

---

[14]This is riskier, as there will be a context shift when moving to the cited article.

[15]Used for training neural nets with sequential input, *e.g.,* transformers like BERT.

a very similar purpose and even structure as OPTs used to represent operations in formulas. For example, the root of the example AMR is a verb with two arguments and a mode modifier (*imperative*) indicating that the statement is a command or request, with one argument for who receives the command (*you*) and what is requested (a *thing* that is the general solution to the provided equation). The AMR representation at left was produced by a neural network after replacing the formula by its formula identifier `EQ:10` in the input text. This allows the network treat the formula as a single node in the AMR. At right we see this AMR formula identifier node replaced by a separately generated OPT with variable and operation type annotations (shown in the middle). Integrated formula+text AMRs were linearized as a token sequence and then converted to dense vectors using Sentence-BERT, and then used to re-rank Math Stack Exchange search results (Mansouri *et al.*, 2022c; Reimers and Gurevych, 2019). One might instead use a graph neural network to embed from the graph structures directly.

## 2.5 Summary

In this chapter we have focused on representing formulas and text in an index. During indexing we create additional information sources, including dictionaries for source contents and formula annotations, and index files. As part of our discussion, we considered the information that formulas carry on their own, and in combination with surrounding text in a document. Related to this, we provided a brief overview of techniques used to annotate formulas with associated text for use in indexing. Symbol Layout Trees (SLTs) and Operator Trees (OPTs) are the most common formula representations, but other representations may be used such as visual-spatial representations like Pyramidal Histograms of Characters (PHOCs), and more might be developed. The integration of text and formulas in index representations is another place where research opportunities exist, perhaps including new ways to incorporate multiple formula representations and text annotations.
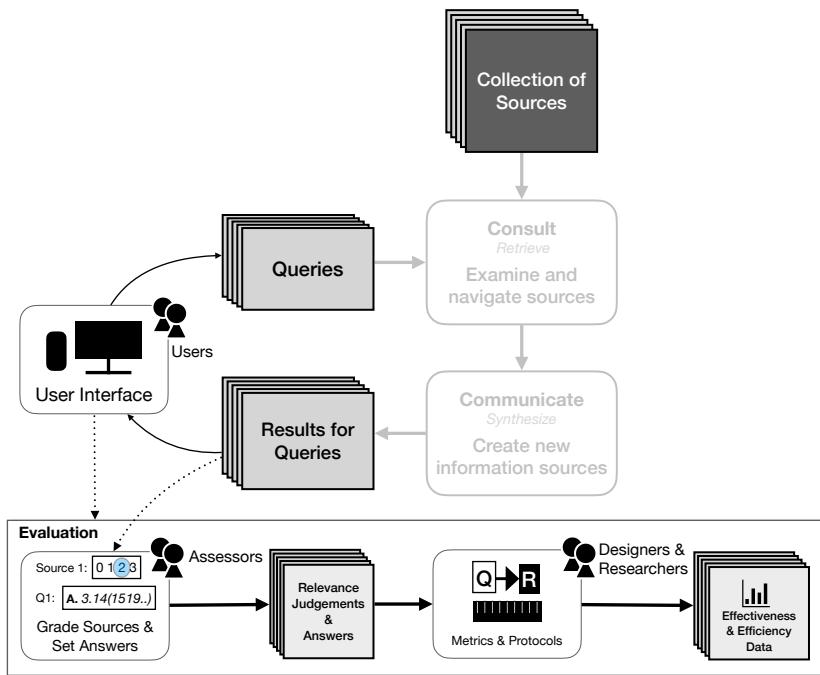
# 3

## Math Retrieval Tasks and Evaluation

Let's now survey retrieval tasks for systems that have been studied in the research literature, along with data and protocols that have been used for system evaluation. There are two main task types that we will discuss: *search tasks* that require a ranked list of sources from a collection, and *question answering* (QA) requiring a single response that may or may not reference a collection of sources.

Figure 3.1 visualizes the people, data, and processes used to create retrieval task data and evaluate search and QA tasks. The figure illustrates an important but easy-to-miss fact: it is *people* and not formal definitions, systems, or algorithms that define retrieval effectiveness, and the target responses that retrieval systems are designed to produce. More specifically, these roles are:

- **Users:** Realistic search queries and questions come from human *users* either directly or from available data (*e.g.,* in query logs).
- **Assessors:** Identify relevant sources for search queries, and define correct question answers.
- **Designers/Researchers:** Collect queries, results, assessments, and set the measures and procedures that quantify performance.

As a result, a chosen evaluation design and the human actors within it

**Figure 3.1:** People, Data, and Methods in Evaluation (expands Figure 1.3). Search tasks rank sources in a specific collection while QA tasks do not, but answers may excerpt or cite sources. System information tasks are in gray/hidden: systems are used to sample sources for grading search relevance, but are mostly unused in QA evaluation.

influence the data produced, along with any observations and conclusions made from this data. It is important to consider this when reviewing evaluation results in the work of others, and when designing our own evaluation frameworks and experiments.

This dependency of system evaluation on people is unavoidable. The tasks retrieval systems perform are motivated by human information tasks, which system designers approximate through observation and then formalization within a system design. Also, large data sets recorded from users and sources authored by people are *human* data. Effectiveness is measured by how well system outputs imitate human responses collected by researchers. For systems utilizing machine learning, desired outputs are obtained by repeatedly playing imitation games scored by

the distance between model outputs and human responses. For these reasons, people determine or constrain nearly every aspect of system design and evaluation.[1]

For search tasks, systems also have a direct role in evaluation aside from producing results, sometimes even for their own evaluation. This is because we normally *pool* sources returned from multiple systems before assessors make relevance judgements. This does introduce bias in evaluation, because only items returned by systems used in pooling are assessed. However, most collections are too large to proceed any other way, and we can choose metrics to mitigate the effects of this bias on evaluation outcomes.

Evaluation-wise, we focus here on the effectiveness of query results as measured offline using *test collections*. A test collection consists of:

1. **Collection of Sources:** sources to be searched for search tasks; optional, usually missing for question answering tasks.
2. **Topics:** queries to run. For search this also includes query information need descriptions and criteria for relevant sources. For QA questions may be accompanied by additional context and/or explanations for how answers are evaluated.
3. **Responses:** includes pooled sources with relevance judgements for search tasks, and question answers for QA tasks.
4. **Protocol:** metrics and methods for producing evaluation data.

In addition to being aware of the roles people play in retrieval evaluation, it is important to remember test collections provide a *sample*, and not all possible queries and results for a task. As such, the data that we collect provides *evidence* for hypotheses (*e.g.,* system $A$ performs better than $B$ for metric $M$), and not proof (Fuhr, 2017). With that said, many valuable things can be learned about system behaviors, evaluation data and frameworks, and even retrieval tasks themselves using test collections. They allow us to address questions using direct observation in addition to observations reported by others.[2]

Briefly returning to our 'jar' model for information tasks in Chapter 1 (Figure 1.2) using test collections for evaluation and related experi-

---

[1]This is a feature, not a bug.

[2]Fun aside: Galileo expressed concern about finding information primarily from book indexes rather than running experiments, see (Duncan, 2021) pp. 9-10.

ments is an important information *synthesis* task that involves *applying* information to *communicate* new information sources such as research papers. Reporting informative evaluation data requires significant effort: careful checking is needed at every step of data collection, measurement, analysis, and reporting. Test collections help ease this burden by providing standard data sets and methods for system comparison.

While discovering new retrieval models and understanding their information use and effectiveness is generally the focus of IR research, for large collections and real-world systems efficiency is also very important.[3] Metrics such as mean query response time ($MRT$, *i.e.,* average seconds/query), query throughput (*i.e.,* average queries/second), and index sizes on disk and in memory are frequently used to evaluate system speed and resource utilization. Efficiency metrics are also required when considering time/space/accuracy tradeoffs discussed earlier. Generally speaking, one should always collect both efficiency and effectiveness metrics, but we will say little about efficiency in this chapter.

## 3.1 Which tasks have been studied?

There are two main math retrieval tasks categories that we will consider: search and question answering (QA). Both have user queries requesting information as input, but return different types of responses. We summarize different search and QA tasks based on the types of responses they require below.

1. **Search:** Ranked list of sources from a provided collection.
   - **Formula search:** ranked formulas.
   - **Math-aware search:** ranked sources/excerpts containing formulas and text (*e.g.,* passages).
2. **Question Answering:** A question answer, sometimes accompanied by a justification or step-by-step solution.
   - **Multiple Choice:** Question includes alternatives that the response selects from.

---

[3]Quoting a former mentor: "Get it right, then make it fast." We add a proviso here: "Still, make it fast and small enough to check behavior quickly."

- **Value:** Answers are numbers or strings, and structures (*e.g.,* lists).
- **Formula:** Formulas or short programs defining how to compute the answer (often paired with value responses).
- **Open response:** Free response with text and formulas.

Some QA tasks also include tables, diagrams and other visual content in their questions, requiring multi-modal processing to generate an answer (*e.g.,* for geometry problems that include a diagram).

Some example formula search queries and results are shown in Example 12. Top-to-bottom, the examples include a formula used directly as a query,[4] a formula query with wildcard symbols that can be replaced by subexpressions, and a *contextualized* formula search query where the context the formula appears is included in the query, and returned formulas include their contexts.

The concrete and wildcard formula queries are symbolic similarity searches, with relevance determined by just formula appearance (SLTs) and/or operations (OPTs). This type of query is motivated by information needs including refinding a previously seen formula in a document collection, or browsing for similar formulas. Wildcards add boolean constraints to queries, as non-wildcard symbols and structures are ideally the same or as similar as possible in the formula returned. Wildcard names can also indicate repetition, *e.g.,* for wildcards **?f**, **?v**, and **?d** in the example.

Contextualized formula queries include the context where a formula appears, incorporating the types of formula-text interactions described in the previous chapter. Here relevance is determined by both the formula and the text within which it appears. For example, two instances of the formula $X^2$ are distinct if in one context $X$ is defined as a number, and in the other it is defined as a set, whereas using symbolic search over SLTs, both formulas have the same representation.[5]

Examples of math-aware search and mathematical question answering tasks are shown in Example 13. We first show an *ad-hoc* math-aware

---

[4] A 'concrete' query

[5] For OPTs, the nodes for the variable and squaring operation may or may not differ, depending upon how the collection is created (see previous chapter).

---

Example 12: Formula Search Tasks

## Formula similarity search

| Query | | Result |
|-------|---|--------|
| $y = \frac{a+bx}{b-x}$ | 1 | $y = \frac{a+bx}{c+x}$ |
| | 2 | $y = a + bx$ |
| | 3 | $y = \frac{a-bx}{c-x}$ |
| | 4 | $y = \frac{a+bx}{x+c}$ |
| | 5 | $g(x) = \frac{x}{x-a}$ |
| | $\vdots$ | $\vdots$ |

## Formula similarity search w. query wildcards (?w)

From NTCIR-10 (Aizawa *et al.*, 2013)

| Query | | Result |
|-------|---|--------|
| $\dfrac{?\mathbf{f}(?\mathbf{v} + ?\mathbf{d}) - ?\mathbf{f}(?\mathbf{v})}{?\mathbf{d}}$ | 1 | $g'(cx) = \lim\limits_{h\to 0} \dfrac{\mathbf{g(cx + h) - g(cx)}}{\mathbf{h}}$ |
| | $\vdots$ | $\vdots$ |

## Contextualized formula search: MSE* answer posts

| Query | Result |
|-------|--------|
| I have the sum $$\sum_{k=0}^{n} \binom{n}{k} k$$ know the result is $n^2 - 1$ but I don't know how you get there. How does one even begin to simplify a sum like this that has binomial coefficients. | 1 $\vdots$ which can be obtained by manipulating the second derivative of $$\sum_{k=0}^{n} \binom{n}{k} z^k$$ and let $z = p/(1 - p)$ $\vdots$ <br> 2 Yes, it is in fact possible to sum this. The answer is $$\sum_{k=0}^{n} \binom{n}{k}\binom{m}{k} = \binom{m+n}{n}$$ assuming that $n \le m$. This comes from the fact that $\vdots$ |

*MSE: Math Stack Exchange

---

## Example 13: Math-Aware Search and Question Answering Tasks

### Ad-hoc retrieval: MSE* answer posts

| Query | Result |
|---|---|
| I have the sum $$\sum_{k=0}^{n} \binom{n}{k} k$$ know the result is $n^2 - 1$ but I don't know how you get there. How does one even begin to simplify a sum like this that has binomial coefficients. | 1 $\vdots$ which can be obtained by manipulating the second derivative of $$\sum_{k=0}^{n} \binom{n}{k} z^k$$ and let $z = p/(1 - p)$ $\vdots$ <br><br> 2  Yes, it is in fact possible to sum this. The answer is $$\sum_{k=0}^{n} \binom{n}{k}\binom{m}{k} = \binom{m+n}{n}$$ assuming that $n \leq m$. This comes from the fact that $\vdots$ |

### Math Q/A: Written response

From ARQMath-3 (Mansouri *et al.*, 2022a)

| Query | Result |
|---|---|
| What does it mean for a matrix to be Hermitian? | A matrix is Hermitian if it is equal to its transpose conjugate. |

### Math word problem Q/A: Equation & value response

From MathQA Amini *et al.*, 2019 and Dolphin18K (Huang *et al.*, 2016)

| Query | Result | |
|---|---|---|
| | Equation | Solution |
| 1. Sarah has 5 pens, David has 3 pens. How many pens do they have? | $x = 5 + 2$ | 8 |
| 2. Find two consecutive integers whose sum is 7. | $x + (x + 1) = 7$ | 3, 4 |

*MSE: Math Stack Exchange

search task with queries that include formulas and text. *Ad-hoc* refers
to the fact that queries can vary greatly, and may include patterns that
are not proper phrases or sentences (*e.g.,* keyword queries, or the query
'$\underline{x^2 + 5 = 30}$ `x value`'). In the example shown, a full question post
from the Community Question Answering (CQA) platform Math Stack
Exchange (MSE) is used as a query, and a collection of MSE answer
posts is searched.

The bottom portion of Example 13 shows two question answering
tasks. The first is an *open response* to an MSE question post, that was
generated using GPT-3 (Mansouri *et al.,* 2022a). The second shows
*math word problems* taken from two test collections. In both cases, the
result should be an answer with two parts: an equation that can be
used to compute the solution, and the solution value, here a number
and a number list. *Multiple choice* questions are also common. These
require choosing from a small number of provided alternatives (*e.g.,* 4
for a fourth alternative *(d) None of the above*). Research-wise, multiple
choice questions allow varying the complexity of information associated
with questions and alternative responses, while constraining system
outputs to an alternative from a small set. As mentioned earlier, in
some collections questions and possible responses also include visual
elements such as tables or diagrams.

## 3.2   Creating topic and response data for test collections

**Selecting test and training topics.**   An important first consideration is
where to collect queries from, and how to select which queries to include.
As our goal is evaluating performance on real-world tasks, it is usually
best if search queries and questions come from real-world users and use
cases. For example, for search tasks topics may come from query logs
or community question-answering websites. Questions on standardized
tests such as the Math SAT are commonly used for question answering.
In some cases, topics generated by the test collection creators are
designed to explore specific scenarios for new features (*e.g.,* wildcards
in formula search queries).

The final topics sets ideally provide a *representative* sample for the
task being evaluated, while including some *diversity* in topics so that

different system capabilities are tested. Diversity is sometimes addressed using separate sub-tasks for a test collection. For math retrieval, criteria to consider include mathematical subjects covered, representations used in queries and responses (*e.g.,* formulas, text, diagrams), and the complexity or mathematical difficulty (*e.g.,* target school grade levels) of queries.

**Train topics, test topics, and cross-validation.**   Normally a test collection divides topics into *training* and *test* topics, so that systems can be compared using the same *test* topics, while being tuned using a separate group of *training* topics. This way, all systems take the same 'test,' without having seen the test search queries/questions previously (*i.e.,* not 'cheating'). This allows us to observe and compare the information and task *generalization* captured in system data structures (*e.g.,* network weights) and algorithms (*e.g.,* for selecting a multiple choice answer) for the same unseen topics. Systems should **never** be tuned on test topics when reporting results for test topics. Published benchmarks for test collections are results for test topics by default.

Training topics are provided for tuning systems. To obtain a more detailed characterization of system behavior using multiple train/test splits, *cross-validation* can be used, and average metric values across splits reported (ideally with standard deviations). All train and test topics may be combined before generating splits for cross-validation, but it must be made clear which topics are used and how train/test splits are produced (*e.g.,* leave-one-out treats every topic as the test sample separately; 5-fold cross validation randomizes the topic order, and then makes 5 equal splits, with each split being rotated as the test split, etc.). While cross-validation provides more robust evaluation measures, it is important to again note that official benchmarks for test collections are computed using *unseen* test topics, and test topics *cannot* be used in training.[6] Unless noted otherwise, for benchmarking test topics are scored just once, without cross-validation.

Test collections sometimes include smaller train/test topic sets used primarily for development, and/or to make use of the collection easier

---

[6]This can be easy to miss amongst multiple data sets versions. Care is needed.

for those new to a task. These can also be safely produced from small
subsets of training topics (*never* test topics), and are very helpful for
fast testing and debugging.

**Responses: Pooled relevance judgments for search tasks.**   For ques-
tion answering tasks, answer data is often compiled from available
sources by those creating the test collection. In contrast, for search
tasks relevance judgements are needed before a test collection can be
released, as they are required to measure effectiveness. The most com-
mon approach is *shared tasks* in which multiple participants run their
systems on provided search topics, and then share the outputs of their
runs for pooling as illustrated in Example 14. After these assessments
have been collected, relevance judgements are used to score participants'
systems, and the relevance judgements produced are included when the
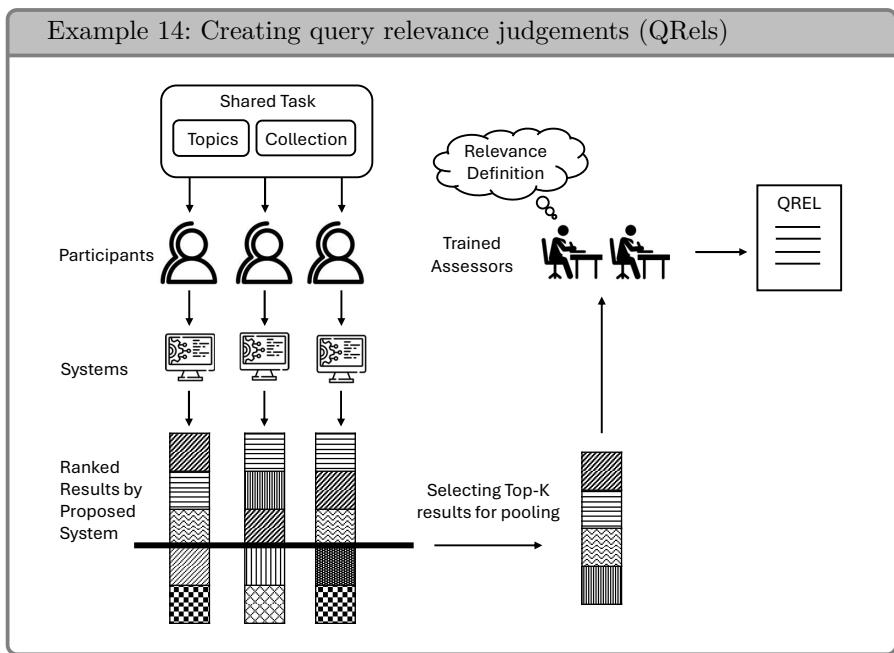final test collection is released.[7]

In a shared task, all participants have their system(s) search the
same collection of sources, and use the same topic queries. Assessors
assign relevance scores to pooled sources, and these assessments are
stored in *qrels* files (quantitative relevance assessments).[8] Shared tasks
are run frequently at conferences including TREC (Text REtrieval
Conference), CLEF (Conference and Labs of the Evaluation Forum),
NTCIR (NII Testbeds and Community for Information access Research),
and FIRE (Forum for Information Retrieval Evaluation).

Example 15 illustrates relevance assessments using a binary scale
(*i.e.,* `1` is relevant, `0` non-relevant), and *graded* relevance where an
ordinal scale of three or more values is used, *e.g.,* `Non-relevant` and
`Low, Medium, High` relevance represented Graded relevances can be
easily converted to binary relevances by thresholding. For our graded
relevance example, we might map `Non-Relevant, Low` (0,1) → 0 (non-
relevant), and `Medium, High` (2,3) → 1 (relevant). We also see an
example of *unknown* relevance for a formula search in the search results
from system A at top. Model A retrieved a formula that was not included
in the pool created during the shared task, and so it is missing in the

---

[7]Ideally, the system *runs* (ranked responses for every topic by each participating
system) are also included in the test collection for later study and comparison.

[8]The most standard format is from TREC (https://trec.nist.gov/data/qrels_eng)

Example 14: Creating query relevance judgements (QRels)

published qrels file. We will come back to this later in the chapter.

Assessing relevance in math search is inherently challenging. As discussed in earlier chapters, a person's mathematical expertise influences their perception of relevance: a highly technical document relevant to an expert might be irrelevant to someone with a basic understanding of math. It is necessary for assessors to have an appropriate mathematical background and to be trained for each search task that they will assess. They should be provided with well-defined relevance definitions, including instructions on how to distinguish between different relevance degrees. It is also a good idea to allow assessors to decline assigning a score when they are very uncertain, or to consult an expert.[9]

Example 16 shows the Turkle assessment interface used for contextualized formulas search in the ARQMath shared tasks. Turkle is used to run assessment tasks online. Each assessor has an account, and is assigned topics to evaluate, with the resulting relevance data compiled automatically by the system. In the example we see a query formula in

---

[9]In our own work, a math professor acted in this role.

---

## Example 15: Relevance Assessments

### Binary relevance and unknown relevance (?)

| Rank | Model A | | Model B | |
|------|---------|------------|---------|------------|
|      | Formula | Relevance | Formula | Relevance |
| 1 | $y = \frac{a+bx}{c+x}$ | **1** | $g(x) = \frac{x}{x-a}$ | 0 |
| 2 | $y = a + bx$ | 0 | $y = \frac{a+bx}{c+x}$ | 1 |
| 3 | $y = \frac{a-bx}{c-x}$ | **?** | $y = \frac{a+bx}{x+c}$ | 1 |
| 4 | $y = \frac{a+bx}{x+c}$ | **1** | $y = \frac{a+x}{b+cx}$ | 1 |
| 5 | $g(x) = \frac{x}{x-a}$ | 0 | $y = a + bx$ | 0 |

### Graded relevance (0-3: Non-Relevant, Low, Medium, High)

| Query | Result | Relevance |
|-------|--------|-----------|
| I have the sum $$\sum_{k=0}^{n} \binom{n}{k} k$$ know the result is $n^2 - 1$ but I don't know how you get there. How does one even begin to simplify a sum like this that has binomial coefficients. | 1    ⋮ which can be obtained by manipulating the second derivative of $$\sum_{k=0}^{n} \binom{n}{k} z^k$$ and let $z = p/(1 - p)$ ⋮ | **(3) High relevance** |
|  | 2    Yes, it is in fact possible to sum this. The answer is $$\sum_{k=0}^{n} \binom{n}{k}\binom{m}{k} = \binom{m + n}{n}$$ assuming that $n \leq m$. This comes from the fact that ⋮ | **(0) Non-relevant** |

its MSE question post on the left, and two formulas in their MSE answer posts taken the assessment pool. Assessors were allowed to view the question threads that queries and results appeared using provided links. This was motivated by the fact that threads are already readily available from MSE posts online, and provide helpful context for assessment. At right we see buttons for the 4-level graded relevance scores, and two additional buttons for system failure (*e.g.,* when a result is unreadable), and when an assessor was uncertain how to rate the result. A box for comments was also included, and was primarily used to explain why assessors selected "System failure" or "Do not know."

Normally measures of agreement between assessors are reported for search test collections. This is done by providing the same set of topics among assessors and comparing their assessments with agreement measures such as Cohen's Kappa coefficient. Properly training assessors can help increase agreement among assessors. For instance, in ARQMath-3 (Mansouri *et al.*, 2022a), for the formula search task, the Cohen's Kappa value increased from 0.21 to 0.52 from the first training to the last (third) training session.

An important practical consideration is the time and effort assessors require to generate answers or judge relevance for search results. For example, in ARQMath the contextualized formula search task had an average assessment time of roughly 35 seconds for *each* formula query and candidate formula. Using the same MSE question/answer posts there was a second answer retrieval task, where assessors had to decide how well an answer addressed a given question. Assessors found this task much more difficult, and average assessment time was nearly twice as long as for formula search: 64 seconds per answer.

**Responses: Question answers for QA tasks.** Unlike search tasks, QA test collections can be created and released without the participation of QA systems (*i.e.,* shared tasks are not needed to create QA test collections).

Answers are created, annotated, and checked in a variety of ways, with differing levels of effort for assessors and test collection creators. For example, LaTeX \boxed{} commands already identify final answers for AMC and AIME questions, while the Amazon Figure Eight annotation

---

**Example 16: Relevance assessment for formula search**

**How to compute this combinatoric sum?**

Thread
I have the sum

$$\sum_{k=0}^{n} \binom{n}{k} k$$

I know the result is $n2^{n-1}$ but I don't know how you get there.
How does one even begin to simplify a sum like this that has
binomial coefficients.

**Answer Post**

Thread
**Title:** Show by counting in two ways that $\sum_{k=0}^{n} \binom{n+k}{k} = \sum_{k=0}^{n} 2^k \binom{n}{k}$?

**Question:**
How does one approach a question like this?

Show by counting in two ways that $\sum_{k=0}^{n} \binom{n+k}{k} = \sum_{k=0}^{n} 2^k \binom{n}{k}$

In general, the fact that $k$ can vary from 0 to $n$ suggests to me that we can consider a universe of $2n$
objects and $2^k$ suggests that subsets are involved.

I have tried a variety of approaches using those ideas, but can't find a nice, simple, interpretation of
$\sum_{k=0}^{n} \binom{n+k}{k}$ specifically

| High |
| Medium |
| Low |
| Not Relevant |

Annotator comment

| System failure |
| Do not know |

Thread
**Title:** How to find the sum $\sum_{k=0}^{n} \binom{n+k}{k}$?

**Question:**
How to find the sum $\sum_{k=0}^{n} \binom{n+k}{k}$?

In my book it is written as a hint that we can use the formula $\binom{n}{r} + \binom{n}{r-1} = \binom{n+1}{r}$ to find
it.But I can't figure out how to apply it here.Any ideas?

| High |
| Medium |
| Low |
| Not Relevant |

Annotator comment

| System failure |
| Do not know |

Turkle assessment interface for ARQMath Formula Search. **Left:**
formula query is highlighted in its question post. **Right:** two question
posts containing the same formula are shown. Assessors consider the
post context when deciding the relevance of each pooled formula, and
could check question threads using the 'Thread' links shown.

platform was used to manually select operations and arguments from provided lists for MathQA. Other collections have assessors generate answers for questions directly (GSM8k), or use machine learning techniques to automatically segment final answer values from available CQA answers (*e.g.,* from answer posted in YahooAnswers using SVMs for Dolphin18K). For AQuA-RAT, crowdsourced workers were used to modify questions (*e.g.,* changing variable values) and provided rationales as needed from available GMAT and GRE test questions.

In addition to final answers, some QA test collections include step-by-step answers or textual annotations (*e.g.,* explanations provided separately from an answer). For machine learning models, some of this textual data associated with *training* topics can be used to improve system answers by providing additional contextual data/information, and used to improve generated explanations for answers or improve responses to queries and comments in interactive discussion (*e.g.,* for large language models).

In all QA datasets, there are also post-processing steps to normalize answer data formatting. The same concerns regarding assessor training, consistency, and effort as mentioned for search tasks apply here as well.

## 3.3   Evaluation metrics and protocols

Consider the top-5 retrieval results from two formula search models in Example 15. Which system would you prefer? We might prefer model 'A', because the first result is relevant, or we might instead prefer model 'B', where more relevant formulae are retrieved. Which result we prefer will depend upon our information need, and will influence which effectiveness metrics should be used to score the rankings.

To be focused in their purpose, in addition to collecting queries, responses, and assessments (*e.g.,* in qrels, or target answer files), test collections also need to define the measures of success to be used, and how they are computed. This way people can use test collections independently, and compare their results in a consistent, meaningful way. We want our measures to be automated for consistency and to avoid error as well. This means that we need computable definitions for 'relevant' sources or 'correct' answers.

As a simple example, suppose that we ask a QA system to 'provide the value of $\pi$', and that we have stored the answer as 3.14159. If a single value is expected, we need to define the required number of digits for a response to be scored as correct. We might accept 3.14, but probably not 3; if more digits of $\pi$ are returned than in our stored answer, we probably don't want to penalize this either.

The differences in answer formats mean that measures of *Accuracy* for math QA systems include normalizations of target responses and provided answers, which include tolerances and constraints to avoid penalizing correct answers with some possible loss in detail (*e.g.,* $\pi$ is irrational, and so has infinite decimal places).

Some protocols for evaluation are more complex, such as the use of *visually distinct* formulas for evaluating formula search in ARQMath, which impacts pooling for assessment and the scoring of search results. To help people using test collections, evaluation scripts are usually often provided for running evaluation protocols automatically.

**Search metrics.** Table 3.1 presents metrics that have been used to evaluate systems for formula search, and math-aware search. Something simple but important to note is that most metrics are defined for a single query, but reported using average values. It is helpful when standard deviations from the mean are also reported, to capture the variation of values for individual queries around the mean. Most common metrics use binary relevance values for their computation. As mentioned earlier, we can binarize graded relevance values to compute these metrics.

The ranking metrics shown provide different ways to evaluate a ranking using qrel data, and there are some trade-offs that tend to occur between these metrics. A classic example is the tradeoff between *recall* and *precision*: the more items we return in search results, the more likely that relevant items will be included, which tends to increase recall metrics. However, returning additional items tends to produce more non-relevant than relevant items, which decreases precision metrics. Intuitively, this is because to retrieve more items, the patterns and/or constraints used to select sources for output needs to be more accepting, allowing more types of sources to be matched.

In theory, we can return the entire collection for all queries to

obtain an average recall ($R$@$hit$) of 100% because all relevant items are included in the returned sources. However, the average query precision ($P$@$hit$) will be close to 0%, because few/none of the returned sources are relevant for a query.

As an another example of a trade-off, systems are often designed to maximize mean reciprocal rank, *i.e.,* the query-averaged value for one over the rank of the first relevant hit ($\frac{1}{r}$), or *precision*@1, *i.e.,* the percentage of queries with a relevant source returned at the first position. Doing this often reduces metrics for effectiveness over a full ranking such as mean average precision ($mAP$) or normalized discounted cumulative gain ($nDCG$). This tends to occur because solely prioritizing high relevance sources at the top of rankings can bias a model away from considering rarer and/or partial relevance signals.

Not all metrics shown here behave as expected, or are always applied or interpreted appropriately.[10] For example both $mAP$ and $nDCG$ ($nDCG$@$hit$) characterize relevance in entire rankings. They are helpful for understanding differences between full rankings, but most real-world users will not consider more than a small number of results returned, and so other measures are more appropriate for this (*e.g., $P$@5*).

This tendency for users to examine few items in a ranking also motivates the logarithmic discount starting at the 3rd returned item for $nDCG$. For example, we can use $nDCG$@5 that will give less credit for the 3rd-5th items returned in a ranking; for $P$@5 all relevant items have the same weight. $nDCG$@5 also differs though, in how relevance is defined: graded relevance is used, and so different relevance ratings produce different scores at each rank. In the end, no metric is *better* on its own – it depends upon what we want to measure. If say we want to know how often relevant items are in the top-5 hits, $P$@5 is simpler to interpret than $nDCG$@5.

To avoid biases in model comparisons, even when there are disputed or clearly incorrect relevance assessments in a search test collection, we always include performance measures using the unaltered qrels file. 'Cleaning up' a qrels file by adding or revising entries is to be avoided,

---

[10]see (Fuhr, 2017) for a critique of mRR (with proposed replacement) and common uses of mAP.

**Table 3.1:** Metrics for Search. By default, returned sources not in relevant source set $R$ are considered non-relevant. $S_k$: first $k$ sources returned. **Note:** reported metrics are typically *averaged* over queries but only mAP and mRR explicitly. Ch. 8 of Croft et. al's textbook provides a helpful overview of these metrics (Croft *et al.*, 2009).

| | Name | Formula | Description |
|---|---|---|---|
| **Binary Relevance** | | | |
| RR | **Reciprocal Rank** | $1/k_f$ | Inverse of rank $k_f$ for the first relevant source |
| mRR | Mean Reciprocal Rank | $\frac{1}{|Q|}\sum_{q \in Q} RR(q)$ | Avg. RR for query set $Q$ |
| R *or* R@hit | **Recall** | $|S \cap R|/|R|$ | Percentage of relevant in returned |
| R@k | Recall at (rank) $k$ | $|S_k \cap R|/|R|$ | Percentage of relevant in first $k$ sources returned |
| P *or* P@hit | **Precision** | $|S \cap R|/|S|$ | Percentage of returned in relevant |
| P@k | Precision at (rank) $k$ | $|S_k \cap R|/k$ | Percentage of first k returned in relevant set |
| AP | Average precision | $\frac{1}{|R|}\sum_{k \in K_r} P@k(S_k)$ | Avg. $P@k$ for relevant documents appearing at ranks $K_r$ in $S$ |
| mAP | Mean Average Precision | $\frac{1}{|Q|}\sum_{q \in Q} AP(q)$ | Avg. AP for query set $Q$ |
| **Graded Relevance** | | | |
| DCG@k | **Discounted Cumulative Gain** | $r_1 + \sum_{i=2}^{k}\frac{r_i}{\log_2 i}$ | Sums relevance scores $r_1$ through $r_k$ for returned sources $S_k$ using log discount from rank 3 on |
| iDCG@k | *Ideal* DCG | | DCG@k for reverse sorted grades of relevant sources $R$, *e.g.*, $(r_1, \ldots, r_5) = (3, 3, 2, 1, 0)$ |
| nDCG@k | Normalized DCG at (rank) $k$ | $DCG@k/iDCG@k$ | Percentage of ideal DCG obtained for returned sources $S_k$ |
| nDCG or nDCG@hit | Normalized DCG | $nDCG@k$ for $k = |S|$ | nDCG for all returned sources $S$ |
| **Scored Sources Only: Includes Non-Relevant Sources $N$** | | | |
| Bpref | **Binary preference** *(w. binary relevance)* | $\frac{1}{|R|}\sum_{k \in K_r}\left(1 - \frac{\min(|S_k \cap N|, |R|)}{|R|}\right)$ | Avg. percentage relevant before non-relevant for relevant sources; treats $|N| = |R|$ to balance classes. |
| M′ | **Prime metric** *(e.g.,* P′@5, nDCG′) | *(as defined for metric M)* | Compute M with $S' = S \cap (R \cup N)$ rather than full ranking $S$ |

as it prevents direct comparison with other published work using the collection: this changes the ideal response, and the specific corrections may be motivated by improving performance for a specific system (which biases what is corrected). It is acceptable to note in a paper where

issues with provided qrels are found, and in some cases to do additional experiments using modified qrels, provided that results on the official qrels are reported first, and that the qrel transformation is unbiased.[11]

**Qrels and items with unknown relevance.**   Qrels provide judgements for pooled sources used in assessment. Unless the collection is extremely small, we do not have the relevance judgements for all collection sources per topic: only pooled sources are evaluated. This means that people using an existing test collection often retrieve sources that are unevaluated, as illustrated for Model A in Example 15, where the 3rd formula retrieved was not included in the assessment pool for the query. This unrated formula will be treated differently, depending upon the metrics that we use. By default, unrated hits are considered non-relevant, *e.g.,* the *Precision@5 (P@5)* for Model A giving the percentage of relevant items within the first 5 returned would be 2/5.

However, some metrics do not use items with unknown relevances, such as the Bpref and *prime* ($'$) metrics shown in Table 3.1. For example, if a 6th formula returned by Model A was graded relevant, then the $P'$@5 would be 3/5; if the 6th formula is non-relevant or no 6th formula was returned, $P'$@5 is still 2/5. Ignoring unevaluated items allows evaluation using only graded items in qrels files, and avoids assuming unpooled sources are non-relevant.[12]

Bpref measures the number of consistent rank *preferences* for relevant vs. non-relevant sources for a query. A relevant source's preference is *consistent* will all other sources at lower ranks; a non-relevant source is *inconsistent* with relevant sources at lower ranks. Bpref treats a relevant sources preferences as entirely inconsistent if the number of non-relevant sources above it is $\geq |R|$, the number of relevant sources. No credit is given for such a source or other relevant sources appearing below it in a ranking. Relevant sources missing in a ranking are also given no credit.

---

[11]Automated qrel changes are preferred, *e.g.,* using topic and result data. Manual 'cherry-picking' of qrel topics for a model family is weak science. However, studying strong/weak topics for a model is good practice; we will discuss this later.

[12]These metrics are also helpful in early design, as models can simply rank graded sources exhaustively. However, metrics obtained this way *cannot* be compared with published systems, because the retrieval step that filters the collection is skipped.

**Question answering metrics.**    The main metrics for question answering tend to be simpler than for retrieval, as shown in Table 3.2. Most evaluations report the percentage of correct responses (with normalizations/tolerances as described earlier) and/or the number of responses that match the target answers in the test collection exactly. Some test collections include measure *perplexity* to characterize system uncertainty in making multiple choice question answer selections, which uses the cross-entropy loss.

For text responses, open response questions, and comparing explanations or step-by-step solutions against those in a test collection, text similarity measures are used at the token, n-gram, and string (*i.e.,* character or token sequence) levels. Note that all string-based measures, including token F1 and BERTScore are affected by the method used to split words into tokens.

The token F1 measure is the harmonic mean for the percentages of target answer tokens in the response (*i.e.,* recall), and response tokens in the answer (*i.e.,* precision). For embedded tokens, a variation called the *BERTScore* (Zhang *et al.*, 2020b) has been used. This is similar to the token F1 score, but makes use of a trained embedding model (*e.g.,* BERT). All target answer and response tokens are embedded using this model, and the highest cosine similarity between each answer and response token is first computed. From these maximum similarities, the average cosine similarity for an answer to response token is used for recall, and the average cosine similarity from a response token to an answer token is used as the precision, and then F1 is computed as before. Embeddings capture token context missing in the token F1 measure, *e.g.,* this can avoid penalizing synonyms, but the token F1 measure is more easily computed and interpreted.

Similarity based on large sequences such as n-grams or full strings can also be used. *BLEU* was originally developed to measure the success of translations by comparing a translation against one or more accepted translations of a sentence (Papineni *et al.*, 2002). It computes the similarity of n-grams (*i.e.,* token sequences of a fixed length, for different values of $n$) between a response and target answer(s), with an additional penalty for short answers. *sBLEU* modifies this to compute similarity at the sentence level, rather than for complete responses. Edit

**Table 3.2:** Metrics for Evaluating Answers to Questions. Metrics computed using answers created by assessor(s) that develop questions. Here the queries $Q$ are questions expecting answers of different types. $a_k$ is the answer given for question $q_k$. $A_k$: probability distribution for all possible answers to question $q_k$ (*e.g.,* multiple choice question). **Note:** Text similarity measures are tokenization-dependent.

| | Name | Formula / Reference | Description |
|---|---|---|---|
| **Correctness and Uncertainty** | | | |
| EM | Exact match | $\frac{1}{|Q|}\sum_{k\in Q}\delta(a_k,q_k)$ | Percentage of answers identical to target answers ($\delta$ returns 0/1) |
| Accuracy | Response accuracy | $\frac{1}{|Q|}\sum_{k\in Q}e(a_k,q_k)$ | Percentage correct answers w. tolerances in function $e$ (returns 0/1) for questions $Q$ |
| Perplexity | Avg. perplexity of alternative selection | $\exp\left(\frac{1}{|Q|}\sum_{q_k\in|Q|}\ell_{CE}(A_k,q_k)\right)$ | Uncertainty for correct answers: no. equally likely alternatives = avg. cross-entropy loss. |
| **Similarity to Text Answers (including Rationales, Step-by-Step)** | | | |
| **Token-Level** | | | |
| Token F1 | Token F1 score | $\frac{2RP}{R+P}$ | Harmonic mean for % model answer tokens in response (*R*ecall), and % response tokens in model answer (*P*recision) |
| BERTScore | Token F1 w. token embedding cos similarity | $\frac{2R_{ac}P_{ac}}{R_{ac}+P_{ac}}$ | Token F1 variant for embeddings. Highest answer to response token pair cos. similarities give avg. token similarity for answer ($R_{ac}$), avg. token similarity for response ($P_{ac}$). |
| **N-gram-Level** | | | |
| BLEU | Bi-Lingual Evaluation Understudy | (Papineni *et al.*, 2002) | [0, 1] score from shared model answer/response n-grams + penalty for short outputs. |
| sBLEU | Sentence-BLUE | | BLEU for individual sentences. |
| **String-Level** | | | |
| Edit distance | String edit distance | (Yu *et al.*, 2016) | Number operations to convert one string to another (*e.g.,* insert, delete, replace). Can be normalized to [0, 1] using string lengths. |

$\delta(a,b)=(a=b)$; $e(a,b)=(normalize(a)=normalize(b))$; $\exp(x)=e^x$
BLEU and edit distance have many variations.

distance considers an entire token sequence, and computes the number of operations from a fixed set needed to transform one to the other (*e.g.,* insert, delete, replace (Yu *et al.*, 2016)).

For the text similarity measures, it is important to realize that they reflect surface structure or correlations learned by a model, and do not directly quantify *semantic* similarity. These metrics are certainly correlated for similar statements, but they capture differences in *information content* between responses and target answers indirectly; what they actually measure is how one string imitates to the other based on tokens, or token embeddings. They are still very useful, but one has to be a bit careful about their interpretation.

**Metric value comparisons and statistical tests.** Imagine that we have two math search systems returning 10 results per query for a test collection with binary relevance grades. We determine that the average $P@10$ is 50% for both systems. However, for the first system half of the first 10 results are always relevant documents from the *qrels* file. For the second system, half of the queries return *no* relevant answers (*i.e., $P@10$* of 0% per query) while the remaining half return only relevant results in the top 10 (*i.e., $P@10$* of 100% per query). While the average $P@10$ scores are identical, we would probably much prefer first system because it is more consistent and avoids missing relevant answers altogether.

To capture variance in our evaluations, such as for rank metrics or differences between target and provided numeric answers to questions, we want to compare *distributions* (*i.e.,* sets) of performance metric values rather than averages. Statistical *hypothesis tests* are used to estimate whether differences in average measures are likely to be stable when running additional queries. They include an estimate for the probability of detecting a difference in error (*i.e.,* a *Type-1 error*) given as the *p-value*. Generally we consider a *p-value* less than either 5% or 1% (*i.e., $p < 0.05$* or *$p < 0.01$*, chosen *before* running an experiment) to be a 'statistically significant' difference suggesting that the averages are *unlikley* to be the same after running additional queries.[13] Note

---

[13]Pro tip: using 'significantly improved' or 'significantly different' without a hypothesis test is a short path to having a research paper rejected.

that hypothesis tests are probabilistic estimates, and not certain answers regarding whether averages are *actually* different in the limit: it is infeasible to run all possible queries to know for certain. Despite this limitation, statistical tests provide a more rigorous and nuanced characterization of differences in metric values than comparing average values directly.

Commonly used statistical hypothesis tests for performance metrics include the standard t-test for comparing two distributions, and the Bonferonni *corrected* t-test when comparing two or more models to a single baseline system. The correction here adjusts computed *p*-values when multiple comparisons are made, because without correction the probability of detecting a difference increases with additional comparisons. Many other tests and comparison types are also used. The selection of a chosen measure or test is motivated by the goal of a comparison, variable data types, and data distribution assumptions (*e.g.,* correlation coefficients, $\chi^2$ ('chi-squared'), and Wilcoxon rank sum tests).

It is also very important when comparing two systems to examine the specific topics where performance differs substantially. These larger differences often help identify specific limitations, patterns of behavior, and information use by models. Equally importantly, this also helps identify bugs in system implementations, including where computed metrics are unusually strong, but *not* because the model is effective.[14]

For search tasks, frameworks like PyTerrier[15] provide easy access to query-specific differences between models, and can be used to compute common statistical tests. Additional helpful evaluation tools include the `trec_eval`[16], `pytrec_eval`[17] and `ranx`.[18] For QA tasks, frameworks such as `nltk`[19] can be used to compute standard text metrics. Standard data matrix tools (*e.g.,* Pandas[20]) and statistical tools (*e.g.,* Scipy stats[21]) can be used to compile descriptive statistics and compute

---

[14]Based on a true story. Or three.

[15]https://pyterrier.readthedocs.io/en/latest

[16]https://github.com/usnistgov/trec_eval

[17]https://github.com/cvangysel/pytrec_eval

[18]https://github.com/AmenRa/ranx

[19]https://www.nltk.org/

[20]https://pandas.pydata.org

[21]https://docs.scipy.org

hypothesis tests.

## 3.4 Formula search

While formula search tasks are defined similarly across different test collections, relevance definitions have evolved over time. There are two basic formula search types, that differ based on differences in representation:

1. **Isolated:** structural similarity of query vs. candidate formulas (SLT and/or OPT), sometimes with optional wildcard symbols
2. **Contextual:** formula relevance depends upon text where query and candidate formulas appear

Both formula search types are visualized in Example 12. Next will describe test collections for formula search task, and explain how relevance definitions differ.

**NTCIR:** NTCIR-10 (Aizawa *et al.*, 2013) was the first shared task on math-aware search, with the formula search task defined as finding similar formulas for a given formula query. In NTCIR-11 (Aizawa *et al.*, 2014), the formula search task was a known-item retrieval task, where for each query there was a specific formula instance in one document to retrieve. NTCIR-12 (Zanibbi *et al.*, 2016a) introduces the Wikipedia Formula Browsing (WFB) task with a similar definition as NTCIR-10; Given a collection of Wikipedia articles, retrieve relevant mathematical formulae for a given mathematical formula query.

NTCIR-10's formula search task used 100,000 technical papers from arXiv (mathematics, physics, and computer science) with 35.5 million formulae. NTCIR-11 used English Wikipedia pages with mathematical formulas as the collection, which was much smaller. NTCIR-12 uses 319,689 articles from English Wikipedia with over 590,000 formulae in the corpus. All three test collections use lab-generated topics. In NTCIR-10, 21 formula queries were chosen by the organizers for arXiv papers, of which 18 queries included wildcards and 3 were concrete queries. NTCIR-11 has 100 queries, with 59 including wildcards, and 41 concrete. Queries were randomly sampled from Wikipedia pages and then modified to include query variables. In the NTCIR-12 Wikipedia Formula Browsing (WFB) task, there were 40 queries, divided into 20

**Table 3.3:** Math Search Tasks and Test Collections. Given a query, search systems consult a collection of sources and communicate ranked sources as output.

| Task & Test Collection | | Collection of Sources | Queries | Results | Effectiveness Metrics |
|---|---|---|---|---|---|
| **Formula Search** | | | | | |
| NTCIR-10 | (2013) | arXiv papers | Organizers$^w$ | F in paper | mAP, P@{5,10,hit} |
| NTCIR-11 | (2014) | Wikipedia | Individual Wiki F (known item) | F in articles | mRR |
| NTCIR-12 | (2016) | Wikipedia$'$ | Organizers$^w$ | F in articles | P@{5,10,15,20}, Bpref |
| ARQMath-1 | (2020) | 2018 MSE As | 2019 MSE Qs | F in MSE As | nDCG$'$, mAP$'$, P$'$@10 |
| ARQMath-2 | (2021) | | 2020 MSE Qs | | |
| ARQMath-3 | (2022) | | 2021 MSE Qs | | |
| AccessMath$^†$ | (2018) | Videos + LaTeX notes | Image region or LaTeX | Image regions or LaTeX form. | R@10, mRR@10 |
| **Math-Aware Search (Ad-Hoc Retrieval)** | | | | | |
| MREC | (2011) | arXiv papers | — — | — — | — — |
| CUMTC | (2015) | arXiv papers (MREC data) | MathOverflow Qs | Papers | mAP |
| NTCIR-10 | (2013) | arXiv papers | Organizers$^w$ | Papers | mAP, P@{5,10,hit} |
| NTCIR-11 | (2014) | arXiv frags. | Organizers$^w$ | Paper frags. | mAP, P@{5,10,hit}, Bpref |
| NTCIR-12 | (2016) | 1. arXiv frags. 2. Wikipedia$'$ | 1. Organizers$^w$ 2. Organizers$^w$ | 1. Paper frags. 2. Wiki articles | P@{5,10,15,20} |
| ARQMath-1 | (2020) | 2018 MSE As | 2019 MSE Qs | F in MSE As | nDCG$'$, mAP$'$, P$'$@10 |
| ARQMath-2 | (2021) | | 2020 MSE Qs | | |
| ARQMath-3 | (2022) | | 2021 MSE Qs | | |
| Cross-Math$^†$ | (2024) | | ARQMath 1-3 Qs (4 languages) | | P$'$@10, nDCG$'$@10 |
| **Reference Retrieval for Theorem Proving** | | | | | |
| NATURALPROOFS (2021) | | Proof references: ProofWiki, Stacks, Textbks. | Theorem | Ref set/sequence in proof | mAP, R@k, edit distance, BLEU |

**F:** Formula; **As**: Answers; **Qs**: Questions; **MSE:** Math Stack Exchange
$^w$**Wildcard symbols** in at least some query formulas
$^†$Cross-modal or cross-language retrieval; **frags.**: Fragments (roughly paragraphs)

concrete queries and 20 wildcard queries The wildcard queries were created by replacing one or more sub-expressions in each concrete formula query with wildcards.

Different pooling processes are used in each NTCIR collection Mansouri *et al.*, 2021a, including no pooling for NTCIR-11 as it was a

known-item retrieval task. Every occurrence of a formula was treated
as a separate source for NTCIR-12 WFB, which led to limited diversity
in the judgement pools after selecting the top-20 instances from each
submitted run. For example, for the query $\beta$ (a short formula consisting
of a single symbol), every formula instance in the pool of formulas to
be judged was $\beta$.

Both NTCIR-10 and -12 test collections used graded relevance (0-2):
Relevant (R), Partially Relevant (P), or Non-relevant (N). For NTCIR-
10, the assessors were mathematicians or math students. They viewed
each formula instance from the judgment pool in isolation, considering
the query-specific scenario and judgment criteria that had been specified
when the query was created. For the NTCIR-12 WFB task, there were
two groups of assessors; each group independently judged each pooled
formula. One group consisted of computer science graduate students,
the other consisted of computer science undergraduates. The pooled
formula instances were shown to the assessors in context (highlighted
in the text where they had appeared in the collection), but assessors
were not asked to interpret the pooled formula in that specific context;
the assessment was to be done based on the pooled formula itself. For
each topic in NTCIR-11, the single Relevant (R) formula instance was
defined as the formula instance that had been used as the formula query.
Note that there may have been other instances of the same or similar
formulas in the collection, but like all instances of other formulas, they
would be scored as Non-relevant (N).

For evaluation protocols, NTCIR-10 and -12 combined the 3-level
judgments from two assessors to form a 5-level "Aggregate" relevance
score. This was done by mapping N to 0, P to 1, and R to 2, and then
summing the two scores. The resulting integer scores ranged from a low
of 0 (both assessors judged N) to a high of 4 (both assessors judged R).
For computing evaluation measures, these 0 to 4 levels are binarized
by treating relevance grades 0, 1, and 2 as non-relevant and treating
3 and 4 as relevant. NTCIR-10 reported P@5, P@10, P@hit (i.e., for
all returned results), and MAP. NTCIR-12 uses P@k for k = {5, 10,
15, 20}. Later, researchers used Bpref (Buckley and Voorhees, 2004)
to avoid penalization for unevaluated formulas. NTCIR-11 uses mean
reciprocal rank (mRR), which is appropriate for single retrieval targets

per topic.

**ARQMath:** Answer Retrieval for Question on Math (ARQMath) introduced a *contextualized* formula search task illustrated in Example 12. The test collection was developed through three years CLEF labs, generating test collections referred to by ARQMath-1 (2020), ARQMath-2 (2021), and ARQMath-3 (2022). ARQMath's collection consists of posts from a math community question-answering website, Math Stack Exchange (MSE), where users can post questions and answers. The benefits of using CQA question posts for queries includes diversity in both subject areas and required mathematical expertise, ranging from simple questions from high school to advanced topics. Formula queries are taken from question posts, and the task is to find relevant formulae inside other question and answer posts. All MSE questions and answers posted from 2010 to 2018, are provided in the collection of sources.

Formula queries for topics were selected from questions posted in 2019 (ARQMath-1), 2020 (ARQMath-2), and 2021 (ARQMath-3). To make topics diverse, ARQMath attached a complexity label to topic formulas, dividing them into low, medium, and high-complexity topics. Additional details on topics and topic selection can be found elsewhere (Mansouri *et al.*, 2021a). Additional training topics are provided in the test collection.

An issue with earlier formula search test collections was some queries had very few unique formulas assessed. To avoid this, ARQMath pooling selects *visually distinct* formulas while pooling to increase diversity for assessed formulas. When two formulas have identical Symbol Layout Trees, they *visually identical*, and otherwise *visually distinct*. The canonicalized SLT representation from Tangent-S (Davila and Zanibbi, 2017a) was used to identify visually distinct formulas when two formulas are parsable, or had identical LaTeX strings otherwise.

For assessment, ARQMath organizers hired undergraduate and graduate students either in mathematics or with strong mathematical backgrounds and trained them for assessment. These assessors were trained with a math professor during three training sessions, after each discussing their choice and aiming to minimize assessment errors. Relevance judgment was defined for these assessors as follows:

*For a formula query, if a search engine retrieved one or more in-*

**Table 3.4:** Relevance scores, ratings, and definitions for ARQMath formula search task.

| Score | Rating | Definition |
|---|---|---|
| 3 | High | Just as good as finding an exact match to the formula query would be |
| 2 | Medium | Useful but not as good as the original formula would be |
| 1 | Low | There is some chance of finding something useful |
| 0 | Not Relevant | Not expected to be useful |

*stances of this retrieved formula, would that have been expected to be useful for the task that the searcher was attempting to accomplish?*

Assessors were presented with formula instances, and asked to decide their relevance by considering whether retrieving either that instance or some other instance of that formula would have been useful, assigning each formula instance in the judgment pool one of four scores as defined in Table 3.4. For example, if the formula query was $\sum \frac{1}{n^{2+\cos n}}$, and the formula instance to be judged is $\sum_{n-1}^{\infty} \frac{1}{n^2}$, the assessors would decide whether finding the second formula rather than the first would be expected to yield good results. To do this, they would consider the content of the question post containing the query (and, optionally, the thread containing that question post) to understand the searcher's actual information need. Thus, the question post fills a role akin to Borlund's simulated work task (Borlund, 2003), although in this case the title, body, and tags from the question post are included in the topic and thus can optionally be used by the retrieval system.

The MSE question post was used in place of the scenarios that had been provided with NTCIR-10 and NTCIR-12 queries; no query-specific relevance judgment criteria were provided. The assessor can also consult the post containing a retrieved formula instance (which may be another question post or an answer post), along with the associated thread, to see if in that case the formula instance would indeed have been a useful basis for a search. Note that the assessment task is not to determine whether a specific post containing the retrieved formula instance is useful, but rather to use that context to estimate how likely useful content would be found if this or other instances of the retrieved

formula were returned by a search engine.

The ARQMath organizers did make one potentially noticeable change to the way this relevance definition was interpreted for ARQMath-2 and -3. ARQMath-1 assessors had been instructed during training that if the query and candidate formulas were the same (in appearance), then the candidate was certainly highly relevant. For ARQMath-2 and -3, the interpretation of 'exact match' was clarified to take the formula semantics and context directly into account, even in the case of identical formulas (so for example, variables of different types would not be considered the same, even if variable names are identical). This means that an exact match with the formula query may in some cases (depending on context) be considered not relevant. On the other hand, formulas that do not share the same appearance or syntax as the query might be considered relevant. This is usually the case where they are both referring to the same concept. For the formula query $\frac{S}{n} \geq \sqrt[n]{P}$ (ARQ-Math query B.277), formula $\frac{1+2+3+...+n}{n} \geq \sqrt[n]{n!}$ has medium relevance. Both formulas are referring to the AM-GM inequality (of Arithmetic and Geometric Means).

For each visually distinct pooled formula, up to five instances of that formula (based on pooling protocols) were then shown to the assessors. Example 16 shows the Turkle[22] interface used for assessment. As shown in the left panel of the figure, the formula query $\sum_{k=0}^{n} \binom{n}{k} k$ is highlighted in yellow. The assessors could use the question context to understand the user's information need. In the right panel, two instances of one visually-distinct formula, $\sum_{k=0}^{n} \binom{n+k}{k}$, are shown in two different posts. For each instance, the assessor could consider the post in which the instance appeared when deciding the relevance degree. The final relevance score for a formula to is the *maximum* relevance score for any judged instance of that formula.

While the official evaluation is done on the visually distinct formulas, ARQMath introduced the use of "Big Qrel Files", where nearly all assessment data is provided, including assessment for each individual formula instance, along with assessor ID. This can be used to study effectiveness of formula search models, under assessment of different

---

[22]https://github.com/hltcoe/turkle

people, and to *change* how final relevance scores are defined (*e.g.,* using average rather than maximum score).

Finally, evaluation is done after removing duplicate instances of visually identical formulas, and then calculating effectiveness measures over the ranked series of visually distinct formulas. This is done by replacing each formula instance with its associated visually distinct formula id, and then removing duplicates starting from the top of the ranking.

ARQMath evaluation measures were different from previous test collections. To avoid issues with unevaluated hits for people using the test collection, the organizers chose the nDCG′ measure (read as "nDCG-prime") introduced by Sakai (Sakai, 2007) as the primary measure. The nDCG measure on which nDCG′ is based is widely used when graded relevance judgments are available. ARQMath also uses two other measures: Mean Average Precision (MAP′), and Precision at 10 (P′@10), after removing the unjudged hits. For MAP′ and P′@10 High+Medium binarization is used, meaning only the high and medium relevant degrees (2 and 3) were considered as relevant.

**AccessMath.** The AccessMath system described in the first chapter (Davila and Zanibbi, 2018) was developed using a small collection of lecture videos and LATEX lecture notes produced for those lecture videos.[23]

## 3.5 Math-aware search

A math-aware search engine supports search in a collection of sources using both formulas and keywords. This task distinguishes itself from isolated and contextualized formula search tasks by representations used for queries, and the associated information needs. Math-aware search tasks range from the simple ad-hoc query "$a^2 + b^2 = c^2$ proof" to full mathematical questions expressed using formulas and text. As with formula search, currently the primary standard test collections for this task are NTCIR and ARQMath. Additional details and comparison of

---

[23]Notes: https://www.cs.rit.edu/~dprl/files/TangentV-data_results.zip, Videos: https://www.cs.rit.edu/~accessmath/am_videos

these collections can be found elsewhere (Mansouri *et al.*, 2021a). Test
collections for math-aware search are summarized in Table 3.3.

**MREC and CUMTC:** The first collection with annotated formu-
las that we are aware of goes back to 2011, when the Mathematical
REtrieval Collection[24] (MREC) (Líška *et al.*, 2011) data set was intro-
duced. MREC consists of 439,423 scientific documents from arXiv with
more than 158 million formulae with MathML annotations. Four years
later, the Cambridge University MathIR Test Collection (CUMTC)
(Stathopoulos and Teufel, 2015) built on MREC, adding 160 test topics
derived from 120 MathOverflow discussion threads. This was one of the
first attempts to use math community question-answering websites for
producing real-world topics rather than topics created by shared task
organizers. CUMTC topics were selected from question excerpts from
120 threads. These threads have at least one citation to the MREC
collection in their accepted answer.[25] The majority of topics (81%) have
only one relevant document, and 17.5% have two relevant documents.

**NTCIR:** NTCIR-10, -11, and -12 used largely the same collections
used for formula search tasks, consisting of arXiv papers and Wikipedia
articles. The retrieval unit for NTCIR-10 was full documents, which
makes assessment challenging. Therefore, in NTCIR-11, the retrieval unit
was reduced to excerpts (roughly paragraphs), resulting in 8,301,578
search units. NTCIR-12 uses the same NTCIR-11 arXiv collection,
and for Wikipedia pages, the retrieval unit is the full article. These
test collections contain different number of topics, growing in each
lab. In NTCIR-10, there are 15 assessed formula+text topics. NTCIR-
11 contains 50 topics, where each has at least one keyword and one
formula. In NTCIR-12, topics were developed for search on two different
collections: arXiv and Wikipedia. 29 arXiv and 30 Wikipedia topics
were assessed. All the topics contained at least one formula, but 5 in the
arXiv and 3 in the Wikipedia set had no keywords. All the NTCIR topics
are lab-generated and only the query is provided with no additional
description of information needs and search scenario. Pooling methods
also differ between the different collections (Mansouri *et al.*, 2021a).

---

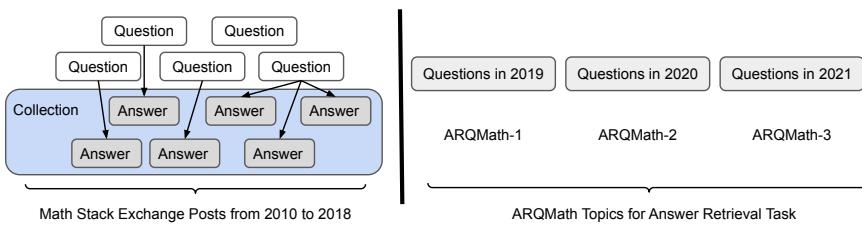[24]https://mir.fi.muni.cz/MREC/index.html
[25]Answer accepted by the user who posted the question.

The assessment process in NTCIR-10 for text+formula searches was similar to the formula search task, with the same assessors. The relevance was decided based on retrieved formulas, not documents due to size and complexity of documents in the collection. For each formula, assessors used a graded 0-2 scale, to represent non-relevant (N), partially relevant (PR) or relevant (R) judgements. Each formula was assessed by one or two assessors.

In NTCIR-11, assessors were shown the title of the topic, the relevance description, and an example hit (if any) as supplementary information. For this collection, relevance was determined using the roughly paragraph-sized retrieval units, rather than individual formulas. To make sure that assessors had enough mathematical background, they were chosen from undergraduate and students in mathematics for the arXiv topics, and computer science for the Wikipedia topics. Relevance levels were the same as for NTCIR-11. Each hit was evaluated by two students, with their judgements combined. The final judgement was decided as R/R and R/PR as relevant (2), PR/PR, R/N, PR/N as partially-relevant (1), and N/N as not relevant (0).

As for the evaluation measures, NTCIR-10 has MAP, P@{5,10, hit} as four basic measures. In NTCIR-11, in addition to MAP and P@{5, 10}, bpref was also included to handle unjudged instances. NTCIR-12 tasks reported P@{5, 10, 15, 20}. In all cases, relevance judgments for formula instances that were missing from the pools (as can happen for P@hit and MAP) were treated as not relevant.

**ARQMath:** ARQMath's main task is retrieving Math Stack Exchange (MSE) answer posts using MSE question containing formulas and text as queries. This was motivated primarily by three things: first, short answer posts are easier to assess (*i.e.,* usually at most a few paragraphs in a post, vs. complete papers), Second, a query log analysis (Mansouri *et al.,* 2019b) that showed that the number of question queries was almost 10% higher for math searches compared to searches in general. Third, the question posts act as both query and information need description. ARMath's topics and collection are built from Math Stack Exchange, as shown in Figure 3.2. All questions and their related answers posted from 2010 to 2018 are provided for training, including roughly 1 million questions and 28 million formulas. ARQMath topics

**Figure 3.2:** ARQMath Answer Retrieval task built on Math Stack Exchange. Questions and their answers posted from 2010 to 2018 are available for training purposes. Topics are selected from new questions from 2019 on, and answers are retrieved the collection of answers from 2010-2018.

are selected from the new questions posted in 2019, 2020, and 2021.

Overall, there are 226 assessed test topics: 77 in ARQMath-1, 71 in ARQMath-2, and 78 in ARQMath-3. A constraint for topic selection is that question posts should contain at least one formula, and 3 categorizations were assigned to candidate topics:

1. **Topic type**: *computation*, *concept* or *proof*
2. **Difficulty**: *low*, *medium*, and *high*
3. **Representation dependency**: *text*, *formulas*, or *both*

These categorizations were used to diversify selected topics.

The pooling process in ARQMath test collections are based on different cuts of participating teams. Each team was allowed to submit up to 5 runs, selecting one run as primary and the others as alternative. Deeper cuts were considered for primary runs. ARQMath provides all system runs used for pooling, which provides a way to study different approaches and understand their behavior in greater detail.

ARQMath assessor for answer retrieval task were also selected from students in mathematics, similar to the formula search task. There were 2-3 training sessions with a math professor to introduce the task and train the assessors in each test collection. Some questions might offer clues as to the level of mathematical knowledge on the part of the person posing the question; others might not. To avoid the need for the assessor to guess about the level of mathematical knowledge available to the person interpreting the answer, we asked assessors to base their judgments on the degree of usefulness for an expert (modeled in this case as a math professor) who might then try to use that answer to

**Table 3.5:** Relevance assessment criteria for the ARQMath Answer Retrieval task.

| SCORE | RATING | DEFINITION |
|---|---|---|
| 3 | High | Sufficient to answer the complete question on its own |
| 2 | Medium | Provides some path towards the solution. This path might come from clarifying the question, or identifying steps towards a solution |
| 1 | Low | Provides information that could be useful for finding or interpreting an answer, or interpreting the question |
| 0 | Not Relevant | Provides no information pertinent to the question or its answers. A post that restates the question without providing any new information is considered non-relevant |

help the person who had asked the original question. Four relevance degrees are considered for this task, defined in Table 3.5. All relevance ratings organized by topic and assessor may be found in a 'big' qrels file available with the test collection.

Finally, for evaluation, the same evaluation measures as the formula search task were used: nDCG′, MAP′ and P′@10. Prime metrics avoids issues discussed earlier for unevaluated hits when using test collections.

**Cross-lingual math information retrieval:** Current math search engines and test collections are primarily developed for the English language, limiting their accessibility and inclusivity. Cross-lingual math information retrieval (CLMIR) is a new task, focusing on retrieving mathematical information across languages. CLMIR has been explored for math-word problems (Tan *et al.*, 2022), where existing datasets were translated into Chinese using online machine translators, and manually refined the translations. CrossMath (Gore *et al.*, 2024) is a novel CLMIR test collection comprising manually translated topics in four languages (Croatian, Czech, Persian, and Spanish). These are the same topics used in the ARQMath Answer Retrieval task. This domain is fairly new and a research gap to be filled. There are approaches proposed for machine translation of mathematical text (Ohri and Schmah, 2021; Petersen *et al.*, 2023) and more to come in the future.

**Other math-aware search collections:** With recent advances in conversational search, these models are explored for math as well.

MathBot (Grossman *et al.*, 2019) is a text-based tutor capable of explaining math concepts providing practice questions and offering feedback to the students. Compared to other domains, there are no standard benchmarks developed for evaluating math conversational search systems. The use of clarifying questions for math has been studied on Math Stack Exchange by extracting comments on math questions (Mansouri and Jahedibashiz, 2023). Tasks such as generating clarifying questions, and detecting ambiguous math queries can be studied on this dataset.

Another interesting problem is theorem proving. NATURALPROOFS (Welleck *et al.*, 2021) for example, is a multi-domain corpus of mathematical statements and their proofs, written in natural mathematical language. The main tasks in this dataset are: 1) *mathematical reference retrieval*: given a theorem, retrieve a set of references that occur in the theorem proof, and 2) *mathematical theorem generation*: generate the sequence of references that occur in a given theorem's proof.

## 3.6    Math question answering

Math question answering test collections require answers to be generated to extracted in response to questions. A summary of available test collections is shown in Table 3.6. Aside from collections requiring open responses (*i.e.,* free text), accuracy is used as the principle metric for performance. Many test collections include math word problems, which are written primarily using words rather than numbers and equations (see Example 13).

Automatically solving math word problems has been an area of interest for artificial intelligence scientists for a long time, back to early work by Bobrow, Feigenbaum, Feldman, and Charniak in the 1960's (Zhang *et al.*, 2020a). The most recent challenge as the time of writing this book is AI Mathematical Olympiad (AIMO).[26] AIMO aims to promote open development of AI models for mathematical reasoning, with the goal of producing a publicly available AI model capable of winning a gold medal at the International Mathematical Olympiad

---

[26]https://aimoprize.com/

(IMO).

Earlier test collections such as Dolphin18K (Huang *et al.*, 2016) focus on arithmetic and algebraic problems, with solutions provided only in the form of equations or final answers. MathQA (Amini *et al.*, 2019) contains 37k English multiple-choice math word problems built on AQuA-RAT (Ling *et al.*, 2017), which annotation and topic selection issues, including wrong solutions and problems that could not be solved without brute-force enumeration of solutions. Despite the corrections to the AQuA-RAT dataset, around 30% of MathQA solutions had inconsistencies.

ASDiv (Academia Sinica Diverse MWP Dataset) (Miao *et al.*, 2020) contains 2.3K math word problems, with each question labeled with a problem type (e.g., Basic arithmetic operations, Aggregative operations) and grade level to show the difficulty level of the problems. This dataset provided different diversity metrics which helped the development of future datasets. This includes GSM8K (Cobbe *et al.*, 2021), which contains 8.5K school math problems created by human problem writers, with 1K problems as the test set. These problems need 2 to 8 steps to get to the solution, using primary calculations with basic math operations to find the answer. This dataset contains questions that need basic knowledge, such as the number of days in a week.

SemEval 2019 Task 10 (Hopkins *et al.*, 2019) provides question sets from MathSAT (Scholastic Achievement Test) practice exams in three categories: Closed Algebra, Open Algebra, and Geometry. Most questions are multiple choice, with some numeric answers. This test collection contains 3860 questions, of which 1092 are used as the test set. The MATH dataset (Hendrycks *et al.*, 2021) also contains multiple choice questions, developed from American high school math competitions including AMC 10/12[27], and AIME[28]. This dataset provides 12,500 problems (7,500 training and 5,000 test). Problems are assigned difficulty levels from 1 to 5, and categorized into seven areas (Prealgebra, Algebra, Number Theory, Counting and Probability, Geometry, Intermediate Algebra, and Precalculus). Answers in this dataset include step-by-step

---

[27]https://maa.org/amc-10-12-information-and-registration/
[28]https://artofproblemsolving.com/wiki/index.php/American_Invitational_Mathematics_Examination

**Table 3.6:** Math Question Answering (QA) Tasks and Test Collections. Given a query, QA systems communicate one answer. Answer formats include written responses, computer programs, values (*e.g.,* numbers, lists), and multiple choice alternatives (*i.e.,* classification). Many queries come from standardized math aptitude tests, such as the SAT.

| Task & Test Collection | | Queries | Results | Effectiveness Metrics |
|---|---|---|---|---|
| **Math Problems** | | | | |
| Dolphin18K | (2016) | Arithm./algebra | Number set | Accuracy° |
| AQuA-RAT | (2017) | MC△ WP† | Rationale + chosen answer | Accuracy, Perplexity, Expert ration. + BLEU |
| SemEval | (2019) | SAT questions | Numeric or chosen answer | Accuracy |
| MathQA | (2019) | MC WP | Rationale + chosen answer | Accuracy |
| ASDiv | (2020) | WP | Equation and value | Accuracy |
| MATH | (2021) | AMC, AIME | Step-by-step soln + final answer | Accuracy |
| GSM8k | (2021) | WP | Step-by-step soln + final answer | Accuracy |
| DROP | (2019) | Paragraph questions | Number, entity, etc. | Accuracy, token F1,EM |
| LILA | (2022) | WP | Numbers or formulas | token F1 |
| AIMO | (2024) | Int'l. Math. Olympiad (IMO) questions | Integers in [0,999] | Accuracy |
| **Math Problems with Graphical/Visual Content** | | | | |
| GeoQA | (2021) | MC geometry problems w. diagram | Chosen answer | Accuracy |
| UniGeo | (2022) | MC w. diagram (GeoQA extension) | Chosen answer | Accuracy, Accuracy@10 |
| MathVista | (2023) | Figure QA, WP, geometry problems, textbook questions, VQA | Chosen answer or numeric | Accuracy |
| Math-Vision | (2024) | Math reasoning with visual elements | Step-by-step soln + final answer | Accuracy |
| TabMWP | (2023) | MC WP w. tabular data | Number or text | Accuracy |
| **Open Response** | | | | |
| ARQMath-3 | (2022) | ARQMath-3 MSE Qs | Written response (formulas + text) | Relevant MSE answers, + token F1, + BERTScore |
| MathDial | (2023) | LLM-simulated student Q. or comment (GSM8k questions) | Responses to Q. and comments | Expert resp. + sBLEU, + BERTScore |

†**Word Problems**; °**Accuracy** = test@1 = % correct ≈ Exact Match (EM)
△**Multiple Choice**

solutions in LaTeX with final answers annotated explicitly, making it suitable for training models that provide explanations for their answers.

Other test collections aim to test mathematical reasoning by systems. The DROP (Discrete Reasoning Over the content of Paragraphs) test collection (Dua *et al.*, 2019) is a reading comprehension task where a paragraph is provided along with questions about the passage. DROP has often been used to evaluate the reasoning capabilities of large language models such as Gemini and PaLM. After DROP, LILA (Mishra *et al.*, 2022) was developed to unify various mathematical reasoning benchmarks. LILA augments 20 existing datasets with solution programs added to answers, along with instructions for producing answers in natural language. Answers are represented as Python strings that prints a number, expression, list, or other data format. For each task, instruction annotations are provided with a clear definition of the task, a prompt providing a short instruction, and examples to help learning by demonstration. The systems' effectiveness is usually measured through token F1-score comparing the models' output against the model answer.

**Questions with graphical/visual content.** While there are many math-word problem with textual data, in recent years there have been attempts to study this area with multimodal data. For example, TabMWP (Lu *et al.*, 2023), contains 38K open-domain grade-level problems that require mathematical reasoning on both textual and tabular data. To solve these problems, systems need to look into the textual data and look up into the given table cells. Visual math question-answering discussed later in this section is also another multimodal math question-answering problem.

Visual aids such as plots, diagrams, and geometric concepts are commonly used in math. While text-based problems have been investigated extensively, visual math question-answering has been explored far less. GeoQA (Chen *et al.*, 2021a) is a geometric question-answering dataset, where for a textual description of the problem and its related diagram, systems should choose the correct answer from multiple choices. UniGeo (Chen *et al.*, 2022a), expanded this dataset, by including the 4,998 calculation problems (from GeoQA) and adding 9,543 proving problems.

While this was an early attempt at this problem, MathVista (Lu *et*

*al.*, 2024) provides 6,141 examples of mathematical and visual tasks and studies several large language models, and large multimodal models.[29] These examples focus on five primary tasks including, figure question answering, geometry problem solving, math word problem, textbook question answering, and visual question answering with two categories of multiple-choice and free-form questions. Math-Vision (Wang *et al.*, 2024) is another recent test collection for studying large multimodal models for mathematical problems with visual contexts, providing 3,040 diverse problems.

**Open response questions.**    ARQMath-3 (Mansouri *et al.*, 2022a) introduced a pilot domain question-answering task using the same topics as the math-aware search task. Evaluation measures for this task used lexical overlap of tokens (viewing answers as a bag of tokens), using $F_1$ score between the system's answer and relevant answers from the answer retrieval task. The maximum $F_1$ score across all the relevant answers is considered the final score. With a similar approach, BERTScore (Zhang *et al.*, 2020b) is also used to measure this overlap.

As AI techniques develop fast, more and more challenging tasks and datasets are being introduced. For example, the MathDial (Macina *et al.*, 2023) dataset brings attention to math tutoring dialogues by connecting an expert annotator, who role-plays a teacher, with an LLM that simulates the student for reasoning problems. This dataset has 3k tutoring conversations grounded in math word problems from GSM8k. It consists of one-to-one interactions between human teachers and a Large Language Model (LLM) that simulates student errors. The dataset is designed to capture the nuances of student-teacher interactions in a tutoring setting, focusing on multistep math problem-solving scenarios. The main task in MathDial is *Tutor Response Generation*, aiming at modeling the teacher in a dialogue by generating follow-up turns to guide the student to solve problems.

---

[29]https://mathvista.github.io/

## 3.7 Summary

In this chapter we have provided an overview of math retrieval tasks and test collections for search and question answering. We have tried to point out a variety of subtle grounding, technical, and data provenance factors that go into the creation of a test collection, including the people involved, how topics are selected, how relevant sources and correct answers are identified, and the metrics and data manipulations chosen for the evaluation protocol.

As part of this, we have shared how the processes and design decisions in test collections have evolved over time. Creating test collections is a large effort, and particularly for areas like mathematical information retrieval that were less actively investigated until recently, it often takes several attempts over years to discover how to build collections that are both well-designed and practical to use by members of the community.[30]

One thing that we hope the reader will take away from all this, is that there is a *lot* that needs to happen before a test collection benchmark metric can be produced using an evaluation script. Further, especially when results aren't as expected, it can be helpful and/or instructive to reflect on the process used to create that collection, and how the evaluation protocol and data collection processes interact and influence those results.

We should also mention that tasks such as math entity linking (Kristianto *et al.*, 2016c) (as discussed in the previous chapter), natural language premise selection (retrieving the premises useful for proving a mathematical statement) (Ferreira and Freitas, 2020b) are related tasks. Another set of tasks closely related to math information retrieval is the extraction of mathematical formulas from documents. One can view this as a step before indexing, but the outcome of this step can greatly impact the quality of retrieval systems. ICDAR Competition on Recognition of On-line Handwritten Mathematical Expressions (CROHME) considers different tasks such as recognizing formulas from online stroke data and image generated from handwritten strokes (Mahdavi *et al.*, 2019; Xie

---

[30]In addition to the providing standard benchmarks, another motivation for test collection development is to grow communities working in an area, and collaborate with them on creating, critiqueing, and improving available test collections.

*et al.*, 2023). While beyond the scope of this book, these tasks play a crucial role in correctly identifying formulas and using them in retrieval models.
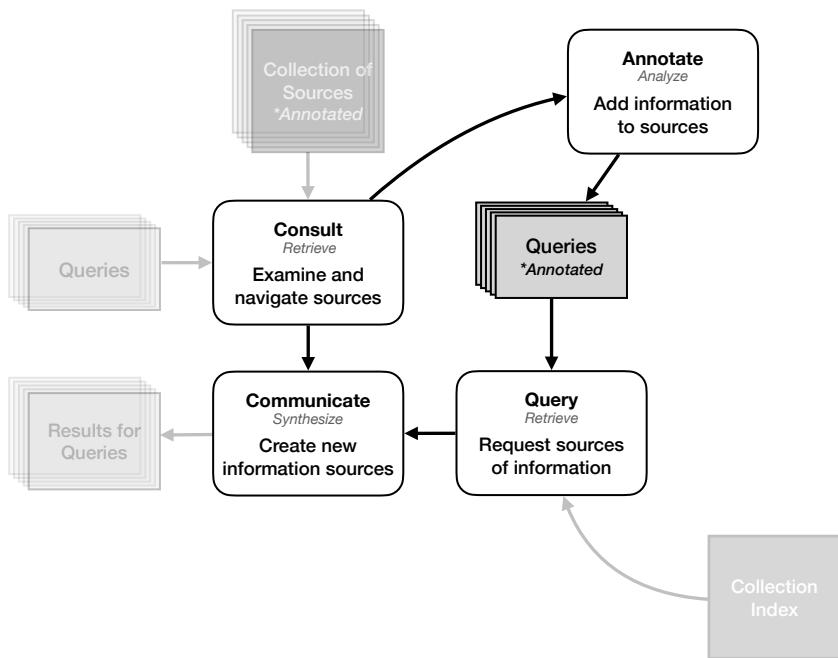
# 4

---

## Formula Search

---

While formulas require context to fully interpret what they represent, there are situations where searching with a formula on its own could be helpful. This includes:

- defining unfamiliar symbols and notation
- re-finding a source using a formula that we remember...almost
- browsing formula variations (*e.g.,* loss functions using cross-entropy loss)
- finding applications of a formula in different domains (*e.g.,* medicine vs. engineering)
- locating premises and theorems in an automated theorem proving system
- finding formulas using subexpressions matches in a document (*i.e.,* `ctrl-f`/filtering for formulas)
- searching formula collections (*e.g.,* formula cards in MathDeck)
- formula auto-completion

These information needs are well-served by formula search.

In Figure 4.1 we illustrate the focus of this chapter in terms of our information task model for retrieval systems. In this discussion, we will assume that formulas have already been indexed using one or more

**Figure 4.1:** After formula patterns have been indexed (see Ch. 2), search queries are consulted and annotated to produce patterns of the same type in the index. The query patterns are used to match formulas in the index, matched formulas are ranked, and then search results are communicated to the user.

representations (*e.g.,* sparse and/or dense: OPT, SLT, visual-spatial, etc.).

In formula search, we first need to annotate the query with the representations used in the available formula indexes. These representations can then be used to query the corresponding index, and a final ranked list of formulas is communicated back to the user.[1] Our focus in this chapter is examining different formula search methods, along with brief notes on the relative effectiveness of methods seen for formula search test collections (*e.g.,* NTCIR, ARQMath). These collections are described in the previous chapter.

---

[1]As discussed in Chapter 1, the presentation of results is important. Unfortunately it has been little studied for formula search.

**Figure 4.2:** Overview of formula retrieval models using different input representations, matching $y = x^2$ against $y = x$.

## 4.1 Formula retrieval models

We have organized approaches to formula search primarily by the type of formula representation used for search. Some of the representations imply sparse vs. dense representations as noted below. A more detailed discussion of formula representations can be found in Chapter 2. The formula representation groups are:

- **Text:** Formula represented as tuple/token sequences for formula text encodings (*e.g.,* LaTeX tokens for individual symbols and operators). Early systems used this with traditional sparse (*i.e.,* inverted index-based) retrieval models.[2]
- **Tree:** Formulas tree representations (mostly SLT and OPT) with retrieval performed using sparse tuple indexing for substructures (*e.g.,* paths, subexpressions), and/or graph alignment techniques.
- **Visual-Spatial:** Formula representation captures formula appearance without the use of writing lines (per SLT) or operation-argument relationships (per OPT).
- **Embedding:** Text, tree, or visual-spatial representations for whole formulas and/or subexpressions are embedded in vector spaces. Nearest-neighbor search using vector similarities are used to find candidates.
- **Other:** Additional approaches including ensembling and learning to rank.

Figure 4.2 gives examples of common representations for the two formulas $y = x^2$ against $y = x$. At left we see text based approaches, where

---

[2]Note (*Error): The linearized tree shown in the figure should only be individual LaTeX tokens, not tuples, which are 'tree-based.'

we have the original L<sup>A</sup>T<sub>E</sub>X strings, and the two lists of symbol pair tokens produced by linearizing the SLT that the L<sup>A</sup>T<sub>E</sub>X represents. Each edge represents an edge in the SLT, with two symbols and a spatial relationship (*e.g.,* horizontally adjacent on the same writing line, or superscripted).

In the middle of Figure 4.2 we have the two operator trees and symbol layout trees at left and right respectively. In tree based approaches subgraphs may be indexed, or whole formulas trees may be compared using methods like tree edit distances and graph alignment. Finally, at right in the figure we see the two formulas represented by points (vectors) in a 3d embedding space; such a representation can be created using text, tree, or other representations, and the specific positions of vectors depend upon training data and the specific learning algorithm used.

Table 4.1 shows the first formula match retrieval using an example search model from each representation category using formula queries of increasing complexity. We see that every models has what can be considered a 'reasonable' first hit. To be useful, it is important for near-identical formulas to appear near the top of a formula similarity search ranking. However, matching identical or nearly-identical formulas is a relatively easy task if a reasonably expressive and correctly implemented representation is used.

Effectiveness-wise, the ability to capture relevance for formulas that are progressively more distinct from the query is what differentiates most formula search models. This is one of the reasons that the ARQ-Math test collections used full rank metrics for ranking formula search systems (*i.e., $nDCG'$*, with $mAP'$ added for comparison), in addition to observing metrics focused on the top of a ranking (*e.g., $P'$@10, $mRR$*).

A summary of formula search models is provided in Table 4.2. In the remainder of this chapter, we will discuss the families of formula search models based on representation type.

## 4.2   Text-based models

It is common practice to use traditional sparse retrieval models for more complex domains such as math. Particularly in the early days of

**Table 4.1:** Rank 1 formula returned by four different formula search engines using different representations, for three queries of increasing structural complexity.

| Model | Repr. | Query | | |
|---|---|---|---|---|
| | | $Cov(x,y)=0$ | $\int_0^1 \frac{\sin^{-1}(x)}{x}$ | $x! = \sqrt{2\pi x} * \left(\frac{x}{e}\right)^x$ |
| **MathDowsers** | Text tokens | $COV(X,Y)=0$ | $J = \int_0^a \frac{\sin^{-1}(x)}{x}\,dx$ | $x! \approx \sqrt{2\pi x}\cdot\left(\frac{x}{e}\right)^x$ |
| **Approach0** | Tree paths (OPT) | $Cov(x,y)=0$ | $\int_0^1 \frac{\sin^{-1}(x)}{x}dx = \frac{\pi}{2}\ln 2$ | $x! \sim \sqrt{2\pi x}\left(\frac{x}{e}\right)^x$ |
| **Tangent-CFT** | Tree Embed. | $\text{cov(y,x)}=0$ | $\int_1^\infty \frac{\sin^2(x)}{x}$ | $n! = \sqrt{2\pi x}\left(\frac{x}{e}\right)^x$ |
| **XY-PHOC** | Visual-spatial | $Cov(x,y)=0$ | $\int_0^1 \frac{\sin(x)dx}{x}$ | $n! = \sqrt{2\pi x}\left(\frac{x}{e}\right)^x$ |

formula search, traditional token-based sparse models such as TF-IDF were used.

An example is one of the earliest formula large-scale formula search engines created for the Digital Library of Mathematical Functions (DLMF) (Miller and Youssef, 2003). The LaTeX representations for formulae are parsed into a tree where nodes correspond to operators and arguments in the formula. This tree is then linearized after normalizing token symbols, which includes converting all symbolic entities to text tokens, and mapping non-alphanumeric characters to alphanumeric strings. For instance, $x^{t-2} = 1$ given as 'x∧{t-2}=1' is converted to the token sequence:

*x, BeginExponent, t, minus, 2, EndExponent, Equal, 1.*

This mapping is applied to both formulae in the collection and input queries. The second normalization is canonical orderings. For commutative operations where argument order is unimportant (*e.g.,* multiplication and addition), a fixed ordering is produced using the alphanumerical ordering of argument tokens. After linearization, DLMF creates an inverted index for use with TF-IDF token-based scoring of queries in the same manner as text.[3]

---

[3]in actual fact, both text and formulas are stored and retrieved from the same index, using a unified token representation.

**Table 4.2:** Formula Search Models. Others refer to representation such as Image or OpenMath. Preprocess refers to the formula representation preprocessing with 'U' indicating unification and 'N' showing normalization.

| Model | Representation | | | | Canonicalization | | Rank | References |
|---|---|---|---|---|---|---|---|---|
| | SLT | OPT | Others | Context | Unif. | Norm. | | |
| **Text** | | | | | | | | |
| DLMF | | | ✓ | | | ✓ | TF-IDF | Miller and Youssef, 2003 |
| ActiveMath | | | ✓ | ✓ | | | Tokens | Libbrecht and Melis, 2006 |
| MathDex | ✓ | | | | ✓ | ✓ | TF-IDF | Miner and Munavalli, 2007 |
| EgoMath | ✓ | | | ✓ | ✓ | ✓ | TF-IDF | Mišutka and Galamboš, 2008 |
| MIaS | ✓ | | | ✓ | ✓ | ✓ | TF-IDF | Sojka and Líška, 2011 |
| LCS | | | ✓ | | ✓ | ✓ | LCS. | Pavan Kumar et al., 2012 |
| **Tree** | | | | | | | | |
| MathWebSearch | | ✓ | | | | ✓ | Paths | Kohlhase and Sucan, 2006 |
| WikiMirs | | ✓ | ✓ | ✓ | ✓ | ✓ | TF-IDF | Hu et al., 2013 |
| SimSearch | ✓ | | | | | | TED | Kamali and Tompa, 2013 |
| MCAT | ✓ | ✓ | | ✓ | ✓ | ✓ | Paths | Kristianto et al., 2016a |
| Tangent-3 | ✓ | | | ✓ | | | Paths | Zanibbi et al., 2016b |
| Tangent-S | ✓ | ✓ | | ✓ | | | Paths | Davila and Zanibbi, 2017a |
| Tangent-L | ✓ | | | ✓ | | | BM25+ | Fraser et al., 2018 |
| Approach0 | | ✓ | | ✓ | ✓ | ✓ | Paths | Zhong and Zanibbi, 2019 |
| MathDowsers | ✓ | | | ✓ | ✓ | ✓ | BM25+ | Ng et al., 2020 |
| Tangent-CFTED | ✓ | ✓ | | | | | TED | Mansouri et al., 2020 |
| **Embedding** | | | | | | | | |
| SMSG5 | ✓ | | ✓ | ✓ | ✓ | | Cosine | Thanda et al., 2016 |
| Formula2vec | | | ✓ | ✓ | | | Cosine | Gao et al., 2017 |
| EqEmb. | ✓ | | | ✓ | | | Cosine | Krstovski and Blei, 2018 |
| Tangent-CFT | ✓ | ✓ | | | ✓ | | Cosine | Mansouri et al., 2019a |
| NTFEM | | | ✓ | | | | Cosine | Dai et al., 2020 |
| Semantic Search | ✓ | | | | | | Cosine | Pfahler and Morik, 2020 |
| Forte | | ✓ | | | | | Cosine | Wang et al., 2021 |
| MathEmb | | ✓ | | ✓ | ✓ | ✓ | Cosine | Song and Chen, 2021 |
| MathBERT | | ✓ | | ✓ | | | Cosine | Peng et al., 2021a |
| MathAMR | | ✓ | | ✓ | | | Cosine | Mansouri et al., 2022d |
| **Visual-Spatial** | | | | | | | | |
| Tangent-V | | | ✓ | | | | Tokens | Davila et al., 2019 |
| XY-PHOC | | | ✓ | | | | Cosine | Avenoso, 2021 |
| EARN | ✓ | ✓ | ✓ | | | | K-NN | Ahmed et al., 2021 |
| **Other** | | | | | | | | |
| TanAPP | ✓ | ✓ | | | ✓ | ✓ | Ens. | Mansouri et al., 2019a |
| Math-L2R | ✓ | ✓ | | | ✓ | ✓ | $SVM^{rank}$ | Mansouri et al., 2021b |
| MathAPP | | ✓ | | ✓ | ✓ | ✓ | Ens. | Peng et al., 2021a |
| FORTEAPP | | ✓ | | | ✓ | ✓ | Ens. | Wang et al., 2021 |

Approaches like DLMF were later extended by applying additional canonicalization steps. For example, EgoMath (Mišutka and Galamboš, 2008) canonicalizes argument ordering, and enumerates variables and constants, using identical symbols to capture variable repetitions.[4] For

---

[4]See Chapter 2 for discussion of symbol enumeration.

constants, formula $74 + a^2 + b^2$ is also indexed as $const + a^{const} + b^{const}$. With variable normalization, formula $a - b$ is also indexed as $id_1 - id_2$. Other normalizations such as removing brackets using distributivity rules are also applied. The goal in these normalizations is to increase recall by making more formulas (ideally, in fact different version of the same formula, or very closely related) have the same token representations.

Aside from sparse retrieval, some other approaches such as using the Longest Common Subsequence (LCS) of a string (Pavan Kumar *et al.*, 2012) have been used to produce a formula similarity score. This model uses tokens from the LaTeX representation of formulae. As before, formulas are canonicalized before applying LCS so that each function, variable, and number is mapped to an unique token, and constants and variables are enumerated.

Another category of text-based models linearized the tree representation of formulae. Math Indexer and Searcher (MIaS) (Sojka and Líška, 2011; Ruzicka *et al.*, 2016) system uses Presentation MathML, encoding each formula as a compact string. For example, for the formula $a + b$, the compact string is encoded as $math\ (mi(a)\ mo(+)\ mi(b))$.

**Effectiveness.** While easier to implement, text-based approaches are generally known to be less effective than other representations for formula search. This is because these models do not capture the hierarchical structure of formulas.

## 4.3 Tree-based models

As discussed in Chapter 2, we normally use graphs to represent structured data, and specifically for formulas, trees to capture a hierarchy of writing lines in SLTs, and a hierarchy of mathematical operations and arguments in OPTs.

Tree-based approaches can be categorized into two main groups: part-based, and full-tree matching. One of the earliest part-based models is MathWebSearch (Kohlhase and Sucan, 2006) that relies on subexpression indexing used originally to unify terms in theorem provers (substitution indexing trees (Graf, 1995a)). Using operator trees, relationships between progressing more concrete formulas are presented by

a series of variable substitutions. A search for expressions with similar operator structures and operands starts from the lowest-precedence operators. Nodes in the substitution indexing tree correspond to expressions with common structures at the top of their operator trees. Moving from the root to the leaves of the substitution tree yields increasingly concrete expressions (*i.e.,* after more variable replacements).

WikiMirs (Hu *et al.*, 2013) creates an OPT-like intermediate presentation from LaTeX strings. To produce the OPT from LaTeX, two template types are used (*opt*: operator, *opr*: operand):

1. **Prefix LaTeX operator**: *opt* $\{opr_1\}$ $\{opr_2\}$...$\{opr_n\}$

2. **Infix operators (*e.g.,* +)**: $opr_1\ opt_1\ opr_2\ opt_2\ ...\ opt_n\ opr_{n+1}$

The OPT-like-tree building process is performed recursively using these templates until no more children can be added. With Template 1, the operator is considered a node with its operands as its children. With Template 2, however, when there is more than one operator the operation order is treated as unknown, and the whole templated expression is mapped to one OPT node. After the tree generation and applying variable and constant enumeration, a term generates tokens for subexpression as well as generalized terms produced by replacing arguments with wildcards ($*$). Tokens are generated using a depth-first traversal.[5] Tokens are then enumerated and used to create an inverted index, which is then searched using TF-IDF match scoring. This system was later extended (Gao *et al.*, 2016), by incorporating text keywords and using more conventional operator trees.

MCAT (Kristianto *et al.*, 2016a) improves part-based retrieval by encoding path and sibling information from symbol layout and operator trees, and uses tuples capturing tree paths as the retrieval units in an inverted index. In addition to the path-based lookup, this model also makes use of a hashing-based formula structure encoding scheme, and text information at three levels of granularity. The first level considers words within a context window of size 10, descriptions, and noun phrases in the same sentence as the formula. The second level considers all the words in the paragraph where the formula appeared. At the third level, the title, and abstract of the document, keywords in the document,

---

[5]This generation is described in more detail in the next chapter.

descriptions of all the formulas, noun phrases, and all words in the document are considered. A formula query combines lookup up in multiple inverted indexes for formula representations and text. This was perhaps the first model to capture surrounding context for formulas in a detailed manner.

**Tangent family models.** Tangent-3 (Zanibbi *et al.*, 2016b) is a two-stage part-based retrieval model. From an SLT, path tuples are generated in the form of $(s_1, s_2, R, \#)$ with ancestor symbol $s_1$, descendant symbol $s_2$, edge label relationship sequence R from s1 to s2, and a count used to capture repetitions ($\#$). These tuples are used to identify an initial set of top-k candidates using a bag-of-words model, scoring by F1 (*i.e.,* harmonic mean of tuples matched on the query and a candidate formula, also known as the *dice coefficient*).

Top candidates based on path tuple F1 are then re-ranked using full-tree matching, aligning the query SLT to each candidate SLT. After alignment, each top-k candidate is scored using the harmonic mean of symbol and relationship recall (the *Maximum Subtree Similarity (MSS)*) and two tie-breakers: symbol precision after unification, and symbol recall without unification. The symbol layout tree representation developed for Tangent-3 provided a basis for subsequent SLT representations (particularly for the Tangent-X family of models). The model introduced a shared container object for matrices, tabular structures, and parenthesized expressions, as well as explicit whitespace, and variable and operation types attached to names (*e.g., N!x* for the number $x$).[6]

A later model combines retrieval over both symbol layout and operator trees (Tangent-S (Davila and Zanibbi, 2017a)). In operator trees, commutative and non-commutative operators have node type (U!) and (O!) for unordered and ordered operations, respectively.

Tangent-L (Fraser *et al.*, 2018) obtains better retrieval results than previous systems due to improve patterns/features and scoring that takes language statistics into account (*e.g.,* through IDF-like terms in the BM25 scoring model). As part of the family of Tangent models, it uses tree traversal to generate math tuples from SLTs, and then

---

[6]See (Zanibbi *et al.*, 2016b) for details.

**Figure 4.3:** Symbol Layout Tree for formula $y_i^j = 1 + x^2$.

tokenizes each pattern to a single term for use in an inverted index. After converting formulae to terms, the BM25+ (Lv and Zhai, 2011) model is used for scoring matches.

For formula feature extraction, consider the formula $y_i^j = 1 + x^2$ with symbol layout tree represented in Figure 4.3. Four types of features are extracted using this representation:

- *Symbol pairs (Tangent-3 tuples)*: For each edge, start and end symbols, along with the edge label. In the example formula, pairs such as $(y, i, \searrow)$. $(y, =, \rightarrow)$ and $(x, 2, \rightarrow)$ are extracted.
- *Terminal symbols (Tangent-3 like)*: Leaf nodes in the symbol layout tree are represented by pairs such as $(i, \varnothing)$, $(j, \varnothing)$, and $(2, \varnothing)$ for the example.
- *Compound symbols*: To capture information about branching in the tree, this feature considers the outgoing edge labels for nodes with more than one outedge; e.g., $(y, \nearrow, \searrow, \rightarrow)$.
- *Augmented locations*: For each of the tuples extracted from the three features above, this feature shows the path from the first symbol in the tree, to the first symbol of that feature. For example, for the tuple $(1, +, \rightarrow)$, tuple $(1, +, \rightarrow, \rightarrow\rightarrow)$ is added to the feature set.[7]

The Tangent-L tuple generator was later re-implemented in the MathDowsers system (Ng *et al.*, 2020; Ng *et al.*, 2021; Kane *et al.*, 2022). The new generator adds additional patterns for repeated symbols, and to apply normalizations. Normalization rules are defined to support

---

[7]A similar approach of representing positions by paths from the root was explored in earlier Tangent models including the original (Stalnaker and Zanibbi, 2015) , but abandoned

operation ('semantic') matches. For example, for commutativity ($A + B$, $B + A$) and symmetry ($A = B$, $B = A$), this model ignores the order of a pair of adjacent symbols. Using a canonical symbol for operator equivalence classes, the model also canonicalizes alternative notations ($A \times B$, $AB$), operator unification ($A \prec B$, $A < B$), and inequality equivalence ($A \leq B$, $B \geq A$). This effectively captures some OPT-type relationships starting from an SLT representation.
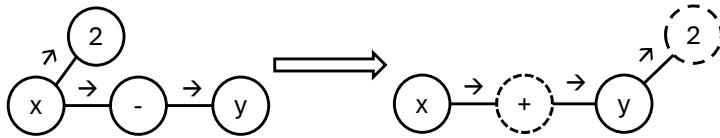
**Approach0.** Approach0 (Zhong and Zanibbi, 2019) uses leaf-root paths in an inverted index within a two-stage model for retrieving operator trees. OPTs are generated by parsing LaTeX with a relatively small but robust expression grammar. To boost recall, variable enumeration is applied using rules in the parser. Candidates are scored based on matches of up to three largest common subtrees identified using dynamic programming. Similarity is scored by a weighted sum of matched leaves (operands) and operators from different common subtrees. In the later version of this system (Zhong *et al.*, 2021), text context is also considered for the search. A textual similarity score is produced using Lucene BM25, and formula structure-based scoring uses the IDF of paths and symbol similarity. These are combined using a linear combination.[8]

**Full tree matching and tree-edit distance (TED).** In addition to tree alignments used for reranking in Tangent-3, full-tree matching via tree-edit distance (TED) has also been used for retrieval. Tree edit distance generalizes string edit distance, defined by the cost (number of edit operations) to convert one tree to the other. The SimSearch (Kamali and Tompa, 2013) model uses tree-edit distance (TED) on SLTs directly as the similarity measure.

In SimSearch, three edit operations are considered: insertion, deletion, and substitution, to find the minimum cost of converting one tree to another. Figure 4.4 shows each of these operations when converting the formula $x^2 - y$ symbol layout tree to $x + y^2$. TED acceleration techniques such as cost-based pruning of candidates and caching subtrees are applied. Operation costs are defined manually, using a set of

---

[8]***NOTE: not verified.

**Figure 4.4:** Converting $x^2 - y$ to $x + y^2$ in symbol layout tree with three edit operations. Node 2 is deleted, node $-$ is substituted with $+$, and node 2 is added as superscript of node $y$.

conditions based on the similarities of node labels, their parent's label, and whether they are leaf nodes. For each condition, there is a cost defined for each edit operation. The final ranking is defined as a reverse edit distance in regard to tree sizes, as Eq.4.1.

$$sim(E_1, E_2) = 1 - \frac{dist(T_1, T_2)}{|T_1| + |T_2|} \tag{4.1}$$

Using a similar approach, Tangent-CFTED (Mansouri *et al.*, 2020) is a re-ranking model using tree-edit distance between SLTs. While in SimSearch the weights for edit operations were set based on heuristics, in this model, weights are learned for each edit operation. Also, this model uses a modified inverse TED as the rank score, as shown in Eq.4.2.[9] Tangent-CFTED uses both symbol layout and operator trees. The final ranking score is the weighted combination of individual ranking scores. This model is named Tangent-CFTED, as it initially re-ranked Tangent-CFT system results (TED standing for tree edit distance).

$$sim(E_1, E_2) = \frac{1}{TED(T_1, T_2) + 1} \tag{4.2}$$

**Effectiveness.**  The strongest models from the most recent ARQMath formula search tasks are tree-based models. Full-tree matching approaches, *e.g.,* using tree-edit distance can be time-consuming and seem to be better suited for re-ranking. Also, these models are stricter than path-based models as they consider the full tree representations. Results on ARQMath-3 formula search task shows the top running systems are tree-based approaches, with Approach0 having the highest

---

[9]1 is added to the denominator to avoid division by zero.

nDCG′ of 0.72, and Tangent-CFTED ranked as second with nDCG′ of 0.69. MathDowsers also provide a high nDCG′ score of 0.64. However, Tangent-CFTED is in fact a two stage model that uses dense vectors for first-stage retrieval, and TED for the second stage. We will discuss dense retrieval models next.
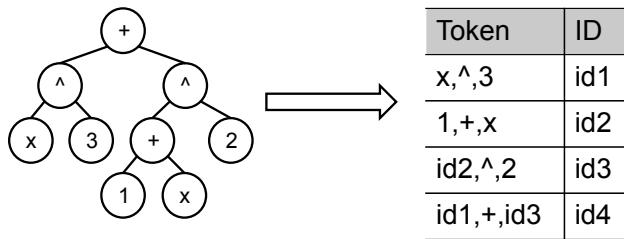
## 4.4 Dense retrieval with formula tree embeddings

As in other domains of information retrieval and natural language processing, researchers working in math IR also turned to embedding models for improved effectiveness and efficiency, and in particular to avoid the types of *vocabulary problems* that traditional sparse model have, and to make greater use of context in patterns used for matching, as discussed earlier in the book.

Models like Word2Vec (Mikolov *et al.*, 2013) were early text embedding models that showed revolutionary results with text problems. These models were then extended to graph data types. The idea was simple; linearizing a graph with different traversals, considering nodes as terms, and then applying Word2Vec. Models like Node2Vec (Grover and Leskovec, 2016) or DeepWalk (Perozzi *et al.*, 2014) use this approach.

SMSG5, the first known embedding model for math formulas, was proposed by the Samsung group as a re-ranking approach (Thanda *et al.*, 2016). In the initial retrieval, formulae in Presentation MathML (SLTs) are linearized and indexed as keywords (using ElasticSearch). For re-ranking, this model uses the doc2vec model (Le and Mikolov, 2014a). Binarized expression trees are embedded into a real-valued vector. Each operator and its operands are treated as tokens for the doc2vec model, using an in-order traversal. If the operand is a subexperssion, the token identifier for the subexpression is used.

Figure 4.5 shows how tokens are extracted for formula $x^3 + (1 + x)^2$. This model uses variable enumeration, replacing variables with identifiers. After training a doc2vec model, the final similarity score between two formulae is calculated using the cosine similarity of their vectors.

Tangent-CFT (Mansouri *et al.*, 2019a) was the first model to use both symbol layout and operator trees for formula embeddings, using

**Figure 4.5:** Binary tree expression for formula $x^3 + (1 + x)^2$ and extracted tokens with SMSG5 model.
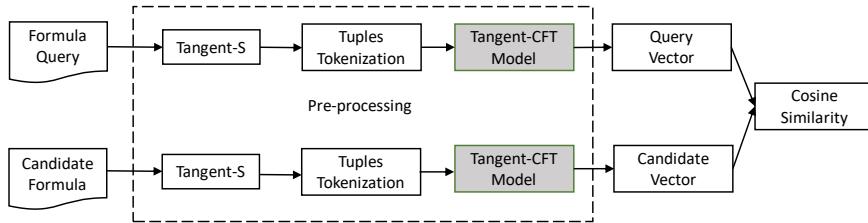
a similar approach to SMSG5. The Tangent-S (Davila and Zanibbi, 2017a) model is used to generates path tuples. For some representations[10] canonicalization is performed using enumeration of variables and constants, and in later versions, operator types. The tuples are then mapped to words, and then tuples are grouped into n-grams. The novelty of this model lies in using an n-gram embedding model, fastText (Bojanowski *et al.*, 2017), that is better suited for queries not seen in the collection as it provides a better sub-formula matching approach. After training a fastText model, vector representations of tuples were averaged to obtain the final representation of a formula.

In the early version of this model, the vectors of different representations (SLT, OPT, Unified SLT) were averaged to get the final vector for a formula. Later, this was converted to combining retrieval results from each representation using a modified Reciprocal Rank (Mansouri *et al.*, 2020). Similar to SMSG5 model, Tangent-CFT uses cosine similarity between vectors for retrieval. A similar approach was applied on other representations such as N-ary tree (Dai *et al.*, 2020) in the N-ary Tree-based Formula Embedding Model (NTFEM).

The approaches seen so far for embedding use linearized tree representations and apply sequence embedding models. As formulas are more naturally represented as trees, graph convolutional neural networks are well-suited to embedding formulas as well. The Semantic Search model (Pfahler and Morik, 2020) generates graph representations from Presentation MathML (SLTs), and also deriving context features from
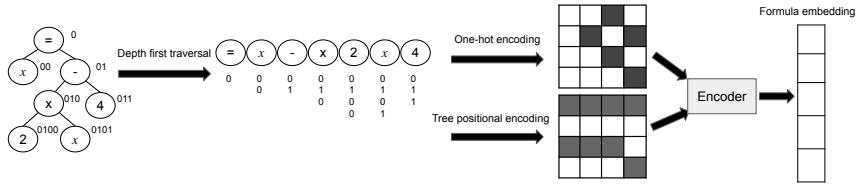
---

[10]*Note: (check claim)

**Figure 4.6:** Retrieval with Tangent-CFT model. Query and candidate formulas are passed to the same pre-processing pipeline to extract their vector representations. The cosine similarity between these two vectors is the similarity score.

tags, attributes, and text. These features were then used to represent nodes as one-hot encoded vectors. To train a graph convolutional neural network model for representation learning, two unsupervised tasks are introduced: 1) a contextual similarity task where labels are generated from the surrounding contexts of mathematical expressions, 2) a self-supervised masking task. Other graph-embedding approaches have considered using symbol layout and operator trees. For example, MathEmb (Song and Chen, 2021) uses operator trees, with Graph Convolution Network, Graph SAmple and aggreGatE (GraphSAGE), and Graph Isomorphism Network.

Encoder-decoder architectures have also been used for formula search. Similar to NLP tasks, reconstruction (also known as the 'fake task') where loss is computed from the ability to decode the original input formula from its embedded vector, can used to train these architectures. After training, only formula encoder is required. FOrmula Representation learning via Tree Embeddings (FORTE) (Wang *et al.*, 2021) uses this architecture by taking in the operator tree on the encoder side, generating the vector representation, and reconstructing the formula in the decoder. On the encoder side, trees are traversed depth-first, and each node is represented by an embedding. To keep track of formula structure, a positional encoding capturing the binary branch path from the root tree is used to represent node position in the tree using a fixed-length vector. The node and tree positional embeddings are then concatenated. The encoding process is shown in Figure 4.7. On the decoder side, this model uses a novel tree beam search generation

**Figure 4.7:** FORTE encoding process for formula $x = 2x - 4$.

algorithm to reconstruct a slightly different version of the input tree, with attached 'end' nodes to each node as its last child.

Different embedding models used formulae textual context when training embeddings for formulae. The formula vector representation in these models is not solely based on the isolated formula and its structure, but the textual context effects the final obtained representation. Early attempts on these models include the Equation embedding model, (Krstovski and Blei, 2018) generating embeddings for words and formulae using textualized formulae. Linearization is done using SLT tuples from the Tangent-3 system. Then Word2Vec embedding model is used within a larger context window size for formulae than words.

A similar idea was adopted in the early days of formula embeddings, using BERT transformer models. In the early days these BERT models (Devlin *et al.*, 2019) were trained on general text. However, later models brought attention to the need for specific tokenizers for math. Note that, fine-tuning BERT-based model for formula search task, is not probably the best approach, as their tokenizers (Byte per encoding or WordPiece), are trained for general text and might not be able to handle formulae correctly.

MathBERT (Peng *et al.*, 2021a), pre-train a BERT model for two tasks: Masked Language Modeling and Context Correspondence Prediction. For formula search, they consider a Masked Substructure Prediction task, with structure representing an operator, its parent node and child nodes as a part of the operator tree. During training, the input to MathBERT model is the formula LaTeX tokens, context and operators as:

[CLS] LaTeX [SEP] Context [SEP] (OPT Nodes)

(a) Abstract Meaning Representation     (b) Operator Tree     (c) Math Abstract Meaning Representation
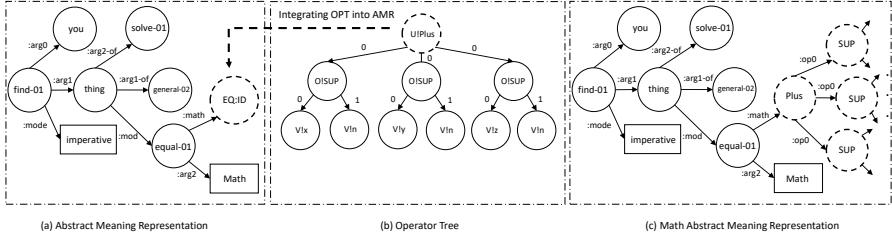
**Figure 4.8:** Example of MathAMR

with [CLS] and [SEP] defined as special tokens. To further incorporate structural information from the operator tree, this model modifies the attention mask matrix, leveraging the edges between nodes in the operator tree.

While MathBERT linearizes formula to obtain a unified representation of text and formulae, MathAMR (Mansouri *et al.*, 2022d) uses Abstract Meaning Representation (AMR) (Banarescu *et al.*, 2013)to have a unified tree representation of both. The process to get the unified tree representation is shown in Figure 4.8. First math formulae are enumerated and replaced with a special token 'EQ:ID' (where ID is the enumerated values). Then using an AMR parser the AMR tree is generated. The the formula is integrated in the AMR by replacing the formula node with the root of operator tree. In the formula operator tree, the edge labels are modified to follow AMR notations. Also special edge label 'math' is added to the set of edge labels in AMR notation to indicate a math formula. Despite having a unified tree representation, the proposed search approach by the authors is linearizing the tree and using a fine-tuning a Sentence-BERT model. With linearization, it is possible that the rich structural representation in the tree is lost. In addition to that, the Sentence-BERT model's tokenizer may not be suitable for this task and handling math formulae.[11]

**Effectiveness.** Experimental results on test collections to date suggest that despite rich contextual features extracted for formulae, these models may be better suited for finding similar or partially-relevant formulae

---

[11]**Note: Repeats material from Ch. 2 (*Revise later)

than very similar or fully relevant. The embedding vectors do well at capturing shared contexts, but not necessarily specific symbols in current models.

Given their ability to match similar formulas well, albeit not rank them in the ideal order, many approaches use embedding-based models to select candidates, and then re-rank the results using similarity scores from the tree-based models. Looking at ARQMath-3 formula search task, Tangent-CFT has a nDCG$'$ of 0.64, and re-ranking with tree-edit distance increases this value to 0.69. A similar pattern is seen for model like MathBERT on NTCIR-12 formula browsing task, where for partial matching the Bpref score is higher than tree-based models with value 0.74, but drops to 0.61 for full matching.

## 4.5 Visual-spatial models

Many math formulas in digital libraries are represented in PDF documents. Even for online resources such as Wikipedia, in many cases, formulae are represented as images rather than Presentation MathML or LaTeX. An interesting example is formulae shown in videos, from technical presentations, classrooms, and tutorials. Some formula search models use visual-spatial representations of formulae, without representation of writing lines or operation hierarchies.

Tangent-V (Davila *et al.*, 2019) is a visual search engine for searching and retrieving mathematical formulas and other graphics in PDF and PNG images. Built on top of the Tangent-S system, Tangent-V utilizes symbol pairs extracted directly from images: for PDFs, symbols are taken directly from the file, and for PNGs, symbols are identified using an open-source OCR system (Davila *et al.*, 2014). Then, this model generates line-of-sight graphs that depict which symbols can see other symbols, storing these symbol pairs with their relative angles in a 2&1/2D representation to capture symbols inside square roots and other containers. During graph generation, the line-of-sight between symbols is tested to determine if two symbols can see each other and if the distance between them falls within a specified threshold.

For search first candidates with shared symbol pairs are retrieved and formulas with large differences in displacement angles and/or symbol

size ratios relative to the query pairs are filtered.[12] The re-ranking step focuses on mapping line-of-sight nodes one-to-one with the query. As for Tangent-S, structural alignment followed by Maximum Subgraph Similarity (MSS) produces final rank scores as a harmonic mean of query node and edge match percentages.

An alternative approach to line-of-sight graphs is a spatial regional embedding of formulae. XY-PHOC (Avenoso, 2021) is a spare retrieval model using simple spatial embeddings of symbols in regions. This extends a one-dimensional spatial encoding previously used for word spotting in handwritten document images, the Pyramidal Histogram of Characters (PHOC) (Sudholt and Fink, 2016),generalized into two-dimensions. The original two-dimensional encoding offers robust spatial embeddings of symbols with lower storage requirements. Scoring is done by cosine similarity of the resulting sparse vectors. Each formula is represented as a bag of symbols. For each symbol, a binary vector of 29 elements is generated, where each element corresponds to a region and value 1 represents the existence of that symbol in that region.

If a symbol appears in a region more than once, it is still represented by a 1.[13] Later investigations found that concentric rectangles were a discrimating region type, and that more detailed representations without large additional space requirements if pyramid region levels are skipped.(Langsenkamp *et al.*, 2024). Unsurprisingly, it also appears that PHOC matching can be improved through the use of part-based matching around symbols rather than using PHOC only for whole formula-level matching (Tucker, 2024). The model is surprisingly effective for formula search despite its simplicity; models that incorporate IDF-like information (*e.g.,* BM25), SPLADE-like token expansions and dense retrieval have not been properly explored with this representation yet, Also, the fact that PHOC provides fixed-length vectors for every formulas suggests this may provide a simple but effective way to combine formulas and other graphics with text in a unified representation.

EARN (Ahmed *et al.*, 2021) is a multimodal embedding model that

---

[12]This informal captures a Boolean query with symbol angles and relative size constraints.

[13]PHOC may be a 2d generalization and/or variation of an unweighted binary independence term model (see (Croft *et al.*, 2009)).

takes advantage of both image and graph representation of formulae. The image encoder uses a formula image (with LATEX rendered as an image) passed to the ResNet (He *et al.*, 2016) model, followed by a Bi-LSTM to extract image embedding. For graph representations, message-passing-based graph encoders are used. The distances from graph-based and image-based embeddings are considered as the features and combined with linear regression similar to the Tangent-S system.

**Effectiveness.**   Visual-spatial models are relatively new to formula search, and their effectiveness present a research opportunity, and should be explored in more detail. The early XY-Phoc provided lower effectiveness compared to other techniques in the ARQMath-3 formula search task for full rank metrics, with nDCG$'$ of 0.47, but was surprisingly competitive in metrics focused on the top of rankings (*e.g.,* $P'$@10). Later refinements including using additional level and rectangular regions produces an $nDCG'$ of 62.3%, and a $P'$@10 of 60.9%, respectively 10% and 8% lower than the best performing tree-based Approach0 model that participated in ARQMath-3, but using a much less complex model without the use of statistical weighting that would improve results (Langsenkamp *et al.*, 2024) (*e.g.,* using BM25).

The multimodal image + SLT dense model EARN has provided higher Bpref scores compared Tangent-S, which is a tree-based model (0.69 vs. 0.64 on NTCIR-12 Wikipedia Formula Browsing task).

## 4.6   Combined representations and formula autocompletion

Ensembling and learning-to-rank approaches have been used in an attempt to combine the benefits of using different representations and dense + sparse retrieval. TanApp (Mansouri *et al.*, 2019a), FORTE-App (Wang *et al.*, 2021), and MathApp (Peng *et al.*, 2021b) are models that consider a linear combination of relevance scores from their embeddings (Tangent-CFT, FORTE, and MathBERT, respectively) with a tree-based approach (Approach0). These ensemble models provide better effectiveness compared to each of the individual models. Models like MathApp and TanApp both had higher Bpref values on NTCIR-12 formula browsing task, compared to their component systems in

| $\int_0^\infty$ | $-\mathrm{dx}$ | $\int_0$ $\quad x$ | $\dfrac{\sin}{x}$ |
|:---:|:---:|:---:|:---:|
| (a) | (b) | (c) | (d) |

**Figure 4.9:** Example of inputting the Equation $\int_0^\infty \frac{sin(x)}{x} dx$ with three symbols when the entry order is a) left-to-right, b) right-to-left, c) alternating left-right outside-in, d) alternating left-right from the middle-out.

isolation.

Using learning to rank models for formula search is an underexplored area. for the formula search task. As mentioned earlier, an early attempt at this was using RankBoost with formula and text features (Gao *et al.*, 2016). A more recent study (Mansouri *et al.*, 2021b) used only formula features, combining similarity features from tree-based and embedding-based models and training an SVM-Rank (Joachims, 2006) model. The input features in this approach include similarity features from path-based, full-tree, and embeddings. An informative outcome of this feature was using SVM-Rank weights to study the importance of features. Based on the findings, full-tree matching features on both symbol layout and operator trees were among the important features.

**Formula autocomplete**   Query auto-completion (QAC) can help users input queries more quickly, and with formulating their queries when they have a specific intent but lack a clear way to express it in words. For text it helps to prevent spelling errors, particularly on devices with small screens. It has been reported that in 2014, when submitting English queries, using QAC by selecting suggested completions saved over 50% of keystrokes for global Yahoo! searchers (Zhang *et al.*, 2015).

Formula auto-completion is employed in search engines like WolframAlpha. This system employs prefix matching for retrieving candidates. Consequently, mathematical expressions that are reordered around commutative operators (e.g., a + b = b + a) or those that utilize a different set of symbols than those in the query will not be presented as candidates.

Despite extensive research in general query autocomplete, the area of formula search remains underexplored. Rohtagi et al. (Rohatgi *et al.*, 2019) proposed an approach that uses LATEX strings and considers three methods: exact matching, prefix matching, and pattern matching. MathDeck (Diaz *et al.*, 2021) uses TangentCFT (Mansouri *et al.*, 2019a) to search a small collection of indexed formulae online as a user inputs a formula, displaying similar formulas.

Both of the previous approaches complete the right part of the query, as the left part is inserted. For math formulae, this is not always the case; for example, when writing fractions or exponentiation. Symbols could be missing from right, left, and middle.

XY-PHOC (Avenoso, 2021) has been used with conjunctive queries, where all symbols in a query must be present in a candidate. This is similar to wildcards in formula queries, where the non-wildcard symbols effectively define a Boolean query/filter. To make this autocompletion work for autocompletion, an additional boolean constraint is added: that returned formulas must contain at least as many symbols as the query. This is the first known model to allow symbols to be inserted in any order for formula autocompletion, because the XY-PHOC is a spatial representation rather than a tree-based one.

To evaluate formula autocomplete, formula search test collections were used. The approach by Rohtagi et al. (Rohatgi *et al.*, 2019) uses partial queries as their approach is for left-to-right autocomplete. An alternative approach is taken by Avenoso (Avenoso, 2021) that compared four different symbol entry orders, illustrated in Figure 4.9. Experiments confirmed the outside-in ordering constraints formulas quickest, raises the rank of target formulas using fewer symbols on average. While retrieval measures such as mean average precision can be used for this task, metrics such as *rsaved* are used in Avenoso's study, representing the target's average reciprocal rank as each symbol is incrementally added for a set of test queries.

## 4.7   Summary

There remain a number of research opportunities for formula search, some of which may have implications for other graphical notations and

multi-modal search in general. These include:

- richer textual annotations and entity linking (*e.g.,* using AMR trees, or other graphs that link text to formula elements directly)
- formula representations: for text/math as well as new visual representations (*e.g.,* PHOC) that may provide a means for representing text and graphics in a single representation (*e.g.,* formulas, tables, chemical diagrams, and other graphics that includes printed characters and symbols).
- improved ensemble techniques for combining results from multiple representations (*e.g.,* via learning-to-rank)
- compact representations: smaller sparse vectors, smaller dense vector spaces
- some work described in the next chapter by the MIaS group started on exploring diversity of math research results. For formula search this has not been investigated further, but might be a useful way to produce strong results in real systems.

# 5

## Multimodal Search, Questions, and Proofs

As previously discussed in Chapter 2, the information conveyed by a formula generally depends heavily upon its surrounding textual context to define variables, operations, and other pertinent information not represented in the formula itself. On their own, formulas primarily convey the operation hierarchy in an expression, and even then assuming that the set of operations and their interactions is known to the reader.

In this chapter we consider three tasks that require using the interaction between formulas and text in sources:

1. math-aware search
2. math question answering
3. theorem proving

Our treatment of theorem proving here is largely focused on applications of dense embeddings and search in the context of supporting premise selection and proof generation.

A common theme across these tasks is the importance of representations used for formulas and text both for original data from sources, and the choice of patterns used to represent mathematical information for computation, whether symbolically in OPTs, inverted index tables, or logical expressions, or as continuous dense vectors embedding the

same chosen patterns. These chosen data structures serve as information sources that constrain and guide subsequent *applications* of information embedded in algorithms used to complete tasks (*e.g.,* to answer a math question from a generated OPT or state representation for sentences in a word problem).

We will also see a recurring role for different forms of annotation of information sources, such as command tokens being inserted in text tokens for transformer inputs, and adding examples of step-by-step answers to questions before a question we want to put to a large language model (LLM). Converting questions in natural language to an OPT or other representation is a form of annotation of its own, adding a new representation for mathematical information in the question. We will return to our information task models from Chapter 1 again later.

## 5.1 Searching with formulas and text

Math-aware search supports queries containing both text and math. Unlike text or formula search where queries and sources have one representation, for math-aware search queries and sources combine text with one or more formula representations (*e.g.,* SLT, OPT, or LaTeX tokens). Because of the different representations (*i.e., modalities*), this is a math-specific variation of *multimodal* information retrieval (Zhu *et al.*, 2024; Shirahama and Grzegorzek, 2016).

A key challenge in multi-modal search is bridging the gap between diverse data formats like text, images, videos, and audio (Bozzon and Fraternali, 2010). There are two main approaches for searching multiple representations: searching modalities separately and combining rank scores to produce the final source ranking, or creating a unified representation for all modalities and searching this directly. Below is a summary of approaches that can be used to combine formula and text search for sparse or dense search indexes.

**Sparse retrieval:**

Formulas are represented using text tokens (*e.g.,* in LaTeX) or token sequences annotated on formulas after traversing nodes of formula trees (*e.g.,* depth-first traversal of OPT or SLT). These formula tokens may represent individual symbols, or tuples, *e.g.,* '+' node

in OPT for $x+1$ as operator-prefix tuple $(+, x, 1)$, or $(x, 1, \rightarrow\rightarrow)$ for the SLT path from x to 1.

1. **Independent:** text and each formula representation have their own inverted index. Query text and formula representations are separated before searches are run, and results are combined to score sources.

2. **Unified:** text and formula tokens belong to one vocabulary, and a 'traditional' inverted index is used to search both together (see Ch. 1). Linearized formula tokens are inserted in source text before indexing, and text/math tokens produced for queries are looked up in the unified inverted index.

**Dense retrieval:**

Formulas and/or subexpressions representations are annotated with vectors in embedding spaces (see Ch. 2). Before embedding, formula representations may be linearized tokens (*e.g.,* LaTeX), trees (*e.g.,* OPT, SLT) or other representations (*e.g.,* PHOC).

- **Independent:** each text granularity (*e.g.,* token vs. sentence) and formula representation have separate embeddings. Query text and formula vectors return the nearest-neighbor vectors in each space, with similarity scores combined to score sources.

- **Unified:** text and formula elements for sources and queries are embedded in the same dense vector space. Multiple space may be used for different granularity (*e.g.,* text+math in passages, vs. individual math/text token embeddings). Vector(s) for sources close to queries in the unified embedding space(s) are used to score sources.

Most systems that combine independent formula and text searches to date use sparse retrieval (*i.e.,* inverted indexes) to produce the initial retrieval results. However some models combine sparse and dense retrieval models. For example, ColBERT (Khattab and Zaharia, 2020b) has been used for text search, and the similarity scores linearly combined with formula search results produced using from Approach0, which was described in the previous chapter.

As a concrete example of combining sparse and dense retrieval models, the MSM team at ARQMath-3 uses Reciprocal Rank Fusion

(RRF) to combine sparse TF-IDF and BM25 retrieval models with a RoBERTa dense retrieval model. The final rank score is produced by EQ. 5.1 to combine retrieval results from a variety of systems, ranging from In this equation, $R$ is a set of rankings, and $r(d)$ is the rank of document $d$. The proposed system achieved nearly the same nDCG$'$ value as the top system for the ARQMath-3 answer retrieval task (0.504 vs. 0.508).
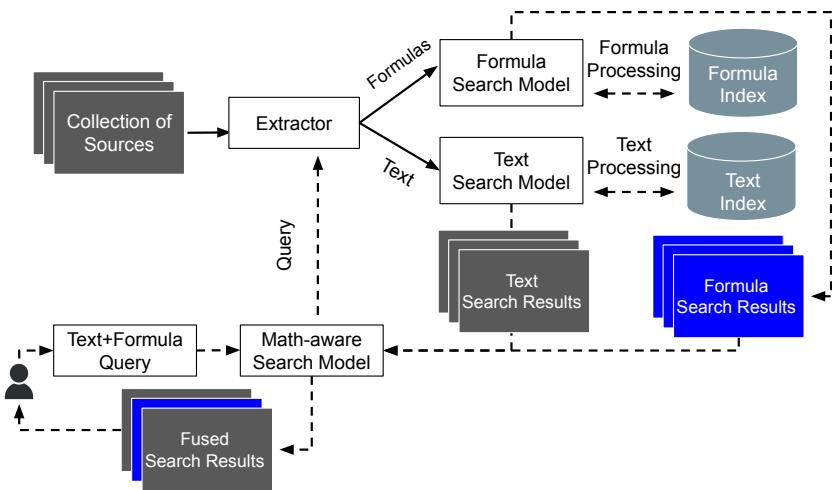
$$RRF(d) = \sum_{r \in R} \frac{1}{60 + r(d)} \qquad (5.1)$$

We review the use of independent and unified text and math representations for search in more detail below.

**Fusing independent formula and text search results.** With several sparse and dense formula search models available, one approach to math-aware search is combining text search with a separate formula search, and then combining the results. Techniques for combining the results include boolean constraints, linearly combining formula and text scores, learning-to-rank, and voting methods.

Figure 5.1 shows a math-aware search model that uses independent searches for text and formulas identify relevant sources (*e.g.,* documents or passages). At indexing time, the *extractor* is used to separate the formulas and text of sources into two separate search indexes, and formulas are annotated with token sequences for sparse models, or embedding vectors for dense models. At query time, the same extractor splits a query into sub-queries for text and formulas, and annotates formulas with a token sequence or vector. The results from both searches are combined into the final score for retrieved sources, and the final result is communicated to the user.

We'll next consider different approaches to fusing independent formula and text searches, starting with Boolean queries used to filter sources without both formula and text token matches. MathWebSearch (Hambasan *et al.*, 2014) represents formulas as tokens for OPT subexpressions stored in a separate inverted index for formulas (implemented in ElasticSearch). Both formula and text inverted indexes are queried separately. However, text search results are used only to implement Boolean

**Figure 5.1:** Fusing independent formula and text search results.

queries of the form $(formula_1 \vee ... \vee formula_n) \wedge (term_1 \vee ... \vee term_n)$, requiring at least one formula and one text token from the query to match a source. Sources without text token matches are removed from the formula rank score list, which produces the final ranking. Despite the simplicity of this approach, the system achieved a $P@5$ of 0.79 for the NTCIR-11 math-aware search task (arXiv papers).

The same type of boolean query/constraint is used to combine formula and text results in the MIaS system (Sojka and Líška, 2011), which uses canonicalized tuples as formula tokens (*e.g.,* with variable unification; implemented in Apache Lucene). Rank scores for sub-queries generated from different combinations of text and formula tokens are used to re-rank the remaining candidates multiple ways, and results from these rankings were interleaved in the final search result to improve the diversity of returned sources (Sojka *et al.*, 2018).

Simple averaging and linear combinations of rank scores have also been used. For the ARQMath Answer Retrieval task three baselines methods were used for answer retrieval:

1. Tangent-S formula search model (linearly combined OPT and SLT scores),
2. TF-IDF text search model, and

3. unweighted average of normalized ($[0, 1]$) Tangent-S and TF-IDF scores.

For Tangent-S, the largest formula (SLT) in the question's title was selected, and if no formula was used in the title, the largest formula in the question body was used. The combined model was more effective then the formula or text model in isolation.

Linearly combining (*i.e.,* scaling and adding) formula and text rank scores was used in the MathDowsers (Ng *et al.*, 2020) system, where Tangent-L results for formula search are combined with $BM25^+$ text search. For text search, a keyword extraction model was used to select tokens in question answer posts for use in text queries. MCAT (Kristianto *et al.*, 2016a) also linearly combines formula and text rank scores, but for multiple formula (OPT, SLT) and text indexes. Text is indexed separately at the paragraph and document (title, abstract, keywords, . . . ) levels. This model was the most effective for participating teams in the NTCIR-12, ArXiv math-aware search task, in part due to the rich variety of formula and text representations.

The WikiMirs system (Gao *et al.*, 2016) uses a learning-to-rank approach to combining formula and text scores. OPT internal nodes for operations (including bracketing/grouping) are converted to sets of concrete and generalized terms (tokens). For example, the formula $(x + 3) \times \frac{a}{b}$, is tokenized into 4 concrete terms, and four generalized terms with wildcards for argument subexpressions:

**Concrete terms** : $\{ (x + 3) \times \frac{a}{b}, \quad (x + 3), \quad \frac{a}{b}, \quad x + 3 \}$
**Generalized terms** : $\{ (*) \times *, \quad\quad (*), \quad\quad \frac{*}{*}, \quad * + * \}$

At query time, query formulas are converted to tuples as shown above, and a set of candidate sources retrieved using inverted indexes for text and formula tokens. These candidates are then re-reanked using RankBoost (Freund *et al.*, 2003) applied to features focused on formulas. This system achieved the highest P@K values among the participating teams in the NTCIR-12 Wikipedia math-aware search task.

The Borda Count has also been used to combine text and formula relevance scores. MaRec (Math answer Recommender) (Gao and Ng, 2023) retrieves answers to math questions. Answers in the collection are classified into topics using text-based Naive Bayes classification (*e.g.,* for

algebra, geometry, etc.). At query time all answers in the sub-collection for the question topic classification are ranked. Answer scoring is done using separate similarity measures for text and formulas.

1. Text similarity: Kullback-Leibler (KL) divergence between question/answer token frequency distributions. Token vocabulary used is chosen with Dirichlet Allocation (LDA (Blei *et al.*, 2003)): terms that characterize topics are selected.[1]

2. Formula similarity: averages tree-edit distance on SLTs and a depth score, based on the sum of leaf-root OPT path lengths matching an answer.

The final ranking score uses the Borda Count, adding the number of answer posts that rank lower than a source in the formula and text rankings.

**Unified representations for formulas and text.**    Combining results from separate indexes for different representations can often produce useful results quickly (especially using sparse retrieval models). However, the text-notation interactions described in Ch. 2 that provide important context for formulas is missing when formulas and text are represented separately. We next consider searching unified representations for formulas and text.

DLMF whose formula retrieval model is described in the previous chapter was the earliest unified sparse retrieval model indexed formula and text tokens together, using a variation of TF-IDF for scoring (Miller and Youssef, 2003). More recently Latent Dirichlet Allocation (LDA) has been used to weight formula and text tokens. In these sparse models, formulas are represented as LaTeX tokens (Yasunaga and Lafferty, 2019) or linearized tree tokens (Thanda *et al.*, 2016), and a single inverted index is used to retrieve and score sources.

In the previous chapter, we saw formula dense retrieval models using the text context around formulas, the majority of which embed text and formulas together using formulas represented by LaTeX tokens. The introduction of the ARQMath answer retrieval task coincided with the

---

[1]Entropy and cross-entropy are closely related to, and even derivable from KL-divergence.

emergence of transformer-based models, and a common technique is using a pre-trained transformer (*e.g.,* BERT variants) that is fine-tuned using pairs of MathSE questions and answers with their associated assessor relevance ratings.

**Using transformers with unified representations and math tokenization.** Reusch (Reusch *et al.*, 2022) studied dense retrieval using Col-BERT and ALBERT models. To fine-tune ALBERT, 1.9M triples containing questions with one relevant and one non-relevant answer were fed to the model. The model attempts to match assessor scores by classifying answers using the learned vector embedding for the *[CLS]* token that starts each token sequence in an answer post. A similar approach was used to fine-tune ColBERT, but using more relevant and non-relevant answers. Somewhat surprisingly, both models proved less effective than sparse retrieval models that participated in the ARQMath-3 shared task.

In subsequent work, the researchers explored how mathematical formulas affect a transfomer model's training (Reusch *et al.*, 2024). They found the transformer models consider formulas when scoring relevance for answers to a given math question. However, a study of the transformer attention weights for formula and text tokens suggests that structural relationships between formula tokens are lost, and that the attention maps do not consider variable overlaps between formulas appearing in questions and answers. This motivates creating transformer-based models that can better capture formula structure and interactions, which are relatively rare in sources used to pretrain transformers.

One approach is adding additional tokens. The MathPredictor model (Jo *et al.*, 2021) extends BERT's tokenizers to support 2,651 new vocabularies. This addresses the BERT WordPiece tokenizer's oversegmention of LATEX commands such as '\overline', which is split into three tokens: {\, over, ##line}. This allows formulas such as $\overline{h}$ (in LATEX expressed as `$\overline h$`) to be correctly tokenized as {$, \overline, h , $ }. The model was then fine-tuned using masked tokens in formulas.

A similar approach is taken in the hybrid MABOWDOR system (Zhong *et al.*, 2023), where the PyA0 toolkit is used for preprocessing mathematical formulas before tokenization by WordPiece. PyA0

canonicalizes math tokens by merging those likely to be semantically identical, such as `\emptyset`, `\empty`, `\varnothing`. 1,000 new math tokens are added to the token vocabulary, and the math tokens are then treated similar to regular text. The final search used combines dense and sparse retrieval, and both independent and unified formula/-text representations. A unified single-vector dense retriever is used for passage-level representations of formulas and text (DPR (Karpukhin *et al.*, 2020)), math is searched independently using Approach0 (combining path-based sparse retrieval with structural alignment), and a dense embedding-augmented sparse retriever (SPLADE (Formal *et al.*, 2021b)) is used to search text independently.

For the MABOWDOR unified representation dense passage retriever, a new pre-training dataset for the math domain was created, Coco-MAE. The pretraining task is Masked Auto-Encoding (MAE), which is similar to masked token pre-training but attempts to decode whole input passages with masked tokens using a decoder. For the final ranking, search results from each component are merged using a convex linear interpolation. This system currently archives the highest $P'$@10 for ARQMath's answer retrieval task.

In an alternative approach, MathBERT (Peng *et al.*, 2021b) (see previous chapter) uses linearized OPT tokens and LATEX formula representations. BERT is pre-trained using Masked Language Modeling, Context Correspondence Prediction (similar to next sentence prediction), and Masked Substructure Prediction (for masked formula tokens). Using an improved tokenization approach for math formulas, this model archives better effectiveness compared to BERT in tasks such as formula topic generation (predicting the topic (tag) associated mathematical formula) (by creating TopicMath-100K dataset on arXiv papers), and formula 'headline' generation (EXEQ-300K (Yuan *et al.*, 2020)), creating a concise description of a formula using the formulas and descriptions in a MathSE question.

Researchers then used MathBERT's language model for retrieval and other 'downstream' tasks. For example, automatic math short-answer grading (Zhang *et al.*, 2022) by classifying grades using an integer ordinal scale from 0 to 4. The MathBERT model is fine-tuned for this task, with pairs of questions and student responses. Additional

information such as the scale of the grade and example answer with grade is also passed to this model for in-context meta-learning.

## 5.2 Math word problem question answering

**Symbolic and logical 'rule-based' approaches.**    The traditional approach to automatically solving math-word problems is rule-based, where a word problem is first converted to some form of symbolic (*e.g.,* logical) data representation which is then passed to an inference algorithm that generates answer values. The STUDENT system (Bobrow, 1964) was a problem-solving system for algebraic problems implemented in LISP that used a sequence of string transformation rules to generate problem instance data representations. While used in the early days, this basic approach remains effective.

For example, the AiFu model (Liu *et al.*, 2019) proposed at the SemEval task uses a translator to convert the word problem into an internal representation using rule-based templates. This representation is based on assertional logic where mathematical objects are formalized as constants, variables, concepts, functions, or relations. For example, *"the integer x equals to 3"*, is transformed to

```
Integer(x), Equal(x,3).
```

This representation is then sent to an encoder to extract a mathematical representation that can be passed to the Satisfiability Modulo Theories (SMT) solver Z3[2] to produce a solution. This encoder is also rule-based, with two templates for unification: one for different concepts and operators, and the other for the new entities. This model had the highest accuracy among participating teams in the SemEval task.

Generating and computing solutions from operator trees (OPT) is also common for automatically solving math word problems, as illustrated in Figure 5.2. These trees are both used to represent the math-word problem:

> *There are 3 boys and 5 girls in a group. Each person wants to buy 9 pencils. How many pencils do they need to buy altogether?*

---

[2]https://github.com/Z3Prover/z3

(a) Operator Tree          (b) Equation Tree

**Figure 5.2:** OPT for '$(3 + 5) \times 9$', and variable-augmented OPT ('equation' tree introducing variable $\chi$) for math-word problem asking for the expression value.

In these approaches, the problem statement in natural language is computed to an OPT, with the $n$ quantities (constants) defined in the question appearing at leaves in the tree. Answers are produced by applying operations bottom-up (*i.e.,* evaluating the expression). In one approach the OPTs use a contrained set of four basic binary operations, with quantities at the leaves. (Roy and Roth, 2015). A similar approach is the ALGES system (Koncel-Kedziorski *et al.*, 2015), where the generated OPT includes a variable for the answer $\chi$ in the tree, attached to an equivalence operator. They call this modification an *equation tree.*

For building the operator trees, it is essential to correctly select relevant variables quantifiers for inclusion in the tree (*e.g.,* $\forall x$, $\exists x$). The approach by Roy and Roth trains an SVM binary classifier to decide which quantifiers are required to answer the math-word problem and adds these to the leaves of the tree. Then another classifier determines the lowest common ancestor (LCA) node between pairs of quantifiers to capture constraints. Alternatively, the ALGES system generates all possible equation trees without dropping any irrelevant quantifiers, and then uses Integer Linear Programming to score the likelihood of each tree using local and global discriminative models.

**Earlier machine learning approaches.** Several statistical machine-learning approaches have been investigated for math-word problems. Like the earlier rule-based approaches, they produce a data representation for a word problem after which the answer is generated from this representation algorithmically. The ARIS system (Hosseini *et al.*,

Sarah had 5 black pens and 3 blue pens. She gave some of her black pens to Jack. Jack has 8 black pens. Sarah has 3 black pens left. **How many black pens did Jack have?**



**Figure 5.3:** Overview of ARIS (Hosseini *et al.*, 2014) approach. States are generated based on sentences, shown in blue boxes. Each state, has one or two containers shown as black boxes.

2014) was developed for arithmetic problems containing addition and subtraction. In this model, steps toward the final solution are represented as a state sequence, with one state per sentence. States are entity triples with attributes, where entities have quantities (*i.e.,* constants or variables), types, and attributes associated with them. These entities are stored as sets in containers associated with subjects and objects in the question.

As an example, consider this math word problem with two subjects with their own quantity containers (Sarah and Jack):

> *Sarah had 5 black pens and 3 blue pens. She gave some of her black pens to Jack. Jack has 8 black pens. Sarah has 3 black pens left. How many black pens did Jack have?*

The ARIS state representation for each sentence is shown in Figure 5.3, which is produced using the Stanford NLP library. $J_0$ shown in the state for the second sentence represents the question quantity, which here is 6 (black pens that Jack originally had). Determining this value depends also on the implicit variable $L_1$, the number of pens that Sarah gave to Jack. After the sentence state representations are generated, constraints in the question are used to compute the question answer (*e.g.,* for $J_0$).

State transitions are identified by verbs in the sentences, classified using an SVM classifier with WordNet-based and other features. Verbs

are categorized into three groups:

1. observation (initialized quantity),
2. positive (quantity increased), and
3. negative (quantity decreased).

For sentences with two containers (*i.e.,* subjects/objects with entities), four other categories are considered:

4. positive transfer (quantity is transferred from the second to the first container),
5. negative transfer (quantity is transferred from the first to the second container),
6. construct (quantity increased in both containers), and
7. destroy (quantity decreased in both containers).

More recently, deep learning approaches are among the strongest approaches used to-date for math-word problems. The earliest attempt, Deep Neural Solver (DNS) (Wang *et al.*, 2017b), used a seq2seq model to translate problem statements into equations (tokenized infix OPT representations) by embedding question text into a vector, and then generating (*decoding*) a formula expression starting from the question vector. This generated expression is then evaluated to produce an answer.

For a given word problem, first numbers appearing in the problem are mapped to variable tokens and stored in a dictionary (*e.g.,* { $(n_1, 3), (n_2, 5), (n_3, 9)$}), and numbers are replaced by their associated variable in the text. This modified text is passed to a seq2seq model to generate an equation over those variables, variables are replaced by their values in the number mapping dictionary, and then the expression value is computed from the underlying OPT. Using the same example as earlier for the earlier rule-based operator tree-based approaches, the architecture of this model is shown in Figure 5.4.

After the DNS model was introduced, there were several attempts to design networks that better capture structural in the question in order to improve generated OPTs. The multi-encoders and decoders model (Shen and Jin, 2020), uses both seq-to-seq (sequential) and graph-based encoder/decoders. Two graph encoders are defined using GraphSage (Hamilton *et al.*, 2017):

1. a dependency parse tree capturing relationships between words in

**Figure 5.4:** Seq2seq model proposed by (Wang *et al.*, 2017b) for math-word problems.

the sentence. Iinitial token embeddings come from the sequence-based encoder.

2. graph-encoder used for numerical comparisons. Nodes in the graph are numbers in the problem with relations $>$ or $\leq$, to capture contraints ignored by models like DNS when mapping values to tokens.

For decoding, sequence-based and tree-based decoders are used. The decoders produce operation and value sequences that correspond to OPT traversals (referred to as abstract syntax trees (AST)) that can be used to directly compute the final answer (*e.g.,* using a simple stack-based algorithm). For the solution equation $(n_1 + n_2) \times n_3$, the sequence decoder will generate a preorder traversal of the OPT in Figure 5.2 with arguments before operations $(n_1 n_2 + n_3 \times)$, and the tree decoder generates a postorder traversal giving operations before arguments $(\times + n_1 n_2 n_3)$. The final output operation sequence is selected by maximum decoder likelihood. This model achieves higher accuracy than the DNS model on the Math23K dataset using 5-fold cross-validation (76.9 vs. 58.1).

**Transformer approaches.**   Solving math word problems automatically has gained increasing attention since the development of generative language models that use transformers, a deep-learning model first proposed in 2017 (Vaswani *et al.*, 2017). Transformers are able to integrate surrounding context in token sequences by consulting all other

tokens embeddings in the input over a series of stages. This produces contextually-enriched embedding vectors for each input token that dramatically improved their usefulness in prediction and generation tasks.

It can be helpful to think about information use and flow to obtain an intuition for the improved effectiveness and flexibility of transformers over earlier embedding models. Thinking back to information tasks in Ch. 1, with a transformer the input tokens act as an initial information source, from which we want to produce a new information source that contains 'contextualized' token vectors. The generated tokens are shaped by token co-occurrence statistics (*i.e.,* a language model) seen in a large *collection* of sources used for training. In terms of information tasks performed by the transformer:

**Retrieve**

- **Consult:** All current token vectors are examined and weighted using a *self-attention mechanism* as input for generating updated vectors. Earlier models did not consult the entire input sequence: more information is used.
- **Query (Designers/Users):** annotated input token streams allow (1) directing parameter learning, and (2) 'programming' outputs in query *prompts*. More flexible with greater expressiveness than earlier 'contextualized' embedding models.

**Analyze**

- **Annotate (Designers/Users):** input annotated with additional tokens, *e.g.,* start/end of sections (*e.g.,* `[CLS]` for first 'classification' token, `[SEP]` ending sentences/sections), responses start (*e.g.,* `[A.]`). *Positional encoding* vectors added to mark input position (*e.g.,* integer enumeration) and/or token relationships. These annotations are more detailed and expressive than for earlier models.
- **Index:** learned network weights represent language statistics in a ('pre-')training corpus, and can be used to produce dense vector indexes.

**Synthesize**

- **Apply:** Transformed token embeddings are generated (*communicated*) over a series of steps, by *applying* network weights

to produce updated token vectors (*e.g.,* in 12 steps). This allows vector transformations to be gradual, using 'depth' to improve generalization (a known property of deep nets). Training tasks apply network weights to make decisions (*e.g.,* guess masked tokens) and then update those weights using backpropogation.

- **Communicate:** the model outputs the contextualized token vectors in a tensor, providing a new information source, along with learned network weights that specify the embedding function represented in the network.

As described in Ch. 2, 'pre-training' transformers on large text corpora is the norm, followed by 'fine-tuning' to produce outputs for specific tasks. 'Pre-training' (*i.e.,* initial language model learning) generally involves imitative games that require predictions after token sequence manipulations (*e.g.,* masked tokens and token re-orderings), while fine-tuning requires replacing the output (decision) layers of the network for the new imitative game(s) of the specific task the network will be used for (*e.g.,* classification or generation tasks).

As for retrieval, we again find here that the token representations used for math with transformers are important. Simple tasks such generating answers to questions about the addition and subtraction of two numbers (e.g., 'what is 52 plus 148?') were explored (*e.g.,* using a T5 a seq-to-seq model (Nogueira *et al.*, 2021)). This was motivated by earlier work where they found that for tasks such as reporting the maximum value of a list of numbers, a BERT model obtained an accuracy of 52% (Wallace *et al.*, 2019). As described earlier, this poor performance was caused partly the default tokenization if WordPiece, which prevented correctly encoding numbers.

The ability of the BERT model to generate answers in response to math word problems provided in the input was explored using the AQuaA-RAT(Ling *et al.*, 2017) rest collection. The focus of this work was exploring fine-tuning approaches (Piękos *et al.*, 2021), including a proposed Neighbor Reasoning Order Prediction (NROP) coherence loss. This considers whether the steps in a rationale for an answer are in their original order, or have been swapped; fine-tuning with this task and loss function improved accuracy by roughly 10% after fine-tuning. This work

inspired subsequent work on improving outputs describing mathematical reasoning, *e.g.,* in producing answers to math word problems.

MWP-BERT (Liang *et al.,* 2022) for example, used the BERT and RoBERTa models and used three families of fine-tuning tasks/objectives.

**Self-supervised (*i.e., data reconstruction/description*)**: for problem descriptions: (1) masked language modeling, (2) number counting (*i.e.,* the number of quantities in a word-problem), and (3) number type grounding (*e.g.,* integer, non-integer).

**Weakly-supervised:** for problems and corresponding answers: (1) answer value type prediction (*i.e.,* discrete or continuous), (2) context-answer type comparison (*i.e.,* predict if problem description quantities are same type as the answer), and (3) number magnitude comparison (*i.e.,* predict the relative size of the answer to the problem statement quantities).

**Fully-supervised:** (1) operation prediction, and (2) tree distance prediction. In operation prediction, the goal is to predict the operator between two quantity nodes in the solution OPT from five operation types (classes): $\{+, -, \times, /, \char`\^\}$. Tree distance prediction estimates differences in the depth of two quantity (numbers) in the solution tree.

It was found that these different training objectives in isolation produce similar accuracy on Math23K dataset, with the fully-supervised objectives providing the highest accuracy. The MWP-BERT is fine-tuned on MathQA dataset and on Math23K (with 5-fold cross-validation) it achieves an accuracy of 82.4, higher than 58.1 of the DNS model. The same pattern can be observed for Math23K test set, where MWP-BERT has an accuracy of 96.2, higher than 53.6 of the DNS model.


## 5.3   Large Language Models (LLMs) & Mathematical Reasoning

The next generation of models for solving math-word problems rely on large language models, which are (roughly) very large transformer models trained on very large amounts of data embedded within a text generation system (*e.g.,* a recurrent neural network). The capability of large language models to solve advanced math questions is one of the important aspects studied when a new LLM is developed, and

math-word problems are commonly studied for the evaluation of LLMs. Open-source models such as Mistral (Jiang *et al.*, 2023) and LLaMA (family) have been used for fine-tuning math-specific language models, helping improve reponsiveness of the model to math-focused *system prompts (messages)*.

Through prompts the model can be given context and instructed on producing outputs, roughly speaking a way to indirectly 'program' outputs as described in our discussion of transformers in the previous section.[3] While the general LLMs have been studied for their capabilities in solving math problems, they are not trained specifically for this task. Therefore, several LLMs, use general models fine-tuned for mathematical tasks. The common approach in domain-specific models is using a base LLM, fine-tuning it with a math-specific dataset (usually a new dataset), and improving the reasoning observed in outputs using different techniques.

**Techniques for improving LLM math question answers.** Here are four common methods for improving answers to math questions given to LLMs vs. directly proving a word problem in a prompt:

**Instruction tuning:** uses a labeled dataset of *(prompt, response)* pairs for additional training of the model, *i.e.,* fine-tuning. This is distinct from prompt engineering, which prompts a fixed LLM multiple times and then a preferred answer is selected; no training is involved. Often done with state-of-the-art language models such as GPT-4 to generate high-quality training data.

**Chain-of-Thought (CoT) prompting (Wei *et al.*, 2022):** prompts include several examples where a desired step-by-step reasoning process is explicitly demonstrated. The goal is to exploit the LLM tendencies to mimic patterns presented in the prompt, in anticipation a similar reasoning step pattern will appear in responses, and help with producing correct answers. COT is disctinct from few-shot prompting (see Figure 5.5), where only questions and final answers are provided within examples included in the prompt.

---

[3]Note that information wise for a user, *prompt=query*: these are requests for generated information sources.

**Figure 5.5:** Standard prompting(left) vs. Chain-of-Thought (right) adopted from (Wei *et al.*, 2022). Grey boxes indicate the model's input.

**Program-of-Thought (PoT) (Chen *et al.*, 2022b):** similar to CoT, where instead of natural language computer programs represent the thought process formally. PoT frees the LLM from having to approximate equations for natural language questions, instead providing explicit computational steps. In a zero-shot setting (*i.e.,* giving no example input and output in the prompt), PoT produces a 12% improvement over CoT for math-word problem datasets.

**Program-Aided Language models (PAL) (Gao *et al.*, 2023):** generates reasoining using interleaved natural language (as in CoT) and programming language statements (per PoT). The final solution produced is a program that is evaluated by an interpreter.

**Math-specific LLMs.**   Most research on math-specific LLMs revolves around improving input token representations and prompts, and expanding datasets for fine-tuning base general models. One of the earliest math-specific LLMs is Minerva (Lewkowycz *et al.*, 2022). This LLM is based on the PaLM (Pathways Language Model) model, and is additional trained on a 118GB dataset of scientific and mathematical data. Minerva can correctly generate LaTeX formulas, which is challenging for general LLMs. The main novelty of this model is correctly handling mathematical formulas; with proper tokenization, formulas such as $E = mc^2$ are processed as `$E=mc^2$`, and not the single token `Emc2`.

To improve answers Minvera uses CoT, prompting with several step-by-step solutions to questions before posing question to be answered. It also uses majority voting to select the most likely final answer, from

multiple generated solutions. From manually analyzed samples, incorrect reasoning, incorrect calculation, and question misunderstanding were among the most frequent issues. For example, for the question:

*If $\sqrt{400} = \sqrt{81} + \sqrt{n}$, then what is the value of n?*

this model makes a mistake in interpreting the question. It infers that $400 = 81 + n$, and then correctly provides 319 as the final answer, but to the wrongly inferred formula.
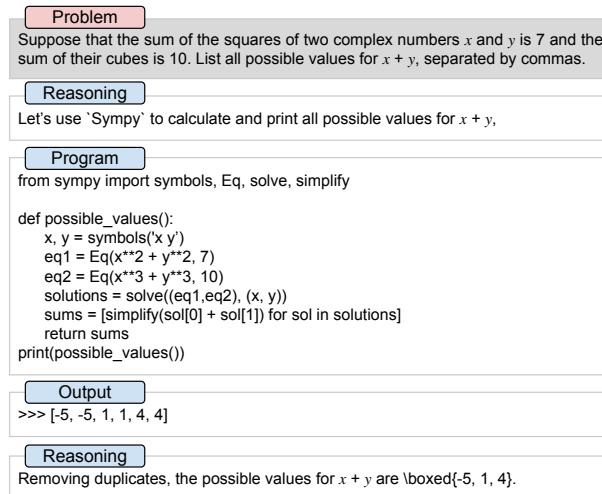
A similar model using CoT is LLEMMA (Azerbayev *et al.*, 2024). The model is built atop LLaMA-Code, and fine-tuned on a new dataset `Proof-Pile-2`. This dataset is a 55 billion token dataset of mathematical and scientific documents, taken from three large collections: arXiv, open-web-math, and algebraic-stack (new in this dataset). By fine-tuning with this dataset, LLEMA provides better effectiveness on tasks such as the ones introduced in GSM8K, and MATH compared to Minerva.

MAmmoTH (Yue *et al.*, 2023) model uses hybrid CoT and PoT for building a new dataset. This model is based on Meta's LLaMA-2, and Code LLaMA and fine-tune these models on their new dataset `MathInstruct` for instruction-tuning. The data is compiled from 13 math rationale datasets (7 existing), focusing on the hybrid use of chain-of-thought (CoT) and program-of-thought (PoT) rationales. While for existing datasets natural text reasoning existed for CoT, this dataset uses GPT-4 to generate programs, and the values generated by programs are verified with ground truth to ensure quality. The dataset contains 260K (instruction, response) pairs with 72% using COT, and 28% using PoT. Figure 5.6 shows an example of CoT and PoT pairs. Analysis by the authors revealed that PoT generally provides better results than CoT, mostly for open-form questions, with algorithmic reasoning being more effective for complex math problems.

ToRA (Tool-integrated Reasoning Agents) (Gou *et al.*, 2024) introduces a new representation for reasoning that interleaves natural language and reasoning using external libraries. For a given math question, ToRA first generates reasoning in natural language, which continues until a program library is better suited (e.g., equation solving). Then, ToRA generates a program from the natural language reasoning,

> **PoT**
> If circles x and y have the same area and circle x has a circumference of 12 π, half of the radius of circle y is : ? Let's write a Python program to solve it.

```
n0 = 12.0
t0 = n0 / 2.0
answer = t0 / 2.0
print(answer)
```

> **CoT**
> There are 3 boxes of cereal. One box holds 14 ounces of cereal. Another box holds half the amount of the first box and 5 ounces less than the third box. How much cereal is there in all 3 cereal boxes?

> First = 14 oz Second = (1/2) * 14 = 7 oz Third = 7 + 5 = 12 oz 14 + 7 + 12 = 33 oz There are 33 ounces of cereal in those 3 boxes. The answer is 33

**Figure 5.6:** Sample instructions and responses from MathInstruct dataset. The grey boxes show the PoT and CoT instructions and the white boxes represent the responses.

> **Problem**
> Suppose that the sum of the squares of two complex numbers $x$ and $y$ is 7 and the sum of their cubes is 10. List all possible values for $x + y$, separated by commas.

> **Reasoning**
> Let's use `Sympy` to calculate and print all possible values for $x + y$,

> **Program**
> ```
> from sympy import symbols, Eq, solve, simplify
>
> def possible_values():
>     x, y = symbols('x y')
>     eq1 = Eq(x**2 + y**2, 7)
>     eq2 = Eq(x**3 + y**3, 10)
>     solutions = solve((eq1,eq2), (x, y))
>     sums = [simplify(sol[0] + sol[1]) for sol in solutions]
>     return sums
> print(possible_values())
> ```

> **Output**
> >>> [-5, -5, 1, 1, 4, 4]

> **Reasoning**
> Removing duplicates, the possible values for $x + y$ are \boxed{-5, 1, 4}.

**Figure 5.7:** ToRA mathematical reasoning approach, interleaving natural language reasoning with program-based tool use. (adapted from (Gou *et al.*, 2024))

and the output further processed for adjustments, sub-task solving, and answer finalization. This process is continued until the final answer is represented using the LaTeX "\boxed{}" command. Figure 5.7 shows an example of this process.[4] To fine-tune LLaMA-2, TORA-CORPUS was introduced by using GPT-4 to synthesize trajectories on the GSM8k and MATH training sets.

---

[4]Note the similarity to the types of answers seen in CQA forums like Math Stack Exchange, where code is part of a larger answer narrative.

**Improving data sets through augmentation.** More techniques are being introduced to expand and/or improve existing datasets, often through techniques that transform existing samples to produce additional ones (*i.e.,* data augmentation). The MetaMathQA dataset is created for fine-tuning MetaMath (Yu *et al.,* 2024). This dataset is based on two previous datasets, GSM8K and MATH. Three main techniques are used for data augmentation. First, is answer augmentation to generate more reasoning paths towards solutions to math questions. For questions, a few-shot CoT promoting is used by appending questions to a few reasoning examples, and fed to LLM for generating its reasoning path and answer. The final answer is compared against the ground-truth, and the wrong answers are filtered out.

A secon data augmentation is question rephrasing using GPT-3.5-Turbo, with few-shot rephrased examples, and the prompt:

> *You are an AI assistant to help me rephrase questions. Follow the given examples.*

To check the quality of rephrasing, the answers generated for the original question and rephrased one were compared, and the accuracy is reported as 76.30% for rephrased questions, comparable to 80.74% for the original questions.

MetaMath uses the idea of backward reasoning to enhance reasoning ability; starting with a given condition and thinking backward to answer the question. To increase the number of questions, in this dataset, a number in question is masked with 'X' and the LLM is asked to predict the value when the answer to the question is provided. The difference between backward and forward reasoning is that backward reasoning starts with the answer, and generates reasoning steps to provide the value of the masked number. After merging all the augmented data, this dataset is then used to fine-tune LLaMaA-2. The experimental results show having augmented data for fine-tuning improves the effectiveness of the fine-tuned model on just the original questions.

While the majority of these models achieve their highest effectiveness with large number of parameters in the model, Orca-Math (Mitra *et al.,* 2024) brings attention to fine-tuning smaller language models for math reasoning. Similar to MetaMath model, first a new dataset is introduced

and then fine-tuning is done over three iterations. Orca-Math-200K dataset is constructed by using previous datasets such MathQA, and GSM8K with a total of 36,217 problems. This dataset was then expanded to create multiple math word problems from each initial problem. To this end, GPT-4-Turbo is used with few-shot examples, and then prompted to first convert the question into a statement by using the answer to the question. Then, from the statement, it creates a new word problem. The solutions for the new problems are also generated by GPT4-Turbo. A total of 120,445 new problems were generated with this approach.

Finally, another augmentation is used within this dataset, for increasing the difficulty level of existing questions. Using two agents *Suggester* and *Editor*, first Suggester proposes methods of increasing the difficulty by suggesting approaches like adding more variables, or increasing certain values. Then, the editor modifies the original question based on the suggester's recommendation. This is done in two rounds, each time increasing the difficulty further. New problems that GPT4-Turbo answers in more than 1800 characters were dropped. 37,157 new problems were collected from this step.

For fine-tuning Mistral-7B, three steps are used for augmentation:
1. instruction-fine-tuning is applied by feeding pairs of questions and answers.
2. Using the teacher-student learning model, for each question, four answers from the first iteration are generated (as student's answers).
3. These answers were then compared against those generated by GPT-4 Turbo (the teacher) as the ground-truth, with correct answers considered as positive, and the wrong were used as negative samples for further fine-tuning.

The correctness of the answers was decided by GPT4-based exact-match prompting. While in the second iteration, the samples were generated with a supervised fine-tuned model, in the last iteration, a similar approach was used for generating samples, by considering two approaches: Direct Preference Optimization (DPO) (Rafailov *et al.*, 2023) and Kahneman-Tversky Optimization (KTO) (Ethayarajh *et al.*, 2024). Both approaches are used for fine-tuning LLMs to align with preferences. With the fine-tuned model, Orca-Math archives pass@1

metric of 86.81, higher than larger models such as MAmmoTH and MetaMath. pass@1 indicates the ratio of problems correctly answered.

Despite the advances in LLMs, their capabilities for math information retrieval need further investigation. Some recent researchers aimed to explore this by comparing the results from LLMs against human-generated and assessed answers.

The use of LLMs as both answer generators and re-rankers have been explored using ARQMath's answer retrieval task. Satpute et al. (Satpute *et al.*, 2024) considered both general LLMs like GPT-4 and math LLMs like ToRA to generate answers for ARQMath topics. These answers were then used as queries to search for relevant answers. This can help by transforming the query into a format more similar to sources in the collection, as the query and collection are expressed as answers. For each question, an embedding (using `BERT_cocomae`) of the generated answer was used to search dense embeddings for ARQMath answers, and cosine similarity was used to find the most similar answer to the generated answer. This approach is currently the state-of-the-art for ARQMath answer retrieval task, with nDCG′ of 0.486 (vs. 0.464 of BERT_cocomae).

In another study, the applications of general LLMs, (LLaMA-2 and Orca-2) were studied for three tasks in math information retrieval: relevance assessment, data augmentation, and re-ranking (point-wise) (Mansouri and Maarefdoust, 2024). This study was done on ARQMath dataset and for each task, an appropriate *system message* is used; for example for assessing the relevance of an answer to a math question, the system prompt was created based on ARQMath assessment protocols as:

> *You are a math professor who will assess the relevance of an answer to a given math question.*

The results of this study revealed that while general LLMs are not yet suited for relevance assessment or re-ranking, the data augmentation from Orca-2 which is useful for expanding the ARQMath training set for use in fine-tuning neural math answer retrieval systems. The data augmentation process is performed by generating additional relevant answers for each training topic using LLMs.

Let's analyze the given series: $\sum_{x=1}^{\infty} \frac{(x-1)}{2^x}$
First, let's rewrite the series as: $\sum_{x=1}^{\infty} \frac{x}{2^x} - \sum_{x=1}^{\infty} \frac{1}{2^x}$
Now, let's focus on the first term: $\sum_{x=1}^{\infty} \frac{x}{2^x}$
This is known as the geometric series with a linear term:

$$\sum_{x=1}^{\infty} rx = \frac{r}{(1-r)^2}$$

In our case, $r = \frac{1}{2}$, so: $\sum_{x=1}^{\infty} \frac{x}{2^x} = \frac{\frac{1}{2}}{\left(1-\frac{1}{2}\right)^2} = 2$
Now, let's focus on the second term: $\sum_{x=1}^{\infty} \frac{1}{2^x}$
This is a geometric series with a common ratio of $\frac{1}{2}$: $\sum_{x=1}^{\infty} \frac{1}{2^x} = \frac{\frac{1}{2}}{1-\frac{1}{2}} = 1$
So, the original series becomes: $2 - 1 = \boxed{1}$
Therefore, the series converges to 1.

**Figure 5.8:** Llama-3-8B-Instruct generated answer for question "can prove it converges but I don't know at what value it converges. $\sum_{x=1}^{\infty} \frac{(x-1)}{2^x}$" from ARQMath-2 Answer Ret. The system prompt is "You are a math professor who will answer math questions, providing step-by-step solutions.".

**LLM limitations.**    To illustrate some limitations of current LLM models, let's consider an ARQMath question. We use the system prompt defined as follows:

> *You are a math professor who will answer math questions, providing step-by-step solutions.*

Using this prompt for Llama-3-8B-Instruct, for math question from ARQMath-2 "can prove it converges but I don't know at what value it converges. $\sum_{x=1}^{\infty} \frac{(x-1)}{2^x}$", leads to the answer shown in Figure 5.8.
    While the given answer is related to the question, there is an error in correctly generating the fourth offset summation, where $rx$ should have been $xr^x$. Moreover, as a "math professor", the splitting of the series as:

$$\sum_{x=1}^{\infty} \frac{(x-1)}{2^x} = \sum_{x=1}^{\infty} \frac{x}{2^x} - \sum_{x=1}^{\infty} \frac{1}{2^x}.$$

is not an obvious step and needs verification/reasoning about convergence of each of those series. So, when we say it is a "step-by-step" solution, then the professor has ignored this crucial step.
    The *system prompt* can guide the responses for LLMs. Changing the prompt for LLaMA-3 to "You are a math professor who will answer

math questions, providing a hint toward the solutions, not the whole solution." for the same ARQMath question, the response is:

> Here's a hint to help you get started:
>
> You can try to recognize this series as a variant of a more well-known series. In particular, consider the geometric series $\sum_{x=1}^{\infty} r^x$, where $|r| < 1$. Can you see how the given series can be rewritten in a similar form?
>
> Also, think about the fact that the series is a telescoping series, meaning that many of the terms cancel out when you sum them up. Can you identify which terms cancel out and how this might help you find the sum?
>
> Take it from here and see if you can prove that the series converges and find its value!

The first part of the hint is a good start for the student. In the sense that it narrows down the scope where the student can focus or get some ideas. However, the second part of the hint is completely wrong; this is not a telescoping series. This part of the hint is very misleading and will not help the person who is asking the query. A *telescoping series* is the one where each term of the series can be written as a difference of two consecutive terms of a sequence. In the sense,

$$\sum_{n=1}^{\infty} (a_{n+1} - a_n) = (a_2 - a_1) - (a_3 - a_2) + \cdots .$$

So, the terms start canceling off and the sum is easy to find. But the series given in the query cannot be written in this form. So, this second part of the hint does not help.

More and more LLMs specific to math continue to be introduced. All these models aim to provide better reasoning strategies. Currently, Chain-of-thought with natural language and Program-of-thought with programs generating outputs are commonly used. Probably in future work, better reasoning techniques will be applied. The current models each need to individually create their dataset for fine-tuning the base LLM model (e.g., LLaMA-2). These datasets are usually created using a pioneer LLM (e.g., GPT-4).

However, recently several researchers have warned against this approach, suggesting that model effectiveness may decrease when AI-generated data is used for training (Shumailov *et al.*, 2024a; Shumailov *et al.*, 2024b).

## 5.4   Applying math-aware search and QA to theorem proving

Automated theorem proving aims to explore the application of computers in proving mathematical theorems. This is actually one of the earliest topics of interest for computer science. In the early 1960s, scientists started exploring the use of machines for theorem proving. for the quantification theory (Davis and Putnam, 1960; Davis *et al.*, 1962) using deduction rules to prove assertions.

Informal theorem proving refers to the way that humans approach proving theorems using reasoning through notation and natural language, likely with some missing details (*e.g.,* assumed definitions) and skipped computational steps, for example. In contrast, formal theorem proving represents theorems in a machine-readable format, making verification by logical rules possible, at the cost of more information being explicitly stated (*e.g.,* all variable types, definitions for all operators, etc.).

Converting informal to formal proof setps is referred to as *Auto-formalization.* The Mizar[5] language is commonly used by mathematicians as a formal language for writing definitions and proofs. Using Mizar, Wang et al. explored converting (translation) of informal LaTeX-written text into formal Mizar language (Wang *et al.*, 2018). This work explored different seq2seq architectures, with LSTM and attention providing the highest BLEU score for this translation.

Researchers have also explored other formalization language, including applying large language models to generate statements in Codex (Chen *et al.*, 2021b) has been studied recently (Wu *et al.*, 2022). This translates statements in natural text into formalized theorems for the interactive proof assistant Isabelle (Wenzel *et al.*, 2008). As an example, the system translated "Prove that there is no function $f$ from the set of

---

[5]https://mizar.uwb.edu.pl/

non-negative integers into itself such that $f(f(n)) = n + 1987$ for every $n$" perfectly to Codex as:

```
theorem
    fixes f :: "nat \<Rightarrow> nat"
    assumes "\<forall> n. f (f n) = n + 1987"
    shows False
```

Another task for theorem proving is retrieving useful lemmas that will help with proving steps, known as *premise selection*. This is a form or search problem, where relevance is defined in terms of suitability for proving a specific conjecture. The problem is defined as (Alama *et al.*, 2014):

> *Given an Automated Theorem Provers and a large number of premises, find premises that are useful to the prover for constructing a proof for a new conjecture.*

DeepMath (Irving *et al.*, 2016) is one of the earliest works to apply deep learning for this task. Conjecture and axiom sequences are embedded separately, concatenated, and then passed to a fully connected neural network for predicting the usefulness of the axiom. Embeddings are at character level for formulas, and word-level for statements defining symbols. The convolutional network model FormulaNet (Wang *et al.*, 2017a) used a similar idea and applied graph neural networks. Using formulas in higher-order logic (Church, 1940), each formula is first parsed into an OPT: internal nodes represent a quantifier or a constant or variable function, and leaf nodes represent variable or constant values. Edges connect a quantifier to all instances of its quantified variables. After creating the tree, a merging step is applied to merge leaf nodes representing the same constant/variable. Finally, a unification technique is used to replace variable names with 'VAR' and function names with 'VARFUNC'. After building the graph, convolution or message passing is applied to get node embeddings. These embeddings are used with max-pooling to form an embedding for the graph.

As theorems are built upon existing mathematical knowledge, a graph representation of mathematical concept statements such as lemma, and definitions is a common approach for this task. One way of building

this graph is to use statements as nodes and ordered edges from node $s_1$ to $s_2$ if there is statement 1 is a premise of statement 2 (Ferreira and Freitas, 2020a). With this definition of a graph, the problem can be viewed as link prediction, for which Deep Graph Convolutional Neural Network (DGCNN) architecture was applied (Zhang *et al.*, 2018). The textual content embedding of each node in this work is encoded using Doc2Vec (Le and Mikolov, 2014b) model with mathematical concepts being encoded as linearized trees, with every sub-expression represented as a token.

For example, the formula $(x + y) \times c$ is represented as a sequence of tokens for subexpressions {'$x$', '$y$', '$(x + y)$', '$(x + y) \times c$'} (similar to what is used to generate subexpression tokens for the WikiMirs search model discussed earlier). The authors later introduced STAtement Representation (STAR) cross-modal representation (Ferreira and Freitas, 2021), treating mathematical formulas as natural language. Each statement is viewed as a combination of words and formulas. However, for embedding, they proposed two separate self-attention layers, one for formulas and one for words. The output of the self-attention layers is then concatenated and passed to a Bi-LSTM to get the final representation of the statement. To find the relatedness score between conjecture and premises, a Siamese neural network was applied.

Generative and large language models have recently been applied for premise selection. LeanDojo (Yang *et al.*, 2023) is an open-source toolkit based on the Lean[6] programming language, that introduces ReProver (Retrieval-Augmented Prover). ReProver is a language model-based prover, augmented with retrieval for selecting premises. Given the initial state of proof, it retrieves a set of useful premises (set at 100) using a Dense Passage Retriever. These premises are concatenated to the initial state and passed to a fine-tuned ByT5 (Xue *et al.*, 2022) model to generate steps toward the proof.

---

[6]https://lean-lang.org/

## 5.5 Summary

We conclude our overview of mathematical information retrieval tasks and models used to address them in this chapter. The tasks discussed were math-aware search (*i.e.,* searching with formulas and text), math question answering, and applications of dense representation and search in mathematical theorem proving.

For math-ware search, a variety of sparse, dense, and sparse-dense hybrid models are used. The models either search text and formula representations independently, or using a single index in which formulas and text have the same representation. Results from searches against multiple indexes (of whatever type) are often combined using some form of rank or rank score weighting (*e.g.,* linear combinations, learning-to-rank algorithms, or reciprocal rank fusion). For both independent and unified formula/text representations, an important question is whether there are more effective representations than existing sequential (*e.g.,* transformer token input sequence) or graph-structured (*e.g.,* AMR tree with embedded OPT) ones; graph representations for formulas+text have been less explored in particular. Improved selections of patterns in these representation for representation in inverted indexes for sparse retrieval, and improved pattern selection and statistical embedding methods for dense vector representations are also important future directions.

Regarding math question answering, while deep learning models have been gaining popularity for solving math-word problems, the input representation used for problems has remained largely the same, based on annotated token sequences in transformers. For formulas in input questions, explicit representations of formula structure using symbol layout and/or operator trees is unexplored. Also, while graphs for dependency parsing are used for text, other graph-based text representations such as Abstract Meaning Representation graphs have yet to be investigated. We envision that in the coming years, researchers will work towards developing better mechanisms for imitating human reasoning better, and for reasoning to reflect human behavior, to insure the interpretability of answers and their rationales.

# 6

---

# **Final Thoughts and Research Opportunities**

---

In this chapter we draw our discussion of mathematical information retrieval to a close. As part of this, we consider directions in which new insights might be found in the future. To do that we come back to the three key concerns we keep returning to throughout the book:

1. the people that IR systems serve and interact with,
2. the information tasks that IR systems perform on behalf of people,
3. how systems and associated processes are evaluated.

We then consider an important question in our current technical climate – from the perspective of an indiviudal, what are the limitations of MIR systems? Finally, we close with some brief comments.

## 6.1   What is next for MathIR?

To try and characterize the space of opportunities in MathIR, we present two views. The first shown in Table 6.1 is information needs for the people that play a role in the creation, use, and evaluation of mathematical information retrieval systems: the users, assessors, and designers/researchers. Improving available information of the types listed, or making these types of information more readily available or easy to use present research opportunities. For example, improved interfaces for users, as-

**Table 6.1:** Information needs for 'the people' involved in mathematical information retrieval. Many of these human information needs motivate research and development in IR generally, in terms of new information that would be useful, and easing access and reuse (*i.e., application*) of existing information.

| Users | Assessors | Designers/Researchers |
|---|---|---|
| *Math info* | *Math info* | *Math info* |
| Query language | Annotation tool use | Useful query mechanisms |
| Query/results interface | Relevance criteria | Evaluation metrics |
| Organizing sources | Assessment strategy | Pooling / protocols |
| Refinding sources | (pragmatics) | Annotation tools/proto-cols |
| Communicating info in sources | | Machine learning tech-niques |
| Stopping criteria (prag-matics) | | Retrieval techniques |
| | | Data augmentation |
| | | Test collections: creation, use, differences |
| | | Data representations (*e.g.,* for formu-las) |
| | | Model behaviors and comparisons |
| | | Implementation: plan-ning and execu-tion |
| | | 'Focus and fences' (prag-matics) |

sessment strategies that improve reliability/agreement without greatly increasing assessment effort, and improved programming libraries and frameworks for machine learning and information retrieval systems for researchers are helpful. Study and innovations providing *new* informa-tion for the needs shown are particularly valuable from the research perspective.

We also present a second view organized around the information tasks in our 'jar' model in Table 6.2, and whether they relate to human-computer interaction, system design and implementation, or evaluation of interaction (front-end) and retrieval (back-end) system components. The majority of items appearing in the table are mentioned elsewhere in the book. Here they are organized in terms of whether they are involved in retrieving information sources, analyzing information sources, or synthesizing new information sources.

At the current moment opportunities in machine learning, includ-

**Table 6.2:** Future Directions for Mathematical Information Retrieval. Opportunities are organized by information task type, as described in Chapter 1.

| Inf. Task | Interaction | Systems | Evaluation |
|---|---|---|---|
| **Retrieve** | | | |
| Query | Search interfaces<br>Multi-modal search (video, docs, other)<br>LLM prompting | Formula autocomplete<br>Text/formula autocomplete | User studies |
| Consult | Formula navigation<br>Formula annotations views<br>Linked math/text | Attention mechanisms<br>Query logs | User studies |
| **Analyze** | | | |
| Annotate | Formula cards (*e.g.,* MathDeck)<br>Formulas in sources | Unified text/math/++ representations (AMR? PHOC?)<br>Formula representations | User studies<br>Test collection creation<br>Identify hard topics, questions |
| Index | Flexible annotation/note organization | Query/source augmentation | User studies<br>Metrics: size, speed |
| **Synthesize** | | | |
| Apply | Reusable formulas (*e.g.,* chips) in LaTeX, code | (Pre)training tasks for model learning<br>Machine learning for embeddings, LTR<br>Model reduction (*e.g.,* for LLMs, teacher-student etc.) | User studies<br>Task metrics, comparisons<br>Cases studies<br>Experiments: within/between models |
| Communicate | Math text and formula authoring tools<br>Shared tasks | Data augmentation<br>Math-focused SERP | User studies<br>Test collection sharing<br>Shared task papers |

ing leveraging and imporving LLMs, improved attention mechanisms for transformers and finding new approaches to distilling contextual

information into vector and other representations for use in search, question answering, and other tasks are both important and widely acknowledged. Some less often discussed but still important research opportunities include improving interfaces for math-aware search and question answering, particularly that flexibly support the frequent organization, reuse, and repurposing of sources that humans undertake when working with mathematical or other information. Another key area is the development of representations for retrieval, including formulas in isolation, but particularly representations that capture text, formulas, and possibly other graphics in a manner that can assist with generating improved representations of context that can then be used as he basis for improved sparse and dense indexing.

## 6.2 Limitations: What MIR cannot provide

In our experience, a frequently expressed concern *and* hope when talking with others about MIR, is that strong systems will remove the need for individuals to understand the math that they use. Related to this, many others have expressed hope that effective math-aware search can accelerate the understanding of mathematics.

Instructors reasonably worry that math problem sets may be completed by issuing prompts to LLMs and questions to search engines, and receiving complete or near-complete answers. This is already often possible for standard math courses, for example when students make use of online resources such as conventional search engines, computational tools like Wolfram Alpha, Community Question Answering (CQA) sites such as Math Stack Exchange, and LLMs like ChatGPT.

On the other hand, students and others who feel challenged, anxious, or even overwhelmed by a mathematical problem might welcome the opportunity to avoid the mental effort and risk of failure associated with working out an answer or proof on their own. Alternatively, studying retrieved solutions to example problems before setting about solving an assignment is often helpful. A student may even reason that by mimicking observed steps, that they will deepen their understanding of the associated topic quickly.

This brings us to a fundamental limitation of technology and re-

sources for education in general, and more broadly, knowledge work. Technology can now certainly be used to store, retrieve, illustrate, and organize information (Teachers of Mathematics, 2011), saving tremendous amounts of human effort, but the benefits for human learning are modest (Cheung and Slavin, 2013). These technologies reside *outside* of a person's mind, and there is evidence that deep understanding comes from effort exerted from *inside* a person's mind. Related to this, to correctly use a search result or question answer, one must understand it first. For example, if a search engine returns a formula that a user cannot understand, what *information* has the user received?

In one psychological study, students who constructed a formula themselves for the area of geometric shapes organized in a lattice were more likely to correctly adapt the formula for new shapes than students who were given the correct formula for the first case (Hallinen *et al.*, 2021). Interestingly, the authors refer to constructing the correct formula as *searching* for the formula; the student must mentally create and compare alternative formulas as they work on the problem. Their study supports the notion that engaged exploration and comparison are needed to deepen understanding – naturally, such engaged thinking takes time.

Further, even where resources or direct answers to mathematical questions can be found, the ability to assess their validity and applicability *requires* understanding of terminology, notation, and associated concepts. As an example, Math CQA sites include a large number of posts and comments seeking clarification of, or adjustment for the context of a question, along with comments upon language and notation used in answers. And so, students trying to complete a (non-trivial) problem set well using MIR tools still need to spend time and effort to evaluate returned answers – the need to understand the pertinent mathematics remains. We do not anticipate that future MIR systems will entirely remove this requirement.

More broadly, for mathematical information differing notations and terminology between subject domains along with differences in assumed mathematical expertise stratify users into a large number of often non-overlapping sub-audiences (even within a single technical discipline). This separation is much more pronounced than for information repre-

sented by only words and targeted at a general audience (e.g., news stories). This may be partly understood as a manifestation of the *vocabulary problem* (i.e., different names for the same thing), but a user's mathematical knowledge, and what is assumed or implied in writing are also contributing factors.

And so, as we understand them, here is a summary of some limitations for mathematical information retrieval systems.

1. Relevance is strongly influenced by searcher expertise and presentation.

2. Retrieval provides sources, not the understanding of them.

3. Deep understanding of information, *e.g.,* for application in new contexts or problems, requires study, exploration, and experimentation – in other words, time.

That this needs saying at all may surprise some readers, and perhaps IR researchers more than anyone. Many other specialized search domains have similar challenges (e.g., law, medicine, chemistry). But our experience has been that mathematics, being both a tremendously powerful tool used for many purposes, and in many cultures a source of both tremendous pride and shame, has sometimes led to expressed hopes and concerns that are at odds with the limitations above.

## 6.3 Some parting thoughts

The time at which this book is being written is an exciting time for mathematical information retrieval. For example, machine learning models (LLMs, other approaches, and hybrid models) recently participated in parallel runs of the International Mathematical Olympiad, and appear to be making great strides. One has a sense that like the game of Go, which few believed computers would be able to play competitively at for many years, we may soon stumble upon the right combination of representations, statistics, and constraints to solve problems and possibly even prove theorems of real mathematical sophistication automatically.

There is another important perspective for the current state of MathIR. There is an opportunity to help people where they are given

their current mathematical expertise, and create systems that help make working with mathematical information just a bit more efficient, and a bit more comfortable. We find this prospect equally exciting, and a challenging goal from human, system, and evaluation perspectives.

# References

Ahmed, S., K. Davila, S. Setlur, and V. Govindaraju. (2021). "Equation Attention Relationship Network (EARN) : A Geometric Deep Metric Framework for Learning Similar Math Expression Embedding". en. In: *2020 25th International Conference on Pattern Recognition (ICPR)*. Milan, Italy: IEEE. 6282–6289. ISBN: 978-1-72818-808-9. DOI: 10.1109/ICPR48806.2021.9412619. URL: https://ieeexplore.ieee.org/document/9412619/ (accessed on 07/26/2021).

Aizawa, A., M. Kohlhase, and I. Ounis. (2013). "NTCIR-10 Math Pilot Task Overview". en. In: *Proc. NTCIR-10*. 8.

Aizawa, A., M. Kohlhase, and I. Ounis. (2014). "NTCIR-11 Math-2 Task Overview". en. In: *Prof. NTCIR-11*. 11.

Alama, J., T. Heskes, D. Kühlwein, E. Tsivtsivadze, and J. Urban. (2014). "Premise Selection for Mathematics by Corpus Analysis and Kernel Methods". *Journal of Automated Reasoning*. 52(2): 191–213. ISSN: 1573-0670. DOI: 10.1007/s10817-013-9286-5. URL: https://doi.org/10.1007/s10817-013-9286-5.

Alexeeva, M., R. Sharp, M. A. Valenzuela-Escárcega, J. Kadowaki, A. Pyarelal, and C. Morrison. (2020). "MathAlign: Linking Formula Identifiers to their Contextual Natural Language Descriptions". English. In: *Proceedings of the Twelfth Language Resources and Evaluation Conference*. Ed. by N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis. Marseille, France: European Language Resources Association. 2204–2212. ISBN: 979-10-95546-34-4. URL: https://aclanthology.org/2020.lrec-1.269.

Almazán, J., A. Gordo, A. Fornés, and E. Valveny. (2014). "Word Spotting and Recognition with Embedded Attributes". *IEEE Trans. Pattern Anal. Mach. Intell.* 36(12): 2552–2566. DOI: 10.1109/TPAMI.2014.2339814. URL: https://doi.org/10.1109/TPAMI.2014.2339814.

Amador, B., M. Langsenkamp, A. Dey, A. K. Shah, and R. Zanibbi. (2023). "Searching the ACL Anthology with Math Formulas and Text". In: *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*. Ed. by H. Chen, W. ( Duh, H. Huang, M. P. Kato, J. Mothe, and B. Poblete. ACM. 3110–3114. DOI: 10.1145/3539618.3591803. URL: https://doi.org/10.1145/3539618.3591803.

Amini, A., S. Gabriel, S. Lin, R. Koncel-Kedziorski, Y. Choi, and H. Hajishirzi. (2019). "MathQA: Towards Interpretable Math Word Problem Solving with Operation-Based Formalisms". In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Ed. by J. Burstein, C. Doran, and T. Solorio. Minneapolis, Minnesota: Association for Computational Linguistics. 2357–2367. DOI: 10.18653/v1/N19-1245. URL: https://aclanthology.org/N19-1245.

Asakura, T., Y. Miyao, and A. Aizawa. (2022). "Building Dataset for Grounding of Formulae — Annotating Coreference Relations Among Math Identifiers". In: *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Ed. by N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, J. Odijk, and S. Piperidis. Marseille, France: European Language Resources Association. 4851–4858. URL: https://aclanthology.org/2022.lrec-1.519.

Asakura, T., Y. Miyao, A. Aizawa, and M. Kohlhase. (2021). "MioGatto: A Math Identifier-oriented Grounding Annotation Tool". In:

Avenoso, R. (2021). "Spatial vs. Graph-Based Formula Retrieval". en: 108.

Avenoso, R., B. Mansouri, and R. Zanibbi. (2021). "XY-PHOC Symbol Location Embeddings for Math Formula Retrieval and Autocompletion". In: *Proceedings of the Working Notes of CLEF 2021 - Conference and Labs of the Evaluation Forum, Bucharest, Romania, September 21st - to - 24th, 2021*. Ed. by G. Faggioli, N. Ferro, A. Joly, M. Maistro, and F. Piroi. Vol. 2936. *CEUR Workshop Proceedings*. CEUR-WS.org. 25–35. URL: https://ceur-ws.org/Vol-2936/paper-02.pdf.

Azerbayev, Z., H. Schoelkopf, K. Paster, M. D. Santos, S. McAleer, A. Q. Jiang, J. Deng, S. Biderman, and S. Welleck. (2024). "Llemma: An Open Language Model For Mathematics". arXiv: 2310.10631 [cs.CL]. URL: https://arxiv.org/abs/2310.10631.

Banarescu, L., C. Bonial, S. Cai, M. Georgescu, K. Griffitt, U. Hermjakob, K. Knight, P. Koehn, M. Palmer, and N. Schneider. (2013). "Abstract meaning representation for sembanking". In: *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*. 178–186.

Bates, M. J. (1989). "The design of browsing and berrypicking techniques for the online search interface". *Online review*.

Belkin, N. J. (1980). "Anomalous states of knowledge as a basis for information retrieval". *Canadian journal of information science*. 5(1): 133–143.

Bender, E. M., T. Gebru, A. McMillan-Major, and S. Shmitchell. (2021). "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" In: *FAccT '21: 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event / Toronto, Canada, March 3-10, 2021*. Ed. by M. C. Elish, W. Isaac, and R. S. Zemel. ACM. 610–623. DOI: 10.1145/3442188.3445922. URL: https://doi.org/10.1145/3442188.3445922.

Blei, D. M., A. Y. Ng, and M. I. Jordan. (2003). "Latent Dirichlet Allocation". *J. Mach. Learn. Res.* 3: 993–1022. URL: http://jmlr.org/papers/v3/blei03a.html.

Bobrow, D. G. (1964). "A question-answering system for high school algebra word problems". In: *Proceedings of the October 27-29, 1964, Fall Joint Computer Conference, Part I. AFIPS '64 (Fall, part I)*. San Francisco, California: Association for Computing Machinery. 591–614. ISBN: 9781450378895. DOI: 10.1145/1464052.1464108. URL: https://doi.org/10.1145/1464052.1464108.

Bojanowski, P., E. Grave, A. Joulin, and T. Mikolov. (2017). "Enriching Word Vectors with Subword Information". en. *Transactions of the Association for Computational Linguistics*. 5(Dec.): 135–146. ISSN: 2307-387X. DOI: 10.1162/tacl_a_00051. URL: https://direct.mit.edu/tacl/article/43387 (accessed on 07/27/2021).

Borlund, P. (2003). "The IIR evaluation model: a framework for evaluation of interactive information retrieval systems". *Inf. Res.* 8(3). URL: http://www.informationr.net/ir/8-3/paper152.html.

Bozzon, A. and P. Fraternali. (2010). "Chapter 8: Multimedia and Multimodal Information Retrieval". In: *Search Computing: Challenges and Directions*. Ed. by S. Ceri and M. Brambilla. Berlin, Heidelberg: Springer Berlin Heidelberg. 135–155. ISBN: 978-3-642-12310-8. DOI: 10.1007/978-3-642-12310-8_8. URL: https://doi.org/10.1007/978-3-642-12310-8_8.

Broder, A. Z. (2002). "A taxonomy of web search". *SIGIR Forum.* 36(2): 3–10. DOI: 10.1145/792550.792552. URL: https://doi.org/10.1145/792550.792552.

Buckley, C. and E. M. Voorhees. (2004). "Retrieval evaluation with incomplete information". en. In: *Proceedings of the 27th annual international conference on Research and development in information retrieval - SIGIR '04*. Sheffield, United Kingdom: ACM Press. 25. ISBN: 978-1-58113-881-8. DOI: 10.1145/1008992.1009000. URL: http://portal.acm.org/citation.cfm?doid=1008992.1009000 (accessed on 07/26/2021).

Chen, J., T. Li, J. Qin, P. Lu, L. Lin, C. Chen, and X. Liang. (2022a). "UniGeo: Unifying Geometry Logical Reasoning via Reformulating Mathematical Expression". In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Ed. by Y. Goldberg, Z. Kozareva, and Y. Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics. 3313–3323. DOI: 10.18653/v1/2022.emnlp-main.218. URL: https://aclanthology.org/2022.emnlp-main.218.

Chen, J., J. Tang, J. Qin, X. Liang, L. Liu, E. Xing, and L. Lin. (2021a). "GeoQA: A Geometric Question Answering Benchmark Towards Multimodal Numerical Reasoning". In: *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Ed. by C. Zong, F. Xia, W. Li, and R. Navigli. Online: Association for Computational Linguistics. 513–523. DOI: 10.18653/v1/2021.findings-acl.46. URL: https://aclanthology.org/2021.findings-acl.46.

Chen, M., J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, *et al.* (2021b). "Evaluating large language models trained on code". *arXiv preprint arXiv:2107.03374*.

Chen, W., X. Ma, X. Wang, and W. W. Cohen. (2022b). "Program of Thoughts Prompting: Disentangling Computation from Reasoning for Numerical Reasoning Tasks". *arXiv preprint arXiv:2211.12588*.

Cheung, A. C. and R. E. Slavin. (2013). "The effectiveness of educational technology applications for enhancing mathematics achievement in K-12 classrooms: A meta-analysis". en. *Educational Research Review*. 9(June): 88–113. ISSN: 1747938X. DOI: 10.1016/j.edurev.2013.01.001. URL: https://linkinghub.elsevier.com/retrieve/pii/S1747938X13000031 (accessed on 10/27/2021).

Church, A. (1940). "A formulation of the simple theory of types". *The journal of symbolic logic*. 5(2): 56–68.

Cobbe, K., V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman. (2021). "Training Verifiers to Solve Math Word Problems". arXiv: 2110.14168 [cs.LG]. URL: https://arxiv.org/abs/2110.14168.

Croft, W. B., D. Metzler, and T. Strohman. (2009). *Search Engines - Information Retrieval in Practice*. Pearson Education. ISBN: 978-0-13-136489-9. URL: http://www.search-engines-book.com/.

Dai, Y., L. Chen, and Z. Zhang. (2020). "An N-ary Tree-based Model for Similarity Evaluation on Mathematical Formulae". en. In: *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. Toronto, ON, Canada: IEEE. 2578–2584. ISBN: 978-1-72818-526-2. DOI: 10.1109/SMC42975.2020.9283495. URL: https://ieeexplore.ieee.org/document/9283495/ (accessed on 07/26/2021).

Davila, K., R. Joshi, S. Setlur, V. Govindaraju, and R. Zanibbi. (2019). "Tangent-V: Math Formula Image Search Using Line-of-Sight Graphs". en. In: *Advances in Information Retrieval*. Ed. by L. Azzopardi, B. Stein, N. Fuhr, P. Mayr, C. Hauff, and D. Hiemstra. Vol. 11437. Cham: Springer International Publishing. 681–695. ISBN: 978-3-030-15711-1 978-3-030-15712-8. DOI: 10.1007/978-3-030-15712-8_44. URL: http://link.springer.com/10.1007/978-3-030-15712-8_44 (accessed on 07/26/2021).

Davila, K., S. Ludi, and R. Zanibbi. (2014). "Using Off-line Features and Synthetic Data for On-line Handwritten Math Symbol Recognition". In: *2014 14th International Conference on Frontiers in Handwriting Recognition*. IEEE. 323–328.

Davila, K., F. Xu, S. Setlur, and V. Govindaraju. (2021). "FCN-LectureNet: Extractive Summarization of Whiteboard and Chalkboard Lecture Videos". *IEEE Access*. 9: 104469–104484. DOI: 10.1109/ACCESS.2021.3099427. URL: https://doi.org/10.1109/ACCESS.2021.3099427.

Davila, K. and R. Zanibbi. (2017a). "Layout and Semantics: Combining Representations for Mathematical Formula Search". en. In: *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Shinjuku Tokyo Japan: ACM. 1165–1168. ISBN: 978-1-4503-5022-8. DOI: 10.1145/3077136.3080748. URL: https://dl.acm.org/doi/10.1145/3077136.3080748 (accessed on 07/26/2021).

Davila, K. and R. Zanibbi. (2017b). "Whiteboard Video Summarization via Spatio-Temporal Conflict Minimization". In: *14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017*. IEEE. 355–362. DOI: 10.1109/ICDAR.2017.66. URL: https://doi.org/10.1109/ICDAR.2017.66.

Davila, K. and R. Zanibbi. (2018). "Visual Search Engine for Handwritten and Typeset Math in Lecture Videos and LATEX Notes". In: *16th International Conference on Frontiers in Handwriting Recognition, ICFHR 2018, Niagara Falls, NY, USA, August 5-8, 2018*. IEEE Computer Society. 50–55. DOI: 10.1109/ICFHR-2018.2018.00018. URL: https://doi.org/10.1109/ICFHR-2018.2018.00018.

Davis, M., G. Logemann, and D. Loveland. (1962). "A machine program for theorem-proving". *Commun. ACM*. 5(7): 394–397. ISSN: 0001-0782. DOI: 10.1145/368273.368557. URL: https://doi.org/10.1145/368273.368557.

Davis, M. and H. Putnam. (1960). "A Computing Procedure for Quantification Theory". *J. ACM*. 7(3): 201–215. ISSN: 0004-5411. DOI: 10.1145/321033.321034. URL: https://doi.org/10.1145/321033.321034.

Devlin, J., M.-W. Chang, K. Lee, and K. Toutanova. (2019). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Ed. by J. Burstein, C. Doran, and T. Solorio. Minneapolis, Minnesota: Association for Computational Linguistics. 4171–4186. DOI: 10.18653/v1/N19-1423. URL: https://aclanthology.org/N19-1423.

Diaz, Y. (2021). "Introducing Handwriting into a Multimodal LaTeX Formula Editor". *MA thesis*. Rochester Institute of Technology.

Diaz, Y., G. Nishizawa, B. Mansouri, K. Davila, and R. Zanibbi. (2021). "The Math-Deck Formula Editor: Interactive Formula Entry Combining LaTeX , Structure Editing, and Search". en. In: *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*. Yokohama Japan: ACM. 1–5. ISBN: 978-1-4503-8095-9. DOI: 10.1145/3411763.3451564. URL: https://dl.acm.org/doi/10.1145/3411763.3451564 (accessed on 07/26/2021).

Dmello, A. (2019a). "Representing Mathematical Concepts Associated With Formulas Using Math Entity Cards". en. *MA thesis*. Rochester Institute of Technology. 167.

Dmello, A. (2019b). *Representing mathematical concepts associated with formulas using math entity cards*. Rochester Institute of Technology.

Dua, D., Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, and M. Gardner. (2019). "DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs". arXiv: 1903.00161 [cs.CL]. URL: https://arxiv.org/abs/1903.00161.

Duncan, D. (2021). *Index, a History of the*. New York: Norton.

Durgin, S., J. Gore, and B. Mansouri. (2024). "MathMex: Search Engine for Math Definitions". In: *Advances in Information Retrieval - 46th European Conference on Information Retrieval, ECIR 2024, Glasgow, UK, March 24-28, 2024, Proceedings, Part V*. Ed. by N. Goharian, N. Tonellotto, Y. He, A. Lipani, G. McDonald, C. Macdonald, and I. Ounis. Vol. 14612. *Lecture Notes in Computer Science*. Springer. 194–199. DOI: 10.1007/978-3-031-56069-9\_17. URL: https://doi.org/10.1007/978-3-031-56069-9%5C_17.

Ethayarajh, K., W. Xu, N. Muennighoff, D. Jurafsky, and D. Kiela. (2024). "KTO: Model Alignment as Prospect Theoretic Optimization". arXiv: 2402.01306 [cs.LG]. URL: https://arxiv.org/abs/2402.01306.

Faggioli, G., L. Dietz, C. L. A. Clarke, G. Demartini, M. Hagen, C. Hauff, N. Kando, E. Kanoulas, M. Potthast, B. Stein, and H. Wachsmuth. (2024). "Who Determines What Is Relevant? Humans or AI? Why Not Both?" *Commun. ACM*. 67(4): 31–34. DOI: 10.1145/3624730. URL: https://doi.org/10.1145/3624730.

Ferreira, D. and A. Freitas. (2020a). "Premise Selection in Natural Language Mathematical Texts". en. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics. 7365–7374. DOI: 10.18653/v1/2020.acl-main.657. URL: https://www.aclweb.org/anthology/2020.acl-main.657 (accessed on 07/27/2021).

Ferreira, D. and A. Freitas. (2021). "STAR: Cross-modal [STA]tement [R]epresentation for selecting relevant mathematical premises". In: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Ed. by P. Merlo, J. Tiedemann, and R. Tsarfaty. Online: Association for Computational Linguistics. 3234–3243. DOI: 10.18653/v1/2021.eacl-main.282. URL: https://aclanthology.org/2021.eacl-main.282.

Ferreira, D. and A. Freitas. (2020b). "Natural Language Premise Selection: Finding Supporting Statements for Mathematical Text". arXiv: 2004.14959 [cs.CL]. URL: https://arxiv.org/abs/2004.14959.

Formal, T., B. Piwowarski, and S. Clinchant. (2021a). "SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking". In: *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*. Ed. by F. Diaz, C. Shah, T. Suel, P. Castells, R. Jones, and T. Sakai. ACM. 2288–2292. DOI: 10.1145/3404835.3463098. URL: https://doi.org/10.1145/3404835.3463098.

Formal, T., B. Piwowarski, and S. Clinchant. (2021b). "SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking". In: *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '21*. Virtual Event, Canada: Association for Computing Machinery. 2288–2292. ISBN: 9781450380379. DOI: 10.1145/3404835.3463098. URL: https://doi.org/10.1145/3404835.3463098.

Fraser, D., A. Kane, and F. W. Tompa. (2018). "Choosing Math Features for BM25 Ranking with Tangent-L". en. In: *Proceedings of the ACM Symposium on Document Engineering 2018*. Halifax NS Canada: ACM. 1–10. ISBN: 978-1-4503-5769-2. DOI: 10.1145/3209280.3209527. URL: https://dl.acm.org/doi/10.1145/3209280.3209527 (accessed on 07/26/2021).

Freund, Y., R. Iyer, R. E. Schapire, and Y. Singer. (2003). "An Efficient Boosting Algorithm for Combining Preferences". *Journal of machine learning research*. 4(Nov): 933–969.

Fuhr, N. (2017). "Some Common Mistakes In IR Evaluation, And How They Can Be Avoided". *SIGIR Forum*. 51(3): 32–41. DOI: 10.1145/3190580.3190586. URL: https://doi.org/10.1145/3190580.3190586.

Gao, L., Z. Jiang, Y. Yin, K. Yuan, Z. Yan, and Z. Tang. (2017). "Preliminary Exploration of Formula Embedding for Mathematical Information Retrieval: can mathematical formulae be embedded like a natural language?" en. *arXiv:1707.05154 [cs]*. Aug. URL: http://arxiv.org/abs/1707.05154 (accessed on 07/26/2021).

Gao, L., K. Yuan, Y. Wang, Z. Jiang, and Z. Tang. (2016). "The Math Retrieval System of ICST for NTCIR-12 MathIR Task". In: *NTCIR*.

Gao, L., A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig. (2023). "PAL: Program-aided Language Models". In: *Proceedings of the 40th International Conference on Machine Learning*. Ed. by A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett. Vol. 202. *Proceedings of Machine Learning Research*. PMLR. 10764–10799. URL: https://proceedings.mlr.press/v202/gao23f.html.

Gao, S. and Y.-K. D. Ng. (2023). "Recommending Answers to Math Questions Based on KL-Divergence and Approximate XML Tree Matching". In: *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region. SIGIR-AP '23*. Beijing, China: Association for Computing Machinery. 21–31. ISBN: 9798400704086. DOI: 10.1145/3624918.3625337. URL: https://doi.org/10.1145/3624918.3625337.

Giacalone, B., G. Paiement, Q. Tucker, and R. Zanibbi. (2024). "Beneath the [MASK]: An Analysis of Structural Query Tokens in ColBERT". In: *Advances in Information Retrieval - 46th European Conference on Information Retrieval, ECIR 2024, Glasgow, UK, March 24-28, 2024, Proceedings, Part III*. Ed. by N. Goharian, N. Tonellotto, Y. He, A. Lipani, G. McDonald, C. Macdonald, and I. Ounis. Vol. 14610. *Lecture Notes in Computer Science*. Springer. 431–439. DOI: 10.1007/978-3-031-56063-7\_35. URL: https://doi.org/10.1007/978-3-031-56063-7%5C_35.

Gore, J., J. Polletta, and B. Mansouri. (2024). "CrossMath: Towards Cross-lingual Math Information Retrieval". In: *The 10th ACM SIGIR / The 14th International Conference on the Theory of Information Retrieval*. URL: https://openreview.net/forum?id=IYcJxgnJ9Q.

Gou, Z., Z. Shao, Y. Gong, Y. Shen, Y. Yang, M. Huang, N. Duan, and W. Chen. (2024). "ToRA: A Tool-Integrated Reasoning Agent for Mathematical Problem Solving". arXiv: 2309.17452 [cs.CL]. URL: https://arxiv.org/abs/2309.17452.

Graf, P. (1995a). "Substitution Tree Indexing". In: *International Conference on Rewriting Techniques and Applications*. Springer. 117–131.

Graf, P. (1995b). "Substitution tree indexing". In: *Rewriting Techniques and Applications*. Ed. by J. Hsiang. Berlin, Heidelberg: Springer Berlin Heidelberg. 117–131. ISBN: 978-3-540-49223-8.

Grossman, J., Z. Lin, H. Sheng, J. T.-Z. Wei, J. J. Williams, and S. Goel. (2019). "MathBot: Transforming online resources for learning math into conversational interactions". *AAAI 2019 Story-Enabled Intelligence.*

Grover, A. and J. Leskovec. (2016). "node2vec: Scalable Feature Learning for Networks". In: *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining.* 855–864.

Ha, J., R. Haralick, and I. Phillips. (1995). "Recursive X-Y cut using bounding boxes of connected components". In: *Proceedings of 3rd International Conference on Document Analysis and Recognition.* Vol. 2. 952–955 vol.2. DOI: 10.1109/ICDAR. 1995.602059.

Hallinen, N. R., L. N. Sprague, K. P. Blair, R. M. Adler, and N. S. Newcombe. (2021). "Finding formulas: Does active search facilitate appropriate generalization?" *Cognitive Research: Principles and Implications.* 6(1): 50. ISSN: 2365-7464. DOI: 10.1186/s41235-021-00316-y. URL: https://doi.org/10.1186/s41235-021-00316-y (accessed on 10/27/2021).

Hambasan, R., M. Kohlhase, and C. Prodescu. (2014). "MathWebSearch at NTCIR-11". en: 6.

Hamilton, W., Z. Ying, and J. Leskovec. (2017). "Inductive Representation Learning on Large Graphs". In: *Advances in Neural Information Processing Systems.* Ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Vol. 30. Curran Associates, Inc. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/5dd9db5e033da9c6fb5ba83c7a7ebea9-Paper.pdf.

He, K., X. Zhang, S. Ren, and J. Sun. (2016). "Deep Residual Learning for Image Recognition". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

Head, A., K. Lo, D. Kang, R. Fok, S. Skjonsberg, D. S. Weld, and M. A. Hearst. (2021). "Augmenting Scientific Papers with Just-in-Time, Position-Sensitive Definitions of Terms and Symbols". en. In: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems.* Yokohama Japan: ACM. 1–18. ISBN: 978-1-4503-8096-6. DOI: 10.1145/3411764.3445648. URL: https://dl.acm.org/doi/10.1145/3411764.3445648 (accessed on 07/26/2021).

Hearst, M. (2009). *Search user interfaces.* Cambridge university press.

Hendrycks, D., C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. (2021). "Measuring Mathematical Problem Solving With the MATH Dataset". arXiv: 2103.03874 [cs.LG]. URL: https://arxiv.org/abs/2103.03874.

Holmes, S., A. Moorhead, R. Bond, H. Zheng, V. Coates, and M. McTear. (2019). "Usability testing of a healthcare chatbot: Can we use conventional methods to assess conversational user interfaces?" In: *Proc. 31st European Conference on Cognitive Ergonomics.* 207–214.

Hopkins, M., R. Le Bras, C. Petrescu-Prahova, G. Stanovsky, H. Hajishirzi, and R. Koncel-Kedziorski. (2019). "SemEval-2019 Task 10: Math Question Answering". In: *Proceedings of the 13th International Workshop on Semantic Evaluation.* Ed. by J. May, E. Shutova, A. Herbelot, X. Zhu, M. Apidianaki, and S. M. Mohammad. Minneapolis, Minnesota, USA: Association for Computational Linguistics. 893–899. DOI: 10.18653/v1/S19-2153. URL: https://aclanthology.org/S19-2153.

Hosseini, M. J., H. Hajishirzi, O. Etzioni, and N. Kushman. (2014). "Learning to Solve Arithmetic Word Problems with Verb Categorization". In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Ed. by A. Moschitti, B. Pang, and W. Daelemans. Doha, Qatar: Association for Computational Linguistics. 523–533. DOI: 10.3115/v1/D14-1058. URL: https://aclanthology.org/D14-1058.

Hu, X., L. Gao, X. Lin, Z. Tang, X. Lin, and J. B. Baker. (2013). "WikiMirs: a mathematical information retrieval system for wikipedia". en. In: *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries - JCDL '13*. Indianapolis, Indiana, USA: ACM Press. 11. ISBN: 978-1-4503-2077-1. DOI: 10.1145/2467696.2467699. URL: http://dl.acm.org/citation.cfm?doid=2467696.2467699 (accessed on 07/26/2021).

Huang, D., S. Shi, C.-Y. Lin, J. Yin, and W.-Y. Ma. (2016). "How well do Computers Solve Math Word Problems? Large-Scale Dataset Construction and Evaluation". In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by K. Erk and N. A. Smith. Berlin, Germany: Association for Computational Linguistics. 887–896. DOI: 10.18653/v1/P16-1084. URL: https://aclanthology.org/P16-1084.

Irving, G., C. Szegedy, A. A. Alemi, N. Een, F. Chollet, and J. Urban. (2016). "Deep-Math - Deep Sequence Models for Premise Selection". In: *Advances in Neural Information Processing Systems*. Ed. by D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett. Vol. 29. Curran Associates, Inc. URL: https://proceedings.neurips.cc/paper_files/paper/2016/file/f197002b9a0853eca5e046d9ca4663d5-Paper.pdf.

Ito, T., T. Matsuzaki, and S. Sato. (2017). "Coreference Resolution on Math Problem Text in Japanese". In: *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Ed. by G. Kondrak and T. Watanabe. Taipei, Taiwan: Asian Federation of Natural Language Processing. 373–377. URL: https://aclanthology.org/I17-2063.

Jiang, A. Q., A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, *et al.* (2023). "Mistral 7B". *arXiv preprint arXiv:2310.06825*.

Jo, H., D. Kang, A. Head, and M. A. Hearst. (2021). "Modeling Mathematical Notation Semantics in Academic Papers". In: *Findings of the Association for Computational Linguistics: EMNLP 2021*. Ed. by M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih. Punta Cana, Dominican Republic: Association for Computational Linguistics. 3102–3115. DOI: 10.18653/v1/2021.findings-emnlp.266. URL: https://aclanthology.org/2021.findings-emnlp.266.

Joachims, T. (2006). "Training Linear SVMs in Linear Time". In: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. 217–226.

Jones, K. S. (1972). "A statistical interpretation of term specificity and its application in retrieval". *J. Documentation*. 28(1): 11–21. URL: https://api.semanticscholar.org/CorpusID:2996187.

Kamali, S. and F. W. Tompa. (2013). "Retrieving documents with mathematical content". en. In: *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*. Dublin Ireland: ACM. 353–362. ISBN: 978-1-4503-2034-4. DOI: 10.1145/2484028.2484083. URL: https://dl.acm.org/doi/10.1145/2484028.2484083 (accessed on 07/26/2021).

Kane, A., Y. K. Ng, and F. W. Tompa. (2022). "Dowsing for Answers to Math Questions: Doing Better with Less". In: *CLEF (Working Notes)*. 40–62.

Karpukhin, V., B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih. (2020). "Dense Passage Retrieval for Open-Domain Question Answering". In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Ed. by B. Webber, T. Cohn, Y. He, and Y. Liu. Online: Association for Computational Linguistics. 6769–6781. DOI: 10.18653/v1/2020.emnlp-main.550. URL: https://aclanthology.org/2020.emnlp-main.550.

Khattab, O. and M. Zaharia. (2020a). "ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT". In: *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*. Ed. by J. X. Huang, Y. Chang, X. Cheng, J. Kamps, V. Murdock, J. Wen, and Y. Liu. ACM. 39–48. DOI: 10.1145/3397271.3401075. URL: https://doi.org/10.1145/3397271.3401075.

Khattab, O. and M. Zaharia. (2020b). "ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT". In: *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '20*. Virtual Event, China: Association for Computing Machinery. 39–48. ISBN: 9781450380164. DOI: 10.1145/3397271.3401075. URL: https://doi.org/10.1145/3397271.3401075.

Kirsch, D. (2010). "Detexify: Recognition of hand-drawn LATEX symbols". *MA thesis*. University of Münster (in German).

Kohlhase, A. and M. Kohlhase. (2007). "*Re* examining the MKM Value Proposition: From Math Web Search to Math Web *Re* Search". In: *Towards Mechanized Mathematical Assistants, 14th Symposium, Calculemus 2007, 6th International Conference, MKM 2007, Hagenberg, Austria, June 27-30, 2007, Proceedings*. Ed. by M. Kauers, M. Kerber, R. Miner, and W. Windsteiger. Vol. 4573. *Lecture Notes in Computer Science*. Springer. 313–326. DOI: 10.1007/978-3-540-73086-6\_25. URL: https://doi.org/10.1007/978-3-540-73086-6%5C_25.

Kohlhase, M. and C. Prodescu. (2013). "MathWebSearch at NTCIR-10". en: 5.

Kohlhase, M. and I. Sucan. (2006). "A Search Engine for Mathematical Formulae". en. In: *Artificial Intelligence and Symbolic Computation*. Ed. by D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, J. Calmet, T. Ida, and D. Wang. Vol. 4120. Berlin, Heidelberg: Springer Berlin Heidelberg. 241–253. ISBN: 978-3-540-39728-1 978-3-540-39730-4. DOI: 10.1007/11856290_21. URL: http://link.springer.com/10.1007/11856290_21 (accessed on 07/26/2021).

Koncel-Kedziorski, R., H. Hajishirzi, A. Sabharwal, O. Etzioni, and S. D. Ang. (2015). "Parsing algebraic word problems into equations". *Transactions of the Association for Computational Linguistics*. 3: 585–597.

Kristianto, G. Y., A. Aizawa, *et al.* (2014). "Extracting Textual Descriptions of Mathematical Expressions in Scientific Papers". *D-Lib Magazine*. 20(11): 9.

Kristianto, G. Y., G. Topic, and A. Aizawa. (2016a). "MCAT Math Retrieval System for NTCIR-12 MathIR Task". en: 8.

Kristianto, G. Y., G. Topic, and A. Aizawa. (2017). "Utilizing dependency relationships between math expressions in math IR". *Inf. Retr. J.* 20(2): 132–167. DOI: 10.1007/S10791-017-9296-8. URL: https://doi.org/10.1007/s10791-017-9296-8.

Kristianto, G. Y., G. Topić, and A. Aizawa. (2016b). "Entity Linking for Mathematical Expressions in Scientific Documents". In: *Digital Libraries: Knowledge, Information, and Data in an Open Access Society.* Ed. by A. Morishima, A. Rauber, and C. L. Liew. Cham: Springer International Publishing. 144–149. ISBN: 978-3-319-49304-6.

Kristianto, G. Y., G. Topić, and A. Aizawa. (2016c). "Entity linking for mathematical expressions in scientific documents". In: *Digital Libraries: Knowledge, Information, and Data in an Open Access Society: 18th International Conference on Asia-Pacific Digital Libraries, ICADL 2016, Tsukuba, Japan, December 7–9, 2016, Proceedings 18.* Springer. 144–149.

Krstovski, K. and D. M. Blei. (2018). "Equation Embeddings". en. *arXiv:1803.09123 [cs, stat].* Mar. URL: http://arxiv.org/abs/1803.09123 (accessed on 07/26/2021).

Lai, V., A. Pouran Ben Veyseh, F. Dernoncourt, and T. Nguyen. (2022). "SemEval 2022 Task 12: Symlink - Linking Mathematical Symbols to their Descriptions". In: *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022).* Ed. by G. Emerson, N. Schluter, G. Stanovsky, R. Kumar, A. Palmer, N. Schneider, S. Singh, and S. Ratan. Seattle, United States: Association for Computational Linguistics. 1671–1678. DOI: 10.18653/v1/2022.semeval-1.230. URL: https://aclanthology.org/2022.semeval-1.230.

Langkilde, I. and K. Knight. (1998). "Generation that Exploits Corpus-Based Statistical Knowledge". In: *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, COLING-ACL '98, August 10-14, 1998, Université de Montréal, Montréal, Quebec, Canada. Proceedings of the Conference.* Ed. by C. Boitet and P. Whitelock. Morgan Kaufmann Publishers / ACL. 704–710. DOI: 10.3115/980845.980963. URL: https://aclanthology.org/P98-1116/.

Langsenkamp, M., B. Amador, and R. Zanibbi. (2024). "A Study of PHOC Spatial Region Configurations for Math Formula Retrieval". *Tech. rep.* No. TBD. arXiv.

Langsenkamp, M., B. Mansouri, and R. Zanibbi. (2022). "Expanding Spatial Regions and Incorporating IDF for PHOC-Based Math Formula Retrieval at ARQMath-3". In: *Proceedings of the Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum, Bologna, Italy, September 5th - to - 8th, 2022.* Ed. by G. Faggioli, N. Ferro, A. Hanbury, and M. Potthast. Vol. 3180. *CEUR Workshop Proceedings.* CEUR-WS.org. 63–82. URL: https://ceur-ws.org/Vol-3180/paper-04.pdf.

Le, Q. and T. Mikolov. (2014a). "Distributed Representations of Sentences and Documents". In: *International conference on machine learning*. PMLR. 1188–1196.

Le, Q. and T. Mikolov. (2014b). "Distributed Representations of Sentences and Documents". In: *Proceedings of the 31st International Conference on Machine Learning*. Ed. by E. P. Xing and T. Jebara. Vol. 32. *Proceedings of Machine Learning Research*. No. 2. Bejing, China: PMLR. 1188–1196. URL: https://proceedings.mlr.press/v32/le14.html.

Lewkowycz, A., A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, Y. Wu, B. Neyshabur, G. Gur-Ari, and V. Misra. (2022). "Solving Quantitative Reasoning Problems with Language Models". In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. Vol. 35. Curran Associates, Inc. 3843–3857. URL: https://proceedings.neurips.cc/paper_files/paper/2022/file/18abbeef8cfe9203fdf9053c9c4fe191-Paper-Conference.pdf.

Liang, Z., J. Zhang, L. Wang, W. Qin, Y. Lan, J. Shao, and X. Zhang. (2022). "MWP-BERT: Numeracy-Augmented Pre-training for Math Word Problem Solving". In: *Findings of the Association for Computational Linguistics: NAACL 2022*. Ed. by M. Carpuat, M.-C. de Marneffe, and I. V. Meza Ruiz. Seattle, United States: Association for Computational Linguistics. 997–1009. DOI: 10.18653/v1/2022.findings-naacl.74. URL: https://aclanthology.org/2022.findings-naacl.74.

Libbrecht, P. and E. Melis. (2006). "Methods to Access and Retrieve Mathematical Content in ActiveMath". en. In: *Mathematical Software - ICMS 2006*. Ed. by D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, A. Iglesias, and N. Takayama. Vol. 4151. Berlin, Heidelberg: Springer Berlin Heidelberg. 331–342. ISBN: 978-3-540-38084-9 978-3-540-38086-3. DOI: 10.1007/11832225_33. URL: http://link.springer.com/10.1007/11832225_33 (accessed on 07/26/2021).

Ling, W., D. Yogatama, C. Dyer, and P. Blunsom. (2017). "Program Induction by Rationale Generation: Learning to Solve and Explain Algebraic Word Problems". In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by R. Barzilay and M.-Y. Kan. Vancouver, Canada: Association for Computational Linguistics. 158–167. DOI: 10.18653/v1/P17-1015. URL: https://aclanthology.org/P17-1015.

Líška, M., P. Sojka, M. Ruzicka, and P. Mravec. (2011). "Web Interface and Collection for Mathematical Retrieval : WebMIaS and MREC". In: *Towards a Digital Mathematics Library*. Ed. by P. Sojka and T. Bouche. Bertinoro, Bertinoro, Italy: Masaryk University Press (Brno, Czech Republic). 77–84.

Liu, Y., K. Ding, and Y. Zhou. (2019). "AiFu at SemEval-2019 Task 10: A Symbolic and Sub-symbolic Integrated System for SAT Math Question Answering". In: *Proceedings of the 13th International Workshop on Semantic Evaluation*. Ed. by J. May, E. Shutova, A. Herbelot, X. Zhu, M. Apidianaki, and S. M. Mohammad. Minneapolis, Minnesota, USA: Association for Computational Linguistics. 900–906. DOI: 10.18653/v1/S19-2154. URL: https://aclanthology.org/S19-2154.

Lu, P., H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao. (2024). "MathVista: Evaluating Mathematical Reasoning of Foundation Models in Visual Contexts". arXiv: 2310.02255 [cs.CV]. URL: https://arxiv.org/abs/2310.02255.

Lu, P., L. Qiu, K.-W. Chang, Y. N. Wu, S.-C. Zhu, T. Rajpurohit, P. Clark, and A. Kalyan. (2023). "Dynamic Prompt Learning via Policy Gradient for Semi-structured Mathematical Reasoning". arXiv: 2209.14610 [cs.LG]. URL: https://arxiv.org/abs/2209.14610.

Lv, Y. and C. Zhai. (2011). "Lower-Bounding Term Frequency Normalization". In: *Proceedings of the 20th ACM international conference on Information and knowledge management.* 7–16.

Macina, J., N. Daheim, S. Chowdhury, T. Sinha, M. Kapur, I. Gurevych, and M. Sachan. (2023). "MathDial: A Dialogue Tutoring Dataset with Rich Pedagogical Properties Grounded in Math Reasoning Problems". In: *Findings of the Association for Computational Linguistics: EMNLP 2023.* Ed. by H. Bouamor, J. Pino, and K. Bali. Singapore: Association for Computational Linguistics. 5602–5621. DOI: 10.18653/v1/2023.findings-emnlp.372. URL: https://aclanthology.org/2023.findings-emnlp.372.

Mahdavi, M., R. Zanibbi, H. Mouchere, C. Viard-Gaudin, and U. Garain. (2019). "ICDAR 2019 CROHME + TFD: Competition on Recognition of Handwritten Mathematical Expressions and Typeset Formula Detection". In: *2019 International Conference on Document Analysis and Recognition (ICDAR).* 1533–1538. DOI: 10.1109/ICDAR.2019.00247.

Mansouri, B. and Z. Jahedibashiz. (2023). "Clarifying Questions in Math Information Retrieval". In: *Proceedings of the 2023 ACM SIGIR International Conference on Theory of Information Retrieval. ICTIR '23.* Taipei, Taiwan: Association for Computing Machinery. 149–158. ISBN: 9798400700736. DOI: 10.1145/3578337.3605123. URL: https://doi.org/10.1145/3578337.3605123.

Mansouri, B. and R. Maarefdoust. (2024). "Using Large Language Models for Math Information Retrieval". In: *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '24.* Washington DC, USA: Association for Computing Machinery. 2693–2697. ISBN: 9798400704314. DOI: 10.1145/3626772.3657907. URL: https://doi.org/10.1145/3626772.3657907.

Mansouri, B., V. Novotný, A. Agarwal, D. W. Oard, and R. Zanibbi. (2022a). "Overview of ARQMath-3 (2022): Third CLEF Lab on Answer Retrieval for Questions on Math". In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction.* Ed. by A. Barrón-Cedeño, G. Da San Martino, M. Degli Esposti, F. Sebastiani, C. Macdonald, G. Pasi, A. Hanbury, M. Potthast, G. Faggioli, and N. Ferro. Cham: Springer International Publishing. 286–310. ISBN: 978-3-031-13643-6.

Mansouri, B., V. Novotný, A. Agarwal, D. W. Oard, and R. Zanibbi. (2022b). "Overview of ARQMath-3 (2022): Third CLEF Lab on Answer Retrieval for Questions on Math (Working Notes Version)". In: *Proceedings of the Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum, Bologna, Italy, September 5th - to - 8th, 2022*. Ed. by G. Faggioli, N. Ferro, A. Hanbury, and M. Potthast. Vol. 3180. *CEUR Workshop Proceedings*. CEUR-WS.org. 1–27. URL: https://ceur-ws.org/Vol-3180/paper-01.pdf.

Mansouri, B., D. W. Oard, and R. Zanibbi. (2020). "DPRL Systems in the CLEF 2020 ARQMath Lab". en: 12.

Mansouri, B., D. W. Oard, A. Agarwal, and R. Zanibbi. (2021a). "Effects of context, complexity, and clustering on evaluation for math formula retrieval". *CoRR*. abs/2111.10504. arXiv: 2111.10504. URL: https://arxiv.org/abs/2111.10504.

Mansouri, B., D. W. Oard, and R. Zanibbi. (2022c). "Contextualized Formula Search Using Math Abstract Meaning Representation". In: *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*. Ed. by M. A. Hasan and L. Xiong. ACM. 4329–4333. DOI: 10.1145/3511808.3557567. URL: https://doi.org/10.1145/3511808.3557567.

Mansouri, B., D. W. Oard, and R. Zanibbi. (2022d). "Contextualized Formula Search Using Math Abstract Meaning Representation". In: *Proceedings of the 31st ACM International Conference on Information & Knowledge Management. CIKM '22*. Atlanta, GA, USA: Association for Computing Machinery. 4329–4333. ISBN: 9781450392365. DOI: 10.1145/3511808.3557567. URL: https://doi.org/10.1145/3511808.3557567.

Mansouri, B., S. Rohatgi, D. W. Oard, J. Wu, C. L. Giles, and R. Zanibbi. (2019a). "Tangent-CFT: An Embedding Model for Mathematical Formulas". en. In: *Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval*. Santa Clara CA USA: ACM. 11–18. ISBN: 978-1-4503-6881-0. DOI: 10.1145/3341981.3344235. URL: https://dl.acm.org/doi/10.1145/3341981.3344235 (accessed on 07/26/2021).

Mansouri, B., R. Zanibbi, and D. W. Oard. (2019b). "Characterizing Searches for Mathematical Concepts". en. In: *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*. Champaign, IL, USA: IEEE. 57–66. ISBN: 978-1-72811-547-4. DOI: 10.1109/JCDL.2019.00019. URL: https://ieeexplore.ieee.org/document/8791164/ (accessed on 07/26/2021).

Mansouri, B., R. Zanibbi, and D. W. Oard. (2021b). "Learning to Rank for Mathematical Formula Retrieval". en. In: *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Virtual Event Canada: ACM. 952–961. ISBN: 978-1-4503-8037-9. DOI: 10.1145/3404835.3462956. URL: https://dl.acm.org/doi/10.1145/3404835.3462956 (accessed on 07/26/2021).

Miao, S.-y., C.-C. Liang, and K.-Y. Su. (2020). "A Diverse Corpus for Evaluating and Developing English Math Word Problem Solvers". In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Ed. by D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault. Online: Association for Computational Linguistics. 975–984. DOI: 10.18653/v1/2020.acl-main.92. URL: https://aclanthology.org/2020.acl-main.92.

Mikolov, T., K. Chen, G. Corrado, and J. Dean. (2013). "Efficient Estimation of Word Representations in Vector Space". *arXiv preprint arXiv:1301.3781*.

Miller, B. R. and A. Youssef. (2003). "Technical Aspects of the Digital Library of Mathematical Functions". en: 19.

Miner, R. and R. Munavalli. (2007). "An Approach to Mathematical Search Through Query Formulation and Data Normalization". In: *International Conference on Mathematical Knowledge Management*. Springer. 342–355.

Mishra, S., M. Finlayson, P. Lu, L. Tang, S. Welleck, C. Baral, T. Rajpurohit, O. Tafjord, A. Sabharwal, P. Clark, and A. Kalyan. (2022). "LILA: A Unified Benchmark for Mathematical Reasoning". In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Ed. by Y. Goldberg, Z. Kozareva, and Y. Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics. 5807–5832. DOI: 10.18653/v1/2022.emnlp-main.392. URL: https://aclanthology.org/2022.emnlp-main.392.

Mišutka, J. and L. Galamboš. (2008). "Extending Full Text Search Engine for Mathematical Content". en: 14.

Mitra, A., H. Khanpour, C. Rosset, and A. Awadallah. (2024). "Orca-Math: Unlocking the potential of SLMs in Grade School Math". arXiv: 2402.14830 [cs.CL]. URL: https://arxiv.org/abs/2402.14830.

Nagy, G. and S. Seth. (1984). "Hierarchical Representation of Optically Scanned Documents". In: *Proc. International Conference on Pattern Recognition*. Vol. 1. Montréal, Canada. 347–349.

Ng, Y. K., D. J. Fraser, B. Kassaie, G. Labahn, M. S. Marzouk, F. W. Tompa, and K. Wang. (2020). "Dowsing for Math Answers with Tangent-L". en: 40.

Ng, Y. K., D. J. Fraser, B. Kassaie, and F. W. Tompa. (2021). "Dowsing for Answers to Math Questions: Ongoing Viability of Traditional MathIR". In: *CLEF (Working Notes)*. 63–81.

Nishizawa, G. (2020). "Visual structure editing of math formulas". *MA thesis*. Rochester Institute of Technology.

Nogueira, R., Z. Jiang, and J. Lin. (2021). "Investigating the Limitations of Transformers with Simple Arithmetic Tasks". arXiv: 2102.13019 [cs.CL]. URL: https://arxiv.org/abs/2102.13019.

Norman, D. A. (1988). *The psychology of everyday things*. Basic books.

Ohri, A. and T. Schmah. (2021). "Machine Translation of Mathematical Text". *IEEE Access*. 9: 38078–38086. DOI: 10.1109/ACCESS.2021.3063715.

Okamoto, M. and B. Miao. (1991). "Recognition of mathematical expressions by using the layout structures of symbols". In: *Prof. International Conference on Document Analysis and Recognition*. Vol. 1. St. Malo, France. 242–250.

Papineni, K., S. Roukos, T. Ward, and W.-J. Zhu. (2002). "Bleu: a Method for Automatic Evaluation of Machine Translation". In: *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Ed. by P. Isabelle, E. Charniak, and D. Lin. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics. 311–318. DOI: 10.3115/1073083.1073135. URL: https://aclanthology.org/P02-1040.

Pavan Kumar, P., A. Agarwal, and C. Bhagvati. (2012). "A Structure Based Approach for Mathematical Expression Retrieval". en. In: *Multi-disciplinary Trends in Artificial Intelligence*. Ed. by D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, C. Sombattheera, N. K. Loi, R. Wankar, and T. Quan. Vol. 7694. Berlin, Heidelberg: Springer Berlin Heidelberg. 23–34. ISBN: 978-3-642-35454-0 978-3-642-35455-7. DOI: 10.1007/978-3-642-35455-7_3. URL: http://link.springer.com/10.1007/978-3-642-35455-7_3 (accessed on 07/26/2021).

Peng, S., K. Yuan, L. Gao, and Z. Tang. (2021a). "MathBERT: A Pre-Trained Model for Mathematical Formula Understanding". en. *arXiv:2105.00377 [cs]*. May. URL: http://arxiv.org/abs/2105.00377 (accessed on 07/26/2021).

Peng, S., K. Yuan, L. Gao, and Z. Tang. (2021b). "Mathbert: A pre-trained model for mathematical formula understanding". *arXiv preprint arXiv:2105.00377*.

Perozzi, B., R. Al-Rfou, and S. Skiena. (2014). "Deepwalk: Online Learning of Social Representations". In: *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 701–710.

Petersen, F., M. Schubotz, A. Greiner-Petter, and B. Gipp. (2023). "Neural Machine Translation for Mathematical Formulae". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by A. Rogers, J. Boyd-Graber, and N. Okazaki. Toronto, Canada: Association for Computational Linguistics. 11534–11550. DOI: 10.18653/v1/2023.acl-long.645. URL: https://aclanthology.org/2023.acl-long.645.

Pfahler, L. and K. Morik. (2020). "Semantic Search in Millions of Equations". en. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Virtual Event CA USA: ACM. 135–143. ISBN: 978-1-4503-7998-4. DOI: 10.1145/3394486.3403056. URL: https://dl.acm.org/doi/10.1145/3394486.3403056 (accessed on 07/26/2021).

Piękos, P., M. Malinowski, and H. Michalewski. (2021). "Measuring and Improving BERT's Mathematical Abilities by Predicting the Order of Reasoning." In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Ed. by C. Zong, F. Xia, W. Li, and R. Navigli. Online: Association for Computational Linguistics. 383–394. DOI: 10.18653/v1/2021.acl-short.49. URL: https://aclanthology.org/2021.acl-short.49.

Pirolli, P. and S. Card. (1999). "Information Foraging". *Psychological Review*. 106(4): 643–675. DOI: 10.1037/0033-295x.106.4.643.

Rafailov, R., A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. (2023). "Direct Preference Optimization: Your Language Model is Secretly a Reward Model". In: *Advances in Neural Information Processing Systems*. Ed. by A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine. Vol. 36. Curran Associates, Inc. 53728–53741. URL: https://proceedings.neurips.cc/paper_files/paper/2023/file/a85b405ed65c6477a4fe8302b5e06ce7-Paper-Conference.pdf.

Reichenbach, M. S., A. Agarwal, and R. Zanibbi. (2014). "Rendering expressions to improve accuracy of relevance assessment for math search". In: *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. 851–854.

Reimers, N. and I. Gurevych. (2019). "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks". en. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics. 3980–3990. DOI: 10.18653/v1/D19-1410. URL: https://www.aclweb.org/anthology/D19-1410 (accessed on 07/27/2021).

Reusch, A., J. Gonsior, C. Hartmann, and W. Lehner. (2024). "Investigating the Usage of Formulae in Mathematical Answer Retrieval". In: *Advances in Information Retrieval*. Ed. by N. Goharian, N. Tonellotto, Y. He, A. Lipani, G. McDonald, C. Macdonald, and I. Ounis. Cham: Springer Nature Switzerland. 247–261. ISBN: 978-3-031-56027-9.

Reusch, A., M. Thiele, and W. Lehner. (2022). "Transformer-Encoder-Based Mathematical Information Retrieval". In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. Ed. by A. Barrón-Cedeño, G. Da San Martino, M. Degli Esposti, F. Sebastiani, C. Macdonald, G. Pasi, A. Hanbury, M. Potthast, G. Faggioli, and N. Ferro. Cham: Springer International Publishing. 175–189. ISBN: 978-3-031-13643-6.

Robertson, S. (2004). "Understanding inverse document frequency: on theoretical arguments for IDF." *J. Documentation*. 60(5): 503–520. DOI: 10.1108/00220410410560582. URL: https://doi.org/10.1108/00220410410560582.

Robertson, S. E. and S. Walker. (1994). "Some Simple Effective Approximations to the 2-Poisson Model for Probabilistic Weighted Retrieval". In: *Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval. Dublin, Ireland, 3-6 July 1994 (Special Issue of the SIGIR Forum)*. Ed. by W. B. Croft and C. J. van Rijsbergen. ACM/Springer. 232–241. DOI: 10.1007/978-1-4471-2099-5\_24. URL: https://doi.org/10.1007/978-1-4471-2099-5%5C_24.

Robertson, S. E. and H. Zaragoza. (2009). "The Probabilistic Relevance Framework: BM25 and Beyond". *Found. Trends Inf. Retr.* 3(4): 333–389. DOI: 10.1561/1500000019. URL: https://doi.org/10.1561/1500000019.

Rohatgi, S., W. Zhong, R. Zanibbi, J. Wu, and C. L. Giles. (2019). "Query Auto Completion for Math Formula Search". arXiv: 1912.04115 [cs.IR].

Roy, S. and D. Roth. (2015). "Solving General Arithmetic Word Problems". In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Ed. by L. Màrquez, C. Callison-Burch, and J. Su. Lisbon, Portugal: Association for Computational Linguistics. 1743–1752. DOI: 10.18653/v1/D15-1202. URL: https://aclanthology.org/D15-1202.

Ruzicka, M., P. Sojka, M. Líška, *et al.* (2016). "Math Indexer and Searcher under the Hood: Fine-tuning Query Expansion and Unification Strategies". In: *Proc. of the 12th NTCIR Conference on Evaluation of Information Access Technologies. Noriko Kando, Tetsuya Sakai, and Mark Sanderson,(Eds.) NII Tokyo*. 331–337.

Sakai, T. (2007). "Alternatives to Bpref". en. In: *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '07.* Amsterdam, The Netherlands: ACM Press. 71. ISBN: 978-1-59593-597-7. DOI: 10.1145/1277741.1277756. URL: http://portal.acm.org/citation.cfm?doid=1277741.1277756 (accessed on 07/26/2021).

Satpute, A., N. Gießing, A. Greiner-Petter, M. Schubotz, O. Teschke, A. Aizawa, and B. Gipp. (2024). "Can LLMs Master Math? Investigating Large Language Models on Math Stack Exchange". In: *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '24.* Washington DC, USA: Association for Computing Machinery. 2316–2320. ISBN: 9798400704314. DOI: 10.1145/3626772.3657945. URL: https://doi.org/10.1145/3626772.3657945.

Schellenberg, T., B. Yuan, and R. Zanibbi. (2012). "Layout-based substitution tree indexing and retrieval for mathematical expressions". In: *Document Recognition and Retrieval XIX, part of the IS&T-SPIE Electronic Imaging Symposium, Burlingame, California, USA, January 25-26, 2012, Proceedings.* Ed. by C. Viard-Gaudin and R. Zanibbi. Vol. 8297. *SPIE Proceedings.* SPIE. 82970I. DOI: 10.1117/12.912502. URL: https://doi.org/10.1117/12.912502.

Schubotz, M., A. Grigorev, M. Leich, H. S. Cohl, N. Meuschke, B. Gipp, A. S. Youssef, and V. Markl. (2016). "Semantification of Identifiers in Mathematics for Better Math Information Retrieval". In: *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '16.* Pisa, Italy: Association for Computing Machinery. 135–144. ISBN: 9781450340694. DOI: 10.1145/2911451.2911503. URL: https://doi.org/10.1145/2911451.2911503.

Shah, A. K., A. Dey, and R. Zanibbi. (2021). "A Math Formula Extraction and Evaluation Framework for PDF Documents". In: *16th International Conference on Document Analysis and Recognition, ICDAR 2021, Lausanne, Switzerland, September 5-10, 2021, Proceedings, Part II.* Ed. by J. Lladós, D. Lopresti, and S. Uchida. Vol. 12822. *Lecture Notes in Computer Science.* Springer. 19–34. DOI: 10.1007/978-3-030-86331-9\_2. URL: https://doi.org/10.1007/978-3-030-86331-9%5C_2.

Shannon, C. E. (1948). "A mathematical theory of communication". *The Bell System Technical Journal.* 27(3): 379–423. DOI: 10.1002/j.1538-7305.1948.tb01338.x.

Shen, Y. and C. Jin. (2020). "Solving Math Word Problems with Multi-Encoders and Multi-Decoders". In: *Proceedings of the 28th International Conference on Computational Linguistics.* Ed. by D. Scott, N. Bel, and C. Zong. Barcelona, Spain (Online): International Committee on Computational Linguistics. 2924–2934. DOI: 10.18653/v1/2020.coling-main.262. URL: https://aclanthology.org/2020.coling-main.262.

Shirahama, K. and M. Grzegorzek. (2016). "Towards large-scale multimedia retrieval enriched by knowledge about human interpretation - Retrospective survey". *Multim. Tools Appl.* 75(1): 297–331. DOI: 10.1007/S11042-014-2292-8. URL: https://doi.org/10.1007/s11042-014-2292-8.

Shumailov, I., Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, and R. Anderson. (2024a). "The Curse of Recursion: Training on Generated Data Makes Models Forget". arXiv: 2305.17493 [cs.LG]. URL: https://arxiv.org/abs/2305.17493.

Shumailov, I., Z. Shumaylov, Y. Zhao, N. Papernot, R. Anderson, and Y. Gal. (2024b). "AI models collapse when trained on recursively generated data". *Nature*. July. DOI: 10.1038/s41586-024-07566-y. URL: https://doi.org/10.1038/s41586-024-07566-y.

Sojka, P. and M. Líška. (2011). "Indexing and Searching Mathematics in Digital Libraries". en. In: *Intelligent Computer Mathematics*. Ed. by J. H. Davenport, W. M. Farmer, J. Urban, and F. Rabe. Vol. 6824. Berlin, Heidelberg: Springer Berlin Heidelberg. 228–243. ISBN: 978-3-642-22672-4 978-3-642-22673-1. DOI: 10.1007/978-3-642-22673-1_16. URL: http://link.springer.com/10.1007/978-3-642-22673-1_16 (accessed on 07/26/2021).

Sojka, P., M. Růžička, and V. Novotný. (2018). "MIaS: Math-Aware Retrieval in Digital Mathematical Libraries". en. In: *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. Torino Italy: ACM. 1923–1926. ISBN: 978-1-4503-6014-2. DOI: 10.1145/3269206.3269233. URL: https://dl.acm.org/doi/10.1145/3269206.3269233 (accessed on 07/27/2021).

Song, Y. and X. Chen. (2021). "Searching for Mathematical Formulas Based on Graph Representation Learning". en. In: *Intelligent Computer Mathematics*. Ed. by F. Kamareddine and C. Sacerdoti Coen. Vol. 12833. Cham: Springer International Publishing. 137–152. ISBN: 978-3-030-81096-2 978-3-030-81097-9. DOI: 10.1007/978-3-030-81097-9_11. URL: https://link.springer.com/10.1007/978-3-030-81097-9_11 (accessed on 07/26/2021).

Stalnaker, D. and R. Zanibbi. (2015). "Math expression retrieval using an inverted index over symbol pairs". In: *Document Recognition and Retrieval XXII, San Francisco, California, USA, February 11-12, 2015*. Ed. by E. K. Ringger and B. Lamiroy. Vol. 9402. *SPIE Proceedings*. SPIE. 940207. DOI: 10.1117/12.2074084. URL: https://doi.org/10.1117/12.2074084.

Stathopoulos, Y., S. Baker, M. Rei, and S. Teufel. (2018). "Variable Typing: Assigning Meaning to Variables in Mathematical Text". In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Ed. by M. Walker, H. Ji, and A. Stent. New Orleans, Louisiana: Association for Computational Linguistics. 303–312. DOI: 10.18653/v1/N18-1028. URL: https://aclanthology.org/N18-1028.

Stathopoulos, Y. and S. Teufel. (2015). "Retrieval of Research-level Mathematical Information Needs: A Test Collection and Technical Terminology Experiment". In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. 334–340.

Sudholt, S. and G. A. Fink. (2016). "PHOCNet: A Deep Convolutional Neural Network for Word Spotting in Handwritten Documents". In: *2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. IEEE Computer Society. 277–282.

Tan, M., L. Wang, L. Jiang, and J. Jiang. (2022). "Investigating Math Word Problems using Pretrained Multilingual Language Models". In: *Proceedings of the 1st Workshop on Mathematical Natural Language Processing (MathNLP)*. Ed. by D. Ferreira, M. Valentino, A. Freitas, S. Welleck, and M. Schubotz. Abu Dhabi, United Arab Emirates (Hybrid): Association for Computational Linguistics. 7–16. DOI: 10.18653/v1/2022.mathnlp-1.2. URL: https://aclanthology.org/2022.mathnlp-1.2.

Tausczik, Y. R., A. Kittur, and R. E. Kraut. (2014). "Collaborative problem solving: a study of MathOverflow". In: *Computer Supported Cooperative Work, CSCW '14, Baltimore, MD, USA, February 15-19, 2014*. Ed. by S. R. Fussell, W. G. Lutters, M. R. Morris, and M. C. Reddy. ACM. 355–367. DOI: 10.1145/2531602.2531690. URL: https://doi.org/10.1145/2531602.2531690.

Teachers of Mathematics, N. C. of. (2011). "Strategic Use of Technology in Teaching and Learning Mathematics - National Council of Teachers of Mathematics". URL: https://www.nctm.org/Standards-and-Positions/Position-Statements/Strategic-Use-of-Technology-in-Teaching-and-Learning-Mathematics/ (accessed on 10/27/2021).

Thanda, A., A. Agarwal, K. Singla, A. Prakash, and A. Gupta. (2016). "A Document Retrieval System for Math Queries". en: 8.

Truong, T.-N., C. Nguyen, R. Zanibbi, H. Mouchére, and M. Nakagawa. (2024). "A survey on handwritten mathematical expression recognition: The rise of encoder-decoder and GNN models". *Pattern Recognition*: (to appear).

Tucker, Q. (2024). "Improving Histogram-Based Mathematical Formula Search with Local Matching and Learned Embeddings". *Tech. rep.* Rochester Institute of Technology.

Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. (2017). "Attention Is All You Need". *CoRR*. abs/1706.03762. arXiv: 1706.03762. URL: http://arxiv.org/abs/1706.03762.

Wallace, E., Y. Wang, S. Li, S. Singh, and M. Gardner. (2019). "Do NLP Models Know Numbers? Probing Numeracy in Embeddings". In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Ed. by K. Inui, J. Jiang, V. Ng, and X. Wan. Hong Kong, China: Association for Computational Linguistics. 5307–5315. DOI: 10.18653/v1/D19-1534. URL: https://aclanthology.org/D19-1534.

Wang, K., J. Pan, W. Shi, Z. Lu, M. Zhan, and H. Li. (2024). "Measuring Multi-modal Mathematical Reasoning with MATH-Vision Dataset". arXiv: 2402.14804 [cs.CV]. URL: https://arxiv.org/abs/2402.14804.

Wang, M., Y. Tang, J. Wang, and J. Deng. (2017a). "Premise Selection for Theorem Proving by Deep Graph Embedding". In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Vol. 30. Curran Associates, Inc. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/18d10dc6e666eab6de9215ae5b3d54df-Paper.pdf.

Wang, Q., C. Kaliszyk, and J. Urban. (2018). "First Experiments with Neural Translation of Informal to Formal Mathematics". In: *Intelligent Computer Mathematics*. Ed. by F. Rabe, W. M. Farmer, G. O. Passmore, and A. Youssef. Cham: Springer International Publishing. 255–270. ISBN: 978-3-319-96812-4.

Wang, X., C. Macdonald, N. Tonellotto, and I. Ounis. (2023). "Reproducibility, Replicability, and Insights into Dense Multi-Representation Retrieval Models: from ColBERT to Col". In: *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*. Ed. by H. Chen, W. ( Duh, H. Huang, M. P. Kato, J. Mothe, and B. Poblete. ACM. 2552–2561. DOI: 10.1145/3539618.3591916. URL: https://doi.org/10.1145/3539618.3591916.

Wang, Y., X. Liu, and S. Shi. (2017b). "Deep Neural Solver for Math Word Problems". In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Ed. by M. Palmer, R. Hwa, and S. Riedel. Copenhagen, Denmark: Association for Computational Linguistics. 845–854. DOI: 10.18653/v1/D17-1088. URL: https://aclanthology.org/D17-1088.

Wang, Z., A. Lan, and R. Baraniuk. (2021). "Mathematical formula representation via tree embeddings". *Online: https://people. umass. edu/˜ andrewlan/papers/preprint-forte. pdf.*

Wangari, K. D. V., R. Zanibbi, and A. Agarwal. (2014). "Discovering real-world use cases for a multimodal math search interface". In: *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. ACM. DOI: 10.1145/2600428.2609481. URL: https://doi.org/10.1145%2F2600428.2609481.

Wei, J., X. Wang, D. Schuurmans, M. Bosma, b. ichter brian, F. Xia, E. Chi, Q. V. Le, and D. Zhou. (2022). "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models". In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. Vol. 35. Curran Associates, Inc. 24824–24837. URL: https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf.

Welleck, S., J. Liu, R. L. Bras, H. Hajishirzi, Y. Choi, and K. Cho. (2021). "Natural-Proofs: Mathematical Theorem Proving in Natural Language". In: *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*. URL: https://openreview.net/forum?id=Jvxa8adr3iY.

Wenzel, M., L. C. Paulson, and T. Nipkow. (2008). "The Isabelle Framework". In: *Theorem Proving in Higher Order Logics*. Ed. by O. A. Mohamed, C. Muñoz, and S. Tahar. Berlin, Heidelberg: Springer Berlin Heidelberg. 33–38. ISBN: 978-3-540-71067-7.

White, R. W. (2016). *Interactions with search systems*. Cambridge University Press.

Wu, Y., A. Q. Jiang, W. Li, M. Rabe, C. Staats, M. Jamnik, and C. Szegedy. (2022). "Autoformalization with Large Language Models". In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. Vol. 35. Curran Associates, Inc. 32353–32368. URL: https://proceedings.neurips.cc/paper_files/paper/2022/file/d0c6bc641a56bebee9d985b937307367-Paper-Conference.pdf.

Xie, Y., H. Mouchère, F. Simistira Liwicki, S. Rakesh, R. Saini, M. Nakagawa, C. T. Nguyen, and T.-N. Truong. (2023). "ICDAR 2023 CROHME: Competition on Recognition of Handwritten Mathematical Expressions". In: *Document Analysis and Recognition - ICDAR 2023*. Ed. by G. A. Fink, R. Jain, K. Kise, and R. Zanibbi. Cham: Springer Nature Switzerland. 553–565. ISBN: 978-3-031-41679-8.

Xue, L., A. Barua, N. Constant, R. Al-Rfou, S. Narang, M. Kale, A. Roberts, and C. Raffel. (2022). "ByT5: Towards a Token-Free Future with Pre-trained Byte-to-Byte Models". *Transactions of the Association for Computational Linguistics*. 10: 291–306. Ed. by B. Roark and A. Nenkova. DOI: 10.1162/tacl_a_00461. URL: https://aclanthology.org/2022.tacl-1.17.

Yang, K., A. Swope, A. Gu, R. Chalamala, P. Song, S. Yu, S. Godil, R. J. Prenger, and A. Anandkumar. (2023). "LeanDojo: Theorem Proving with Retrieval-Augmented Language Models". In: *Advances in Neural Information Processing Systems*. Ed. by A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine. Vol. 36. Curran Associates, Inc. 21573–21612. URL: https://proceedings.neurips.cc/paper_files/paper/2023/file/4441469427094f8873d0fecb0c4e1cee-Paper-Datasets_and_Benchmarks.pdf.

Yasunaga, M. and J. D. Lafferty. (2019). "TopicEq: A Joint Topic and Mathematical Equation Model for Scientific Texts". en. *Proceedings of the AAAI Conference on Artificial Intelligence*. 33(July): 7394–7401. ISSN: 2374-3468, 2159-5399. DOI: 10.1609/aaai.v33i01.33017394. URL: https://wvvw.aaai.org/ojs/index.php/AAAI/article/view/4728 (accessed on 07/26/2021).

Yu, L., W. Jiang, H. Shi, J. Yu, Z. Liu, Y. Zhang, J. T. Kwok, Z. Li, A. Weller, and W. Liu. (2024). "MetaMath: Bootstrap Your Own Mathematical Questions for Large Language Models". arXiv: 2309.12284 [cs.CL]. URL: https://arxiv.org/abs/2309.12284.

Yu, M., G. Li, D. Deng, and J. Feng. (2016). "String similarity search and join: a survey". *Frontiers Comput. Sci.* 10(3): 399–417. DOI: 10.1007/S11704-015-5900-5. URL: https://doi.org/10.1007/s11704-015-5900-5.

Yuan, K., D. He, Z. Jiang, L. Gao, Z. Tang, and C. L. Giles. (2020). "Automatic Generation of Headlines for Online Math Questions". *Proceedings of the AAAI Conference on Artificial Intelligence*. 34(05): 9490–9497. DOI: 10.1609/aaai.v34i05.6493. URL: https://ojs.aaai.org/index.php/AAAI/article/view/6493.

Yue, X., X. Qu, G. Zhang, Y. Fu, W. Huang, H. Sun, Y. Su, and W. Chen. (2023). "MAmmoTH: Building Math Generalist Models through Hybrid Instruction Tuning". arXiv: 2309.05653 [cs.CL]. URL: https://arxiv.org/abs/2309.05653.

Zanibbi, R., A. Aizawa, and M. Kohlhase. (2016a). "NTCIR-12 MathIR Task Overview". en. In: *Prof. NTCIR-12*. 10.

Zanibbi, R. and D. Blostein. (2012). "Recognition and retrieval of mathematical expressions". en. *International Journal on Document Analysis and Recognition (IJDAR)*. 15(4): 331–357. ISSN: 1433-2833, 1433-2825. DOI: 10.1007/s10032-011-0174-4. URL: http://link.springer.com/10.1007/s10032-011-0174-4 (accessed on 07/26/2021).

Zanibbi, R., K. Davila, A. Kane, and F. W. Tompa. (2016b). "Multi-Stage Math Formula Search: Using Appearance-Based Similarity Metrics at Scale". en. In: *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. Pisa Italy: ACM. 145–154. ISBN: 978-1-4503-4069-4. DOI: 10.1145/2911451.2911512. URL: https://dl.acm.org/doi/10.1145/2911451.2911512 (accessed on 07/26/2021).

Zanibbi, R. and L. Yu. (2011). "Math Spotting: Retrieving Math in Technical Documents Using Handwritten Query Images". In: *2011 International Conference on Document Analysis and Recognition, ICDAR 2011, Beijing, China, September 18-21, 2011.* IEEE Computer Society. 446–451. DOI: 10.1109/ICDAR.2011.96. URL: https://doi.org/10.1109/ICDAR.2011.96.

Zeng, H., C. Luo, B. Jin, S. M. Sarwar, T. Wei, and H. Zamani. (2024). "Scalable and Effective Generative Information Retrieval". In: *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*. Ed. by T. Chua, C. Ngo, R. Kumar, H. W. Lauw, and R. K. Lee. ACM. 1441–1452. DOI: 10.1145/3589334.3645477. URL: https://doi.org/10.1145/3589334.3645477.

Zhang, A., A. Goyal, W. Kong, H. Deng, A. Dong, Y. Chang, C. A. Gunter, and J. Han. (2015). "adaQAC: Adaptive Query Auto-Completion via Implicit Negative Feedback". In: *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*. 143–152.

Zhang, D., L. Wang, L. Zhang, B. T. Dai, and H. T. Shen. (2020a). "The Gap of Semantic Parsing: A Survey on Automatic Math Word Problem Solvers". *IEEE Trans. Pattern Anal. Mach. Intell.* 42(9): 2287–2305. DOI: 10.1109/TPAMI.2019.2914054. URL: https://doi.org/10.1109/TPAMI.2019.2914054.

Zhang, M., S. Baral, N. Heffernan, and A. Lan. (2022). "Automatic Short Math Answer Grading via In-context Meta-learning". In: *Proceedings from the 15th International Conference on Educational Data Mining*.

Zhang, M., Z. Cui, M. Neumann, and Y. Chen. (2018). "An End-to-End Deep Learning Architecture for Graph Classification". *Proceedings of the AAAI Conference on Artificial Intelligence*. 32(1). DOI: 10.1609/aaai.v32i1.11782. URL: https://ojs.aaai.org/index.php/AAAI/article/view/11782.

Zhang, T., V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi. (2020b). "BERTScore: Evaluating Text Generation with BERT". arXiv: 1904.09675 [cs.CL]. URL: https://arxiv.org/abs/1904.09675.

Zhao, J., M.-Y. Kan, and Y. L. Theng. (2008). "Math information retrieval: user requirements and prototype implementation". en. In: *Proceedings of the 8th ACM/IEEE-CS joint conference on Digital libraries - JCDL '08*. Pittsburgh PA, PA, USA: ACM Press. 187. ISBN: 978-1-59593-998-2. DOI: 10.1145/1378889.1378921. URL: http://portal.acm.org/citation.cfm?doid=1378889.1378921 (accessed on 07/26/2021).

Zhong, W., S.-C. Lin, J.-H. Yang, and J. Lin. (2023). "One Blade for One Purpose: Advancing Math Information Retrieval using Hybrid Search". In: *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '23*. Taipei, Taiwan: Association for Computing Machinery. 141–151. ISBN: 9781450394086. DOI: 10.1145/3539618.3591746. URL: https://doi.org/10.1145/3539618.3591746.

Zhong, W. and R. Zanibbi. (2019). "Structural Similarity Search for Formulas Using Leaf-Root Paths in Operator Subtrees". en. In: *Advances in Information Retrieval*. Ed. by L. Azzopardi, B. Stein, N. Fuhr, P. Mayr, C. Hauff, and D. Hiemstra. Vol. 11437. Cham: Springer International Publishing. 116–129. ISBN: 978-3-030-15711-1 978-3-030-15712-8. DOI: 10.1007/978-3-030-15712-8_8. URL: http://link.springer.com/10.1007/978-3-030-15712-8_8 (accessed on 07/26/2021).

Zhong, W., X. Zhang, J. Xin, R. Zanibbi, and J. Lin. (2021). "Approach Zero and Anserini at the CLEF-2021 ARQMath Track: Applying Substructure Search and BM25 on Operator Tree Path Tokens". *Proc. CLEF 2021 (CEUR Working Notes)*.

Zhu, L., C. Zheng, W. Guan, J. Li, Y. Yang, and H. T. Shen. (2024). "Multi-Modal Hashing for Efficient Multimedia Retrieval: A Survey". *IEEE Trans. Knowl. Data Eng.* 36(1): 239–260. DOI: 10.1109/TKDE.2023.3282921. URL: https://doi.org/10.1109/TKDE.2023.3282921.

Zobel, J. and A. Moffat. (2006). "Inverted Files for Text Search Engines". *ACM Comput. Surv.* 38(2): 6–es. ISSN: 0360-0300. DOI: 10.1145/1132956.1132959. URL: https://doi.org/10.1145/1132956.1132959.