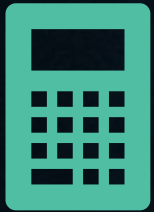# MATH INFORMATION RETRIEVAL

David Na

# WHAT IS MATH INFORMATION RETRIEVAL?

Math IR is "concerned with finding information in documents that include mathematics".

Document similarity and retrieval is calculated considering any math keywords and/or formulae found in the collection and query.
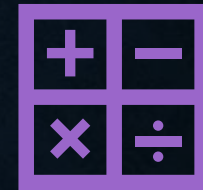
It's the same as standard Information Retrieval but instead of querying just strings in documents, you query using some mathematical context.

Math IR is essentially a performing "standard" IR using Math concepts

**First Math-aware search engines developed in 2000:**

[NIST Digital Library of Mathematical Functions (DLMF)](#)

**Examples:**

Given a mathematical concept in keywords and/or formulae, find the technical papers that use similar mathematical models.

Another example is you can browse for documents containing a given formula.
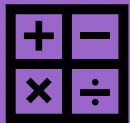
# DIFFICULTIES OF MATH IR

Mathematical notation is dialectic, some variables represent similar concepts but are named differently depending on the subject area

Operations can be defined in different ways depending on context. For instance, a line over a variable can be a division simple or a Boolean negation depending on context.

Math formulae can look very similar but still be very different.

Most of the time Math IR queries are combinations of keywords and formulae.

# TYPES OF MATH IR SYSTEMS

## Formula Retrieval

- Related to computation
- Example: How do you compute this formula?

## Text Answer Retrieval

- Related to a Mathematical Concept
- Example: What is a Hamiltonian Circuit?

## Formula and Text Retrieval

- Related to Proofs
- Example: Prove $a$ is divisible by $b$

# FORMULA RETRIEVAL

**TEXT ANSWER RETRIEVAL**

# TWO MAIN MATH IR CONFERENCES

## NTCIR – NII Testbeds and Community for Information Access Research

- Run by National Institute of Informatics (Tokyo, Japan)
- Math IR was only the focus in workshops 10, 11 and 12
- Held from 2014-2016

## ARQMath - Answer Retrieval for Questions on Math

- Hosted by Rochester Institute of Technology and the NSF
- ARQMath has had 3 iterations
- More Recently held in 2019-2021

# NTCIR

- NTCIR 10-12 (2014-2016): Math IR
  - 3 Main Tasks:
    - Formula Search: Search given a formula
    - Formula + Text Search: Search documents given keywords and formula
    - Open Information Retrieval: Search the collection given a text query
  - NTCIR-12 focused on the Wikipedia dataset
    - Wikipedia Formula Browsing: Focusing on formula search
    - Formula Similarity Search: Finding formulas related to the query
  - Two Datasets: Wikipedia was introduced in NTCIR 12
    - 100,000 scientific articles from arXiv and 35,000 Wikipedia articles
    - All HTML files
    - Includes 30-60M math formulae
  - Given a set of queries, the Math IR system must return a ranked list of search results from both datasets.
    - Queries will be a set of formulae and keywords and you must search within the document collection

https://www.cs.rit.edu/~rlaz/files/ntcir12-mathir.pdf

# ARQMATH

- ARQMath 1, 2 and 3 (2019-2021):

  - Focused directly on Answer Retrieval

  - Dataset: Question + Answer posts from Stack Exchange

    - Data from 2010 to 2020

    - 28 Million mathematical formulas in LaTeX

  - Task 1: Finding relevant answers to a mathematical question

    - Given a math question as a query, search all answers and return relevant answers

  - Task 2: Search based on a formula

    - Given a question with an identified formula as a query, return relevant formulas

  - Task 3: Open Domain Question Answering (ARQMath 3)

    - Like Task 1, but now the question can come from anywhere (even written by hand). Task 1 used questions from the Stack Exchange dataset as input. Task 3 can take as input any query.

http://sigir.org/wp-content/uploads/2020/12/p06.pdf

# REPRESENTING MATH FORMULAS

**LaTeX**
- Standard
- Commonly used in PDFs

**MathML**
- XML Format
- HTML <math> tag

**Tangent-CFT**
- Symbol Layout Trees (SLT)
- Operator Trees (OPTs)

# LATEX

Reference: http://www.malinc.se/math/latex/basiccodeen.php

# MATHML

# TOKENIZING MATHEMATICAL FORMULAE



(a) Formula  (b) Symbol Layout Tree  (c) Operator Tree

Reference: https://clgiles.ist.psu.edu/pubs/IDTIR2019.pdf

# TOKENIZING MATHEMATICAL FORMULAE – TANGENT-L AND SLT



Original Formula:

$$y_i^j = 1 + x^2$$

Converted Symbol Layout Tree (SLT):

Tangent-L's Math Tokens:

**Symbol Pairs** — For each edge, start and end symbols and edge label

$(y, j, \nearrow)$ $(y, =, \rightarrow)$ $(y, i, \searrow)$ $(=, 1, \rightarrow)$

$(1, +, \rightarrow)$ $(+, x, \rightarrow)$ $(x, 2, \nearrow)$

**Terminal symbols** — List of symbols with no outedges

$(j, \triangle)$ $(i, \triangle)$ $(2, \triangle)$

**Compound symbols** — List of outedge labels for nodes with more than one outedge

$(y, \nearrow \searrow \rightarrow)$

**Augmented locations** — For each token of the above three types, that token together with the SLT path to the token's (first) symbol from the root node

$(y, j, \nearrow, \emptyset)$ $(y, =, \rightarrow, \emptyset)$ $(y, i, \searrow, \emptyset)$ $(=, 1, \rightarrow, \rightarrow)$

$(1, +, \rightarrow, \rightarrow\rightarrow)$ $(+, x, \rightarrow, \rightarrow\rightarrow\rightarrow)$ $(x, 2, \nearrow, \rightarrow\rightarrow\rightarrow\rightarrow)$

$(j, \triangle, \nearrow)$ $(i, \triangle, \searrow)$ $(2, \triangle, \rightarrow\rightarrow\rightarrow\rightarrow\nearrow)$ $(y, \nearrow\searrow\rightarrow, \emptyset)$

A detailed animation provided here: https://www.scg.uwaterloo.ca/brushsearch/tangent-l/

# BASELINE MATH IR SYSTEM – TANGENT-CFT

# INSPIRATION: MATHDOWSERS



Reference: https://cs.uwaterloo.ca/~yk2ng/MathDowsers-ARQMath/ & https://ceur-ws.org/Vol-3180/paper-03.pdf

# MY MATH IR SYSTEM

- GitHub: https://github.com/davidna22/math-IR-ARQMath-CompuBERT

- Overview:
    - My math IR system was created following the ARQMath workshop. My primary goal was to create an IR system capable of vectorizing q+a pairs from the ARQMath dataset, in order to be able to take in a math query and return a ranked list of documents matching that query. I chose to develop my system following Task 1 of the ARQMath workshop.

- Dataset:
    - Math Stack Exchange posts from 2010-2021
    - Data preprocessing was provided by ARQMath and found on their github
    - I implemented their preprocessing, added some others that I found online in papers

- Architecture:
    - I use a variation of a model called CompuBERT (Novotny et al.) that I found very interesting. It utilizes pre-trained Transformer models such as BERT to create feature embeddings for question answer pairs. Using those feature embeddings, I can compute cosine similarity scores and return a ranked list of documents. Most of my code was taken from their open-source implementation (https://github.com/MIR-MU/CompuBERT) as it was one of the first transformer-encoder models and easiest to implement.

# COMPUBERT



**Q:** *Can anyone explain ...*   **A_relevant:** *Consider $c^2 = a^2 + b^2$...*   **A_irrelevant:** *Ask elsewhere ...*

$$\text{minimize} \quad \sum_{i=1}^{|Qs|} \sum_{j=1}^{|A_i|} \left| (1 - \cos(q_i, a_{ij})) - \text{dist}_{\exp}(q_i, a_{ij}) \right|$$

Reference: https://ceur-ws.org/Vol-2696/paper_235.pdf

The Model relies on the feature vector embeddings generated by the BERT model and uses cosine similarity to find the best answers per query

The overall results of the this first transformer model was poor relative to other submissions, so part of my experiment was to test how it will perform with a newer pre-trained model

Since CompuBERT, other Transformer models have been implemented to success. This (given my lack of experience) was the easiest to implement and thus, a good starting point for me.

# TRANSFORMER-ENCODER AND DECODER MDOELS



**Figure 1:** Overview of our approach for Task 1 - Mathematical Answer Retrieval including examples for training and evaluation data.

Reference: https://ceur-ws.org/Vol-3180/paper-07.pdf

# MY MODEL

Math-aware ALBERT - https://huggingface.co/AnReu/math_albert

ALBERT for ARQMath 3

MPNet – Trained by Microsoft

MathBERT

Embeddings are calculated for all corpus documents and stored on disk

Query Embeddings are passed and Cosine Similarity Score is calculated

Ranked List of Documents is returned to the User

# RESULTS

| | pid | doc | score |
|---|---|---|---|
| 1 | tensor(13090, device='cuda:0') | For p \geq 1 the generalized mean defines a norm, because it is the \el... | tensor(0.9593, device='cuda:0') |
| 2 | tensor(43157, device='cuda:0') | Well, actually it is well known that \\A\\ = \\A\\_{1\rightarrow\infty} = \m... | tensor(0.9572, device='cuda:0') |
| 3 | tensor(27479, device='cuda:0') | Another approach and somehow an answer to your second question: Thi... | tensor(0.9536, device='cuda:0') |
| 4 | tensor(8717, device='cuda:0') | Hint . The Fourier transform preserves norms on L^2 , i.e. \\f\\_2 = \\\ha... | tensor(0.9531, device='cuda:0') |
| 5 | tensor(45930, device='cuda:0') | Comment converted to answer as suggested by Jonas Teuwen: Hint: I... | tensor(0.9527, device='cuda:0') |
| 6 | tensor(27213, device='cuda:0') | If we can show that A doesn't increase the 1-norm, i.e., \\Ax\\_1\leq\\x\\... | tensor(0.9519, device='cuda:0') |
| 7 | tensor(29167, device='cuda:0') | The other answers cover how to solve this with Lagrange multipliers, but ... | tensor(0.9515, device='cuda:0') |
| 8 | tensor(36557, device='cuda:0') | The rough idea is to show a series of inequalities: \int\|fgh\|\leq\|fg\|... | tensor(0.9511, device='cuda:0') |
| 9 | tensor(13837, device='cuda:0') | Since I could not find a formula for an expectation value of a generalized... | tensor(0.9507, device='cuda:0') |
| 10 | tensor(13595, device='cuda:0') | If we look at the discrete group \mathbb{Z} , we have L^1(\mathbb{Z})=\... | tensor(0.9503, device='cuda:0') |
| 11 | tensor(39592, device='cuda:0') | Suppose that fg\in L^1 , but there is no C so that \\fg\\_{L^1}\le C\\g\\_{... | tensor(0.9498, device='cuda:0') |
| 12 | tensor(44480, device='cuda:0') | The usual definition of the operator norm is \\A\\_{\mathrm{op}} = \sup_{... | tensor(0.9494, device='cuda:0') |
| 13 | tensor(38730, device='cuda:0') | Using the normalization of the Fourier Transform shown above, we get ... | tensor(0.9479, device='cuda:0') |
| 14 | tensor(14897, device='cuda:0') | I'm not sure if this is what you're asking: A norm on a vector space is in... | tensor(0.9479, device='cuda:0') |
| 15 | tensor(32114, device='cuda:0') | I won't include all details since you don't seem to want that... Up to min... | tensor(0.9470, device='cuda:0') |
| 16 | tensor(44250, device='cuda:0') | If 1 \leq p \leq q \lt \infty then \\x\\_{q} \leq \\x\\_{p} and clearly \\x\\_p \g... | tensor(0.9467, device='cuda:0') |
| 17 | tensor(7486, device='cuda:0') | Using the spectral theorem to obtain an orthogonal basis of eigenvector... | tensor(0.9454, device='cuda:0') |
| 18 | tensor(23948, device='cuda:0') | You ought to specify what space a_i and u are coming from, and what ... | tensor(0.9448, device='cuda:0') |
| 19 | tensor(7860, device='cuda:0') | Well there are slightly weaker condition that improves the result you cite.... | tensor(0.9448, device='cuda:0') |
| 20 | tensor(43757, device='cuda:0') | In general, \varphi(f(x)) = f(x) h(x) is impossible, because the right side ... | tensor(0.9447, device='cuda:0') |

| | MPNet | Math-Aware Albert | Albert for Math 3 | MathBERT |
|---|---|---|---|---|
| f1 Scores | 0.699 | 0.626 | 0.691 | 0.644 |
| Precision | 0.544 | 0.532 | 0.54 | 0.527 |
| Accuracy | 0.546 | 0.51 | 0.538 | 0.506 |

# TO-DO LIST

- Trying different Transformer Models
  - There have been new iterations of CompuBERT that have had high performance. An example is found in this paper, Transformer-Encoder and Decoder Models for Questions on Math (https://ceur-ws.org/Vol-3180/paper-07.pdf).

- Training on a larger subsample of the dataset
  - Training on the entire dataset would have taken me about a full month of constant training given the size of the corpus and the size of the models I chose.
  - Bunch of limitations such as GPU RAM and compute power that I would need to pay for via cloud to train effectively

- Applying different transformations in text preprocessing
  - Math text preprocessing is very different from standard IR, so I would love to do more research and apply different preprocessing techniques

- Attempting Task 2 and 3 of the ARQMath competition
  - Task 2 is very different as I would need to fully implement a math-only textual embedding. Most likely requiring the use of symbol trees (Tangent-CFT) and something called Formula2Vec.

# REFERENCES

- NTCIR-12 MathIR Task Overview
  - https://www.cs.rit.edu/~rlaz/files/ntcir12-mathir.pdf
- NTCIR Website
  - https://ntcir-math.nii.ac.jp/introduction/
- ARQMath: A New Benchmark for Math-Aware CQA and Math Formula Retrieval
  - http://sigir.org/wp-content/uploads/2020/12/p06.pdf
- ARQMath Website
  - https://www.cs.rit.edu/~dprl/ARQMath/2021/
- Tangent-CFT: An Embedding Model for Mathematical Formulas
  - https://clgiles.ist.psu.edu/pubs/IDTIR2019.pdf
- Math ML Docs
  - https://www.w3.org/TR/mathml-core/
- Latex Docs
  - http://www.malinc.se/math/latex/basiccodeen.php
- CompuBERT
  - https://github.com/MIR-MU/CompuBERT
- Three is Better than One: Ensembling Math Information Retrieval Systems
  - https://ceur-ws.org/Vol-2696/paper_235.pdf
- Transformer-Encoder and Decoder Models for Questions on Math
  - https://ceur-ws.org/Vol-3180/paper-07.pdf