

Hunter T_EXNet | CSCI 350 Project Proposal

Alex Taradachuck Ralph “Blake” Venté

March 2020

Abstract

The OpenAI request for research spawned research in using Visual Neural Models for translating mathematical expressions into markup code. Since the debut, progress has been largely incremental. Singh 2018, Wang 2019, and Bender 2019 all attribute properties of their respect models to the biases in the corpus and its pre-processing. We present a new dataset and a pipeline for harvesting examples from source files hoping to minimize these biases.

1 Previous Work

1.1 Debut

The prospect of accurately transcribing mathematical expression into a markup representation is enticing because it opens the doors for bringing new life to old mathematical texts or those for which the source code is unavailable.

A Harvard project called “What you get is what you see: A visual markup decompiler” by Deng, Kanervisto, and Rush documents strategies for machine translation of mathematical notation using an attention-based encoder-decoder neural model. Notably, the researchers were interested in the influence of markup alone on the efficacy of a model — without providing explicit information about underlying grammars [1, p. 1]. The previous state of the art required exhaustively documenting the grammar to allow a conventional character recognition model to translate the expression.

Introduced in 2016 by deng2016you, the `im2latex100k` dataset has facilitated research by reducing the overhead of heavy pre-processing. Using this dataset, other solutions to IM2LATEX have seen substantial progress. Recently, it has come under question as a possible *bottleneck* for progress. We discuss this in greater detail in Section 2.2.

1.2 im2latex Problem Specification

Before delving into how architectures were designed to solve the problem, let us first define it. Formally, this problem is a sequence-to-sequence image captioning task. A neural model learns an abstract encoding and its respective decoding into the original sequence [1, p. 1]. At inference time, when presented a single image, the trained model predicts one token at a time, aiming to reconstruct the original sequence. All the while, the model receives as input the output of previous steps. This is accomplished using Long-Short Term Memory, which we outline in 1.3.

1.3 Architectural Designs

Since its inception, solutions to IM2LATEX have used neural models. In fact, OpenAI’s Request for Research recommended the use of a neural model and even specified the sequence-to-sequence attention mechanism that makes these models suitable for the problem, but there have been a wide variety in the particular architectures of the neural models.

Deng, Kanervisto, and Rush Genthial and Sauvestre Wang and Liu Bender, Haurilet, Roitberg, et al.

2 Corpus Harvesting

2.1 Challenges

We harvested examples from Cornell’s arXiv archive. As the arXiv strictly forbids scraping, we found a patron who properly obtained copies and uploaded them to archive.org. We wrote a script to extract the source code files, and notably, we use the `pandoc` document parser.

On this note, \LaTeX and its subset, \TeX pose an unique set of challenges compared to other steps of the pre-processing pipeline. \LaTeX is a Turing-Complete language, equivalent in power to general purpose programming languages. From the theoretical standpoint, it may contain loops and recursion, classifying it as a recursively enumerable language. As such, no \LaTeX parser is guaranteed to terminate, and regular expressions alone lack the expressive power to extract all cases of a particular pattern. Practically, we mostly needed to contend with macro expansion. This means that the original scripts of [1] required substantial changes before being suitable for our purposes.

2.2 Innovations

To start, we follow closely in the footsteps of Singh. After preprocessing steps on `im2latex100k`, Singh arrived upon a dataset of 93,784 compiling examples, which he regarded as “too small for good generalization.” As such, he mined additional examples from the KDD Cup 2003 dataset, resulting in about additional 50,000 examples [5, p. 8].

Practically, `pandoc` handled the files from the corpus gracefully, but due to unlinked files and mismatched encodings, we set `timeout 5` to terminate the program after 5 seconds of stalling. `pandoc` also facilitates extracting expressions. Without it, we were presented with the challenging prospect of filtering out text not containing mathematical expressions, but `pandoc` generates an abstract syntax tree and encases the mathematical expressions in the more “specific” expression `\((.*)\)` compared to the expression `$(.*)$` which resulted in false matches that contained only text. One such false match that was eliminated is provided in Figure 2. As a fall-back, we rely on Singh’s script to test for the presence of at least one $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ command (from the command vocabulary, Σ) before writing out the final formula code. This has the added benefit of removing examples that are strictly numbers and those that are otherwise linear strings of alphanumeric characters.

Most importantly, `pandoc` expands all the macros in the source documents. Without this step, papers in the dataset that relied heavily on macros to shorten the length of the written expressions contributed very few formulas to the resulting collection because they caused definition-not-found errors and were eliminated in later steps. In [1] the processed formulas were discarded implicitly because they caused such errors, but with the macros expanding, documents containing macros reclaimed their representation in the final set, eliminating one possible source of bias in the dataset.

From Deng, Kanervisto, and Rush, we also reiterate that several $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ input sequences map to the same output image, so we use the same normalization with Khan Academy’s $\text{K}^{\text{A}}\text{T}_{\text{E}}\text{X}$ as they did in their scripts. Like `pandoc`, $\text{K}^{\text{A}}\text{T}_{\text{E}}\text{X}$ parses the input sequences into an abstract syntax tree and rebuilds the expressions without disrupting semantics.

Bender, Haurilet, Roitberg, et al., p. 5 note that their `IM2LATEX` model struggled very short sequences, and performed better on those that were slightly longer. They attribute this to an imbalance in the dataset: very short formulas are underrepresented. In fact, this may be traced back to the scripts written by [1], which discarded examples shorter than 40 bytes in length. We halved the minimum length of an example.

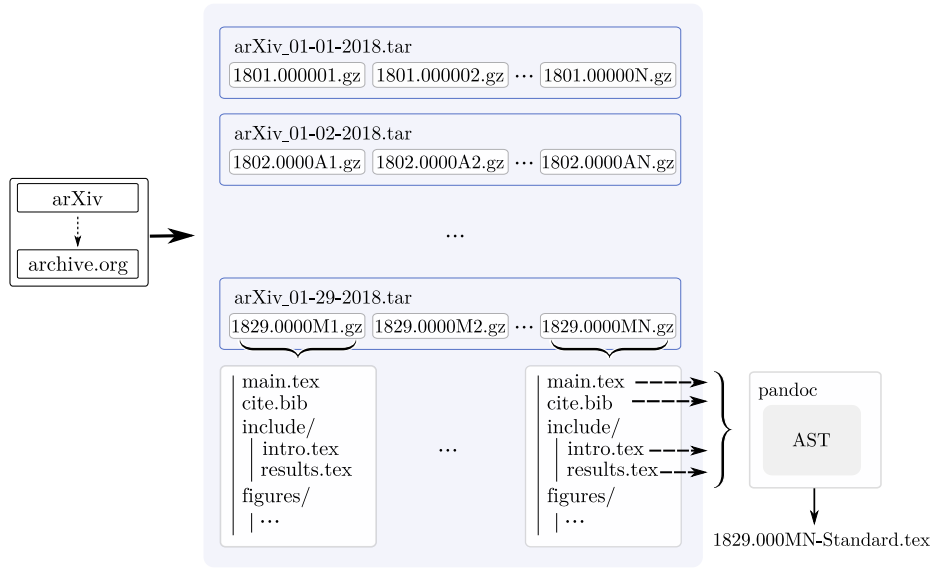


Figure 1: Examples flow from archive.org and are stored on our disks, where we extract the `tar` files first by day, then by submission, into the document folders. Then `pandoc` pulls in all of the dependencies of this particular document and writes out one stand-alone document with macros expanded. Then we match the regular expression `\((.*?)\)` that `pandoc` uses to wrap mathematical expressions. As a general tool `pandoc` can convert between other markup language in the same way, making it useful for standardizing data from disparate sources.

2.3 Product

In all, we mined 1 month of `tar` files, extracted 15,200 `tex` files, expanded macros 7,200 unique files and generated data 500,000 rendering examples. We decided to randomly subsample about 40% of that as our final candidate set, after deduplicating, we were left with 170,00 images. Our formula files can be found at <https://github.com/rvente/TeXNet.ai/Dataset>. A graphical view of our data pipeline can be seen in Figure 1.

3 Network Architecture

Our network has its basis in Singh’s research. The architecture has several essential features that we enumerate with brief explanations and proceed with detailed remarks on their implementation in our architecture.

Convolutional Neural Network (CNN) Layers in our neural model learn convolutions on image space that extracts meaningful information into a compressed internal embedding, that is, a vector which has a dense encoding of the contents of the input image.

Recurrent Neural Network (RNN) At time t , the network’s layer n takes as input the output of layer n at time $t - 1$. To train such a network using backpropagation, it is “unrolled” so as to emulate a feed-forward network, allowing the updates to propagate. This allows for reasoning over time because information from prior time-steps can factor into current time steps.

Long-Short Term Memory (LSTM) Loosely speaking, the LSTM cells in the neural network “learn what to forget and when to forget it.” This aspect of the architecture allows a particular subset of the previous layer’s output to be factored in dynamically based on what the LSTM cells have learned is important.

Attention mechanism Attention mechanisms allow for the network to mask way, or “ignore” portions of the input image. The network can now use information about past predictions to “keep track of where it’s looking,” allowing for focused steps in the prediction process.

3.1 Recurrent Neural Networks

The issue with an RNN alone is that memory of previous inputs quickly decay when moving from time-step to time-step.

3.2 Long-Short Term Memory Configuration

To solve these issues, researchers use need a network to selectively add and remove inputs into the current state [6]. Dating back to 1997, this LSTM architecture solves this using gates to control the flow of data during training and inference [6]. It takes 3 inputs, c_{t-1} memory at time $t - 1$,

4 Evaluation

For evaluation, we use Bilingual Evaluation Understudy (BLEU), the emerging standard among neural network solutions to the IM2LATEX problem. BLEU is an automated metric for assessing machine translation accuracy [7, p. 1]. It is computed by comparing a target translation (ground truth) to the predicted sequence.

We report 3 metrics: first we compute global BLEU score across all examples and then again with 2 splits. The first split is performance by length in tokens the target sequence.

From Papineni, Roukos, Ward, et al., BLEU score has a penalty for translations that are too long and those which are too short. The former case is handled implicitly, and the latter requires an explicit brevity penalty, BP computed as

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{1-r/c} & \text{otherwise} \end{cases} . \quad (1)$$

with r as the length of of the ground-truth sentence and c as the length of the candidate sentence. Then, we have

$$\text{BLEU} = \text{BP} \exp \left(\sum_{i=1}^n w_i \log p_i \right) , \quad (2)$$

where p_n is geometric mean of n -gram precisions from 1 to N and w_n is a parameter weight, such that $\sum_i w_i = 1$. Papineni, Roukos, Ward, et al. used $w_n = 1/N$. BLEU score has gained widespread use because it is predictive of human judgements.

5 Future work

References

- [DKR16] Yuntian Deng, Anssi Kanervisto, and Alexander M Rush. “What you get is what you see: A visual markup decompiler”. In: *arXiv preprint arXiv:1609.04938* 10 (2016), pp. 32–37.

- [GS16] Guillaume Genthial and Romain Sauvestre. “Image to Latex”. In: (2016).
- [WL19] Zelun Wang and Jyh-Charn Liu. “Translating Mathematical Formula Images to LaTeX Sequences Using Deep Neural Networks with Sequence-level Training”. In: *arXiv preprint arXiv:1908.11415* (2019).
- [Ben+19] Sidney Bender, Monica Haurilet, Alina Roitberg, et al. “Learning Fine-Grained Image Representations for Mathematical Expression Recognition”. In: *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*. Vol. 1. IEEE. 2019, pp. 56–61.
- [Sin18] Sumeet S Singh. “Teaching Machines to Code: Neural Markup Generation with Visual Attention”. In: *arXiv preprint arXiv:1802.05415* (2018).
- [Zha+19] Aston Zhang, Zachary C Lipton, Mu Li, et al. “Dive into Deep Learning”. In: *Unpublished draft. Retrieved 3* (2019), p. 319.
- [Pap+02] Kishore Papineni, Salim Roukos, Todd Ward, et al. “BLEU: a method for automatic evaluation of machine translation”. In: *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics. 2002, pp. 311–318.
- [Ben+18] Carolina Benedetti, Rafael S. González D’León, Christopher R. H. Hanusa, et al. *A combinatorial model for computing volumes of flow polytopes*. 2018. arXiv: 1801.07684 [math.CO].

A Qualitative Analysis

A.1 Reducing Plain Text Matches

The source code

```
1132 The net flow vector  $(1,1,0,\ldots,0)$  has previously been
      considered for the complete graph by Corteel, Kim, and Mészáros \cite{CKM}. They used the Lidskii formula~\eqref{eq:lidskiiivol} and constant term identities to derive the
      following product formula for the volume of  $\mathcal{F}_{K_{n+1}}(1,1,0,\ldots,0)$ .
1133 It would be of interest to rederive this result using the
      refined Kostant constants.
1134 \begin{theorem}[\cite{Theorem 1.1}{CKM}] \label{thm:volkn11}
1135 Let  $n \geq 1$ . For the complete graph  $K_{n+1}$ ,
      matched the regular expression  $$(*.*?)$$  with the string
      $. It would be of interest to rederive this result using the
      refined Kostant constants. \begin{theorem}[\cite{Theorem
      1.1}{CKM}] \label{thm:volkn11} Let $
```

Figure 2: From arXiv:1801.07684 Benedetti, D’León, Hanusa, et al.

This is a sample from a particular document whose source code induced a false match Math encased in dollar-signs has been deprecated in \LaTeX , but remains in wide use as an artifact from \TeX . It was deprecated to facilitate parsing: unlike $\backslash($ which indicates that it opens to the right, the dollar-sign has no such indication. The reason the previous scripts were matching these plain-text strings was because matching all strings that start and end in dollar-signs doesn’t preclude text from matching. In this case, the $\$$ of the first line matched with the first $\$$ of the 4th line.