

MathWriting: A Dataset For Handwritten Mathematical Expression Recognition

Philippe Gervais*[†]
 pgervais@acm.org
 Inceptive
 Switzerland

Anastasiia Fadeeva[†]
 fadeich@google.com
 Google DeepMind
 Zurich, Switzerland

Andrii Maksai
 amaksai@google.com
 Google DeepMind
 Zurich, Switzerland

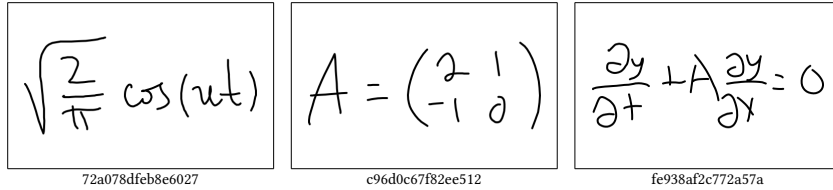


Figure 1: Three examples of HME from MathWriting. More examples can be found in Appendix J. Each ink is accompanied by a unique identifier that matches a corresponding filename in the dataset.

Abstract

Recognition of handwritten mathematical expressions allows to transfer scientific notes into their digital form. It facilitates the sharing, searching, and preservation of scientific information. We introduce MathWriting, the largest online handwritten mathematical expression dataset to date. It consists of **230k human-written samples** and an additional **400k synthetic ones**. This dataset can also be used in its rendered form for offline HME recognition. One MathWriting sample consists of a formula written on a touch screen and a corresponding \LaTeX expression. We also provide a normalized version of \LaTeX expression to simplify the recognition task and enhance the result quality. We provide baseline performance of standard models like OCR and CTC Transformer as well as Vision-Language Models like PaLI on the dataset. The dataset together with an example colab is accessible on Github.

CCS Concepts

• **Applied computing** → **Online handwriting recognition**; *Optical character recognition*.

Keywords

Mathematical Expression Recognition, Online Handwriting Recognition, Digital Ink

*Work performed while employed at Google

[†]Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'17, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
 ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

ACM Reference Format:

Philippe Gervais, Anastasiia Fadeeva, and Andrii Maksai. 2025. MathWriting: A Dataset For Handwritten Mathematical Expression Recognition. In . ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

MathWriting dataset (2.9 GB):

https://storage.googleapis.com/mathwriting_data/mathwriting-2024.tgz

Partial dataset (1.5 MB):

https://storage.googleapis.com/mathwriting_data/mathwriting-2024-excerpt.tgz

Online text recognition interprets handwritten input on devices with touchscreens like smartphones and tablets. This process analyzes the handwriting as a sequence of strokes, rather than a static image [14]. The recognition of mathematical expressions (MEs) presents a significant challenge, and has historically received less research attention compared to character and word recognition [29]. The fundamental differences between mathematical expression (ME) recognition and text recognition create obstacles to the direct application of improvements developed in one domain to the other. While mathematical expressions (MEs) utilize many of the same symbols as standard text, they are distinguished by a more rigid and two-dimensional structural organization, see Figure 1. Where text can be treated to some extent as a one-dimensional problem amenable to sequence modeling, MEs require two-dimensional analysis due to the importance of symbol spatial positioning. It is also different from symbol segmentation or object detection because the output of a recognizer has to contain the relationship between symbols, serialized in some form (\LaTeX , a graph, InkML, etc.). Handwritten mathematical expressions (HMEs), like handwritten text, are more challenging to recognize than their printed counterparts due to increased variability in writing and data scarcity.

Handwritten data acquisition is expensive due to the need for manual input with specialized hardware like touchscreens and digital pens. The publication of the MathWriting dataset aims to address

the data requirements for mathematical expression recognition research. Samples include a large number of human-written inks, as well as synthetic ones. MathWriting can readily be used with other online datasets like CROHME [33] or Detexify [9] - we publish the data in InkML format to facilitate this. It can also be used for offline ME recognition by rasterizing the inks, using code provided on the Github page¹.

MathWriting is the largest set of online HME published so far - both human-written and synthetic. It significantly expands the set of symbols covered by CROHME [33], enabling more sophisticated recognition capabilities. Since inks can be rasterized, MathWriting can also be seen as larger than existing offline HME datasets [8, 34, 1]. For these reasons we introduce a new benchmark, applicable to both online and offline ME recognition.

This work's main contributions are:

- a large dataset of Handwritten Mathematical Expressions under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International².
- \LaTeX ground truth expressions in normalized form to simplify training and to make evaluation more robust.
- Evaluation of different models like CTC Transformer and PaLI on the dataset to show what recognition quality could be achieved with the provided data.

The paper focuses on the high-level description of the dataset: creation process, postprocessing, train/test split, ground truth normalization, statistics, and a general discussion of the dataset content to help practitioners understand what can and cannot be achieved with it. All the low-level technical information like file formats can be found in the `readme.md` file present at the root of the dataset archive linked above. We also provide code examples on Github¹, to show how to read the various files, process and rasterize the inks, and tokenize the \LaTeX ground truth.

2 Related work

Mathematical expression recognition offers valuable tools for educational software, including features like Math Notes on iPhone [3]. As tablets such as the Google Pixel, iPad, and Remarkable gain widespread adoption, the demand for effective handwriting tools is growing.

There are multiple datasets, such as Aida Calculus Math Handwriting Recognition Dataset [1] and IM2LATEX-100K [28], that focus on image-based (offline) math recognition. The Aida dataset consists solely of synthetically generated handwritten math expressions. IM2LATEX-100K has become a key benchmark for assessing performance in typeset math formula recognition. There is also a dataset for typeset math formula detection - TFD-ICDAR 2019 [21] that was part of Competition on Recognition of Handwritten Mathematical Expressions (CROHME) 2019. The most widely used dataset for evaluating online math recognition systems originates from the CROHME competition [33]. We establish in Section 4.4 that the MathWriting dataset surpasses CROHME with a richer vocabulary and more unique formulas.

The highest-performing model in the CROHME 2019 competition [21] was an LSTM-based encoder-decoder, achieving its result

through training on an augmented dataset. This highlights the critical role of training data volume and methods for synthetic data generation. The MyScript model [23], a CTC BLSTM architecture, achieved third place in the competition. Similarly, Section 5.2 reveals that a CTC Transformer yields strong results on the MathWriting dataset. We also present results obtained with the encoder-decoder model PaLI [6].

3 Dataset Creation

MathWriting dataset primarily consists of \LaTeX expressions from Wikipedia, more details about the acquisition of expressions are provided in Appendix B. These expressions were used for both ink collection from human contributors Section 3.1 as well as synthetic data generation Section 3.2. We did a very limited filtering of very noisy human-written examples (described in Appendix C).

3.1 Ink Collection

Inks were obtained from human contributors through an in-house Android app. Participants agreed to the standard Google terms of use and privacy policy. The task consisted in copying a rendered mathematical expression (prompt) shown on the device's screen using either a digital pen or a finger on a touch screen. Mathematical expressions used as prompt were first obtained in \LaTeX format, then rendered into a bitmap through the \LaTeX compiler (see Appendix A for the template used). 95% of MathWriting expressions were obtained from Wikipedia. The remaining ones were generated to cover underrepresented cases in Wikipedia, like isolated letters with nested sub/superscripts or complicated fractions (see Section B). Contributors were hired internally at Google. 6 collection campaigns were run between 2016 and 2019, each lasting between 2 to 3 weeks. Collected data contains only inks and labels, so no personally identifiable information is present in the dataset. Offensive content is highly unlikely because \LaTeX expressions were taken from Wikipedia and we conducted a filtering of noisy data (described in Appendix C).

3.2 Synthetic Samples and Isolated Symbols

We created synthetic samples in order to further increase the label diversity for training. This also enabled compensating for limitations of the human collection like the maximum length of the expressions, which were limited by the size of the screen they were written on. We used \LaTeX expressions from Wikipedia that were not used in the data collection. The resulting synthetic data has a 90th percentile of expression length of 68 characters, compared to 51 in train. This is especially important as deep neural nets often fail to generalize to inputs longer than their training data [2, 31]. Using synthetic long inks together with the original human-written inks can help to eliminate that problem as shown in [30, 24]. The synthesis technique is as follows: starting from a raw \LaTeX mathematical expression, we computed a DVI file using the \LaTeX compiler, from which we extracted bounding boxes. We then used those bounding boxes to place handwritten individual symbols, resulting in a complete expression. See Figure 2 for an example of extracted bounding boxes and the resulting synthetic example. We make the bounding box data public to facilitate custom synthetic dataset creation.

¹<https://github.com/google-research/google-research/tree/master/mathwriting>

²<https://creativecommons.org/licenses/by-nc-sa/4.0/>

$$((p + q) + (p - q)) / 2 = q$$

Figure 2: An example of a synthetic ink created from bounding boxes with label $((p+q)+(p-q))/2=q$

Inks for individual symbols are all from the `symbols` split. They have been manually extracted from inks in `train`. For each symbol, we manually selected strokes corresponding to it for 20-30 distinct occurrences in `train`, and used that information to generate a set of individual inks. Similar synthesis techniques have been used by [33] with inks, [8] and [1] with raster images.

The synthetic inks can incorporate symbols from multiple participants which can result in statistical inconsistency. However, even human handwriting displays occasional stylistic variations within one sample (see Appendix J.1 ink `02229a0c174d8dbe`). We argue that having more diversity in synthetic dataset is valuable for training robust recognition models. We also show in ablation experiments (Section 5.3) that synthetic data improves the quality of recognition.

Another important distinction between the synthetic and human-written inks is the stroke order. For synthetic inks, stroke order follows the order of the bounding boxes in the DVI file, which can be different from the usual order of writing for mathematical expressions. However, the writing order within a given symbol is consistent with human writing.

3.3 Dataset Split

MathWriting is composed of five different sets of samples, which we call 'splits': `train`, `valid`, `test`, `symbols`, and `synthetic`. The splits `train`, `valid` and `test` consist only of human-written examples. The split `symbols` is provided for synthetic data generation and is not used in training. The split of human-written samples between `train`, `valid` and `test` was partially done based on writers, partially based on labels. More details are provided in Appendix K. Experiments have shown that a more important factor than the handwriting style was whether the *label* had already been seen during training. This fact is also supported by research in the area of compositional generalization [18]. In the published version, `valid` has a 55% (8.5k samples) intersection with `train` based on unique normalized labels, and `test` has an 8% intersection (647 samples). We chose to have a low intersection between `train` and `test` in order to correctly measure generalization of trained models to unseen labels.

3.4 Label Normalization

All samples in the dataset come with two labels: the \LaTeX expression that was used during the data collection (annotation `label` in the InkML files), and a normalized version of it meant for model training, which is free from a few sources of confusions for an ML model (annotation `normalizedLabel`). An example with original

and normalized labels is provided in Figure 3. Label normalization covers three main categories (details are provided in Appendix L):

- variations used in print that can't be reproduced in handwriting - e.g. bold, italic - or that haven't been reproduced consistently by contributors.
- non-uniqueness of the \LaTeX syntax. e.g. `\frac{1}{2}` and `1\over 2` are equivalent.
- visual variations that can be reproduced in handwriting but can't reliably be inferred by a model. This includes size modifiers like `\left`, `\right`.

We provide the raw labels to make it possible to experiment with alternative normalization schemes, which could lead to better outcomes for different applications.

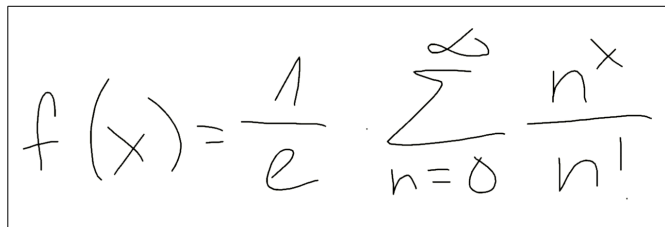
3.4.1 Limitations of normalization. The normalization process is purely syntactic, and can not cover cases where the meaning of the expression has to be taken into account. For example, a lot of expressions from Wikipedia use `cos` instead of `\cos`. It is often clear to a human reader whether the sequence of characters `c,o,s` represents the `\cos` command or simply three letters. However, this can not be reliably inferred by a syntactic parser, for example in `tacos` vs `ta\cos`. An alternative would be to update the raw labels, which we didn't do because we wanted to keep the information that was used during the collection as untouched as possible. Similarly, cases like `10^{-1}` usually mean `\{10\}^{-1}`, though they render exactly the same. We made the choice to normalize to the former because it's the only option with a purely syntactic normalizer. It's also better than not removing these extra braces because it gives more consistent label structures, which simplifies the model training problem.

4 Dataset Statistics

In this section we describe the key characteristics of MathWriting and compare it to CROHME23 [33]. In Table 1 we provide the information about the volume of the dataset splits both in terms of examples (inks) and unique labels.

Table 1: Statistics on different subsets of MathWriting dataset.

	<code>train</code>	<code>synthetic</code>	<code>valid</code>	<code>test</code>
# distinct inks	230k	396k	16k	8k
# distinct labels	53k	396k	8k	4k



51b364eb9ba2185a

Figure 3: An example from the train split, with its labels: Raw: $f(x)=\frac{1}{e}\cdot \sum_{n=0}^{\infty}\frac{x^n}{n!}$
Normalized: $f(x)=\frac{1}{e}\cdot\sum_{n=0}^{\infty}\frac{n^x}{n!}$

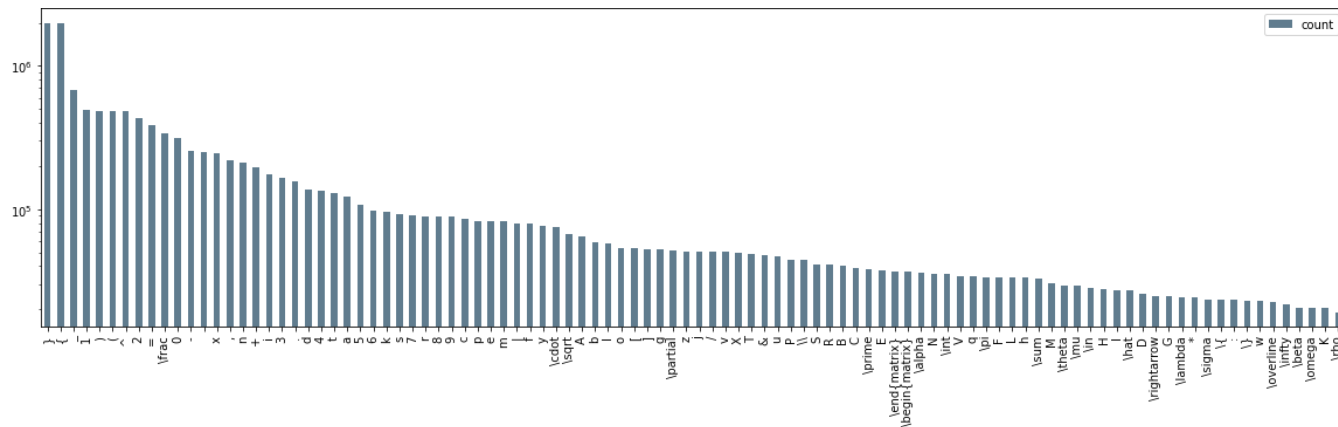


Figure 4: Histogram of the top-100 most frequent tokens in MathWriting.

4.1 Label Statistics

MathWriting contains 457k unique labels after normalization (see Section 3.4). From Table 1 we see that most unique expressions are covered by the synthetic portion of the dataset. However, the absolute number of unique expressions in human-written part is still high – 61k. This underlines the importance of synthetic data as it allows models to see a much bigger variety of expressions. It is important to note that the synthetic split has essentially no repeated expressions. On the other hand, in real data multiple different writings of the same expression are quite common (see Figure 12 in Appendix F). This fact allows us to separately evaluate model’s quality on expressions that were observed during training and that those that hadn’t. As seen in Table 2 the biggest intersection in expressions is between valid and train. The minimal overlap between test and train splits is beneficial for assessing a model’s ability to generalize to expressions that were not seen in train.

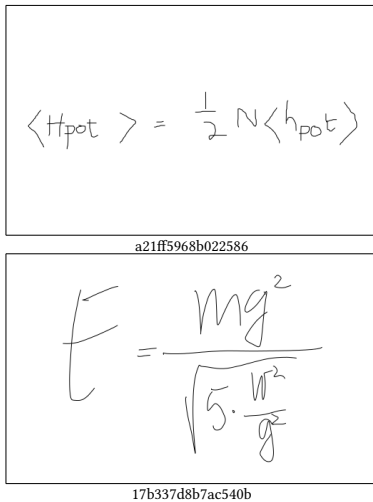
Table 2: Counts of unique labels shared between MathWriting splits

	train	synthetic	valid	test
train	-	0	3.6k	355
synthetic	0	-	0	0
valid	3.6k	0	-	239
test	355	0	239	-

Table 3: Ink statistics for MathWriting.

	10th percentile	median	90th percentile
# strokes	5	14	39
# points	131	350	1069
writing time (sec)	1.88	6.03	16.42
aspect ratio	1.32	3.53	9.85

The median length of expressions in characters is 26 which is comparable to one of the most popular English recognition datasets IAMonDB [20] which has median of 29 characters. However, it is important to note that \LaTeX expressions have tokens that span multiple characters like `\frac`. The median length of expressions in tokens (provided in Appendix I) is 17, thus making training a model on tokens rather than characters easier due to shorter target lengths [32, 25]. We want to emphasize that MathWriting can be used with a different tokenization scheme and token vocabulary from what we propose in Appendix I. In Figure 4 we show the number of occurrences for the most frequent tokens. Tokens `{` and `}` are by far the most frequent as they are integral to the \LaTeX syntax.



**Figure 5: Top: an ink with very low sampling rate (9.4 points per second)
Bottom: an ink with very high sampling rate (260 points per second)**

4.2 Ink Statistics

Each ink in MathWriting dataset is a sequence of strokes $I = [s_0, \dots, s_n]$, each stroke s_j consisting of points. A point is represented as a triplet (x, y, t) where x and y are coordinates on the screen and t is a timestamp. In Table 3 we provide statistics on number of strokes, points, and duration of writing. It’s important to note that as inks were collected on different devices, the absolute coordinate values can vary a lot. In human-written data the time information t always starts from 0 but it is not always the case in the synthetic split. Different samples often have different sampling rates (number of points written in one second) due to the use of different devices (see Figure 5). More details in Section 4.3. Consequently, the same ink written on two different devices can result in inks with a different number of points. For human-written inks, the sampling rate is consistent between strokes, but it is not the case for synthetic ones. In order to accommodate a model and make sequences shorter, inks can be resampled in time.

4.3 Devices Used

Around 150 distinct device types have been used by contributors. In most cases inks were written on smartphones using a finger on a touchscreen. However, there are cases where tablets with styluses were used. The main device used in this case is Google Pixelbook, which accounted for 51k inks total (see Table 8, Appendix F). Out of all device types, 37 contributed more than 1000 inks. Note that writing on a touchscreen with a finger or a stylus results in different low-level artifacts. All devices were running the same Android application for ink collection, regardless of whether their operating system was Android or ChromeOS.

Table 4: Counts of inks, distinct labels and distinct tokens used in MathWriting and CROHME23. The single token present in CROHME23 but not in MathWriting is the literal dollar sign $\backslash\$$.

	MathWriting	CROHME23	Common
Inks	650k	164k	0
Labels	457k	102k	47k
Vocab	254	105	104

Table 5: Count of human-written and synthetic inks for MathWriting and CROHME23. Human-written inks represent 38% of the total for MathWriting, and 10% for CROHME23.

	MathWriting	CROHME23
human	253k	17k
synthetic	396k	147k

4.4 Comparison With CROHME23

In this section we compare main dataset statistics of MathWriting and CROHME23 [33] as it is a popular publicly available dataset for HME recognition. In terms of overall size, MathWriting has nearly 3.9 times as many samples and 4.5 times as many distinct labels after normalization, see Table 4. A significant number of labels can be found in both datasets (47k), but the majority is dataset-specific. This suggests that combining both datasets during training could yield improved HME recognition quality. MathWriting has more human-written inks than CROHME23 as seen in Table 5, and contains a much larger variety of tokens. It has 254 distinct tokens including all Latin capital letters and almost the entire Greek alphabet. It also contains matrices, which are not included in CROHME23. Therefore, more scientific fields like quantum mechanics, differential calculus, and linear algebra can be represented using MathWriting.

5 Experiments

5.1 Evaluation setup

We propose the following evaluation setup based on MathWriting for the quality of handwriting math expression recognition.

- **evaluation samples:** the test split of MathWriting.
- **metric:** character error rate (CER) [22], where a "character" is a \LaTeX token as defined by the code in Appendix M.

We provide a reference implementation of the evaluation metric at the Github page ¹. We propose the use of CER as a metric to make results comparable to other recognition tasks like text recognition [26, 17], and the use of \LaTeX tokens instead of ASCII characters so that an error on a single non-latin letter (e.g. $\backslash\alpha$ recognized as a) counts as one instead of many.

5.2 Baseline Recognition Models

Table 6 compares different recognition methods discussed in Section 2. All models are trained exclusively on the MathWriting dataset (train and synthetic), except for the OCR API that was trained on other datasets as well. The following models represent

Table 6: Recognition results for different models. The evaluation metric is reported on both the valid and test splits.

Model	Time	Params	CER ↓	valid		test		
				EM ↑	≤ 1 dist ↑	CER ↓	EM ↑	≤ 1 dist ↑
OCR [12]	no	-	6.50	64	76	7.17	53	68
CTC Transformer [10]	yes	35M	4.52 (0.08)	71 (0.46)	81 (0.3)	5.49 (0.05)	60 (0.42)	72 (0.42)
PaLI [6]	yes	700M	4.47 (0.08)	76 (0.25)	83 (0.11)	5.95 (0.06)	64 (0.35)	73 (0.22)
PaLIGemma [4]	yes	3B	3.95 (0.04)	80 (0.04)	86 (0.03)	5.97 (0.08)	69 (0.46)	77 (0.26)

Ground Truth $\psi_2(\Omega_2^{\Omega_2}) \frac{1}{z}(z+1) \theta \in [q_i, q_i+1]$

Predictions $\psi_1(\Omega_2^{\boxed{\Omega_2}}) \frac{1}{2}(\boxed{2}+1) \theta \in [q_i, q_i+\boxed{1}]$

Figure 6: Examples of recognition mistakes from the CTC Transformer model. We observe similar mistakes from the other models.

different approaches to handwriting recognition – offline [17], on-line [5] and mixed [10]. Through fine-tuning, we establish a benchmark for the quality attainable with MathWriting data.

OCR. This is a publicly available Document AI OCR API [12], which processes bitmap images. It has been trained partly on samples from MathWriting. We sent inks rendered with black ink on a white background and searched for optimal image size and stroke width to get the best evaluation result from the model. We also fine-tuned TrOCR model following methodology outlined in [27] and observed only 25% Exact Match on valid split. Similar behavior was observed with TrOCR [19] on CROHME dataset [15] and the possible explanation is that the text tokenizer used by TrOCR is not well adapted to handle LaTeX expressions.

CTC Transformer. This model is a transformer base with a Connectionist Temporal Classification loss on top (CTC) [13]. It contains 11 transformer layers with an embedding size of 512. We used swish activation function and dropout of 0.15 as those parameters performed best on valid. We train with an Adam optimizer, learning rate of 1e-3, batch size 256 for 100k steps. One training run took 4 hours on 4 TPU v2. We trained from scratch and exclusively on MathWriting (train and synthetic). The model is similar to [5], replacing LSTM layers by Transformer layers and not using any external language model on top.

VLM. We fine-tuned a large Vision-Language Model PaLI [6] with encoder-decoder architecture on MathWriting (train and synthetic). We used the representation proposed in [10] where an ink is represented as both a sequence of points (similar to CTC Transformer) and its rasterized version (similar to OCR). We train three models with different data shuffling for 200k steps with batch size 128, learning rate 0.3 and dropout 0.2. One training run took 14 hours on 16 TPU v5p. Models were fine-tuned exclusively on train and synthetic MathWriting data. Overall, it took 2 TPU v2 days and 28 TPU v5p days to run the experiments.

We also fine-tuned a publicly available PaLIGemma model [4] that utilizes a decoder-only LLM – Gemma [11]. This model has recently shown strong performance on captioning and VQA tasks after fine-tuning. The model is trained on image input with 448px resolution. We leverage the image representation with speed information from [10], that demonstrated superior performance compared to black-on-white rendering. We’ve trained the model with learning rate 1e-4 and batch size 512 on 64 TPU v5p over 36 hours.

Table 6 shows the evaluation comparison between the four models. The OCR model has no information about the order of writing and speed (offline recognition), which explains its lower performance than methods that take time information into account (online recognition). The other methods – PaLI, PaLIGemma and CTC Transformer perform significantly better than OCR. These results show that our dataset can be used to train classical recognition models like CTC transformer as well as more recent architectures like VLMs.

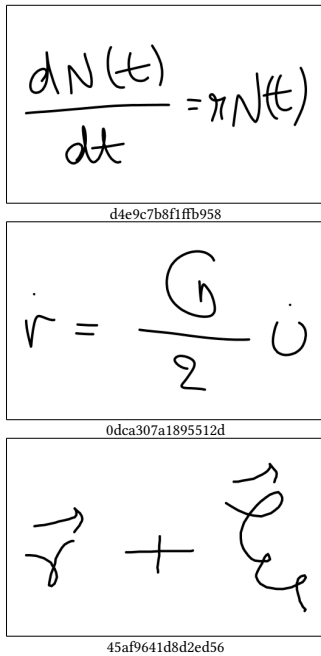
Figure 6 shows examples of model mistakes. Two of the main causes of mistakes are confusing similar-looking characters like “z” and “2”, and errors in the structural arrangement of the characters, for instance not placing a sub-expression in a subscript or superscript.

5.3 Synthetic data ablation

In Section 3.2 we described the synthetic dataset generation process. This section examines the empirical impact of synthetic data on the model’s performance metrics. We have opted for a CTC Transformer model for this purpose as it allows for faster experimentation. Table 7 shows that there is a 10% decrease in CER on the test dataset when synthetic data is used in training, and close to 5% on the valid dataset. These results demonstrate the substantial benefit of synthetic data that we provide for recognition training.

Table 7: Ablation of synthetic data on CTC transformer model.

synthetic	CER eval ↓	CER test ↓
with	4.52 (0.08)	5.49 (0.05)
without	4.64 (0.07)	6.2 (0.08)

**Figure 7: Three ways of writing a lowercase 'r'.**

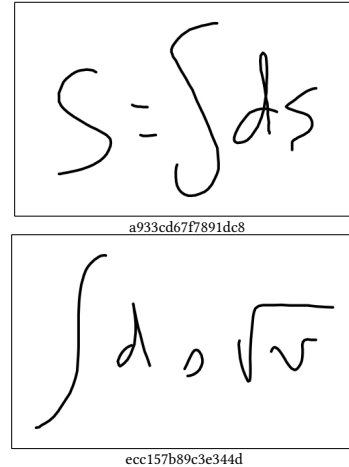
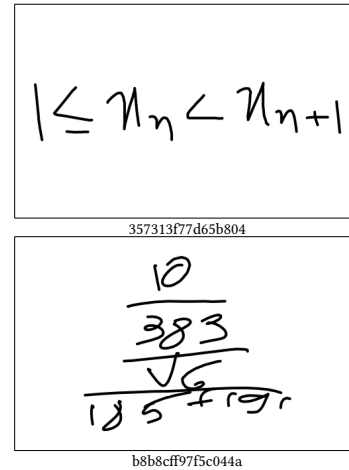
6 Discussion

6.1 Differences in Writing Style

The number of contributors was large enough that a variety of writing styles are represented in the dataset. An example for different ways of writing letter 'r' can be seen in Figure 7. Additional examples are provided in Figure 8. Similar though less obvious differences exist for other letters. Style differences also show through writing order (example – Figure 14, Appendix G).

6.2 Recognition Challenges

MathWriting presents some inherent recognition challenges, which are typical of handwritten representations. For example, it's not really possible to distinguish these pairs from the ink alone: $\frac{\underline{a}}{b}$ vs $\frac{a}{\overline{b}}$, and $\overline{\omega}$ vs ϖ . We'd like to point out that these ambiguities are not an issue for humans in practice, because they rely on contextual information to disambiguate: a particular writing idiosyncrasy, consistency with nearby expressions, knowledge of the scientific domain, etc. See Figures 9 and 10 for more examples.

**Figure 8: Two ways of writing lowercase s.****Figure 9: Top: character ambiguity. Is it $1 \leq x_n < x_{n+1}$ or $1 \leq n_\eta < n_{\eta+1}$? Bottom: what is the fraction nesting order?**

6.3 Dataset Applications and Future Work

Mathwriting can be used to train recognizers for a large variety of scientific fields, and is also large enough to enable synthesis of mathematical expressions. Combining it with other large datasets like CROHME23 would increase the variety of samples even further, both in terms of writing style and number of expressions, likely improving the performance of a model.

Bounding box information for synthetic samples together with individual symbols are provided to enable experimentation with synthetic ink generation. Synthetic samples were generated through the straightforward process of pasting inks of individual symbols (symbols) exactly where bounding boxes were located. This gives synthetic samples a very regular structure, see Figure 2. It is possible to improve this process by modifying the location, size or orientation of bounding boxes prior to generating the synthetic inks. This would soften \LaTeX 's rigid structure and make synthetic data closer to human handwriting. Another application of these

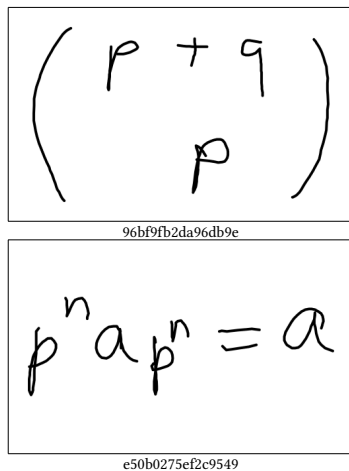


Figure 10:

Top: $\binom{p+q}{p}$ or 2-element matrix? **Bottom:** $p^n a p^n = a$ or $p^n a p^n = a$?

bounding boxes would be to bootstrap a recognizer that would also return character segmentation information. This kind of output is critical for some UI features - for example, editing a handwritten expression.

MathWriting can also be improved by varying the label normalization. Changing it can have different benefits depending on the application, as mentioned above. We provide the source \LaTeX string for that reason. Another possible improvement in recognition can come from additional contextual information, for instance the scientific field [7] that can be added post-hoc. Combining recognizers with a language model [5] trained on a large set of mathematical expressions would be a step in a similar direction.

7 Limitations

A single sample in MathWriting dataset has one handwritten \LaTeX formula, see Figure 3. As a result, models that are trained on this dataset would perform poorly on complete handwritten documents, such as the IAMonDo dataset [16]. Also, as the dataset contains only \LaTeX expressions, it is unlikely that models trained on it will accurately recognize handwritten text in English or other languages. As shown in Figure 4, some \LaTeX tokens are way more frequent than others. Some infrequent tokens like $\backslash ni$ could be hard to recognise.

8 Conclusion

We introduced MathWriting, the largest dataset of online handwritten mathematical expressions to date, together with the experimental results of three different types of models. We hope this dataset will help advance research in both online and offline mathematical expression recognition. Additionally, we invite data practitioners to build on the dataset. We intentionally chose a file format for MathWriting close to the one used by CROHME to facilitate their combined use. We also provided original or intermediate representations (raw \LaTeX strings, bounding boxes) to enable experimentation with the data itself, and suggested a few directions (Section 6.3).

Acknowledgements

We warmly thank all the contributors without whom this dataset would not exist. We thank Henry Rowley and Thomas Deselaers for their contribution to organizing the data collection and supporting this effort for many years. We thank Ashish Myles and MH Johnson for related contributions that resulted in important data and model improvements. We thank Vojta Letal and Pedro Gonnet for their contributions to the CTC Transformer model as well as to synthetic data. We thank Peter Garst and Jonathan Baccash for their contribution to the label normalizer. We thank Blagoj Mitrevski and Henry Rowley for their useful suggestions regarding the text of the paper. We thank our product counsels Janel Thamkul and Rachel Stigler for their legal advice.

References

- [1] [n. d.] Aida calculus math handwriting recognition dataset. <https://www.kaggle.com/datasets/aidapearson/ocr-data>. ()
- [2] Cem Anil et al. 2022. Exploring length generalization in large language models. (2022). arXiv: 2207.04901 [cs. CL].
- [3] Apple. 2023. Solve math with math notes in calculator on iphone. <https://support.apple.com/guide/iphone/solve-math-with-math-notes-iph46efa613a/ios>.
- [4] Lucas Beyer et al. 2024. Paligemma: a versatile 3b vlm for transfer. (2024). <https://arxiv.org/abs/2407.07726> arXiv: 2407.07726 [cs. CV].
- [5] Victor Carbune et al. 2020. Fast multi-language lstm-based online handwriting recognition. (2020). arXiv: 1902.10525 [cs. CL].
- [6] Xi Chen et al. 2023. Pali: a jointly-scaled multilingual language-image model. (2023). arXiv: 2209.06794 [cs. CV].
- [7] Mark Collier, Rodolphe Jenatton, Efi Kokiopoulou, and Jesse Berent. 2022. Transfer and marginalize: explaining away label noise with privileged information. In *International Conference on Machine Learning*. <https://api.semanticscholar.org/CorpusID:246996751>.
- [8] Yuntian Deng, Anssi Kanervisto, Jeffrey Ling, and Alexander M. Rush. 2017. Image-to-markup generation with coarse-to-fine attention. (2017). arXiv: 1609.04938 [cs. CV].
- [9] [n. d.] Detexify data. <https://github.com/kirel/detexify-data>. ()
- [10] Anastasiia Fadeeva, Philippe Schlattner, Andrii Maksai, Mark Collier, Efi Kokiopoulou, Jesse Berent, and Claudiu Musat. 2024. Representing online handwriting for recognition in large vision-language models. (2024). arXiv: 2402.15307 [cs. CV].
- [11] Gemma Team. 2024. Gemma: open models based on gemini research and technology. (2024). <https://arxiv.org/abs/2403.08295> arXiv: 2403.08295 [cs. CL].
- [12] Google Cloud. 2023. Detect handwriting in image. https://cloud.google.com/document-ai/docs/enterprise-document-ocr#ocr_add_ons.
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning (ICML '06)*. Association for Computing Machinery, Pittsburgh, Pennsylvania, USA, 369–376. ISBN: 1595933832. doi:10.1145/1143844.1143891.
- [14] Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and Jürgen Schmidhuber. 2009. A novel connectionist system for unconstrained handwriting recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31, 855–868. <https://api.semanticscholar.org/CorpusID:14635907>.
- [15] [SW]. How to improve accuracy (Model on TrOCR) version 1.0, 2022. URL: <https://github.com/huggingface/transformers/issues/16458>.
- [16] Emanuel Indermühle, Marcus Liwicki, and Horst Bunke. 2010. Iamondo-database: an online handwritten document database with non-uniform contents. In (June 2010), 97–104. doi:10.1145/1815330.1815343.
- [17] Dmitrijs Kass and Ekta Vats. 2022. Attentionhtr: handwritten text recognition based on attention encoder-decoder networks. (2022). arXiv: 2201.09390 [cs. CV].
- [18] Daniel Keyzers et al. 2020. Measuring compositional generalization: a comprehensive method on realistic data. (2020). arXiv: 1912.09713 [cs. LG].
- [19] Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. 2022. Trocr: transformer-based optical character recognition with pre-trained models. (2022). <https://arxiv.org/abs/2109.10282> arXiv: 2109.10282 [cs. CL].
- [20] M. Liwicki and H. Bunke. 2005. IAM-OnDB-an on-line english sentence database acquired from handwritten text on a whiteboard. In *ICDAR'05*. IEEE.
- [21] Mahshad Mahdavi, Richard Zanibbi, Harold Mouchère, Christian Viard-Gaudin, and Utpal Garain. 2019. Icdar 2019 crohme + tfd: competition on recognition

- of handwritten mathematical expressions and typeset formula detection. *2019 International Conference on Document Analysis and Recognition (ICDAR)*, 1533–1538. <https://api.semanticscholar.org/CorpusID:201611288>.
- [22] Johannes Michael, Roger Labahn, Tobias Grüning, and Jochen Zöllner. 2019. Evaluating sequence-to-sequence models for handwritten text recognition. (2019). arXiv: 1903.07377 [cs. CV].
- [23] Harold Mouchère, Richard Zanibbi, Utpal Garain, and Christian Viard-Gaudin. 2016. Advancing the state of the art for handwritten math recognition: the crome competitions, 2011–2014. *Int. J. Doc. Anal. Recognit.*, 19, 2, (June 2016), 173–189. doi:10.1007/s10032-016-0263-5.
- [24] Arun Narayanan, Rohit Prabhavalkar, Chung-Cheng Chiu, David Rybach, Tara N. Sainath, and Trevor Strohman. 2019. Recognizing long-form speech using streaming end-to-end models. (2019). arXiv: 1910.11455 [eess. AS].
- [25] Masato Neishi and Naoki Yoshinaga. 2019. On the relation between position information and sentence length in neural machine translation. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*. Association for Computational Linguistics, Hong Kong, China, (Nov. 2019), 328–338. doi:10.18653/v1/K19-1031.
- [26] George Retsinas, Giorgos Sfikas, Basilis Gatos, and Christophoros Nikou. 2024. Best practices for a handwritten text recognition system. (2024). arXiv: 2404.11339 [cs. CV].
- [27] [SW] Niels Rogge, "Transformers Tutorials" version 1.0, Sept. 2020. doi:10.5281/zenodo.1234, URL: <https://github.com/NielsRogge/Transformers-Tutorials>.
- [28] Shahrukh Khan. 2021. Im2latex-100k. <https://www.kaggle.com/datasets/shahrukhkhan/im2latex100k>.
- [29] May Mowaffaq AL-Taee, Sonia Ben Hassen Neji, and Mondher Frikha. 2020. Handwritten recognition: a survey. *2020 IEEE 4th International Conference on Image Processing, Applications and Systems (IPAS)*, 199–205. <https://api.semanticscholar.org/CorpusID:231823977>.
- [30] 2023. *Dss: synthesizing long digital ink using data augmentation, style encoding and split generation. Document Analysis and Recognition - ICDAR 2023*. Springer Nature Switzerland, 217–235. ISBN: 9783031416859. doi:10.1007/978-3-031-41685-9_14.
- [31] Dusan Varis and Ondřej Bojar. 2021. Sequence length is a domain: length-based overfitting in transformer models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. doi:10.18653/v1/2021.emnlp-main.650.
- [32] Dusan Varis and Ondřej Bojar. 2021. Sequence length is a domain: length-based overfitting in transformer models. In (Jan. 2021), 8246–8257. doi:10.18653/v1/2021.emnlp-main.650.
- [33] 2023. *Icdar 2023 crome: competition on recognition of handwritten mathematical expressions*. (Aug. 2023), 553–565. ISBN: 978-3-031-41678-1. doi:10.1007/978-3-031-41679-8_33.
- [34] Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai. 2022. Syntax-aware network for handwritten mathematical expression recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4553–4562.

Appendix

A \LaTeX template for label rendering

All the packages and definitions that are required to compile all the normalized and raw labels:

```
\usepackage{amsmath}
\usepackage{amsfonts}
\usepackage{amssymb}
\newcommand{\R}{\mathbb{R}}
\newcommand{\C}{\mathbb{C}}
\newcommand{\Q}{\mathbb{Q}}
\newcommand{\Z}{\mathbb{Z}}
\newcommand{\N}{\mathbb{N}}
```

B Acquisition of \LaTeX Expressions

The labels we publish mostly come from Wikipedia (95% of all samples have labels from Wikipedia). A small part were generated, to cover deeply nested fractions, number-heavy expressions, and isolated letters with nested superscripts and subscripts, which are rare in Wikipedia.

The extraction process from Wikipedia followed these steps:

- download an XML Wikipedia dump which provides Wikipedia's raw textual content. `enwiki-20231101-pages-articles.xml` was used for synthetic samples, older dumps for human-written ones
- extract all \LaTeX expressions from that file. This gives the list of all mathematical expressions in \LaTeX notation from Wikipedia
- keep those which could be compiled using the packages listed in Appendix A. Wikipedia contains a significant number of expressions that are not accepted by the \LaTeX compiler, because of syntax errors or other reasons
- keep only those which can be processed by our normalizer which only supports a subset of all \LaTeX commands and structures

For expressions used for synthesis, the following extra steps were performed:

- keep only the expressions whose normalized form contains more than a single \LaTeX token. Example: `\alpha` is rejected but `\alpha^{2}` is kept. This step is useful to eliminate trivial expressions that wouldn't add any useful information
- de-duplicate expressions based on their normalized form. e.g. `\frac{1}{2}` and `\frac{1}{2}` normalize to the same thing, we kept only one of them in raw form
- restrict the list of expressions to the same set of tokens used in the train split: if the normalized form of an expression contained at least one token that was not also present somewhere in train, it was discarded.

C Postprocessing of MathWriting dataset

We applied no postprocessing to the collected inks other than dropping entirely those that were completely unreadable or had stray marks. Inks are provided in their original form, as they were recorded with the collection app. What we *did not do* was to discard samples that were very hard to read or ambiguous, because we believe this type of sample to be essential in training a high-quality model.

Some cleanup was performed on the labels (ground truths). The goal was to make the dataset better suited to training an ML model, and eliminate unavoidable issues that occurred during the collection. After training some initial models, we manually reviewed samples for which they performed poorly. This helped identify a lot of unusable inks (near-blank, lots of stray strokes, scribbles, etc.), and a lot of ink/label discrepancies. A fairly common occurrence was a contributor forgetting to copy part of the prompted expression. We adjusted the label to what was actually written unless the ink contained a partially-drawn symbol, in which case we discarded the sample entirely. In this process we eliminated or fixed around 20k samples.

The most important postprocessing step was to normalize the labels: there are many different ways to write a mathematical expression in \LaTeX format that will render to images that are equivalent in handwritten form. We applied a series of transformations to eliminate as many variations as possible while retaining the same semantic. This greatly improved the performance of models and

made their evaluation more precise. We publish both the normalized and raw (unnormalized) labels, to enable people to experiment with other normalization procedures.

This normalization is similar to what [8] did, but pushed further because of the specifics of handwritten MEs. See Section 3.4 for more detail.

D Dataset split

The `valid` and `test` splits are the result of multiple operations performed between 2016 and 2019. The first split operation, performed on the data available in 2016, was based on the contributor id: any given contributor's samples would not appear in more than one split (either `train`, `valid`, `test`). This is common practice for handwriting recognition systems, to test how the recognizer performs on unseen handwriting styles.

Experiments then showed that a more important factor than the handwriting style was whether the *label* had already been seen during training. Subsequent data collection campaigns focused on increasing label variety, and new samples were added to `valid` and `test`, this time split by label: a given normalized mathematical expression would not appear in more than one split.

E Label Normalization

E.1 Syntactic Variations

There are several ways to change a \LaTeX string without changing the rendered output significantly. The normalization we implemented does the following:

- all unnecessary space is dropped
- all command arguments are consistently put in curly braces
- superscripts and subscripts are put in curly braces and their order is normalized. e.g. a^{2_1} becomes a_{1}^{2} .
- redundant braces are dropped
- infix commands are replaced by their prefix versions. e.g. $\frac{}{}{\over}$ is replaced by $\frac{}{}{\over}$
- a lot of synonyms are collapsed. e.g. \leq and \leq , \rightarrow and \rightarrow , etc. Some of the synonyms are only synonyms in handwriting. For example \star and $*$ are different in print (5-prong and 6-prong stars), but the difference was not expressed in handwriting by our contributors.
- functions commands like \sin are replaced by the sequence of letters of the function name (e.g. \sin is replaced by \sin). This reduces the output vocabulary, and eliminates a source of confusion because we found that \LaTeX expressions from Wikipedia come with a mix of function commands and sequences of letters.
- expansion of abbreviations. e.g. \cdots , \ldots , etc. have been replaced by the corresponding sequence of characters.
- matrix environments are normalized to use only the 'matrix' environment surrounded by the proper delimiters like brackets or parentheses.
- $\binom{}{}{\binom}$ is turned into a 2-element column matrix. Expressions from Wikipedia did not use those consistently, so we made the choice to normalize \binom away.

Device type	Ink
Google PixelBook	51k
Google Nexus 5X	28k
Coolpad Mega 2.5D	14k
OnePlus One	13k
Google Nexus 5	11k
Google Nexus 6	11k
Google Nexus 6P	11k
Coolpad Mega 3	8k
LG Optimus L9	8k
Galaxy Grand Duos	7k
Google Pixel XL	6k
Samsung Galaxy S7	5k

Table 8: Top-12 devices used, with the number of samples obtained from each device. The bias towards Google devices simply reflects the conditions in which inks were collected.

E.2 Differences Between Print And Handwriting

The following characteristics can not be represented in handwriting and have been normalized away:

- color
- accurate spacing: e.g. \sim , $\backslash\text{quad}$.
- font style and size: e.g. $\backslash\text{mathrm}$, $\backslash\text{mathit}$, $\backslash\text{mathbf}$, $\backslash\text{scriptstyle}$.

There are others that can be represented in handwriting, but that are not consistent enough in MathWriting to be preserved:

- font families: Fraktur, Calligraphic. In practice, only Blackboard ($\backslash\text{mathbb}$) has been written consistently enough by contributors that we were able to keep it: $\backslash\text{mathcal}$ and $\backslash\text{mathfrak}$ are dropped.
- some variations like \rightarrow and \rightarrow .
- some character variations. e.g. ϱ , ε
- size modifiers like \left , \right , \big . Similarly, variable-width diacritics like $\widehat{}$.

F Additional dataset statistics

In this section we show additional graphs that illustrate dataset statistics that are described in Section 4. The frequencies of normalized \LaTeX expressions are presented in Figure 12. Figure 13 illustrates the distribution of sampling rates within human-written data. Results of resampling points in time are presented in Figure ??.

G Variety of Writing Styles

In this section we provide additional examples of differences in the writing order of fractions – Figure 14. These examples show that MathWriting dataset contains a variety of writing styles.

H Sources of Noise

The result of any task performed by humans will contain mistakes, and MathWriting is no exception. We've done our best to remove most of the mistakes, but we know that some remain.

Sub/superscript are put in braces, \over is replaced by \frac

Raw: $\overline{hu^2} + \frac{1}{2} k_{ap} g_{zh}^2$

Normalized: $\overline{hu^{\{2\}}} + \frac{\{1\}}{\{2\}} k_{\{ap\}} g_{\{z\}} h^{\{2\}}$

Subscripts are put before superscripts, extra space is dropped

Raw: $\int^{a_{-a}} f(x) dx = 0$

Normalized: $\int_{-a}^a f(x) dx = 0$

Single quotes are replaced by a superscript

Raw: $f'(\overline{x})$

Normalized: $f^{\{\prime\}}(\overline{x})$

Text formatting commands like \rm are dropped

Raw: $\sim A_{\{0\}} = \frac{ND}{\{\sigma_{\{rm\ as\}} + \sigma_{\{rm\ es\}}\}} \sim$

Normalized: $A_{\{0\}} = \frac{ND}{\{\sigma_{\{as\}} + \sigma_{\{es\}}\}}$

Matrix environments with delimiters like bmatrix are replaced by matrix surrounded by delimiters
 Commands like \cos are replaced by the series of letters

Raw: $\begin{bmatrix} -\sin t \\ \cos t \end{bmatrix}$

Normalized: $[\begin{matrix} -sint \\ cost \end{matrix}]$

Delimiter size modifiers like \big are dropped

Raw: $\big(\frac{a}{N}\big)$

Normalized: $(\frac{a}{N})$

Figure 11: Examples of expression normalization. See Section 3.4 for details.

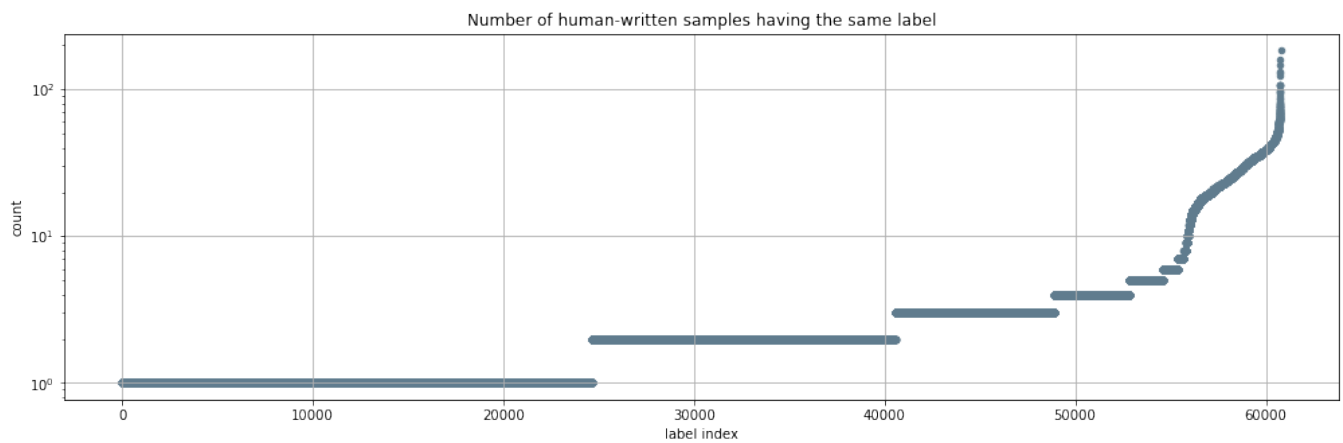


Figure 12: Counts of inks corresponding to the same normalized expression, ordered by increasing count. Each position on the horizontal corresponds to a unique normalized expression. Almost 5k unique expressions have been written 10 times or more by contributors.

Stray strokes. These do not carry any meaning and should be ignored by any recognizer. Since they also appear in real applications, there could be some benefit in having some in the dataset

to teach the model about them. That said, it being usually easier to add noise rather than to remove it, we made the choice of

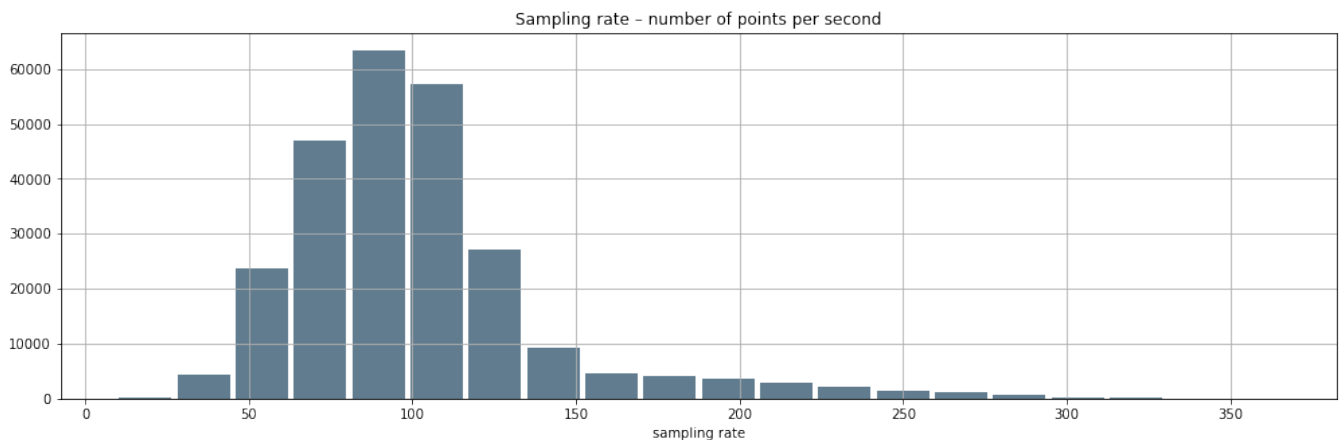


Figure 13: Histogram of sampling rates in human-written data of MathWriting dataset.

discarding as many inks containing stray strokes as possible. Not all inks with stray strokes have been found and removed though (e.g. `train/9e64be65cb874902.inkml` that was discovered post-publication). The fraction of inks containing stray strokes is significantly lower than 1%, and should not be an issue for training a model.

Incorrect ground truth. Contributors did not always copy the prompt perfectly, leading to a variety of differences. In most of the cases we spotted, we were able to fix the label to match what had actually been written. A short manual review once the dataset was in its final state showed the rate of incorrect ground truth to be between 1% and 2%. Most of the mistakes are very minor, usually a single token added, missing or incorrect. Errors here also come from ambiguities or misuse of the \LaTeX notation: expressions coming from Wikipedia contain some misuse like using `\Sigma` where `\sum` was more appropriate, `\triangle` instead of `\Delta`, `\triangledown` instead of `\nabla`, `\begin{matrix}` instead of `\binom`, and also some handwriting-specific ambiguities like `\dagger` vs `\top` vs `T`. There are also some instances where reference numbers or extra punctuation are included.

Aggressive normalization. While the above sources of noise are unavoidable, normalization is a postprocessing operation that can in theory be tweaked to perfection. In practice, it's a compromise between reducing accidental ambiguities (i.e. removing synonyms), and removing information. Examples: we made the choice of treating `\binom` as a synonym for a 2-element matrix. While it does improve recognition accuracy by making the problem easier, it also moves the burden of distinguishing between the two cases to downstream steps in the recognition pipeline. Similar things can be said about removing all commands that indicate that their content is text instead of math (e.g. `\mbox`), dropping size modifiers, rewriting function commands (e.g. `\sin`, `\cos`), etc. Using a different normalization could prove beneficial depending on the context the recognizer is used in practice. However, for the purpose of a benchmark any reasonable compromise is adequate.

I Tokens

Using the above code to compute tokens, the set of all samples in the dataset (human-written, synthetic, from all splits) contain the following after normalization:

- Syntactic tokens: `_` `^` `&` `\` `space`
- Latin letters and numbers: `a-z` `A-Z` `0-9`
- Blackboard capital letters `\mathbb{A}`-`\mathbb{Z}` `\mathbb`
- Latin punctuation and symbols: `,` `;` `!` `?` `.` `()` `[]` `\{` `\}` `*` `/` `+` `-` `_` `&` `\#` `\%` `\backslash` `slash`
- Greek letters: `\alpha` `\beta` `\delta` `\Delta` `\epsilon` `\eta` `\chi` `\gamma` `\Gamma` `\iota` `\kappa` `\lambda` `\Lambda` `\nu` `\mu` `\omega` `\Omega` `\phi` `\Phi` `\pi` `\Pi` `\psi` `\Psi` `\rho` `\sigma` `\Sigma` `\tau` `\theta` `\Theta` `\upsilon` `\Upsilon` `\varphi` `\varpi` `\varsigma` `\vartheta` `\xi` `\Xi` `\zeta`
- Mathematical constructs: `\frac` `\sqrt` `\prod` `\sum` `\iint` `\int` `\oint`
- Diacritics and modifiers - Note the absence of the single-quote character, which is normalized to `^{\prime}`:
`\hat` `\tilde` `\vec` `\overline` `\underline` `\prime` `\dot` `\not`
- Matrix environment: `\begin{matrix}` `\end{matrix}`
- Delimiters: `\angle` `\rangle` `\lceil` `\rceil` `\lfloor` `\rfloor` `\|`
- Comparisons: `\ge` `\gg` `\le` `\ll` `<>`
- Equality, approximations: `=` `\approx` `\cong` `\equiv` `\ne` `\propto` `\sim` `\simeq`
- Set theory: `\in` `\ni` `\notin` `\sqsubseteq` `\subseteq` `\subset` `\supseteq` `\subsetneq` `\supset` `\supseteq` `\emptyset`
- Operators: `\times` `\bigcap` `\bigcirc` `\bigcup` `\bigoplus` `\bigvee` `\bigwedge` `\cap` `\cup` `\div` `\mp` `\odot` `\ominus` `\oplus` `\otimes` `\pm` `\vee` `\wedge`
- Arrows: `\hookrightarrow` `\leftarrow` `\leftrightarrow` `\Leftrightarrow` `\longrightarrow` `\mapsto` `\rightarrow` `\Rrightarrow` `\rightleftarpoons` `\iff`
- Dots: `\bullet` `\cdot` `\circ`
- Other symbols: `\aleph` `\angle` `\dagger` `\exists` `\forall` `\forall` `\bar` `\infty` `\nabla` `\nabla` `\neg` `\partial` `\perp` `\top` `\triangle` `\triangleleft` `\triangleleft` `\triangleleft` `\vdash` `\Vdash` `\vdots`

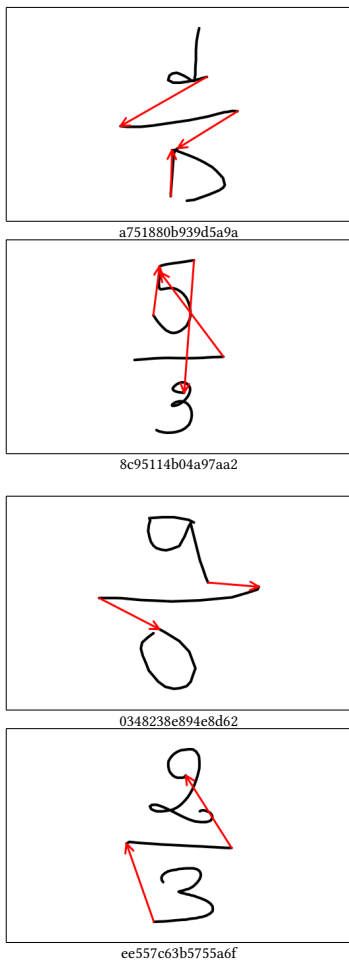
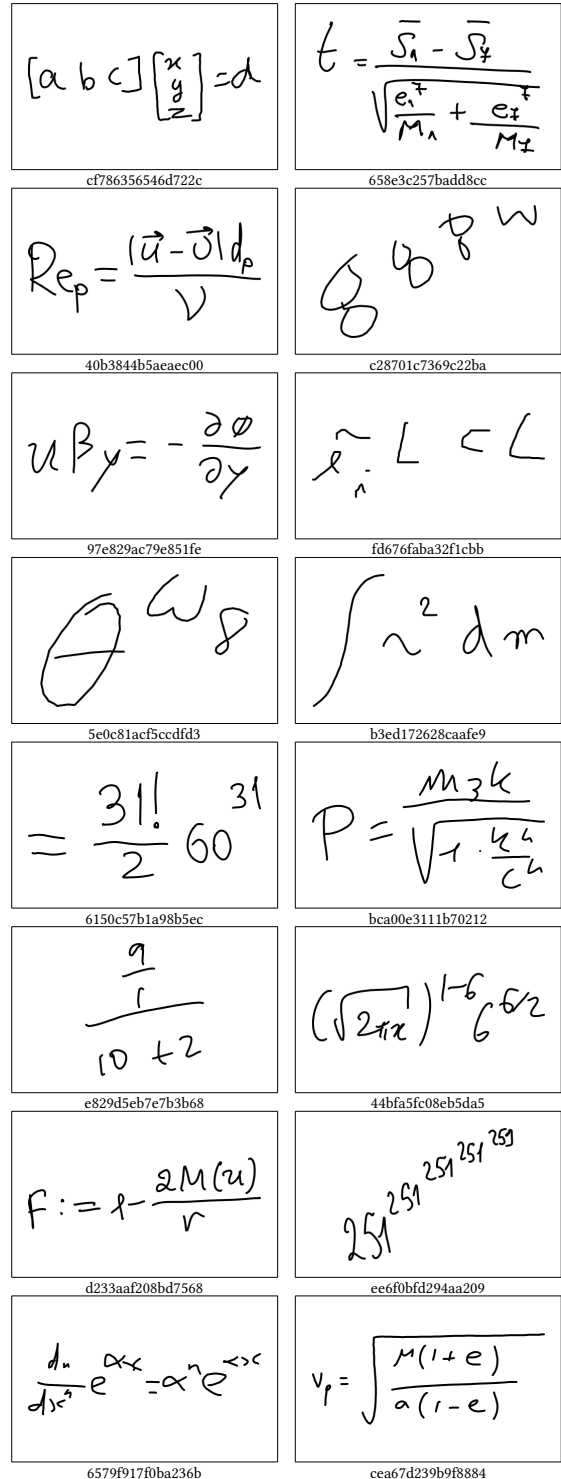


Figure 14: Examples of various writing orders found in the training set. Red arrows show the movement of the pen between strokes. The top one: most common writing order (top-down, fraction bar drawn left-to-right), second from the top: fraction bar written first, third from the top: fraction bar drawn right-to-left, forth from the top: fraction written bottom-up.

J Examples of inks

This section shows a few examples of rendered inks, so that the reader can get a feel for the kind of data that is in MathWriting. All samples are from the training set. They have been manually picked to show a variety of sizes, characters and structures.

J.1 Human-Written Samples



$$\frac{\frac{\sqrt{6}}{10} - 456}{\frac{\sqrt{6}}{61} / 7}$$

bb4bca53d0e6336d

$$\binom{n}{k}$$

b14bca3fc2d2819a

$$L = \sqrt{1 + [f'(x)]^2}$$

2409d2feaa79b9d7

$$\frac{\partial u}{\partial \bar{z}} - i z \frac{\partial u}{\partial z} = \varphi'(z)$$

51486f88b789d6d

$$d \propto \sqrt{2 \cdot k \cdot R \cdot h}$$

02229a0c174d8dbe

$$\tilde{C}_7$$

478e10a15203fa3a

$$W \begin{pmatrix} \psi \\ z \end{pmatrix}$$

ccb15825579a096b

$$\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$$

7a4b95285de0caf0

$$A = (A \cup B) - B$$

ac25b4d053596ded

$$\int_{-\infty}^{\infty} x^n$$

7d597c52bf8bdd1e

$$\left(\frac{246}{1297}\right)^6 \cdot 363^4 + 424$$

acc5f6620fad1ce8

$$t \equiv \sqrt{1 - 1/\beta^2}$$

88c3551d373c72e5

$$\theta_b^m = \sum_{a,b} \theta^a$$

5009d32d32f80324

$$\forall b_1, \dots, b_n \notin T(s)$$

068de3aad90c403c

$$\sum_{n=1}^{\infty} \frac{1}{|z_n|^{p+1}}$$

a3cf115524f0c55b

$$\bar{E} = 0, \bar{\xi}_i \bar{\xi}_j = \frac{1}{3} \bar{\xi}^i \bar{\xi}^j$$

355f5df56a16913a

$$\theta = n \times 137.508^\circ$$

d9b2ce7aa3495888

$$\frac{\partial(\bar{r}_w \cdot r_h)}{\partial x_h} = 0$$

adceb80fdad9f9fe

$$264-175-365-445$$

25892f7caecac8c36

$$(P_j + iQ_j) \Delta^{-1/2}$$

41e1261a951c6f33

$$F_{out} \alpha_{out} = \eta F_{in} \alpha_{in}$$

02a7f7f172671fb4

$$\deg_P f = \deg_P g$$

0d848d4b170d36b9

$$2 m \varepsilon \frac{d^3 p d^3 r}{h^3}$$

2adc4f10d42b641c

$$v_e = c \sqrt{2\eta - \eta^2}$$

408f904038dcbba0

J.2 Synthetic Samples: Expressions from Wikipedia

$$\infty, \beta, 1$$

133829b5a10b783f

$$\gamma = 1/(\sigma \mu)$$

11623165e9bab0e4

$$L_2 \otimes I_5$$

5ca38d17bf2bea0a

$$u(x_1, \dots, x_s) \cdot \sum_i^s c_i$$

089650cc894c024b

$$Q_{mn}^* = \int v_m^* u_n dt$$

8e88b75cb5f03bf4

$$\frac{1}{(4\pi t)^{3/2}} e^{-\frac{|x|^2}{4t}}$$

5c1573b41e762307

$$E(r) \equiv 1 - \int d^d k \epsilon(r, k)$$

68e9560a27a093c8

$$\begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix}$$

daab3bae071f8bc0

$$|1\rangle_A |1\rangle_B$$

77602abaea39b774

$$Z_q$$

eb817bcbfd11df18

$$n^3 < N_1 N_2 N_3$$

3bd062b33ea5db6f

$$\frac{\partial}{\partial t} dr_s = \frac{1}{2} g^{ij} v_{pj} dv_p$$

51ef5122de326151

$$c_v = \frac{nR}{\gamma - 1}$$

f2a613fc323df342

$$x y^2$$

cbd19abc03ba4098

$$\bigoplus_{n \in \mathbb{Z}} \text{Hom}_n(A, B)$$

1879aa5c882b445d

$$\begin{pmatrix} f_1(x, y, y', y'', \dots, y^{(n)}) \\ f_2(x, y, y', y'', \dots, y^{(n)}) \\ \vdots \\ f_m(x, y, y', y'', \dots, y^{(n)}) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

094c7e52a3f0934d

$$ds^2 = (R \cos \phi)^2 d\lambda^2 + R^2 d\phi^2$$

0772aeaac09d3415

$$r = \begin{pmatrix} \frac{\partial A}{\partial x} & \frac{\partial A}{\partial y} \end{pmatrix} + \frac{1}{c\pi v} e^{i\pi k} (v\gamma) \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}, \quad \pi = \frac{e}{\pi} \nabla(\ln r)$$

20b4ebf292cfa8d1

$$f(x_1, x_2, x_3, \dots, x_n; u(x_1, x_2, x_3, \dots, x_n)) \delta'(u), D^2(u), D^3(u), \dots, D^n(u) = 0$$

96613ae167f35f8a

$$\frac{2\sqrt{3} - 1\sqrt{5} + (2 - \sqrt{3} - 1)(2 + \sqrt{6}) - (2\sqrt{2 + \sqrt{3}} + \sqrt{10 - 2\sqrt{6}})}{4}$$

4b1f5165e3698343

$$E = -\nabla\phi - \frac{dA}{dt} - \nabla \frac{\partial \mathcal{U}}{\partial t} = -\nabla \left(p + \frac{\partial \mathcal{U}}{\partial t} \right) - \frac{\partial A}{\partial t}$$

1cd654228d7ca6bb

$$c_{solid, p} = \sqrt{\frac{k + \frac{4}{3}G}{p}} = \sqrt{\frac{E(1-v)}{p(1+v)(1-2v)}}$$

6f86778996d5f514

$$\int |\Psi|^2 dx dy \approx A + 2Re \left[\frac{f(z)}{z} \int_{-\infty}^{\infty} e^{ikx^2/2\pi} dx \int_{-\infty}^{\infty} e^{iky^2/2\pi} dy \right]$$

fc00050933165b70

$$m x \frac{d^2 x}{dt^2} = m \frac{d[x(dx/dt)]}{dt} - m \left(\frac{dx}{dt} \right)^2$$

685bd5676bda74ce

$$A^{(1)} = \begin{pmatrix} 4 & 2 & 1 \\ 0 & 6 & 5 \\ 0 & 5 & 2 \end{pmatrix}, \quad P^{(1)} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

60553a301aaf84a3

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & \omega & \omega^2 & \omega^3 & \omega^4 \\ 1 & \omega^2 & \omega^4 & \omega & \omega^3 \\ 1 & \omega^3 & \omega & \omega^4 & \omega^2 \\ 1 & \omega^4 & \omega^3 & \omega^2 & \omega \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad \omega = e^{-\frac{j2\pi}{5}}$$

b79dfced7f5b5cb4

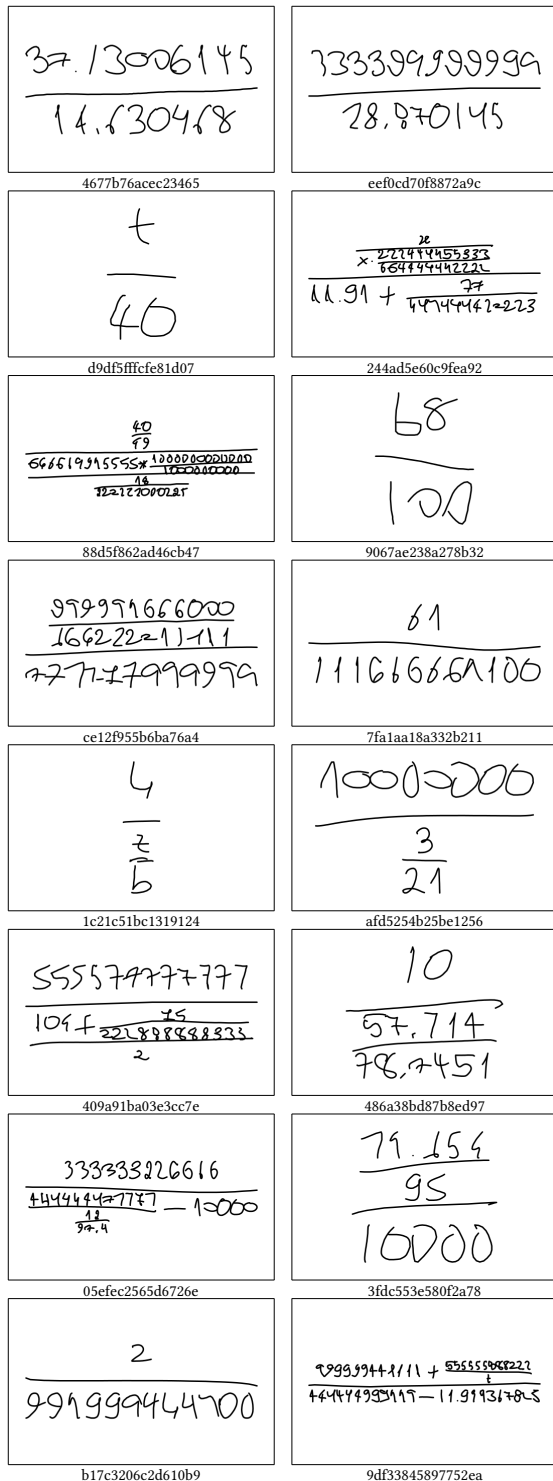
$$d\phi_n(p_1, p_2, \dots, p_n) = (2\pi)^{4n} \delta^4 \left(p - \sum_{i=1}^n p_i \right) \prod_{i=1}^n \frac{d^3 p_i}{(2\pi)^3 E_i}$$

b2aeaf7a0fd30ed6

$$y = -(g/2)t^2 + v t \sin \alpha + h \quad (2)$$

254dbc2b3843dcf8

J.2.1 Synthetic Samples: Generated Fractions.



K Dataset split

The valid and test splits are the result of multiple operations performed between 2016 and 2019. The first split operation, performed on the data available in 2016, was based on the contributor id: any given contributor's samples would not appear in more than one

split (either train, valid, test). This is common practice for handwriting recognition systems, to test how the recognizer performs on unseen handwriting styles.

Experiments then showed that a more important factor than the handwriting style was whether the *label* had already been seen during training. Subsequent data collection campaigns focused on increasing label variety, and new samples were added to valid and test, this time split by label: a given normalized mathematical expression would not appear in more than one split.

L Label Normalization

L.1 Syntactic Variations

There are several ways to change a \LaTeX string without changing the rendered output significantly. The normalization we implemented does the following:

- all unnecessary space is dropped
- all command arguments are consistently put in curly braces
- superscripts and subscripts are put in curly braces and their order is normalized. e.g. a^2_1 becomes a_{1}^{2} .
- redundant braces are dropped
- infix commands are replaced by their prefix versions. e.g. \over is replaced by \frac
- a lot of synonyms are collapsed. e.g. \leq and \leq , \rightarrow and \rightarrow , etc. Some of the synonyms are only synonyms in handwriting. For example \star and $*$ are different in print (5-prong and 6-prong stars), but the difference was not expressed in handwriting by our contributors.
- functions commands like \sin are replaced by the sequence of letters of the function name (e.g. \sin is replaced by \sin). This reduces the output vocabulary, and eliminates a source of confusion because we found that \LaTeX expressions from Wikipedia come with a mix of function commands and sequences of letters.
- expansion of abbreviations. e.g. \cdots , \ldots , etc. have been replaced by the corresponding sequence of characters.
- matrix environments are normalized to use only the 'matrix' environment surrounded by the proper delimiters like brackets or parentheses.
- \binom is turned into a 2-element column matrix. Expressions from Wikipedia did not use those consistently, so we made the choice to normalize \binom away.

L.2 Differences Between Print And Handwriting

The following characteristics can not be represented in handwriting and have been normalized away:

- color
- accurate spacing: e.g. \sim , \quad .
- font style and size: e.g. \mathrm , \mathit , \mathbf , \scriptstyle .

There are others that can be represented in handwriting, but that are not consistent enough in MathWriting to be preserved:

- font families: Fraktur, Calligraphic. In practice, only Blackboard (\mathbb) has been written consistently enough by contributors that we were able to keep it: \mathcal and \mathfrak are dropped.

Sub/superscript are put in braces, \over is replaced by \frac	
Raw:	<code>\overline{hu^2}+{1 \over 2}{k_{ap}g_zh^2}</code>
Normalized:	<code>\overline{hu^{2}}+\frac{1}{2}k_{ap}g_{z}h^{2}</code>

Subscripts are put before superscripts, extra space is dropped	
Raw:	<code>\int^a_{-a}f(x) dx=0</code>
Normalized:	<code>\int_{-a}^af(x)dx=0</code>

Single quotes are replaced by a superscript	
Raw:	<code>f'(\overline{x})</code>
Normalized:	<code>f^{\prime}(\overline{x})</code>

Text formatting commands like \rm are dropped	
Raw:	<code>\sim A_{\emptyset}=\frac{ND}{\sigma_{\rm as}+\sigma_{\rm es}}\sim</code>
Normalized:	<code>A_{\emptyset}=\frac{ND}{\sigma_{as}+\sigma_{es}}</code>

Matrix environments with delimiters like bmatrix are replaced by matrix surrounded by delimiters	
Commands like \cos are replaced by the series of letters	
Raw:	<code>\begin{bmatrix} -\sin t \ \ \ \cos t \end{bmatrix}</code>
Normalized:	<code>[\begin{matrix}-sint\ \ \ cost\end{matrix}]</code>

Delimiter size modifiers like \big are dropped	
Raw:	<code>\big(\tfrac{a}{N}\big)</code>
Normalized:	<code>(\frac{a}{N})</code>

Figure 15: Examples of expression normalization. See Section 3.4 for details.

- some variations like `\rightarrow` \rightarrow and `\longrightarrow` \longrightarrow .
- some character variations. e.g. `\varrho`, `\varepsilon`
- size modifiers like `\left`, `\right`, `\big`. Similarly, variable-width diacritics like `\widehat`.

`s = s[len(tokens[-1]):]`

`return tokens`

Received 23 February 2025

M Tokenization Code

Python code used in this work to tokenize \LaTeX mathematical expressions.

```
%\begin{verbatim}
import re

_COMMAND_RE = re.compile(
    r'\\(\mathbb{[a-zA-Z]}|begin{[a-z]+}|end{[a-z]+}|operatorname\*[a-zA-Z]+|\.)')

def tokenize_expression(s: str) -> list[str]:
    tokens = []
    while s:
        if s[0] == '\\':
            tokens.append(_COMMAND_RE.match(s).group(0))
        else:
            tokens.append(s[0])
```