

---

# MathWriting: A Dataset For Handwritten Mathematical Expression Recognition

---

Philippe Gervais\*  
pgervais@acm.org

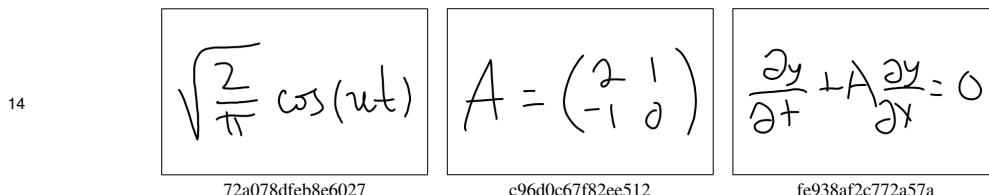
Asya Fadeeva  
Google  
fadeich@google.com

Andrii Maksai  
Google  
amaksai@google.com

## Abstract

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

## 13 1 Introduction



15 Three examples of HME from MathWriting. More examples can be found in Appendix K. Each ink  
16 is accompanied by a unique identifier that matches a corresponding filename in the dataset.

### 17 MathWriting dataset (2.9 GB):

18 [https://storage.googleapis.com/mathwriting\\_data/mathwriting-2024.tgz](https://storage.googleapis.com/mathwriting_data/mathwriting-2024.tgz)

### 19 Associated code:

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

21 Online *text* recognition models have improved a lot over the past years, because of improvements  
22 in model structure [1, 2, 3] and also because of an increase in the amount of training data [4, 5, 6].  
23 Mathematical expression (ME) recognition is a challenging task that has received less attention  
24 than regular recognition of words and characters [7]. ME recognition is different from regular text  
25 recognition in a number of interesting ways which can prevent improvements from transferring  
26 from one to the other. Though MEs share with text most of their symbols, they follow a more rigid  
27 structure which is also two-dimensional, see Figure 1. Where text can be treated to some extent as a  
28 one-dimensional problem amenable to sequence modeling, MEs cannot because the relative position

---

\*Work performed while employed at Google

29 of symbols in space is meaningful. It is also different from symbol segmentation or object detection  
30 because the output of a recognizer has to contain the relationship between symbols, serialized in  
31 some form ( $\LaTeX$ , a graph, InkML, etc.). Similarly to the case of text, *handwritten* MEs (**HME**) are  
32 more difficult to recognize than *printed* ones as they are more ambiguous and less training data is  
33 available.

34 Handwritten data is costly to obtain as it must be written by hand, which is compounded in the case  
35 of online representation (**ink**) by the necessity to use dedicated hardware (touchscreen, digital pen,  
36 etc.). By publishing the MathWriting dataset, we hope to alleviate some of the needs for data for  
37 research purposes. Samples include a large number of human-written inks, as well as synthetic ones.  
38 MathWriting can readily be used with other online datasets like CROHME [8] or Detexify [9] - we  
39 publish the data in InkML format to facilitate this. It can also be used for offline ME recognition  
40 simply by rasterizing the inks, using code provided on the Github page<sup>2</sup>.

41 MathWriting is the largest set of online HME published so far - both human-written and synthetic.  
42 It significantly expands the set of symbols covered by CROHME [8], enabling more sophisticated  
43 recognition capabilities. Since inks can be rasterized, MathWriting can also be seen as larger  
44 than existing offline HME datasets [10, 11, 12]. For these reasons we introduce a new benchmark,  
45 applicable to both online and offline ME recognition.

46 This work’s main contributions are:

- 47 • a large dataset of Handwritten Mathematical Expressions under the Creative Commons  
48 Attribution-NonCommercial-ShareAlike 4.0 International<sup>3</sup>.
- 49 •  $\LaTeX$  ground truth expressions in normalized form to simplify training and to make evalua-  
50 tion more robust.
- 51 • Evaluation of different models like CTC Transformer and PaLI on the dataset to show what  
52 recognition quality could be achieved with the provided data.

53 The paper focuses on the high-level description of the dataset: creation process, postprocessing,  
54 train/test split, ground truth normalization, statistics, and a general discussion of the dataset content  
55 to help practitioners understand what can and cannot be achieved with it. All the low-level technical  
56 information like file formats can be found in the `readme.md` file present at the root of the dataset  
57 archive linked above. We also provide code examples on Github<sup>2</sup>, to show how to read the various  
58 files, process and rasterize the inks, and tokenize the  $\LaTeX$  ground truth.

## 59 2 Dataset Creation

60 MathWriting dataset primarily consists of  $\LaTeX$  expressions from Wikipedia, more details about the  
61 acquisition of expressions are provided in Appendix B. These expressions were used for both ink  
62 collection from human contributors Section 2.1 as well as synthetic data generation Section 2.2. We  
63 did a very limited filtering of very noisy human-written examples (described in Appendix C).

### 64 2.1 Ink Collection

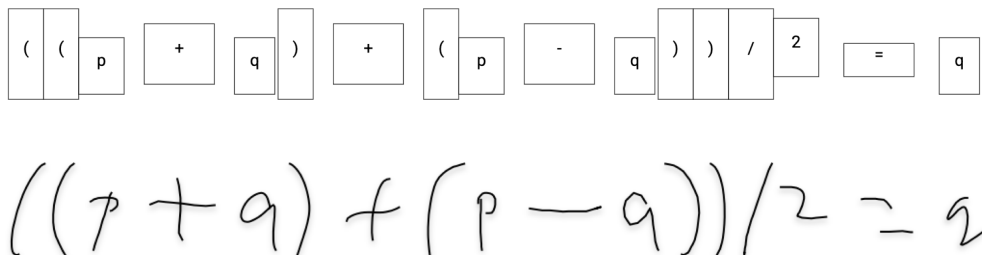
65 Inks were obtained from human contributors through an in-house Android app. Participants agreed  
66 to the standard Google terms of use and privacy policy. The task consisted in copying a rendered  
67 mathematical expression (prompt) shown on the device’s screen using either a digital pen or a finger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

<sup>2</sup><https://github.com/google-research/google-research/tree/master/mathwriting>

<sup>3</sup><https://creativecommons.org/licenses/by-nc-sa/4.0/>

77 **2.2 Synthetic Samples and Isolated Symbols**

78 We created synthetic samples in order to further increase the label diversity for training. This  
79 also enabled compensating for limitations of the human collection like the maximum length of the  
80 expressions, which were limited by the size of the screen they were written on. We used  $\LaTeX$   
81 expressions from Wikipedia that were not used in the data collection. The resulting synthetic  
82 data has a 90th percentile of expression length of 68 characters, compared to 51 in `train`. This is  
83 especially important as deep neural nets often fail to generalize to inputs longer than their training  
84 data [13, 14]. Using synthetic long inks together with the original human-written inks can help to  
85 eliminate that problem as shown in [15, 16]. The synthesis technique is as follows: starting from a  
86 raw  $\LaTeX$  mathematical expression, we computed a DVI file using the  $\LaTeX$  compiler, from which  
87 we extracted bounding boxes. We then used those bounding boxes to place handwritten individual  
88 symbols, resulting in a complete expression. See Figure 1 for an example of extracted bounding  
89 boxes and the resulting synthetic example.



**Figure 1:** An example of a synthetic ink created from bounding boxes with label  $((p+q)+(p-q))/2=q$

90 Inks for individual symbols are all from the `symbols` split. They have been manually extracted  
91 from inks in `train`. For each symbol that we wanted to support, we manually selected strokes  
92 corresponding to it for 20-30 distinct occurrences in `train`, and used that information to generate a  
93 set of individual inks. Similar synthesis techniques have been used by [8] with inks, [10] and [12]  
94 with raster images.

95 A significant difference between synthetic and human-written inks is the stroke order. For synthetic  
96 inks, stroke order follows the order of the bounding boxes in the DVI file, which can be different  
97 from the usual order of writing for mathematical expressions. However, the writing order within a  
98 given symbol is consistent with human writing.

99 **2.3 Dataset Split**

100 MathWriting is composed of five different sets of samples, which we call 'splits': `train`, `valid`,  
101 `test`, `symbols`, and `synthetic`. The splits `train`, `valid` and `test` consist only of human-written  
102 examples. The split `symbols` is provided for synthetic data generation and is not used in training.  
103 The split of human-written samples between `train`, `valid` and `test` was partially done based on  
104 writers, partially based on labels. More details are provided in Appendix D. Experiments have shown  
105 that a more important factor than the handwriting style was whether the *label* had already been seen  
106 during training. This fact is also supported by research in the area of compositional generalization  
107 [17]. In the published version, `valid` has a 55% (8.5k samples) intersection with `train` based on  
108 unique normalized labels, and `test` has an 8% intersection (647 samples). We chose to have a low  
109 intersection between `train` and `test` in order to correctly measure generalization of trained models  
110 to unseen labels.

111 **2.4 Label Normalization**

112 All samples in the dataset come with two labels: the  $\LaTeX$  expression that was used during the  
113 data collection (annotation `label` in the InkML files), and a normalized version of it meant for  
114 model training, which is free from a few sources of confusions for an ML model (annotation  
115 `normalizedLabel`). An example with original and normalized labels is provided in Figure 2. Label  
116 normalization covers three main categories (details are provided in Appendix E):

51b364eb9ba2185a

**Figure 2:** An example from the train split, with its labels:

Raw:  $f(x) = \frac{1}{e} \cdot \sum_{n=0}^{\infty} \frac{n^x}{n!}$

Normalized:  $f(x) = \frac{1}{e} \cdot \sum_{n=0}^{\infty} \frac{n^{\{x\}}}{n!}$

- 117 • variations used in print that can't be reproduced in handwriting - e.g. bold, italic - or that
- 118 haven't been reproduced consistently by contributors.
- 119 • non-uniqueness of the L<sup>A</sup>T<sub>E</sub>X syntax. e.g. `\frac{1}{2}` and `1\over 2` are equivalent.
- 120 • visual variations that can be reproduced in handwriting but can't reliably be inferred by a model.
- 121 This includes size modifiers like `\left`, `\right`.

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

### 124 2.4.1 Limitations of normalization

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

## 135 3 Dataset Statistics

136 In this section we describe the key characteristics of MathWriting and compare it to CROHME23  
 137 [8]. In Table 1 we provide the information about the volume of the dataset splits both in terms of  
 138 examples (inks) and unique labels.

**Table 1:** Statistics on different subsets of MathWriting dataset.

	train	synthetic	valid	test
# distinct inks	230k	396k	16k	8k
# distinct labels	53k	396k	8k	4k

### 139 3.1 Label Statistics

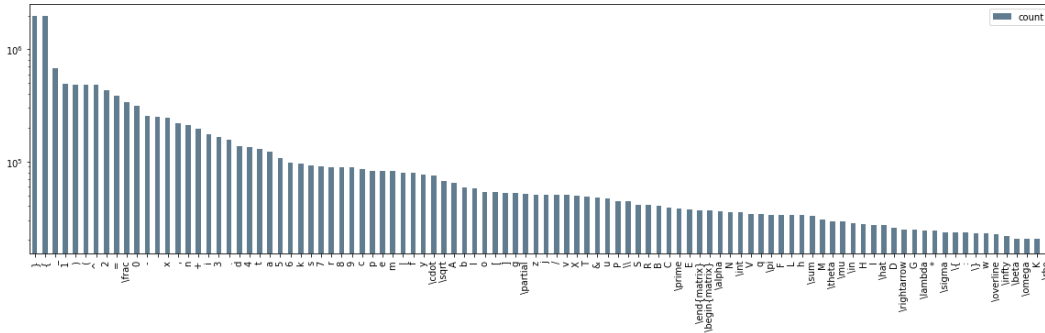
140 MathWriting contains 457k unique labels after normalization (see Section 2.4). From Table 1 we  
 141 see that most unique expressions are covered by the synthetic portion of the dataset. However, the  
 142 absolute number of unique expressions in human-written part is still high – 61k. This underlines  
 143 the importance of synthetic data as it allows models to see a much bigger variety of expressions. It  
 144 is important to note that the synthetic split has essentially no repeated expressions. On the other  
 145 hand, in real data multiple different writings of the same expression are quite common (see Figure 11  
 146 in Appendix F). This fact allows us to separately evaluate model's quality on expressions that were

147 observed during training and that those that hadn't. As seen in Table 2 the biggest intersection in  
 148 expressions is between `valid` and `train`. The minimal overlap between `test` and `train` splits is  
 149 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	-

150 The median length of expressions in characters is 26 which is comparable to one of the most popular  
 151 English recognition datasets IAMonDB [18] which has median of 29 characters. However, it is  
 152 important to note that  $\LaTeX$  expressions have tokens that span multiple characters like `\frac`. The  
 153 median length of expressions in tokens (provided in Appendix J) is 17, thus making training a model  
 154 on tokens rather than characters easier due to shorter target lengths [19, 20]. We want to emphasize  
 155 that MathWriting can be used with a different tokenization scheme and token vocabulary from what  
 156 we propose in Appendix J. In Figure 3 we show the number of occurrences for the most frequent  
 157 tokens. Tokens `{` and `}` are by far the most frequent as they are integral to the  $\LaTeX$  syntax.



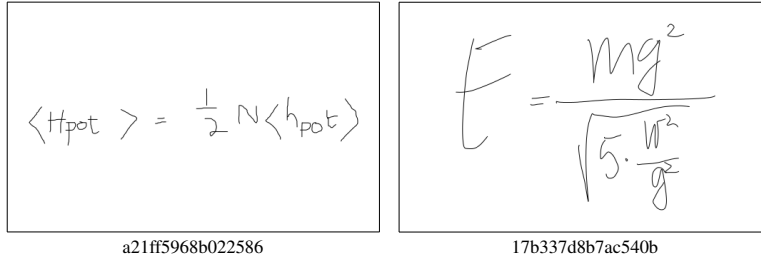
**Figure 3:** Histogram of the top-100 most frequent tokens in MathWriting.

### 158 3.2 Ink Statistics

159 Each ink in MathWriting dataset is a sequence of strokes  $I = [s_0, \dots, s_n]$ , each stroke  $s_i$  consisting  
 160 of points. A point is represented as a triplet  $(x, y, t)$  where  $x$  and  $y$  are coordinates on the screen and  $t$   
 161 is a timestamp. In Table 3 we provide statistics on number of strokes, points, and duration of writing.  
 162 It's important to note that as inks were collected on different devices, the absolute coordinate values  
 163 can vary a lot. In human-written data the time information  $t$  always starts from 0 but it is not always  
 164 the case in the `synthetic` split. Different samples often have different sampling rates (number of  
 165 points written in one second) due to the use of different devices (see Figure 4). More details in  
 166 Section 3.3. Consequently, the same ink written on two different devices can result in inks with a  
 167 different number of points. For human-written inks, the sampling rate is consistent between strokes,  
 168 but it is not the case for synthetic ones. In order to accommodate a model and make sequences shorter,  
 169 inks can be resampled in time (see example in Figure 13, Appendix F).

**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



**Figure 4:** Left: an ink with very low sampling rate (9.4 points per second)  
 Right: an ink with very high sampling rate (260 points per second)

**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  $\$$ .

	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

### 170 3.3 Devices Used

171 Around 150 distinct device types have been used by contributors. In most cases inks were written on  
 172 smartphones using a finger on a touchscreen. However, there are cases where tablets with styluses  
 173 were used. The main device used in this case is Google Pixelbook, which accounted for 51k inks total  
 174 (see Table 7, Appendix F). Out of all device types, 37 contributed more than 1000 inks. Note that  
 175 writing on a touchscreen with a finger or a stylus results in different low-level artifacts. All devices  
 176 were running the same Android application for ink collection, regardless of whether their operating  
 177 system was Android or ChromeOS.

### 178 3.4 Comparison With CROHME23

179 In this section we compare main dataset statistics of MathWriting and CROHME23 [8] as it is a  
 180 popular publicly available dataset for HME recognition. In terms of overall size, MathWriting has  
 181 nearly 3.9 times as many samples and 4.5 times as many distinct labels after normalization, see  
 182 Table 4. A significant number of labels can be found in both datasets (47k), but the majority is  
 183 dataset-specific. This suggests that combining both datasets during training could yield improved  
 184 HME recognition quality. MathWriting has more human-written inks than CROHME23 as seen in  
 185 Table 5, and contains a much larger variety of tokens. It has 254 distinct tokens including all Latin  
 186 capital letters and almost the entire Greek alphabet. It also contains matrices, which are not included  
 187 in CROHME23. Therefore, more scientific fields like quantum mechanics, differential calculus, and  
 188 linear algebra can be represented using MathWriting.

## 189 4 Experiments

### 190 4.1 Evaluation setup

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

- 193 • **evaluation samples:** the test split of MathWriting.
- 194 • **metric:** character error rate (CER) [21], where a "character" is a  $\LaTeX$  token as defined by  
 195 the code in Appendix I.

196 We provide a reference implementation of the evaluation metric at the Github page <sup>2</sup>. We propose the  
 197 use of CER as a metric to make results comparable to other recognition tasks like text recognition [22,  
 198 23], and the use of  $\LaTeX$  tokens instead of ASCII characters so that an error on a single non-latin  
 199 letter (e.g.  $\backslash\alpha$  recognized as a) counts as one instead of many.

## 200 4.2 Baseline Recognition Models

201 In Table 6 we provide results for different models. All models are trained exclusively on the  
 202 MathWriting dataset (`train` and `synthetic`), except for the OCR API that was trained on other  
 203 datasets as well. The following models represent different approaches to handwriting recognition –  
 204 offline [23], online [24] and mixed [25].

205 **OCR** This is a publicly available Document AI OCR API [26], which processes bitmap images. It  
 206 has been trained partly on samples from MathWriting. We sent inks rendered with black ink on a  
 207 white background and searched for optimal image size and stroke width to get the best evaluation  
 208 result from the model.

209 **CTC Transformer** This model is a transformer base with a Connectionist Temporal Classification  
 210 loss on top (CTC) [27]. It contains 11 transformer layers with an embedding size of 512. We used  
 211 swish activation function and dropout of 0.15 as those parameters performed best on `valid`. We  
 212 train with an Adam optimizer, learning rate of 1e-3, batch size 256 for 100k steps. One training run  
 213 took 4 hours on 4 TPU v2. We trained from scratch and exclusively on MathWriting (`train` and  
 214 `synthetic`). The model is similar to [24], replacing LSTM layers by Transformer layers and not  
 215 using any external language model on top.

216 **VLM** We fine-tuned a large Vision-Language Model PaLI [28] on MathWriting (`train` and  
 217 `synthetic`). We used the representation proposed in [25] where an ink is represented as both a  
 218 sequence of points (similar to CTC Transformer) and its rasterized version (similar to OCR). We train  
 219 three models with different data shuffling for 200k steps with batch size 128, learning rate 0.3 and  
 220 dropout 0.2. One training run took 14 hours on 16 TPU v5p. Models were finetuned exclusively on  
 221 `train` and `synthetic` MathWriting data. Overall, it took 2 TPU v2 days and 28 TPU v5p days to  
 222 run the experiments.

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

Model	Input	Parameters	CER on <code>valid</code>	CER on <code>test</code>
OCR [26]	Image	-	6.50	7.17
CTC Transformer [25]	Ink	35M	4.52 (0.08)	5.49 (0.05)
PaLI [25]	Image+Ink	700M	4.47 (0.08)	5.95 (0.06)

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

229 Figure 5 shows examples of model mistakes. Two of the main causes of mistakes are confusing  
 230 similar-looking characters like “z” and “2”, and errors in the structural arrangement of the characters,  
 231 for instance not placing a sub-expression in a subscript or superscript.

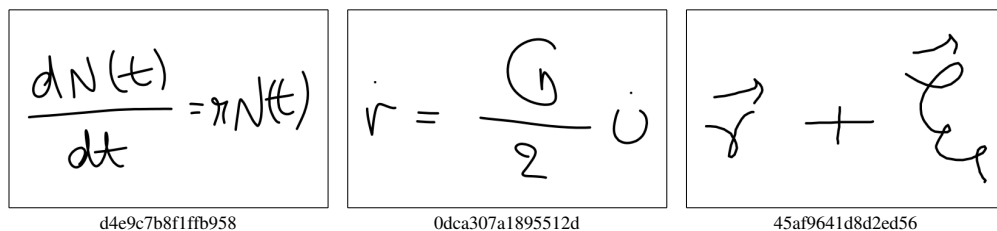
## 232 5 Discussion

### 233 5.1 Differences in Writing Style

234 The number of contributors was large enough that a variety of writing styles are represented in the  
 235 dataset. An example for different ways of writing letter ‘r’ can be seen in Figure 6. Additional

	$\psi_2(\Omega_2^{\Omega_2})$	$\frac{1}{2}(z+1)$	$\theta \in [q_i, q_i+1]$
Ground Truth	$\psi_1(\Omega_2^{\Omega_2})$	$\frac{1}{2}(z+1)$	$\theta \in [q_i, q_i+1]$
Predictions	$\psi_1(\Omega_2^{\boxed{\Omega_2}})$	$\frac{1}{2}(\boxed{2}+1)$	$\theta \in [q_i, q_i+\boxed{1}]$

**Figure 5:** Examples of recognition mistakes from the CTC Transformer model. We observe similar mistakes from the other models.



**Figure 6:** Three ways of writing a lowercase 'r'.

236 examples are provided in Figure 7. Similar though less obvious differences exist for other letters.  
 237 Style differences also show through writing order (example – Figure 14, Appendix G).

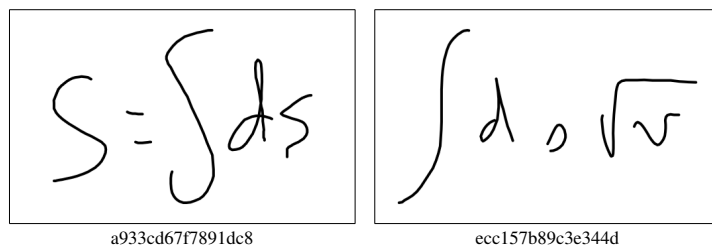
238 **5.2 Recognition Challenges**

239 MathWriting presents some inherent recognition challenges, which are typical of handwritten rep-  
 240 resentations. For example, it's not really possible to distinguish these pairs from the ink alone:  
 241  $\frac{\underline{a}}{b}$  vs  $\frac{a}{\overline{b}}$ , and  $\overline{\omega}$  vs  $\varpi$ .  
 242 We'd like to point out that these ambiguities are not an issue for humans in practice, because they  
 243 rely on contextual information to disambiguate: a particular writing idiosyncrasy, consistency with  
 244 nearby expressions, knowledge of the scientific domain, etc. See Figures 8 and 9 for more examples.

245 **5.3 Dataset Applications and Future Work**

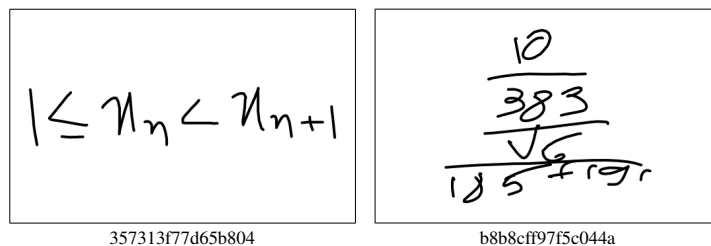
246 Mathwriting can be used to train recognizers for a large variety of scientific fields, and is also large  
 247 enough to enable synthesis of mathematical expressions. Combining it with other large datasets like  
 248 CROHME23 would increase the variety of samples even further, both in terms of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

250 Bounding box information for synthetic samples together with individual symbols are provided to  
 251 enable experimentation with synthetic ink generation. Synthetic samples were generated through the  
 252 straightforward process of pasting inks of individual symbols (`symbols`) exactly where bounding  
 253 boxes were located. This gives synthetic samples a very regular structure, see Figure 1. It is possible

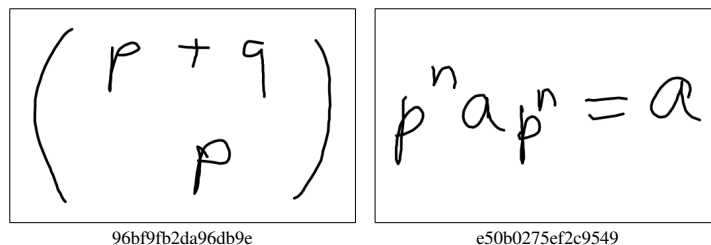


**Figure 7:** Two ways of writing lowercase s.





**Figure 8:** Left: character ambiguity. Is it  $1 \leq x_n < x_{n+1}$  or  $1 \leq n_\eta < n_{\eta+1}$ ? Right: what is the fraction nesting order?



**Figure 9:** Left: `\binom` or 2-element matrix? Right:  $p^n a p^n = a$  or  $p^n a_p n = a$ ?

254 to improve this process by modifying the location, size or orientation of bounding boxes prior to  
 255 generating the synthetic inks. This would soften  $\LaTeX$ 's rigid structure and make synthetic data  
 256 closer to human handwriting. Another application of these bounding boxes would be to bootstrap a  
 257 recognizer that would also return character segmentation information. This kind of output is critical  
 258 for some UI features - for example, editing a handwritten expression.

259 MathWriting can also be improved by varying the label normalization. Changing it can have different  
 260 benefits depending on the application, as mentioned above. We provide the source  $\LaTeX$  string  
 261 for that reason. Another possible improvement in recognition can come from additional contextual  
 262 information, for instance the scientific field [29] that can be added post-hoc. Combining recognizers  
 263 with a language model [24] trained on a large set of mathematical expressions would be a step in a  
 264 similar direction.

## 265 6 Limitations

266 A single sample in MathWriting dataset has one handwritten  $\LaTeX$  formula, see Figure 2. As a result,  
 267 models that are trained on this dataset would probably perform poorly on complete handwritten  
 268 documents, such as the IAMonDo dataset [30]. Also, as the dataset contains only  $\LaTeX$  expressions,  
 269 it is unlikely that models trained on it will accurately recognize handwritten text in English or other  
 270 languages. As shown in Figure 3, some  $\LaTeX$  tokens are way more frequent than others. Some  
 271 infrequent tokens like `\ni` could be hard to recognise.

## 272 7 Conclusion

273 We introduced MathWriting, the largest dataset of online handwritten mathematical expressions  
 274 to date, together with the experimental results of three different types of models. We hope this  
 275 dataset will help advance research in both online and offline mathematical expression recognition.  
 276 Additionally, we invite data practitioners to build on the dataset. We intentionally chose a file format  
 277 for MathWriting close to the one used by CROHME to facilitate their combined use. We also  
 278 provided original or intermediate representations (raw  $\LaTeX$  strings, bounding boxes) to enable  
 279 experimentation with the data itself, and suggested a few directions (Section 5.3).

## 280 References

- 281 [1] Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and Jürgen  
282 Schmidhuber. A novel connectionist system for unconstrained handwriting recognition. *IEEE*  
283 *Transactions on Pattern Analysis and Machine Intelligence*, 31:855–868, 2009.
- 284 [2] Johannes Michael, Roger Labahn, Tobias Grüning, and Jochen Zöllner. Evaluating sequence-to-  
285 sequence models for handwritten text recognition. *2019 International Conference on Document*  
286 *Analysis and Recognition (ICDAR)*, pages 1286–1293, 2019.
- 287 [3] Yutian Liu, Wenjun Ke, and Jianguo Wei. Attention guidance mechanism for handwritten  
288 mathematical expression recognition, 2024.
- 289 [4] Emre Aksan, Fabrizio Pece, and Otmar Hilliges. Deepwriting: Making digital ink editable via  
290 deep generative modeling, 2018.
- 291 [5] Harold Mouchère, Christian Viard-Gaudin, Richard Zanibbi, and Utpal Garain. Icfhr 2014  
292 competition on recognition of on-line handwritten mathematical expressions (crohme 2014).  
293 volume 2014, 09 2014.
- 294 [6] Harold Mouchère, Christian Viard-Gaudin, Richard Zanibbi, and U. Garain. Icfhr2016 crohme:  
295 Competition on recognition of online handwritten mathematical expressions. In *2016 15th*  
296 *International Conference on Frontiers in Handwriting Recognition (ICFHR)*, pages 607–612,  
297 2016.
- 298 [7] May Mowaffaq AL-Taee, Sonia Ben Hassen Neji, and Mondher Frikha. Handwritten recognition:  
299 A survey. *2020 IEEE 4th International Conference on Image Processing, Applications and*  
300 *Systems (IPAS)*, pages 199–205, 2020.
- 301 [8] Yejing Xie, Harold Mouchère, Foteini Liwicki, Sumit Rakesh, Rajkumar Saini, Masaki Nak-  
302 agawa, Cuong Nguyen, and Thanh-Nghia Truong. *ICDAR 2023 CROHME: Competition on*  
303 *Recognition of Handwritten Mathematical Expressions*, pages 553–565. 08 2023.
- 304 [9] Detexify data. <https://github.com/kirel/detexify-data>.
- 305 [10] Yuntian Deng, Anssi Kanervisto, Jeffrey Ling, and Alexander M. Rush. Image-to-markup  
306 generation with coarse-to-fine attention, 2017.
- 307 [11] Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai.  
308 Syntax-aware network for handwritten mathematical expression recognition. In *Proceedings*  
309 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4553–4562,  
310 2022.
- 311 [12] Aida calculus math handwriting recognition dataset. [https://www.kaggle.com/datasets/](https://www.kaggle.com/datasets/aidapearson/ocr-data)  
312 [aidapearson/ocr-data](https://www.kaggle.com/datasets/aidapearson/ocr-data).
- 313 [13] Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh,  
314 Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length general-  
315 ization in large language models, 2022.
- 316 [14] Dusan Varis and Ondřej Bojar. Sequence length is a domain: Length-based overfitting in  
317 transformer models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural*  
318 *Language Processing*. Association for Computational Linguistics, 2021.
- 319 [15] Aleksandr Timofeev, Anastasiia Fadeeva, Andrei Afonin, Claudiu Musat, and Andrii Maksai.  
320 *DSS: Synthesizing Long Digital Ink Using Data Augmentation, Style Encoding and Split*  
321 *Generation*, page 217–235. Springer Nature Switzerland, 2023.
- 322 [16] Arun Narayanan, Rohit Prabhavalkar, Chung-Cheng Chiu, David Rybach, Tara N. Sainath, and  
323 Trevor Strohman. Recognizing long-form speech using streaming end-to-end models, 2019.
- 324 [17] Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashu-  
325 bin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov,  
326 Xiao Wang, Marc van Zee, and Olivier Bousquet. Measuring compositional generalization: A  
327 comprehensive method on realistic data, 2020.

- 328 [18] M. Liwicki and H. Bunke. IAM-OnDB-an on-line english sentence database acquired from  
329 handwritten text on a whiteboard. In *ICDAR'05*. IEEE, 2005.
- 330 [19] Dusan Varis and Ondřej Bojar. Sequence length is a domain: Length-based overfitting in  
331 transformer models. pages 8246–8257, 01 2021.
- 332 [20] Masato Neishi and Naoki Yoshinaga. On the relation between position information and sentence  
333 length in neural machine translation. In *Proceedings of the 23rd Conference on Computational  
334 Natural Language Learning (CoNLL)*, pages 328–338, Hong Kong, China, November 2019.  
335 Association for Computational Linguistics.
- 336 [21] Johannes Michael, Roger Labahn, Tobias Grüning, and Jochen Zöllner. Evaluating sequence-to-  
337 sequence models for handwritten text recognition, 2019.
- 338 [22] George Retsinas, Giorgos Sfikas, Basilis Gatos, and Christophoros Nikou. Best practices for a  
339 handwritten text recognition system, 2024.
- 340 [23] Dmitrijs Kass and Ekta Vats. Attentionhtr: Handwritten text recognition based on attention  
341 encoder-decoder networks, 2022.
- 342 [24] Victor Carbune, Pedro Gonnet, Thomas Deselaers, Henry A. Rowley, Alexander Daryin, Marcos  
343 Calvo, Li-Lun Wang, Daniel Keysers, Sandro Feuz, and Philippe Gervais. Fast multi-language  
344 lstm-based online handwriting recognition, 2020.
- 345 [25] Anastasiia Fadeeva, Philippe Schlattner, Andrii Maksai, Mark Collier, Efi Kokiopoulou, Jesse  
346 Berent, and Claudiu Musat. Representing online handwriting for recognition in large vision-  
347 language models, 2024.
- 348 [26] Google Cloud. Detect handwriting in image, 2023.
- 349 [27] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist  
350 temporal classification: Labelling unsegmented sequence data with recurrent neural networks.  
351 In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, page  
352 369–376, New York, NY, USA, 2006. Association for Computing Machinery.
- 353 [28] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz,  
354 Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan  
355 Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish  
356 Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol  
357 Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and  
358 Radu Soricut. Pali: A jointly-scaled multilingual language-image model, 2023.
- 359 [29] Mark Collier, Rodolphe Jenatton, Efi Kokiopoulou, and Jesse Berent. Transfer and marginalize:  
360 Explaining away label noise with privileged information. In *International Conference on  
361 Machine Learning*, 2022.
- 362 [30] Emanuel Indermühle, Marcus Liwicki, and Horst Bunke. Iamondo-database: An online  
363 handwritten document database with non-uniform contents. pages 97–104, 06 2010.

## 364 Checklist

- 365 1. For all authors...
- 366 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
367 contributions and scope? **[Yes]** We claim to provide a large dataset of HME – see  
368 Section 1, together with normalized labels – the process is described in Section 2.4 and  
369 example is provide in Figure 2. Experimental results on this dataset are presented in  
370 Section 4.2.
- 371 (b) Did you describe the limitations of your work? **[Yes]** We described general limitations  
372 of MathWriting dataset in Section 6, limitations of label normalization in Section 2.4.1,  
373 recognition challenges of mathematical expressions in Section 5.2 and sources of noise  
374 in the dataset in Section H.

- 375 (c) Did you discuss any potential negative societal impacts of your work? [N/A] The type  
376 of the dataset we are publishing is not new, there are similar datasets like CROHME23  
377 [8]. Given the widespread use of handwriting recognition, we don't see any potential  
378 negative impacts of our work.
- 379 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
380 them? [Yes] All the participants were payed the minimum hourly rate as discussed in  
381 Section 2.1. The dataset doesn't include any personal information about contributors  
382 apart from their handwriting samples that they agreed to share.
- 383 2. If you are including theoretical results...
- 384 (a) Did you state the full set of assumptions of all theoretical results? [N/A] Our paper  
385 doesn't include any theoretical results.
- 386 (b) Did you include complete proofs of all theoretical results? [N/A] Our paper doesn't  
387 include any theoretical results.
- 388 3. If you ran experiments (e.g. for benchmarks)...
- 389 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
390 mental results (either in the supplemental material or as a URL)? [No] All experiments  
391 in this paper are conducted using publicly available datasets. We provide code in  
392 Github for ink rasterisation, CER computation and expression tokenization. The Visual-  
393 Language Model PaLI used in the experiments is non-open-sourced, so full results  
394 from Table 6 cannot be reproduced.
- 395 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
396 chosen)? [Yes] We specify the training details like number of parameters, learning rate,  
397 dropout, training data, etc. for the models – CTC transformer and PaLI in Section 4.2.
- 398 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
399 ments multiple times)? [Yes] In Section 4.2 we report average and variance of three  
400 training runs with different shuffling of training data.
- 401 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
402 of GPUs, internal cluster, or cloud provider)? [Yes] In Section 4.2 we provide the total  
403 number of TPU days it took to run our experiments.
- 404 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 405 (a) If your work uses existing assets, did you cite the creators? [Yes] We used pretrained  
406 PaLI model for finetuning and cited [28].
- 407 (b) Did you mention the license of the assets? [N/A] As the PaLI model is non-open-  
408 sourced there is no license for it.
- 409 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
410 URLs to the dataset and code are provided in Section 1.
- 411 (d) Did you discuss whether and how consent was obtained from people whose data you're  
412 using/curating? [Yes] We discussed the conditions of data collection in Section 2.1.
- 413 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
414 information or offensive content? [Yes] We mention in Section 2.1 that there is no  
415 personally identifiable information present in the dataset and offensive content is highly  
416 unlikely given the nature of the dataset.
- 417 5. If you used crowdsourcing or conducted research with human subjects...
- 418 (a) Did you include the full text of instructions given to participants and screenshots, if  
419 applicable? [Yes] We describe the instructions of the data campaigns in Section 2.1 as  
420 they are quite simple – to write a rendered expression provided on the screen.
- 421 (b) Did you describe any potential participant risks, with links to Institutional Review  
422 Board (IRB) approvals, if applicable? [N/A] The paper does not involve research with  
423 human subjects.
- 424 (c) Did you include the estimated hourly wage paid to participants and the total amount  
425 spent on participant compensation? [Yes] We write about the employment of partic-  
426 ipants in Section 2.1, we don't disclose the exact amount of compensation as it is  
427 confidential information.

## 428 Appendix

### 429 A $\LaTeX$ template for label rendering

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

```
431 \usepackage{amsmath}
432 \usepackage{amsfonts}
433 \usepackage{amssymb}
434 \newcommand{\R}{\mathbb{R}}
435 \newcommand{\C}{\mathbb{C}}
436 \newcommand{\Q}{\mathbb{Q}}
437 \newcommand{\Z}{\mathbb{Z}}
438 \newcommand{\N}{\mathbb{N}}
```

### 439 B Acquisition of $\LaTeX$ Expressions

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

443 The extraction process from Wikipedia followed these steps:

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

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

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

### 463 C Postprocessing of MathWriting dataset

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

469 Some cleanup was performed on the labels (ground truths). The goal was to make the dataset better  
470 suited to training an ML model, and eliminate unavoidable issues that occurred during the collection.  
471 After training some initial models, we manually reviewed samples for which they performed poorly.  
472 This helped identify a lot of unusable inks (near-blank, lots of stray strokes, scribbles, etc.), and

473 a lot of ink/label discrepancies. A fairly common occurrence was a contributor forgetting to copy  
474 part of the prompted expression. We adjusted the label to what was actually written unless the ink  
475 contained a partially-drawn symbol, in which case we discarded the sample entirely. In this process  
476 we eliminated or fixed around 20k samples.

477 The most important postprocessing step was to normalize the labels: there are many different ways  
478 to write a mathematical expression in  $\text{\LaTeX}$  format that will render to images that are equivalent in  
479 handwritten form. We applied a series of transformations to eliminate as many variations as possible  
480 while retaining the same semantic. This greatly improved the performance of models and made their  
481 evaluation more precise. We publish both the normalized and raw (unnormalized) labels, to enable  
482 people to experiment with other normalization procedures.

483 This normalization is similar to what [10] did, but pushed further because of the specifics of hand-  
484 written MEs. See Section 2.4 for more detail.

## 485 **D Dataset split**

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

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

## 495 **E Label Normalization**

### 496 **E.1 Syntactic Variations**

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

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

---

Sub/superscript are put in braces, <code>\over</code> is replaced by <code>\frac</code>
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 <code>\rm</code> are dropped
Raw: <code>\tilde{A}_{0}=\frac{ND}{\sigma_{\rm as}+\sigma_{\rm es}}</code>
Normalized: <code>A_{0}=\frac{ND}{\sigma_{as}+\sigma_{es}}</code>

---

Matrix environments with delimiters like <code>bmatrix</code> are replaced by <code>matrix</code> surrounded by delimiters
Commands like <code>\cos</code> 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 <code>\big</code> are dropped
Raw: <code>\big(\tfrac{a}{N}\big)</code>
Normalized: <code>(\frac{a}{N})</code>

---

**Figure 10:** Examples of expression normalization. See Section 2.4 for details.

## 519 E.2 Differences Between Print And Handwriting

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

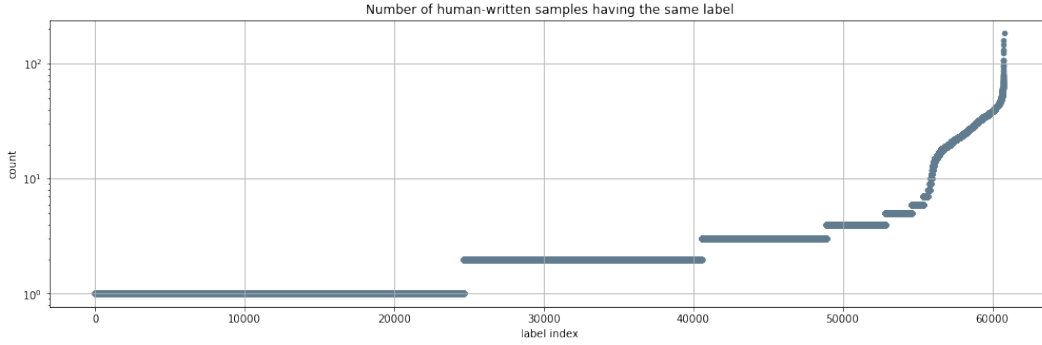
- 521
- color
  - 522 • accurate spacing: e.g. `\sim`, `\quad`.
  - 523 • font style and size: e.g. `\mathrm`, `\mathit`, `\mathbf`, `\scriptstyle`.

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

- 526
- font families: Fraktur, Calligraphic. In practice, only Blackboard (`\mathbb`) has been  
527 written consistently enough by contributors that we were able to keep it: `\mathcal` and  
528 `\mathfrak` are dropped.
  - 529 • some variations like `\rightarrow`  $\rightarrow$  and `\longrightarrow`  $\longrightarrow$ .
  - 530 • some character variations. e.g. `\varrho`, `\varepsilon`
  - 531 • size modifiers like `\left`, `\right`, `\big`. Similarly, variable-width diacritics like  
532 `\widehat`.

## 533 F Additional dataset statistics

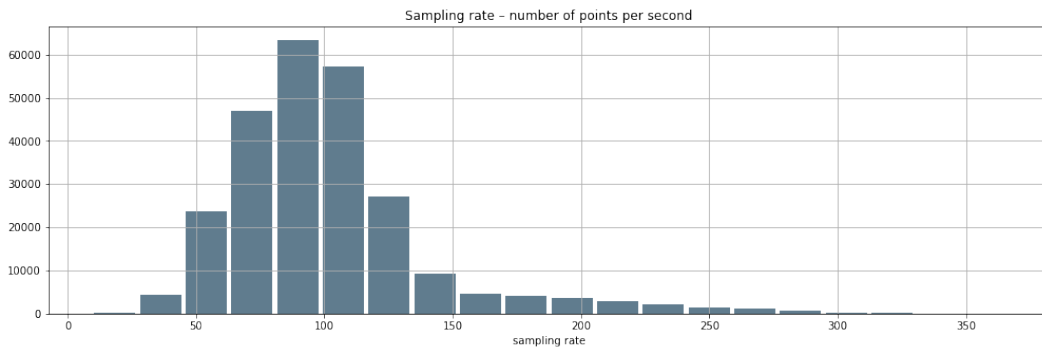
534 In this section we show additional graphs that illustrate dataset statistics that are described in Section 3.  
535 The frequencies of normalized L<sup>A</sup>T<sub>E</sub>X expressions are presented in Figure 11. Figure 12 illustrates the



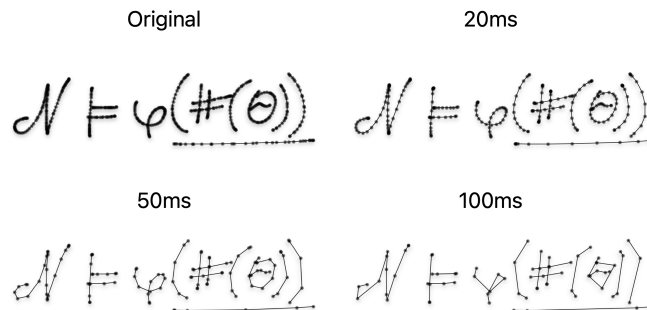
**Figure 11:** Counts of inks corresponding to the same normalized expression, ordered by increasing count. Each position on the horizontal corresponds to a unique normalized expression. Almost 5k unique expressions have been written 10 times or more by contributors.

536 distribution of sampling rates within human-written data. Results of resampling points in time are  
 537 presented in Figure 13.

538 with time resampling are given on Figure 13.



**Figure 12:** Histogram of sampling rates in human-written data of MathWriting dataset.



**Figure 13:** Examples of time resampling with different time periods. Larger periods result in shorter sequences of points.

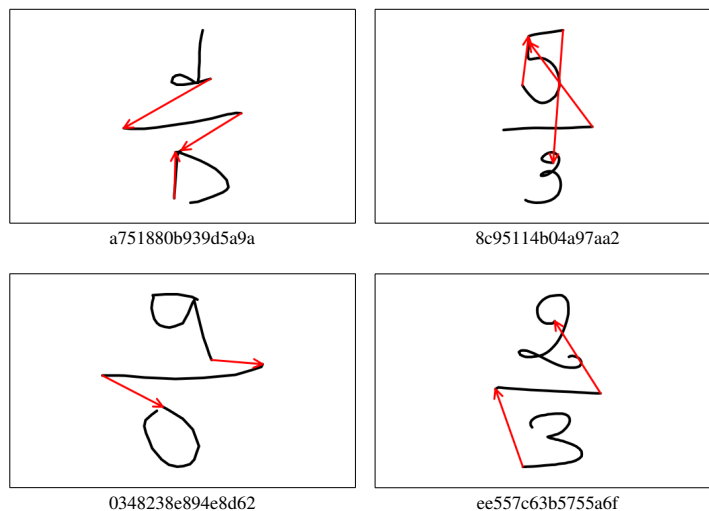
539 **G Variety of Writing Styles**

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



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 7:** 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.



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

## 542 H Sources of Noise

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

545 **Stray strokes** These do not carry any meaning and should be ignored by any recognizer. Since  
 546 they also appear in real applications, there could be some benefit in having some in the dataset to  
 547 teach the model about them. That said, it being usually easier to add noise rather than to remove it,  
 548 we made the choice of discarding as many inks containing stray strokes as possible. Not all inks with  
 549 stray strokes have been found and removed though (e.g. `train/9e64be65cb874902.inkml` that  
 550 was discovered post-publication). The fraction of inks containing stray strokes is significantly lower  
 551 than 1%, and should not be an issue for training a model.

552 **Incorrect ground truth** Contributors did not always copy the prompt perfectly, leading to a variety  
 553 of differences. In most of the cases we spotted, we were able to fix the label to match what had  
 554 actually been written. A short manual review once the dataset was in its final state showed the rate  
 555 of incorrect ground truth to be between 1% and 2%. Most of the mistakes are very minor, usually  
 556 a single token added, missing or incorrect. Errors here also come from ambiguities or misuse of

557 the L<sup>A</sup>T<sub>E</sub>X notation: expressions coming from Wikipedia contain some misuse like using `\Sigma`  
 558 where `\sum` was more appropriate, `\triangle` instead of `\Delta`, `\triangledown` instead of  
 559 `\nabla`, `\begin{matrix}\end{matrix}` instead of `\binom`, and also some handwriting-specific  
 560 ambiguities like `\dagger` vs `\top` vs `T`. There are also some instances where reference numbers or  
 561 extra punctuation are included.

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

## 572 I Tokenization Code

573 Python code used in this work to tokenize L<sup>A</sup>T<sub>E</sub>X mathematical expressions.

```
574 import re
575
576 _COMMAND_RE = re.compile(
577     r'\\(\mathbb{[a-zA-Z]}|begin{[a-z]+}|end{[a-z]+}|operatorname*|[a-zA-Z]+|.)')
578
579 def tokenize_expression(s: str) -> list[str]:
580     tokens = []
581     while s:
582         if s[0] == '\\':
583             tokens.append(_COMMAND_RE.match(s).group(0))
584         else:
585             tokens.append(s[0])
586
587         s = s[len(tokens[-1]):]
588
589     return tokens
```

## 590 J Tokens

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

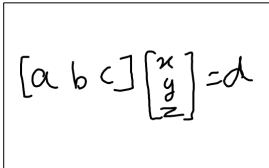
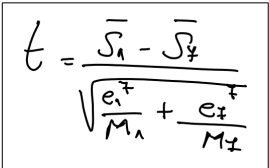
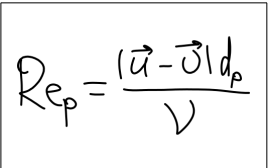
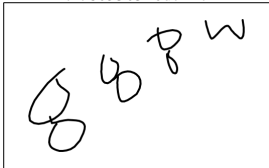
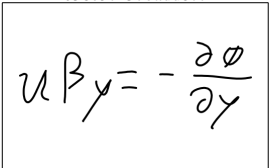
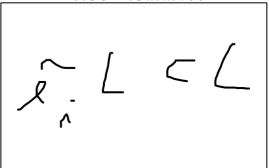
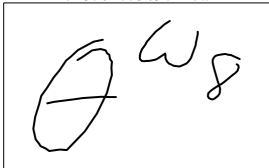
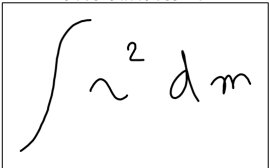
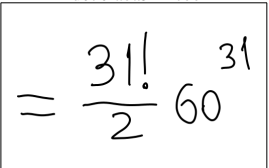
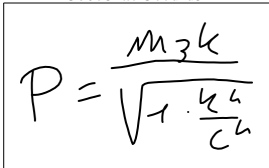
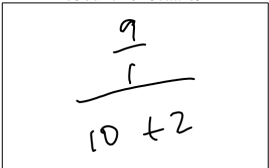
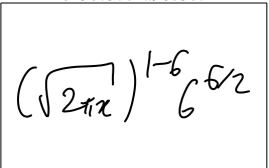
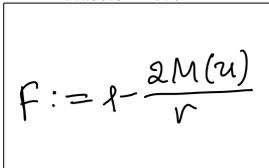
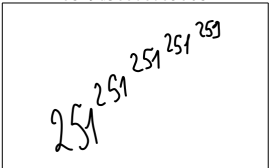
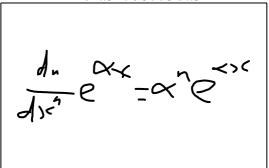
- 593 • Syntactic tokens: `_` `^` `{` `}` `&` `\\` `space`
- 594 • Latin letters and numbers: `a-z` `A-Z` `0-9`
- 595 • Blackboard capital letters `\mathbb{A}`-`\mathbb{Z}` `\mathbb`
- 596 • Latin punctuation and symbols: `,` `;` `!` `?` `.` `()` `[]` `\{` `\}` `*` `/` `+` `-` `_` `&` `\#` `\%` `|` `\backslash`
- 597 • Greek letters: `\alpha` `\beta` `\delta` `\Delta` `\epsilon` `\eta` `\chi` `\gamma` `\Gamma`  
 598 `\iota` `\kappa` `\lambda` `\Lambda` `\nu` `\mu` `\omega` `\Omega` `\phi` `\Phi` `\pi` `\Pi` `\psi`  
 599 `\Psi` `\rho` `\sigma` `\Sigma` `\tau` `\theta` `\Theta` `\upsilon` `\Upsilon` `\varphi` `\varpi`  
 600 `\varsigma` `\vartheta` `\xi` `\Xi` `\zeta`
- 601 • Mathematical constructs: `\frac` `\sqrt` `\prod` `\sum` `\iint` `\int` `\oint`
- 602 • Diacritics and modifiers - Note the absence of the single-quote character, which is normalized  
 603 to `~{\prime}`:  
 604 `\hat` `\tilde` `\vec` `\overline` `\underline` `\prime` `\dot` `\not`
- 605 • Matrix environment: `\begin{matrix}` `\end{matrix}`

- 606 • Delimiters: `\langle \rangle \lceil \rceil \lfloor \rfloor \|`
- 607 • Comparisons: `\ge \gg \le \ll \lt`
- 608 • Equality, approximations: `= \approx \cong \equiv \ne \propto \sim \simeq`
- 609 • Set theory: `\in \ni \notin \sqsubseteq \subset \subseteq \subsetneq \supset \supseteq \emptyset`
- 610
- 611 • Operators: `\times \bigcap \bigcirc \bigcup \bigoplus \bigvee \bigwedge \cap \cup \div \mp \odot \ominus \oplus \otimes \pm \vee \wedge`
- 612
- 613 • Arrows: `\hookrightarrow \leftarrow \leftrightarrows \Leftrightarrow \longrightarrow \mapsto \rightarrow \Rightarrow \rightleftarrows \iff`
- 614
- 615 • Dots: `\bullet \cdot \circ`
- 616 • Other symbols: `\aleph \angle \dagger \exists \forall \hbar \infty \models \nabla \neg \partial \perp \top \triangle \triangleleft \triangleq \vdash \Vdash \vdots`
- 617
- 618

## 619 K Examples of inks

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

### 623 K.1 Human-Written Samples

624	 cf786356546d722c	 658e3c257badd8cc	 40b3844b5aeac00
625	 c28701c7369c22ba	 97e829ac79e851fe	 fd676faba32f1cbb
626	 5e0c81acf5ccdf3	 b3ed172628caafe9	 6150c57b1a98b5ec
627	 bca00e3111b70212	 e829d5eb7e7b3b68	 44bfa5fc08eb5da5
628	 d233aaf208bd7568	 ee6f0bfd294aa209	 6579f917f0ba236b

629

$$v_f = \frac{M(1+e)}{a(1-e)}$$

cea67d239b9f8884

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

bb4bca53d0e6336d

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

b14bca3fc2d2819a

630

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

2409d2feaa79b9d7

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

51486ff88b789d6d

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

02229a0c174d8dbe

631

$$\tilde{C}_7$$

478e10a15203fa3a

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

ccb15825579a096b

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

7a4b95285de0caf0

632

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

ac25b4d053596ded

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

7d597c52bf8bdd1e

$$\left(\frac{24}{\sqrt{39}}\right)^6 \cdot 36 \cdot 3^4 + 404$$

acc5f6620fad1ce8

633

$$t \equiv \sqrt{1 - 1/5^2}$$

88c3551d373c72e5

$$\theta_b^m = 5_{ab}^m \theta^a$$

5009d32d32f80324

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

068de3aad90c403c

634

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

a3cf115524f0c55b

$$\bar{E} = 0, \bar{E}_i \bar{E}_j = \frac{1}{3} \bar{E}^2 \delta_{ij}$$

355f5df56a16913a

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

d9b2ce7aa3495888

635

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

adceb80fdadf9f6e

$$264 - 175 - 365 - 445$$

25892f7caecac8c36

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

41e1261a951c6f33

636

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

02a7f7f172671fb4

$$\deg_P f = \deg_P g$$

0d848d4b170d36b9

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

2adc4f10d42b641c

637

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

408f904038dcbbba0

638 **K.2 Synthetic Samples: Expressions from Wikipedia**

639  $\infty, \beta, \eta = 1/(\sigma\mu), L_2 \otimes I_S$

133829b5a10b783f

11623165e9bab0e4

5ca38d17bf2bea0a

640  $u(x_1, \dots, x_s) \cdot \sum_i^s c_i, \alpha_{mn}^* = \int v_m^* u_n d\tau, \frac{1}{(4\pi t)^{\frac{n}{2}}} e^{-\frac{|x|^2}{4t}}$

089650cc894c024b

8e88b75cb5f03bf4

5c1573b41e762307

641  $E(r) \equiv r - \int dr' k(r, r'), \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix}, |1\rangle_A |1\rangle_B$

68e9560a27a093c8

daab3bae071f8bc0

77602abaea39b774

642  $Z_q, n^3 < N_1 N_2 N_3, \frac{\partial}{\partial t} dr_s = \frac{1}{2} g^{pq} v_{p2} dr_p$

eb817bcbfd11df18

3bd062b33ea5db6f

51ef5122de326151

643  $c_v = \frac{nR}{\gamma - 1}, xy^2, \bigoplus_{n \in \mathbb{Z}} \text{Hom}_n(A, B)$

f2a613fc323df342

cbd19abc03ba4098

1879aa5c882b445d

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

094c7e52a3f0934d

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

1cd654228d7ca6bb

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

fc00050933165b70

647

$$A^{(1)} = \begin{pmatrix} 4 & 2 & 1 \\ 0 & 6 & 8.5 \\ 0 & 5 & \frac{2}{3} \end{pmatrix}, \quad P^{(1)} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

60553a301aaf84a3

648

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

b2aeaf7a0fd30ed6

649

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

254dbc2b3843dcf8

650

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

0772aeaac09d3415

651

$$r = \begin{pmatrix} \frac{\partial A}{\partial x} & \frac{\partial A}{\partial y} \\ \frac{\partial A}{\partial x} & \frac{\partial A}{\partial y} \end{pmatrix} + \frac{1}{c\pi v} e^{i\psi} \kappa_0(\psi) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \Pi = \frac{e}{c\pi} \nabla(\psi r)$$

20b4ebf292cfa8d1

652

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

96613ae167f35f8a

653

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

4b1f5165e3698343

654

$$c_{solid, r} = \sqrt{\frac{k + \frac{4}{3}g}{\rho}} = \sqrt{\frac{E(1-\nu)}{\rho(1+\nu)(1-2\nu)'}}$$

6f86778996d5f514

655

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

685bd5676bda74ce

656

$$\begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & \omega & \omega^2 & \omega^3 & \omega^4 \\ 1 & \omega^2 & \omega^4 & \omega & \omega^3 \\ 1 & \omega^3 & \omega & \omega^4 & \omega^2 \\ 1 & \omega^4 & \omega^3 & \omega^2 & \omega \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad \omega = e^{-\frac{j2\pi}{5}}$$

b79dfccd7f5b5cb4

657 **K.2.1 Synthetic Samples: Generated Fractions**

658

$$\frac{37.13006145}{14.630468}$$

4677b76acec23465

$$\frac{333399999999}{28.870145}$$

eef0cd70f8872a9c

$$\frac{t}{40}$$

d9df5ffcf81d07

659

$$\frac{\begin{array}{r} 26 \\ \times 22744445333 \\ \hline 6644444222 \\ 77 \\ \hline 11.91 + 4444442=22.3 \end{array}}{77}$$

244ad5e60c9fea92

$$\frac{\begin{array}{r} 60 \\ 79 \\ \hline 66619915555 \times 12000000000 \\ \hline 12212200012 \end{array}}{12212200012}$$

88d5f862ad46cb47

$$\frac{68}{100}$$

9067ae238a278b32

660

$$\frac{\begin{array}{r} 999991666000 \\ 1662222111 \\ \hline 777777999999 \end{array}}{1662222111}$$

ce12f955b6ba76a4

$$\frac{61}{11166666100}$$

7fa1aa18a332b211

$$\frac{4}{\frac{7}{b}}$$

1c21c51bc1319124

661

$$\frac{10000000}{\frac{3}{21}}$$

afd5254b25be1256

$$\frac{\begin{array}{r} 5557777777 \\ 104 + \frac{75}{2228888888} \\ \hline 2 \end{array}}{2}$$

409a91ba03e3cc7e

$$\frac{10}{\frac{57.714}{78.7451}}$$

486a38bd87b8ed97

662

$$\frac{\begin{array}{r} 33333226616 \\ +444447777 - 1000 \\ \hline 11 \\ 32.4 \end{array}}{11}$$

05efec2565d6726e

$$\frac{\begin{array}{r} 79.654 \\ 95 \\ \hline 10000 \end{array}}{95}$$

3fdc553e580f2a78

$$\frac{2}{999999444700}$$

b17c3206c2d610b9

663

$$\frac{\begin{array}{r} 999991111 + 5555588822 \\ \hline 44444999119 - 11.91931785 \end{array}}{11}$$

9df33845897752ea