



COLLEGE OF
ENGINEERING, ARCHITECTURE AND TECHNOLOGY

Deep Learning for Recognizing and Converting Handwritten Scientific Equations into LaTeX

Deep Learning (ECEN 5060)

Haya Monawwar and Sungjoo Chung (Group 04)

May 06 2025

Department of Electrical and Computer Engineering

Oklahoma State University – Stillwater Campus

Outline

- Introduction
- Problem
- Dataset
- Methodology
- Results
- Limitations
- Conclusion
- Future Work

Introduction

- As engineers, we often want to fast-track our handwritten notes into digital format for various purposes. **But can we digitalize handwritten equations directly?**
- ...Hence, the motivation for this project.

What if we could extend the problem and get our hand-written equations converted to LaTeX version directly?

Problem

Develop a machine learning based system that accurately recognizes and converts handwritten mathematical expressions/scientific equations into LaTeX format using deep learning techniques.

Dataset

- CROHME 2019 (Competition on Recognition of Handwritten Mathematical Expressions)
- Contains approximately 12,000 HMEs from previous CROHME competitions (2014-2019)
- Format:
 - Greyscale Images: Sized 1000x1000 pixels with 5 pixels of padding.
 - SymLG annotations: SymLG is a structured representation of HMEs that captures both individual symbols and their spatial relationships.
- Annotations are used to create a vocabulary split.
- Train/test/validation split: 8,835 / 2,186 / 1,147 (images and captions each)

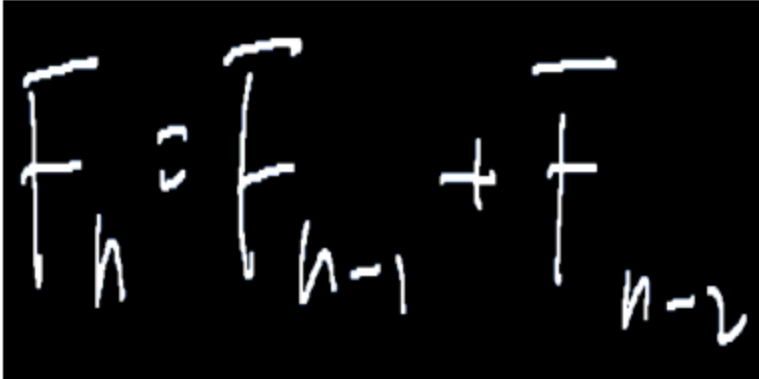
Methodology

Pre-processing

- Tokenizing / creating a unified vocabulary
- Implemented CROHMEDataset to load images and their corresponding tokenized labels.
- Data augmentation:
 - Resizing to 100 x 100
 - Random rotation
 - Color jitter by varying brightness and contrast
 - Normalizing with mean = standard deviation = 0.5
 - Converted images to PyTorch tensors and normalized to [0,1].

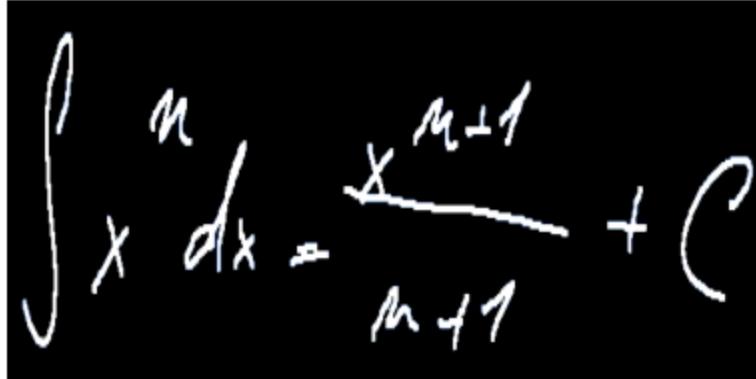
Methodology

GT: $F_{\{n\}} = F_{\{n-1\}} + F_{\{n-2\}}$



A handwritten equation on a black background showing the Fibonacci recurrence relation: $F_n = F_{n-1} + F_{n-2}$. The variables are written in a stylized, slightly cursive font.

GT: $\int x^{\{n\}} dx = \frac{x^{\{n+1\}}}{\{n+1\}} + C$



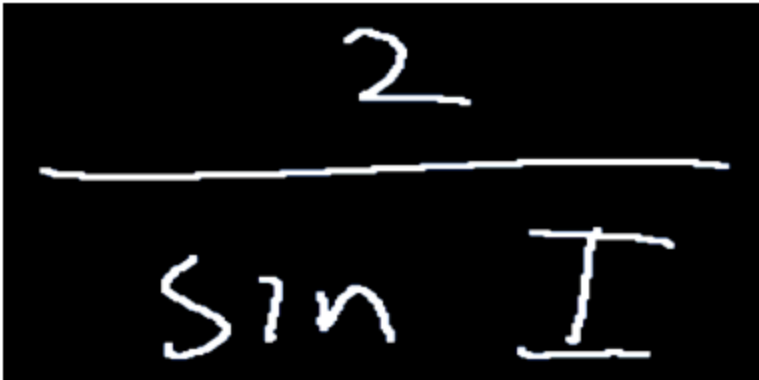
A handwritten equation on a black background showing the power rule for integration: $\int x^n dx = \frac{x^{n+1}}{n+1} + C$. The variables are written in a stylized, slightly cursive font.

GT: $c_{\{t\}}$



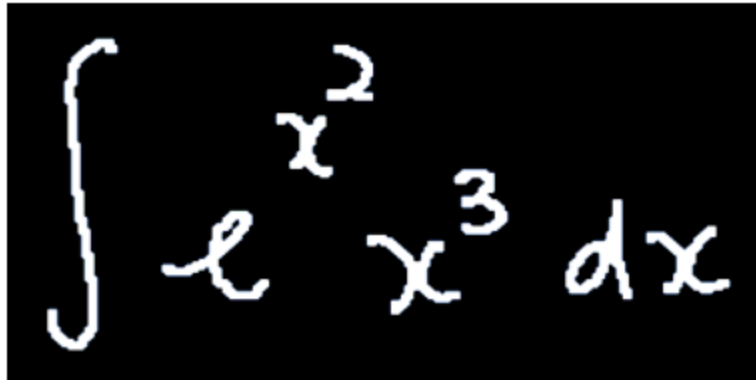
A handwritten expression on a black background showing the variable c_t . The 'c' is a simple capital letter, and the 't' is a simple lowercase letter.

GT: $\frac{\{2\}}{\{\sin I\}}$



A handwritten equation on a black background showing the fraction $\frac{2}{\sin I}$. The '2' is a simple digit, and 'sin I' is written in a simple font.

GT: $\int e^{\{x^2\}} x^{\{3\}} dx$



A handwritten equation on a black background showing the integral $\int e^{x^2} x^3 dx$. The variables are written in a stylized, slightly cursive font.

GT: $([R_{\{y\}}])$



A handwritten expression on a black background showing the variable $([R_y])$. The 'R' is a simple capital letter, and the 'y' is a simple lowercase letter.

Methodology

Baseline Model

- Using a CNN-RNN
- Pre-trained ResNet18
- Freezing the early layers of the ResNet to reduce computation time.
- Adding Batch Normalization and Dropout layers for regularization and preventing overfitting.
- Enhancing the LSTM layers to improve temporal sequence handling.

```
def __init__(self, num_classes):
    super(CRNN, self).__init__()
    resnet = models.resnet18(pretrained=False)
    state_dict = torch.load("/kaggle/input/resnet18/resnet18-f37072fd.pth")
    resnet.load_state_dict(state_dict)

    # Freeze first few layers to avoid overfitting
    for param in list(resnet.children())[:5]:
        for p in param.parameters():
            p.requires_grad = False

    # Extract CNN layers up to layer3
    self.cnn = nn.Sequential(
        *list(resnet.children())[:3],          # Keep until layer3
        nn.BatchNorm2d(256),                  # Add BatchNorm after last conv layer
        nn.Dropout2d(p=0.3)                   # Dropout after batchnorm
    )

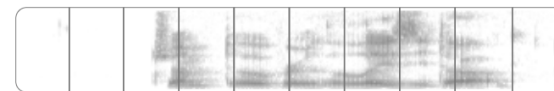
    # Bidirectional LSTM with dropout
    self.rnn = nn.LSTM(
        input_size=256,
        hidden_size=256,
        num_layers=2,
        dropout=0.5,                          # Dropout between LSTM layers
        bidirectional=True,
        batch_first=True
    )

    self.dropout_fc = nn.Dropout(p=0.3)       # Dropout before final classification
    self.fc = nn.Linear(512, num_classes)

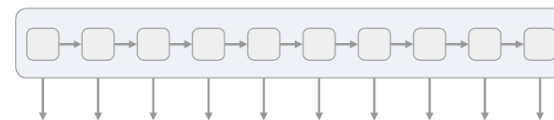
def forward(self, x):
    x = self.cnn(x) # (B, 256, H, W)
    x = nn.functional.adaptive_avg_pool2d(x, (1, x.size(3))) # (B, 256, 1, W)
    x = x.squeeze(2) # (B, 256, W)
    x = x.permute(0, 2, 1) # (B, W, 256)
    x, _ = self.rnn(x)
    x = self.fc(x)
    return x
```


Loss function

- CTC Loss [2]
 - Enables models to learn alignments between input sequences and output label sequences without requiring explicit frame-level annotations.
 - It allows for flexible spacing and repeated characters, while handling intermediate blank tokens.



We start with an input sequence, like a spectrogram of audio.



The input is fed into an RNN, for example.

h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o	o
€	€	€	€	€	€	€	€	€	€

The network gives $p_t(a | X)$, a distribution over the outputs $\{h, e, l, o, \epsilon\}$ for each input step.

h	e	€	l	l	€	l	l	o	o
h	h	e	l	l	€	€	l	€	o
€	e	€	l	l	€	€	l	o	o

With the per time-step output distribution, we compute the probability of different sequences:

h	e	l	l	o
e	l	l	o	
h	e	l	o	

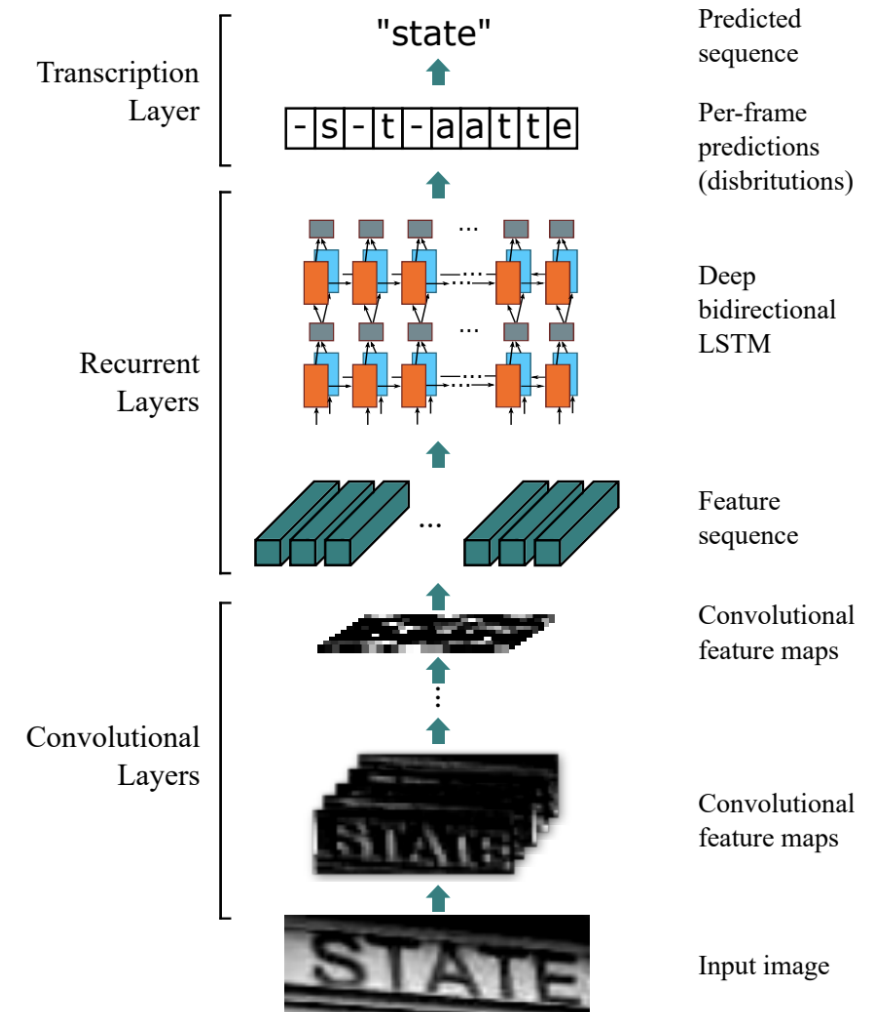
By marginalizing over alignments, we get a distribution over outputs

Evaluation Metric

- **Levenshtein distance:** a string metric that measures how many **single-character edits** are needed to change one string into another. These edits include:
 - Insertion
 - Deletion
 - Substitution
- Character Error Rate (CER): $\text{Levenshtein distance} / \text{Total number of characters}$
- Word Error Rate (WER): $\text{Levenshtein distance (words)} / \text{Total number of words}$

Methodology – 1st iteration

- CRNN framework from [1]
 - Convolution layer: extract features from images
 - Recurrent layer: makes prediction for each frame of the feature sequence (bi-directional)
 - Transcription layer: translates per-frame predictions into labeled sequence
- Feature maps are "sliced" to keep the height, helps distinguish stacked symbols like fractions or superscripts



Methodology – 1st iteration results

- Image is resized to 32 x 300 (unlike paper, which uses varying width)
- Train and Validation WER: 19.98% and 21.66%
- However, predictions were spammed with "Right", inflating the accuracy
- Indicators such as "Right", "Below", "Inside"... are unnecessary due to the CRNN framework

5 Best (Lowest Error) Predictions

Err=6
GT: z Right = Right a Right + Right b Right ...
PR: Right Right Right Right Right

Err=6
GT: r Right \geq Right 1
PR: Right Right

Err=7
GT: t Right - Right 6
PR: Right NoRel Right

Err=8
GT: r Right (Right x Right)
PR: Right Right Right Right

Err=9
GT: m Right v
PR: - - \sqrt - - -

$z = a + b$ $r \geq 1$ $t - 6$ \sqrt{x} $m v$

Methodology – 2nd iteration

- Modified the vocabulary to exclude structural / relationship indicators
- Model's predictions were still very bad even after hyperparameter optimization (Train and validation WER: 43.29% and 37.14%)

Err=84
GT:\frac { 2 G M r - 2 r ^ { 3 } } { \pm r \sqrt { r ... }
PR:}

Err=88
GT:R _ { o } = \frac { \{ \frac { \beta + 1 ... } { \beta + 1 ... }
PR: { + +

Err=90
GT:\sum \limits _ { k = 1 } ^ { N } a _ { n ... }
PR: + 3 3 x {

Err=90
GT:\frac { 2 } { \{ \sqrt { 3 } - 1 \} \times \dots }
PR: { +

Err=2
GT:d s
PR: _ _

Err=2
GT:- 1
PR: { {

Err=2
GT:1 m
PR: { {

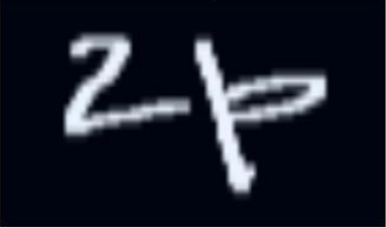
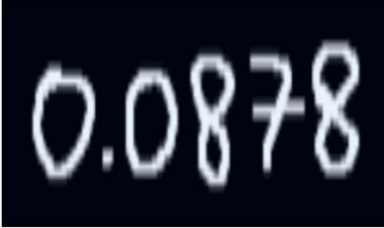
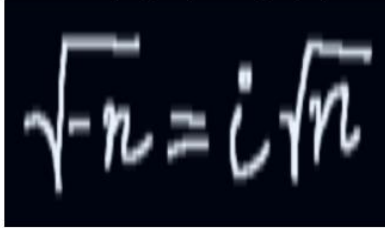

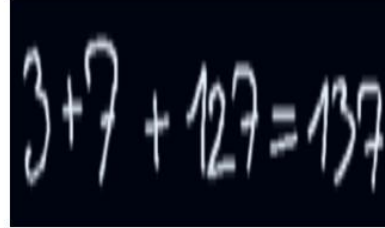
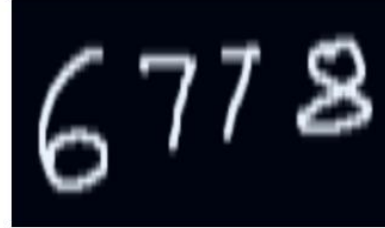
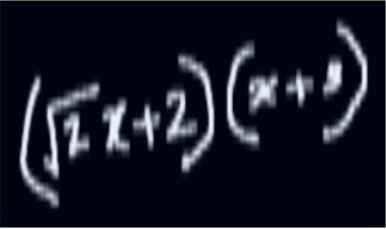





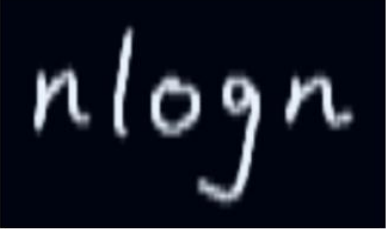
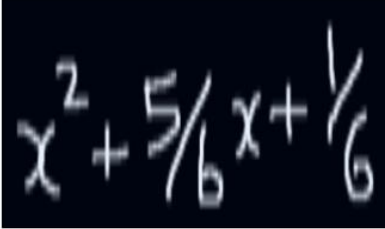
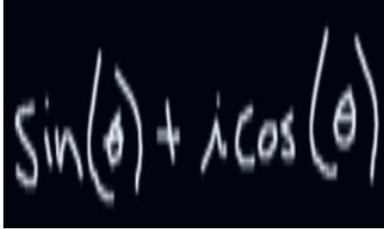
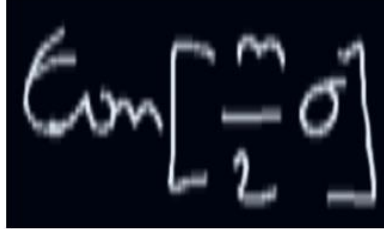
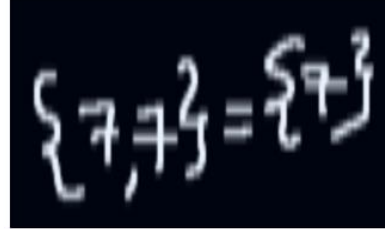

Err=2
GT:6 3
PR: {

Methodology – 3rd iteration

- Greedy decoding:
 - At each timestep, selects the single most probable token
 - Does not consider the **long-term sequence likelihood**. A locally optimal token may lead to poor global predictions.
- Beam decoding:
 - Maintains the top- k most likely sequences at each timestep (e.g., top 5 partial predictions).
 - Common math structures (fractions, subscripts, braces) can be preserved due to global consistency.
- Due to the nature of CTC loss, beam decoding is especially effective in the proposed framework

Methodology – 3rd iteration results

Prediction Samples with Labels

<p>Err=0 GT:2 p PR:2 p</p> 	<p>Err=0 GT:0.0878 PR:0.0878</p> 	<p>Err=0 GT:\sqrt{-n} = i\sqrt{n} PR:\sqrt{-n} = i\sqrt{n}</p> 	<p>Err=0 GT:\cos(m y) PR:\cos(m y)</p> 	<p>Err=0 GT:3 + 7 + 127 = 137 PR:3 + 7 + 127 = 137</p> 	<p>Err=1 GT:6778 PR:678</p> 
<p>Err=1 GT:(\sqrt{2} x + 2)(x + 3) PR:(\sqrt{2} x + 2)(x +)</p> 	<p>Err=1 GT:X \times X \times X PR:x \times X \times X</p> 	<p>Err=2 GT:\log PR:\log e</p> 	<p>Err=2 GT:\frac{1}{c} PR:\frac{1}{c}</p> 	<p>Err=3 GT:\mu \geq 0 PR:M \geq 0</p> 	<p>Err=3 GT:1.0737 + 1.2227 PR:1.0737 + 1.27</p> 
<p>Err=4 GT:n \log n PR:n(0 n)</p> 	<p>Err=5 GT:x^2 + 5/6 x + 1/6 PR:x^2 + 5 {b} x + 1 {6}</p> 	<p>Err=6 GT:\sin(\theta) + i \cos(\theta) PR:\sin(s) + i \cos(\theta)</p> 	<p>Err=7 GT:\tan[\frac{n}{2} \sigma] PR:\tan[\frac{2}{2} 0]</p> 	<p>Err=10 GT:\{7, 7\} = \{7\} PR:z 7.73 = \sum^{\{73\}}</p> 	<p>Err=20 GT:and one goes down from PR:a n (\tan(g \cos d \cos \beta n 0 n))</p> 

Methodology

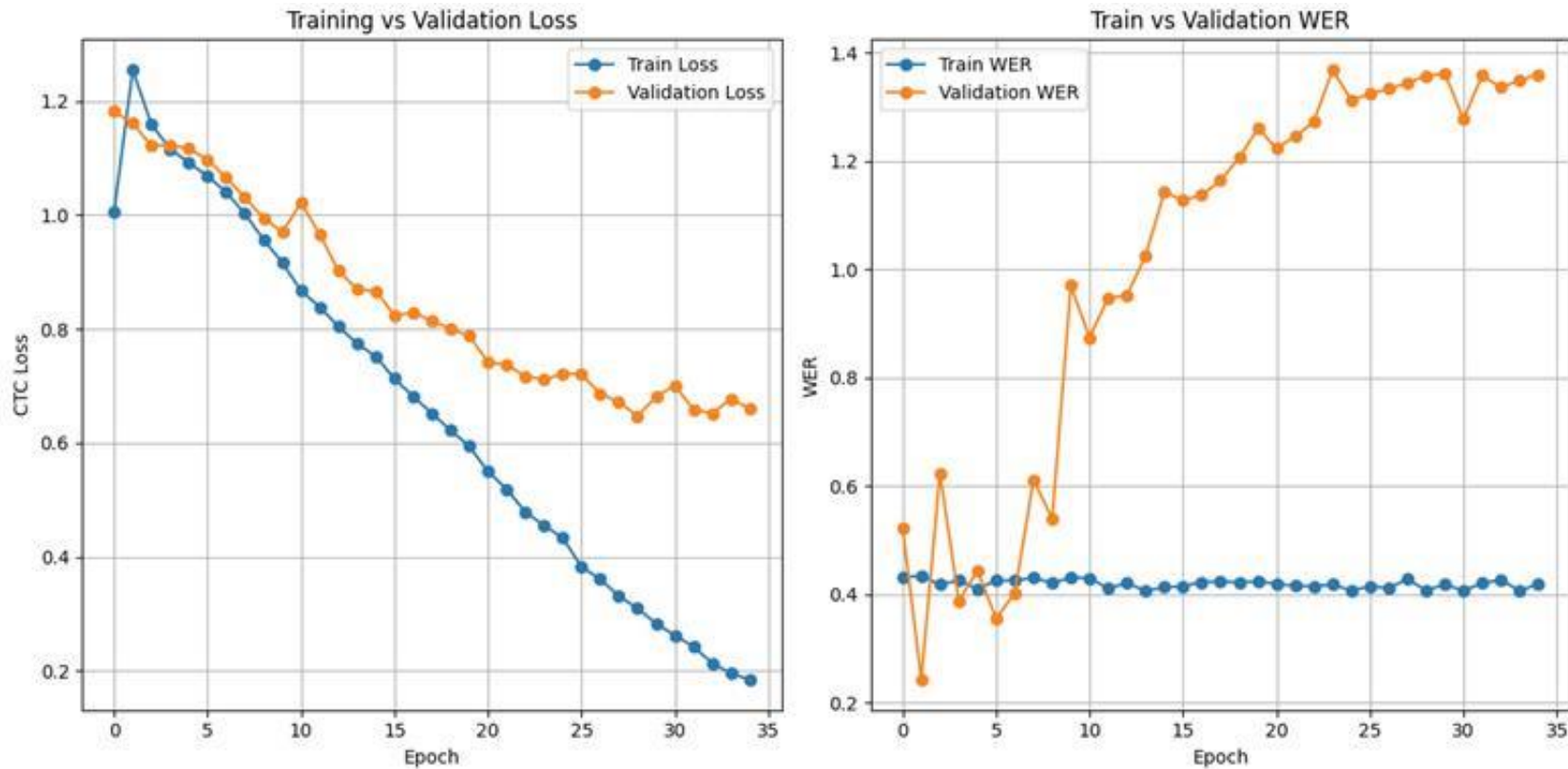
Post-processing

GroundTruth	Prediction	Error
k N	k N	1
1 2	1 2	1
P a	p { a }	6
1 9	7 9	2
2 6	2 6	1
1 m	1 m	1
N m	N	1
H z	H { z }	5
k g	k	1
m v	m v	1

Notice something?

Final Results

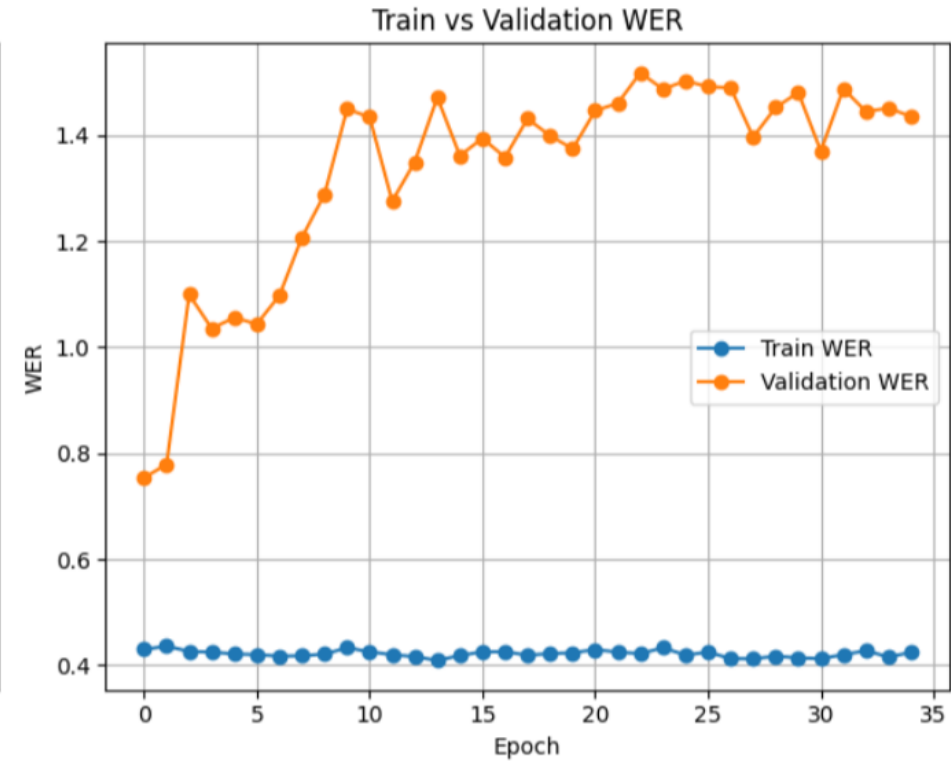
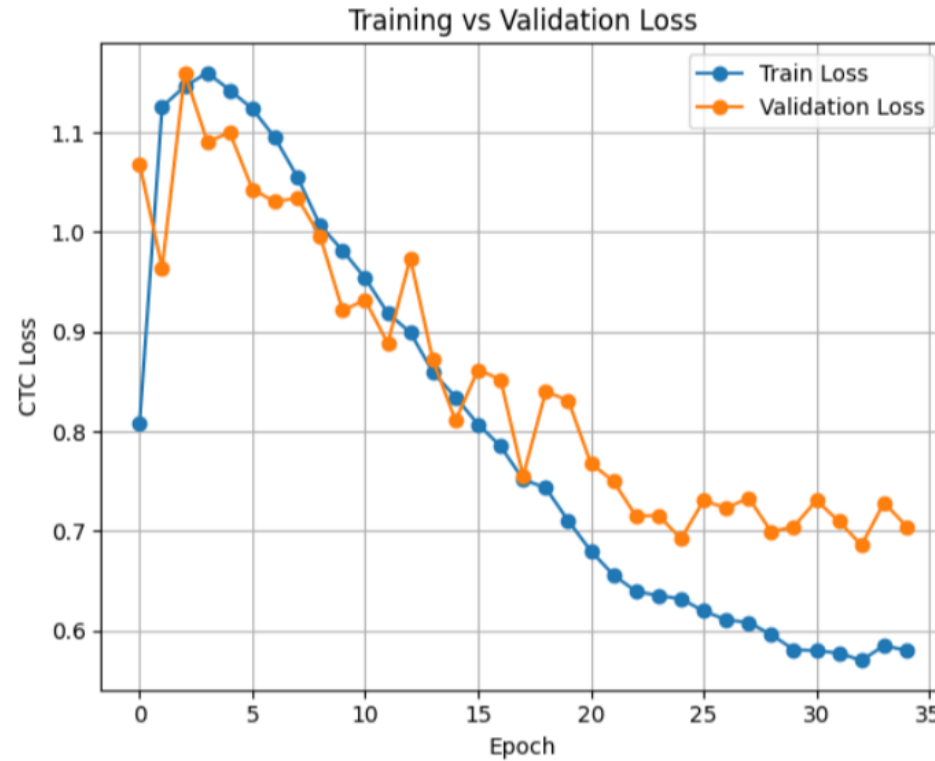
Training loop/Intermediate Results



Baseline model – Iteration 1

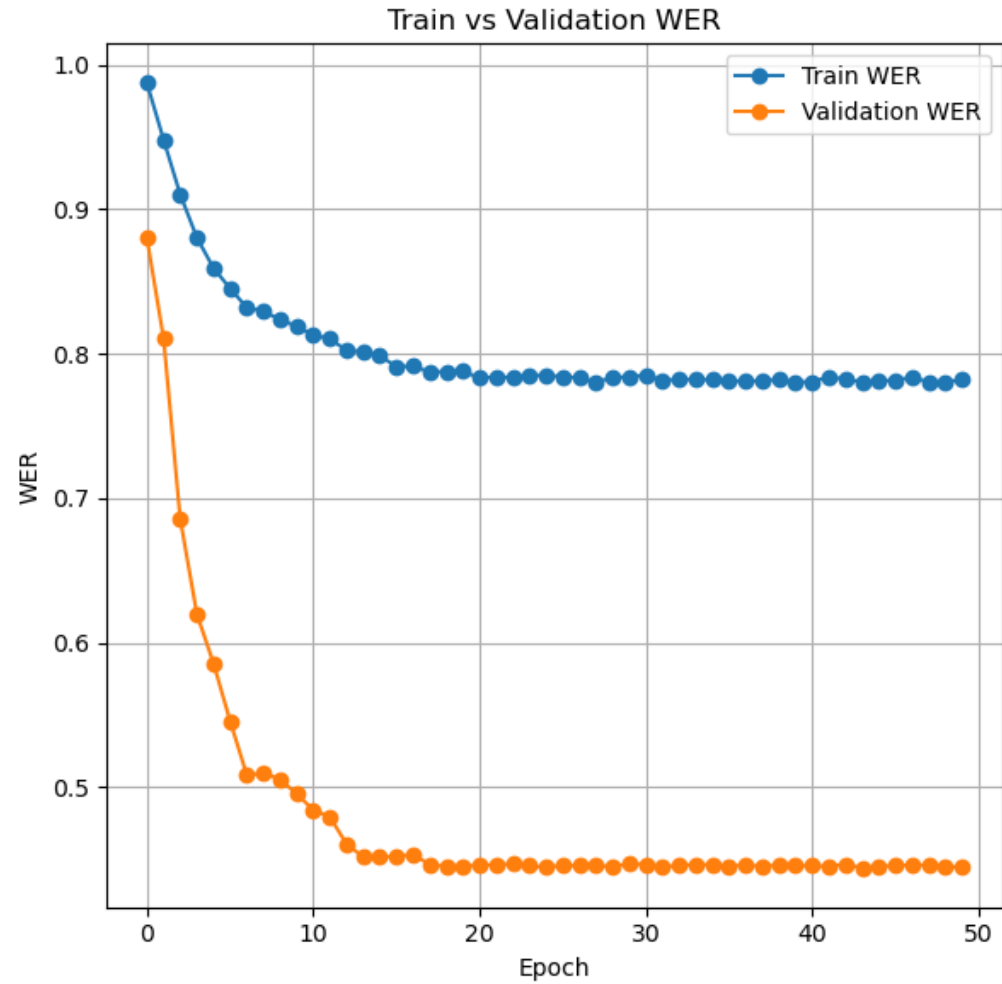
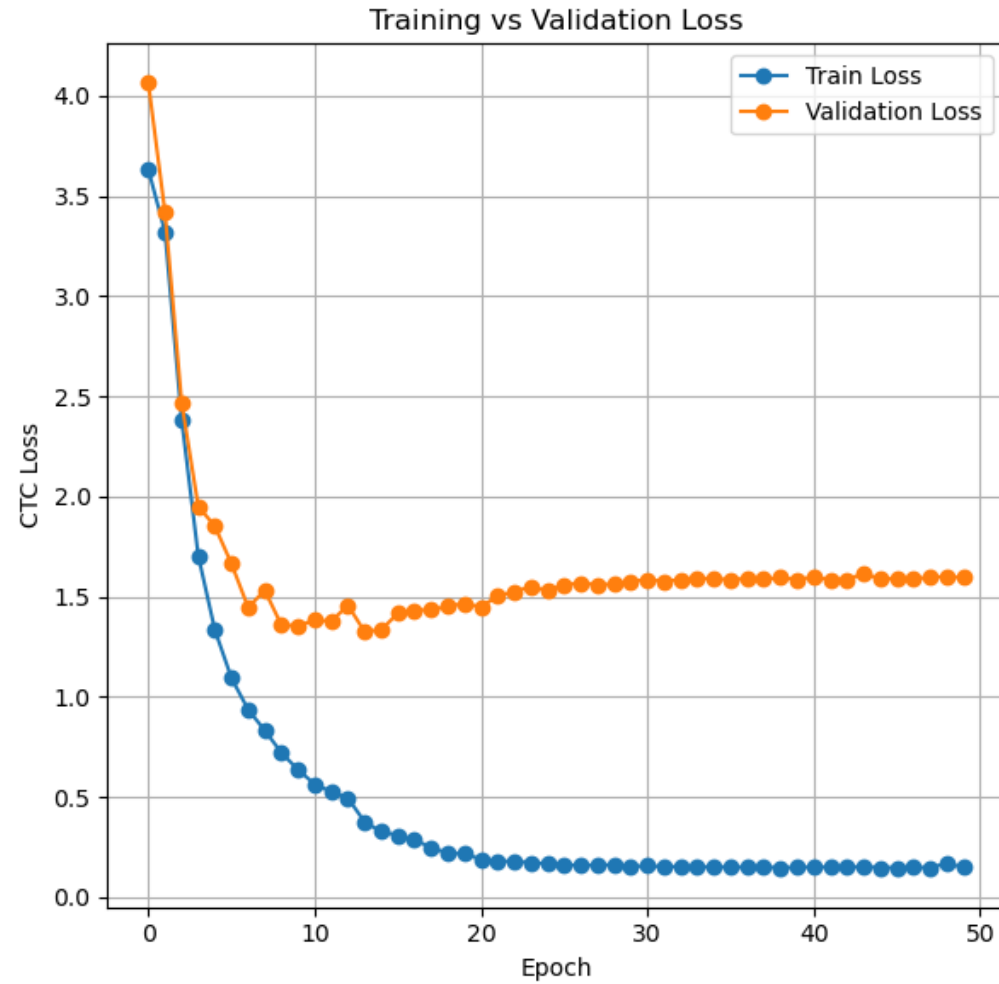
Final Results

Training loop/Intermediate Results



Baseline model – Iteration 2

Final Results



Final Model Architecture and Parameters

Component	Hyperparameter	Value / Setting
Input	Image Size	(3, 32, 300)
CNN Backbone	Architecture	VGG16
RNN	Type	2-layer Bidirectional LSTM
	Hidden Size	256
	Dropout (between layers)	0.5
Fully Connected Layer	Input / Output Size	512 (256 × 2) / 110 tokens
	Dropout	0.3
Loss	Loss Function	CTC
Optimizer	Type	Adam
	Learning Rate	5e-4
Scheduler	Type	ReduceLROnPlateau
	Factor and Patience	0.5 (halves LR on plateau), 2
Decoding	Beam Width	10
Training	Epochs	50

Post Processing

- Due to the garbage token insertion, the WER were inflated
- After preprocessing, we split the predictions by length, and compare the original and cleaned WER
 - Split was done by finding the median expression length value
 - We see a significant improvement in overall WER, especially for short expressions

Expression Length	Original WER	Cleaned WER
Short (≤ 14 tokens)	0.7639	0.2385
Long (> 14 tokens)	0.5344	0.3936

Limitations

- CTC loss suffers with long sequences, as evident from WER
- CTC assumes the output sequence is *monotonically aligned* with the input (left-to-right), which doesn't always hold for math expressions
- The model treats output as a flat sequence of tokens, thus it can't handle hierarchical or spatial relationships
- Even with beam search, the model makes decisions at the token level, which means it can still miss globally optimal expression sequences.

Conclusion

- Developed a deep learning pipeline to recognize and convert handwritten mathematical expressions from images into LaTeX format
- Implemented a CRNN-based model with CTC loss for sequence prediction.
- Iteratively refined the system by adjusting model architecture, vocabulary, and decoding strategies (greedy → beam search).
- Integrated post-processing to remove garbage tokens, significantly improving WER rates, especially for shorter sequences

Future Work

- Attempt a Transformer based model – does it out perform our work?
- Experimenting with deeper trained models.
- Incorporate additional data augmentation strategies, such as elastic distortions or random cropping, to make the model more robust to various writing styles and formats.
- Integrate the proposed mechanism into a mobile/web application to enable wide-spread use.

References

1. Baoguang Shi, Xiang Bai, and Cong Yao. "An End-to-End Train- able Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition". In: IEEE Transactions on Pattern Analysis and Machine Intelligence 39.11 (2016), pp. 2298– 2304.
2. Hannun, "Sequence Modeling with CTC", Distill, 2017.

THANK YOU

Questions?

