

# Image-to-Markup Generation with Coarse-to-Fine Attention

Yuntian Deng, Anssi Kanervisto, Jeffrey Ling, Alexander M. Rush

OpenNMT Workshop Paris, Mar 2018



**HARVARD**  
John A. Paulson  
School of Engineering  
and Applied Sciences



# Outline

1 Introduction: Image-to-Markup Generation

2 Dataset: IM2LATEX-100K

3 Model

4 Experiments

5 Conclusions & Future Work

## Multimodal Generation

Real text is not disembodied. It always appears in context... As soon as we begin to consider the generation of text in context, we immediately have to countenance issues of **typography** and **orthography** (for the written form) and **prosody** (for the spoken form)... This is perhaps most obvious in the case of systems that **generate both text and graphics** and attempt to combine these in sensible ways.

Dale et al. [1998]

# Image to Text

- Natural OCR [Shi et al., 2016, Lee and Osindero, 2016, Mishra et al., 2012, Wang et al., 2012]



cocacola

- Image Captioning [Xu et al., 2015, Karpathy and Fei-Fei, 2015, Vinyals et al., 2015]



A man in street  
racer armor is  
examining the tire  
of another racers  
motor bike

# Image to Text

- Natural OCR [Shi et al., 2016, Lee and Osindero, 2016, Mishra et al., 2012, Wang et al., 2012]



cocacola

- Image Captioning [Xu et al., 2015, Karpathy and Fei-Fei, 2015, Vinyals et al., 2015]



A man in street  
racer armor is  
examining the tire  
of another racers  
motor bike



# IM2LATEX-100K

$$A_0^3(\alpha' \rightarrow 0) = 2g_d \varepsilon_\lambda^{(1)} \varepsilon_\mu^{(2)} \varepsilon_\nu^{(3)} \left\{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_3^\mu - p_1^\mu) + \eta^{\mu\nu} (p_2^\lambda - p_3^\lambda) \right\}.$$

A\_{0}^{(3)}(\alpha' \rightarrow 0) = 2 g\_d \varepsilon\_\lambda^{(1)} \varepsilon\_\mu^{(2)} \varepsilon\_\nu^{(3)} \left\{ \eta^{\lambda\mu} (p\_1^\nu - p\_2^\nu) + \eta^{\lambda\nu} (p\_3^\mu - p\_1^\mu) + \eta^{\mu\nu} (p\_2^\lambda - p\_3^\lambda) \right\}

# IM2LATEX-100K

$$\left\{ \begin{array}{l} \delta_\epsilon B \sim \epsilon F, \\ \delta_\epsilon F \sim \partial \epsilon + \epsilon B, \end{array} \right.$$

```
\left\{ \begin{array}{l} \delta_\epsilon B \sim \epsilon F, \\ \delta_\epsilon F \sim \partial \epsilon + \epsilon B, \end{array} \right.
```

# IM2LATEX-100K

$$\int_{\mathcal{L}_{d-1}^d} f(H) d\nu_{d-1}(H) = c_3 \int_{\mathcal{L}_2^A} \int_{\mathcal{L}_{d-1}^L} f(H) [H, A]^2 d\nu_{d-1}^L(H) d\nu_2^A(L).$$

$$\begin{aligned} & \int \limits_{\mathcal{L}_{d-1}^d} \{ \{ \mathcal{L} \}^d - (d-1) \} f(H) d\nu_{d-1}(H) = c_3 \int \limits_{\mathcal{L}_2^A} \{ \{ \mathcal{L} \}^A - 2 \} \\ & \int \limits_{\mathcal{L}_{d-1}^L} \{ \{ \mathcal{L} \}^L - (d-1) \} f(H) [H, A]^2 d\nu_{d-1}^L(H) d\nu_2^A(L). \end{aligned}$$

# IM2LATEX-100K

$$J = \begin{pmatrix} \alpha^t & \tilde{f}_2 \\ f_1 & \tilde{A} \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & L \end{pmatrix} \begin{pmatrix} \alpha & \tilde{f}_1 \\ f_2 & A \end{pmatrix} = \begin{pmatrix} \tilde{f}_2 L f_2 & \tilde{f}_2 L A \\ \tilde{A} L f_2 & \tilde{A} L A \end{pmatrix}$$

```
J = \left( \begin{array}{cc} c c & \alpha ^{ t } & \tilde{f} _ { 2 } \\ f _ { 1 } & \tilde{A} \end{array} \right) \left( \begin{array}{cc} 0 & 0 \\ 0 & L \end{array} \right) \left( \begin{array}{cc} \alpha & \tilde{f} _ { 1 } \\ f _ { 2 } & A \end{array} \right) = \left( \begin{array}{cc} \tilde{f} _ { 2 } L f _ { 2 } & \tilde{f} _ { 2 } L A \\ \tilde{A} L f _ { 2 } & \tilde{A} L A \end{array} \right)
```

# IM2LATEX-100K

$$\lambda_{n,1}^{(2)} = \frac{\partial \overline{H}_0}{\partial q_{n,0}} , \lambda_{n,j_n}^{(2)} = \frac{\partial \overline{H}_0}{\partial q_{n,j_n-1}} - \mu_{n,j_n-1} , \quad j_n = 2, 3, \dots, m_n - 1 .$$

\lambda\_{n,1}^{(2)} = \frac{\partial \overline{H}\_0}{\partial q\_{n,0}} \{ \lambda\_{n,j\_n}^{(2)} \} = \frac{\partial \overline{H}\_0}{\partial q\_{n,j\_n-1}} - \mu\_{n,j\_n-1} , \quad j\_n = 2, 3, \dots, m\_n - 1 .

# IM2LATEX-100K

$$(P_{ll'} - K_{ll'})\phi'(z_q)|\chi> = 0$$

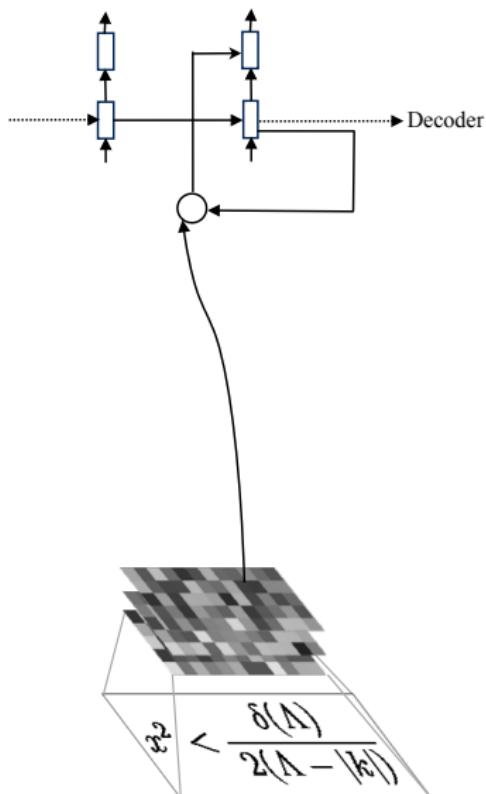
$$(P_{\{ll'\}} - K_{\{ll'\}})\phi' (z_{\{q\}})|\chi> = 0$$

# IM2LATEX-100K

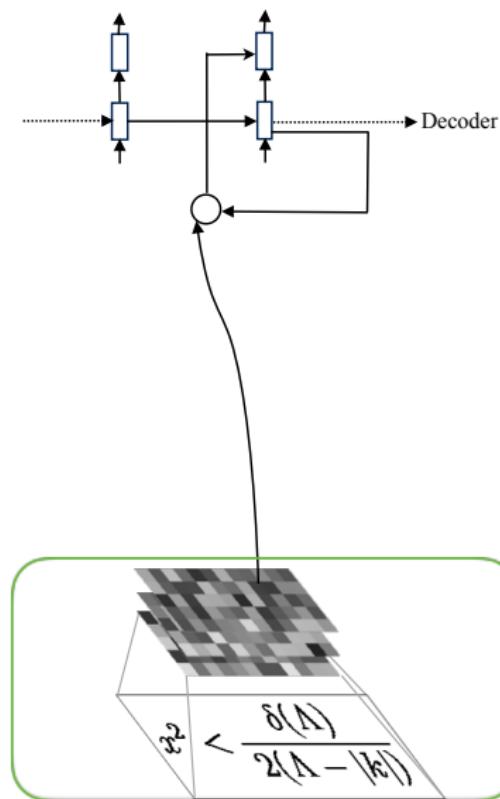
#	img size	median #char	min #char	max #char
103,556	1654×2339	98	38	997

- Originally developed for OpenAI requests for research
- LaTeX sources of arXiv papers on high energy physics from 2003 KDD cup [Gehrke et al., 2003]
- Extracted with regular expressions
- Rendered in a vanilla LaTeX environment

## Attention-based Image Captioning (Xu et al. 2015)

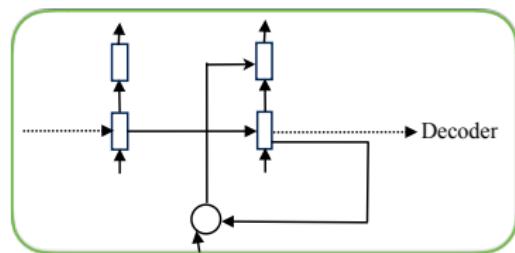


## Attention-based Image Captioning (Xu et al. 2015)



- Encoder: CNN
- Decoder: RNN with attention
- Objective: maximize log-likelihood

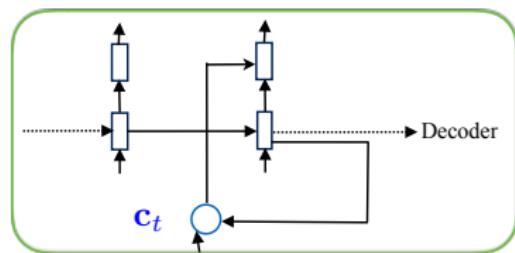
## Attention-based Image Captioning (Xu et al. 2015)



- Encoder: CNN
- Decoder: RNN with attention
- Objective: maximize log-likelihood

$$x^2 < \frac{\delta(\Lambda)}{2(\Lambda - |k|)}$$

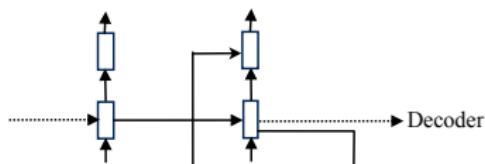
## Attention-based Image Captioning (Xu et al. 2015)



- Encoder: CNN
- Decoder: RNN with attention
- Objective: maximize log-likelihood

A 3D plot showing a function of two variables. The vertical axis is labeled  $x^2$ . The horizontal axes are labeled  $\delta(\Lambda)$  and  $2(\Lambda - |k|)$ . The surface shows a local minimum at the origin.

## Attention-based Image Captioning (Xu et al. 2015)

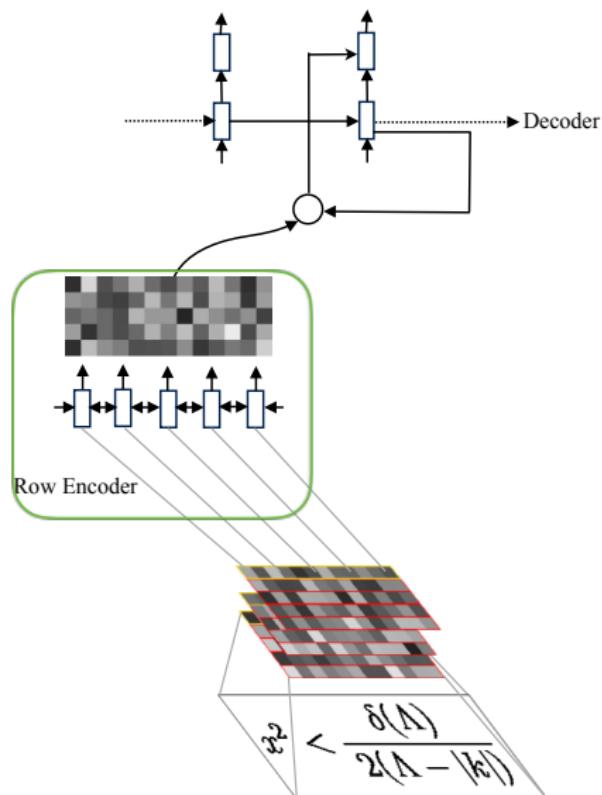


- Encoder: CNN
- Decoder: RNN with attention
- Objective: maximize log-likelihood

A diagram showing a triangular search space with a grid of points. Below the triangle, there is a mathematical expression:

$$x^2 < \frac{\delta(\Lambda)}{2(\Lambda - |k|)}$$

# Model Extensions



- Row Encoder: RNN over each row of feature map
- Parameters shared across rows
- Row embeddings to initialize RNN

## Attention

$r = \left\{ \frac{\sqrt{Q} - \sqrt{3}}{3} \right\} \cdot \{$

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

## Attention

$r = \{ \frac{\sqrt{Q} - \{ 3 \}}{l} \} \cdot \{$

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

## Attention

$$r = \frac{\sqrt{Q_3}}{l} \sin\left(\frac{l}{\sqrt{Q_3}} u\right),$$

## Coarse-to-Fine Attention

$$r = \frac{\sqrt{Q_3}}{l} \sin\left(\frac{l}{\sqrt{Q_3}} u\right),$$

## Coarse-to-Fine Attention

$$r = \frac{\sqrt{Q_3}}{l} \sin\left(\frac{l}{\sqrt{Q_3}} u\right),$$

## Coarse-to-Fine Attention

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

## Coarse-to-Fine Attention

$r \in \frac{1}{\sqrt{Q_3}} \cdot \{ \dots \} \cdot \{ \}$

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

## Coarse-to-Fine Attention

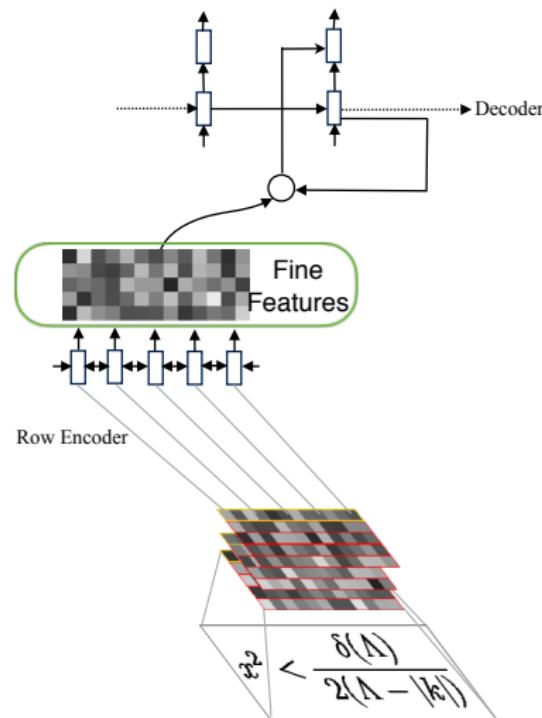
$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

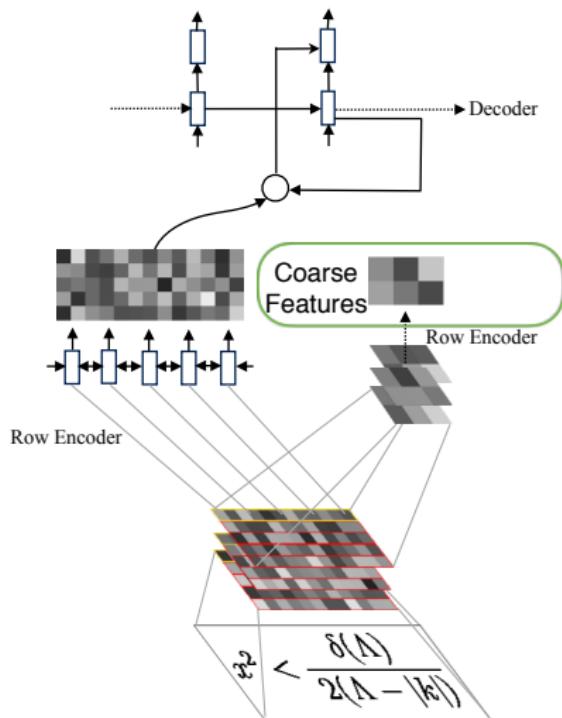
## Coarse-to-Fine Attention

$$r = \frac{\sqrt{Q_3}}{l} \sin \left( \frac{l}{\sqrt{Q_3}} u \right),$$

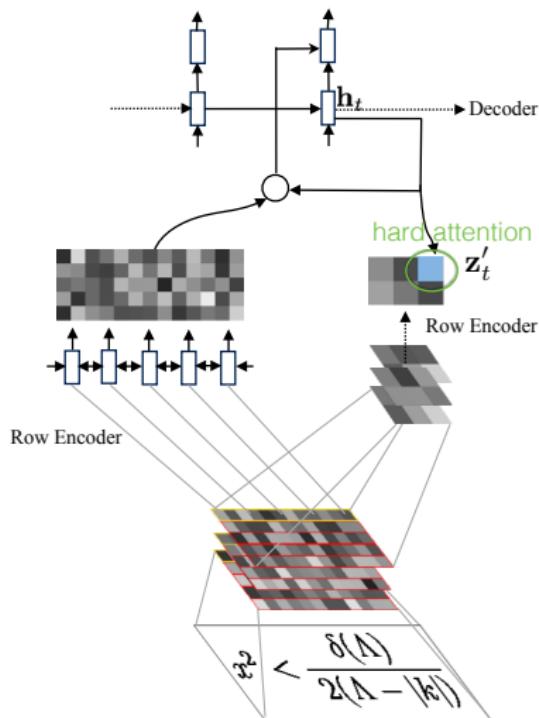
# Coarse-to-Fine Attention



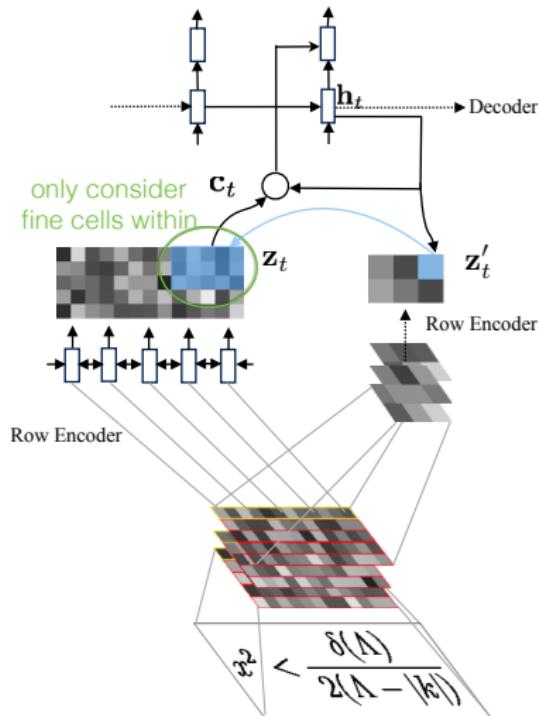
# Coarse-to-Fine Attention



# Coarse-to-Fine Attention

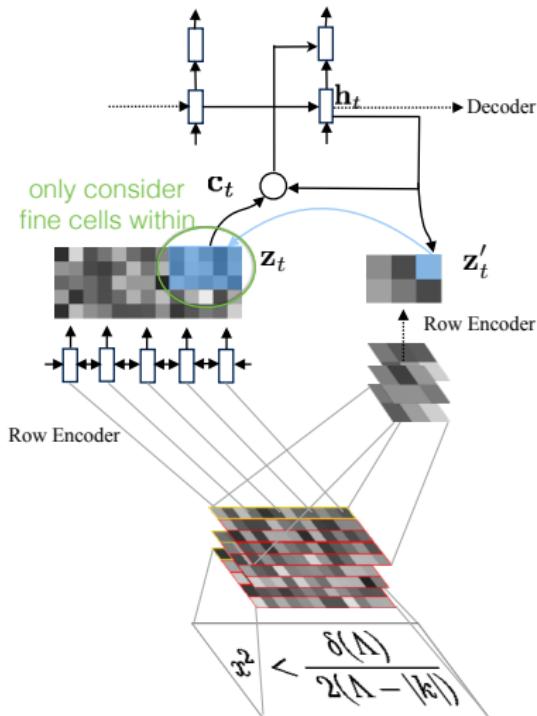


# Coarse-to-Fine Attention



$$p(z_t) = \sum_{z'_t} p(z'_t)p(z_t|z'_t)$$

# Coarse-to-Fine Attention



$$p(z_t) = \sum_{z'_t} p(z'_t)p(z_t|z'_t)$$

## Coarse-to-Fine Variants

- REINFORCE: hard attention [xu et al., 2015] to select a **single** coarse cell, the presented model
- SPARSEMAX: use sparse activation function Sparsemax [Martins and Astudillo, 2016] instead of Softmax to select **multiple** coarse cells

## Experiment Details

- Tokenization & Normalization:

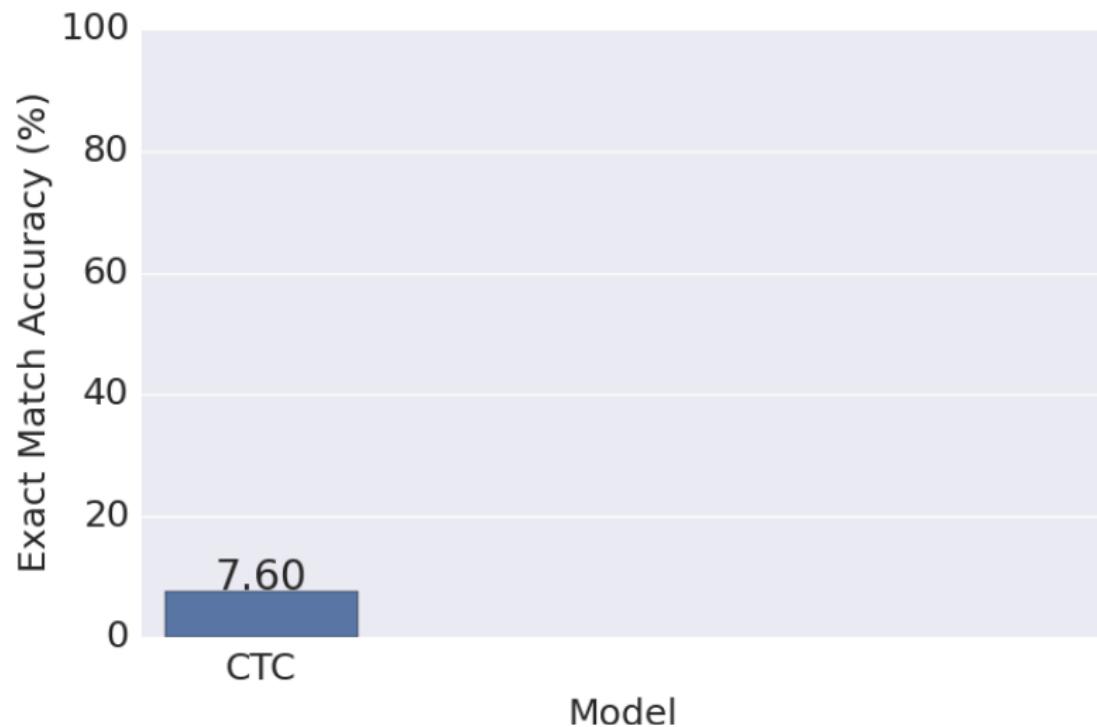
$P_{\{11'\}}^1 - K^2_{\{11\}}$



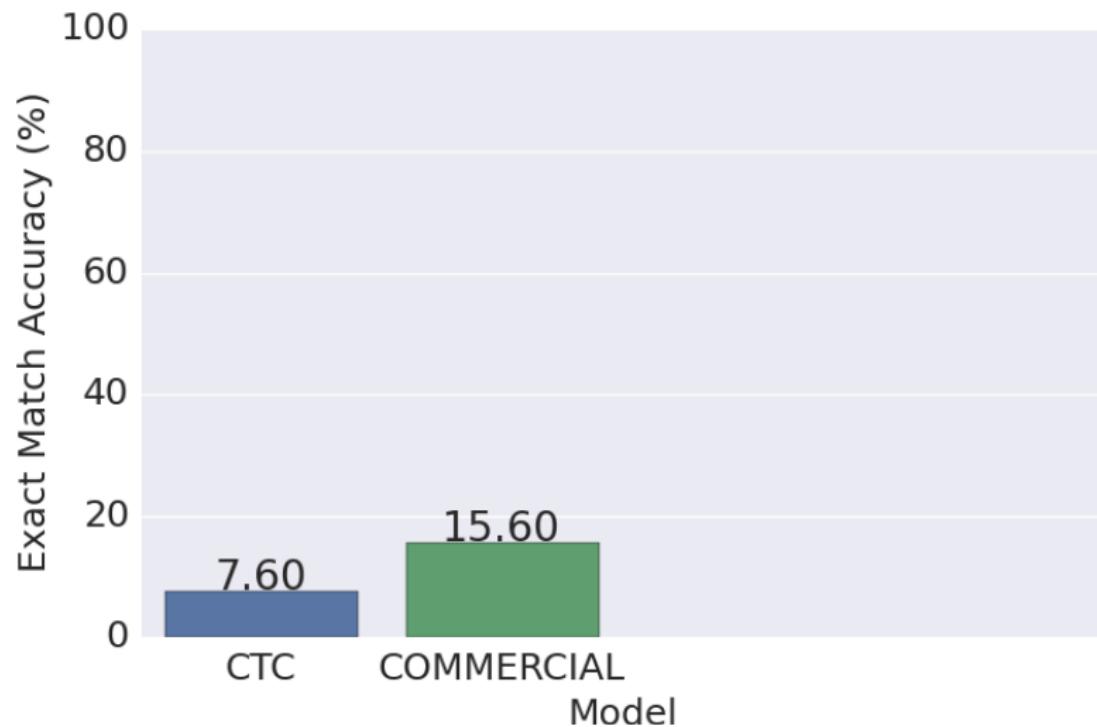
$P_{\{11'\}}^1 - K^2_{\{11\}}$

- Evaluation: exact image match accuracy (rendered prediction versus original image)
- Implementation: Torch [Collobert et al., 2011], based on OpenNMT [Klein et al., 2017]

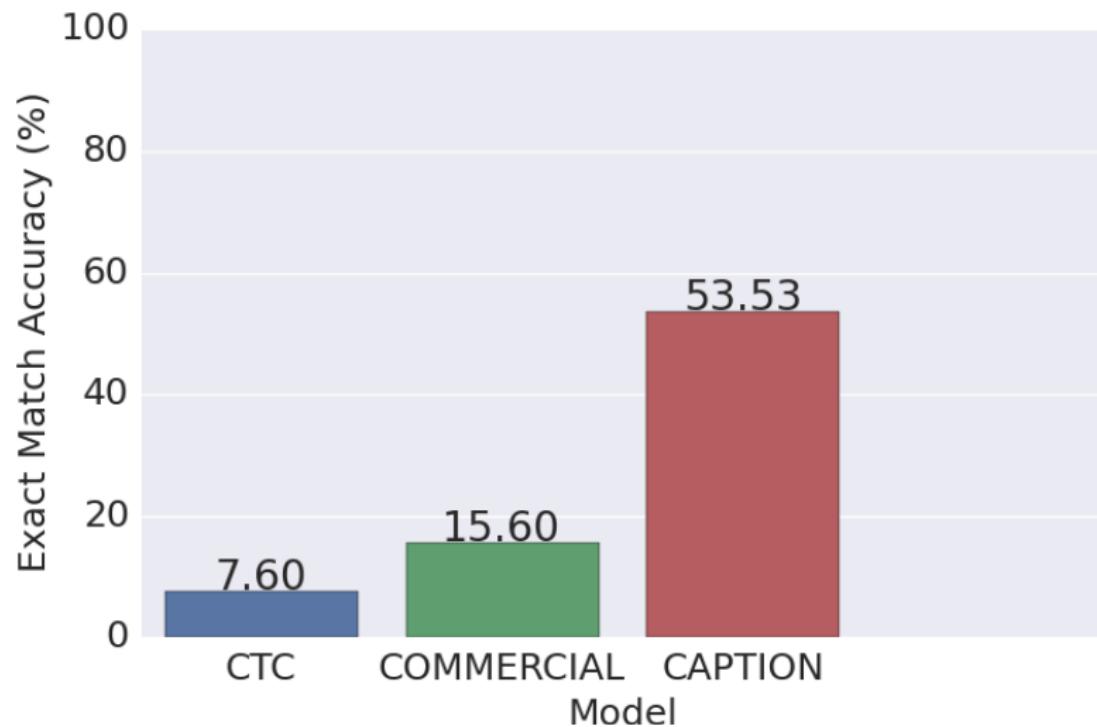
## Baseline Results



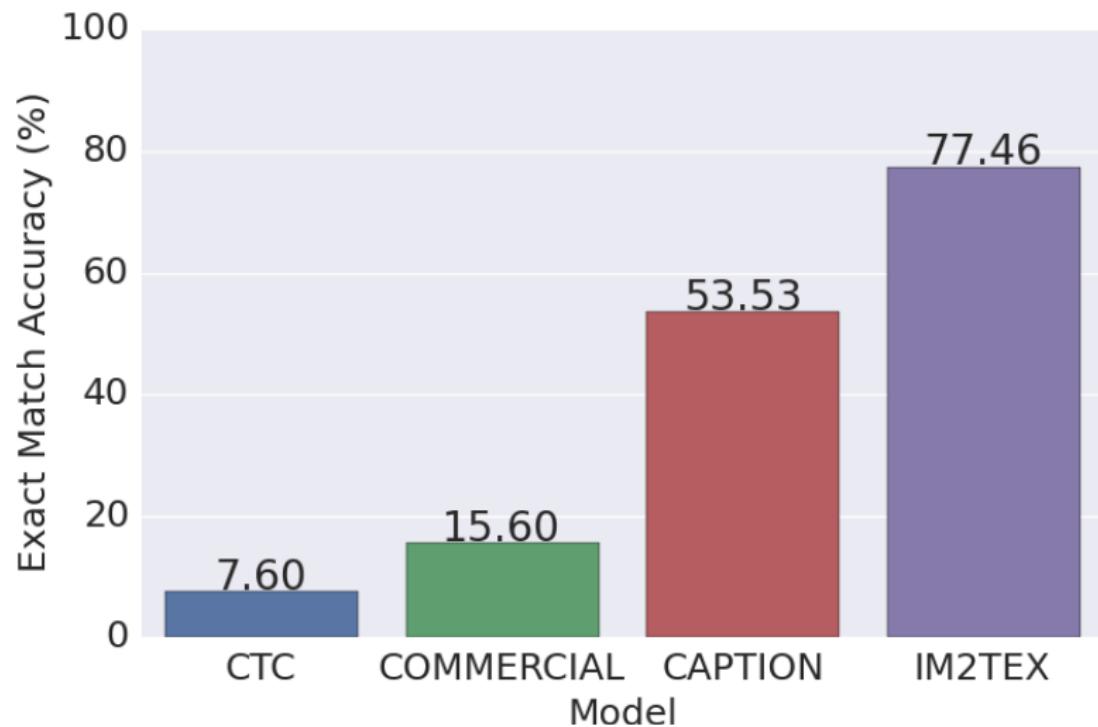
## Baseline Results



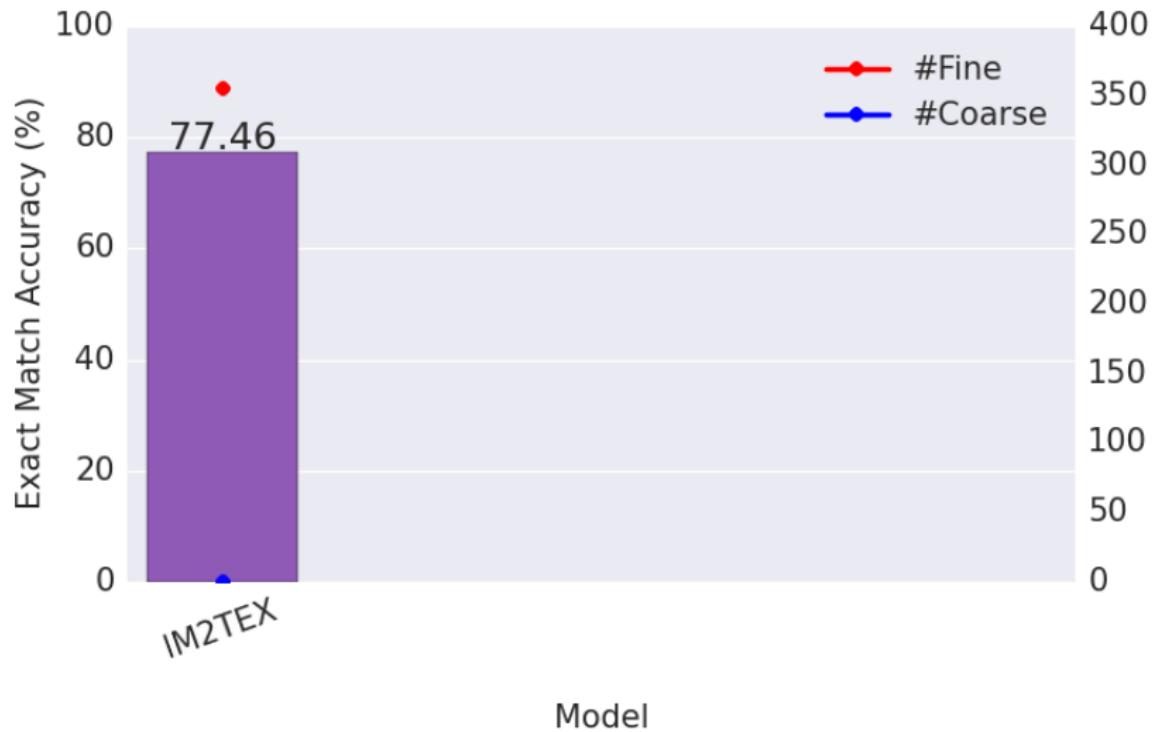
## Baseline Results



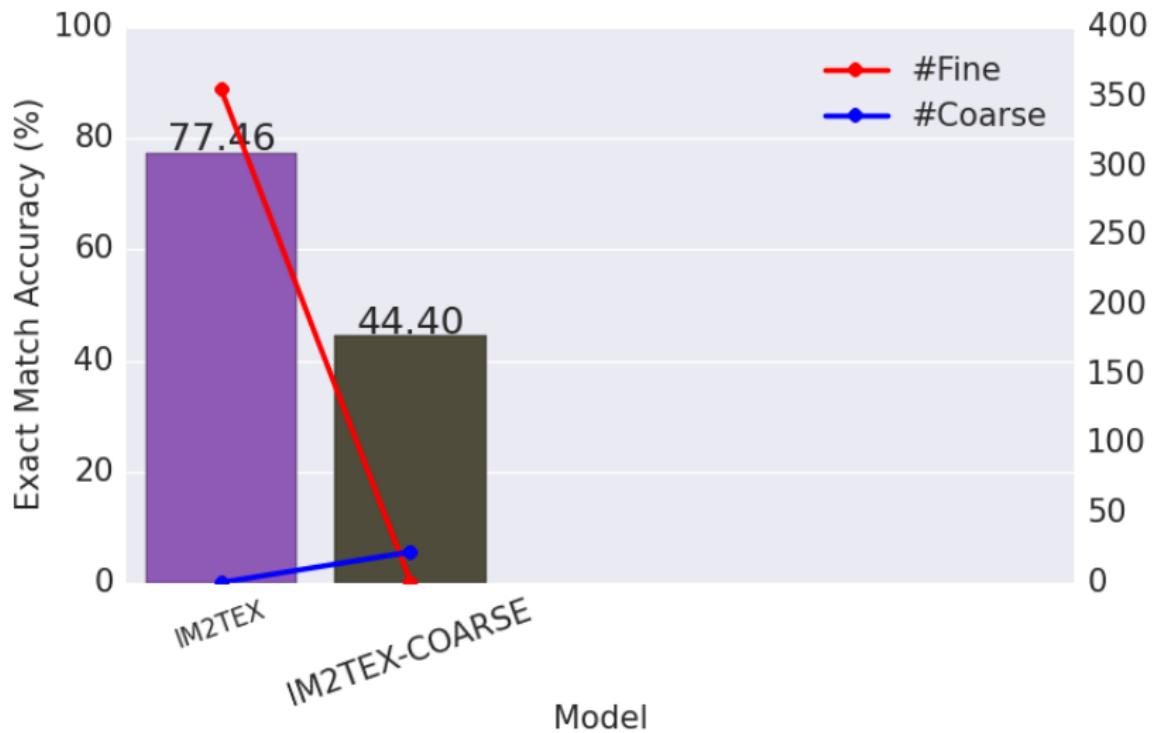
## Baseline Results



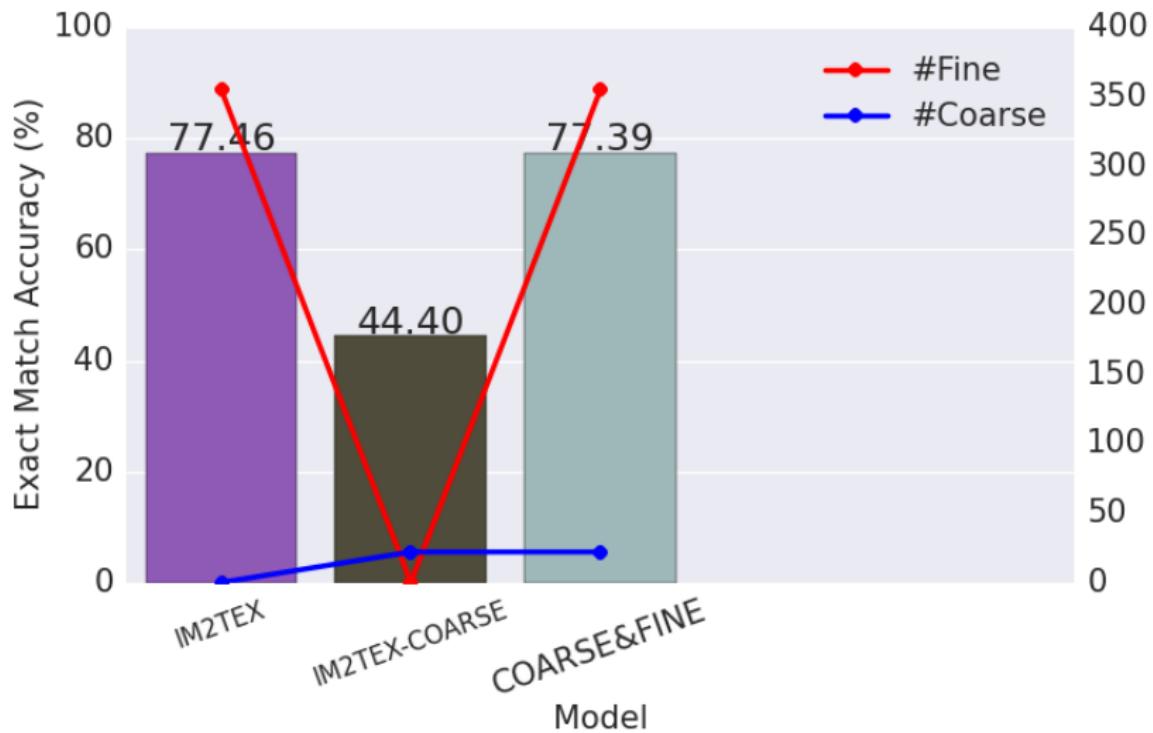
## Main Results



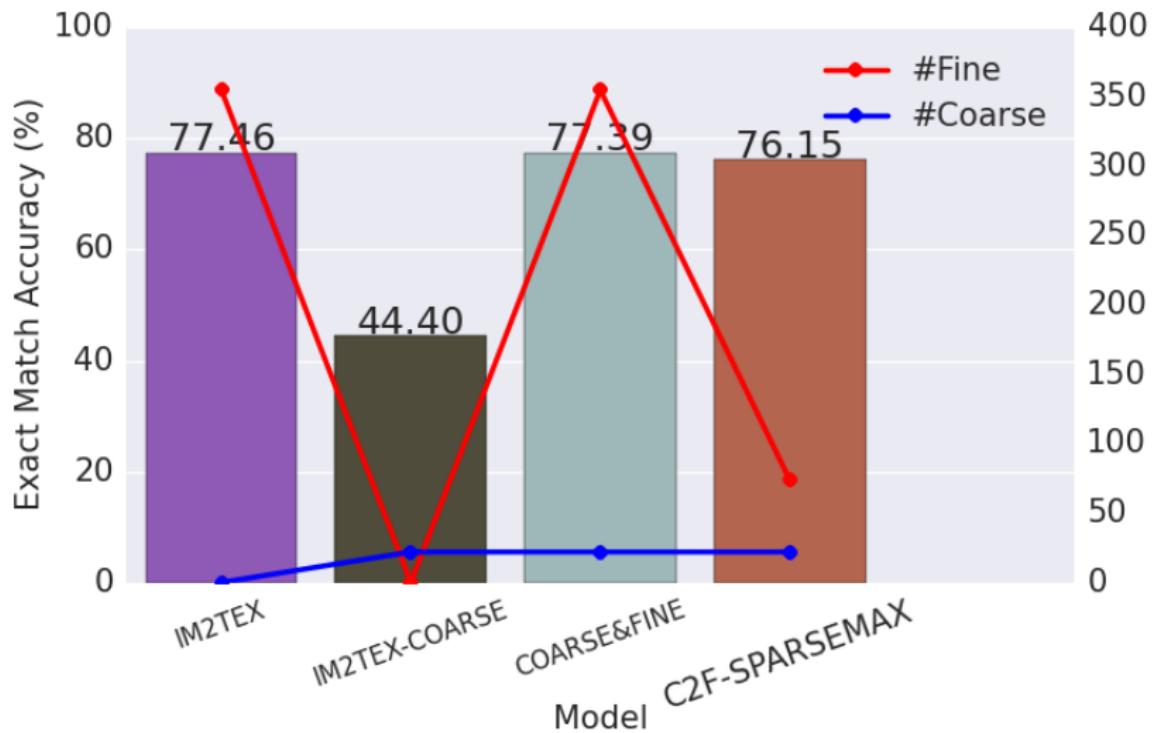
## Main Results



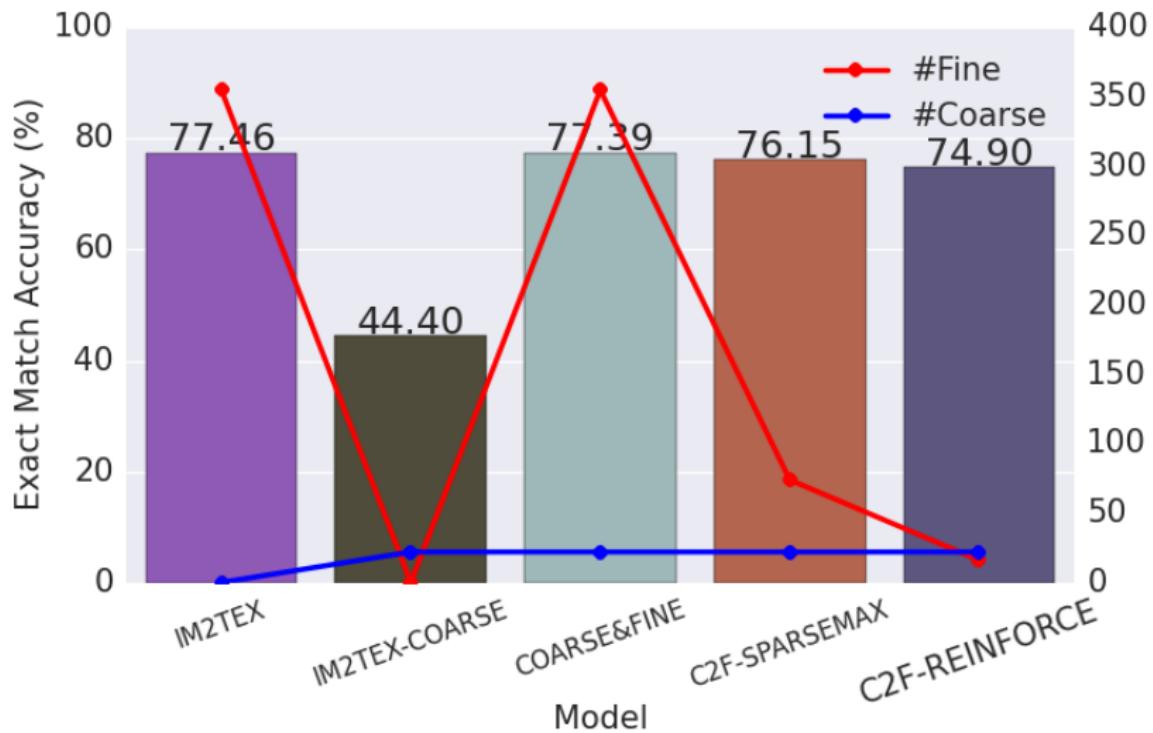
## Main Results



## Main Results



# Main Results



## Qualitative Results

$$Z = \sum_{\text{spins}} \prod_{\text{cubes}} W(a|e, f, g|b, c, d|h),$$

$$\{\Psi \circ \mu, f\} = (\overline{X}_i f) (Y^i \Psi) \circ \mu,$$

$$U_n(\theta, \phi) = \begin{pmatrix} \cos(\theta/2) & -e^{-in\phi} \sin(\theta/2) \\ \sin(\theta/2) e^{in\phi} & \cos(\theta/2) \end{pmatrix}$$

$$\sin \frac{\pi \alpha' s}{2} + \sin \frac{\pi \alpha' t}{2} + \sin \frac{\pi \alpha' u}{2} = -\frac{\pi^3}{16} \alpha'^3 stu + o(\alpha'^5),$$

$$Y(T, U) = \int_{\mathcal{F}} \frac{d^2 \tau}{\Im \tau} \Gamma_{2,2}(T, U) \left( -6 \left[ \overline{\Omega}_{\textcolor{red}{2}} - \frac{1}{8\pi \Im \tau} \right] \frac{\overline{\Omega}}{\bar{\eta}^{24}} - \frac{\bar{j}}{8} + 126 \right) ,$$

## Handwritten Formulas

- Synthetic 100K handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune on CROHME 13, 14 and got competitive results

## Handwritten Formulas

- Synthetic 100K handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
  - Finetune on CROHME 13, 14 and got competitive results

$A_0^3(\alpha' \rightarrow 0) = 2g_d\,\varepsilon_{\lambda}^{(1)}\varepsilon_{\mu}^{(2)}\varepsilon_{\nu}^{(3)}\left\{\eta^{\lambda\mu}\left(p_1^{\nu}-p_2^{\nu}\right)+\eta^{\lambda\nu}\left(p_3^{\mu}-p_1^{\mu}\right)+\eta^{\mu\nu}\left(p_2^{\lambda}-p_3^{\lambda}\right)\right\}.$ $\text{A}_{\{0\}}^{\{3\}}(\text{alpha}^{\prime }\rightarrow 0)=2 g_{\{d\}}\backslash \backslash \varepsilon \! \! \! \varepsilon ^{(1)}_{\lambda }\varepsilon \! \! \! \varepsilon ^{(2)}_{\mu }\varepsilon \! \! \! \varepsilon ^{(3)}_{\nu }\left\{\eta ^{\lambda \mu }\left(p_1^{\nu }-p_2^{\nu }\right)+\eta ^{\lambda \nu }\left(p_3^{\mu }-p_1^{\mu }\right)+\eta ^{\mu \nu }\left(p_2^{\lambda }-p_3^{\lambda }\right)\right\}.$	$\begin{cases} \delta_\epsilon B & \sim \epsilon F\,, \\ \delta_\epsilon F & \sim \partial \epsilon + \epsilon B\,, \end{cases}$
$\text{A}_{\{0\}}^{\{3\}}(\text{alpha}^{\prime }\rightarrow 0)=2 g_{\{d\}}\backslash \backslash \varepsilon \! \! \! \varepsilon ^{(1)}_{\lambda }\varepsilon \! \! \! \varepsilon ^{(2)}_{\mu }\varepsilon \! \! \! \varepsilon ^{(3)}_{\nu }\left\{\eta ^{\lambda \mu }\left(p_1^{\nu }-p_2^{\nu }\right)+\eta ^{\lambda \nu }\left(p_3^{\mu }-p_1^{\mu }\right)+\eta ^{\mu \nu }\left(p_2^{\lambda }-p_3^{\lambda }\right)\right\}.$	$\text{Left}\{\text{begin}[array]{lcl}\text{delta}_{\_}\text{epsilon} & \text{B} & \text{\&} \text{\sim} \text{\&} \text{epsilon}\\ \text{\&} \text{\delta}_{\epsilon} \text{F} & \text{\sim} & \text{\partial} \epsilon + \epsilon \text{B}\end{array}\text{right}.$
$\int\limits_{\mathcal{L}_{d-1}^d}f(H)d\nu_{d-1}(H)=c_3\int\limits_{\mathcal{L}_2^d}\int\limits_{\mathcal{L}_{d-1}^L}f(H)[H,A]^2d\nu_{d-1}^L(H)d\nu_2^A(L).$ $\text{Int}\text{ }\text{limits}_{\_}\{(\text{cal L})^{\{d\}}\text{-(d-1)}\}\text{f(H)}\text{d}\backslash \text{nu}_{\_}\{d-1\}\text{(H)}=\text{c}_{\{3\}}\text{ }\text{int}\text{ }\text{limits}_{\_}\{(\text{cal L})^{\{A\}}\text{-(2)}\}\text{ }\text{int}\text{ }\text{limits}_{\_}\{(\text{cal L})^{\{L\}}\text{-(d-1)}\}\text{f(H)}\text{[H,A]}^{\{2\}}\text{d}\backslash \text{nu}_{\_}\{d-1\}\text{^}\{L\}\text{(H)}\text{d}\backslash \text{nu}_{\_}\{2\}\text{^}\{A\}\text{(L)}.$	$J=\left(\begin{array}{cc}\alpha^t&\bar f_2\\f_1&\bar A\end{array}\right)\left(\begin{array}{cc}0&0\\0&L\end{array}\right)\left(\begin{array}{cc}\alpha&\bar f_1\\f_2&A\end{array}\right)=\left(\begin{array}{cc}\bar f_2Lf_2&\bar f_2LA\\\bar ALf_2&\bar ALA\end{array}\right)$ $\text{J}=\text{Left}\{\text{begin}[array]{cc}\text{alpha}^t & \bar f_2 \\ f_1 & \bar A\end{array}\text{right}\}\text{ }\text{left}\{\text{begin}[array]{cc}0 & 0 \\ 0 & L\end{array}\text{right}\}\text{ }\text{left}\{\text{begin}[array]{cc}\alpha & \bar f_1 \\ f_2 & A\end{array}\text{right}\}=\text{Left}\{\text{begin}[array]{cc}\bar f_2Lf_2 & \bar f_2LA \\ \bar ALf_2 & \bar ALA\end{array}\text{right}\}$

# Handwritten Formulas

- Synthetic 100K handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune on CROHME 13, 14 and got competitive results

$A_0^3(\alpha' \rightarrow 0) = 2g_d \varepsilon_A^{(1)} \varepsilon_V^{(2)} \varepsilon_V^{(3)} \{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_3^\nu - p_1^\nu) + \eta^{\lambda\nu} (p_2^\nu - p_3^\nu) \}$ <p><math display="block">(A_{\{0\}}^3 \{ \alpha \wedge (\prime) \rightarrow 0 \ }) = 2g_d d \varepsilon_A^{(1)} \varepsilon_V^{(2)} \varepsilon_V^{(3)} \{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_3^\nu - p_1^\nu) + \eta^{\lambda\nu} (p_2^\nu - p_3^\nu) \}</math></p> <p><math display="block">\left( \begin{array}{l} \delta_\epsilon B \sim \epsilon F, \\ \delta_\epsilon F \sim \partial \epsilon + \epsilon B, \end{array} \right)</math></p>	$\left( \begin{array}{l} \delta_\epsilon B \sim \epsilon F, \\ \delta_\epsilon F \sim \partial \epsilon + \epsilon B, \end{array} \right)$
$\int_{L_{d-1}^A} f(H) d\nu_{d-1}(H) = c_3 \int \int f(H)  H, A ^2 d\nu_{d-1}^L(H) d\nu_d^A(L)$ <p><math display="block">\int \int \int f(H)  H, A ^2 d\nu_{d-1}^L(H) d\nu_d^A(L) d\nu_{d-1}^A(L)</math></p> <p><math display="block">\int \int \int f(H)  H, A ^2 d\nu_{d-1}^L(H) d\nu_d^A(L) d\nu_{d-1}^A(L)</math></p>	$J = \begin{pmatrix} \alpha & f_2 \\ f_1 & \tilde{A} \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha & f_1 \\ f_2 & A \end{pmatrix} = \begin{pmatrix} f_2 L f_2 & f_2 L A \\ \tilde{A} L f_2 & \tilde{A} L A \end{pmatrix}$ <p><math display="block">J = \begin{pmatrix} \alpha &amp; f_2 \\ f_1 &amp; \tilde{A} \end{pmatrix} \begin{pmatrix} 0 &amp; 0 \\ 0 &amp; 1 \end{pmatrix} \begin{pmatrix} \alpha &amp; f_1 \\ f_2 &amp; A \end{pmatrix} = \begin{pmatrix} f_2 L f_2 &amp; f_2 L A \\ \tilde{A} L f_2 &amp; \tilde{A} L A \end{pmatrix}</math></p> <p><math display="block">J = \begin{pmatrix} \alpha &amp; f_2 \\ f_1 &amp; \tilde{A} \end{pmatrix} \begin{pmatrix} 0 &amp; 0 \\ 0 &amp; 1 \end{pmatrix} \begin{pmatrix} \alpha &amp; f_1 \\ f_2 &amp; A \end{pmatrix} = \begin{pmatrix} f_2 L f_2 &amp; f_2 L A \\ \tilde{A} L f_2 &amp; \tilde{A} L A \end{pmatrix}</math></p>

## Handwritten Formulas

- Synthetic 100K handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
  - Finetune on CROHME 13, 14 and got competitive results

$$A_0^3(\alpha' \rightarrow 0) = 2g_d \varepsilon_A^{(1)} \varepsilon_V^{(2)} \varepsilon_V^{(3)} \{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_2^\mu - p_1^\mu) + \eta^{\mu\nu} (p_2^\lambda - p_3^\lambda) \}$$

$$\begin{cases} \delta_\epsilon B \sim \epsilon F, \\ \delta_\epsilon F \sim \partial F + \epsilon B. \end{cases}$$

$$\left. \begin{array}{l} \delta_{\epsilon} B \sim \epsilon F \\ \delta_{\epsilon} F \sim \partial \epsilon + \epsilon B \end{array} \right\}$$

$$\int_{\mathcal{L}_{d-1}^A} f(H) d\nu_{d-1}(H) = c_3 \int \int f(H) [H, A]^2 d\nu_{d-1}^L(H) d\nu_2^A(L)$$

$$J = \begin{pmatrix} \alpha^t & f_2 \\ 0 & \tilde{f}_2 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} \alpha & f_1 \\ 0 & \tilde{f}_1 \end{pmatrix} = \begin{pmatrix} f_2 L f_2 & f_2 L A \\ \tilde{f}_2 L \tilde{f}_2 & \tilde{f}_2 L A \end{pmatrix}$$

```
\int \limits_{\{(\text{cal L})^{\wedge}(\text{d})\}_{-}(\text{d}-1)} \text{f}(\text{H}) \text{d}\nu_{-}(\text{d}-1)(\text{H}) = c_{-}\{3\} \int \limits_{\{(\text{cal L})^{\wedge}(\text{A})\}_{-}(\text{2})} \text{f}(\text{H}) \text{d}\nu_{-}(\text{d}-1)^{\wedge}(\text{L})(\text{H}) \text{d}\nu_{-}(\text{2})^{\wedge}(\text{A})(\text{L}).
```

```
J=\left(\begin{array}{cc}\alpha & \tilde{f}_1 \\ \tilde{f}_2 & 0\end{array}\right) = \left(\begin{array}{cc}\alpha & \tilde{f}_1 \\ \tilde{f}_2 & 0\end{array}\right)\left(\begin{array}{cc}0 & 0 \\ 0 & 1\end{array}\right) = \left(\begin{array}{cc}\alpha & \tilde{f}_1 \\ \tilde{f}_2 & 0\end{array}\right)
```

$$\lambda_{n,j}^{(2)} = \frac{\partial \bar{H}_0}{\partial q_{m_0}}, \quad \lambda_{n,j,n}^{(2)} = \frac{\partial \bar{H}_0}{\partial q_{m_{j-1}}} - \mu_{n,j,n-1}, \quad j_n = 2, 3, \dots, m_n - 1$$

$$(P_{W'} - K_{W'})\phi'(z_0)|x> = 0$$

```
\lambda_{n,1}^{(2)}=\frac{\partial\overline{H}_0}{\partial q_{n,0}}, \lambda_{n,j,n}^{(2)}=\frac{\partial\overline{H}_0}{\partial q_{n,j-1}}-\mu_{n,j-1}, j=2,3,\dots,m-1.
```

$$(P_{\{q\}} - K_{\{q\}}) \varphi'(z_{\{q\}}) \chi \geq 0$$

## Conclusions & Future Work

- The constructed dataset IM2LATEX-100K is rich structured and challenging
- A case study of multi-modal document recognition/generation
- Coarse-to-fine attention can be applied to other tasks

# References |

- R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A matlab-like environment for machine learning. In BigLearn, NIPS Workshop, number EPFL-CONF-192376, 2011.
- R. Dale, D. Scott, and B. Di Eugenio. Introduction to the special issue on natural language generation. Computational Linguistics, 24(3):346–353, 1998.
- J. Gehrke, P. Ginsparg, and J. Kleinberg. Overview of the 2003 kdd cup. ACM SIGKDD Explorations Newsletter, 5(2):149–151, 2003.
- A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In Proceedings of the 23rd international conference on Machine learning, pages 369–376. ACM, 2006.
- A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3128–3137, 2015.
- D. Kirsch. Detexify: Erkennung handgemalter LaTeX-symbole. PhD thesis, Diploma thesis, Westfälische Wilhelms-Universität Münster, 10 2010.[Online]. Available: <http://danielkirs.ch/thesis.pdf>, 2010.
- G. Klein, Y. Kim, Y. Deng, J. Senellart, and A. M. Rush. Opennmt: Open-source toolkit for neural machine translation. arXiv preprint arXiv:1701.02810, 2017.
- C.-Y. Lee and S. Osindero. Recursive recurrent nets with attention modeling for ocr in the wild. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2231–2239, 2016.
- A. Martins and R. Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In International Conference on Machine Learning, pages 1614–1623, 2016.
- A. Mishra, K. Alahari, and C. Jawahar. Scene text recognition using higher order language priors. In BMVC 2012-23rd British Machine Vision Conference. BMVA, 2012.

## References II

- B. Shi, X. Bai, and C. Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2016.
- M. Suzuki, F. Tamari, R. Fukuda, S. Uchida, and T. Kanahori. Infty: an integrated ocr system for mathematical documents. In *Proceedings of the 2003 ACM symposium on Document engineering*, pages 95–104. ACM, 2003.
- O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164, 2015.
- T. Wang, D. J. Wu, A. Coates, and A. Y. Ng. End-to-end text recognition with convolutional neural networks. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 3304–3308. IEEE, 2012.
- K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*, pages 2048–2057, 2015.

## Q & A

- More visualizations:

<http://lstm.seas.harvard.edu/latex/>

- Source code (part of OpenNMT-lua and OpenNMT-py):

<http://opennmt.net/OpenNMT/applications/#image-to-text>

<http://opennmt.net/OpenNMT-py/im2text.html>