

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

Yuntian Deng¹ Anssi Kanervisto² Jeffrey Ling¹
Alexander M. Rush¹

¹Harvard University

²University of Eastern Finland

- 1 Introduction: Image-to-Markup Generation
- 2 Dataset: IM2LATEX-100K
- 3 Model
- 4 Experiments
- 5 Conclusions & Future Work

*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]

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



cocacola

- 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

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

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

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

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

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

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

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

$J = \left(\begin{array}{cc} \alpha^t & \tilde{f}_2 \\ f_1 & \tilde{A} \end{array} \right) \left(\begin{array}{cc} 0 & 0 \\ 0 & L \end{array} \right) \left(\begin{array}{cc} \alpha & \tilde{f}_1 \\ f_2 & A \end{array} \right) = \left(\begin{array}{cc} \tilde{f}_2 L f_2 & \tilde{f}_2 L A \\ \tilde{A} L f_2 & \tilde{A} L A \end{array} \right)$

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

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

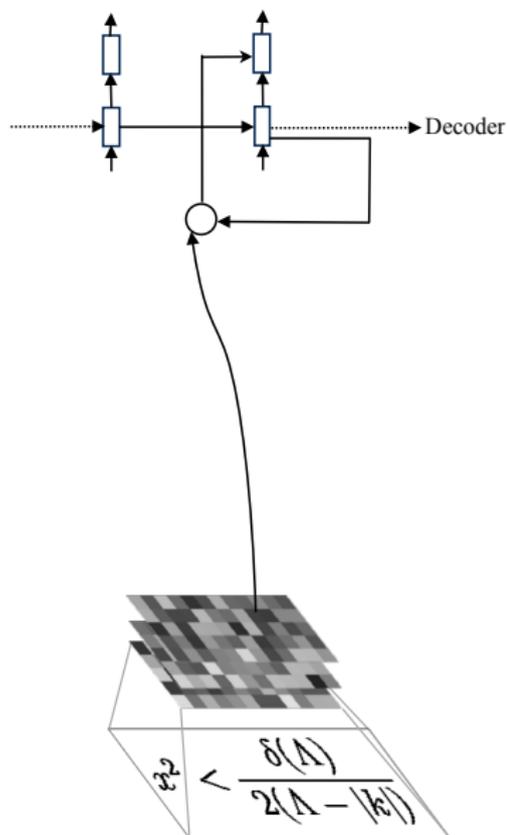
$$(P_{w'} - K_{w'})\phi'(z_q)|_{\chi} \geq 0$$

$$(P_{\{l\}} - K_{\{l\}})\phi'(z_q)|_{\chi} \geq 0$$

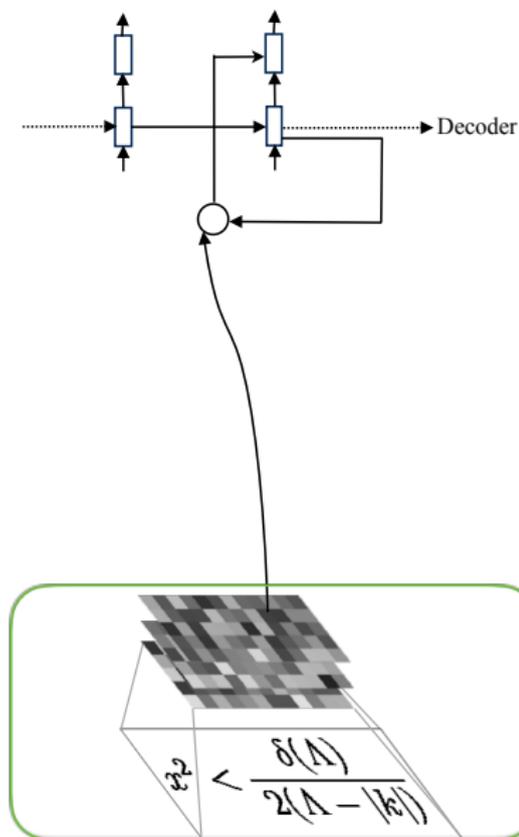
#	img size	median #char	min #char	max #char
103,556	1654×2339	98	38	997

- 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)

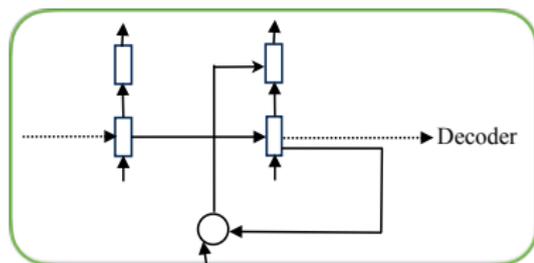


Attention-based Image Captioning (Xu et al. 2015)

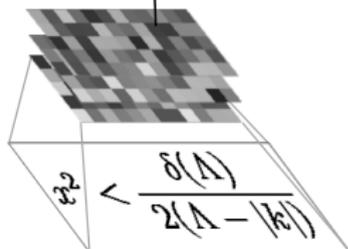


- Encoder: CNN

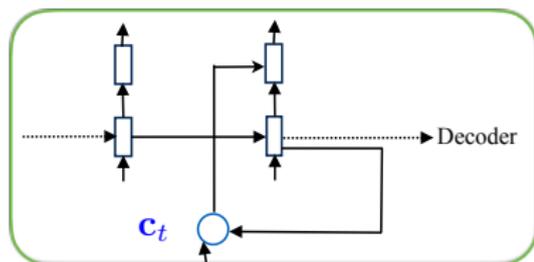
Attention-based Image Captioning (Xu et al. 2015)



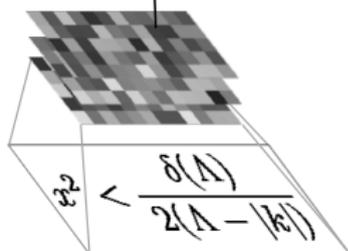
- Encoder: CNN
- Decoder: RNN with attention



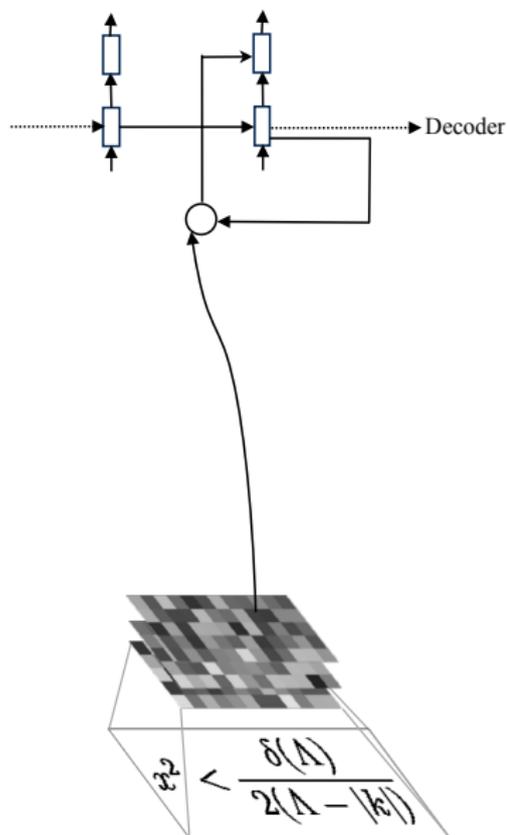
Attention-based Image Captioning (Xu et al. 2015)



- Encoder: CNN
- Decoder: RNN with attention

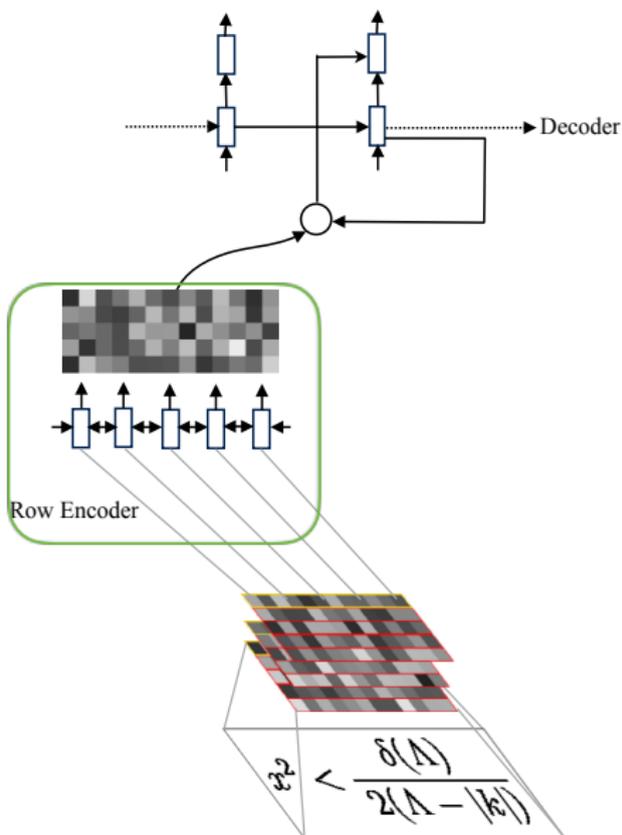


Attention-based Image Captioning (Xu et al. 2015)



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

Model Extensions



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

Attention

`r = { \frac {\sqrt{Q_3}}{l} }`

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

Attention

$r = \left\{ \frac{\sqrt{Q_3}}{l} \right\}$

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

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

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

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

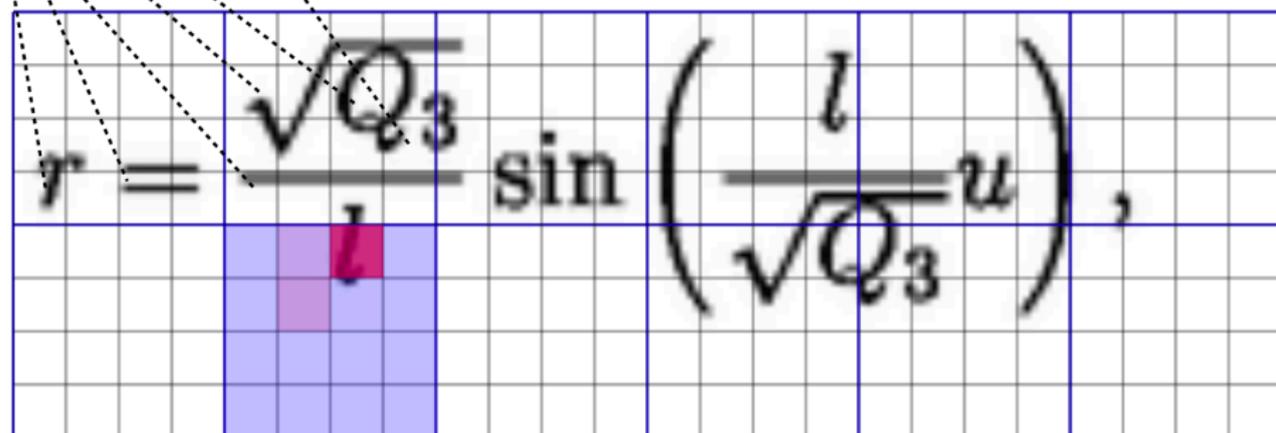
Coarse-to-Fine Attention

$$r = \frac{\sqrt{Q_3}}{l}$$

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

Coarse-to-Fine Attention

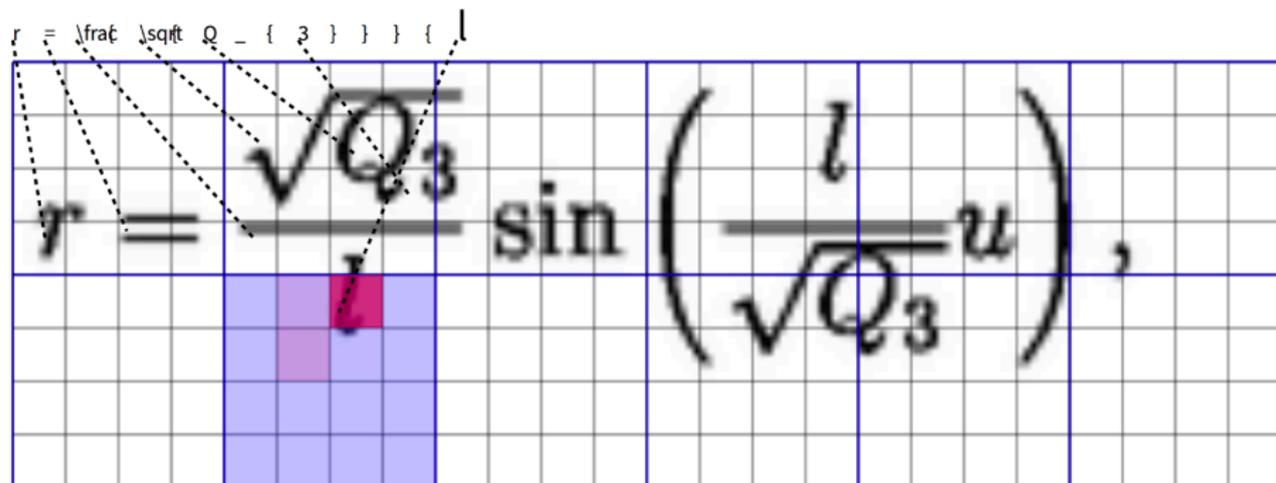
$$r = \frac{\sqrt{Q_3}}{l}$$



The diagram shows a grid with a 3x3 region highlighted in light blue. The center cell of this region is highlighted in red and contains the letter 'l'. Dotted lines connect the 'l' in the red cell to the 'l' in the denominator of the equation above and to the 'l' in the numerator of the sine function's argument.

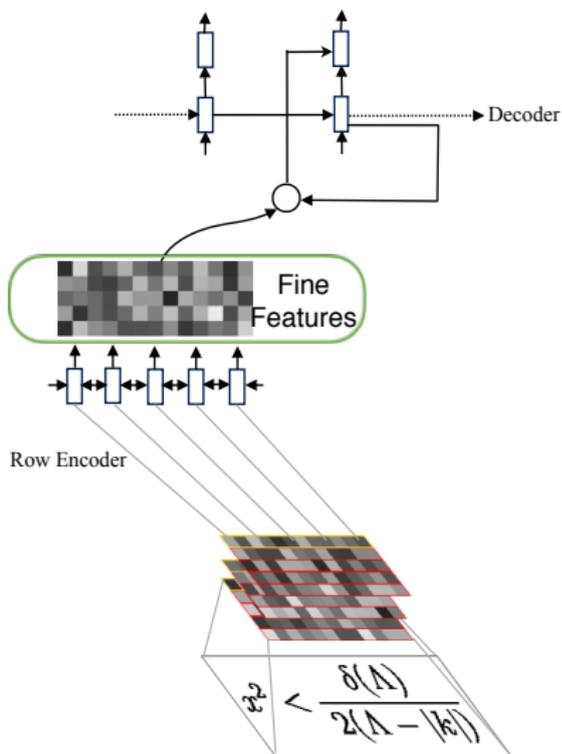
$$r = \frac{\sqrt{Q_3}}{l} \sin\left(\frac{l}{\sqrt{Q_3}}u\right),$$

Coarse-to-Fine Attention

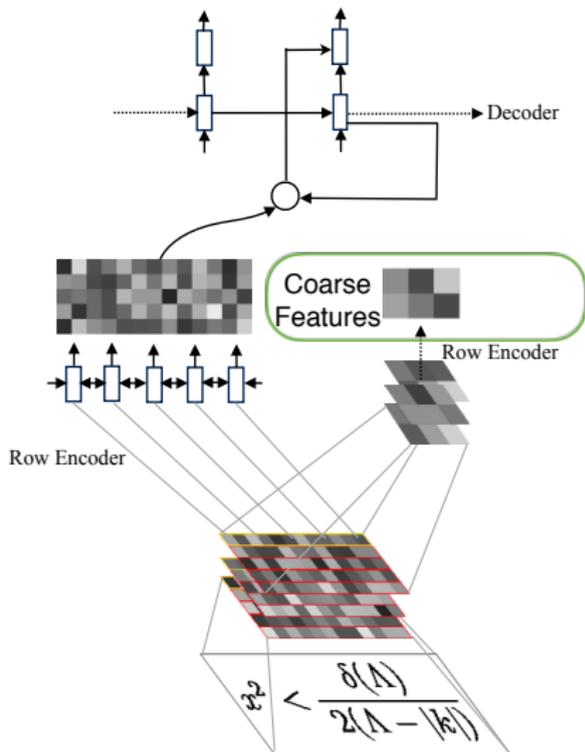
$$r = \frac{\sqrt{Q_3}}{l} \sin\left(\frac{l}{\sqrt{Q_3}}u\right),$$


The diagram illustrates the Coarse-to-Fine Attention mechanism. A grid of colored squares (purple and red) is shown below the equation. Dotted lines connect the variables in the equation to the grid: 'r' points to the red square, 'sqrt(Q3)' points to the purple squares, 'l' points to the red square, and 'u' points to the red square.

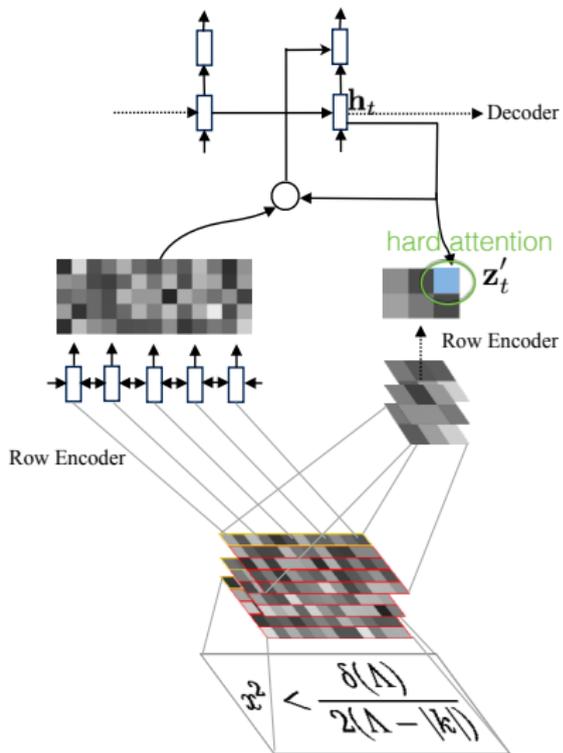
Coarse-to-Fine Attention



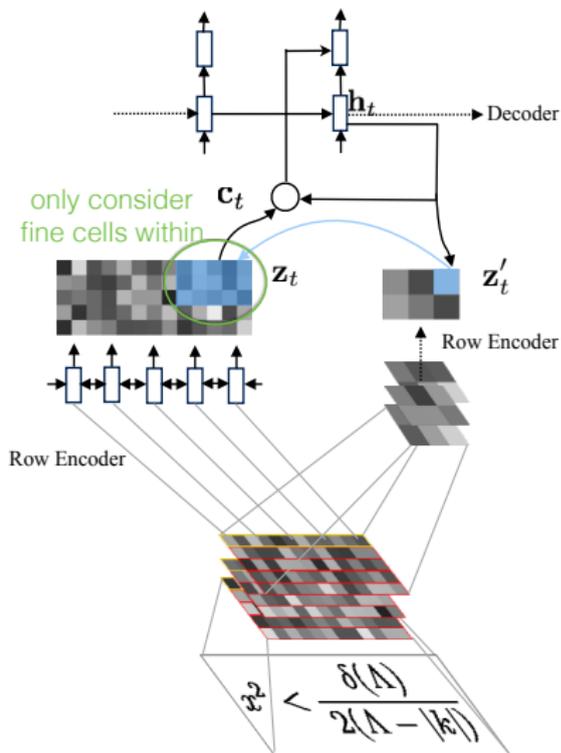
Coarse-to-Fine Attention



Coarse-to-Fine Attention

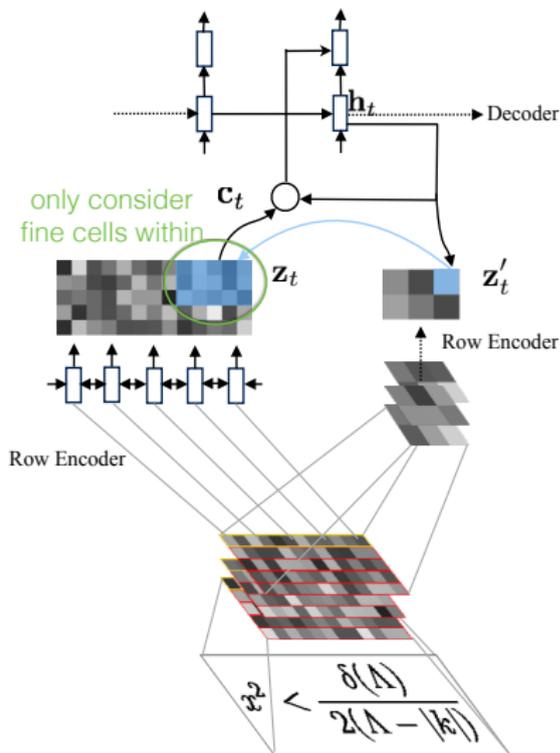


Coarse-to-Fine Attention



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

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

- Tokenization & Normalization:

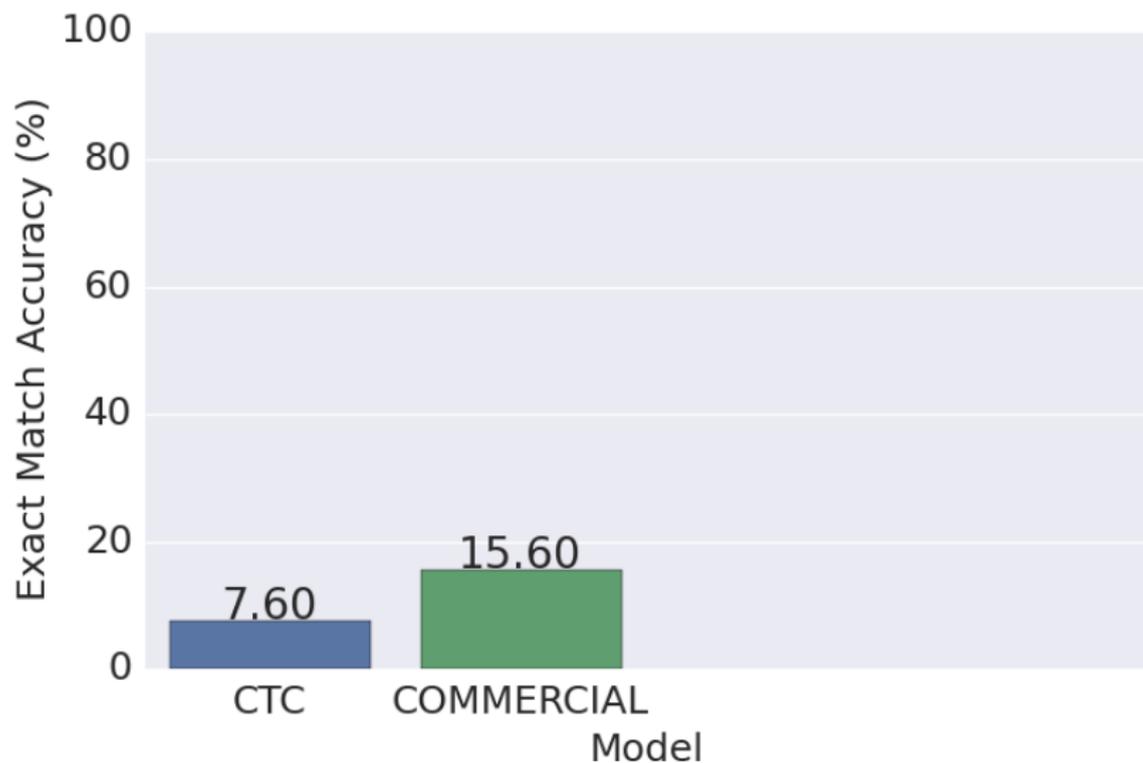
$$P_{\{11'\}}^{1-K^2_{\{11\}}}$$
$$\Downarrow$$
$$P_{\{11\}}^{\{\prime\}}^{\{1\}} - K_{\{11\}}^{\{2\}}$$

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

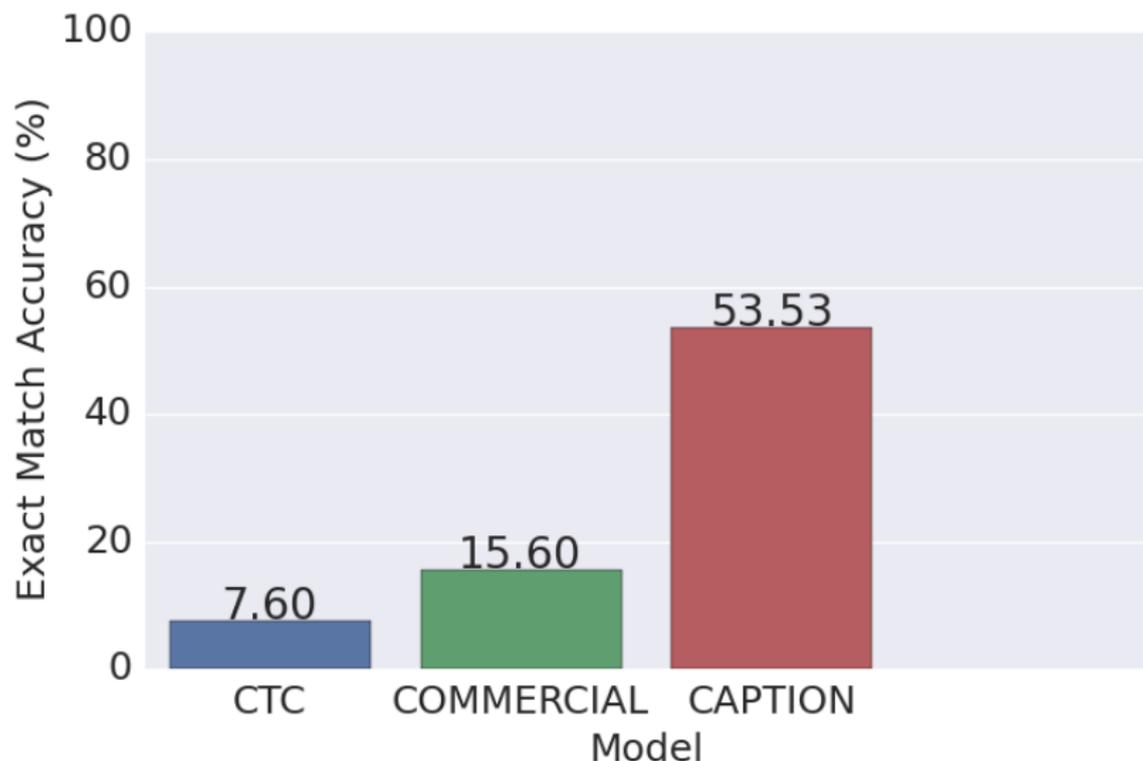
Baseline Results



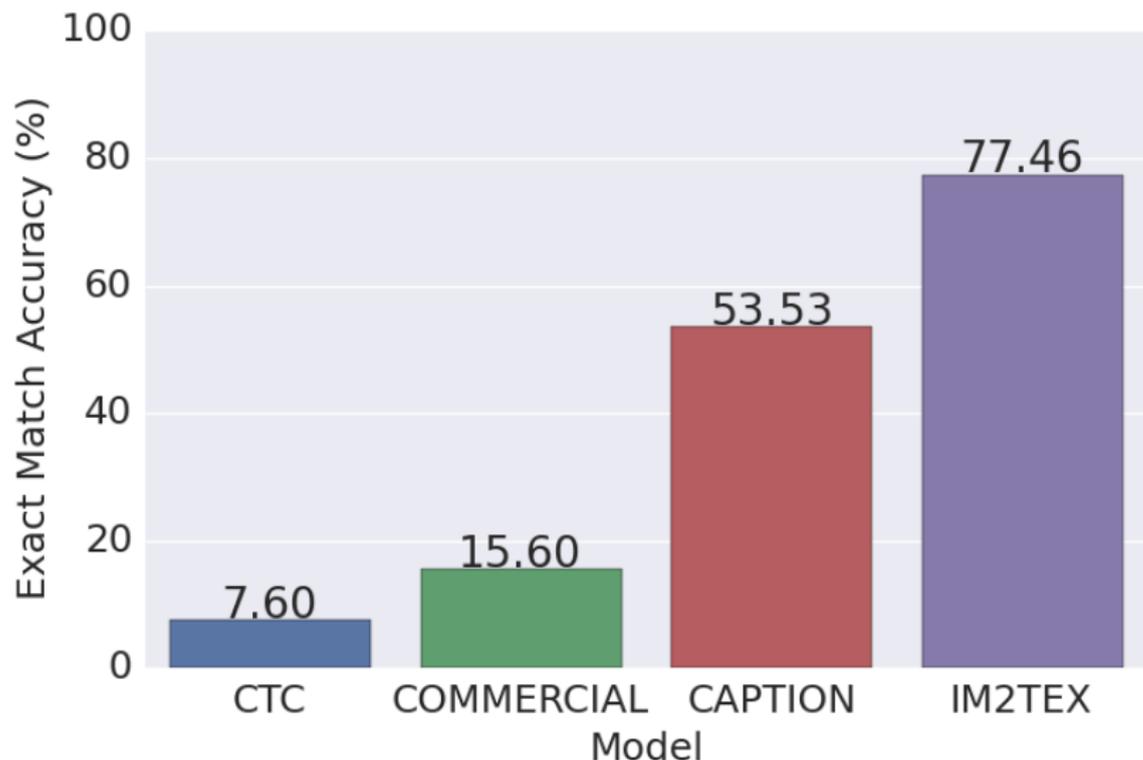
Baseline Results



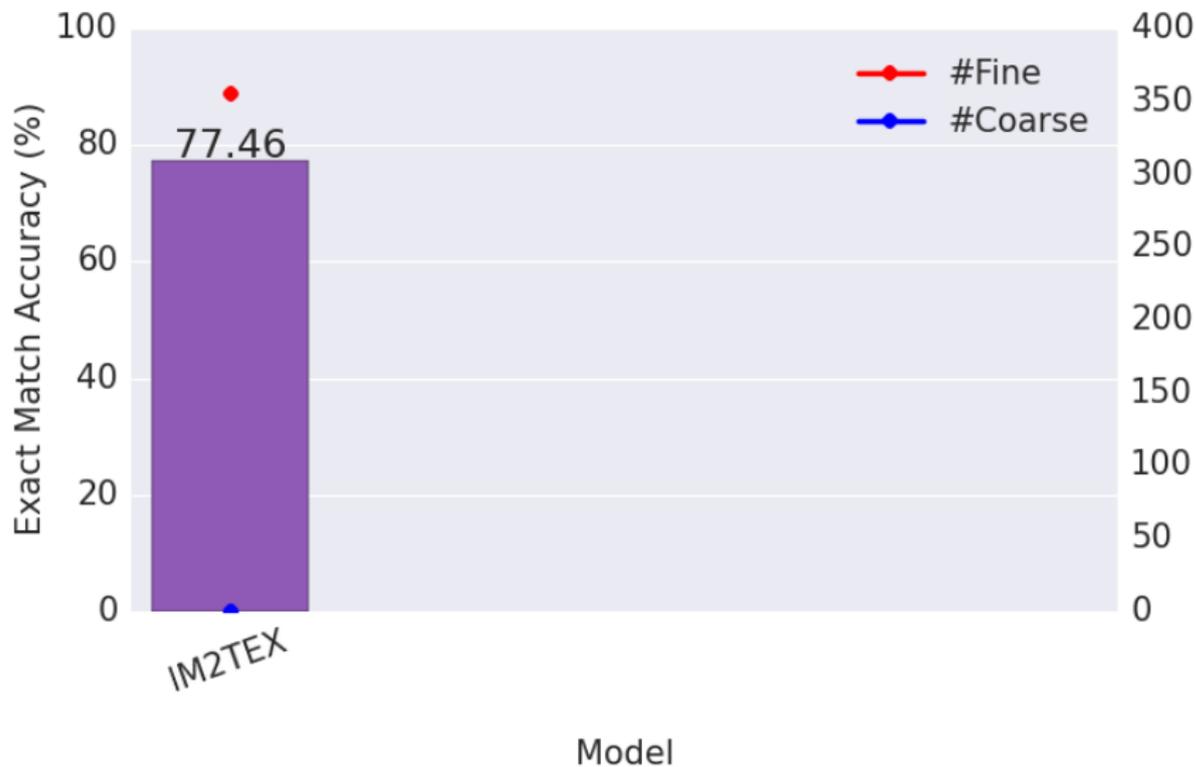
Baseline Results



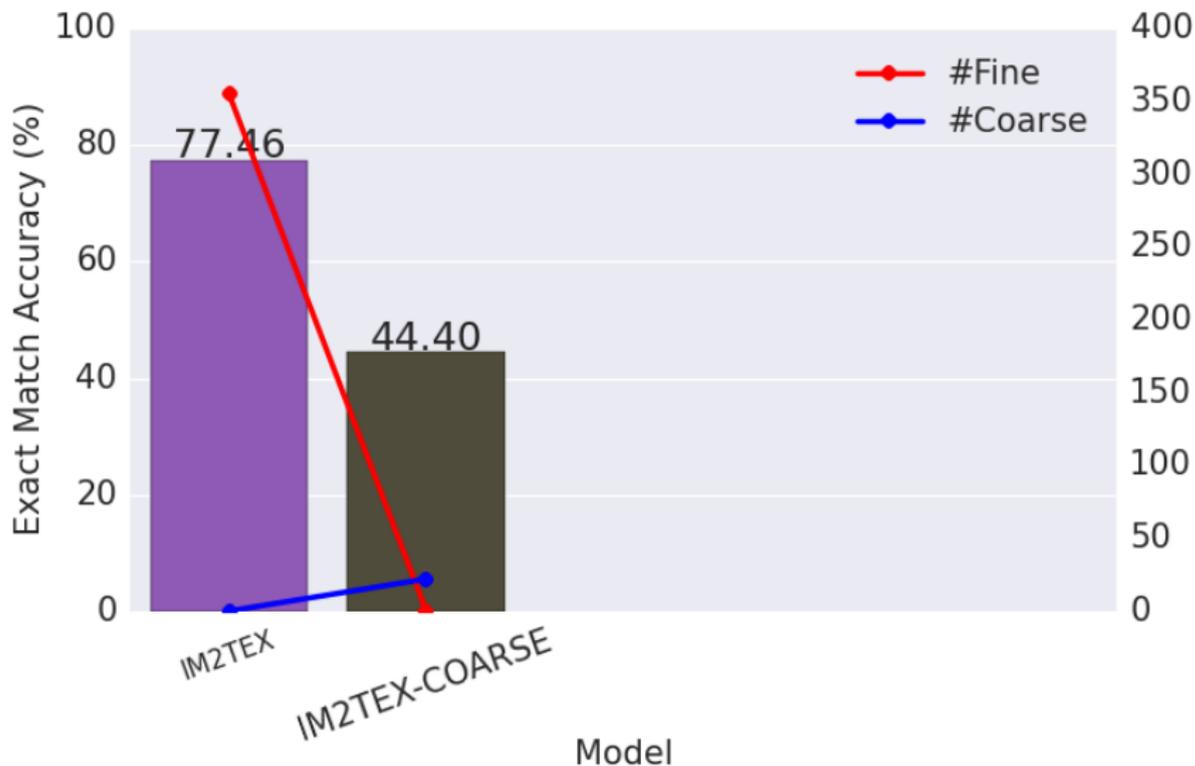
Baseline Results



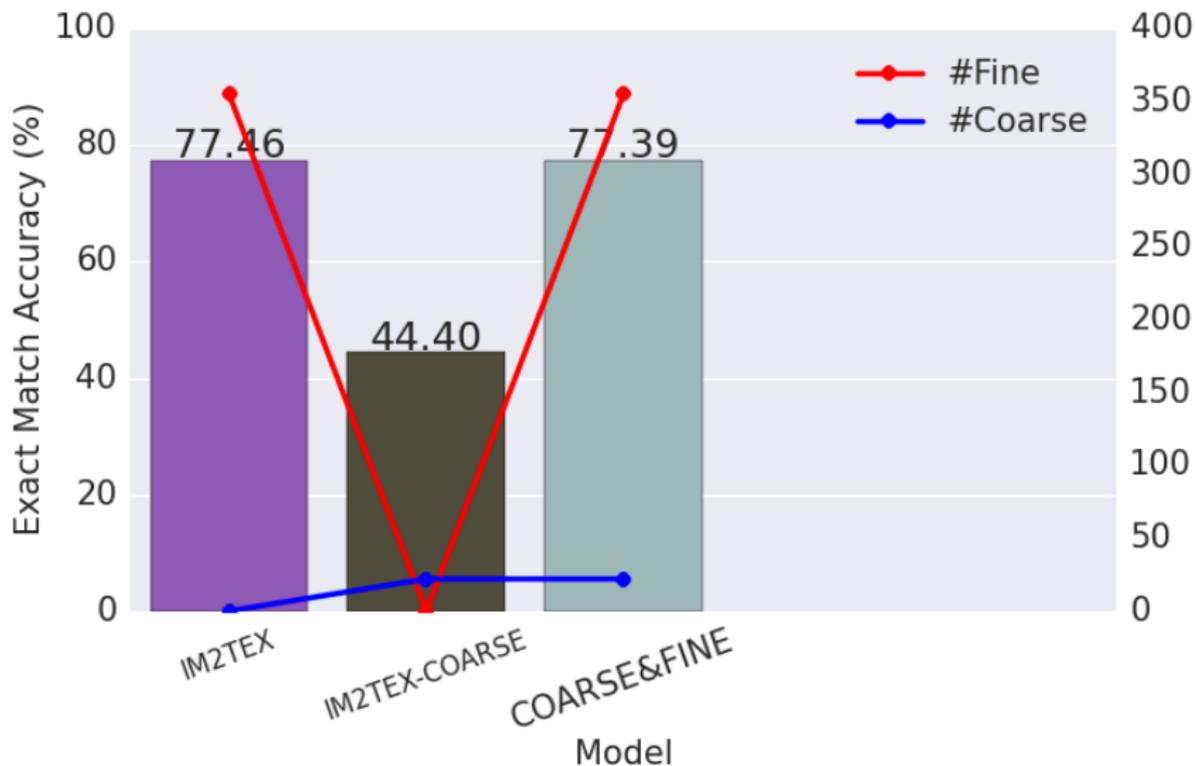
Main Results



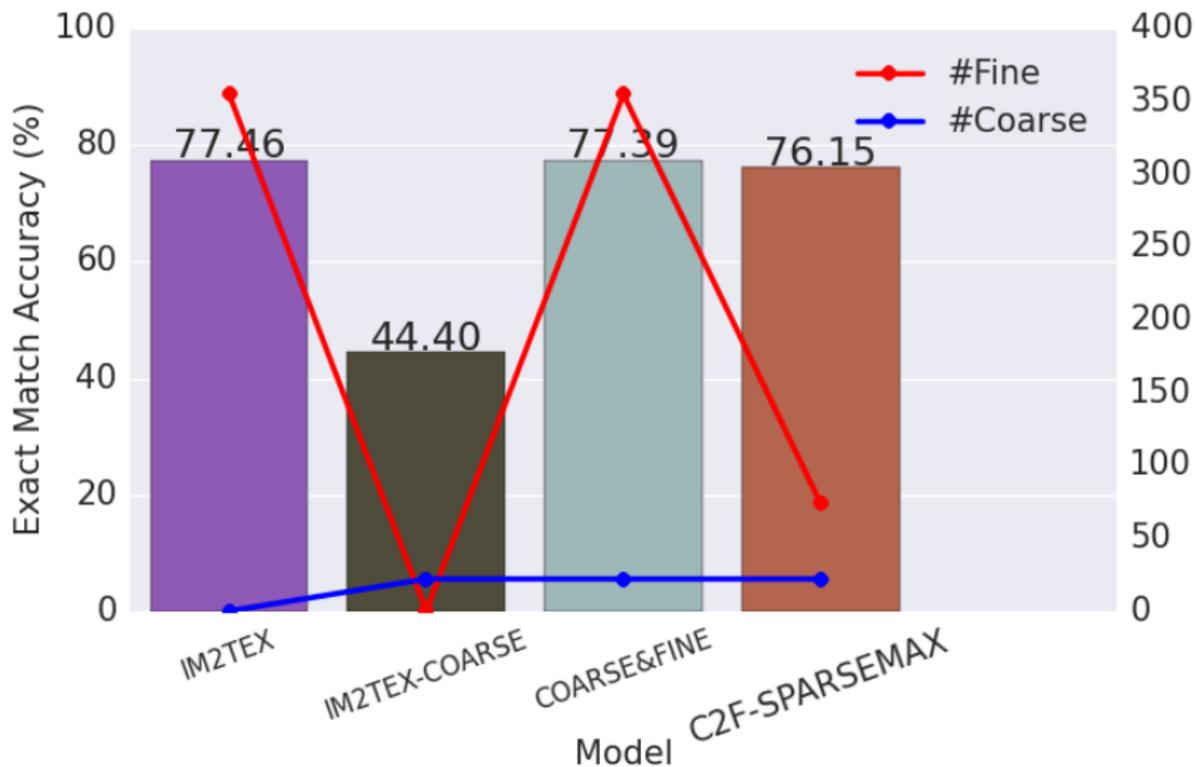
Main Results



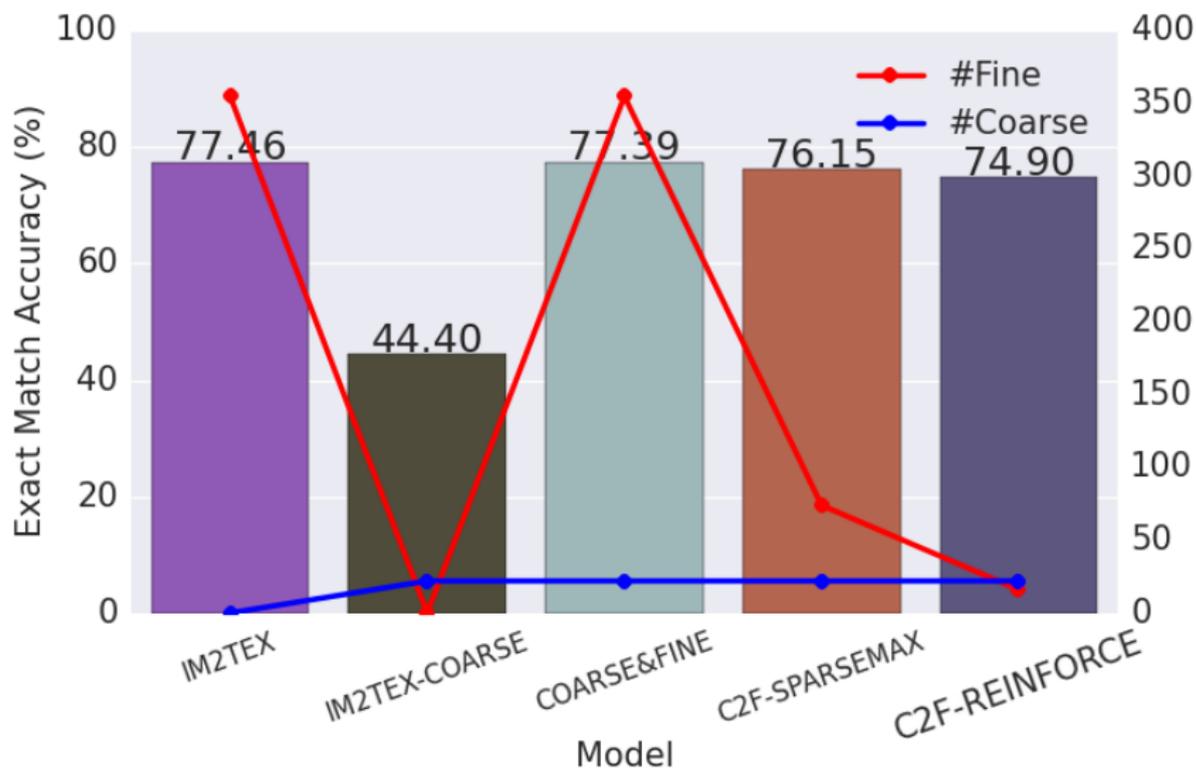
Main Results



Main Results



Main Results



Qualitative Results

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

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

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

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

$$Y(T, U) = \int_{\mathcal{F}} \frac{d^2 \tau}{\Im \tau} \Gamma_{2,2}(T, U) \left(-6 \left[\overline{\Omega}_2 - \frac{1}{8\pi \Im \tau} \right] \frac{\overline{\Omega}}{\eta^{24}} - \frac{\overline{j}}{8} + 126 \right),$$

Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

$$A_3^3(\alpha' \rightarrow 0) = 2g_d \varepsilon_\lambda^{(1)} \varepsilon_\mu^{(2)} \varepsilon_\nu^{(3)} \left\{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_3^\mu - p_1^\mu) + \eta^{\mu\nu} (p_2^\lambda - p_3^\lambda) \right\}.$$

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

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

$$\left(\begin{array}{c} \delta_\varepsilon B \\ \delta_\varepsilon F \end{array} \right) \sim \left(\begin{array}{c} \varepsilon F \\ \partial \varepsilon + \varepsilon B \end{array} \right)$$

$$\int_{C_{d-1}^2} f(H) d\nu_{d-1}(H) = c_3 \int_{C_2^2} \int_{C_{d-1}^2} f(H)[H, A]^2 d\nu_{d-1}^L(H) d\nu_2^A(L).$$

$$\int_{C_{d-1}^2} f(H) d\nu_{d-1}(H) = c_3 \int_{C_2^2} \int_{C_{d-1}^2} f(H)[H, A]^2 d\nu_{d-1}^L(H) d\nu_2^A(L).$$

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

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

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

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

$$(P_W - K_W) \phi'(z_q) | \chi \rangle = 0$$

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

Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

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

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

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

$$\left(\begin{array}{c} \delta_\epsilon B \\ \delta_\epsilon F \end{array} \right) = \left(\begin{array}{cc} \epsilon & F \\ \partial\epsilon + \epsilon B & \end{array} \right)$$

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

$$\int \lim_{\epsilon \rightarrow 0} \frac{d}{d\epsilon} f(H) d\nu_{d-1}(H) = c_3 \int \lim_{\epsilon \rightarrow 0} \frac{d}{d\epsilon} f(H) d\nu_{d-1}^L(H) d\nu_{d-1}^A(L)$$

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

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

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

$$(P_{ij} - K_{ij}) \phi'(z_{ij}) | \chi \rangle = 0$$

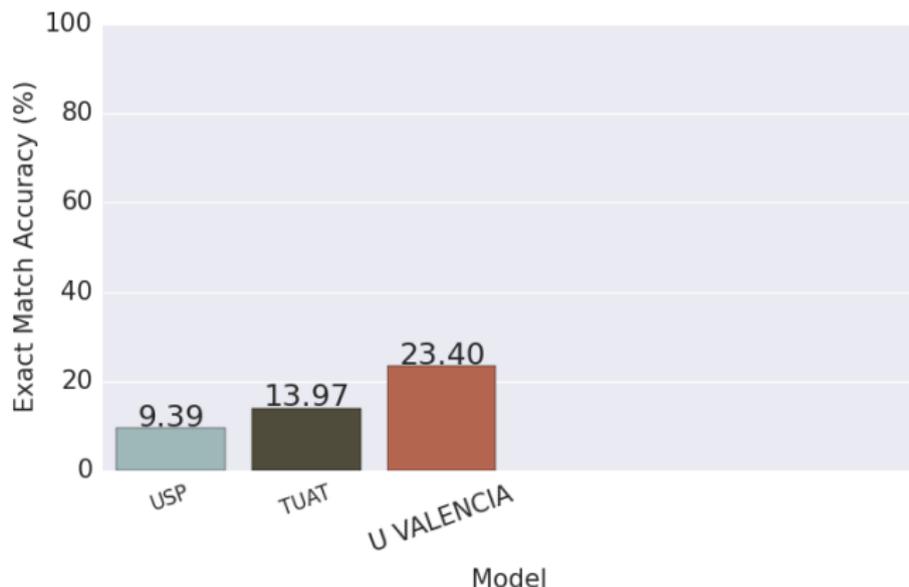
$$\lambda_{n,1}^{(2)} = \frac{\partial \overline{H}_0}{\partial q_{n,0}}, \quad \lambda_{n,j_n}^{(2)} = \frac{\partial \overline{H}_0}{\partial m_{j_n-1}} - \rho_{n,j_n-1}, \quad j_n = 2, 3, \dots, m_n - 1$$

$$(P_{ij} - K_{ij}) \phi'(z_{ij}) | \chi \rangle = 0$$

Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

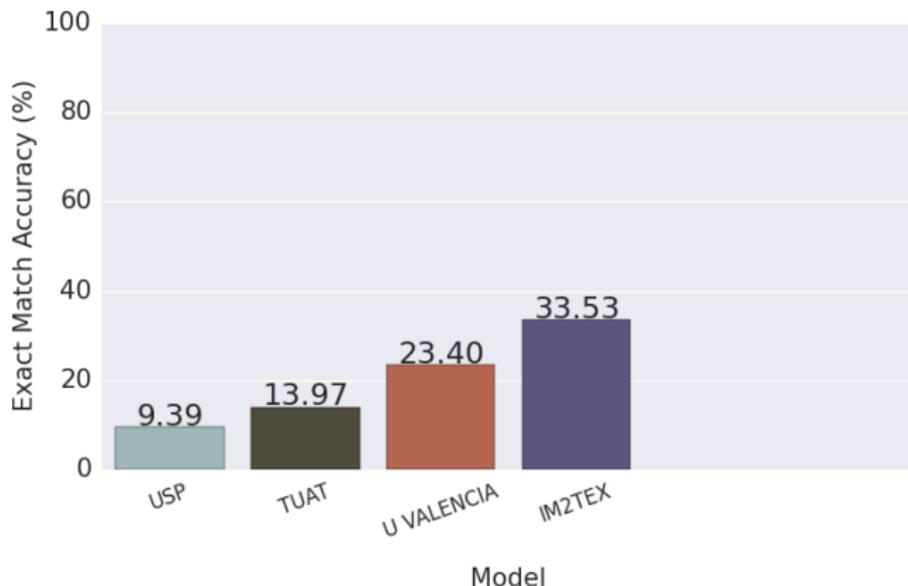
CROHME 13



Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

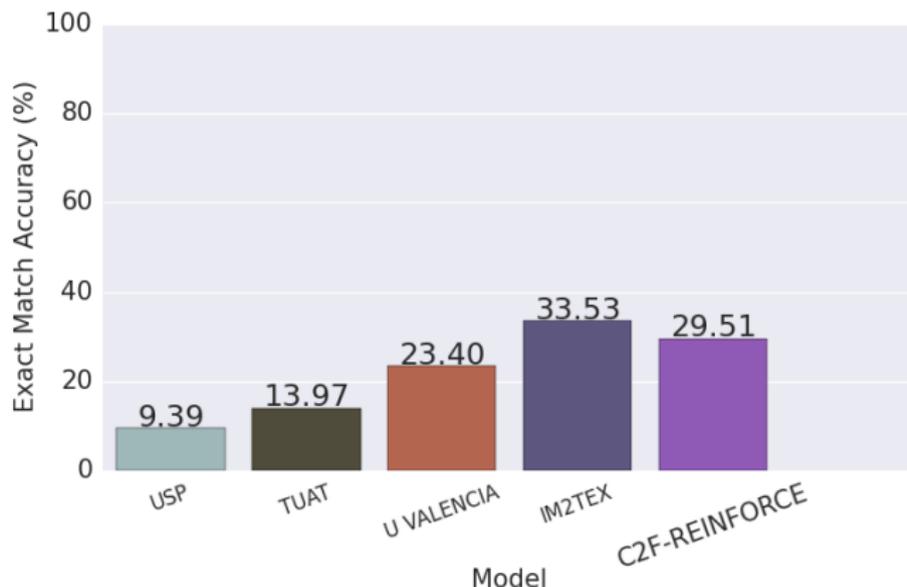
CROHME 13



Handwritten Formulas

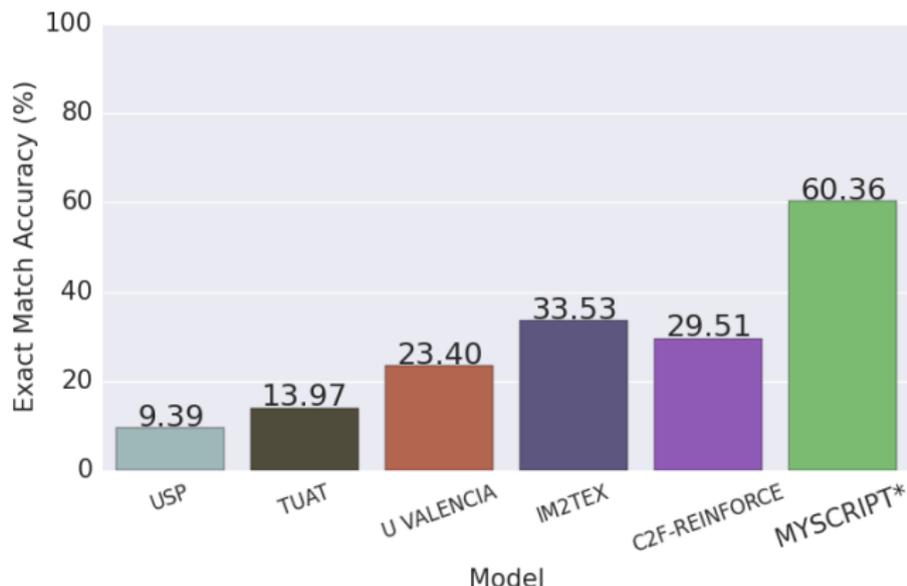
- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

CROHME 13



Handwritten Formulas

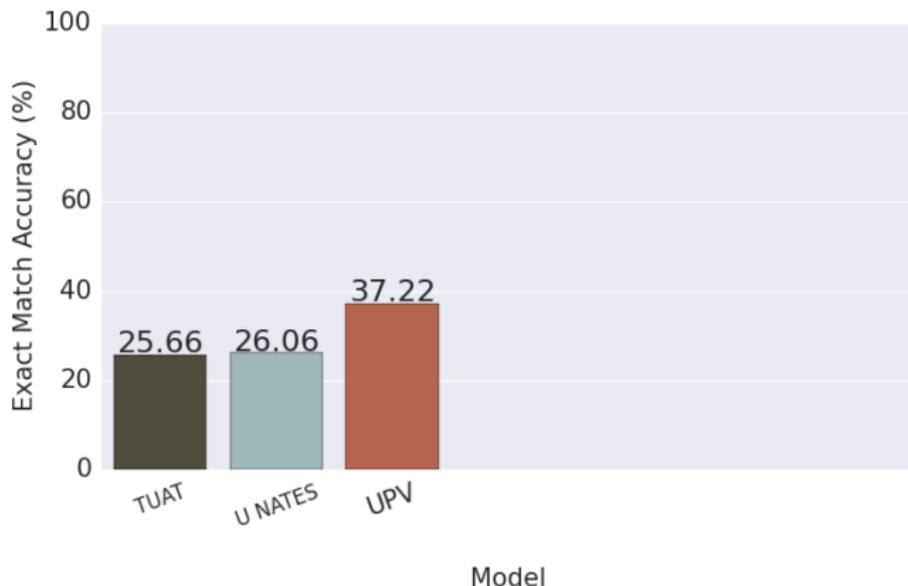
- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)
CROHME 13 (*uses private in-domain handwritten training data)



Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

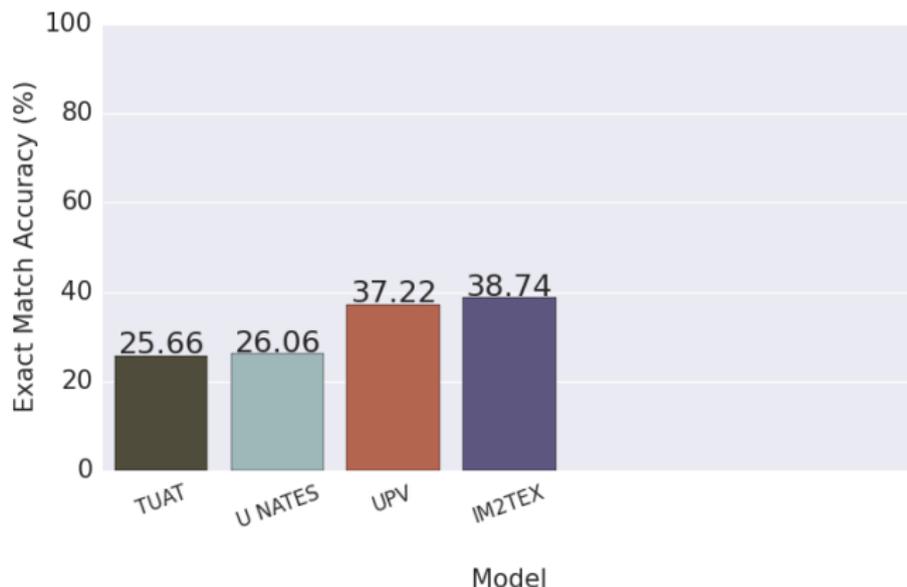
CROHME 14



Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

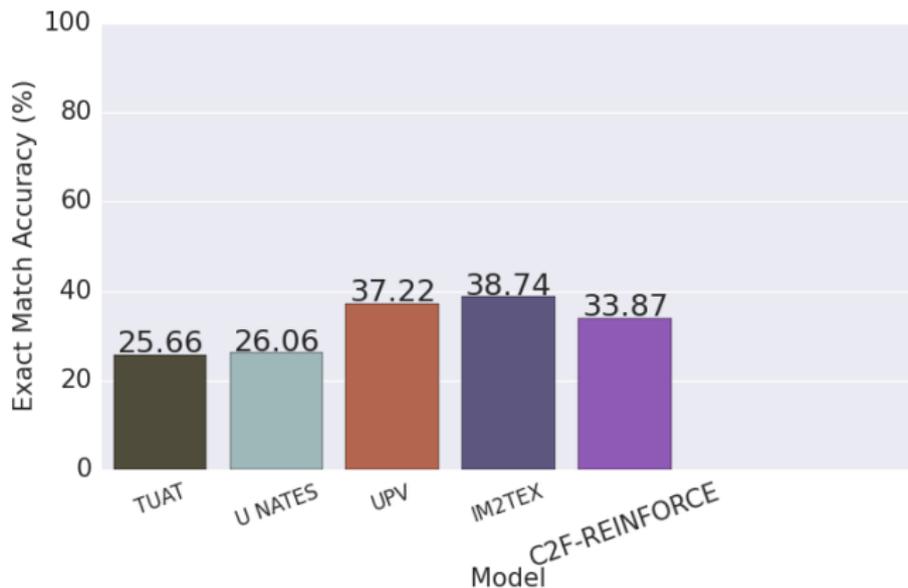
CROHME 14



Handwritten Formulas

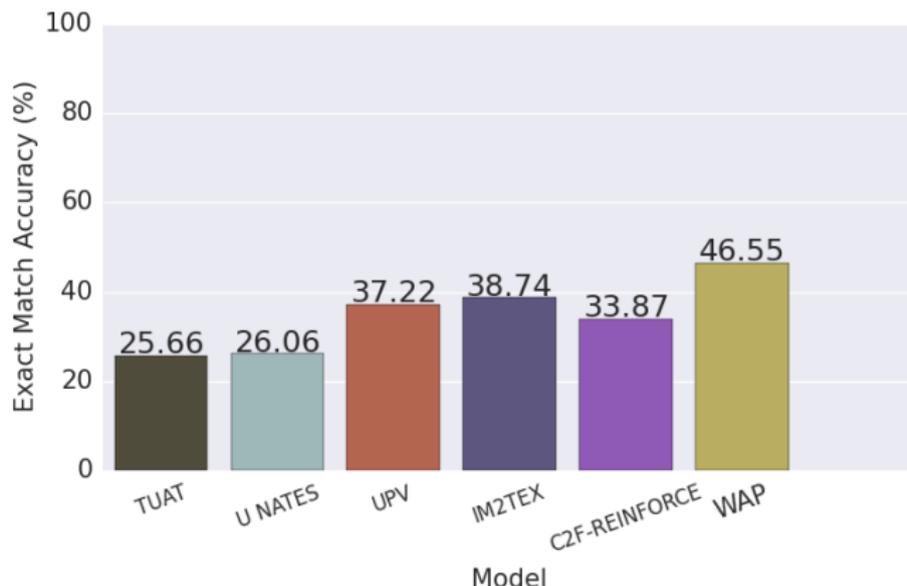
- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)

CROHME 14



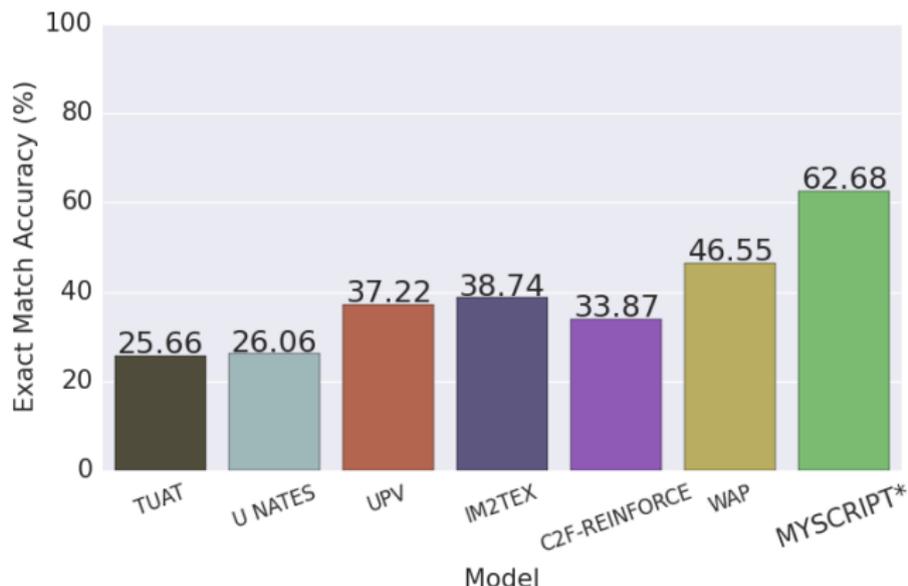
Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)
CROHME 14 (WAP: Zhang et al. [2017])



Handwritten Formulas

- Synthetic handwritten formulas by using handwritten characters [Kirsch, 2010] as font, used for pretraining
- Finetune and evaluate on CROHME 13 and 14 (8K training set)
CROHME 14 (*uses private in-domain handwritten training data)



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.
- 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. Infnty: 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.
- J. Zhang, J. Du, S. Zhang, D. Liu, Y. Hu, J. Hu, S. Wei, and L. Dai. Watch, attend and parse: An end-to-end neural network based approach to handwritten mathematical expression recognition. *Pattern Recognition*, 2017.

- More visualizations:
<http://lstm.seas.harvard.edu/latex/>
- Source code (part of OpenNMT):
<http://opennmt.net/OpenNMT/applications/>