



# Image Captioning to Assist the Visually Impaired

PIERRE DOGNIN, IGOR MELNYK, YOUSSEF  
MROUEH, INKIT PADHI, MATTIA RIGOTTI,  
JERRET ROSS, YAIR SCHIFF\*

*CVPR VizWiz Grand Challenge Workshop  
June 14, 2020*

\*All authors contributed equally

**Pierre Dognin**



**Igor Melnyk**



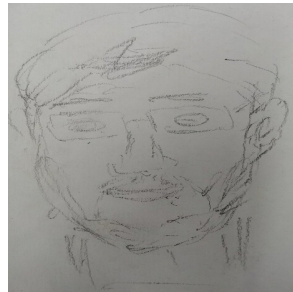
**Youssef Mroueh**



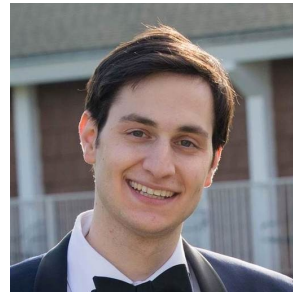
**Inkit Padhi**



**Mattia Rigotti**



**Jerret Ross**

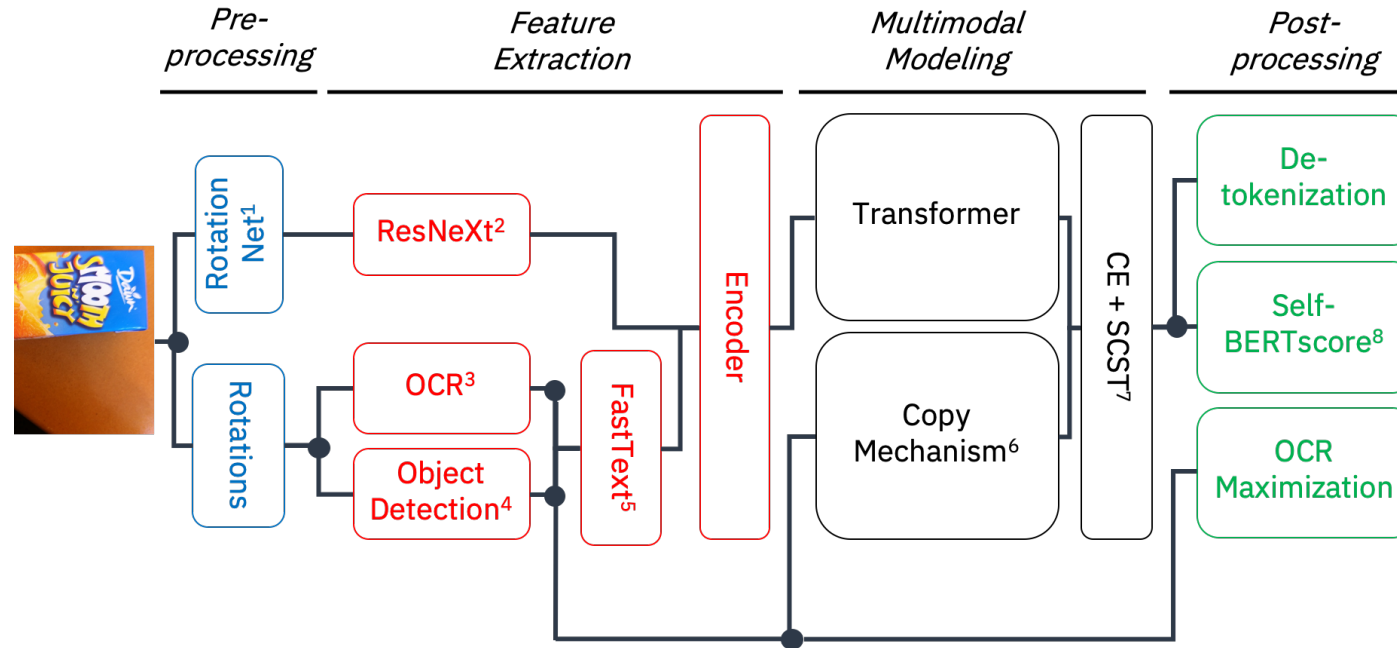


**Yair Schiff**



**Richard Young**

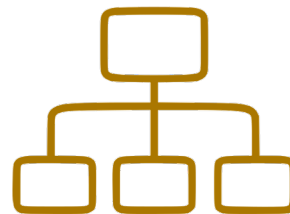
# MODEL:



# MAIN IDEAS:



Open Source



Ensemble



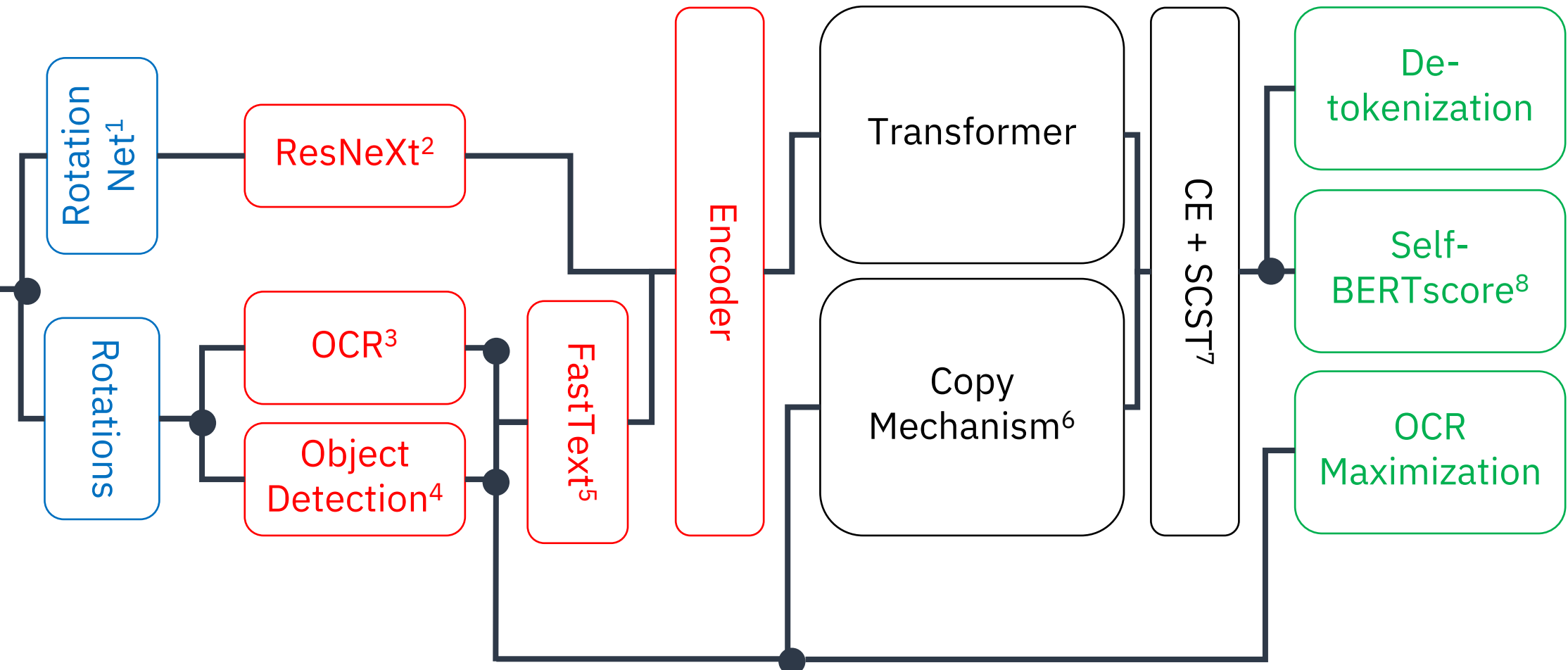
Pre/Post-process

*Pre-processing*

*Feature Extraction*

*Multimodal Modeling*

*Post-processing*



1. Spyros Gidaris, Praveer Singh, and Nikos Komodakis. "Unsupervised Representation Learning by Predicting Image Rotations". In: CoRRabs/1803.07728(2018). arXiv: 1803.07728. URL: <http://arxiv.org/abs/1803.07728>
2. Saining Xie et al. "Aggregated Residual Transformations for Deep Neural Networks". In: arXiv preprint arXiv: 1611.05431(2016)
3. Jeonghun Baek et al. "What Is Wrong With Scene Text Recognition Model Comparisons? Dataset and Model Analysis". In: International Conference on Computer Vision (ICCV). to appear. 2019. published.
4. Youngmin Baek et al. "Character Region Awareness for Text Detection". In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019, pp. 9365–9374.
5. Mingxing Tan, Ruoming Pang, and Quoc V Le. "Efficientdet: Scalable and efficient object detection". In: arXiv preprint arXiv: 1911.09070 (2019).

6. Piotr Bojanowski et al. "Enriching Word Vectors with Subword Information". In: Transactions of the Association for Computational Linguistics5 (2017), pp. 135–146.issn: 2307-387X
7. Jiatao Gu et al. "Incorporating Copying Mechanism in Sequence-to-Sequence Learning". In: CoRRabs/1603.06393 (2016). arXiv: 1603.06393. URL: <http://arxiv.org/abs/1603.06393>.
8. Steven J. Rennie et al. "Self-critical Sequence Training for Image Captioning". In: CoRRabs/1612.00563 (2016). arXiv: 1612.00563. URL: <http://arxiv.org/abs/1612.00563>.
9. Tianyi Zhang et al. "BERTScore: Evaluating Text Generation with BERT". In: CoRRabs/1904.09675 (2019). arXiv: 1904.09675. URL: <http://arxiv.org/abs/1904.09675>.

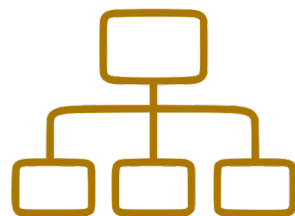
Copy mechanism included	Transformer Architectures	Test-dev CIDEr Score
✗	<i>Same</i>	76.6
✗	<i>Varying</i>	78.0
✓	<i>Same</i>	75.1
✓	<i>Varying</i>	75.6

**Total Ensemble** (90 models, with post-processing)

**80.38**



**Open Source**



**Ensemble**



**Pre/Post-process**



## Open Source

*Leveraging existing implementation for OCR and Object Detection modules made incorporation into pipeline more seamless*



## Ensemble

Jeonghun Baek et al. “What Is Wrong With Scene Text Recognition Model Comparisons? Dataset and Model Analysis”. In: International Conference on Computer Vision (ICCV). to appear. 2019. published.



## Pre/Post-process

Youngmin Baek et al. “Character Region Awareness for Text Detection”. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019, pp. 9365–9374.

Mingxing Tan, Ruoming Pang, and Quoc V Le. “Efficientdet: Scalable and efficient object detection”. In: arXiv preprint arXiv: 1911.09070 (2019).



**Open Source**



**Ensemble**

*Ensembling varying  
Transformer architectures and  
combining copy / non-copy  
mechanism injected caption  
diversity*



**Pre/Post-process**





**Open Source**



**Ensemble**



**Pre/Post-process**

*Combining de-tokenization, Self Bert scoring, and OCR-maximization provided additional gains*



# Live image captioning with text-to-speech



**Richard Young**



**Thank  
you!**