



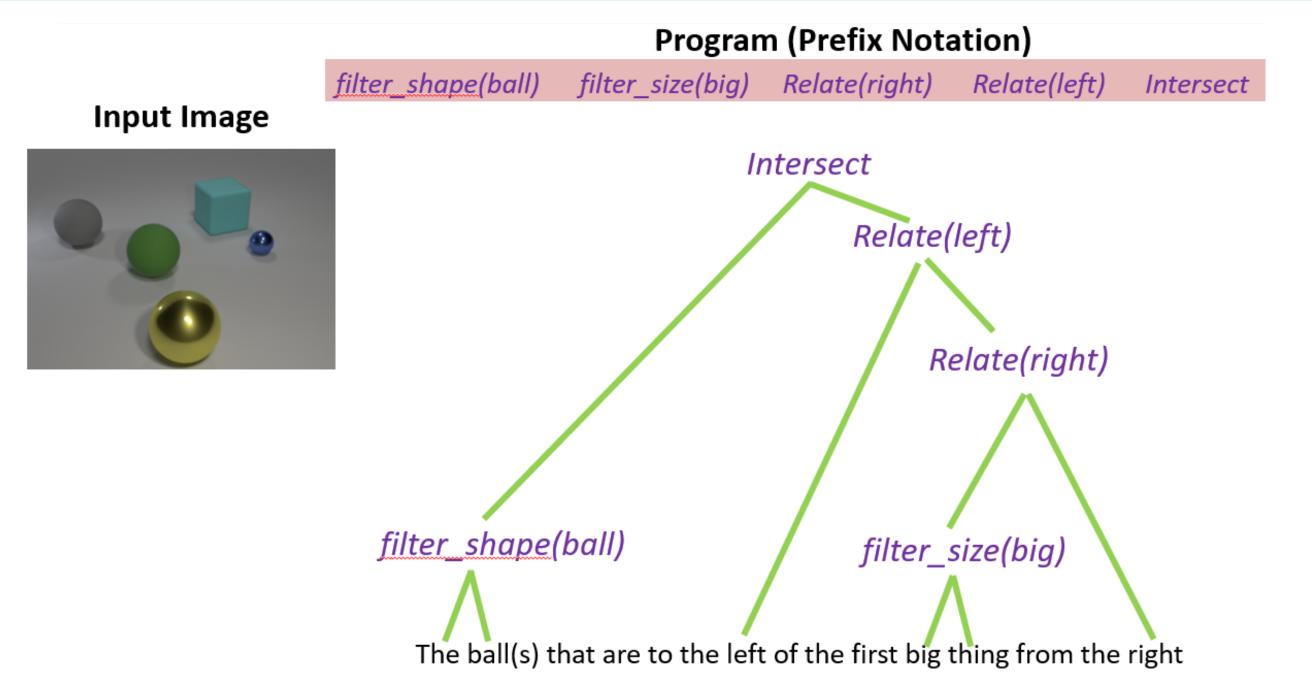
# Robust Visual Reasoning via Language Guided Neural Module Networks

Arjun R Akula<sup>1</sup> Varun Jampani<sup>2</sup> Soravit(Beer) Changpinyo<sup>2</sup> Song-Chun Zhu<sup>3</sup>



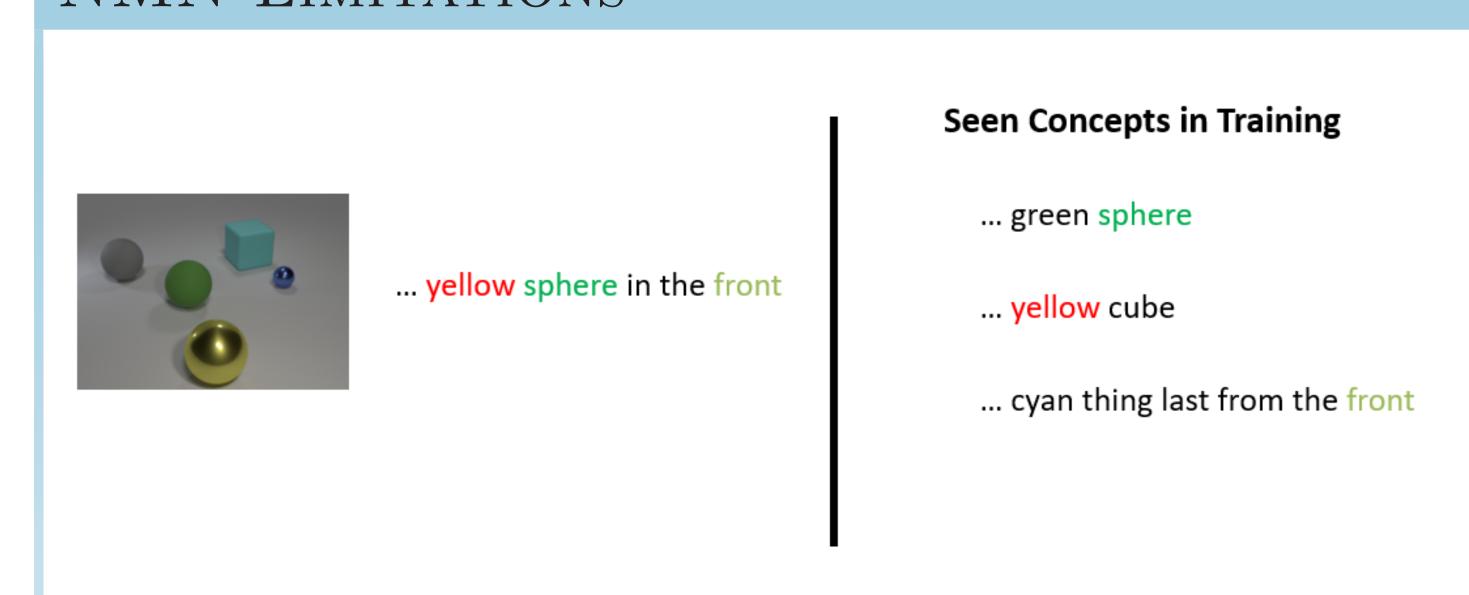
<sup>1</sup>UCLA, <sup>2</sup>Google Research, <sup>3</sup>BIGAI; Peking University; Tsinghua University, <sup>3</sup>Mila; McGill University

### NEURAL MODULE NETWORK (NMN)

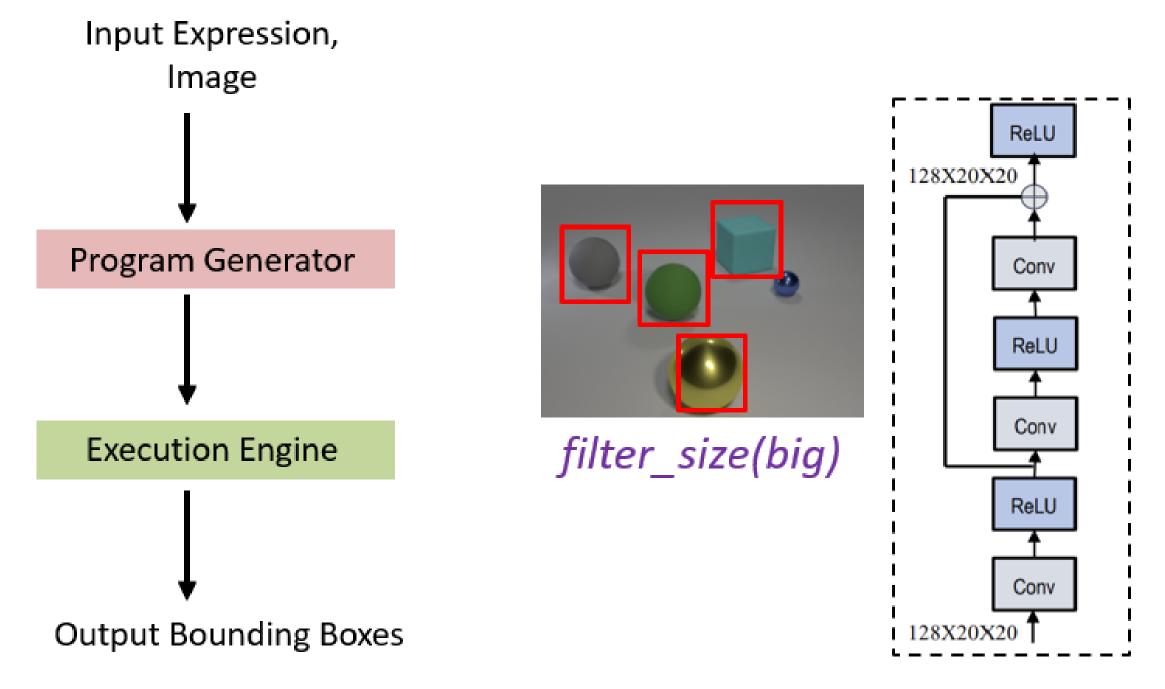


NMNs use an explicit modular reasoning process and provide interpretable output. NMNs first extract logical expression/program from input expression. Then composes a neural network to ground referring expressions.

## NMN LIMITATIONS

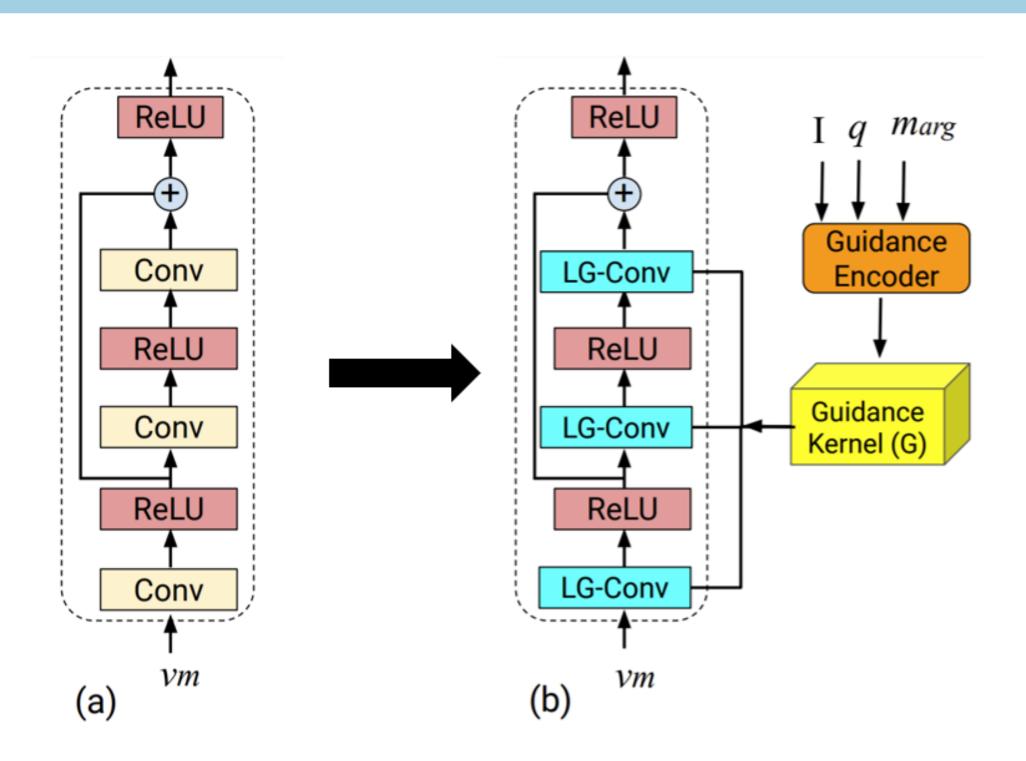


Limitation 1: NMNs are less effective in generalizing to unseen but known language constructs



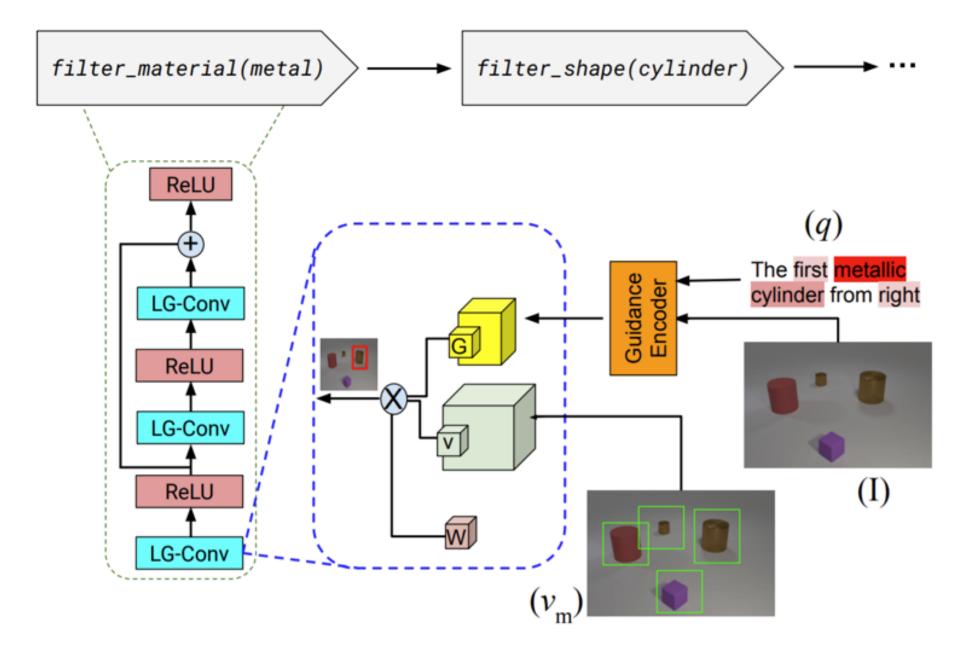
Limitation 2: NMNs use shallow, in-direct language guidance

#### CONTENT ADAPTIVE CONVOLUTIONS

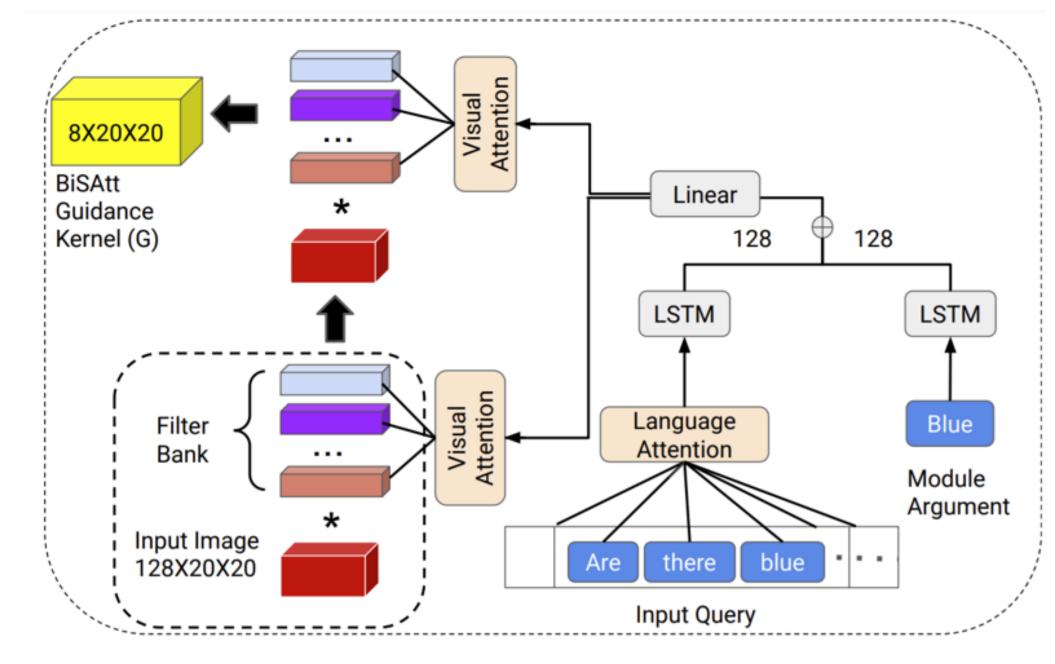


We replace the standard convolution layers with content adaptive convolution layers LG-Conv.

#### BI-SALIENT ATTENTION GUIDANCE ENCODER

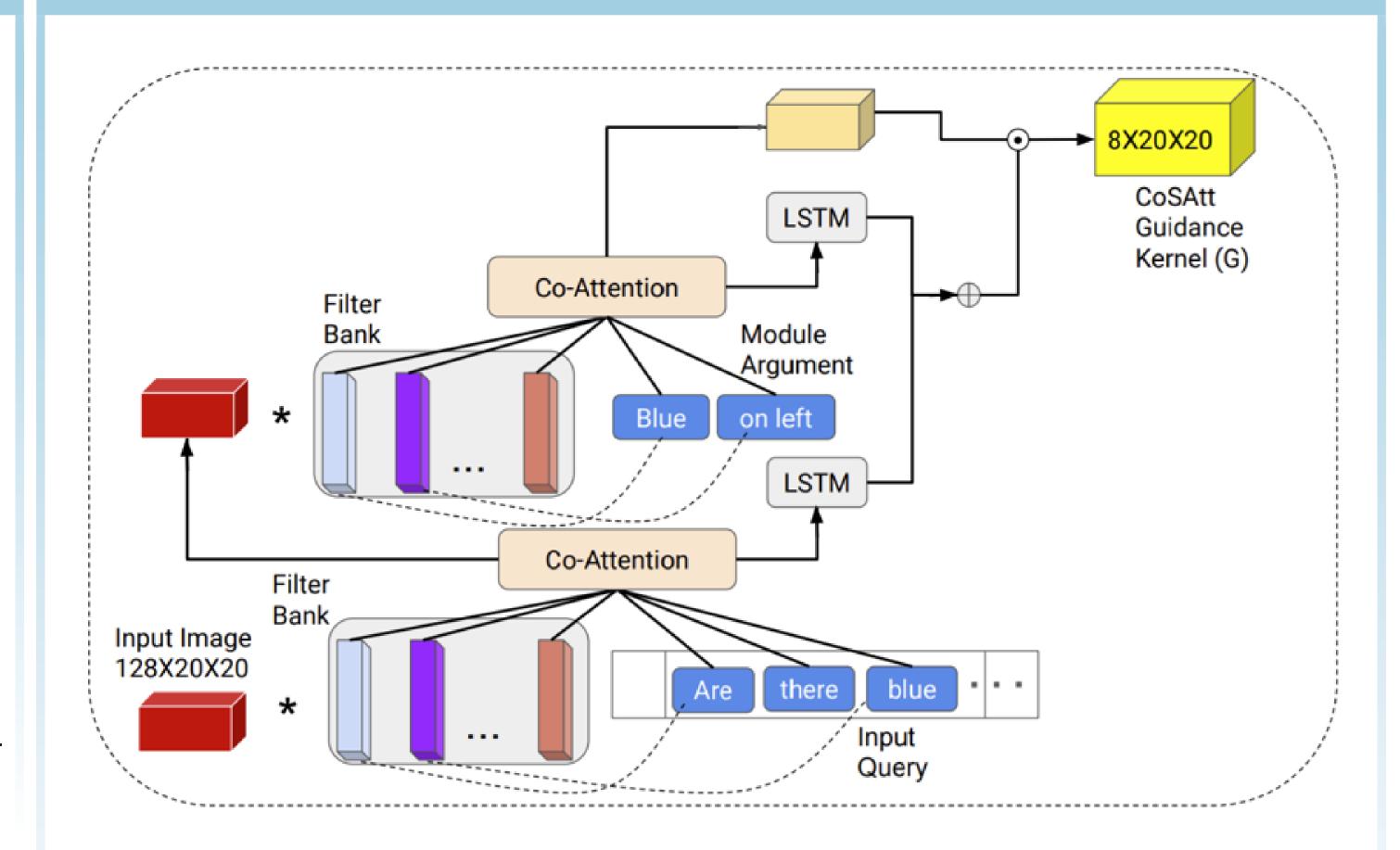


The spatial convolution W is adapted at each pixel in the visual input using the guidance kernel G.



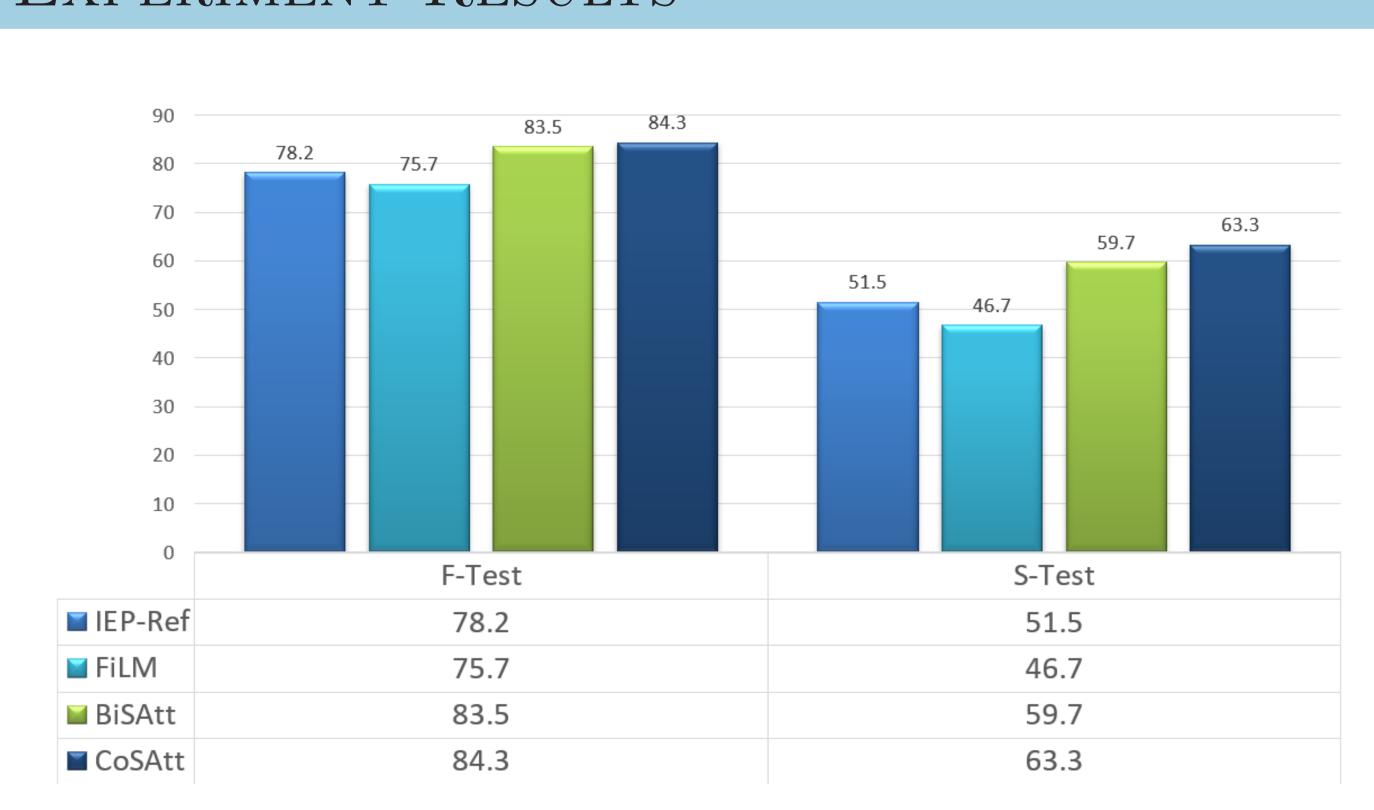
we first encode text inputs and then use it to learn a set of adaptive weights to linearly combine the basis filters.

#### Co-Salient Attention Guidance Encoder



CoSAtt Encoder jointly attends over the input image and text inputs (early fusion) to identify co-salient regions and relationships in visual and language features that are contextually associated with each other.

### EXPERIMENT RESULTS



Our model improves the accuracy on CLEVR-Ref+ test splits by 9.8% on S-Ref and 4.7% on F-Ref, compared with the current state-of-the-art method IEP-Ref.

Relatively more improvements with CoSAtt encoder compared to BiSAtt encoder shows that early fusion of image and text features facilitate in generating more robust guidance kernel.