

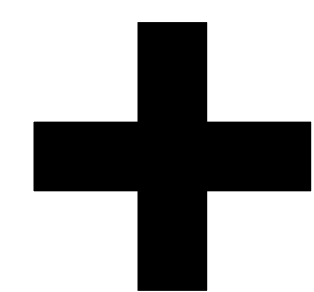
Decoupled Box Proposal and Featurization with Ultrafine-Grained Semantic Labels

Improve Image Captioning and Visual Question Answering

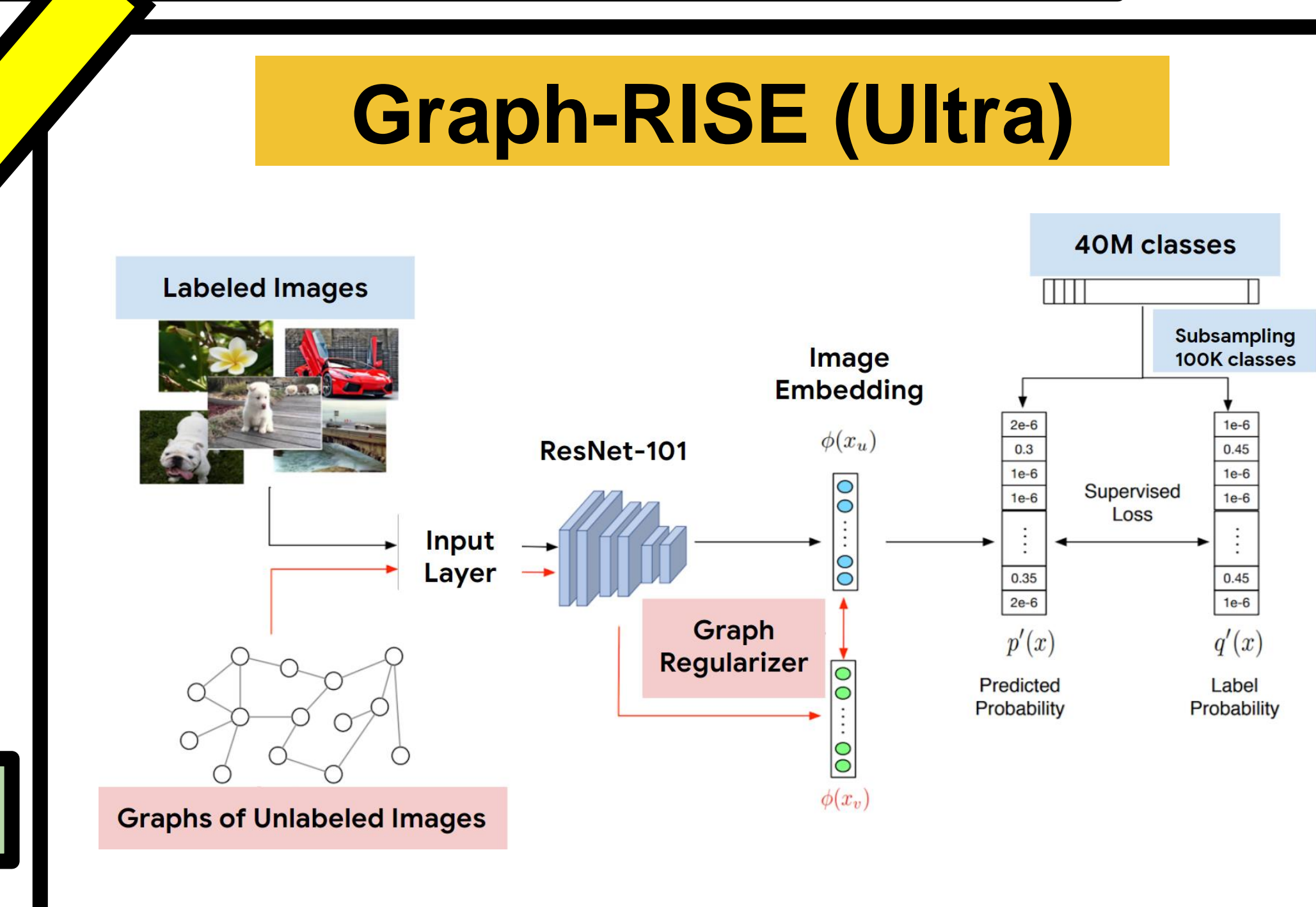
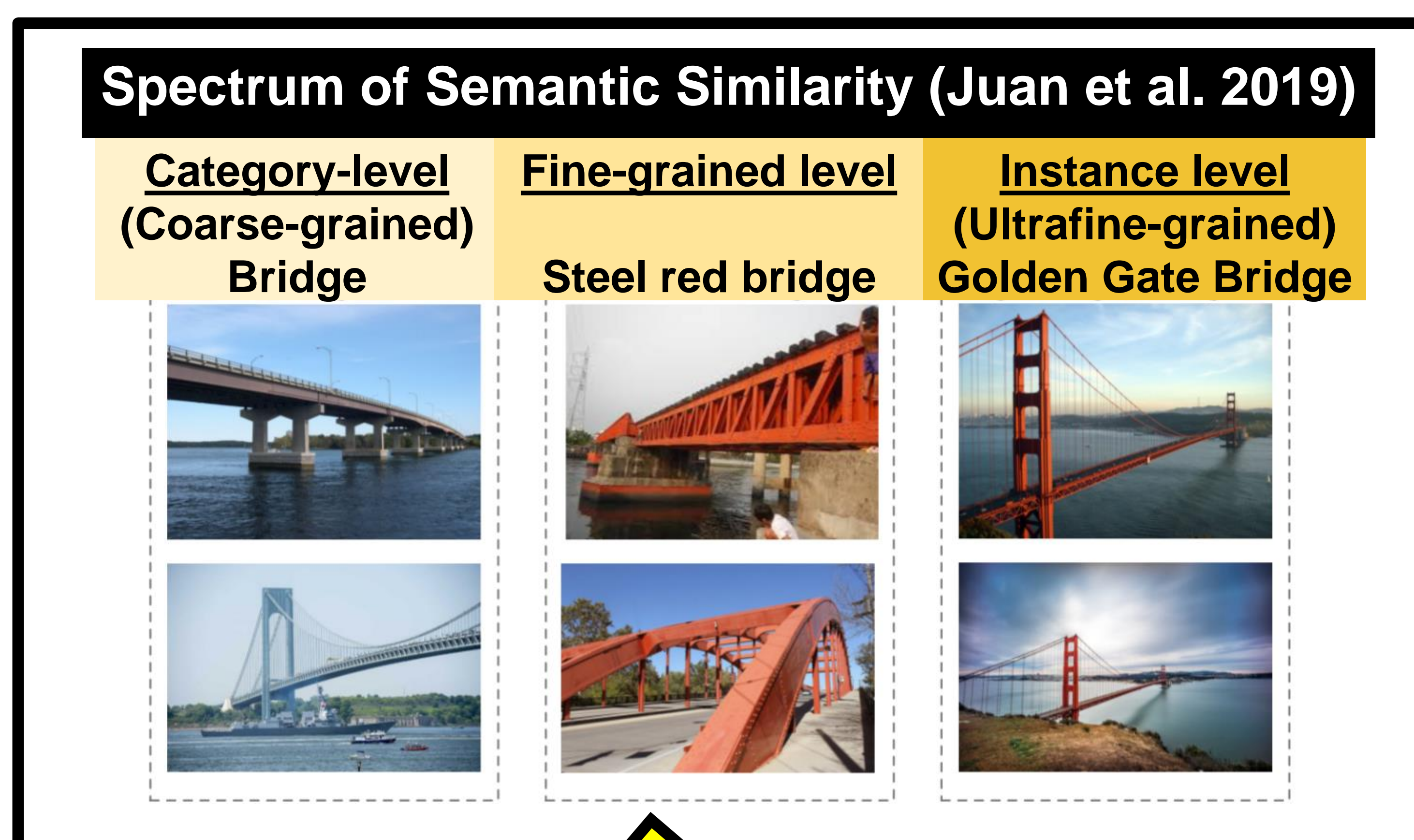
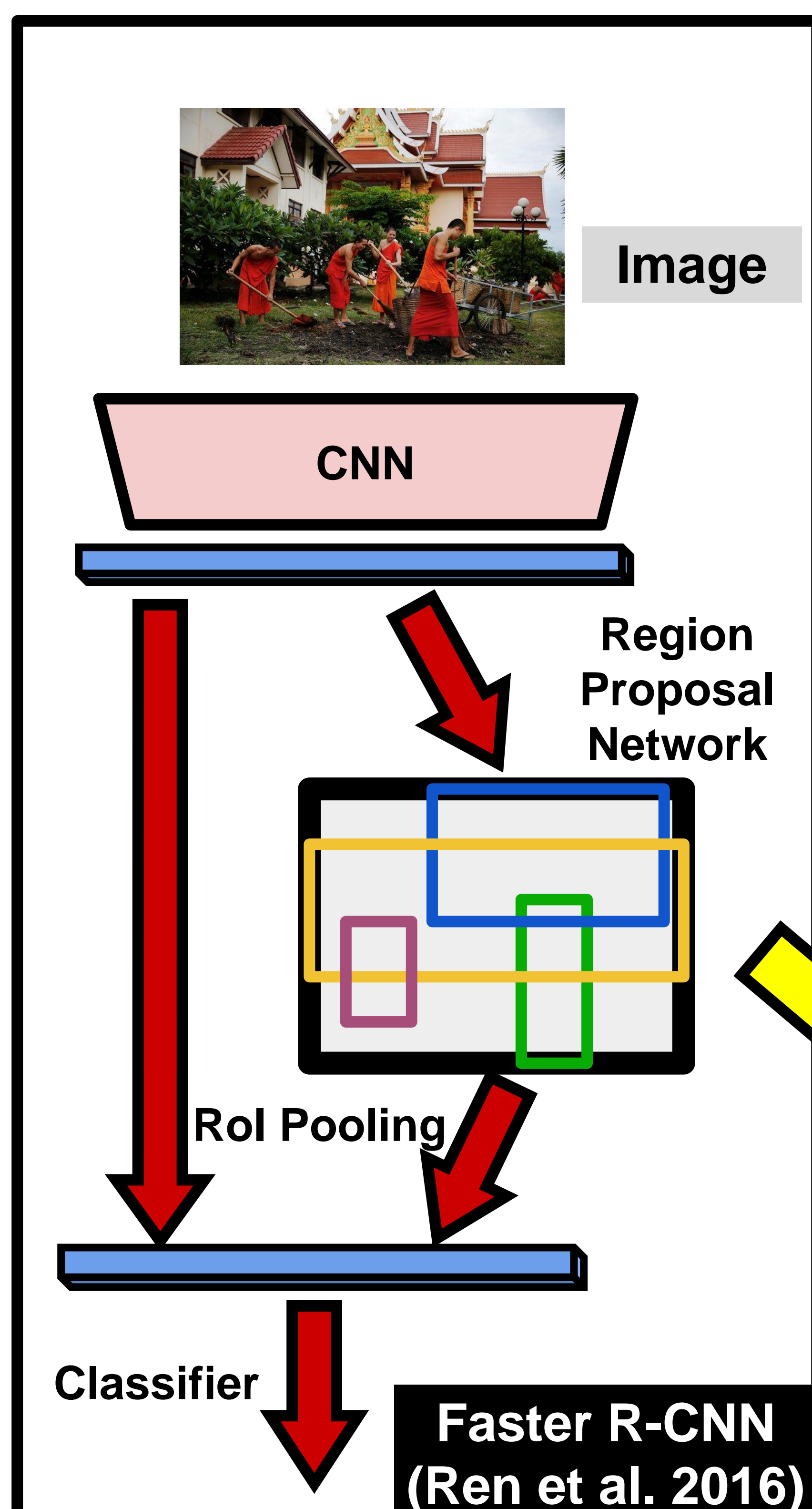
Soravit Changpinyo, Bo Pang, Piyush Sharma, Radu Soricut



Decoupled Box Proposal



Featurization with Ultrafine-Grained Semantic Labels



Vision & Language Models

Image Captioning on Conceptual Captions (Sharma et al. 2018)



Ground-truth
FRCNN (Visual Genome)
Ultra

monks clean a garden at a temple .
a woman walks through the streets .
monks walking in front of a temple

black sesame seeds on a white background
a pile of dried flowers
black chia seeds on a white background

VQA on VizWiz (Gurari et al. 2018)



Q: How much money is this?
A: 1 dollar
A: 20



Q: What is this?
A: beer
A: bbq sauce