

Deep Variation-structured Reinforcement Learning (VRL) for Visual Relationship & Attribute Detection

Xiaodan Liang, Lisa Lee, Eric Xing

Machine Learning Department, Carnegie Mellon University

Motivation

Computers still struggle to understand the interdependency of objects in the scene as a whole.

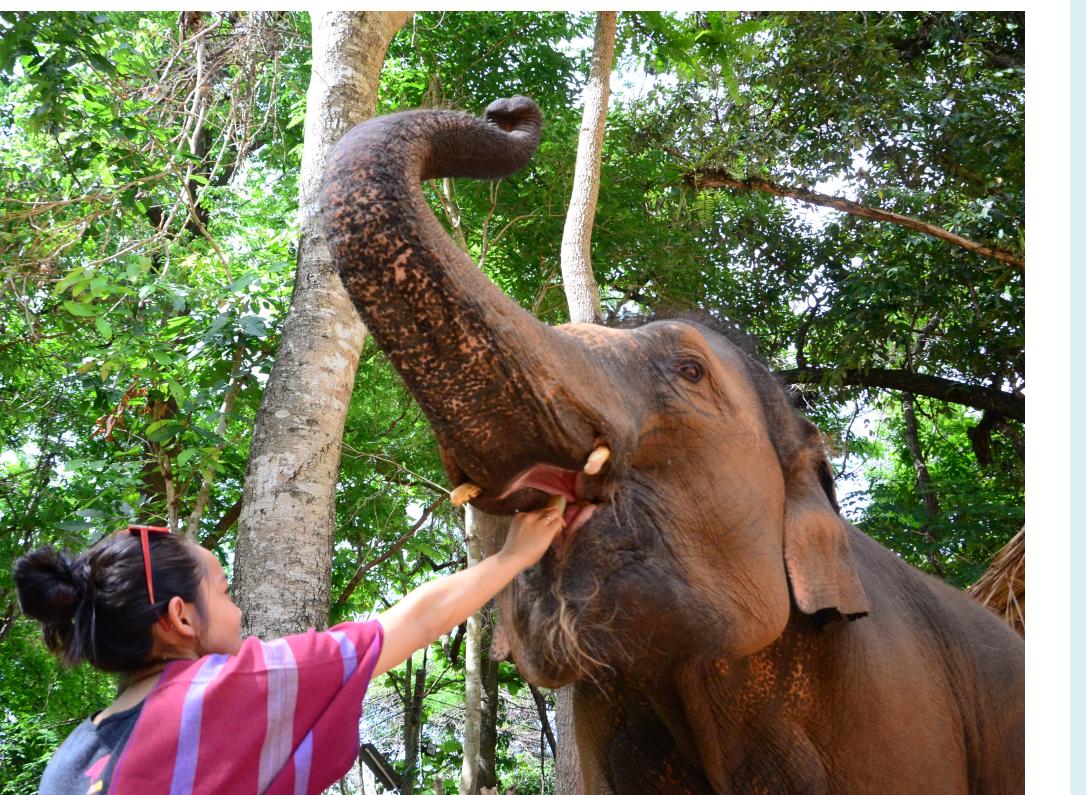
Existing methods ignore **global context cues** capturing the interactions among different object instances.

Relationship $(c, p, c') \in \mathcal{E}^P$

girl → feeding → elephant

Attribute modifier $(c, a) \in \mathcal{E}^A$

elephant → brown



Our Contributions

Use a novel **variation-structured traversal scheme** over the semantic action graph that dynamically constructs small action sets for each step.

Make **sequential decisions** using a deep RL framework, incorporating global context cues and semantic embeddings in the state vector.

Achieve **state-of-the-art** results on Visual Genome & Visual Relationship Detection datasets.

Method	Phr. R@100	Phr. R@50	Rel. R@100	Rel. R@50	Attr. R@100	Attr. R@50
Joint CNN+RPN	2.52	2.44	2.37	2.23	9.77	8.35
Lu et al. (2016)	10.23	9.55	7.96	6.01	-	-
Our VRL	16.09	14.36	13.34	12.57	26.43	24.87

Results for relationship phrase detection (Phr.), relationship detection (Rel.), and attribute detection (Attr.) on the Visual Genome dataset.

R@N is short for Recall@N.

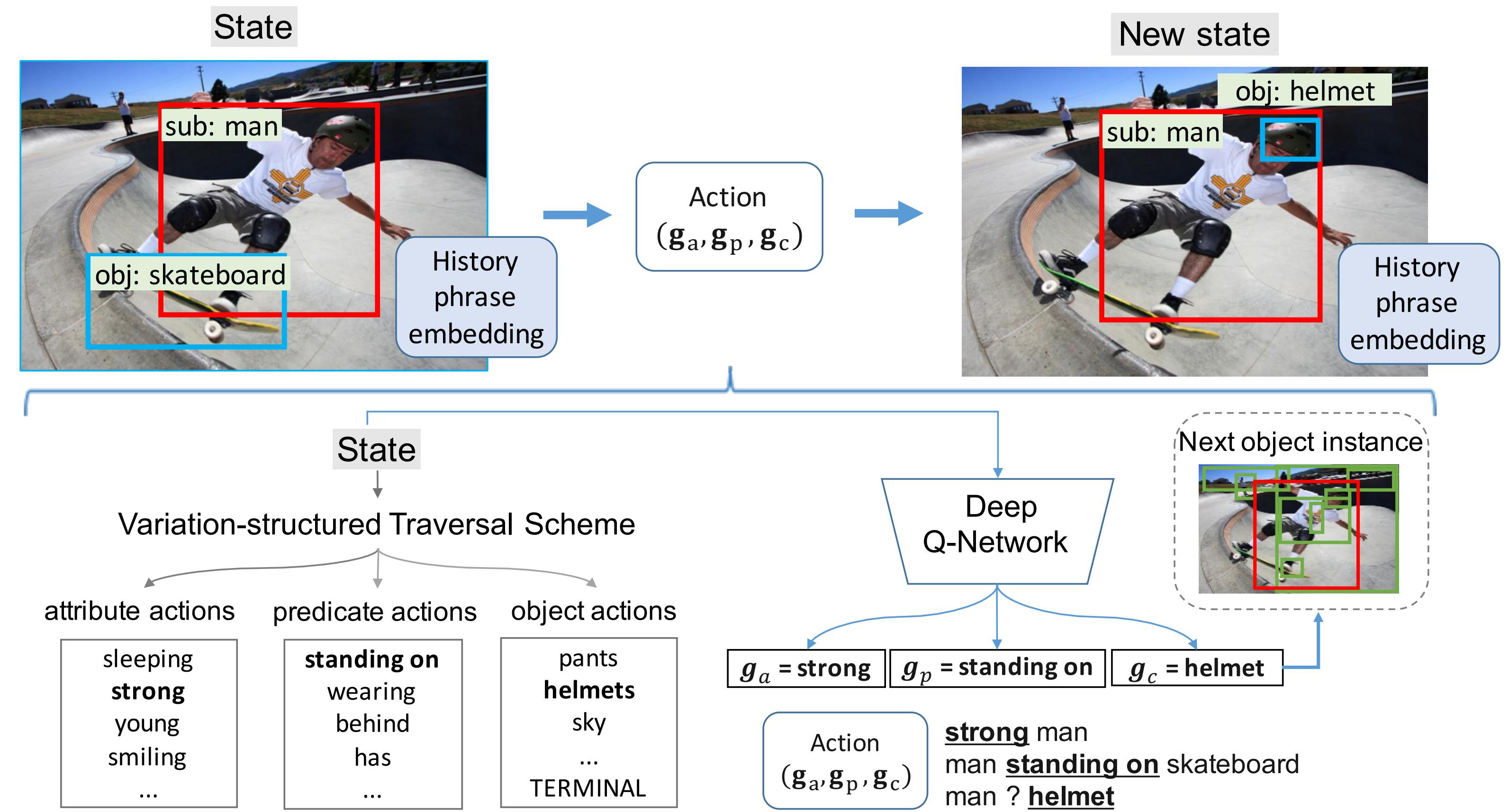
Methods

ϵ -greedy actions: Given the current **subject** and **object** instances, choose:

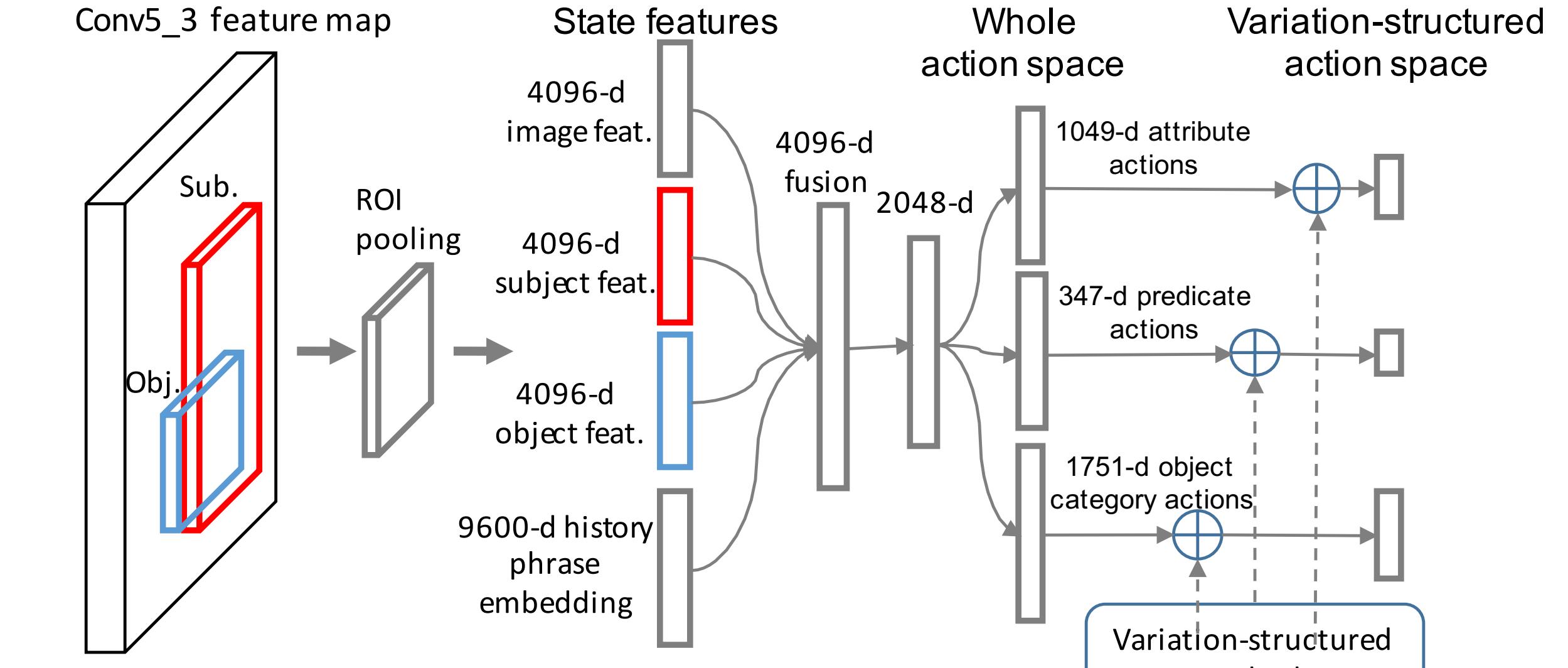
1. An **attribute** g_a describing the subject
2. A **predicate** g_p relating the subject and object
3. The **next object instance** g_c .

Q-learning update given actions (g_a, g_p, g_c) , state transition $f \rightarrow f'$, rewards $(\mathcal{R}_a, \mathcal{R}_p, \mathcal{R}_c)$:

$$\begin{aligned}\theta_a^{(t+1)} &= \theta_a^{(t)} + \alpha \left(\mathcal{R}_a + \lambda \max_{g_{a'}} Q(f', g_{a'}; \theta_a^{(t)}) - Q(f, g_a; \theta_a^{(t)}) \right) \nabla_{\theta_a^{(t)}} Q(f, g_a; \theta_a^{(t)}) \\ \theta_p^{(t+1)} &= \theta_p^{(t)} + \alpha \left(\mathcal{R}_p + \lambda \max_{g_{p'}} Q(f', g_{p'}; \theta_p^{(t)}) - Q(f, g_p; \theta_p^{(t)}) \right) \nabla_{\theta_p^{(t)}} Q(f, g_p; \theta_p^{(t)}) \\ \theta_c^{(t+1)} &= \theta_c^{(t)} + \alpha \left(\mathcal{R}_c + \lambda \max_{g_{c'}} Q(f', g_{c'}; \theta_c^{(t)}) - Q(f, g_c; \theta_c^{(t)}) \right) \nabla_{\theta_c^{(t)}} Q(f, g_c; \theta_c^{(t)})\end{aligned}$$



Network Architecture



Experimental Results

Comparison vs. VRD from Lu et al. (2016)



Example results generated by our VRL on the Visual Genome dataset.

