

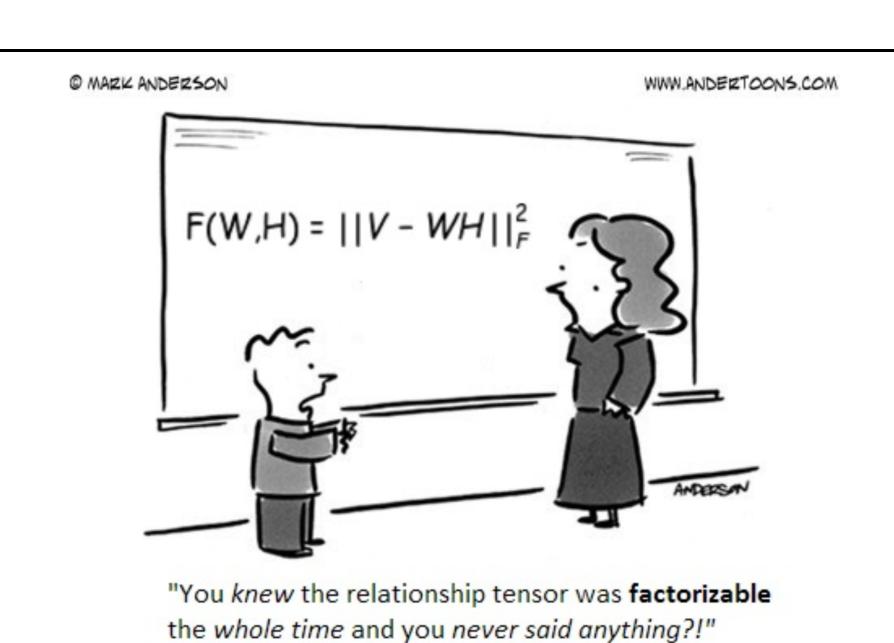
Tensorize, Factorize and Regularize: Robust Visual Relationship Learning

Seong Jae Hwang Sathya N. Ravi Zirui Tao Hyunwoo J. Kim Maxwell D. Collins Vikas Singh http://pages.cs.wisc.edu/~sjh



OBJECTIVES

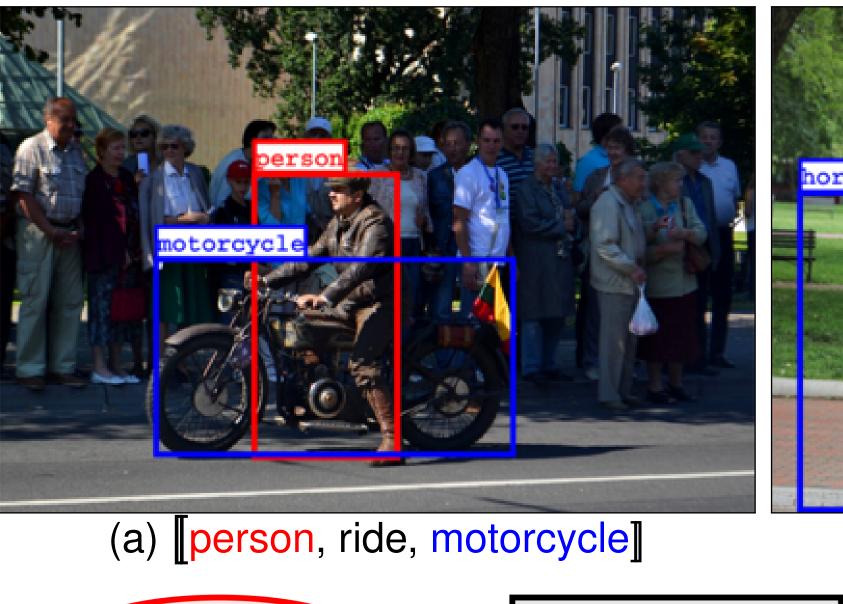
Regularize the Scene Graph learning deep network via multi-relational tensor factorization for robust visual relationship learning.

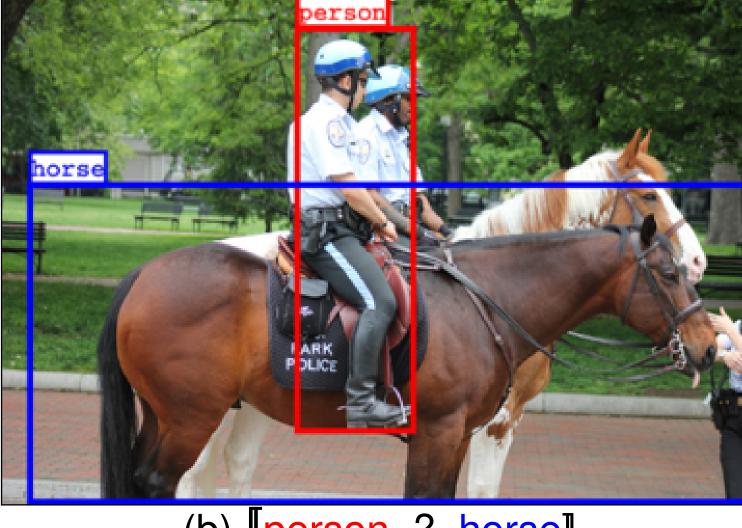


PROBLEM: VISUAL RELATIONSHIP DETECTION

- ▶ Given:
- 1. Image with Object i and Object j of interest
- 2. **Bounding boxes** and **features** (i.e., Faster R-CNN)
- 3. Object labels (ground truths or detection results)
- ▶ Goal:

Predict the visual relationship of the objects: [Object i, ?, Object j $| \cdot |$ corresponding k^{th} predicate.





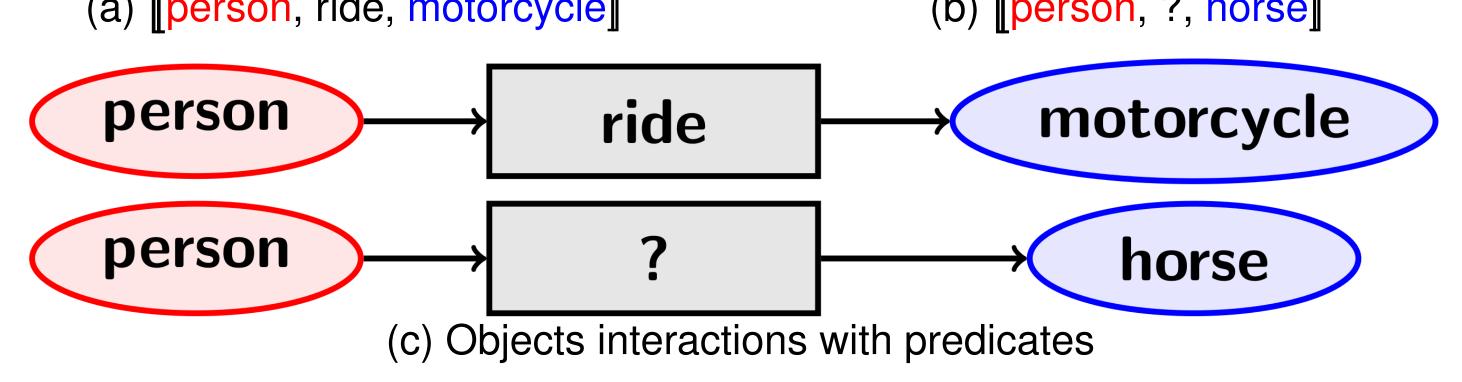
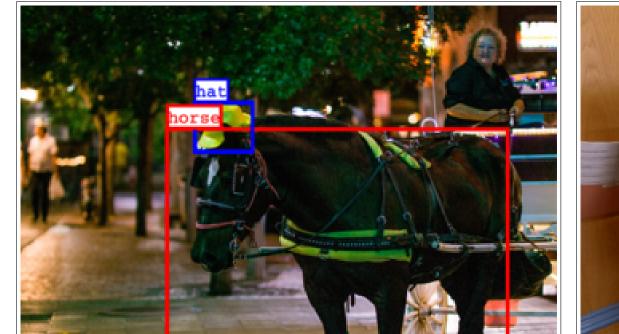


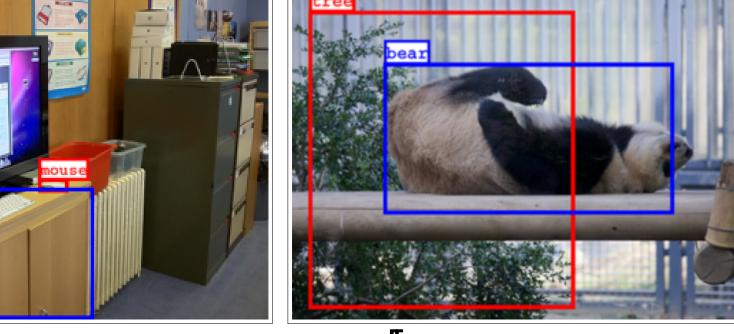
Figure: (a): A relationship instance in a training set. (b): An unknown relationship to predict. (c): The interactions of the objects (i.e., motorcycle and horse are both 'ridable') can be used to infer the correct relationship.

CHALLENGES IN SEMANTIC INFERENCE

- ▶ Large observation space:
- $\Rightarrow N^2M$ possible combinations
- Sparse observations:
- > Visual Genome has 1M relationship instances
- \Rightarrow But observed only \sim 2% of possible combinations
- ➤ Zero-shot learning:
- > Inferring cases unobserved in train set
- \Rightarrow ~98% of possible cases *NOT* observed in train set
- Zero-shot Learning: Unobserved Visual Relationship Detection







mouse, on, cabinet **'ee**, behind, <mark>bear</mark>] horse, wear, hat Figure: The *unobserved* observed relationships are potentially much harder to detect.

STEP 1: TENSORIZE THE VISUAL RELATIONSHIPS

- ▶ Multi-relational tensor $X \in \mathbb{R}^{n \times n \times m}$ given n object categories and *m* possible predicates
- $\rightarrow X(i,j,k)$: number of [Object i, Predicate k, Object j] in train set
- $ightharpoonup \sim$ 2% is non-zero \Rightarrow extremely sparse tensor

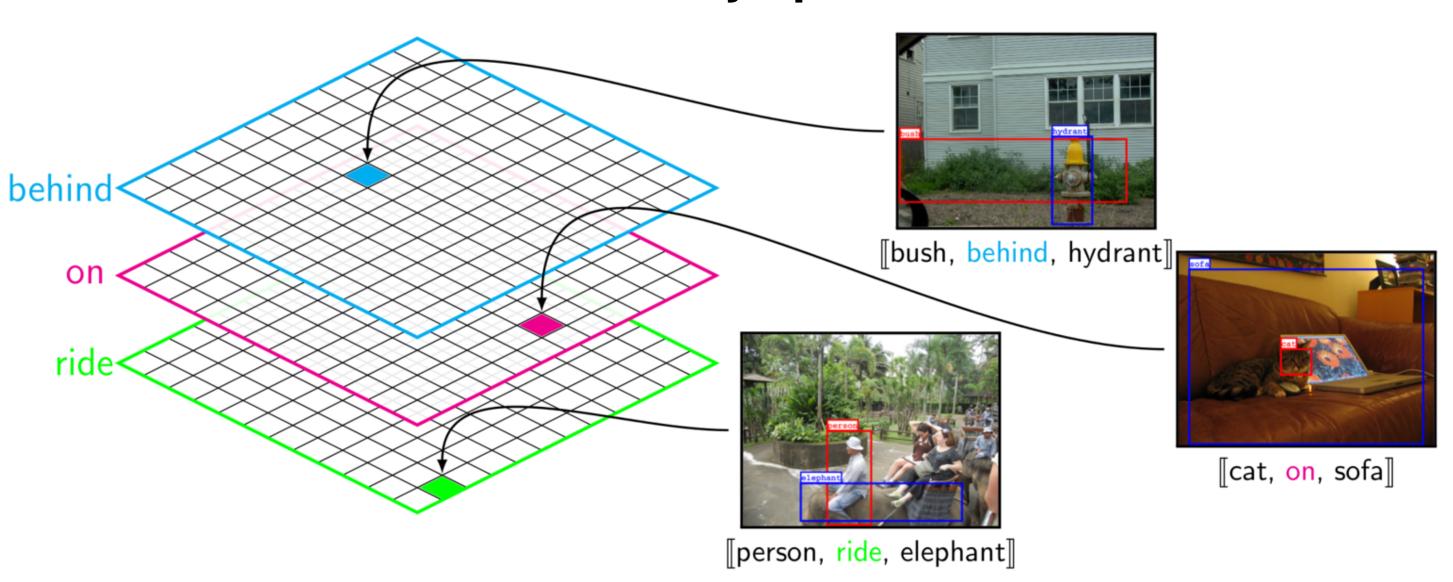


Figure: Each 'slice' X_k encodes possible relationships involving its

STEP 2: FACTORIZE THE RELATIONAL TENSOR

Multi-relational Tensor Factorization

- ► Based on the multi-relational tensor X from train set, derive
- . Common latent representation of objects $A \in \mathbb{R}^{n \times r}$
- 2. **Relationship-specific** factor matrix $R_k \in \mathbb{R}^{r \times r}$ for each relationship $k \in \{1, \ldots, m\}$
- such that $X_k \approx AR_kA^T$

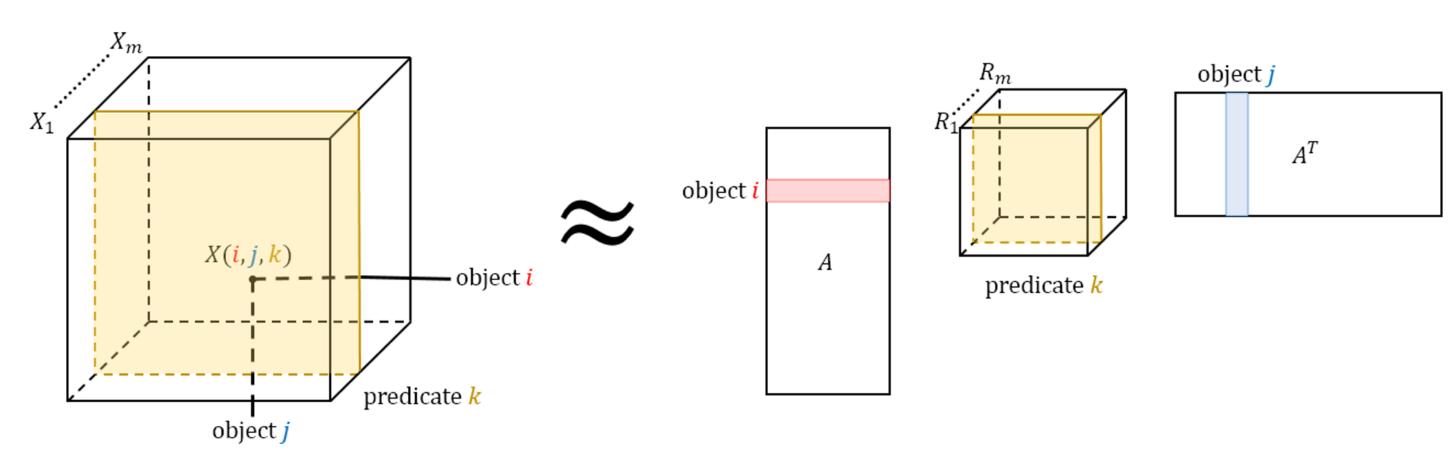


Figure: Multi-relational tensor factorization $X_k \approx AR_kA^T$ for $k \in \{1, ..., m\}$.

MULTI-RELATIONAL TENSOR FACTORIZATION

$$\min_{A,R_k} \sum_{k=1}^{m} ||X_k - AR_k A^T||_F^2$$
 (

. **4**th-order term A: Use auxiliary variables to decouple A and A^T

$$\min_{A,R_k} \sum_{t=1}^m ||X_k - B_k A^T||_F^2 \quad \text{s.t.} \quad B_k = AR_k.$$

2. Low-rank Initialization via SVD: For m = 1, $U\Sigma V^T = X$. Initialize with the "basin of attraction" (Luo et al., Tu et al.):

$$A = V\Sigma^{1/2}, \quad B = U\Sigma^{1/2}.$$

3. Restrict degenerate cases: $A' = AP^{-T}$ and $B' = BP^{-T}$ for any invertible P. Normalize A and B_k :

$$\lambda_p \sum_{k=1}^m ||A^T A - B_k^T B_k||_F^2.$$
 (3)

FINAL FORMULATION

$$\min_{A,R_k,B_k} \sum_{k=1}^m ||X_k - B_k A^T||_F^2 + \gamma \sum_{k=1}^m ||B_k - A R_k||_F^2 + \lambda_p \sum_{k=1}^m ||A^T A - B_k^T B_k||_F^2$$

ALGORITHM: ALTERNATING BLOCK COORDINATE DESCENT

Algorithm 1 Alternating Block Coordinate Descent on (4)

- 1: Given: $X \in \mathbb{R}^{n \times n \times m}$, $X_k := X(:,:,k)$, rank r > 0
- Low-rank Initialization:
- 3: $\overline{X} \leftarrow \sum_{k=1}^m X_k$
- 4: $\overline{U}\overline{\Sigma}\overline{V}^T \leftarrow \text{SVD}(\overline{X},r)$
- 5: $A \leftarrow \overline{V} \overline{\Sigma}^{1/2}$
- 6: for k = 1, ..., m do
- $B_k \leftarrow \overline{U}\overline{\Sigma}^{1/2}$
- $R_k \leftarrow (A^T A)^{-1} (A^T X_k A) (A^T A)^{-1}$
- end for
- ! Iterative descent method:
- while Convergence criteria not met do
- 12: $A \leftarrow \text{gradient descent on (4) w.r.t. } A$
- for k = 1, ..., m do
 - $B_k \leftarrow \text{gradient descent on (4) w.r.t. } B$
- $R_k \leftarrow (A^T A)^{-1} (A^T B_k)$
- end for
- end while
- 18: Output: $A \in \mathbb{R}^{n \times r}$, $B_k \in \mathbb{R}^{n \times r}$, $R_k \in \mathbb{R}^{r \times r}$ for $\forall k$

EXPERIMENT 1: PREDICATE DETECTION

Scene Graph dataset:

- ► 5000 training (< 1% unique relationships), 1000 test images
- ightharpoonup n = 100 object categories, m = 70 predicates
- ► Given: Test images with FastRNN bboxes / labels
- ► Goal: Predict [Object, Predicate, Object] with 3 tasks

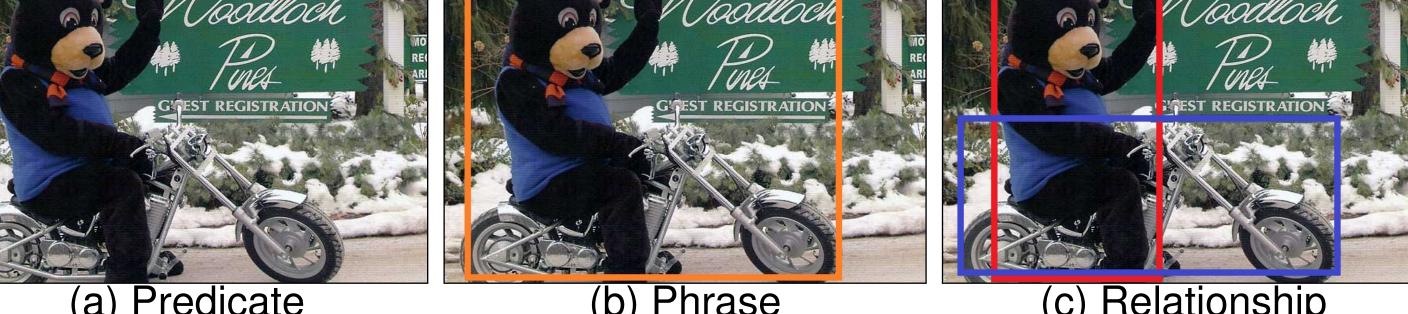
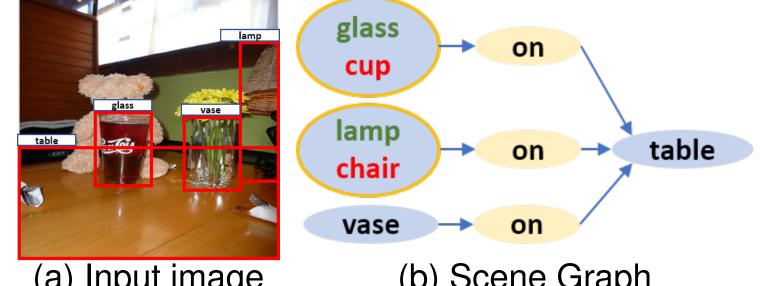


Figure: Detection Task Conditions: (a) Predicate (easy): does not require bounding boxes. (b) Phrase (moderate): requires relationship bounding box e) containing both objects. (c) **Relationship** (hard): requires individual bounding boxes (red/blue).

EXPERIMENT 2: SCENE GRAPH DETECTION

Visual Genome dataset:

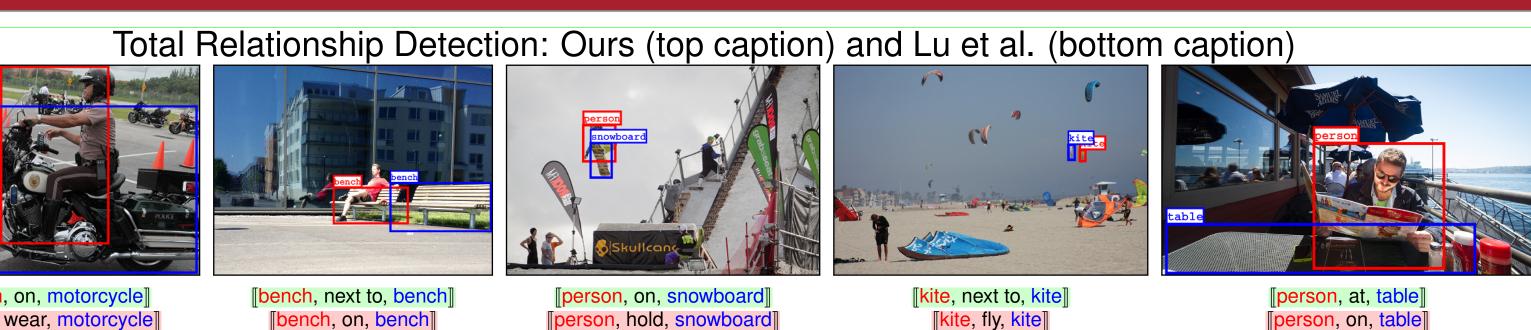
- ► 108,077 rels. (< 2% unique relationships), 70% training, 30% test
- ightharpoonup n = 150 object categories, m = 50 predicates
- ► Given: Test images with Region Proposal Network bboxes / labels
- ► Goal: Predict Scene Graphs with 3 standards

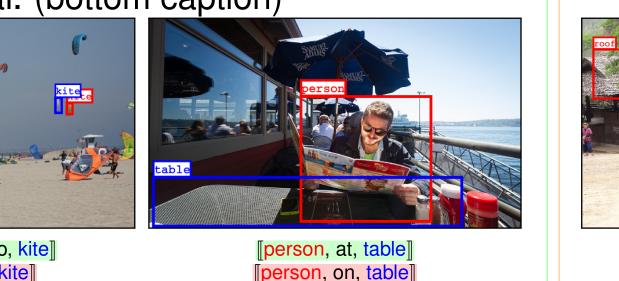


Prediction Tasks Pred. Obj. B-box Predict Predicate (**PredCls**) Classify SG (**SgCls**) Generate SG (**SgGen**)

Figure: Scene graph detection tasks. Check marks indicate required prediction components. The tasks become incrementally more demanding from top (PredCls) to bottom (SgGen).

EXPERIMENT 1: RESULTS











on Scene Graph dataset using our algorithn

(top caption) and Lu et al. (bottom caption). The correct and incorrect predictions are highlighted in green and red respectively.

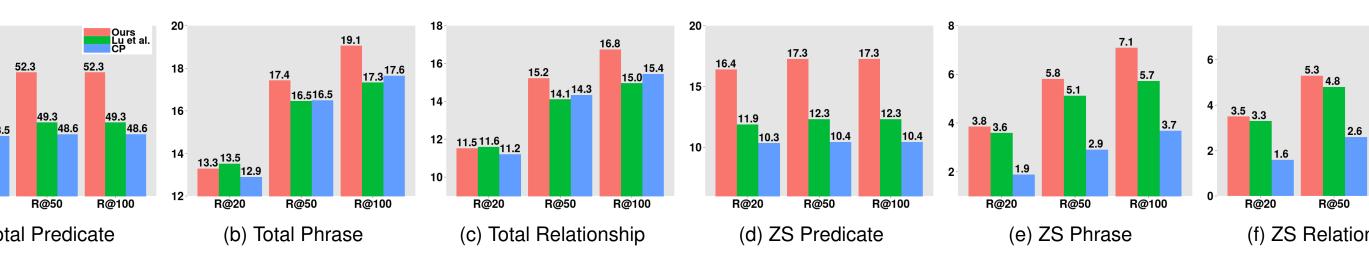
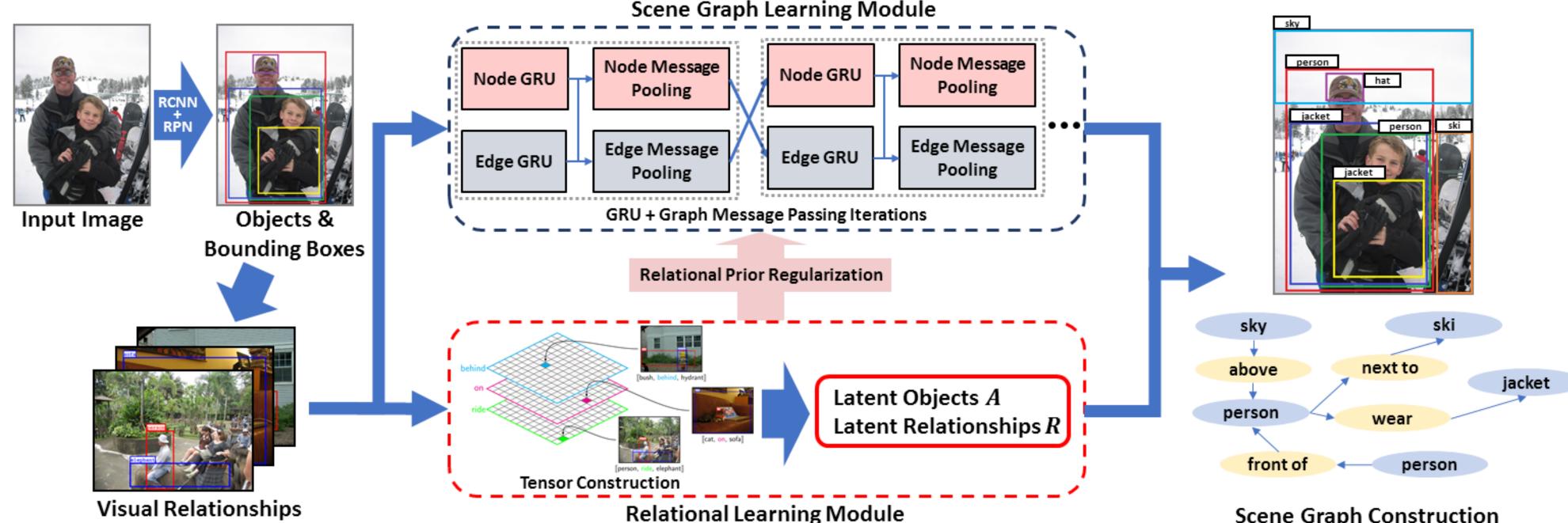


Figure: Left: Examples of semantically confusing relationships of the ground truth (top caption), ours (middle caption) and Lu et al. (bottom caption). The left three relationships (green box) have varying detection results, but they can all be considered to be semantically correct. However, the right two relationships (orange box) have disagreeing interpretations of the relationships. Right: Visual relationship detection results on Scene Graph.

STEP 3: REGULARIZE THE SCENE GRAPH LEARNING PIPELINE

Scene Graph detection: Predicate detection ⇔ Object detection

- ▶ Module 1: Scene Graph Module (Xu et al.) \Rightarrow Faster-RCNN \rightarrow GRU \rightarrow Message Pooling \rightarrow Exchange Expressive but highly dataset dependent (i.e., prone to outliers and mislabels)
- ► Module 2: **Relational Learning Module** (Ours) ⇒ Faster-RCNN → Multi-relational Tensor Factorization Robust but does not infer objects by construction
- Pre-trained and fixed during training



Balance the learning modules

- Bring the best of both worlds!
- Regularize with global prior from RL module

PROBABILISTIC GLOBAL PRIOR INJECTION Given (i) SG module prediction k_{SG}^* , (ii) RL module prediction k_{RI}^* , and (iii) a 'y-or-1' filter

 $D(k_{BI}^*(i), \theta) = k_{BI}^*(i)$ with probability θ , we obtain the regularized prediction k^* as follows:

 $oldsymbol{k}^* = oldsymbol{k}_{SG}^* \odot oldsymbol{D}(oldsymbol{k}_{BL}^*, heta)$ where \odot is a Hadamard product.

Result: Increase underestimated scores (i.e., rare relationships) or decrease overestimated scores. High θ , more frequent influence of prior k_{BI}^* .

Training the pipeline:

- 1. 'Warm start' with *only* SG for 100K iters ($\theta = 0$)
- 2. **Regularize** k_{SG}^* with k_{BI}^* to obtain k^* as 5 with $\theta = 0.2$ for 50K iters (RL module is fixed)

EXPERIMENT 2: RESULTS

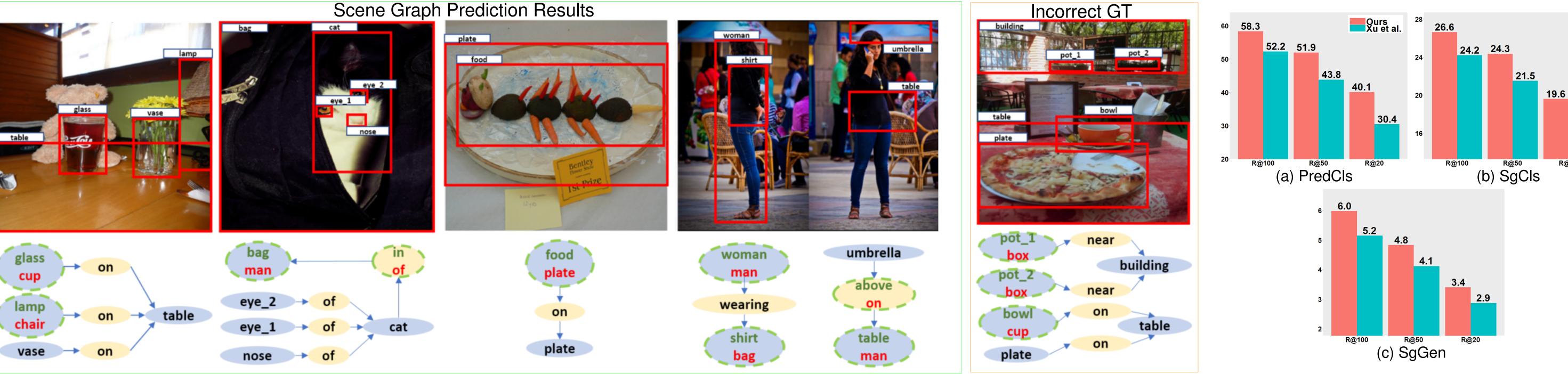


Figure: Left: For each column, the predicted objects (blue ellipses) and their relationships (yellow ellipses) are constructed as a scene graph of its top image. The bounding boxes labels reflect our prediction results. For difficult predictions (green dashed boundary) where our model has correctly predicted (top green) and while Xu et al. has misclassified (bottom red) are shown. The rightmost column is an example of a case where our model provides more accurate predictions (pot and bowl) than those of the ground truth (box and cup). Right: Scene graph detection task results on Visual Genome.

ACKNOWLEDGMENT

SJH was supported by a University of Wisconsin CIBM fellowship (5T15LM007359-14). We acknowledge support from NIH R01 AG040396 (VS), NSF CCF 1320755 (VS), NSF CAREER award 1252725 (VS), UW ADRC AG033514, UW ICTR 1UL1RR025011, Waisman Core grant P30 HD003352-45 and UW CPCP AI117924 (SNR).

Conference on Computer Vision and Pattern Recognition (CVPR) 2018