
Image Retrieval with Mixed Initiative and Multimodal Feedback, BMVC '18

CS688 Paper Presentation

2018/11/22

20174315 Chiwan Song

KAIST

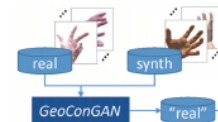
The KAIST logo consists of the letters "KAIST" in a bold, blue, sans-serif font. Below the text is a light blue, horizontal oval shape that serves as a shadow or base for the letters.

GANerated Hands for Real-Time 3D Hand Tracking from Monocular RGB, Mueller et al., CVPR '18

Conclusion & Summary

- Presents a more robust model for occlusions
- Presents
 - a data set similar to the real hand domain
 - a model that can create the data set
- Demonstrates these benefits in the evaluation
 - particularly in difficult occlusion scenarios.
- Summary
 - Real-time full 3D hand tracking from single monocular RGB video.
 - Technical Novelties

1)



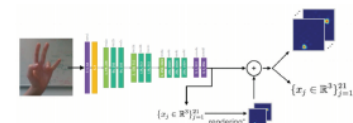
New GAN for **geometrically consistent unpaired image-to-image translation**

2)



Novel **enhanced RGB dataset** with **3D hand joint annotations** (>260k frames)

3)



CNN with projection layer for tightly coupled regression of **2D and 3D joint locations**

Table of contents

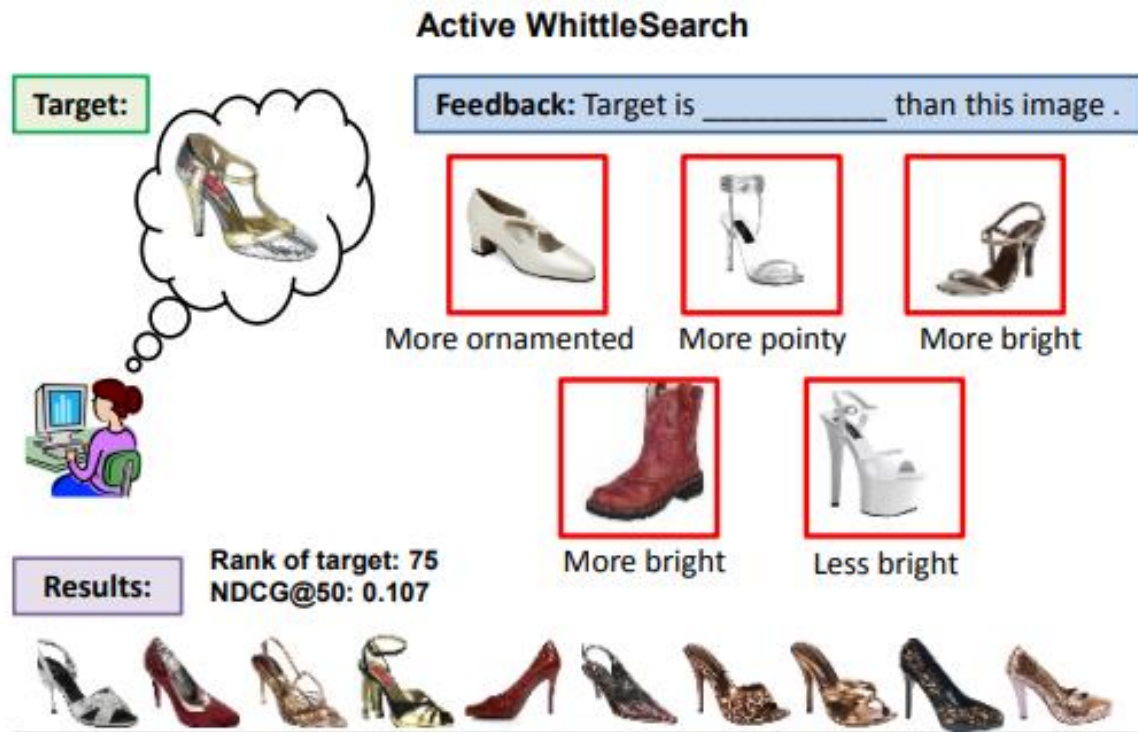
- Introduction
- Backgrounds
- Main idea
- Experiment & Result



Introduction

Giving user's mental concept to a system

- The users can give language-based guidance to the system



Giving user's mental concept to a system

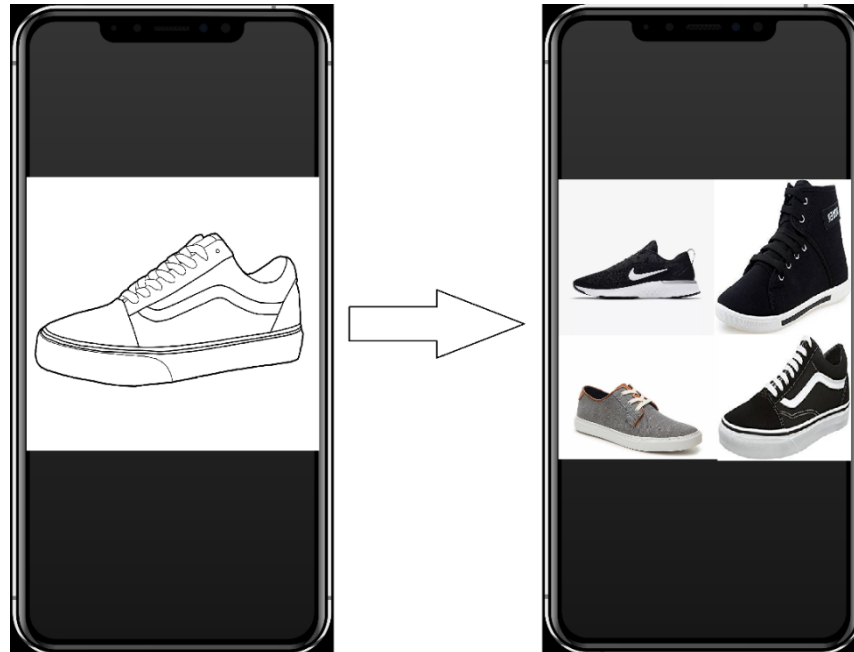
- The system asks to users to derive what users want



Figure 6. Using the user's feedback on the left, we retrieve the images on the right at the top of the results list.

Giving user's mental concept to a system

- The users can provide visual cues for what they are looking for
 - i.e. Sketch based image retrieval



Problems in previous work

- The interaction was driven by user **or** system, but **not both**
- The single modality of user input - language or visual

Main contribution of the paper

- A new framework that the interaction can be driven **either** the user **or** system
 - By agent trained with **Reinforcement Learning**
- The **dynamic decision** of interaction methods
 - Sketch, free-form attribute feedback, or answering attribute-based question

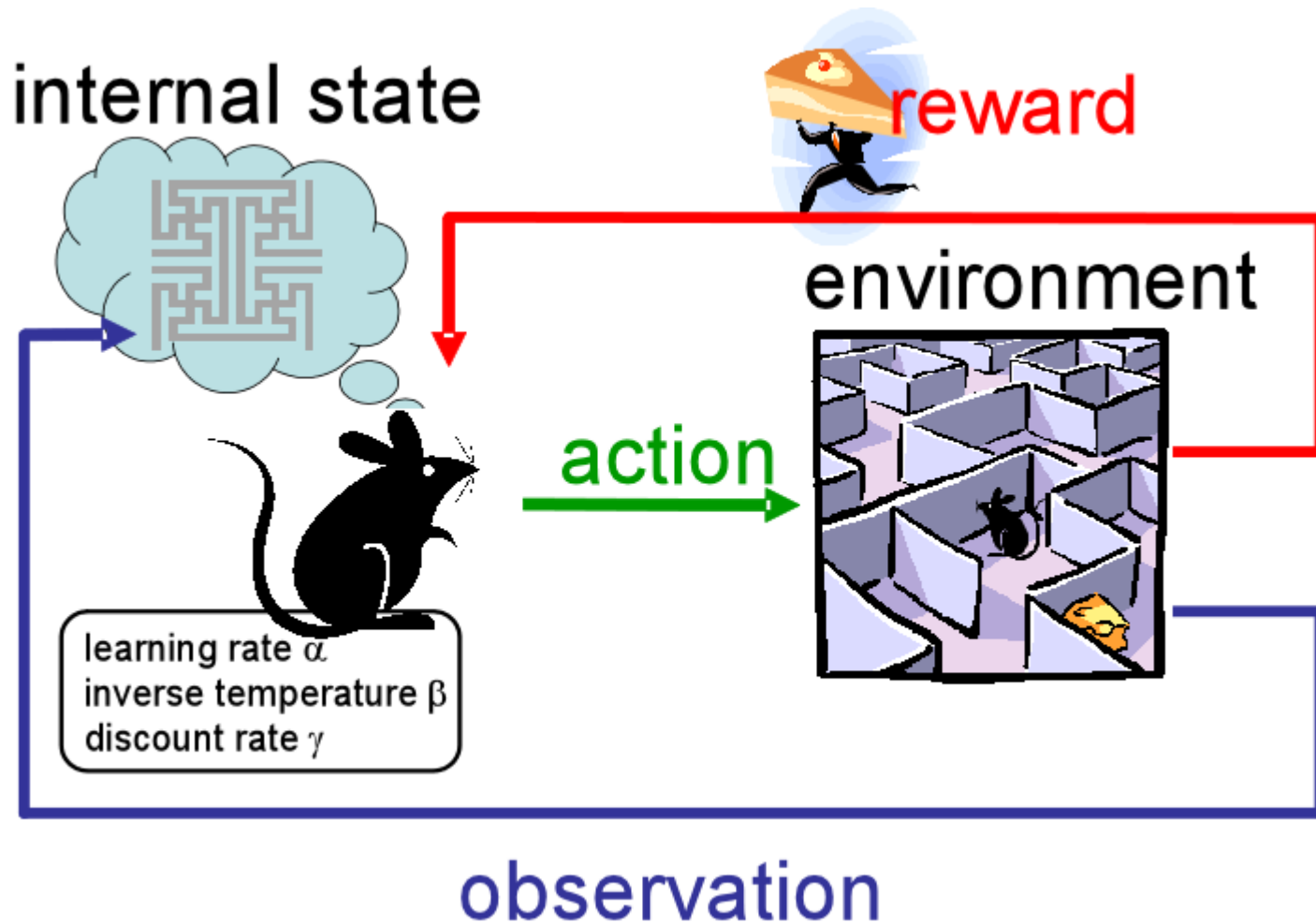


Backgrounds

Concept of Reinforcement Learning

- Inspired by the psychology of behavior
- Consist of **agent, environment, state, action, and reward**
- The **agent** changes **state** by doing **action**
- The **environment** gives **reward** to **agent** according to the **action**

Concept of Reinforcement Learning



Properties of Reinforcement Learning

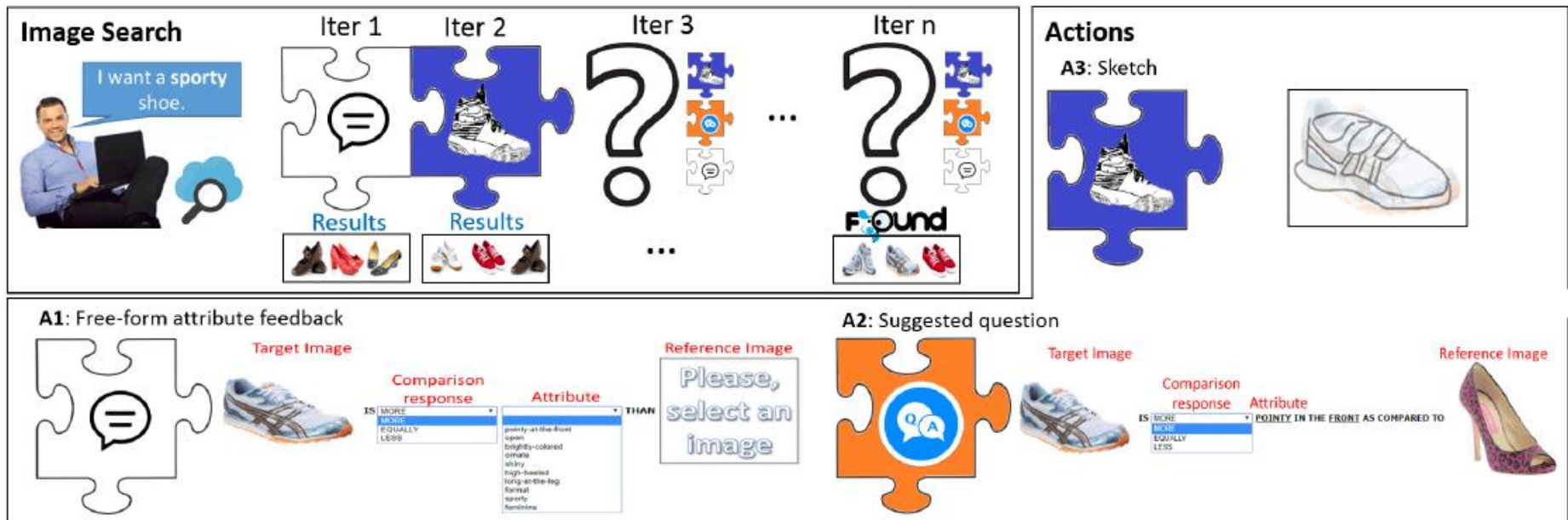
- The **agent** is trained for maximizing the **reward**
- The **agent's action** effects its next input
- No supervisor present



Main idea

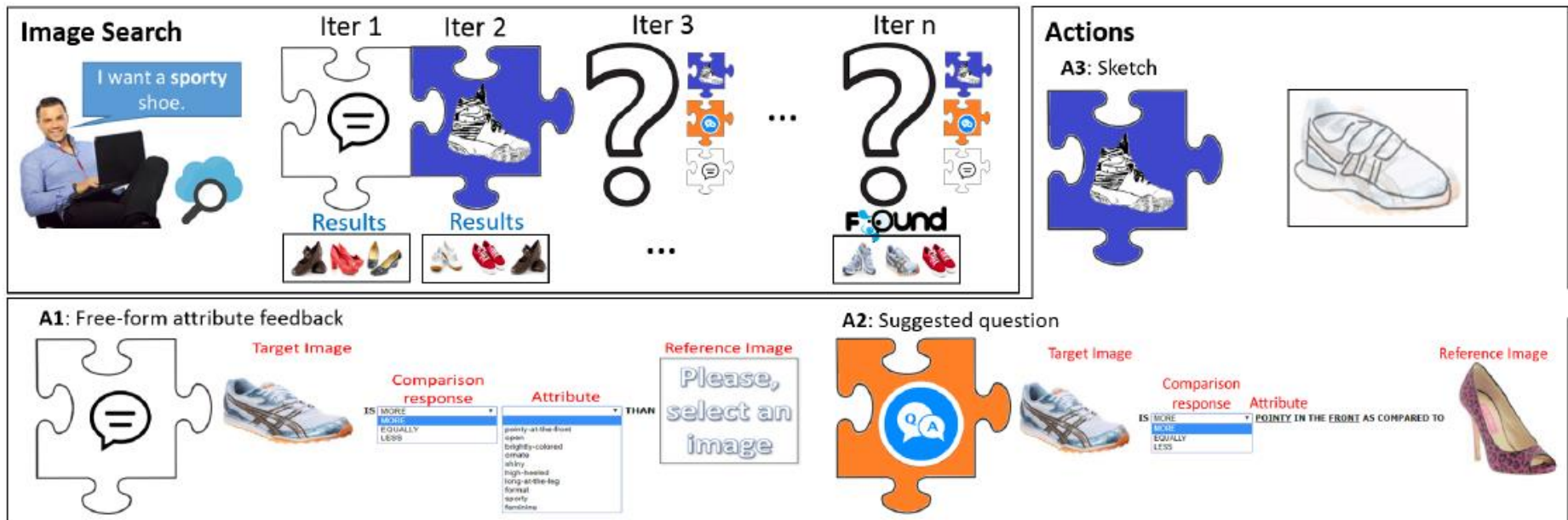
Overall process of image retrieval

- The user initiates a search with images or simple query



Overall process of image retrieval

- The system chooses the **best feedback method** for retrieving the target image




Refining the image retrieval


- Use (attribute, reference image, comparison response) triplet or a sketch
- The feedback is used to update the next top K results

(a) Attribute feedback


Target Image




IS

Comparison response	Attribute	THAN
<input type="button" value="MORE"/>	<input type="button" value="pointy-at-the-front"/>	
<input type="button" value="MORE"/>	<input type="button" value="open"/>	
<input type="button" value="EQUALLY"/>	<input type="button" value="brightly-colored"/>	
<input type="button" value="LESS"/>	<input type="button" value="ornate"/>	
	<input type="button" value="shiny"/>	
	<input type="button" value="high-heeled"/>	
	<input type="button" value="long-at-the-leg"/>	
	<input type="button" value="formal"/>	
	<input type="button" value="sporty"/>	
	<input type="button" value="feminine"/>	

(b) Sketch feedback



Current Top Results



Refining the image retrieval

- Sketch feedback is converted to photographs by GAN



Figure 6: Sample sketch-to-photo colored images for Pubfig (columns 1-3), Shoes (columns 4-6), and Scenes (columns 7-9). Each column denotes a different category.

Training the agent with RL

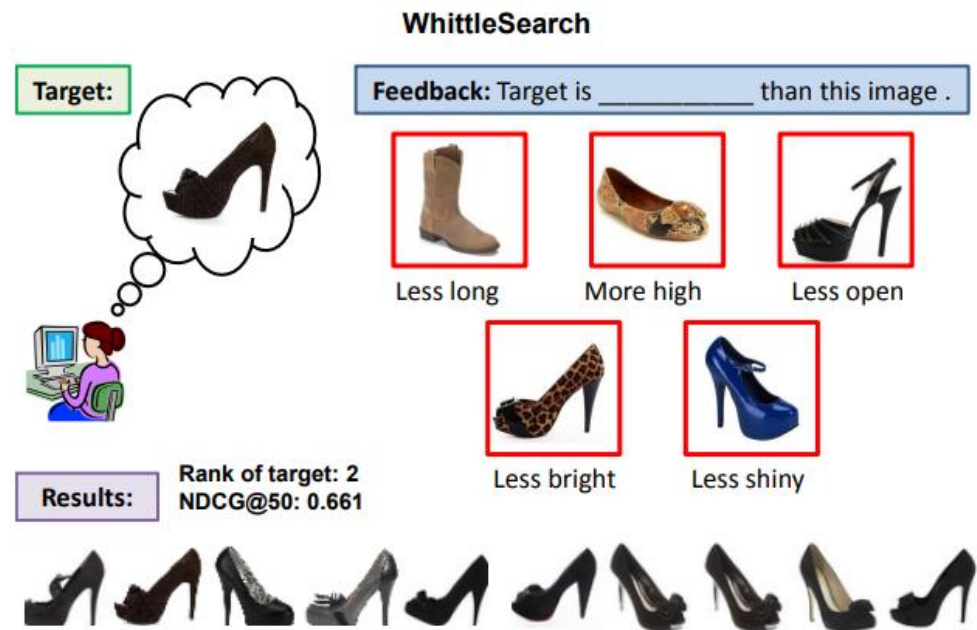
- State - (History of top 1 image, 5 closest neighbor, 5 furthest neighbor, previous action)
- Action - Selecting interaction method



Experiment & Results

Experiment setup

- Experiment on 3 dataset
 - Pubfig, Scenes, Shoes
- Comparing with 3 baselines
 - Whittle search (WS)



Experiment setup

- Experiment on 3 dataset
 - Pubfig, Scenes, Shoes
- Comparing with 3 baselines
 - Whittle search (WS)
 - Pivot round robin (PRR)



Experiment setup

- Experiment on 3 dataset
 - Pubfig, Scenes, Shoes
- Comparing with 3 baselines
 - Whittle search (WS)
 - Pivot round robin (PRR)
 - Sketch retrieval + PRR (SK_PRR)
- Experimented by simulated and real users

Results

- Percentile rank plots for each dataset

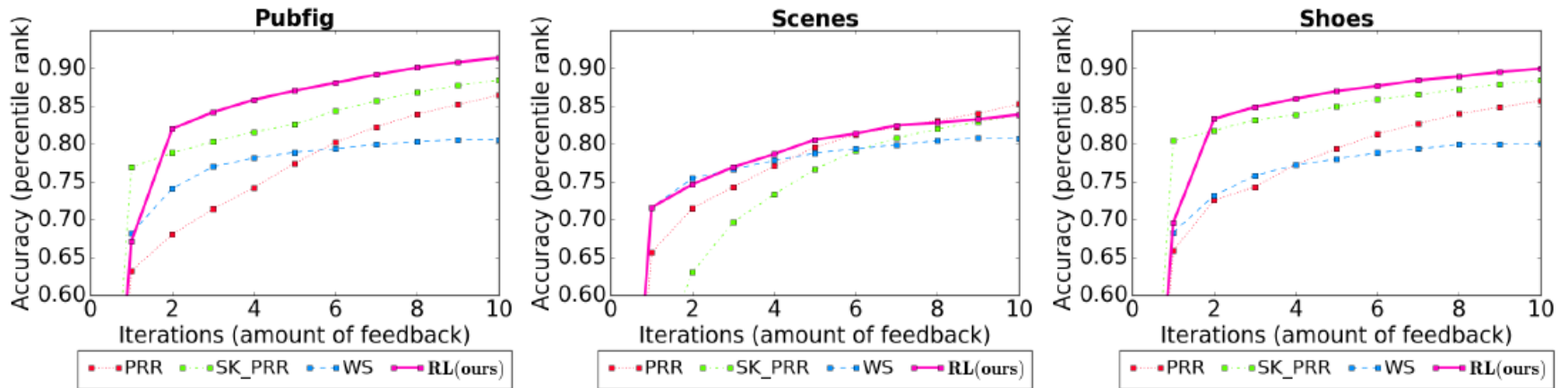


Figure 4: Percentile rank plots for Pubfig, Scenes, and Shoes. Our mixed-initiative RL agent outperforms the other baselines on Pubfig and Shoes, and performs competitively for Scenes.

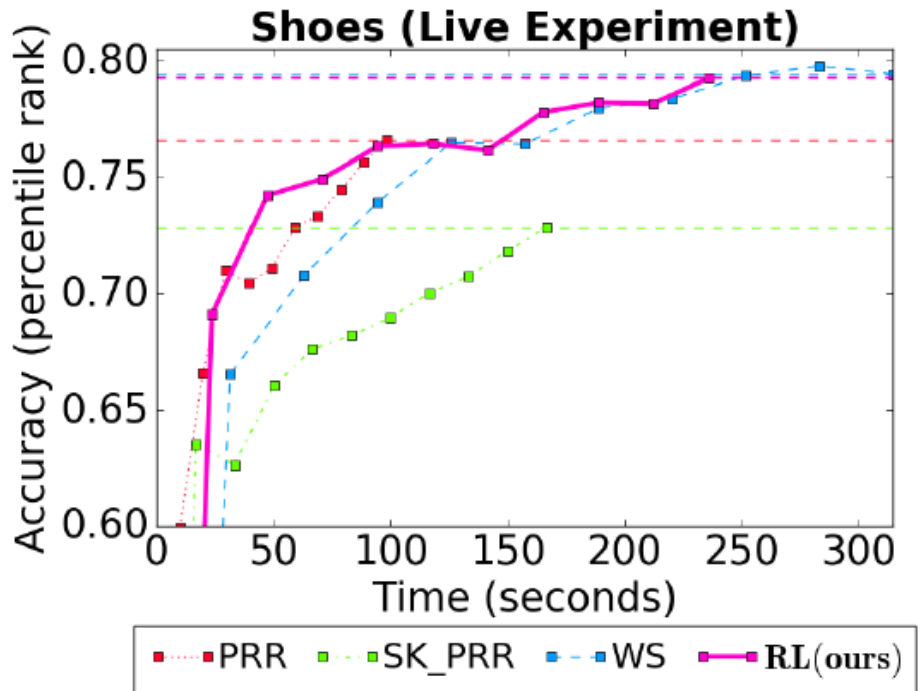
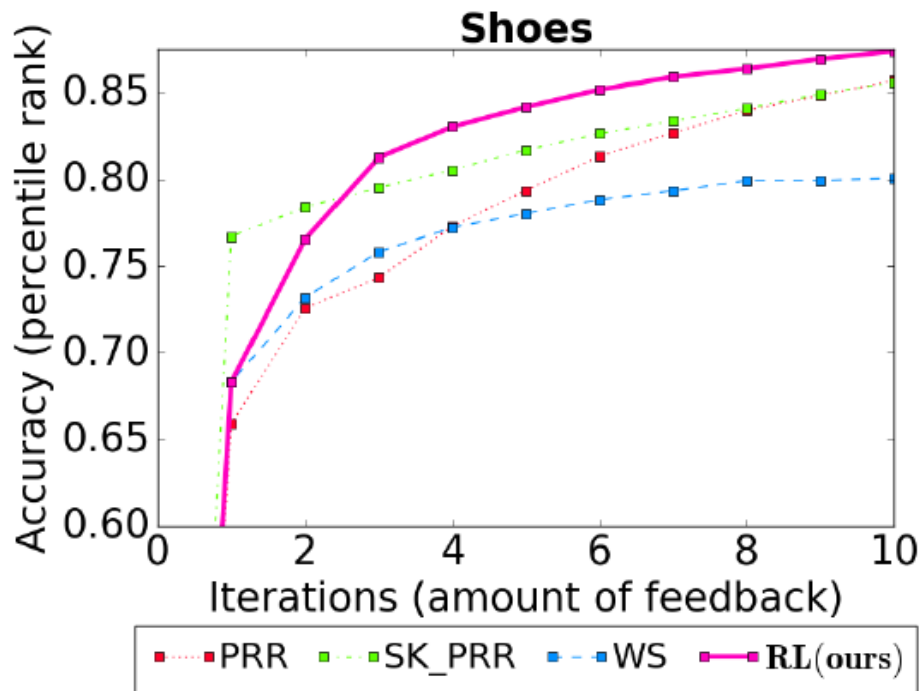
Results

- AUC for percentile rank curves

	PRR [16]	WS [19]	SK_PRR [16, 54]	RL (ours)
Pubfig	0.729	0.737	0.789	0.810
Scenes	0.741	0.741	0.699	0.754
Shoes	0.745	0.731	0.806	0.810
avg	0.738	0.736	0.764	0.791

Results

- Live experiment result





Thank you!