**Project A: Figure Annotation Interface**

**Overview**

The goal of this project is to build a web-based figure annotation user interface (UI). The UI will support annotation of compound figure segmentation and semantic information extraction. The figures are from US design patents. The UI has landing page including a text search box and an advanced search box that allows users to search figures by specifying multiple conditions. The advanced search box also allows registered users to specify figures that have been assigned, assigned to the user, and assigned but not annotated. When a user clicks an item in the search engine result page, an annotation page will show up and displays two panels. The panel on the left displays the original compound figure and basic metadata. The panel on the right displays the segmented figures and automatically extracted information, including objects and viewpoints. Because the segmentation and information extraction were performed automatically, the program may make mistakes. The user can mark whether the segmentation and information extraction results were correct or not. If the results are wrong, the user can make corrections. For example, if a compound figure contains 2 subfigures, but was mistakenly segmented to 3, the user needs to specify the correct number of subfigures contained in the original figure. If the caption of the original figure contains the object "high heel shoe" but this phrase was not extracted, the user should be able to type the correct object name in a text box. If the caption of the original figure contains the viewpoint "top view" but this phrase was not extracted, the user should be able to type the correct viewpoint (if any) in a text box. The website should allow an admin login and a regular user login. An admin can assign annotation tasks to users.

The UI search engine will implement other features such as user registration, log-in, reCAPTCHA, spell-check, auto-complete, customized list, like button, etc. More details will be specified in milestone specifications.

**Data**

The students who do this project will be provided a dataset containing the following parts.

1. Metadata of 1400 compound figures (aka original figures) and around 2500 segmented figures, including patent ID, figure ID, figure caption, figure file name, segmented figure file name, automatically extracted object and viewpoints, etc.
2. The PNG files of compound figures and segmented figures.

**Infrastructure**

The students will build the web applications on a virtual machine in the CCI academic environment. The search function should be powered by Elasticsearch. The user information, metadata and annotation results should be stored in MySQL.