*Yujie Lu*
*Muge Chen*
*Yanqiong Chen*

Google ML Winter Camp
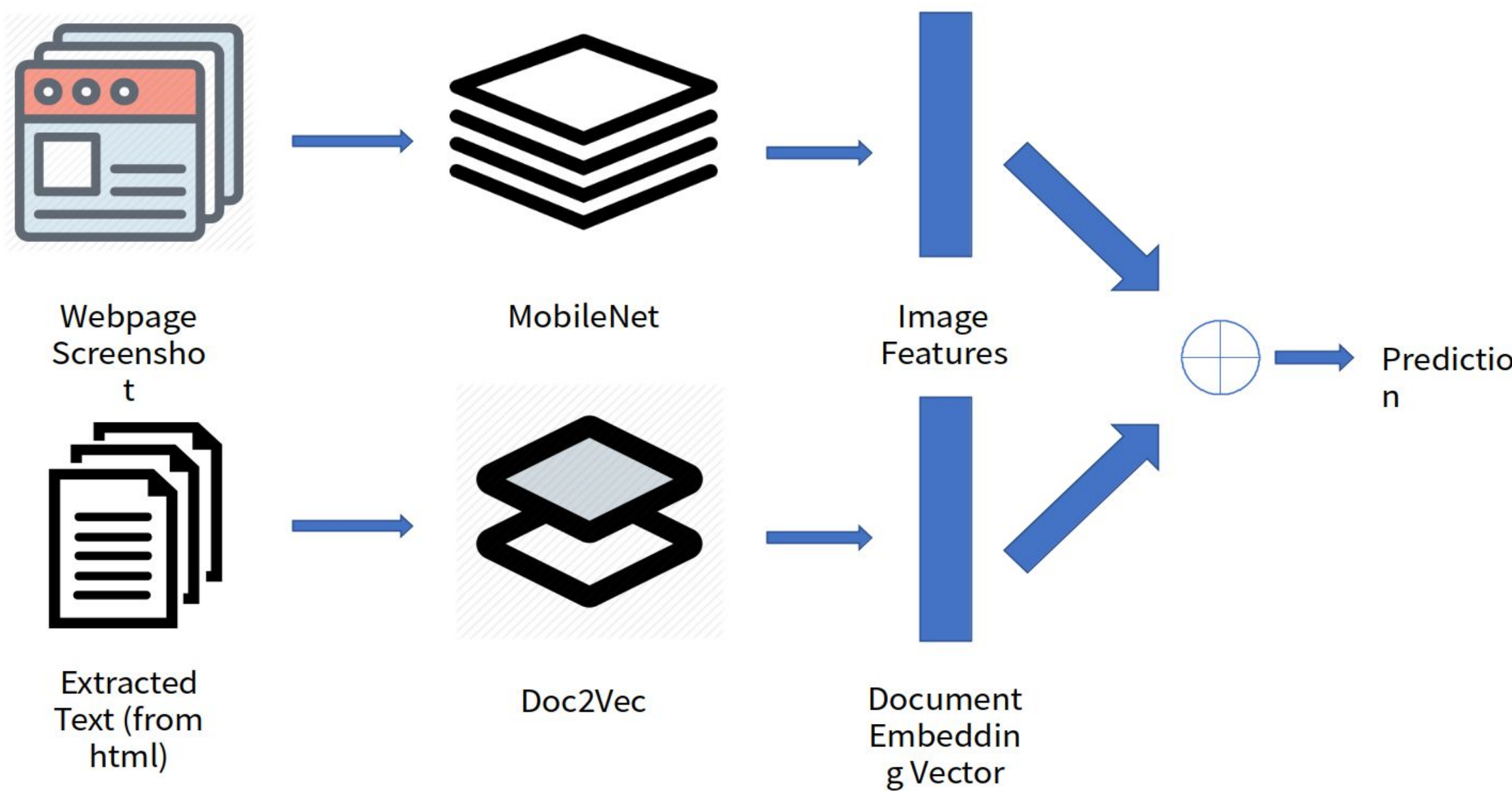谷歌机器学习应用冬令营

# Webpage Wizard
## - A Useful Model for Webpage Classification

• Introduction

In the world of machine learning, most of the application domain suffers from not having sufficient labeled data whereas unlabeled data is available cheaply.

In our project, we try to do webpage classification on a dataset that **only a small amount of data is labeled.**

• Pipeline



Webpage Screenshot → MobileNet → Image Features

Extracted Text (from html) → Doc2Vec → Document Embedding Vector

⊕ → Prediction

• Implementation Detail

**See github link**

https://github.com/foreseeable/Aelous

• Applications

① Our Demo

```
In [54]: def Predict(img_num, img_clf, show_img=True):
             file_name = 'render' + str(img_num) + '.png'
             this_x = process_train(os.path.join(input_dir, file_name), show_img=show_img)
             t1 = np.ndarray([1,224,224,3])
             t1[0] = this_x
             max_id = np.argmax(img_clf.predict(t1))
             print('Predicting category of ' + url[img_num])
             print('This render is most likely to be ' + cate[max_id])

         Predict(24, img_clf, True)
```



Predicting category of https://brobible.com/life/article/54952/
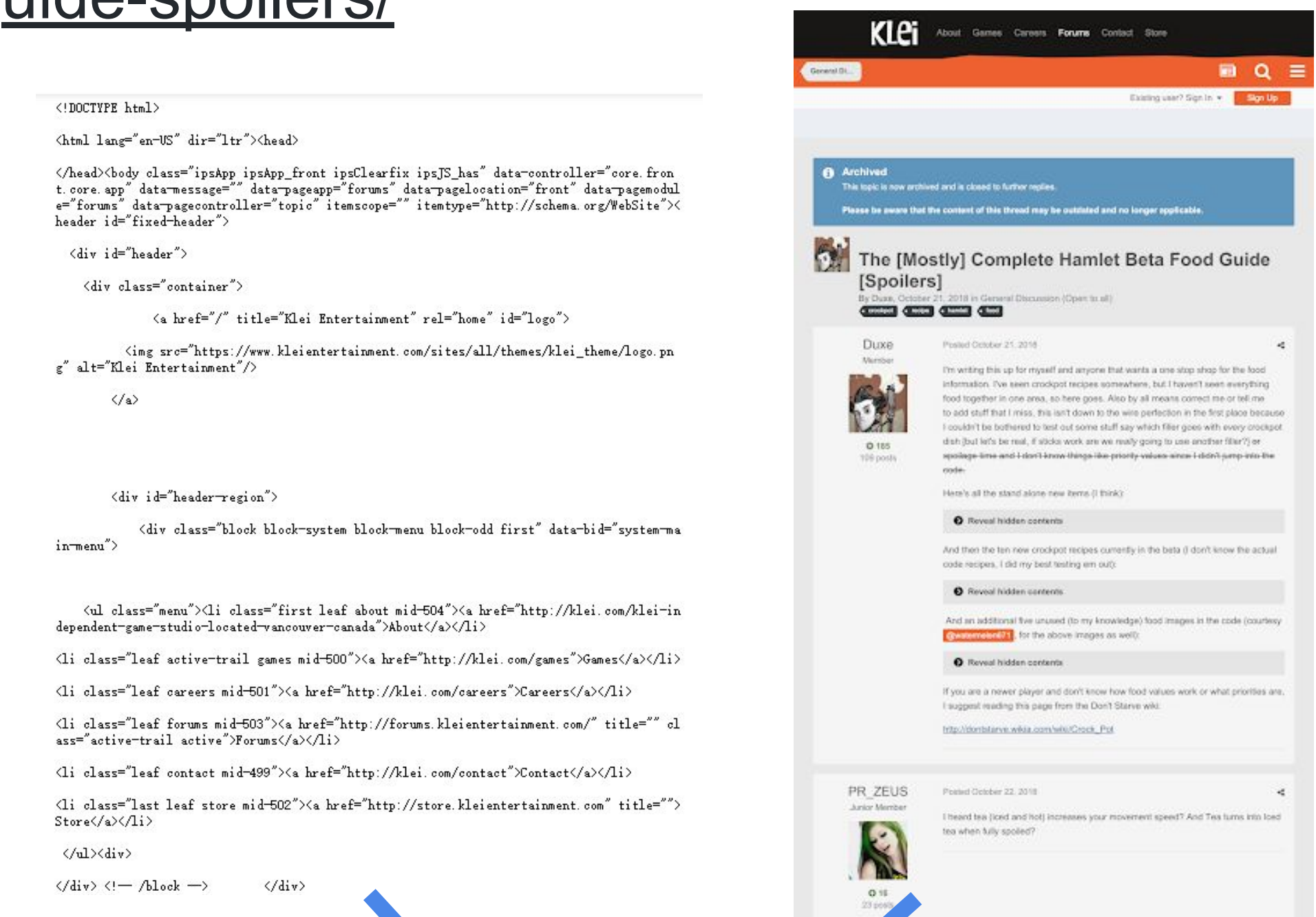This render is most likely to be article

• Dataset:

10K data: url + html + screenshot
2590 with label **is_entity**
800 with label **is_entity** & **category**
**10 categories:**

- media_introduction
- others
- location
- social_media_profile
- encyclopedia
- qa_forum
- shopping_item
- list
- media_player
- article

example:

https://forums.kleientertainment.com/forums/topic/97192-the-mostly-complete-hamlet-beta-food-guide-spoilers/



qa_forum

② Other Promising Applications

- Assist building an extension for Chrome which can beautify the UI with different strategies according to the category of the webpage. Similar to switch omega.

- Help the browser collecting information about which category the user visits most frequently to decide what ADs to present.

投 票 区 域