# Computing Environment for Course Project

Chen Zhao

# CodaLab

# CodaLab

# CodaLab

- CodaLab set up:

  - Create your group account (both CodaLab and class)

  - Install command-line interface (Recommended)

  - Submit!

# CodaLab

- Run official evaluation locally

  - Use the evaluation script we provided

  - Run the trained model on dev set

- Submit your model to the leaderboards

  - Docker submission

  - We test the model on private test set, then update the leaderboard

# Docker

- A modern platform for AI applications

  - Automate the deployment of software applications inside **containers**

  - Allow users to package an application with all of its dependencies into the container

  - Applications run on **different OS** via a layer of indirection

# Docker

# Docker

- We will provide a basic docker image

    - All dependencies for the baseline model

    - Include other dependencies in your code

    - We also provide some practice for using Docker

- You can customize your docker image from scratch(not recommended)

# Course Project

- Overview

- Input/Output

- Dataset

- Baseline

- Evaluation

# Course Project - Overview

- Incremental

  - We have several sub-tasks as homework (Deep Learning, Sequence Tagging..)

  - Final project will benefit from the previous sub-task code/results

  - Sub-task will use docker for submission

    - We will provide some hands on practice for getting used to Docker in sub-tasks

  - Final course project will use both docker and CodaLab

# Course Project - I/O

- Example I/O

  - Input: {112; At its premiere, the librettist of this opera portrayed a character who asks for a glass of wine with his dying wish; False; WAIT}

  - Output: {The Magic Flute; BUZZ }

- Input:

  - Index (character position)

  - Question text up to this position

  - Is_new_sent

  - Opponent Action (probably BLANK)

- Output

# Course Project - I/O

- Input:

- Output:

  - Top-1 guess

  - Action (buzz or wait)

# Course Project - Dataset

- Data locates at  https://github.com/Pinafore/qb#downloading-data

- Training/dev/test

  - Training (questions, correct answer w/wiki mapping)

  - Two evaluation datasets

    - Held-out data (from different tournaments)

    - Adversarial questions (more challenging)

      - Avoid clues that are easy for computers to answer

# Course Project - Baseline

- Baseline system (2017 NIPS competition)

  - Guesser

    - IR system (TF-IDF) keyword-matching

    - 31862 QB questions, 6991 entities

    - Given a query(question), compare to previously asked questions, sort the answer entities based on similarity

    - TOP1 Accuracy 0.55

  - Buzzer

# Course Project - Baseline

- Baseline system (2017 NIPS competition)

  - Guesser

  - Buzzer

    - Based on the guesser score

    - Threshold

      - Top-score exceeds the threshold(0.3)

# Course Project - Evaluation

- Evaluation

  - Accuracy (Based on top-1 guess)

  - Head-to-head competition

    - Evaluate on same output, but like real competition

    - Check the answer of model that buzz first, if not correct, check the other model's answer at the end

    - 10 points for correct answer, -5 for incorrect if not at the end of question