# Crawling Hupu's NBA Game Comments?

## Overview

This article explains how to use Scrapy to crawl and scrape Hupu's hottest comments (高亮评论) after NBA games, along with player statistics and basic game information such as scores.

## Background

I've always wanted to fine-tune an LLM to mimic Hupu JRs' commenting style, which is typically concise, harsh, and direct in both criticism and praise. These comments are closely tied to game statistics and the specific players being discussed (for example, LeBron James haters have particular patterns when mocking him).
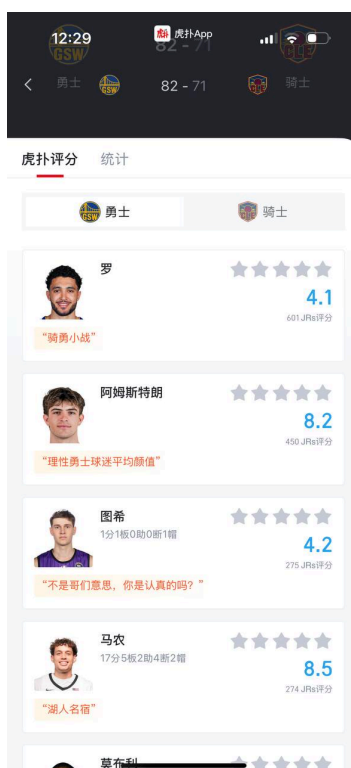
To build an effective training dataset for the LLM, we need to collect game information, player statistics, and the comments about individual players.

## Method

### Context

Hupu can be accessed in two ways: through the mobile app or through the website (www.hupu.com). While Hupu's website provides a comprehensive database of game statistics, it lacks the game review and highlighted comment sections. The comment section I'm targeting appears under "虎扑评分" (as shown in the screenshot) and is only available in the Hupu mobile app.
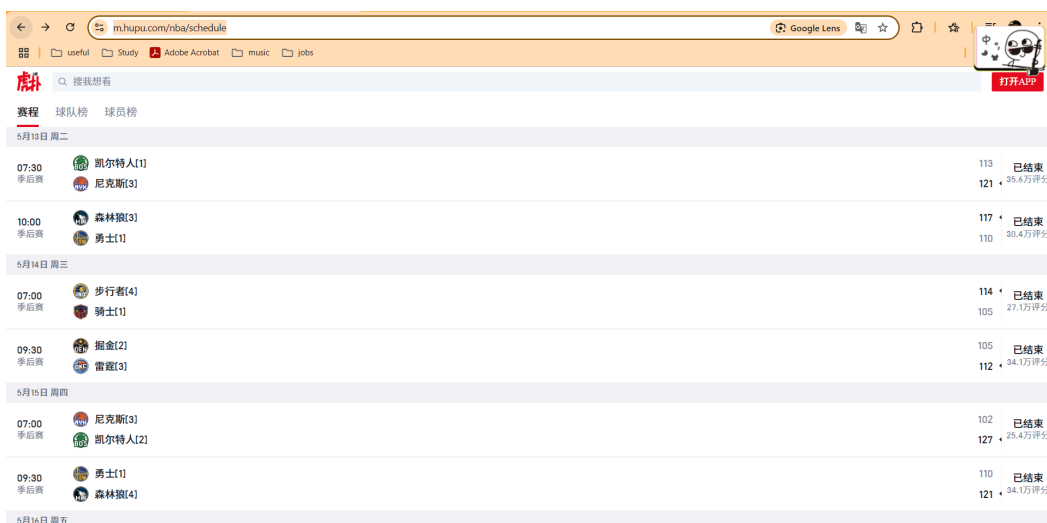
Consequently, we need to crawl data from the Hupu app, which is more challenging than scraping the website since it's typically harder to locate the API endpoints that the app uses.
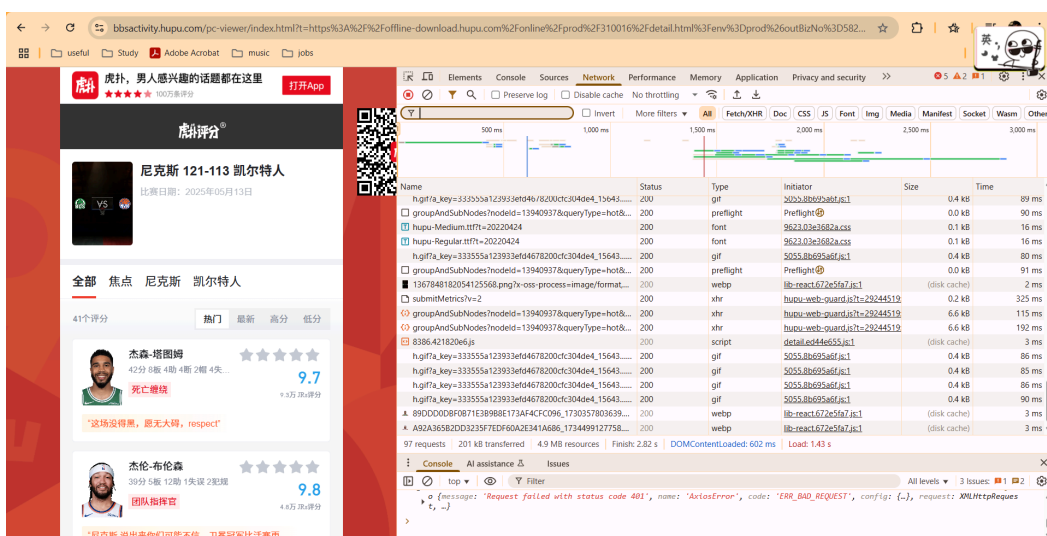
Fortunately, there is a website version of Hupu's mobile app ([m.hupu.com](m.hupu.com)). This allows us to identify the endpoints where Hupu fetches the data.

## Crawling Process

- First, we navigate to the NBA schedule page ([https://m.hupu.com/nba/schedule](https://m.hupu.com/nba/schedule)). Note that this page only displays games within a limited time range, with most games hidden.
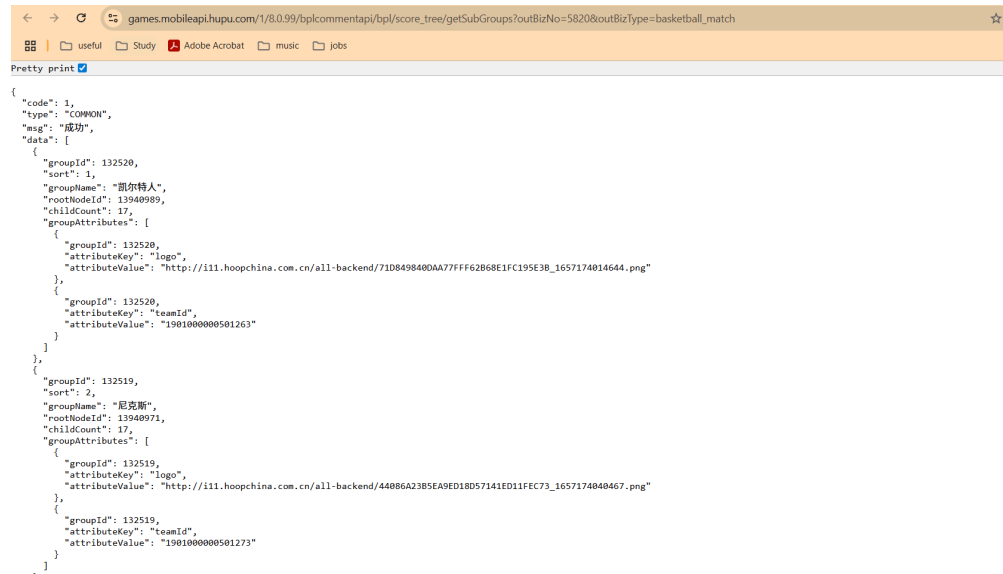


- Next, we select any game, open the developer tools, and examine the network tab to analyze the traffic:



- Several notable XML requests appear:

    - [https://games.mobileapi.hupu.com/1/8.0.99/bplcommentapi/bpl/score_tree/getSubGroups?outBizNo=5820&outBizType=basketball_match](https://games.mobileapi.hupu.com/1/8.0.99/bplcommentapi/bpl/score_tree/getSubGroups?outBizNo=5820&outBizType=basketball_match)

Based on observation, the **outBizNo** parameter appears to be a numbering system corresponding to the **groupID**. We will find out what **groupID** is in the next step.



- https://games.mobileapi.hupu.com/1/8.2.99/bplcommentapi/bff/bpl/score_tree/groupAndSubNodes?nodeId=13940937&queryType=hot&page=1&pageSize=10
  By inserting the groupID into the "nodeID=<groupID>" parameter, we can retrieve a JSON response containing information about a specific game (see attached screenshot). Note that this only shows one team's stats, meaning each game has two different groupIDs—one for each team.

  From this JSON response, we can extract:
  1. Game scores for each team
  2. Individual player statistics
  3. Each player's individual **bizID** (which becomes crucial later)
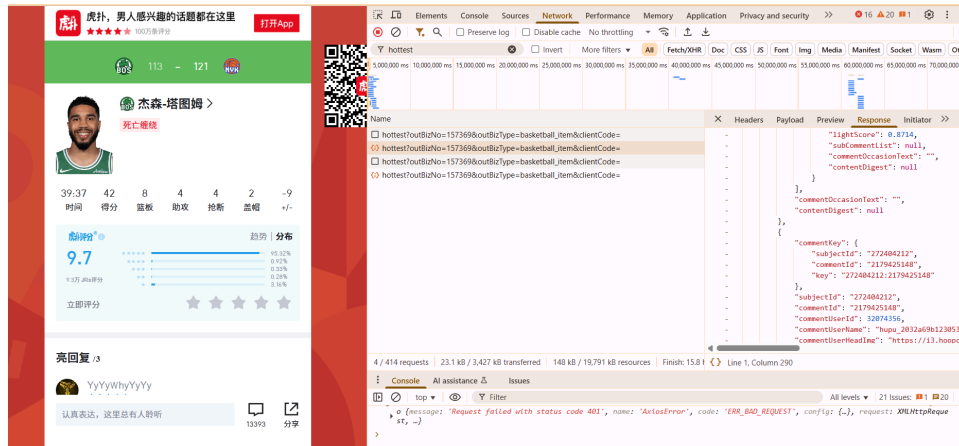
  However, this page only displays one top comment. Since I want to collect the top 3 comments for each player, I need to find something that display all the comments.

- https://games.mobileapi.hupu.com/1/8.2.99/bplcommentapi/bpl/comment/list/primarySingleRow/hottest?outBizNo=157369&outBizType=basketball_item&clientCode=
  After clicking into each player's individual tab (see screenshot), I observed this XML request appearing in the network traffic. The **outBizNo** here corresponds to the **bizID** of the player.



- The response (see screenshot below) contains all the comments from Hupu users about this specific player in this game, with each comment accessible through the "commentContent" parameter.

games.mobileapi.hupu.com/1/8.2.99/bplcommentapi/bpl/comment/list/primarySingleRow/hottest?outBizNo=157369&outBizType=basketball_item&clientCode=

Pretty print ✅

        "score": 10,
        "hasLight": false,
        "hasBlack": false,
        "publishTime": 1747102346338,
        "chosenTags": null,
        "ipLocation": "安徽",
        "subCommentCount": 51,
        "descendantCount": 93,
        "lightScore": 0.9991,
        "subCommentList": [
            {
                "commentKey": {
                    "subjectId": "272404212",
                    "commentId": "2179190340",
                    "key": "272404212:2179190340"
                },
                "subjectId": "272404212",
                "commentId": "2179190340",
                "commentUserId": 111663660,
                "commentUserName": "无名者杜哥",
                "commentUserHeadImg": "https://i2.hoopchina.com.cn/user/111663660/1730709037993.jpg",
                "commentContent": "搏至无憾了，伤病远离",
                "commentContentImages": null,
                "commentAncillaryContents": null,
                "parentCommentCanSee": true,
                "parentCommentDeleteFlag": false,
                "parentCommentUserId": 19682197,
                "parentCommentId": "2179520173",
                "parentCommentUserName": "YyYyWhyYyYy",
                "parentCommentUserHeadImg": "https://i3.hoopchina.com.cn/user/197/19682197/19682197.jpg",
                "parentCommentContent": "这场没得黑，恩无大碍，respect",
                "parentCommentContentImages": [],
                "parentCommentAncillaryContents": [],
                "commentDate": "05-13",
                "lightCount": 2613,
                "blackCount": 1,
                "score": 10,
                "hasLight": false,
                "hasBlack": false,
                "publishTime": 1747102363030,
                "chosenTags": null,
                "ipLocation": "山西",
                "subCommentCount": 13,

- Now, we have all the necessary responses needed to gather the information. We need to design a crawling script to generate all the **groupID**s of basketball matches, get the game stats and the individual **bizID** of the player, request for the json response, retrieve the top 3 comments of the player and finally plug them all into a .csv file.

The actual implementation using scrapy can be found here: https://github.com/alexgoexercise/hupu-comments-crawler.git
Have fun crawling!