

Web Scrapping Using IOT-Based Raspberry PI

Bhavya Shah
Department of Computer
Engineering and Techonlogy
Dr. Vishwanath karad MIT
World Peace University
Pune, Maharastra
bhavyashah16@outlook.com

Lakshya Upadhyaya
Department of Computer
Engineering and Techonlogy
Dr. Vishwanath karad MIT
World Peace University
Pune, Maharastra
lakshyau15@gmail.com

Gautam Sharma
Department of Computer
Engineering and Techonlogy
Dr. Vishwanath karad MIT
World Peace University
Pune, Maharastra
sharma.gautam0905@gmail.com

Samarth Patel
Department of Computer
Engineering and Techonlogy
Dr. Vishwanath karad MIT
World Peace University
Pune, Maharastra
samarthpatel@gmail.com

Dr. Vitthal Gutte
Department of Computer
Engineering and Techonlogy
Dr. Vishwanath karad MIT
World Peace University
Pune, Maharastra
vitthalgutte2014@gmail.com

Abstract—The proliferation of machine learning (ML) models has further driven the demand of massive and high-quality datasets, yet alternative datasets collection approaches may be cost- and resource-intensive. The current paper presents a novel, low-power Internet of Things (IoT) system to automatically scrape websites dynamically to generate ML-ready datasets, based on a limited resource edge-device. It is implemented on a Raspberry Pi edge node which gathers information automatically on dynamic websites using python packages Selenium and BeautifulSoup. This piece of work also features an active health monitoring device to facilitate resilience on prolonged working hours. The temperature of the CPU is being continuously measured and relayed over the lightweight MQTT protocol to the ThingSpeak cloud to be visualized and analyzed in real time and proves the effectiveness of passive thermal control provided by a heat sink. This two-fold design is an example of successful data-gathering of structured data types that can be trained in ML pipelines, at the same time the collection equipment is not damaged. The proposed framework confirms the possibility of low-cost edge devices to support more advanced data acquisition, and demonstrates a scalable, low-energy consuming, and allowable provisioning of researchers and developers to perform the development of custom datasets.

Index—Internet of Things (IoT), Edge Computing, Web Scrapping, Raspberry Pi, Dataset Creation, Machine Learning (ML), Thermal Management, MQTT

I. INTRODUCTION

Artificial Intelligence (AI) and Machine Learning (ML) have been widely used in a sudden manner that has mostly defined the phenomenon of the fourth industrial revolution. The models largely depend on the quantity and quality of the data used in training to be successful. This has led to the capacity to acquire succinct large domain-specific datasets becoming a problem that is challenging and has to be solved. Scraping data, especially at the web scale, can be highly expensive in terms of computing power and cost, which

can impede researchers and students to join the discussion in more traditional ways of scraping data.

This paper acknowledges the need to have an easier and more affordable means of not data acquisition in the broader framework of 'IoT' and helps! The technology of edge computing may meet those requirements. This paper presented a low power network that used a Raspberry Pi as a fully independent edge node that performs web scraping and that can be accessed via the internet. The system will be capable of extracting structured information on the web and capturing the data with the aid of Python and packages Selenium and BeautifulSoup that can be utilized in an ML pipeline.

However, the use of intensive apps in devices with limited resources is posing an essential operating issue to the limelight — thermal management. The constant operation of applications programs may lead to higher temperature of the CPU and lower performance and degradation of the hardware. In order to deal with this challenge, our framework employs a real-time health monitoring element. The temperature of the device is transferred by the MQTT protocol to the IoT cloud platform ThingSpeak, which makes it remotely available and gives the opportunity to confirm the influence of using an integrated passive heat sink.

The paper had the following significant contributions:

- 1) the architecture of an IoT system with dual use, which will enable data collection and device maintenance;
- 2) the proof of the creation of ML datasets using a low-cost edge device;
- 3) integration of a MQTT based monitoring system to facilitate operational stability.

II. LITERATURE REVIEW

Overview:

The high pace of Internet of Things (IoT) evolution has boosted the research of lightweight communication, edge-computing

system, and system-level management. The most common protocols such as MQTT, CoAP, DPWS, and XMPP are also briefly discussed including the explanation of QoS capabilities of MQTT, small header size, and publish/subscribe pattern, and absence of discovery of MQTT is explained [3]. These security analyses consist of authentication, authorization based on ACLs, the payload encryption, tradeoff of TLS/DTLS, and the link-layer encryption (LLSec) in which hardware acceleration is applied [3].

Raspberry pi as an Edge device:

Raspberry boards are presented as low-power, low-cost SBCs that could be utilized in edge processing to reduce bandwidth bandwidth of the cloud-based design [4]. Pi-based image processing, service deployment, PaaS clusters, and real-time lightning are several works that are taken into consideration. Real time measurements of a Pi3B+ over time show that there is constant CPU load (around 25 per cent to 30 per cent) and slight rise in temperature during real time signal processing, which can also be regarded as viability to flash-flood forecasting. Comparative studies of the different Pi models indicate that Pi 4B can do more than Zero -W can do little.

High-performance processors:

Thermal control:

Thermal Assist Unit (TAU) is an on-chip micro sensing, dual-threshold registration, and interrupt controller that offers an ability to control the temperature fine-grainedly [5]. Together with the instruction cache throttling and sleep modes TAU has an optimum 1 2 derating of 7 -C temperature at its peak performance, better than naive clock-scaling. Results of silicon result show process corner correction of -1 o C and power variance.

IoT-Based Digital Signage:

To deliver remote, context-sensitive advertisements and real-time information (traffic, weather), the signage system is driven by a Raspberry instead of microcontroller boards, a Node.js server, web APIs, and optional DHT11 sensors [2]. The option of remote control, by a web control panel, enables to plan and update the content in real time.

Application Domains (Agriculture, Transportation, Healthcare):

Pi-enabled IoT are deployed in surveyed literature to introduce such applications as intelligent transportation (driver aid, parking, road-condition monitoring, fleet tracking) [1], smart agriculture (greenhouse climate control, irrigation, plant fertilizer, etc.), and health care (vital-sign monitoring, non-invasive glucose monitoring, elderly support). The protocols are lightweight and have edge compute functionality and powerful thermal or security enforcement across these regions, which make up a single research environment of scalable IoT solutions.

III. SYSTEM ARCHITECTURE & METHODOLOGY

A. Hardware Architecture:

The hardware design of our proposed framework is a stand-alone and self-contained platform, with Raspberry Pi single-board computer at its centre. The entire hardware setup is shown in Fig. 1.

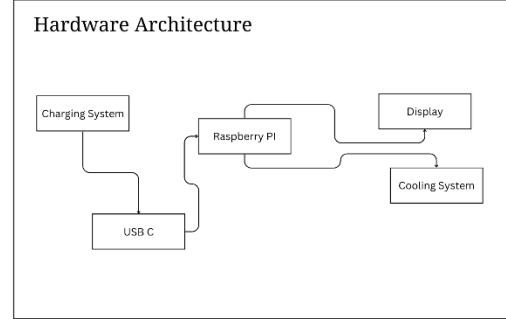


Fig 1. Hardware Architecture

The hardware is based on a Raspberry Pi, used as a hub of the central processing. A stable Charging System is used to power the system through USB-C so that it operates continuously. The Pi has two main peripherals that are controlled by the Pi: a Display, which is used as a user interface and as a real-time status indicator, and an active, GPIO-controlled Cooling System (e.g. a fan) to allow the Pi to dynamically manage thermal effects during high computational loads. This hard system is a small, all-in-one system with the ability to conduct independent data collection.

B. Software Architecture:

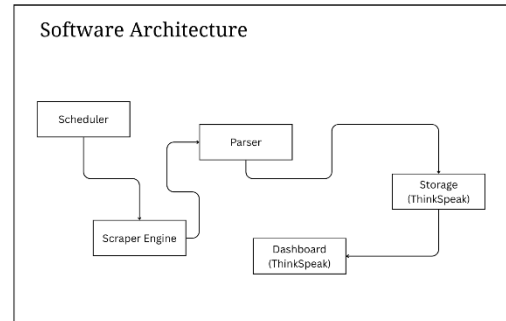


Fig 2. Software Architecture

The software architecture is presented in the form of a pipeline that is modular to automate the process of data acquisition and visualization as shown in Fig. 2. Scheduler activates the workflow, and the Scraper Engine is triggered by a Scheduler at the specific intervals. Access to a target site and fetching of the raw HTML content are done by the Scraper Engine, which is created using Selenium and BeautifulSoup. This raw data is

then forwarded to the Parser module whose functions is to process the HTML, get the data fields of interest, and organize the information into a clean format. Lastly, structured data is sent to the ThingSpeak cloud platform that can be used in two ways: as a Storage solution to log time-series data and as a Dashboard to get real-time visualization and analyze the data received.

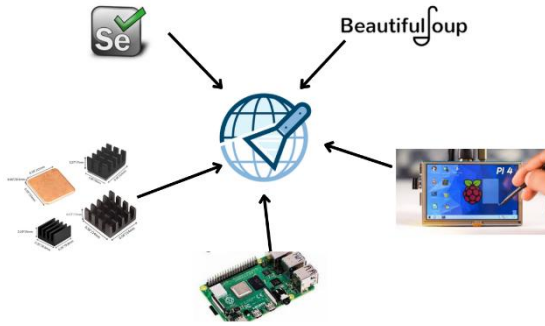


Fig 3. Graphical Abstract

C. Data Acquisition Module:

This module works on the removal of information on the target site. It is written in Python 3 and utilizes two important libraries:

1. **Selenium:** It is used to automate a web browser (Chromium). This is necessary in the management of modern and dynamic websites that use JavaScript to load pages. The script loads the target URL, executes activities such as scrolling or clicking, and gets the fully loaded page source.
2. **BeautifulSoup:** When the page source is retrieved with the help of Selenium, BeautifulSoup is applied to the HTML document to extract the HTML. It offers an effective means of traversing the parse tree to pick out particular data points by tagging, classes, and ID.

The obtained data (e.g., product names, prices, reviews) is processed, cleaned, and arranged into a Pandas DataFrame, which is further stored in the local disk as a CSV file. This format can be used across all types of machine learning systems such as Scikit-learn and TensorFlow.

D. Health Monitoring Module:

1. This will be a parallel module that operates as a different process or thread to the scraping task. Its only role is to provide the stability of the Raspberry Pi.
2. **Temperature Sensing:** This is a Python script that reads the internal temperature of the CPU of the Raspberry Pi periodically out of the system files under `/sys/class/thermal/thermal_zone0/temp`.

3. **MQTT Transmission:** The temperature value is delivered to a particular subject on a public or private MQTT server. The lightweight Paho-MQTT client library is used to do this. Publication is done after a predetermined time (e.g. after every 60 seconds) in order to prevent congestion on the network.
4. **Cloud Visualization** The ThingSpeak platform will be set up to subscribe to the MQTT topic. ThingSpeak is an IoT analytics service where it is possible to record and visualize data streams in real-time. An open channel is created and a field of temperature is put, which then automatically graphs the incoming data during the time

IV. IMPLEMENTATION & RESULTS

The framework was implemented and tested to validate its two primary functions. The target for the web scraping task was a mock e-commerce site with dynamically loaded product listings.

A. Data Acquisition Results:

The scraping program was run over a time period of one hour. The system was able to navigate and identify the data required fields in several pages and put them together in a structured CSV file. A very small sample of the resulting dataset is presented in Table I, and it indicates that it is well-organized and can be used in the purpose of the ML.

Title	Description	Developer(s)	Publisher(s)	Release	Genre(s)	Mode(s)
30XX	Following c Battery	stap	Batterystap	Windows V	Rogue-lite	Single-player , multiplayer
.hack//G.U.	This is part of the list of	Nintendo	Switch	games.		
Achilles: Le	Players con	Dark Point	Dark Point	Windows, P	Action role-	Single-player
Ad Infini	tun Players con	Hekate	Nacon	14-Sep-23	Survival hor	Single-player
Advance W	Advance W	WayForward	Nintendo	21-Apr-23	Turn-based	Single-player , multiplayer
AEW Fight	f AEW Fight	F Yuke's	THQ Nordic	29-Jun-23	Sports	Single-player , multiplayer
AFL 23	AFL 23 is a	2 Big Ant	Stur Nacon	Windows, P	Sports	Single-player , multiplayer
After Us	Players con	Piccolo Stur	Private Divi	WW : May	Action-adve	Single-player
Against the	The game is	Eremita Gar	Hooded Hoi	Windows V	City-buildin	Single-player
Age of Won	At the end c	Triumph St	Paradox Int	02-May-23	4X turn-bas	Single-player , multiplayer
Aimlabs	Aimlabs, for	State Space	State Space	16-Jun-23	Shooter	Single-player
Akka Arrh	It was relea	Llamasoft	Atari	February 21	Action	Single-player
Night Sprin	Alan Wake	Remedy Ent	Epic Games	27-Oct-23	Survival hor	Single-player

Fig 4. Sample of Scraped Data

B. Dashboard Visualization Result:

In order to show a real-world implementation of the data acquisition framework in action, the system was implemented to scrape an open video game database. Its aim was to gather data on new game titles and their developers that have been released recently and create a dataset to analyze them. Once the data was gathered and organized, a primary analysis was conducted in order to define the most prolific developers in the data.

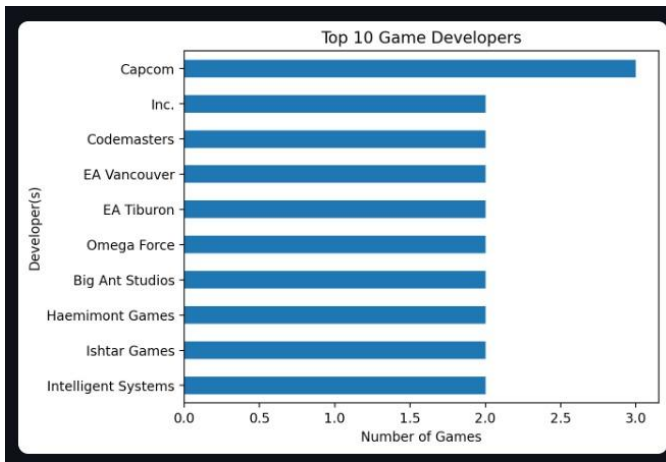


Fig 5. Bar Chart of Top 10 Game Developers from Scraped Dataset

V. CONCLUSION

This article developed a dual-purpose low-power IoT solution that was designed, implemented, and tested to create machine learning datasets. Our IoT system is based on a Raspberry Pi, an edge node, demonstrating that it can be possible to engage in advanced web scraping and data shaping with limited infrastructure needs. Also, a health monitor utilizing MQTT in the real-time, context which can offer significant safeguards and allow the device to continue throughout extended periods of the task can be implemented online. This paper suggests a low-power, scalable and low-cost solution to formalizing a significant step in the machine learning life cycle.

REFERENCES

- [1] Hosny, Khalid M., et al. "Internet of things applications using Raspberry-Pi: a survey." *International Journal of Electrical & Computer Engineering* (2088-8708) 13.1 (2023).
- [2] Alase, Swapnil, and Vaibhavi Chinchur. "IoT based digital signage board using Raspberry Pi 3." *International Research Journal of Engineering and Technology* 4.05 (2017): 310-313.
- [3] S. Katsikeas et al., "Lightweight & secure industrial IoT communications via the MQ telemetry transport protocol," 2017 IEEE Symposium on Computers and Communications (ISCC), Heraklion, Greece, 2017, pp. 1193-1200, doi: 10.1109/ISCC.2017.8024687. keywords: {Protocols; Encryption; Wireless sensor networks; Payloads; Communication system security; Wireless communication; Industrial Internet of Things; IIoT; MQTT; information security; confidentiality; authentication; secure communication; Wireless Sensor Networks; WSN; Industrial Networks},
- [4] Hassan, Aslinda, et al. "Performance evaluation of raspberry pi as an iot edge signal processing device for a real-time flash flood forecasting system." *International Journal of Advanced Computer Science and Applications* 13.10 (2022).
- [5] H. Sanchez et al., "Thermal management system for high performance PowerPC/sup TM/ microprocessors," *Proceedings IEEE COMPCON 97. Digest of Papers*, San Jose, CA, USA, 1997, pp. 325-330, doi: 10.1109/COMPCON.1997.584744. keywords: {Thermal management; Energy management; Power system management; Portable computers; Process design; Power dissipation; Costs; Microprocessors; Thermal sensors; Logic},
- [6] Thomas Zachariah, Noah Klugman, and Prabal Dutta. 2023. ThingSpeak in the Wild: Exploring 38K Visualizations of IoT Data. In *Proceedings of the 20th ACM Conference on Embedded Networked Sensor Systems (SenSys '22)*. Association for Computing Machinery, New York, NY, USA, 1035–1040. <https://doi.org/10.1145/3560905.3567766> M. Young, *The Technical Writer's Handbook*. Mill Valley, CA: University Science, 1989.
- [7] Thapliyal, Aditya, and C. Kumar. "Development of data acquisition console and web server using Raspberry Pi for marine platforms." *International Journal of Information Technology and Computer Science* 8 (2016): 46-53.
- [8] Linghe Kong, Jinlin Tan, Junqin Huang, Guihai Chen, Shuaitian Wang, Xi Jin, Peng Zeng, Muhammad Khan, and Sajal K. Das. 2022. Edge-computing-driven Internet of Things: A Survey. *ACM Comput. Surv.* 55, 8, Article 174 (August 2023), 41 pages. <https://doi.org/10.1145/3555308>
- [9] Nettikadan, David, and Subodh Raj. "Smart community monitoring system using Thingspeak IoT platform." *International Journal of Applied Engineering Research* 13.17 (2018): 13402-13408.
- [10] K. M P and D. N R, "Crop Prediction Based on Influencing Parameters for Different States in India- The Data Mining Approach," 2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2021, pp. 1785-1791, doi: 10.1109/ICICCS51141.2021.9432247. keywords: {Temperature dependence; Temperature; Vegetation; Soil; Prediction algorithms; Agriculture; Classification algorithms; crop; climate; pesticides; fertilizers; groundwater; soil; Agriculture},