



# Design Document

Hybrid IRCTC Bot-Detection And Train Search Web App

Prepared by: Subradeep Das

11th August 2025

# Contents

<b>1</b>	<b>Executive Summary</b>	<b>1</b>
<b>2</b>	<b>Scope and Goals</b>	<b>2</b>
2.1	Purpose . . . . .	2
2.2	Primary goals . . . . .	2
<b>3</b>	<b>High-level Architecture</b>	<b>3</b>
3.1	Overview . . . . .	3
<b>4</b>	<b>Component Design</b>	<b>4</b>
4.1	Frontend (Browser) . . . . .	4
4.1.1	Responsibilities . . . . .	4
4.1.2	Behavioral payload (example) . . . . .	4
4.1.3	UI mockups (selected) . . . . .	5
4.2	Backend: API Gateway (Flask) . . . . .	5
4.2.1	Primary endpoints . . . . .	5
4.2.2	Login sequence . . . . .	5
4.3	Behavioral Classifier (TensorFlow) . . . . .	6
4.3.1	Input / Output . . . . .	6
4.3.2	Suggested lightweight architecture . . . . .	6
4.3.3	Preprocessing . . . . .	6
4.4	Ticket Confirmation Predictor (Random Forest) . . . . .	6
4.4.1	Features . . . . .	6
4.4.2	Prediction contract . . . . .	6
<b>5</b>	<b>Data Model and Database Design</b>	<b>7</b>
5.1	Key entities . . . . .	7
5.2	Recommended schemas (Postgres) . . . . .	7
5.3	Indexes and performance . . . . .	8
<b>6</b>	<b>Data Flow Diagrams</b>	<b>9</b>
<b>7</b>	<b>API &amp; Integration Contracts</b>	<b>10</b>
7.1	POST /api/login . . . . .	10
7.2	Admin endpoints . . . . .	10
<b>8</b>	<b>Model Training, Retraining and MLOps</b>	<b>11</b>
8.1	Datasets . . . . .	11
8.2	Retraining policy . . . . .	11
8.3	Model artifacts and metadata . . . . .	11

<b>9 Deployment and Operations</b>	<b>12</b>
9.1 Containerization . . . . .	12
9.2 Suggested process . . . . .	12
<b>10 Security and Privacy</b>	<b>13</b>
10.1 Authentication and secrets . . . . .	13
10.2 Privacy . . . . .	13
<b>11 Performance and Scalability</b>	<b>14</b>
11.1 Targets . . . . .	14
11.2 Optimization suggestions . . . . .	14
<b>12 Testing and Validation</b>	<b>15</b>
12.1 Test types . . . . .	15
<b>13 Admin Dashboard and Monitoring</b>	<b>16</b>
13.1 Dashboard features . . . . .	16
13.2 Monitoring . . . . .	16
<b>14 Risks and Mitigations</b>	<b>17</b>
<b>15 Appendices</b>	<b>18</b>
15.1 Operational files and snippets . . . . .	18
15.1.1 Requirements (example) . . . . .	18
15.1.2 Useful artifacts (included images) . . . . .	18
15.2 Sample graphs . . . . .	19

# List of Figures

3.1	System workflow overview . . . . .	3
4.1	Frontend screens and flows . . . . .	5
6.1	Detailed data flow showing feature aggregation and model inference . . . . .	9
15.1	Results with predictions . . . . .	19
15.2	Demand heatmap — example 1 . . . . .	19

# **Chapter 1**

## **Executive Summary**

This document describes the system architecture, components, data flows, API contracts, database design, model training and retraining policy, deployment plan, security considerations, and testing strategy for the Hybrid IRCTC Bot-Detection and Train Search Web App.

# **Chapter 2**

## **Scope and Goals**

### **2.1 Purpose**

Provide a detailed, implementable blueprint for the system so developers and reviewers can understand how each module is built and integrated.

### **2.2 Primary goals**

- Reliable bot detection combining CAPTCHA and behavioral analysis.
- Fast and informative train search with probability estimates for confirmation and punctuality.
- Clear retraining and model governance process.
- Practical deployment and operational guidance with emphasis on security and privacy.

# Chapter 3

## High-level Architecture

### 3.1 Overview

The system is split into:

- Frontend: Login (behavior capture + CAPTCHA), Train Search UI, Admin Dashboard.
- Backend: API gateway (authentication, feature aggregation, inference orchestration).
- Inference services: Behavioral classifier (TensorFlow) and Ticket probability predictor (Random Forest).
- Persistence: PostgreSQL (recommended), model artifacts and logs on disk/object storage.
- Background workers: Celery + Redis for retraining, batch jobs, and reporting.
- Edge: Nginx for TLS termination, static asset serving, and basic rate-limiting.

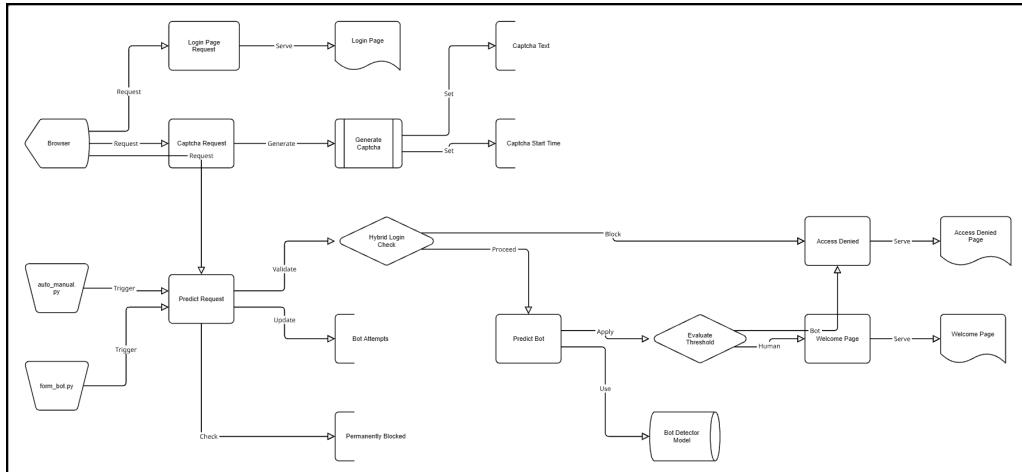


Figure 3.1: System workflow overview

# Chapter 4

## Component Design

### 4.1 Frontend (Browser)

#### 4.1.1 Responsibilities

- Collect behavioral telemetry (mouse, keyboard, focus/blur, timing) on login page.
- Display image CAPTCHA and validate token from server.
- Provide train search UI with station suggestions, calendar, and result list showing probabilities and metrics.
- Expose admin pages for authorized users.

#### 4.1.2 Behavioral payload (example)

```
1  {
2      "session_id": "uuid-v4",
3      "events": [
4          {"type": "mousemove", "t": 120, "x": 234, "y": 105},
5          {"type": "keydown", "t": 350, "key": "a"},
6          {"type": "keyup", "t": 370, "key": "a"}
7      ],
8      "summary": {
9          "mouse_total_distance": 523.2,
10         "avg_mouse_velocity": 120.3,
11         "typing_cpm": 210,
12         "focus_changes": 2,
13         "captcha_time_ms": 2350
14     }
15 }
```

### 4.1.3 UI mockups (selected)

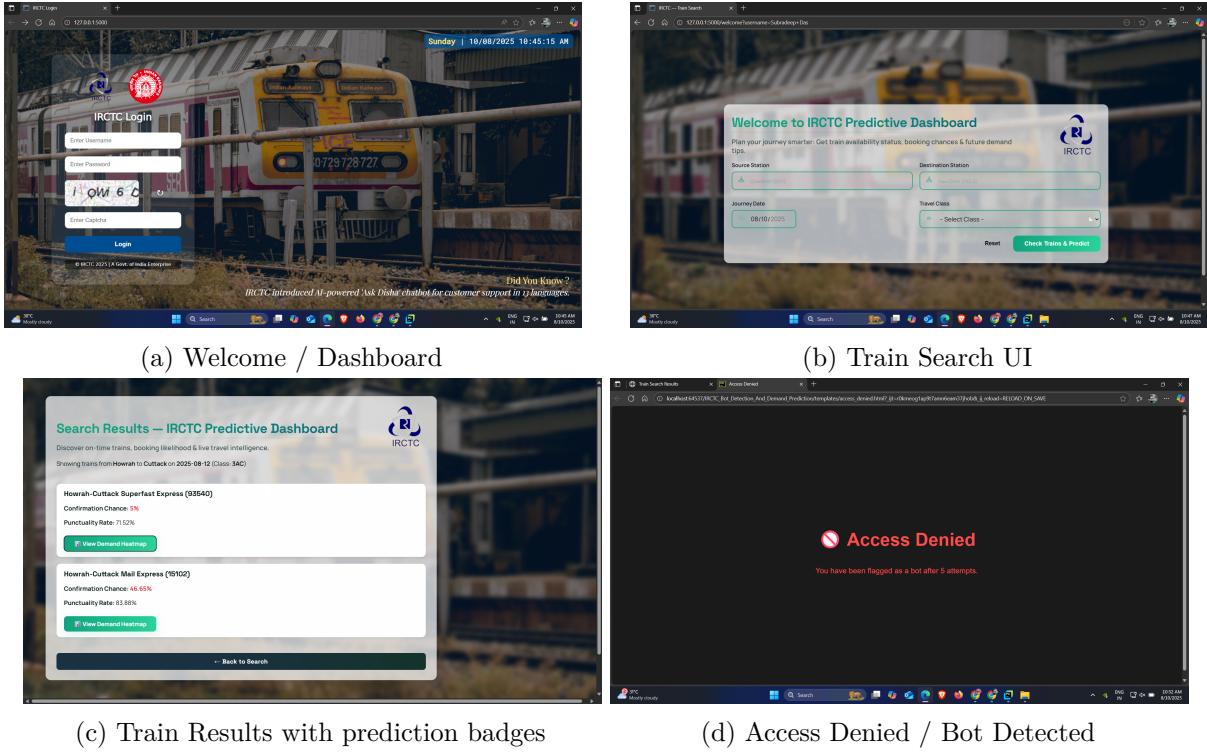


Figure 4.1: Frontend screens and flows

## 4.2 Backend: API Gateway (Flask)

### 4.2.1 Primary endpoints

Endpoint	Description / Payload
POST /api/login	<code>{email, password, captchaToken, behavioral} → returns {status, label, confidence, session}</code>
GET /api/search	Query params: source, dest, date, class → returns list of trains + prediction objects
POST /api/retrain	Admin-only trigger; optional force.
GET /api/admin/logs	Admin-only paginated logs.

### 4.2.2 Login sequence

1. Verify CAPTCHA server-side.
2. Aggregate behavioral events into features.
3. Call behavioral classifier for label and confidence.
4. Log attempt (JSONB) and decide flow: allow or deny (with human-review fallback).

## 4.3 Behavioral Classifier (TensorFlow)

### 4.3.1 Input / Output

Input: aggregated session-level features (and optional short sequences). Output: class probabilities for {human, bot, human-like-bot}.

### 4.3.2 Suggested lightweight architecture

Dense network for low-latency inference:

Input (N) → Dense(128, ReLU) → Dropout(0.3) → Dense(64, ReLU) → Dense(3, Softmax)

### 4.3.3 Preprocessing

- Winsorize extreme outliers.
- Z-score normalization using stored mean/std.
- Persist scaler and feature list alongside model.

## 4.4 Ticket Confirmation Predictor (Random Forest)

### 4.4.1 Features

days\_before\_journey, day\_of\_week, month, past\_avg\_waitlist, past\_punctuality, class, quota, encoded station features, train\_no.

### 4.4.2 Prediction contract

Return calibrated probability and model metadata:

```
1 { "train_no": "12345", "confirmation_prob": 0.71, "model_version": "  
rf_v1.2" }
```

# Chapter 5

## Data Model and Database Design

### 5.1 Key entities

- **users** — registered users and roles.
- **login\_attempts** — behavioral payloads, IP, UA, label, model\_version.
- **trains** — train metadata and historical stats.
- **predictions** — runtime predictions with features and model metadata.

### 5.2 Recommended schemas (Postgres)

```
1 CREATE TABLE users (
2     user_id SERIAL PRIMARY KEY,
3     name TEXT,
4     email TEXT UNIQUE NOT NULL,
5     password_hash TEXT NOT NULL,
6     role TEXT DEFAULT 'user',
7     created_at TIMESTAMP DEFAULT now()
8 );
9
10 CREATE TABLE login_attempts (
11     attempt_id SERIAL PRIMARY KEY,
12     user_id INTEGER REFERENCES users(user_id),
13     timestamp TIMESTAMP DEFAULT now(),
14     ip TEXT,
15     user_agent TEXT,
16     features_json JSONB,
17     label TEXT,
18     model_version TEXT
19 );
20
21 CREATE TABLE trains (
22     train_no TEXT PRIMARY KEY,
23     train_name TEXT,
24     source TEXT,
25     destination TEXT,
26     departure_time TIME,
27     arrival_time TIME,
28     punctuality_rate REAL
29 );
30
```

```
31 CREATE TABLE predictions (
32   prediction_id SERIAL PRIMARY KEY,
33   train_no TEXT REFERENCES trains(train_no),
34   timestamp TIMESTAMP DEFAULT now(),
35   features_json JSONB,
36   predicted_prob REAL,
37   model_version TEXT
38 );
```

### 5.3 Indexes and performance

Index on `login_attempts(timestamp)`, partial index for denied attempts, and indexes on `predictions(train_no)` and `predictions(timestamp)`.

# Chapter 6

## Data Flow Diagrams

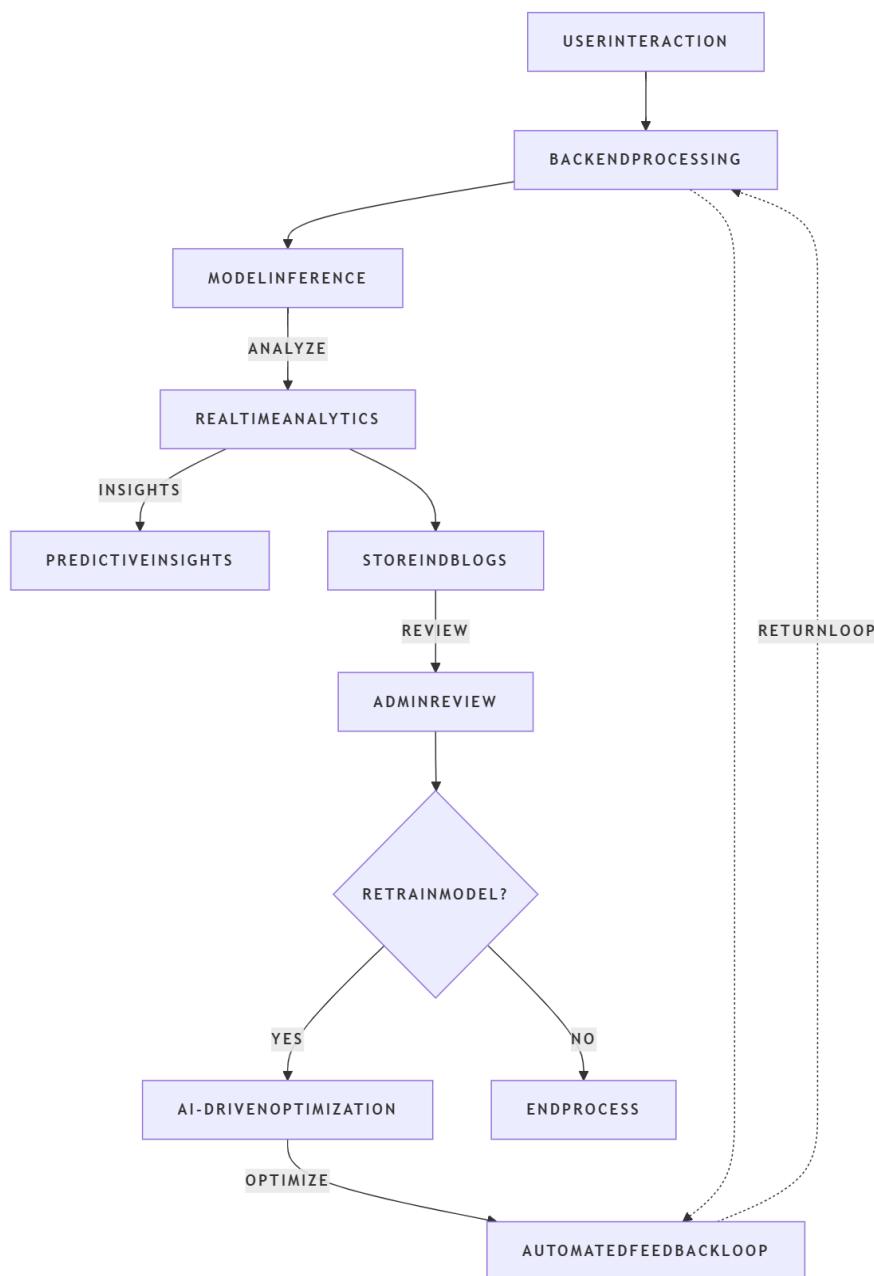


Figure 6.1: Detailed data flow showing feature aggregation and model inference

# Chapter 7

## API & Integration Contracts

### 7.1 POST /api/login

**Request:**

```
1 {  
2   "email": "user@example.com",  
3   "password": "strongpassword",  
4   "captcha_token": "...",  
5   "behavioral": { ... }  
6 }
```

**Response (allowed):**

```
1 { "status": "ok", "label": "human", "confidence": 0.98, "session": "sid  
- ..." }
```

**Response (denied):**

```
1 { "status": "denied", "label": "bot", "confidence": 0.92, "message": "  
Access denied" }
```

### 7.2 Admin endpoints

Require role-based authentication and IP allowlist: POST /api/retrain, GET /api/admin/logs, GET /api/models.

# Chapter 8

# Model Training, Retraining and MLOps

## 8.1 Datasets

- Behavioral dataset with aggregated features and labels.
- Historical ticket dataset with waitlist, booking, punctuality records.

## 8.2 Retraining policy

- Auto-trigger retrain when new login data reaches a threshold (e.g., +10 new entries since last train) and minimum totals satisfied.
- Promote a new model only if validation accuracy improves and validation loss decreases compared to the current production model.
- Archive prior models and log retraining metadata (version, samples, metrics).

## 8.3 Model artifacts and metadata

Store: model\_file (HDF5 or SavedModel), scaler (joblib), metadata.json (version, date, features, metrics).

# Chapter 9

# Deployment and Operations

## 9.1 Containerization

Use Docker to package application and inference components; orchestrate with docker-compose for dev and a container platform for prod.

## 9.2 Suggested process

- Nginx as reverse proxy for TLS and gzip/static assets.
- Gunicorn (WSGI) for Flask app; persistent TF serving or in-process model loaded at startup for low latency.
- Redis for caching and Celery broker.
- Postgres for structured persistence.

# **Chapter 10**

## **Security and Privacy**

### **10.1 Authentication and secrets**

Passwords stored with bcrypt/Argon2; sessions use secure cookies; admin endpoints restricted.

### **10.2 Privacy**

Avoid storing raw typed input. Persist aggregated features instead of raw event dumps whenever possible. Hash any PII before logs.

# **Chapter 11**

# **Performance and Scalability**

## **11.1 Targets**

- Behavioral model inference: target under 500 ms.
- Train search latency: under 2 seconds for baseline concurrency.

## **11.2 Optimization suggestions**

Warm models at service start, reuse inference process, cache common queries, tune Gunicorn worker count and DB pools.

## Chapter 12

# Testing and Validation

### 12.1 Test types

- Unit tests for feature aggregator and model wrapper.
- Integration tests for full login → search flow.
- Load tests for login and search endpoints, fuzz tests for malformed telemetry.

## **Chapter 13**

# **Admin Dashboard and Monitoring**

### **13.1 Dashboard features**

Model version and metrics, denied attempts summary, retrain controls, CSV export for audit.

### **13.2 Monitoring**

Export metrics to Prometheus and set alerts for latency spikes, inference errors, and model regression events.

## Chapter 14

# Risks and Mitigations

Risk	Mitigation
Behavioral spoofing	Multi-signal approach (CAPTCHA + behavior), continual retraining, threshold tuning
False positives blocking humans	Conservative initial thresholds, human-review flow, clear logging and rollback
Model drift	Time-aware retraining and monitoring for metric degradation
Data privacy breaches	Minimize retention of raw events, hashing, and access controls

# Chapter 15

## Appendices

### 15.1 Operational files and snippets

#### 15.1.1 Requirements (example)

```
1 Flask
2 gunicorn
3 tensorflow
4 scikit-learn
5 joblib
6 pandas
7 numpy
8 celery
9 redis
10 psycopg2-binary
11 bcrypt
```

#### 15.1.2 Useful artifacts (included images)

- Workflow.png — system workflow overview (used in architecture).
- Dataflow.png / Dataflow1.png — data flow diagrams.
- train\_search.png, train\_results.png, welcome\_page.png — frontend mockups.
- demand\_graph1.png, demand\_graph2.png — sample demand graphs used in analysis.
- access\_denied.png — access-denied page mockup.

## 15.2 Sample graphs

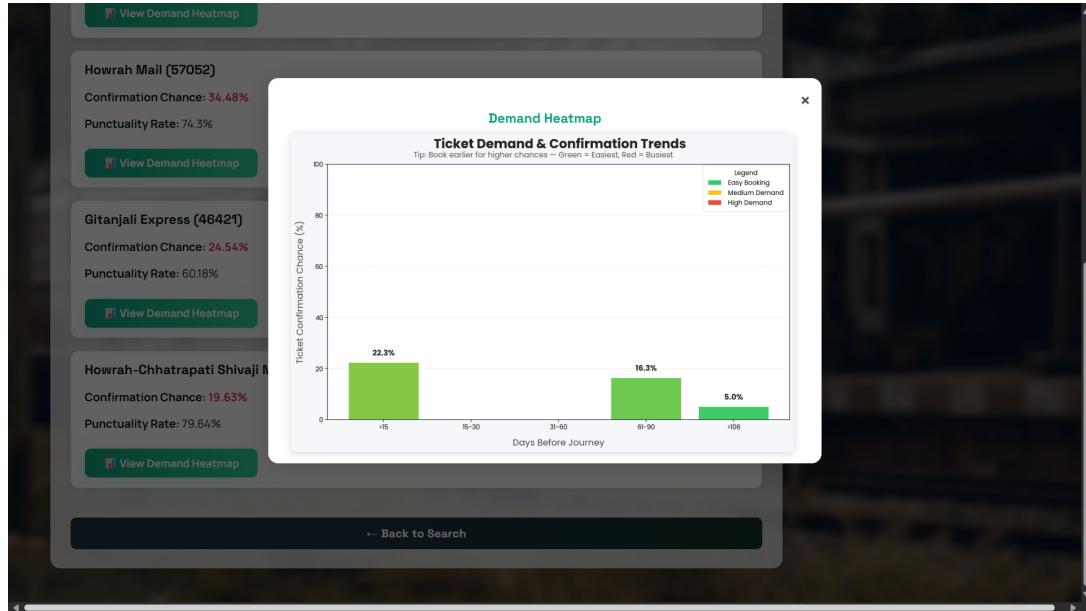


Figure 15.1: Results with predictions

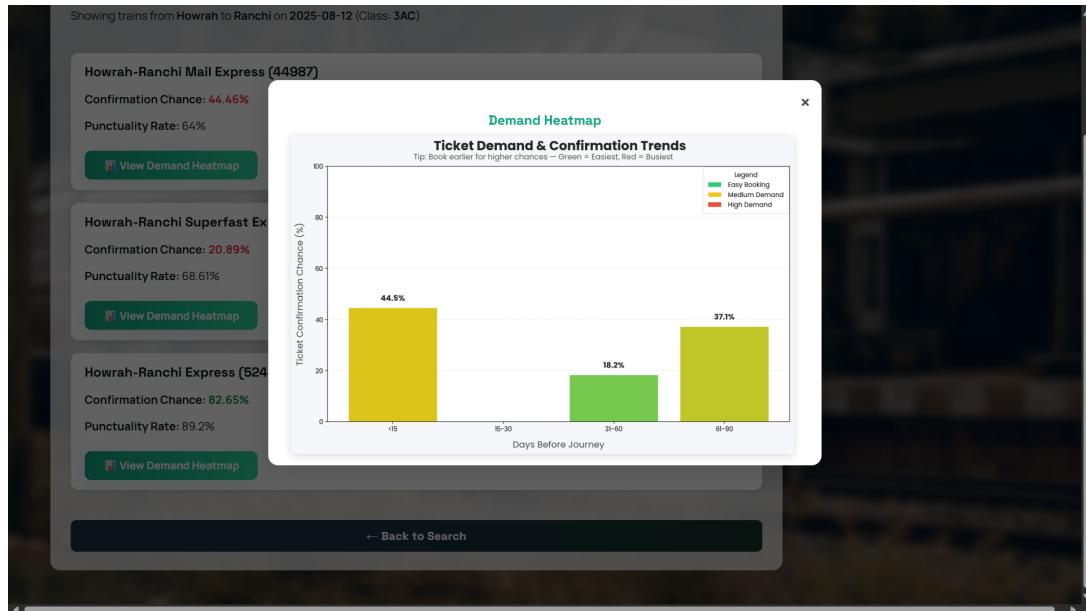


Figure 15.2: Demand heatmap — example 1

# Notes

This document captures the architecture and component design required to implement the Hybrid IRCTC Bot-Detection and Train Search Web App.